

IFAC



WARSZAWA 1969

INTERNATIONAL FEDERATION
OF AUTOMATIC CONTROL

Translations of papers from Russian into English

Fourth Congress of the International
Federation of Automatic Control
Warszawa 16–21 June 1969

Vol. II

TECHNICAL
SESSIONS
28 – 46



Organized by
Naczelna Organizacja Techniczna w Polsce

INTERNATIONAL FEDERATION OF AUTOMATIC CONTROL

Translations of papers from Russian into English

Vol. II

TECHNICAL SESSIONS 28 – 46

**FOURTH CONGRESS OF THE INTERNATIONAL
FEDERATION OF AUTOMATIC CONTROL
WARSZAWA 16 – 21 JUNE 1969**



**Organized by
Naczelna Organizacja Techniczna w Polsce**

C o n t e n t s

Paper No		Page
28.3	USSR - A.A.Pervozvansky - Decentralization Principle in Optimization of Complex Systems.....	5
28.5	USSR - A.I.Kukhtenko - On a Complex Systems Control Theory.....	17
29.1	USSR - I.V.Pranghishvily, V.V.Ignatushchenko - Construction of Checking and Diagnostic Procedures for General-Purpose Uniform Arrays....	35
29.2	USSR - A.F.Volkov, I.N.Vasil'ev, V.A.Vedeshenkov, V.A.Petrov - On Design of Tests for Digital Devices with Delays.....	52
29.3	USSR - P.P.Parkhomenko - The Theory of Questionnaires and Problems of Technological Diagnostics.....	66
29.4	USSR - A.L.Garcavi, V.B.Gogolevsky, V.P.Grabovezky - Effect of Monitoring Periodicity on Reliability of Restorable Devices.....	79
32.1	USSR - M.D.Klimovitsky, O.S.Kozhinsky, R.V.Lyambakh V.V.Naumchenko, A.B.Chelustkin - Digital Slab Tracking and Production Logging System for continuous Hot-Strip Mill.....	92
34.4	USSR - V.M.Kuntsevich, Yu.N.Chekhovoi - Fundamentals of Nonlinear Control Systems with the Pulse-Frequency and Pulse-width Modulation..	103

35.1	USSR	- A.Ya.Lerner, A.I.Teiman - On Optimal Resources Allocation.....	119
35.2	USSR	- V.Avdiysky, A.Voronov, S.Lovetsky - On Stock Control Theory.....	135
35.3	USSR	- Oleg.G.Tchebotarev - Resource Allocation in Multi-Project Based on Aggregation of the Project Networks.....	146
35.4	USSR	- V.N.Burkov - Optimal Project Control.....	158
35.5	USSR	- M.K.Badunachvili, D.I.Golenko, S.S.Naumov - - Some Questions of the testing and Construction Principles of an Optimum Multilevel Control Structure in Systems with a Specific Objective Function.....	166
36.1	USSR	- A.P.Shorygin - Electrochemical Transducers, Comparative Properties, Basic Characteristics and Fields of Application.....	175
39.1	USSR	- E.L.Suchanov, V.S.Shvidki, B.I.Kitaev, Ju.G.Yaroschenko, Ju.N.Ovchinnikov, V.G.Lisienko - Rational Algorithm of Controlling the Thermal Condition of Blast Furnace Using Computers.....	187
39.6	PL	- R.Górecki - Time Sub-Optimum Control of the Work of Cranes with Special Regard to Its Realization in Practice.....	204
41.3	USSR	- Ye.P.Popov, Ye.I.Khlypalo - An Extension of the Harmonic Linearization Technique.....	222
42.4	USSR	- A.P.Kopelovich, A.A.Belostotsky, B.A.Vlasjuk V.M.Khrupkin, G.I.Nikitin - Control Systems and Algorithms for a "Steel-Rolled Products" Manufacturing Complex of a Steel Works.....	232
43.1	USSR	- B.N.Petrov, N.P.Kolpakova, V.A.Vasilyev, A.I.Pavlenko - Considering Synthesis of Lifting Reentry Vehicle Control System Structures in Atmospheric Maneuver.....	249

43.3	USSR	- V.M.Ponomarev, V.I.Gorodezky - Optimal Parametric Control for The Re-Entry Space Vehicle	260
43.4	USSR	- A.G.Vlasov, E.I.Mitroshin, I.S.Ukolov - Stochastic Optimization of Spaceship Reentry Control in Atmosphere.....	272
43.5	USSR	- D.E.Okhotsimski, A.P.Bukharkina, Yu.F.Golubiev - Atmosphere Re-Entry Control Problem..	290
43.6	USSR	- Y.P.Plotnikov - Stochastic Problems of Missile Dynamics.....	303
43.7	USSR	- A.Ya.Andrienko - Statistical Synthesis of Optimal Pulse Control Systems with Regard to System's Structure Constraints.....	321
43.8	USSR	- A.A.Lebedev, M.N.Krasilshchikov, V.V.Malishov - Optimal Control System / For Stationary Artificial Circumterrestrial Satellite Orbit..	331
46.6	USSR	- R.Tavast, L.Mytus - An Adaptive Man-Computer Control System For a Chemical Plant.....	348

DECENTRALIZATION PRINCIPLE IN OPTIMIZATION OF COMPLEX SYSTEMS

PERVOZVANSKY A.A.

Polytechnical Institute, Leningrad, USSR.

A paper consists of two parts. In the first part there is considered the possibility of reduction of optimal planning problem for systems of reasonably general structure to the series of extremum problems for separate elements. An influence of reciprocal supplies on the efficiency of the work of every element is studied and the simple local approximation for these functions is proposed. On this basis in the second part an influence of supplies irregularity during planning period on the average efficiency of the work of an element is analyzed with provision for the possibility of optimization of inventory policy. It leads to certain conclusions on the necessity of correction of the statement of a planning problem for a system as a whole.

I. A problem of a planning of a work of a system combined of interactive elements is considered. A structure of interaction may be defined by the incidence matrix of connections graph, i.e. for every i it is possible to define a set of indices j of the inputs $J_1(i)$ or a set of output indices $J_2(i)$. Define also an element $n+1$ (graph node), which has no arcs coming out of it. A performance of every element i characterized by its output P_i distributed to connected elements

$$P_i = \sum_{j \in J_2(i)} P_{ij}'' \quad (I.1)$$

such that

$$\{u_i; P_i; P_{ji}'\} \in R_i \quad (I.2)$$

where u_i - a vector characterizing an intensity of the processes in an element itself, P_{ji}' , $j \in J_1(i)$, - vectors characterizing a part of other elements output consumed by the element i , R_i - some closed sets.

Note that in some cases it is convenient to separate the constraints into two groups: system constraints characterized by the output consumption levels of system elements and own (local) constraints characterizing the possibilities of a given element only.

It is assumed that a performance of a system as a whole is governed by an extremum principle expressed in additive form

$$F = \sum_{i \in J(n+1)} F_i(P''_{i,n+1}) \quad (I.3)$$

and by the conditions of balance

$$P'_{ij} = P''_{ij}, \quad (I.4)$$

Extremum problem just formulated will be called optimal deterministic planning for a system. It is a mathematical programming problem and it represents a structural generalization of classical schemes given by L.V. Kantorovich and D. Gale which are used in econometrics. It may be interpreted as a problem of a choice of output and distribution levels in the system combined of series of partially connected productive units (shops, enterprises, firms and so on). Although a time factor does not occur explicitly in a model, it is possible to suggest that given constraints determine a behaviour of each element of a system at a certain finite interval of time (planning subperiod). This interval may not be identical for different elements. It makes possible to consider so called determined perspective planning models although in this case one productive unit must be represented by several elements in a scheme corresponding to its state at successive planning subperiods.

After construction of a planning problem it may be solved at least in principle, by any known method of mathematical programming. However, special features of a model as a complex system impose essential limitations on the choice of a method. These limitations are of two types:

- a) informational, connected with difficulties of collection and keeping of the data in one place (center).
- b) computing, caused by practical inability of modern computers to operate with problems of very big volume.

For this reason, in this paper an attention is paid exclusively to the methods which use a decentralization principle. It means that the main problem must be separated into n extremal

subproblems, independent in a sense that direct information about constraints on every element of a system is not required for their coordination.

In another treatment of the principle an additional requirement is put forward, namely, that a subproblem coordination may be done not in a single center but in several ones, each center possessing partial information only.

Mathematical treatment of an idea of decentralization has arisen (although in a limited form) on a basis of works by Arrow, Hurwicz, Uzawa on gradient schemes of search of extremum ^{1,2}, then got a new impulse when the decomposition algorithm was created by Dantzig and Wolfe ³, and now is widely discussed, mainly in connection with different schemes of so called block programming (^{4,8} and others).

It is possible to mark out (omitting some essential mathematical details) two different approaches on subproblem forming.

The first one has a general character and is based on utilization of a main theorem by Kuhn and Tucker on equivalence (with certain limitations) of a given problem and a problem of a search of a saddle point of Lagrange function

$$\varphi = \sum_{i \in J_1(n+1)} F_i(P''_{i,n+1}) + \sum_i \sum_j \lambda_{ij} (P''_{ij} - P'_{ij}). \quad (I.5)$$

at the conditions (I.1-I.2). Then the problem is separated into n subproblems for the elements

$$f_i = \max \left[F_i(P''_{i,n+1}) + \sum_{j \in J_2(i)} \lambda_{ij} P''_{ij} - \sum_{j \in J_1(i)} \lambda_{ji} P'_{ji} \mid P = \sum P''_{ij}; \{u_i; P'_{ji}; P_i\} \in R_i \right] \quad (I.6)$$

where "prices" λ_{ij} are considered as parameters, and one problem for a "center", which is a problem of optimal coordination of "prices"

$$\varphi_0 = \min \left\{ \varphi = \sum_{i=1}^n f_i \left[\lambda_{ij}, j \in J_2(i); \lambda_{ji}, j \in J_1(i) \right] \right\} \quad (I.7)$$

This scheme may be called a scheme of "buy-sell" of intermediate products, the subproblems are treated as those of a profit maximization for every element, the prices on intermediate products are coordinated by the "center".

In the second approach it is assumed that every element of a system participates in forming of final products (or, fi-

nal expenses), i.e. a set $J_1(n+1)$ includes all $i=1, 2, \dots, n$. Then, considering the levels of distributed products as parameters, one comes to n subproblems of the type

$$F_i = \max [F_i(P_{i,n+1}'') / \sum_{ij} P_{ij}'' = P_i; \{u_i, P_i, P_{ji}'\} \in R_i] \quad (I.8)$$

and one problem for a "center" which is a problem of optimal balance of intermediate products:

$$F = \max \left\{ F = \sum_i F_i \left[P_{ji}', j \in J_1(i); P_{ij}'', j \in J_2(i) \right] \right\} \quad (I.9)$$

This scheme may be called a scheme of resources distribution, or a scheme of "optimal balance". In this case an optimization of final products value is made inside of every element, the levels of reciprocal supplies are given by the "center".

The partitions described are purely formal. In fact, there is no way to specify the "prices" (objectively conditioned estimates in L.V. Kantorovich's terminology) or the levels of intermediate production, which give an optimal coordination of subproblems, except the construction of a solution of a problem as a whole. However, a possibility of utilization of iterative procedures, i.e. the procedures of gradual coordination, makes this approach efficient enough.

Almost all of the known procedures are the versions of generalized gradient descent method or the method of feasible directions applied to the problem of a "center". In application to (I.7) the main idea is reduced to the following: an arbitrary set of "prices" is chosen by a "center" (in practical problems, the choice, naturally, is dictated by some practical considerations), and each element is informed about the choice. Then a problem of a type (I.6) is solved for every element, and this solution provides at given $\lambda_{ij} = \lambda_{ij}^0$ optimal values of activities u_i , supplies P_{ij}' , $j \in J_2(i)$, and external (for the element in question) consumption. In addition to it, a tendency of change of purpose function f_i due to small deviations of the "prices" from the level λ_{ij}^0 is revealed, i.e. a local approximation $f_i(\lambda_{ij})$ in some small vicinity λ_{ij}^0 is constructed.

Just the local approximation is sent to a "center", where on the basis of this information a tendency of a change of

function φ as a whole is revealed, that makes it possible to find necessary direction of a change of "price" system. So far as the step in the chosen direction is accomplished, and a set is got of new values of $\lambda_{ij} = \lambda_{ij}^1$, a procedure, naturally, can be repeated. A choice of the way of local approximation and size of a step in the direction of decrease is a speciality of a particular method.

Note only some general features arising in a case when original problem, and hence - subproblems, are formalized as linear programming problems. Here a local approximation of the behaviour of functions f_i , which coincides exactly with a real behaviour in certain finite vicinity, is as follow

$$\hat{f}_i = f_i(\lambda_{ij}^0) + \min_{s \in S_i^0} \left\{ \sum_{j \in J_2(i)} \Delta \lambda_{ij} P_{ij,s} - \sum_{j \in J_1(i)} \Delta \lambda_{ji} P'_{ji,s} \right\} \quad (I.10)$$

where S_i^0 - a set of optimal basises of a problem (I.7) at $\lambda_{ij} = \lambda_{ij}^0$, and $P_{ij,s}, P'_{ji,s}$ - optimal values of the variables P_{ij}'' , corresponding to these basises. Hence, in particular, it is clear that the functions $f_i(\lambda_{ij})$, generally speaking, are not differentiable. If the solution (I.7) is unique at $\lambda_{ij} = \lambda_{ij}^0$, then a gradient $f_i(\lambda_{ij})$ exists at the point and may be readily found as soon as an optimal solution is known. If all the problems of a type (I.7) have a unique solution then a calculation of a gradient of a function φ as a whole is reduced to calculation of debalances on connections:

$$\left. \frac{\partial \varphi}{\partial \lambda_{ij}} \right|_{\lambda_{ij} = \lambda_{ij}^0} = P_{ij, \text{opt}}'' - P_{ij, \text{opt}}' \quad (I.11)$$

Noted feature essentially facilitates an organization of a gradient descent. Moreover, it then follows a principally important result: a determination of the direction of a change of "prices" may be done without participation of a united center, but by coordination of the results of planning of directly connected elements. At the same time, it is obvious that a condition of uniqueness is not always fulfilled, it is certainly violated at the extremum point of a function φ , a consequence of this, being an absence of strict convergence of gradient descent at constant coefficients of proportionality Γ . However, there are some ways, a decrease of coefficients with a grow of its iteration number, for example, that permit to avoid the troubles

noted. Moreover, when solving practical problems a requirement of strict convergence is not very essential, it would be more important to increase a speed of motion, that can be done, for example by utilization of finite-step schemes with a use of simplified ways of local approximation¹⁰. Note also, that all that was said before may be completely applied to the second version of decentralization principle. Here the schemes of generalized gradient descent may also be used for the solution of the problem of a "center", but in a space of variables P'_{ij}, P''_{ji} . The functions F_i are also piecewise linear if original problem is a problem of linear programming. Their local approximation may be constructed through the use of optimal basis solutions of the problems dual to (I.8). If these solutions are unique, i.e. the problems (I.8) at given values of parameters are not degenerated, then it is possible to find the gradients of functions F_i with the help of optimal dual variables (objectively conditioned estimates). In the vicinity of an optimum only piecewise linear locally-exact approximation of the following type is permitted

$$\hat{F}_i = F_i [P''_{ij, \text{opt}}; P'_{ji, \text{opt}}] + \min_{\ell \in \ell_i} \left\{ \sum_j \lambda'_{ji, \ell} \Delta P'_{ji} - \sum_j \lambda''_{ij, \ell} \Delta P''_{ij} \right\} \quad (\text{I.12})$$

where ℓ_i denote a set of optimal bases of a dual problem and $\lambda'_{ji, \ell}, \lambda''_{ij, \ell}$ - corresponding values of dual variables. Iterative procedures can be done in a manner described earlier for the first version of separation of a problem, this time, however, the conditions of a balance on connections are fulfilled exactly in every iteration, that makes its realization more complicated but increases a validity of constructed approximation, every of which represents a suboptimal plan of a problem as whole.

Note that a plan (program) is usually treated as a specification of final production and reciprocal supplies levels for every element of a system, these levels specified integrally at the planning period. However, due to inevitable fluctuations of conditions an exact realization of the plan turns out to be impossible, therefore in carrying out a control, i.e. in realization of a plan, a question of compensation of deviations from the program naturally arises. Since there are many plan indices, it is important to estimate the influence of their individual

deviations on the efficiency of a system. In practice, this estimate is done subjectively. At the same time, decentralization principle permits not only to find optimal levels during construction of a plan, but also to estimate an efficiency of small deviations from these levels, because in the course of construction of a plan a local approximation of dependency of local purpose functions from external parameters is also constructed. It is essential that even in a small vicinity of optimal plan this approximation is nonlinear. This point predetermines impossibility of use of such "linear" criteria as summed value of deviations calculated in prices that do not depend on the level or direction of the deviations.

2. Let us now proceed to the analysis of influence of some fluctuative factors on the efficiency of realization of the program in the course of its accomplishment and of a back influence of these factors on a scheme of construction of a plan itself.

It is obvious that designation of a reciprocal supplies program integrally at the planning period does not specify a distribution of these supplies during this period. This uncertainty, peculiar to the method of planning, may lead (and practically leads) to a rise of irregularity in the level of supplies.

Consider certain element from a system and study a change of its final productions due to deviations in the level of supplies to this element from other elements. Divide a planning period in subperiods and enumerate them with the indices k ($k = 1, 2, \dots, K$). Assume that the resources of this element are not changed during a whole period and that its performance is governed by local extremal principle, formulated earlier, coordinated integrally with extremal principle for a system as a whole. Then at every subperiod a change of efficiency with respect to final production output is given by formula analogous to (I.12), assuming that the requirements on supplies from a given element are fulfilled and the supplies to the element have small deviations from prescribed level. This piecewise linear dependence may be efficiently approximated by means of a system of coefficients (marginal values of a problem, following a terminology of ^{II}) characterizing a change of function at increase or decrease of every of the components of supplied

production from optimal program.

Omitting indices, characterizing a number of an element and its suppliers, a local approximation may be given as follows ¹⁰:

$$\bar{F}_k = F_k(p_{k,opt}) + \sum_z \min(\lambda_z^+ \Delta p_{zk}; \lambda_z^- \Delta p_{zk}) = F_k(p_{k,opt}) + \sum_z \Delta F_{zk}^{(2.1)}$$

where λ_z^+, λ_z^- - accordingly right and left partial derivatives of purpose function with respect to z -th resource. It is possible to show that $\bar{F} \leq \hat{F}$. An advantage of an approximation (2.1) consists in a possibility to consider an influence of changes of every component separately. A difference among the components demonstrates itself in a difference of marginal values only. In what follows, the indices of resource components $p_{z,k}$ and these of corresponding changes $\Delta F_{z,k}$ of purpose function will be omitted and following notations introduced:

$$\Delta p_{z,k} \equiv v_k; \Delta F_{z,k} \equiv \psi(v_k) = \begin{cases} \lambda^+ v_k, & v_k \geq 0, \\ \lambda^- v_k, & v_k \leq 0. \end{cases} \quad (2.2)$$

Note also that always $\lambda^- \geq \lambda^+ \geq 0$. The values $v_k, k=1,2,\dots,K$, characterize a deviation of a quantity of a resource component by given element of a system from planned level, which is the $1/K$ -th part of planned supplies for a period as a whole. We may also introduce the values η_k , characterizing the deviations of real supplies during subperiod k from a given planned level. If an accumulation of resource surpluses is possible, then a choice of the value v_k is restricted by the quantity stored at a warehouse at the beginning of a subperiod and

$$v_k \leq y_k, \quad (2.3)$$

$$y_{k+1} = y_k + \eta_k - v_k, \quad k=1,2,\dots,K. \quad (2.4)$$

In this way a general problem of supplies irregularity is reduced to a classical onedimensional inventory problem (a particular case - water resource control was considered by Karlin and Gessford in ¹²). Both an analysis of a system as a whole and utilization of simplified local approximation (2.1) give a general character to this problem.

We shall now describe the concrete results of analysis with a use of optimal many-stage policy at the whole period and

the simplest policy $v_{k+1} = \eta_k$. Under these conditions, assume that η_k deviates from zero (planned level) with equal probability on the value $\pm a$, an average being equal to zero.

Then if $S_t(y)$ is a value of average loss due to irregularity at t subperiods to go to the end of a period, at the initial stock level y and optimal policy, then

$$S_t(y) = \min_{v \leq y} [-\phi(v) + \int S_{t-1}(y + \eta - v) w(\eta) d\eta],$$

$$S_1(y) = \min_{v \leq y} [-\phi(v)]. \quad (2.5)$$

It is possible to show that an optimal policy is as follows:

$$v_t = \begin{cases} 0 & ; y \geq 0 \\ y & ; y \leq 0 \end{cases}, \quad t = 2, 3, \dots; \quad v_t = y. \quad (2.6)$$

Then

$$S_t(y) = \begin{cases} \frac{1}{2} S_{t-1}(y+a) + \frac{1}{2} S_{t-1}(y-a); & y \geq 0 \\ -\lambda^- y + \frac{1}{2} S_{t-1}(a) + S_{t-1}(-a); & y \leq 0 \end{cases} \quad (2.7)$$

$$S_1(y) = \begin{cases} -\lambda^+ y; & y \geq 0 \\ -\lambda^- y; & y \leq 0 \end{cases}$$

A use of the simplest policy leads to

$$\bar{S}_t(y) = -\phi(y) + (t-1) \frac{\lambda^- - \lambda^+}{2} a. \quad (2.8)$$

Diagrams of the function $S_t(y)$ at $t=1, 2, \dots, 6$ are given at fig. 1.

It is now possible to make the following conclusions:

1. Supplies fluctuations with an average equal to zero lead to an increase of an average loss which generally speaking is not equal to zero.

2. In absence of initial stock the losses at the simplest policy are equal to

$$(n-1) \frac{\lambda^- - \lambda^+}{2} a$$

i.e. they grow proportionally to the number of subperiods.

3. At the presence of initial stock and under optimal use of accumulations it is possible to reduce average losses. Even at the absence of initial stock, average losses at $n=6$ are reduced approximately twice compared with the results of the simplest policy.

It is possible to show also that the probability of an ab

sence of losses essentially depends on the quantity of initial stock level y_0 and is equal to $(1/2)^{K-N}$, if $N = \frac{y_0}{\alpha} \leq K$.

The described effect of an influence of supplies' irregularity on the efficiency of a system is important itself, but it reflects one side of a problem only. In fact, in the course of analysis it was implicitly assumed that optimization problem for given element of a system may be solved not only at the planned level of supplies but also at the fluctuations near the level. In other words, it was assumed, that at the presence of fluctuations it is possible to satisfy the restriction on planned level of supplies from given element. Generally speaking, it is not always true: at the deviations of supplies it is often impossible to fulfil the requirements on orders and hence in the case of decentralized current planning a necessity arises to commensurate optimally the deviations from the levels required.

A use of the local functionals in a form (I.12) does not give a solution of this problem, because the functions $F_i(\rho)$ are defined only on the region of existence of the solutions. This problem becomes especially critical for a structure of a technological chain type 7, I³ where the elements are sequentially connected and the last element alone gives a final production with definite estimates of its components. At the same time a scheme of decentralization constructed on the basis of dynamical programming may be used for the analysis of such a structure. Here, generally speaking, it is necessary only to construct the dependence of the efficiencies of each element on given resources at small deviations of those from the optimally planned level. In a paper I³ there was shown that locally-exact approximation of a purpose function for the p -th element of technological chain may be represented as a piecewise linear function

$$\hat{f}_p(v) = f_{p, \text{opt}} + \min_{l \in l_p^0} \lambda_p v \quad (2.9)$$

where v - a column-matrix of the deviations from the level optimal level of resources supplied to the p -th element from its predecessor in technological chain. The dependence of each type of products on v has a similar character. With the aid of simplified approximation of a type (2.1) it may be shown that the dependence of deviations from optimal plan for any product of the p -th element on small deviations of any component of a resource consumed from the $(p-1)$ -th element may be re-

presented by the function of a type $\psi(v)$ (look at 2.2). It leads to the effect of accumulation of an influence of fluctuations caused by the irregularity of reciprocal supplies along technological chain if a necessary level of initial reserve for each element is not provided in a chain. In fact, the fluctuations of supplies to the first element of a chain, lead to impossibility to maintain an average level of supplies from the first element to the second one, and so on. Due to this factor a real average levels of supplies to the following elements may considerably differ from the planned ones, and the estimates of the efficiency of deviations justified only in the vicinity of a plan will become incorrect. Therefore an optimal planning itself may be effective only if it takes into account necessity of expenditures on forming of the reserves that give an opportunity to localize an influence of irregularities.

It is natural that creation and maintenance of reserves seems to be not economical from the point of view of a purpose of a system. However, it is necessary and it has the same meaning as an introduction of redundancy into technical systems composed from elements that are not completely reliable. It is worthwhile to stress that it is not an effect of fluctuations of external conditions which is discussed (a necessity of reservation with respect to external conditions is generally accepted) but it is an effect of internal fluctuations stimulated by the fact that the system's planning is incompletely deterministic.

1. Arrow K.J., Hurwicz L., Uzawa H., Studies in linear and non-linear programming, Stanf., Calif., 1958.
2. Marschak T., Centralization and decentralization in economic organizations, Econometrica, v.27, 3 /July 1959/.
3. Dantzig G.B., Linear programming and extensions, Princet., N.-J., 1963.
4. Гольштейн Е.Г., Юдин Д.Б., Новые направления в линейном программировании, "Сов.радио", М., 1966.
5. Волконский В.А., Оптимальное планирование в условиях большой размерности, Экономика и математические методы, 1965, т.1, № 2.
6. Корнаи Н., Липтак Т., Планирование на двух уровнях, в сб. "Применение математики в экономических исследованиях", т.3, "Мысль", 1965.

7. Первозванская Т.Н., Первозванский А.А., Распределение ресурсов между многими предприятиями, Экономика и математические методы, 1966, т.2, № 5.
8. Каценеленбойген А.И., Овсиенко Ю.В., Фраерман Е.Ю., Методические вопросы оптимального планирования в социалистической экономике, ЦЭМИ АН СССР, М., 1966.
9. Шор Н.З., Принцип обобщённого градиентного спуска в блочном программировании, Кибернетика, № 3, 1967.
10. Первозванская Т.Н., Первозванский А.А., Алгоритм поиска оптимального распределения централизованных ресурсов. Изв. АН СССР, Техническая кибернетика, 1966, № 3.
11. Mills H.D., Marginal values of matrix games and Linear programs, "Linear inequalities and related systems", Princet., N.-J., 1956.
12. Arrow K.J., Karlin S., Scarf H., Studies in the Theory of Inventory and Production, Stanf., Calif., 1958.
13. Первозванская Т.Н., Метод приближения в пространстве целевых функций при решении задач на "узкие места", ЖВМ и МФ, т.7, № 3, 1967.

ON A COMPLEX SYSTEMS CONTROL THEORY

A.I. Kukhtenko

Institute of Cybernetics of the Academy of Sciences
of the Ukrainian Soviet Socialist Republic

Kiev

USSR

It was already indicated in a number of publications [1-4] that the study of complex systems must be made, as to necessity by using different levels of an abstract description. Depending on a system designation, either a theoretical informational, or logical mathematical, or dynamical, or at last, a heuristical treatment of the problem may be used. But in reality one must use in most cases a few different levels of abstract description simultaneously.

Without dwelling here repeatedly upon a characteristic of the term "a complex control system", given in the paper [3], we shall mark the fact that the necessity of a description of a complex systems control at a few abstract levels simultaneously compels us to search for those mathematical means that enabling us to make it. However, attempts of applying for this purpose the well-known methods of the automatic control theory or, in general, the dynamic systems theory, the finite automatic theory, Shannon's information theory, etc. demonstrate an evident groundlessness of each of them for this aim. It may be affirmed at least, with respect to a

state of things existing at present. Each branch of the whole scientific trend, connected with the control problem has been developing independently, and only recently the contacts between them are outlined. This report that is a summary statement of a part of a more complete work, prepared on this topic by the author for print is just dedicated to a description of the known things in this field and to a discussion on some possible ways of an investigation of the complex systems control.

§ 1. On an unique conception in the finite automata
theory and the dynamic control theory

It is quite natural that, before "A General Control Theory" is created which will enable us to study in detail behaviour of a complex control system, different investigators try to solve a simpler problem, and try to create a method covering simultaneously, at least, only two of a number of possible abstract treatments of problems. The authors of papers [5, 6], for example, make efforts of uniting the dynamic theory methods and those of the information theory. From this point of view, the opinions stated in work [7] are quite interesting, the author has shown that if one uses some ideas of the abstract algebra, profound analogies existing between the finite automata theory and the dynamic control theory can be revealed. If one represents the Milley's finite automaton, as one makes it usually [8], in terms of five quantities

$$M = \{X, Z, S, f_Z, f_S\},$$

where X and Z are input and output alphabets of the automaton, respectively,

S_y

is a quantities set determining the automaton state,

$$Z_y = f_Z(S_y, X_y)$$

is a characteristic function, by virtue of which the automaton output quantities are determined, if the input quantities and its state are known,

$$S_{y+1} = f_S(S_y, X_y)$$

is a characteristic function, by virtue of which the automaton state in the $y+1$ th tact is determined, if the input quantities and the automaton state in the y th tact are known,

then, as it is shown in [7], for a discrete time scale by setting that X, Z, S are representable through the Abelian (commutative) groups, the system (i.e., the automaton) is quite additive, if and only if such homomorphisms [9] exist

$$A: S \rightarrow S; B: X \rightarrow S; C: S \rightarrow Z; D: X \rightarrow Z;$$

that means that for all $S \in S, X \in X$ the correlations

$$S_{y+1} = f_S(S_y, X_y) = AS_y + BX_y$$

$$Z_y = f_Z(S_y, X_y) = CS_y + DX_y$$

take place, and may be written in the more expanded form

$$S_{y+1} = f_S(S_y, X_1, \dots, X_n) = A^n S_y + \sum_{m=1}^n A^{n-m} B X_{n-m+1}$$

$$Z_y = f_Z(S_y, X_1, \dots, X_n) = CA^{n-1} S_y + \sum_{m=1}^n CA^{n-m} B X_{n-m} + D X_n$$

If now one introduces the notations

$$\Phi(n) = CA^{n-1}$$

$$h(m) = \begin{cases} D & \text{for } m=0 \\ CA^{m-1} B & \text{for } m>0 \end{cases}$$

the expression for can be written in this form

$$Z_y = f_Z(S_y, X_1, \dots, X_n) = \Phi(n) S_y + \sum_{m=0}^{n-1} h(m) X_{n-m}.$$

Substituting, at last the continuous time for the discrete one and passing from a groups homomorphism to a vectorial space homomorphism, the last expression may be given a form, well-known in a linear theory of the dynamic systems control

$$[10] \quad Z(t) = \Phi(t-t_0)x(t_0) + \int_{t_0}^t h(t-\xi)U(\xi)d\xi$$

where $Z(t)$ is a P -dimensional vector characterizing the system output,

$x(t_0)$ is h -dimensional vector characterizing the system state at the moment to t_0 ,

$\Phi(t) = (P \times h)$ is a Z -dimensional vector characterizing the system output,

is a matrix, the i th column of which represents a system reaction at the moment t ,

$h(t)$ is a pulse transient function of the system.

The possibility of this type of transitions from the finite automata theory correlations to the linear dynamic theory correlations reveals to us a close connection, which enables us to speak of an unique conception for these two branches of scientific knowledges earlier independently developed. However, it seems to be necessary to make one more step forward along the path of uniting logical and dynamical treatments of complex systems control theory problems. The section below is dedicated to a discussion of this possibility.

§ 2. On logical dynamical systems

It is not difficult to realize that, besides demonstrating of the fact itself of existence of a close connection, revealed between the finite automata theory and the linear dynamical systems theory, for a factual study of complex control systems, which consist simultaneously of logical and dynamical links united immanently as a whole, the construction of an absolutely new theory is required, consequently mathematical means are also required for its completion. The principal difficulty, arising in studying logical dynamical systems (we shall name them so, for the sake of brevity) consists in the fact that it is necessary to find such a supple language, which should enable us to operate in an equally convenient way both with mathematical analysis ordinary variables and logical ones. While creating this type of language one may apparently go along different paths. From this point of view, these things are worthy of attention: an n -functions language [11], a calculus of operators for Boolean functions [12], a continuous logics language [13] and others. In particular, for the same purposes in the work [14], the concept of an hybrid function $G(x_1, \dots, x_n)$ is introduced that represents a product of the ordinary function of real variable $F(x_1, \dots, x_n)$ and the function of logical variables $f(x_j, \dots, x_m)$, i.e.

$$G(x_1, \dots, x_n) = F(x_1, \dots, x_n) f(x_j, \dots, x_m).$$

The logical function $f(x_j, \dots, x_m)$ may be a predicate, a formula or a quantifier, but it may assume only two values: 1 (verity) and 0 (falsehood), however, the moments, at which this takes place, depend in a complicated way on a variables value of different nature.

They can be:

- 1) predicates depending on a real variable function,
- 2) partly real variables and partly logical variables,
- 3) only logical variables not depending on real variables.

All this creates various possibilities for describing logical dynamical systems. Unfortunately, the fact itself of hybrid function introduction does not mean yet that we have desirable language already.

We need further investigations due to a rules development, by means of which necessary operations with hybrid functions should be made (to differentiate them, to integrate them and so on), therefore the problem of creation of a mathematical tool, fitting for the study of logical-dynamical systems, remains still unsolved. There are already works, in which a tool based on the hybrid function concept [15] is developed, but there are also critical works [16], in which the possibility of such a type are denied. Now it may be noted only that G. von Neumann's prediction about a necessity of a continuous and finite mathematics methods merging (i.e. of merging the mathematical analysis methods basing on the fact of variables continuity and the mathematical logics methods operating with discrete logical variables) comes true, and the needs of a logical dynamical systems theory construction will apparently impel many investigators to go at this aim more hurriedly.

It may be also noted that no less difficulties arise in studying systems, which require a simultaneous use of any two other abstract levels of complex systems description, for example, of the logical and the heuristic ones. Still greater difficulties arise, if one uses simultaneously not two, but three

levels of the systems abstract description, for example, the informational, the logical and the heuristic ones. A creation of mathematical means permitting to realize investigations in such cases and in some more complicated ones is exactly the principal scope of the complex systems general control theory.

§ 3. On a multidimensionality problem in the complex systems theory

The above-mentioned difficulties are not the unique ones, which one meets while studying systems of similar type.

Essential difficulties arise while using one definite level of an abstract description too, if the system consists of many elements or subsystems, interconnected with each other.

By the way, "the multidimensionality curse" (by Richard Bellman's figurative expression [17]) is overcome equally with difficulty in using any of the levels of abstract description. Thus, for example, the problems are solved easily and elegantly in the finite automata theory with a low dimensionality of systems, and the difficulties increase essentially, as soon as a dimensionality of the system studied grows. Therefore a wish appears, quite naturally, to find paths to overcome difficulties conditioned by a system multidimensionality, which should be, for example, suitable simultaneously for the study of multidimensional dynamical and logical systems of a great dimensionality. In our opinion, one can struggle against a multidimensionality and a multirelationship of complex systems not only by a new mathematical methods research (to what many authors limited themselves on the whole) but by using physical and engineering informations concerning the systems studied. Thus, the fact may

be used that the system (it would be of technical, economical, biological or social character) consists of big groups of elements of the same type. If this is the case, we may use a mathematical tool elaborated for the multiphase liquids description, as I.M. Gelfand and M.L. Tsetlin [18] proposed to make it for studying systems of this type. The second possible way of overcoming the difficulties, connected with the problem multidimensionality and considered in [3], is based on physical representations too. Namely: when the system is symmetrical in that or other sense, we may essentially simplify the investigation of a complex system of great dimensionality having both the dynamical and the logical treatment of problems. In these cases the groups theory methods [19] and, more exactly, groups representations theory methods [20] may be used, applied so widely in the quantum physics, the quantum chemistry and in the modern theory of elementary particles [21 - 23] .

By the way, one should not think that the question is indispensable of a symmetry of purely geometrical character. By no means. The symmetry properties, hidden more deeply, are the most important ones for the investigation problem of great dimensionality systems. Thus, the symmetry property manifests, for many and quite diverse dynamical systems, that the Lagrangian (or the Hamiltonian) remains for them invariant with respect to linear transformations of coordinates. If the presence of symmetry is stated in that or other way (v.g. it may be revealed directly by applying a matricial form of equations), a formal apparatus of the groups representations theory enables us to replace an initial problem of great dimensionality by some problems of the same type, but of a considerably smaller

dimensionality (for example, by 10 or 100 times smaller). As a matter of fact the corresponding initial matrix of a greater dimensionality can be reduced, by recipes quite determined and conditioned by the type of an existing symmetry, to a blocking diagonal form

$$\left\| \begin{array}{ccccccc} \Gamma_{ij}^{(1)}(\rho) & 0 & 0 & \dots & 0 \\ 0 & \Gamma_{ij}^{(2)}(\rho) & 0 & \dots & 0 \\ 0 & 0 & \Gamma_{ij}^{(3)}(\rho) & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & \Gamma_{ij}^{(n)}(\rho) \end{array} \right\|$$

where each of the submatrices standing on the main diagonal is already irreducible, i.e. it can be represented once more through matrices of a smaller dimensionality. The groups representation theory enables us to answer, what a dimensionality of the submatrices $\Gamma_{ij}(\rho)$ may be, what their number is, in how many different ways such a type of representation of a matrix of a greater dimensionality by means of matrices of a low dimensionality may be realized and, in particular it gives criteria to judge, whether a further decomposition of the matrices $\Gamma_{ij}(\rho)$ into matrices of some a smaller dimensionality and so on is or is not possible. In a concise summary it is impossible to describe all those procedures, which one must by the way, to fulfil practically, viz to divide groups elements into classes, to find the character for each class, i.e. a trace of the submatrix

$S_{\rho} \|\Gamma_{ij}(\rho)\|$, to determine a canonical basis for an initial multidimensional vectorial space of the given problem and so forth.

We can get familiar with all necessary concepts and theorems (equivalence, and homomorphism of groups, Shoor's lemmas, Lagrange's theorems and so on) and with the techniques of their ap-

plication through the aforementioned books dedicated to an exposition of the groups theory and the groups representation theory 18 - 23 . So far as an investigator himself has often at his disposal a possibility to form a structure of the studied complex system (for example, while controlling economical systems and industrial objects), one has, in choosing the structure as a symmetrical one, a possibility to apply groups representation theory ideas due to which the latter theory acquires a highly great practical importance.

To explain the essence of the matter we shall consider elementarily a simple example, v.g. vibrations of some mechanical system (see more detailed 4).

Suppose that the equations written in a matricial form have the shape:

$$M\ddot{X} + KX = 0,$$

where M is a diagonal matrix of masses,

K is a rigities matrix,

X is a variables column-vector.

$$K = \begin{vmatrix} \frac{3}{2}K_2 & 0 & 0 & -\frac{\sqrt{3}}{2}K_2 & \frac{\sqrt{3}}{2}K_2 \\ 0 & \frac{3}{2}K_2 & -K_2 & \frac{K_2}{2} & \frac{K_2}{2} \\ 0 & -K_2 & \sqrt{3}K_1 + K_2 & 0 & 0 \\ -\frac{\sqrt{3}}{2}K_2 & \frac{K_2}{2} & 0 & \sqrt{3}K_1 + K_2 & 0 \\ \frac{\sqrt{3}}{2}K_2 & \frac{K_2}{2} & 0 & 0 & \sqrt{3}K_1 + K_2 \end{vmatrix}, \quad X = \begin{vmatrix} x \\ y \\ y_1 \\ y_2 \\ y_3 \end{vmatrix}$$

$$M = \begin{vmatrix} m_1 & 0 & 0 & 0 & 0 \\ 0 & m_1 & 0 & 0 & 0 \\ 0 & 0 & m & 0 & 0 \\ 0 & 0 & 0 & m & 0 \\ 0 & 0 & 0 & 0 & m \end{vmatrix}$$

Thus, it is necessary to investigate vibrations of a system, having matrices of dimensionality 5. The problem consists in the fact that, using symmetry property, one should be able to solve, instead of an initial problem, a problem of smaller dimensionality. Applying the known methodics of the groups representations theory [18 - 23] and the fact that the vibrational system used is symmetrical (as it manifests in the matrices symmetry), the new variables $h_{11}, h_{12}, h_{21}, h_{22}, h_{23}$, which are related to the old variables x, y, y_1, y_2, y_3 through such correlations, so that the problem of dimensionality 5×5 splits into two identical problems of dimensionality 2×2 and one problem of dimensionality 1×1 . Without expounding here all details of the transformations based on the groups representation theory ideas, we write out final correlations connecting the old and the new variables

$$\begin{aligned} y_1 &= \frac{1}{\sqrt{3}} \sum_{\alpha=1}^3 h_{\alpha 1}, & x &= \frac{1}{\sqrt{2}} (h_{12} + h_{22}), \\ y_2 &= \frac{1}{\sqrt{3}} \sum_{\alpha=1}^3 e^{\frac{2\pi i \alpha}{3}} h_{\alpha 1}, & y &= -\frac{i}{\sqrt{2}} (h_{12} - h_{22}), \\ y_3 &= \frac{1}{\sqrt{3}} \sum_{\alpha=1}^3 e^{-\frac{2\pi i \alpha}{3}} h_{\alpha 1}, \end{aligned}$$

The same correlations may be written in a matrixial form as follows:

$$X = RH,$$

where the matrix R and the column-vector H have the shape

$$\left\| \begin{array}{cc} 0 & \frac{1}{\sqrt{2}} \\ 0 & \frac{i}{\sqrt{2}} \\ \frac{1}{\sqrt{3}} & 0 \\ \frac{1}{\sqrt{3}} e^{\frac{2\pi i}{3}} & \frac{1}{\sqrt{3}} \\ \frac{1}{\sqrt{3}} e^{-\frac{2\pi i}{3}} & 0 \end{array} \right\| \begin{array}{cc} 0 & \frac{1}{\sqrt{3}} \\ 0 & \frac{1}{\sqrt{3}} e^{\frac{2\pi i}{3}} \\ \frac{1}{\sqrt{3}} & 0 \\ \frac{1}{\sqrt{3}} e^{\frac{2\pi i}{3}} & 0 \\ \frac{1}{\sqrt{3}} e^{-\frac{2\pi i}{3}} & 0 \end{array} \left\|, \quad H = \left\| \begin{array}{c} h_{11} \\ h_{12} \\ h_{21} \\ h_{22} \\ h_{23} \end{array} \right\|$$

The feedback of the new variables with the old ones in matrix shape is written in the form of equation

$$H = R^{-1}X,$$

where R^{-1} is a matrix, inverse to the matrix R . Substituting the linear transformation $X = RH$ for the initial equation and multiplying the left side by the matrix R^{-1} , we reduce the initial system of differential equations to a blocking diagonal form:

$$R^{-1}MR\ddot{H} + R^{-1}KRH = 0$$

where $R^{-1}MR$ is a masses diagonal matrix of the shape

$$R^{-1}MR = \begin{vmatrix} m & 0 & 0 & 0 & 0 \\ 0 & m_1 & 0 & 0 & 0 \\ 0 & 0 & m & 0 & 0 \\ 0 & 0 & 0 & m_1 & 0 \\ 0 & 0 & 0 & 0 & m \end{vmatrix}$$

and $R^{-1}KR$ is a blocking diagonal rigidities matrix of the shape

$$R^{-1}KR = \begin{vmatrix} \sqrt{3} K_1 + K_2 & \frac{\sqrt{3} K_2 l}{\sqrt{2}} & 0 & 0 & 0 \\ -\frac{\sqrt{3} K_2 l}{\sqrt{2}} & \frac{3}{2} K_2 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{3} K_1 + K_2 & \frac{\sqrt{3} K_2 l}{\sqrt{2}} & 0 \\ 0 & 0 & \frac{\sqrt{3} K_2 l}{\sqrt{2}} & \frac{3}{2} K_2 & 0 \\ 0 & 0 & 0 & 0 & \sqrt{3} K_1 + K_2 \end{vmatrix}$$

Thus, on the basis of a quite determined procedure, we have come to a blocking diagonal shape of matrices. One makes this not on the basis of a method of attempts, but according to quite definite algorithms in the method given, if the initial system is symmetrical. The differential equations system in terms of the new variables has the form:

$$\left. \begin{aligned}
 m \ddot{h}_{11} + (\sqrt{3} K_1 + K_2) h_{11} + i \sqrt{\frac{3}{2}} K_2 h_{12} &= 0 \\
 m_1 \ddot{h}_{12} + \frac{3}{2} K_2 h_{12} - i \sqrt{\frac{3}{2}} K_2 h_{11} &= 0
 \end{aligned} \right\} \\
 \left. \begin{aligned}
 m \ddot{h}_{21} + (\sqrt{3} K_1 + K_2) h_{21} - i \sqrt{\frac{3}{2}} K_2 h_{22} &= 0 \\
 m_1 \ddot{h}_{22} + \frac{3}{2} K_2 h_{22} + i \sqrt{\frac{3}{2}} K_2 h_{21} &= 0
 \end{aligned} \right\} \\
 m \ddot{h}_{31} + (\sqrt{3} K_1 + K_2) h_{31} = 0$$

Thus, one can see clearly that the matter comes, instead of solving the problem of dimensionality 5×5 , to a twofold solution of the same problem of dimensionality 2×2 and to a separate solution of an equation of dimensionality 1×1 .

We take an interest in the complex systems control problem, and the method just exposed gives a possibility to investigate a controlled member taken separately, having a symmetrical structure, if the control system is asymmetrical. Therefore, in these cases, the stability, invariance, optimality problems and others are independent ones, each of which is solved on the basis of using the controlled member symmetry properties. Thus, the stability problem may be solved by means of an application of a decomposition method, according to which all the system is disintegrated into subsystems. The controlled member described by a symmetrical matrix enters the first of these subsystems, and the control system enters the second subsystem. At first one investigates a stability of each of the subsystems. Moreover, it is known that one uses, while studying the stability of a controlled member taken separately, their symmetry properties. Then, determining the subsystem stability and on the basis, for example, of the known Baley's theorem, the asymptotic stability

conditions can be determined for the system as a whole. The optimal control problems are particularly interesting for controlled members of this type, when a common use of the decomposition method and of the groups representation theory ideas enables us to solve highly difficult multidimensional problems of optimal control. The factual demonstration of a similar type of possibilities requires a separate report.

§4. On universality of electronic digital computers and on a complex systems control processes simulation

The principal particularity of the complex systems theory, as of the Cybernetics in general, becomes apparent in the fact that it must give possibilities to study objects of any nature (technical, economical, biological, social etc.) at a corresponding abstract level. Therefore it is necessary to take care not so much of a creation of " A general systems theory", from which particular theories as special cases should result (the linear dynamic systems theory, the information theory, the Markovian processes theory and so on), which [24] call for at times, as of a creation of many branches of a theory for different levels of a systems abstract description. A consideration of the problem at a singular abstraction level enables us to receive answers only for a certain group of questions, and in order to receive answers for other questions it is necessary to make an investigation at an abstract description level of other systems. To reach a maximally possible completeness of an information it is necessary to study the same system at all abstraction levels suitable for the case

given. A path of mathematical simulation is the most expedient for this purpose, and it is practically accessible. Although one usually looks at an electronic digital computer as a high-speed computer, but it is in reality a universal device, which besides a rapid computation, can make the processing of alphabetical or other symbols, transform an information into a form that we need, draw "conclusions", "deductions" and so on. All this enables us to study by means of an electronic digital computer not only informational processes in a complex systems control, but logical, dynamical and heuristic treatment. Nowadays special abstract languages are created (SIMSCRIPT, SIMPAC and others, [25]) permitting to economize the time and efforts due to a programming and a process itself of simulation. Therefore just this path of using the electronic digital computer for a simulation is now principal while developing really complex control systems. One often undertakes even the creation of special scientific centers intended exceptionally for the simulation purposes of a complex control system developed. As an example, one may mention a scientific center, created especially for developing an automatic system, controlling air traffic over the Western Europa territory [26] . To characterize a real complexity of control systems of this type one may give data about programs needed for a control by means of electronic digital computer and other technical means (radars, communication equipment etc.) of aircraft flows (with total number 300 - 600) over a territory of about 2000 km [27] . One needs efforts of 250 programmers during two years only for a development of such programs, and the total quantity of commands reaches 2.5 millions. In spite of the fact that the cost of scientific simulation centers of this type is sufficiently high, the

economical expediency of their creation is doubtless while developing really complex control systems, and one follows this path in many cases both while solving technical, economical or defensive problems and fulfilling of great social investigations. Therefore the problem of developing a general theory of algorithms transformation (acad. V.M. Glushkov [28]) and of a complex systems theory in general becomes (some aspects of which were considered above) most significant, because so powerful means of mathematical modelling can be used completely only in a case of its presence. Of course, only some complex systems theory general questions are marked in this brief survey, and its numerous important aspects, which are ^{due} to a complex systems reliability problematics [29 - 30] and their efficiency [31, 32] etc.

References

1. A General Systems Theory. Publishing House " The World", Moscow, 1966.
2. Acad. A.I.Berg, J.L. Cherniak. Information and Control. Publishing House " The Economics", Moscow, 1966.
3. A.I. Kukhtenko. Principal trends of a development of the complex systems control theory. In the collection "Complex Control Theory", Issue IV, Publishing House " The Scientific Thought", Kiev, 1968.
4. A.I. Kukhtenko. Principal Problems of the Complex Systems Control Theory. Transactions of the Seminar "Complex Control Systems". Issue I, Kiev, 1968. Publishing House of the Institute of Cybernetics of the Ukrainian Soviet Socialist Republic.

5. A.A. Krassovsky. Continuous Dynamic Systems Entropy Change. News of the Academy of Sciences of the USSR. "Technical Cybernetics", N 5, 1964.
6. R.L. Stratonovich. Theory Information Extremal Problems and a Dynamic Programming. News of the Academy of Sciences of the USSR. "Technical Cybernetics", N 5, 1967.
7. M. Arbib. A Common Framework for Automata Theory and Control Theory. I Soc. Industr. and Appl.Math., 1965, A.3. N 2.
8. A. Gill. Introduction into a Finite Automata Theory. Publishing House "The Science", 1966.
9. R.Fort, A.Kauffmann, M.Denis-Papin. A Modern Mathematics . Publishing House "The World", 1966.
10. V.M. Brown. An Analysis of Linear Systems Invariants in Time. Publishing House "Machine-Building", 1966.
11. V.L. Rvachev. Logics Algebra Geometrical Applications. Publishing House "Technics", Kiev, 1967.
12. I. Richalet. Calcul Opérationnel Booléen, "L'onde électrique", 1965, N 439, v. 43, p. 63.
13. S.A. Ginsburg. A Continuous Logics and its Applications. AIT, N 2, 1967.
14. R. Terno. Hybrid Functions Are a New Method of Complex Systems Description. News of the Academy of Sciences of the USSR. "The Technical Cybernetics", N 6, 1965.
15. M.M. Gerdov. One Numerical Analysis Method of Complex Systems. News of the Academy of Sciences of the USSR. "The Technical Cybernetics", N 6, 1967.
16. G.P. Sbrodov. Logical Relations and Complex Control Systems Description. Proceedings of the Conference on the Economical Cybernetics, Batumi, 1966.
17. R. Bellman. A Trend of Mathematical Investigations in the Nonlinear Circuits Theory. The Collection "Mathematics", 8, 15. Publishing House of Foreign Literature, 1964.
18. I.M. Gelfand, M.L. Tsetlin. On Some Complex Systems Control Methods. "Successes of Mathematical Sciences", v. 16, Issue I, 1962.

19. A.G. Kurosh. A Groups Theory. Publishing House "The Science", 1967.
20. F.D. Murangan. A Groups Representation Theory. Publishing House of the Foreign Literature, 1950.
21. E. Vigner. A Groups Theory and its Application in the Quantum Mechanics. Publishing House of the Foreign Literature, 1947.
22. V. Heine. A Groups Theory in the Quantum Mechanics. Publishing House of the Foreign Literature, 1963.
23. G.I. Liubarsky. A Groups Theory and its Application in the Physics. Publishing House of the Physical and Mathematical Literature, 1958.
24. M. Mesarovic. A General Systems Theory. In the Collection "A General Systems Theory". Publishing House "The World", 1966.
25. G. Markovits et al. Simscript. An Algorithmical Language for a Modelling. Publishing House. "The Soviet Radio", 1966.
26. Control. II, N 105, 1967.
27. T.B. Steel. The Development of Very Large Programs. "Proc. IFI. P. Congr. N 1, vol. 1, 1965.
28. V.M. Glushkov. Perspectives of the Automation of a Computers Design. The Herald of the Academy of Sciences of the USSR, N 4, 1967. Publishing House "The Thought".
29. I.P. Popchev. Great Systems Structure. Principal Dependencies and Description Methodics. "The Technical Thought", N 1, 1967.
30. I. Popchev, M. Panova. A Determination of a Great Systems Efficiency Index. News of the Institute of the Technical Cybernetics, v. VI, 1967.
31. B.S. Fleischmann. On a Complex Systems Vitality. News of the Academy of Sciences of the USSR. "The Technical Cybernetics", N 5, 1966.
32. B.S. Fleischmann. Complex Systems Efficiency Statistical Limits. The Collection "Applied Problems of the Technical Cybernetics". Publishing House "Soviet Radio", 1966.

CONSTRUCTION OF CHECKING AND DIAGNOSTIC PROCEDURES FOR
GENERAL-PURPOSE UNIFORM ARRAYS

I.V. Pranghishvily

V.V. Ignatushchenko

Institute of Automatics and Telemechanics
Moscow, U.S.S.R.Introduction

Uniform arrays represent a new class of networks consisting of regularly interconnected identical cells capable to be preprogrammed to perform the desired functions. Investigations carried out in the Institute of Automatics and Telemechanics, and also in some other organizations in the U.S.S.R. and abroad showed that uniform programmable general-purpose arrays are among the most effective and promising vehicles for implementation of digital computers and control systems based on integrated circuits, especially LSI-based. Both uniform programmable microelectronic arrays and functions realized by them are advantageous in that they give general-purpose, flexible, economic, reliable and "hardy", easy-and-cheap-to-fabricate, high-capacity devices of a high standardization level.

We are going to show in this paper that regularity of uniform arrays leads to considerable simplification in the construction of checking and diagnostic procedures, and to drastic reduction in the time necessary for checking and diagnosis of arrays, irrespective of the function specified.

It is assumed here that any sequential function realized by a uniform array, is in order if all functional elements (cells) of the array used to perform it are in order. Then for any function realized by a uniform array, checking and failure diagnosis are reduced to testing of the array cells. Since all cells of a uniform array are identical and have identical interconnections, the tests that make the checking and diagnostic procedure for a uniform array as a whole, are sets of tests forming the

procedure for one array cell.

Here we discuss those arrays described in Ref.¹⁾ whose functional element is a finite automaton with k inputs, q outputs, and memory that serves only to set up the cell. p of k inputs correspond to the interconnections, i.e. they are identified with the outputs of other elements of the array. We shall call them further "p-inputs". The remaining $k-p$ inputs are used for control, and setting up signals are fed along them from outside into the element.

Given a (k,q) -terminal network (functional element of the array), one can construct for each of its outputs a minimal (elementary) checking procedure Π_{\min} by the use of the technique of Refs.²³⁾. Let the minimal procedure Π_{\min} for one output include n tests:

$$\Pi_{\min} = t_1 \cdot t_2 \dots t_i \dots t_n \quad (1)$$

In general, a uniform array has external connections through peripheral cells, and through control inputs. For instance, Fig.1 shows that a uniform array $m_1 \times m_2$ whose cells have four inputs and four outputs (i.e. a signal may be transmitted through a cell in four directions), has $2m_1 + 2m_2$ external p-inputs and the same number of external outputs.

Since each p-input and output of an array functional element is identified with a certain direction of signal transmission, it is convenient to regard an uniform array functional element as a set of $p=q$ identical elementary transmitting channels that interact inside it. Consequently the functional element of the array shown in Fig.1 is regarded as a set of four elementary channels.

From this point of view, elementary procedure Π_{\min} serves to test one of all channels with due regard to other channels. It is assumed here that the output of a cell elementary channel F_i is the input of a similar channel F_j of a neighbouring cell.

Let us give definitions for tests that are fed into adjacent cells a and b connected by the i -th channel.

Let cells a have at its inputs signals

$(\epsilon_1 \dots \epsilon_i \dots \epsilon_p, \epsilon_{p+1}^H \dots \epsilon_K^H)^a$ that correspond to the test t_i from Π min, where ϵ_i assumes the values 0 or 1; ϵ_i is the signal at the i -th p -input (i.e. the signal at the input of i -th elementary channel), ϵ_{p+1}^H is the signal at the $(p+1)$ -th control (set-up) input. If cell a is in order, the signals $(\epsilon_{out1} \dots \epsilon_{outi} \dots \epsilon_{outp})^a$ will appear at its outputs. If the i -th channel of cell a is out of order, or the cell a as a whole is out of order, the signal $\bar{\epsilon}_{outi}$ will appear at the i -th output.

We shall call the test $(\epsilon_1 \dots \epsilon_i \dots \epsilon_{p+1}^H \dots \epsilon_K^H)^b$ fed into the cell b , compatible in the i -th direction (channel) with the test $(\epsilon_1 \dots \epsilon_i \dots \epsilon_{p+1}^H \dots \epsilon_K^H)^a$ fed into cell a , if value of ϵ_{outi}^a at the i -th output of a good cell a coincides with signal value at the input of the i -th channel cell b in the test $(\epsilon_1 \dots \epsilon_i \dots \epsilon_{p+1}^H \dots \epsilon_K^H)^b$ i.e. $\epsilon_{outi}^a = \epsilon_i^b$. It is obvious that if $\epsilon_{outi}^a \neq \epsilon_i^b$, the test $(\epsilon_1 \dots \epsilon_i \dots \epsilon_{p+1}^H \dots \epsilon_K^H)^b$ may not be fed into the cell b at all.

Let for instance the test $(\epsilon_1 \epsilon_2 \epsilon_3)^a = 001$ be fed into the cell a of the bidirectional uniform array shown in Fig.2, and let the signal $\epsilon_{out}^a = 0$ appear at the horizontal output of the good cell a . Then the cell b may be fed only with such test which is compatible with $(\epsilon_1 \epsilon_2 \epsilon_3)^a$, i.e. has at the horizontal input the signal value $\epsilon_3^b = 0$. If tests t_i and t_j are compatible with each other, we shall call these two tests intercompatible.

Fig.2 shows a special case where compatible tests for cells a and c that are neighbours along vertical line coincide $((\epsilon_1 \epsilon_2 \epsilon_3)^a = (\epsilon_1 \epsilon_2 \epsilon_3)^c = 001)$. We shall call such tests self-compatible.

We shall call the compatible test $(\epsilon_1 \dots \epsilon_i \dots \epsilon_{p+1}^H \dots \epsilon_K^H)^b$, fed into the cell b , conjugate in the i -th direction (channel) with the test $(\epsilon_1 \dots \epsilon_i \dots \epsilon_{p+1}^H \dots \epsilon_K^H)^a$

fed into the cell a, if change of signal value from $\sigma_{out i}^a$ to $\bar{\sigma}_{out j}^a$ at the i-th output of the faulty cell a makes the cell b to change its i-th output from $\sigma_{out i}^b$ to $\bar{\sigma}_{out i}^b$.

If each cell, following the cell a in the i-th direction is fed with a test that is conjugate with the test fed into the preceding cell it becomes possible to obtain at array external output signals $\sigma_{out i}^a$ or $\bar{\sigma}_{out i}^a$ corresponding to the good and faulty states of the cell a.

The sets of tests $[(\sigma_1 \dots \sigma_k^a)^a, (\sigma_1 \dots \sigma_k^b)^b, \dots, (\sigma_1 \dots \sigma_k^r)^r]$, where each test is conjugate in i-th direction with the preceding one, and the last with the first, will be called the test period (r).

In a special case conjugate tests may coincide. Such tests will be called self-conjugate, their period contains one test, (r=1).

If a period contains two tests (r=2), these tests will be called interconjugate.

The functional element of a uniform structure may be set up so that the same signal $\sigma_{out i}$ appears at several outputs corresponding to different channels (rather than at one output), i.e. the output signal of the cell is reproduced, and $\sigma_{out i} = \sigma_{out j}$.

Let a cell a, which is set up so that $\sigma_{out i}^a = \sigma_{out j}^a$ be fed with test t_i from Π_{min} . If its neighbour in the j-th direction (a cell c) is fed with such test t , that change of the cell c j-th input signal from

$\sigma_{out i}^c = \sigma_{out j}^c = \sigma_j^c$ to $\bar{\sigma}_{out i}^c = \bar{\sigma}_{out j}^c = \bar{\sigma}_j^c$, results in change of the cell c output signal value from

$\sigma_{out j}^c$ to $\bar{\sigma}_{out j}^c$, we shall call the test t conjugate with the test t_i in intersecting directions (i and j).

The problem of checking uniform array cells (detection of the fact of existence of faulty cells^(x)) may

(x) It is implied here and henceforth that there is one or more non-compensating faults.

be formulated then in terms of checking array channels with tests from Π_{\min} . In so doing one should feed conjugate tests into the cells situated in between checked ones and external outputs. The most convenient tests are those that make cells to perform connecting functions.

The problem of diagnosis consists in finding two coordinates of a faulty cell and is reduced then to the checking of the cell in two intersecting directions (channels).

The less is the number of tests in a period, the less is the time (number of clock cycles) necessary to check array cells along the given direction. Indeed, if a test is self-conjugate, it can simultaneously (during one clock cycle) check all cells along an i -th direction (from first to the last that has external i -output). For instance, 1-channels of all first row cells in the array shown in Fig.1, may be checked during one clock cycle. If, in addition, the given test is self-compatible in other channels all i -th channels of all array cells may be checked simultaneously during one clock cycle.

If the given test t_i from Π_{\min} is interconjugate (in an i -th channel) with a test \bar{t} from Π_{\min} , then the test t_i can during one cycle check half of the cells (to the accuracy of the greater integer) and the test \bar{t} checks the rest of them. At the next cycle the tests must be interchanged. In the case of a test \bar{t} (interconjugate with the given one) not being included in the minimal procedure Π_{\min} , the given test t_i needs two clock cycles to check the i -th channels.

Thus the time for checking array cells with the minimal procedure Π_{\min} depends on the number of tests in periods constructed for each test in Π_{\min} .

For each test t_i from Π_{\min} the conjugated in an i -th channel tests are found as follows:

a) The values of $(\epsilon_1 \dots \epsilon_i \dots \epsilon_k) = t_i$ are substituted into the truth table or algebraic expression of the function F_i (i.e. output function of the i -th channel of the cell), and $G_{out\ i} = F_i$ is found. If $G_{out\ i} = \epsilon_i$

the given test is selfcompatible;

b) if the test t_i is selfcompatible, it is examined for being selfconjugate, to this end the values of $(\bar{c}_1 \dots \bar{c}_i \dots \bar{c}_k)$ are substituted into the truth table or algebraic expression of F_i . If $\bar{c}_{out\ i}$, is the result, the given test is self-conjugate;

c) If the given test is not selfconjugate (in an i -th channel), one should pick out from the truth table those combinations of input signals $(\bar{c}_1 \dots \bar{c}_i \dots \bar{c}_k)' = t'$, which have $\bar{c}'_i = \bar{c}_{out\ i}$; these combinations correspond to tests that are compatible with the given test t_i .

d) from the tests of t' type those tests t'' are chosen, which have $\bar{c}_{out\ i}'' = \bar{c}_i$. Then the given test t_i will have intercompatible tests of t'' type.

e) the tests interconjugate with the given one (t_i), are picked out from tests of t'' type. For this purpose the values of $(\bar{c}_1 \dots \bar{c}_i \dots \bar{c}_k)''$, and $(\bar{c}_1 \dots \bar{c}_i \dots \bar{c}_k)$ are substituted (similar to the point b)) into the truth table or algebraic expression of F_i . This corresponds to $\bar{c}_i'' = \bar{c}_{out\ i}$ and $\bar{c}_i = \bar{c}_{out\ i}''$. If in doing so the value of F_i changes to the inverse one (i.e. there are $\bar{c}_{out\ i}''$ and $\bar{c}_{out\ i}$, respectively), the corresponding test is interconjugate with t_i .

f) if no one test is found that is interconjugate with the given one, then from the tests of t' type compatible with t_i , the tests of t''' type having $\bar{c}_{out\ i}''' \neq \bar{c}_i$ are chosen. Now the tests of t''' type are looked for, which are conjugate with the test of t'' type and have $\bar{c}_{out\ i}''' = \bar{c}_i$. If they are not found, search goes on among tests of t''' type until a test is found that is conjugate with the previous one and has $\bar{c}_{out\ i}^{(2)} = \bar{c}_i$, i.e. until a period of conjugated tests is found. For uniform array cells known to authors, the test period does not exceed three tests.

In order the given checking test t_i to be also diagnostic one, it is necessary to find out a test t , which is conjugate with the given one in a direction intersect-

ing the i -th direction. Simultaneously with feeding the given test t_i into cells connected by the i -th channel, the neighbouring cells are fed with the conjugate test \bar{t} . If it is necessary to diagnose e.g. cells situated in the 1st row of the array shown in Fig.1, the checking tests from \bar{t} are fed into the cells of the 1 row providing for transmission of "in order/out of order" signal to the external output of the channel 1. At the same time the cells of all other rows are fed with tests that are conjugate with the given one in the 2nd channel. This provides for the transmission of "out of order" signals from faulty cells of the 1 row to vertical outputs of the lower row cells.

It is clear now that those checking and diagnostic tests are most effective that are selfconjugate in two intersecting channels (directions). One such test enables all appropriate channels of array cells to be examined simultaneously (during one clock cycle). It serves not only to detect the presence of a failure in the array, but to find out coordinates of faulty cells too.

The above discussion leads to the following algorithm for construction of uniform array checking and diagnostic procedures (Method 1):

1. By the use of the technique of Refs.2,3, a set of minimal checking procedures for the channels of an array functional element is found. For each test from the minimal procedure, the minimal conjugate test period (including the given test) is found.

2. Of all minimal checking procedures for the array cell that one is chosen which contains more selfconjugate (in two intersecting directions, or channels) tests. All array cells are checked and diagnosed simultaneously (during one clock cycle) with each of such tests.

3. For each test t_i which belongs to the elementary procedure (found in accordance with the point 2.) and is selfconjugate only in the i -th channel, a test \bar{t} is found which is conjugate with the given one in a direction which intersects the direction 1. A group of cells whose

i -th channels are set up so that $G_{out i} = G_{out j}$ is isolated. For example, the array shown in Fig.1 has cells of one row (connected with channels 1) isolated and set up to perform $G_{out 1} = G_{out 2}$. The isolated group of cells is fed with the test t_i (selfconjugate in the i th channel), and the rest of the array is fed with the tests t . If the isolated group has some faulty cells, the "out of order" signals will appear at the external output of the i -th channel and at the corresponding external outputs of the j -th channels.

4. Checking and diagnosing of the isolated group of the i -th channels with the rest of tests (included in the procedure found in accordance with point 2) is performed similar to point 3.

Points 3 and 4 are performed separately for each group of cells having connected i -channels. In Fig.1, for instance, each row of cells is examined.

The rest of channels is examined in a similar manner.

The total number of uniform array checking and diagnostic procedures (Mcd. by Method 1) and the time necessary for them (number of clock cycles (N)) is related to the number of array functional elements (elementary channels) as follows:

$$(M_{c.d.} = N) \leq l \cdot \left\{ \frac{p}{2} \right\} + (l' + l'' + l''' \cdot 2) / (m_1 + m_2 + \dots + m_p) \quad (2)$$

where: l is the number of tests in the minimal procedure min that are selfconjugate in two crossing directions;

$\left\{ \frac{p}{2} \right\}$ is the number of pairs of intersecting channels (directions) rounded off to the greater integer (e.g. for the Fig.1 array, where $p=4$, $\left\{ \frac{p}{2} \right\} = 2$);

l' is the number of tests from \prod min that are selfconjugate ($r=1$) only in i -th channel (direction);

l'' is the number of tests from \prod min that are interconjugate ($r=2$) in i -th channel;

$\{n_i\}$ is the number of tests from Π min that are conjugate in i -th channel and have period $r=2$;

$m_1, m_2, \dots, m_i, \dots, m_p$ is the number of groups of connected elementary channels of 1st, 2nd, i -th, \dots , p -th directions (e.g. the array of Fig.1 has m_1 groups (rows) of connected channels of direction 1, m_2 groups (columns) of connected channels of direction 2).

Uniform array testing time may be made independent of the number of array elements. For this purpose, it is necessary to check all groups (rows, columns) of cells or a certain part of them with conjugate (in i -th channel) tests rather than one group. If some of the i -th channels have at their external outputs "out of order" signals, the cells are diagnosed by Method I to find second coordinates of faulty cells. In accordance with this, the following method for checking and diagnosis of uniform array (Method II) is proposed:

1. The point 1 of Method 1 is to be performed.
2. An elementary procedure is chosen which belongs to the set of array cell minimal procedures and has the greatest number of pairs of intercompatible (in j -th channel) tests selfconjugate in i -th channel. If there are several such procedures, that one is chosen which has the greatest number of selfconjugate (in two intersecting directions) tests.
3. The i -th channels of all uniform array cells are checked with each pair of the above mentioned (point 2) tests during two cycles. All other channels of array cells in this case, are fed with set-up signals providing for compatibility of their inputs and outputs. All channels are checked in a similar manner.
4. During $\{\frac{p}{2}\}$ cycles all channels of all cells are checked with selfconjugate (in two intersecting directions) tests.

5. All i -th channels of all array cells are checked then with the rest of tests during r cycles. Thereupon the rest of channels is checked.

The number of checking procedures (M'_c) and checking time (number of cycles (N') for uniform array (irrespective of number of cells) are defined in Method II as follows:

$$(M'_c = N') \leq \ell \cdot \left\{ \frac{p}{2} \right\} + (\ell' + \ell'' + \ell''' \cdot z + z') \cdot p \quad (3)$$

where r' is the total number of auxiliary tests t_{aux} that are fed into the elements situated in between m - and $(m+r)$ -th groups of the i -th channels checked with the given test t_i . If, at the same time, t_i is selfcompatible in any other (j -th) channel, but the i -th one, or is intercompatible in the sense of point 2, tests t_{aux} are looked for that are not included in Π_{min} , but are compatible with t_i in j -th channel. Thus, tests t_{aux} serve to separate cells, if t_i may not be fed into neighbouring groups of elements whose i -th channels are not connected together. Only one coordinate of each faulty cell is found during N' clock cycles. In order to find out the second coordinates of m -th group faulty cells, a checking and diagnostic procedure is used that corresponds to checking (with the same test t_i) i -th channels of m -th group by Method I. Thus the total number of uniform array checking and diagnostic procedures ($M'_{c,d}$) by Method II is as follows:

$$M'_{c,d} = M_{c,d} + M'_c - \ell \cdot \left\{ \frac{p}{2} \right\} \quad (4)$$

Only $M'_c + S$ procedures are needed to diagnose S faults, since each fault needs only one procedure from $M_{c,d}$. One additional clock cycle is needed to find out the second coordinate of each faulty cell. Thus S faults are diagnosed in $N' + S$ cycles.

Examples of Procedures for Uniform Arrays. Majority Array.

The majority array of Ref.⁴ capable of performing an arbitrary switching function consists of three-input cells. (Fig.3). Only cells along edges have their horizontal and vertical inputs and outputs led outside. The central (control) inputs to all array cells are external. Each cell has one bifurcating output, i.e. $\bar{G}_{out i} = \bar{G}_{out j}$.

By the use of the technique of Refs.²³, a set of minimal procedures is found for checking the cell performing the function $\beta = AB + AC + BC$ (Fig.4):

$$\begin{aligned} \prod \min (G_1=A, G_2=B, G_3=C) = & (010, 101, 011, 001) \vee \\ & (010, 001, 011, 110) \vee \\ & \vee (101, 001, 100, 110) \vee (010, 100, 011, 110) \vee \\ & \vee (010, 101, 100, 110) \vee (001, 011, 101, 101) \end{aligned} \quad (5)$$

According to the points 1 and 2 of Method I, the minimal procedure

$$\prod \min (ABC) = 010, 101, 011, 001 \quad (6)$$

is picked out, because it has two tests selfconjugate in two channels (010, 101). Indeed, for $A=0, B=1$ and $C=0$ we get $\bar{G}_{out} = f=0$, i.e. $\bar{G}_1 = \bar{G}_3 = \bar{G}_{out}$, and the test 010 is selfcompatible in the channels A and C. Since alteration of signal value at the input of A, or C results in $f=1 = \bar{G}_{out}$, the test 010 is also selfconjugate in these channels.

The 4x4 majority array (Fig.5.a) is checked and diagnosed with the test 010 in one cycle T_1 . To do this the external inputs of extreme array cells A and C are fed with zero signals, and the central inputs of all cells are fed with 1 signals. If any cell is out of order (e.g. in the 3rd row and 2nd column), it will generate 1 instead of 0, and all cells to the right and lower the faulty one will change their outputs from 0 to 1. Thus both coordinates of the faulty element will be determined. In a similar manner all cells will be diagnosed with the test 101 during cycle T_2 (Fig.5,b).

The test 011, belonging to procedure (6), is self-

conjugate in the horizontal channel (C), because $f = \overline{C_{out}} = 1 = C$ at $A=0$, $B=C=1$, and $f = \overline{C_{out}} = 0$, if C changes to 0. According to point 3 of Method I, the test 011 is fed into all elements of the first row (Fig.5,6). All other cells are fed with the selfconjugate test 101 which is conjugate with 011 in the vertical channel. Thus, the test 011 checks and diagnoses all cells of the 1-st row during one cycle (T_3), all cells of the 2nd row during T_4 , 3rd T_5 and 4th T_6 (Fig.5,2). If, for instance, cell (2,2) is faulty, it will generate 0 instead of 1 at the output. This will result in generation of 0 (instead of 1) at outputs of all cells to the right and below the faulty one, and both coordinates of the faulty cell will be determined.

The last test from (6) (001) is selfconjugate in the vertical channel (A). According to the point 3 of Method I, it checks and diagnoses all cells of the 1-st column during one cycle (T_7) (Fig.5,8). All other cells are fed with selfconjugate test 010 which is conjugate with 001 in the horizontal channel. The cells of 2nd, 3rd and 4th columns are checked and diagnosed during cycles T_8 , T_9 , and T_{10} , respectively. Thus, all cells of the 4x4 majority array are checked and diagnosed with the tests of procedure (6) during 10 cycles.

Now let us construct a checking procedures M_C^I for the same array by the use of Method II. According to point 2 of Method II, the procedure

$$\prod \min (ABC) = (010, 001, 011, 110) \quad (7)$$

from the set (5) of majority cell minimal procedures is chosen, since it contains two intercompatible (in the horizontal channel) tests (001 and 110). Each of them is selfconjugate in the vertical channel.

During the cycle T_1 the test 001 checks all cells of odd columns (see hatched cells in Fig.6,a), and the test 110 checks cells of even columns. In the cycle T_2 the tests are interchanged (Fig.6,б).

During cycle T_3 all cells are checked with test 010 (Fig.5,a). During cycle T_4 cells of odd rows are checked with selfconjugate (in horizontal channel) test 011 (Fig. 6,6). All cells of even rows are fed at the same time with the auxiliary test 100 which is intercompatible with 011, but is not included in (7). The cells of even rows are checked during cycle T_5 with 011 (Fig.6,2). If the cell (2,2), for instance, is out of order, the corresponding signal (shown in brackets in Fig.6,2) will appear at the output of the 2nd row when it is checked with the test 011. In order to find out the second coordinate of a faulty cell it is necessary to diagnose the array with a test, corresponding to checking cells of the 2-nd row with the same test 011 by Method I (Fig.5,2). Thus, the majority array is diagnosed during 5 clock cycles, and the time for diagnose by Method II is independent of array size.

The Uniform Array Capable of Performing an Arbitrary Sequential Function

Let us construct a checking procedure M_c^I (Method II) for a four-directional array of arbitrary size (Fig.1). Fig.7 shows the functional diagram of one of four channels of the array cell, where A,B,C,D are external inputs serving to control (set up) the memory cells (the flip-flops TP5 and TP4), X is working input (or p-input), F is the channel output. The channel performs the logical function

$$F = X \oplus Y_{TP5} + Y_{TP4} \quad (8)$$

where \oplus is the "sum by moduls two" function.

It is assumed in the procedure construction for (8), that during all cycles T of channel F testing $A(t)=0$, i.e. the flip-flop reset signals are absent. It is also assumed that the inputs C and D get pulse signals i.e., their duration is less than T. The signal at a flip-flop output should be constant (either $Y_{TP}=0$, or $Y_{TP}=1$). Under these restrictions, and with signal duration at the B,C, and D inputs disregarded, one can determine Y_{TP5} and Y_{TP4}

as follows:

$$Y_{TP5} = BC, Y_{TP4} = BD \quad (9)$$

Substituting (9) into (8), one gets:

$$F = [X \oplus (BC)] + BD \quad (10)$$

The procedure

$$\prod_{\min}(XBCD) = (1110, 0110, 1100, 0011, 0101) \quad (11)$$

is one of the elementary procedures for circuits implementing the function (10). To see this, look through the Table 1, where the tests 1 and 2 are interconjugate, and the tests 3 and 4 are conjugate.

Table 1

No.No.	X	B	C	D	F
1	1	1	1	0	0
2	0	1	1	0	1
3	1	1	0	0	1
4	0	0	1	1	0
5	0	1	0	1	1

As Fig.8 shows, the array functional element is an aggregate of four similar channels (F_1, F_2, F_3 , and F_4).

Thus, it is described by the system of four equations.

$$\left. \begin{aligned} F_1(i,j) &= \underbrace{[F_1(i,j-1) \cdot F_2(i-1,j) \cdot F_3(i,j+1) \cdot F_4(i+1,j)]}_X \oplus Y_{TP5} + Y_{TP4} \\ F_2(i,j) &= X \oplus Y_{TP5} + Y_{TP2} \\ F_3(i,j) &= X \oplus Y_{TP5} + Y_{TP3} \\ F_4(i,j) &= X \oplus Y_{TP5} + Y_{TP4} \end{aligned} \right\} \quad (12)$$

Test 3 is selfconjugate in any pair of intersecting channels. The array testing consists in testing all elementary channels F of all cells with the tests of Table 1, hand in hand with testing conjunction in X of (12).

According to points 3 and 6 of Method II, the F_1 channels of all odd cells in each row are fed with the first test (1110), and even cells get the test 0110 which is interconjugate with the first one (See Fig.9a.).

All other channels (F_2, F_3, F_4) of each cell get test 0111 (compatible with 1110), or 111 (compatible with 0110) which correspond to channel blocking. If any of F_1 channels is faulty, the "out of order" signal will emerge at the horizontal output of corresponding row.

It should be noted that during the cycle T_1 in the even cells of each row conjunction in X is tested with the test $(F_1 F_2 F_3 F_4)^X = (0111)^X$, i.e. with one of five tests serving to check the four-input "AND" gate: $(0111, 1011, 1110, 1111, 1101)^X$. At the expiration of cycle T_1 (and each successive cycle) all cells are fed with the signal $A=1$ (at $B=C=D=0$) that results in resetting the cell flip-flops into initial state, i.e. $YTP=0$.

During cycle T_2 the test 1 checks channels F_1 of even cells in each row and test 2 -- of odd ones. Channels F_2, F_3 , and F_4 are tested with the same tests in a similar manner during cycles T_3 through T_8 . During the cycle T_9 the channels F_1 , and F_2 of all cells are checked with the selfconjugate (in a pair of intersecting channels) test 1100 (Fig.9,5); the rest of the channels (F_3, F_4) is fed with the auxiliary test 1101. During cycles T_{10} , and T_{11} the channels F_1 and F_2 are checked with the selfconjugate (in one channel) test 4 from Table 1 (Fig.9,6). During cycles T_{12}, T_{13}, T_{14} the other pair of channels (F_3 and F_4) is checked with same tests 3 and 4.

There is no one period of conjugated tests for the last test (0101) in Table 1, since input signals $B=D=1$ block the channel, i.e. the output of the cell becomes independent of its input. Such being the case, a test is used which is conjugate (without a period) and intercompatible with the given one in the i -th channel, and has conjugation period in the j -th channel. Thus, the "out of order" signal from the i -th channel will appear at the external output of the j -th channel. When, during

the cycle T_{15} , test 0101 is fed into channels F_1 , of odd cells in odd rows (Fig.9,2), the channels F_2 of even cells in all rows are fed with tests 1110 and 0110 (1-st and 2-nd in Table 1) which are interconjugate (in any channel, including F_2). Test 1110 is conjugate and intercompatible in the channel F_1 with the given test 0101. The odd cells in even rows are fed with any test which is intercompatible with the test 0110 in the channel. The "out of order" signal for the channel F_1 (0 instead of 1) will appear at the external output of the channel F_2 . Channels $F_1 \div F_4$ of all other cells are checked in a similar manner with test 0101 during cycles T_{16} through T_{30} .

The checking time may be reduced if the array channels are tested with tests 1, 2, and 5, as shown in Fig.9,2. During a single cycle the test 0101 is fed into the channels F_1 , of odd cells in odd rows, and the channels F_2 of even cells in odd and even rows are fed with interconjugate tests 1110 and 0110. Then the whole of the array, irrespective of its size, may be tested during 22 clock cycles. It is clear that at the same time all X-generating cell networks are tested.

It stands to reason that the checking time N' does not change if inputs B, C , and D_1 through D_4 of each cell are not external (as it was supposed before), but are connected to coordinate control buses as shown in Fig.10 for one row.

Summary

The uniform array checking and fault diagnosis (and consequently checking and diagnosis of any finite automaton (sequential function) realized by the array) are drastically simplified as compared with checking of non-uniform circuits, since they are reduced to checking and diagnosis of identical cells having regular interconnections. The number of checking and diagnostic procedures by the methods discussed above and time of array testing (number

of working cycles) are nearly or completely independent of the number of cells (array size), i.e. of complexity of the function realised by the array.

References

1. Прангшвили И.В. и др. "Микроэлектроника и однородные структуры для построения логических и вычислительных устройств", "Наука", Москва, 1966.
(I.V.Pranghishvily et al. "Microelectronics and Uniform Arrays for Logic and Computation", "Nauka" Publishers, Moscow, 1966. (In Russian)).
2. Чегис И.А., Яблонский С.В. "Логические способы контроля электрических схем", Труды математического института им. Стеклова, № 51, 1958.
(I.A.Cheghis, S.V.Yablonsky, "Logical Checking of Electrical Networks". Trudy Matem.Inst.in. Steklova, No.51, 1958). (In Russian).
3. Карибский В.В. и др. "Техническая диагностика комбинационных устройств" в сб. "Абстрактная и структурная теория релейных устройств", "Наука", 1966.
(V.V.Karibsky et al. in "Abstract and Structural Theory of Relay Devices", "Nauka" Publishers, 1966). (In Russian).
4. R.H.Canaday. "Two-Dimensional Iterative Logic".
Proc.of the 1965 FJCC, AFIPS, Vol.27, Part I.

ON DESIGN OF TESTS FOR DIGITAL
DEVICES WITH DELAYSVolkov A.F., Vasil'ev I.N.,
Vedeshenkov V.A., Petrov V.A.Institute of Automatics and Telemechanics
(Technological Cybernetics), Moscow

U S S R

INTRODUCTION

A direct application of the conventional technique of tests construction¹ by tables of faults function leads to labor consuming cumbersome methods². However, when the number and nature of faults that exist simultaneously and the action of contactless elements is known to be unidirectional, more effective modular techniques can be used to synthesize unitary tests.

We intend to discuss development and modification of tests for sequential units without memory elements such as flip-flops and feedbacks. These were chosen for analysis because:

- a) a major portion of a control unit circuits satisfy the requirements of such a model;
- b) the transitional nature of the model analysed (from combinational units to units with memory elements) makes it possible to resort to tests construction techniques used for combinational elements^{3,4} and allow for specific features of sequential units;
- c) if memory elements are included, the algorithm will have to be altered rather than invalidated.

The functioning of an element without fixation of the internal state can be described by the equation

$$W_t = f(u_1, u_2, \dots, u_m)_{t-2}, \quad (1)$$

where W is a variable at the output of the element;

u_1, \dots, u_m are the corresponding variables at the inputs of the element;

τ a delay between changes in u_1, \dots, u_m and W .

Any possible fault will change the logic or time response of the element and can only be detected if there is such a set of variables u_1^0, \dots, u_m^0 and such a time instant t^0 for which the test function turns to unity

$$P_{t-\tau} = W(u_1^0, \dots, u_m^0)_{t-\tau} \oplus W^*(u_1^0, \dots, u_m^0)_{t-\tau} = 1, \quad (2)$$

where \oplus is a symbol of addition for modulo 2,

W^* is an element operation equation for a given fault. The tests of separate elements are best developed experimentally by introducing faults into the elements and obtaining test tables.

The element checked is assumed to be connected to an output terminal of the unit through a group of elements whose operation is described as

$$y_t = F(z_{1,t-\tau_1}, \dots, W_{t-\tau}, \dots, z_{n,t-\tau_n})_{t-\tau_0}. \quad (3)$$

Then the function of the test for detection of faults in this element for output y can be represented as an intersection of two functions:

$$R_{t-\tau_0} = P_{t-\tau-\tau_0} \cdot \varphi_{w,t-\tau_0}^y. \quad (4)$$

The function of the test $P_{t-\tau-\tau_0}$ from (2) defines the conditions for detection of the fault when the output of the faulty element is being checked.

The function $\varphi_{w,t-\tau_0}^y$ determines the conditions for transfer of the control point from the output W of the element to the output y which is linked to the

output W directly or through a number of intermediate elements that are connected arbitrarily. This is the function of the output y sensitivity to changes in the variable W ⁴. Formally the sensitivity function is determined as

$$\varphi_{w,t-\tau_0}^y = F(z_1, t-\tau_1, \dots, W_{t-\tau}^0, \dots, z_n, t-\tau_n)_{t-\tau_0} \oplus F(z_1, t-\tau_1, \dots, \overline{W_{t-\tau}^0}, \dots, z_n, t-\tau_n)_{t-\tau_0}, \quad (5)$$

where

$$W_{t-\tau}^0 = W(u_1^0, \dots, u_m^0)_{t-\tau}.$$

By equation (5) one can find the sensitivity functions for elements of any type. E.g. for elements OR (NOR), AND (NAND), the sensitivity function for the i input will be (for n input elements)

$$\begin{aligned} \varphi_i^{OR} &= \varphi_i^{NOR} = \overline{x_1} \cdot \overline{x_2} \cdot \dots \cdot \overline{x_{i-1}} \cdot \overline{x_{i+1}} \cdot \dots \cdot \overline{x_n}, \\ \varphi_i^{AND} &= \varphi_i^{NAND} = x_1 \cdot x_2 \cdot \dots \cdot x_{i-1} \cdot x_{i+1} \cdot \dots \cdot x_n. \end{aligned} \quad (6)$$

For a follower, inverter, the delay line $\varphi \equiv 1$.

As follows from the formula a solution to the equation enables to find such values of variables at other inputs of the element (circuit) at which the values of y depend unambiguously on the values of W :

$$y_t = W_{t-\tau-\tau_0}.$$

or

$$y_t = \overline{W_{t-\tau-\tau_0}}.$$

We will say that there is a sensitive path between the terminals W and y if changes in the value of W leads necessarily to changes in the value of y . Eq. (5) determines the conditions for the existence of a sensitive path between W and y . A solution to Eq. (5) would yield these conditions, but this may appear too cumbersome when the size of the circuit increases.

Therefore we undertook to analyze the relation between the structure of the device which implements equation (3) and properties of the sensitivity function. That analysis has

shown that the sensitivity function of a sequence of elements connected in series is the product of the sensitivity functions of all the elements in the sequence. Therefore, the conditions for existence of a sensitive path through a sequence of elements connected in series can be achieved by a consistent unification of conditions for existence of separate parts of the path which are associated with each element of the succession.

If there are several paths for a signal from the terminal W to propagate to the terminal Y along, the sensitivity path chosen may prove impossible to realize. Therefore an essential part of construction of a unitary test suggested here is an analysis of conditions for realizability of the sensitive path chosen.

11. A modular technique for development of a unitary test in case of sequential device with delay elements

The algorithm suggested employs local characteristics of separate elements making up a sequential unit: tests tables, operational equations, magnitudes of delay, sensitivity functions, tables of links to neighbouring elements through inputs and outputs. Its steps are:

1. Ranking of the circuit.
2. Making up a primary tests table.
3. Finding the sensitive path.
4. Finding and analysis of parallel paths.
5. "Widening" of the sensitive path.
6. Making input variables consistent.
7. Making the times when values of variables change consistent.
8. Minimization of the test total duration.

Each step will be discussed for the case of a unit whose block diagramme is shown in the Figure. The unit is made of conjunction (\wedge_i) disjunction (\vee_j), inverters (J_{n_k}) and delay, whose magnitude (τ) is taken arbitrarily and shown in relative units.

1. The first part of the algorithm is to denote the

row chosen is rewritten in the row numbered $(a_i + \tau_i)$ the secondary table while shifting the value of the output y_j by τ_j rows downwards. If y_j is not an output terminal of the unit, then we proceed to construct the sensitive path by finding in the primary table the sensitivity function of the element which follows y_j .

Table 2 shows how a version of the sensitive path is constructed for the second row of Table 1.

TABLE 2

No. of terminal of time unit	; ; ; ; ; ; ; ; ; ; ; ; ; ; ;															Note
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	
1				1												2
2																(a)
3																
4																
5					0	0										6
6																14
7																
8																
9																
10					1				0	1						0
11																
12																
1				$\frac{1}{1^x}$												2
2				$\frac{1}{1^x}$												(b)
3																
4		0^x						1^x								
5		$\frac{0}{0^x}$				$\frac{0}{0^x}$		$\frac{0}{0^x}$								6
6		$\frac{0}{0^x}$		1^x		$\frac{0}{0^x}$		1^x								10
7				$\frac{1}{1^x}$		1^x		$\frac{0}{0^x}$	1^x	1^x						
8				$\frac{1}{1^x}$		$\frac{1}{1^x}$		$\frac{0}{0^x}$	$\frac{1}{1^x}$	$\frac{0}{0^x}$						
9						$\frac{1}{1^x}$			$\frac{1}{1^x}$	$\frac{0}{0^x}$				1^x		
10						$\frac{1}{1^x}$			$\frac{0}{0^x}$	$\frac{0}{0^x}$				$\frac{0}{0^x}$		14
11														$\frac{0}{0^x}$		0
12														$\frac{0}{0^x}$		

4. Let us say that between terminals i and j of the circuit there are loops (groups of elements connected in parallel) if one can construct two or more paths along which the signal propagates from i to j and differing by at least one element. Such terminals in the circuit of Figure are 8 and 13, 9 and 15. A linking element will be the one whose output terminal coincides with the final point of the loop. For each branch of a loop one can determine:

a) whether the number of signal ℓ_i inversions is even (the phase of signal propagation);

b) total delay time of signal T_i , the index i relating to that branch for which a sensitive path was constructed during the third step of the algorithm.

If $\ell_1 = \ell_2$, the two branches of the loop will be termed synphase, in the opposite case, counterphase. After a loop has been found with a sensitive path passing through a branch, the difference $(T_1 - T_2)$ of the summed signal delays has to be computed for the branches. If the difference is zero, the phase relations of the signals transmitted $(\ell_1 - \ell_2)$ has to be found. If the branches of the loop are in phase we have to find the function $\varphi_{1,2}$ of the link element sensitivity to a simultaneous synphase change of its two input variables. If $\varphi_{1,2} \neq 0$ the sensitive path constructed has to be adjusted. For linking elements of the types AND (NAND), OR (NOR) this reduces to the value the variable obtained at the output of the first branch being rewritten to the column which corresponds to the variable at the output of the second branch.

If the branches are counterphase, or $\varphi_{1,2} \equiv 0$, it is necessary to prevent the signal from transmission through the second branch. For this purpose each branch should contain fixation elements or elements that make it possible to fix at the output of the given branch those values which are determined by the sensitivity function of the linking element.

At $|T_1 - T_2| = \Delta T > 0$ the branches can be assumed non-shunting and the third step of the algorithm has to be

repeated for the sensitivity path, crossing the second branch of the loop.

Table 2b differs from Table 2a in that records have been added that describe the sensitivity path through terminals 8, 9, 11, 13 with time $T_2 = 9 < T_1 = 11$. The records which describe the first and the second sensitive path are underlined once or twice, respectively. Besides, 1 in the cell of the 11th column and 10th row is replaced by 0 taken from the 8th row to simplify the stage of consistency. (The numbers with asterisks will be explained below).

5. Changes in the quantities τ_i of delay elements caused by faults will make the correct value of the output variable appear at the terminal checked at time t which does not coincide with the desired time t_0 . When the output variable $f(t)$ values at time $(t_0 - 1)$ are fixed and that value is compared with the quantity $f(t_0)$, we may check whether the magnitudes τ_i are correct for the delay elements on the given sensitive path.

A test which characterises the sensitive path constructed can be modified so that not only a qualitative reply (yes or no) to the question, whether the delay elements are in order or not, can be obtained, but also we will be able to quantitatively estimate the total change of delays. To do this we will "expand" the sensitive path.

Assume that the number of time units d is the maximal permissible deviation of the total delay. Then starting with the row numbered $(a_i + \tau_i)$ we will repeat the records which correspond to sensitivity functions of separate elements d times upwards and d times downwards; in the lower rows the columns will have the values of y_i^0 and in the upper lines the values of $\underline{y_i^0}$. The columns of input variables which are related to d rows above the a_i -th row will be left blank, while in the same columns of the lower rows we will write the values of variables from the a_i -th row.

In Table 2b this construction was completed for a sensitivity path with a shorter delay time at $d = 1$. (The

values added at this stage are marked with asterisks (*).

The values of variables which expand the sensitive path fix the time t_0 when the values of variables at the output of each delay elements change (the boundaries between y_i^0 and y_i^1) and maintain the constant values of variables that describe the sensitive path during the interval $(t_0 - d, t_0 + d)$. As a result when τ_i of one of delay elements on the sensitive path changes by the magnitude $|\pm \delta \tau_i| < d$ the time t of the output value change will be shifted with respect to time t_0 to by the same value $\pm \delta \tau_i$. A comparison of the value

$f(t_0 + d)$ with the value which corresponds to the operative unit gives an answer to the question: "Do the changes in the delay magnitude exceed the permissible limits?" If the values coincide, the answer is no, other wise, yes, they do.

6. The sensitive path which results from the previous parts of the algorithm depends on the values of input variable elements which form the path. Some of these variables are outputs of other elements and thus cannot be fed from outside. The sixth part of the algorithm is to make consistent the values of input variable unit which are essential to feed the appropriate values of variables to the input of elements which form the sensitive path. Since for most logical elements there are no unambiguous relations between input and output variables the process of specification can lead to inconsistent result and a trial of a certain number of versions will be required until a satisfactory result is obtained.

The values of all variables that are required for a sensitive path are made consistent starting from the element which corresponds to a filled column with the greatest number in the row with the greatest number of time units. After input variables have been made consistent for remaining elements of the same row, in the decreasing order of their numbers, one can proceed to the elements of the nearest row which describes the value of at least one variable with the number $(n+1, \dots, n+m)$. Thus at the stage of

making variables consistent the secondary table is looked through from the right to the left and upwards until there are no values of variables with the numbers $n+1, \dots, n+m$, which are not determined by the values of variables with the numbers $1, 2, \dots, n$.

To choose one possible version of making the variables consistent one can resort to a relatively simple technique at which the chosen input variables coincide with the values of the same variables in the nearest upper or lower rows of the secondary table. Table 3 represents the results of making the variables consistent for the data of Table 2b whose numbers of rows are given in the first column in parenthesis.

TABLE 3

No. of terms of min: time unit	1 : 2 : 3 : 4 : 5 : 6 : 7 : 8 : 9 : 10 : 11 : 12 : 13 : 14 : 15														
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1				0											
2				0											
3				0											
4(1)				1											
5(2)				1			1								
6(3)	1			1			1								
7(4)	1	0					1								
8(5)	1	0			0		0								
9(6)	1	0			1		0	1							
10(7)		0			1		0	0	1	1					
11(8)					1			0	1	0					
12(9)								0	1	0		1			
13(10)									0	0		0			
14(11)													0		
15(12)														0	

7. Among the columns of the secondary table let us find those that contain 0 and 1 with a group of blanks between

them $x/$ and fill out the half of cells closer to 0 with zeroes and the half closer to 1 with ones. To avoid contradictions, the blanks will be filled in this order:

a) fill out one of the intervals by the above rule; the next interval will be the one which is in the column with the greatest number;

b) make the input variables consistent.

8. The secondary table resulting from the making the values of input variables consistent and times when these are fed contains a great number of blank cells. Their content is of no importance for correct realization of the part of the test which is related to the sensitive path chosen. Therefore blank cells of a secondary table can be filled for various purposes in particular to minimize the total length of the test when separate parts of the test are linked. Efficiency of various minimization techniques depends on the structure of a unit which feeds the sets of a test. Therefore we will not deal with minimization in this paper.

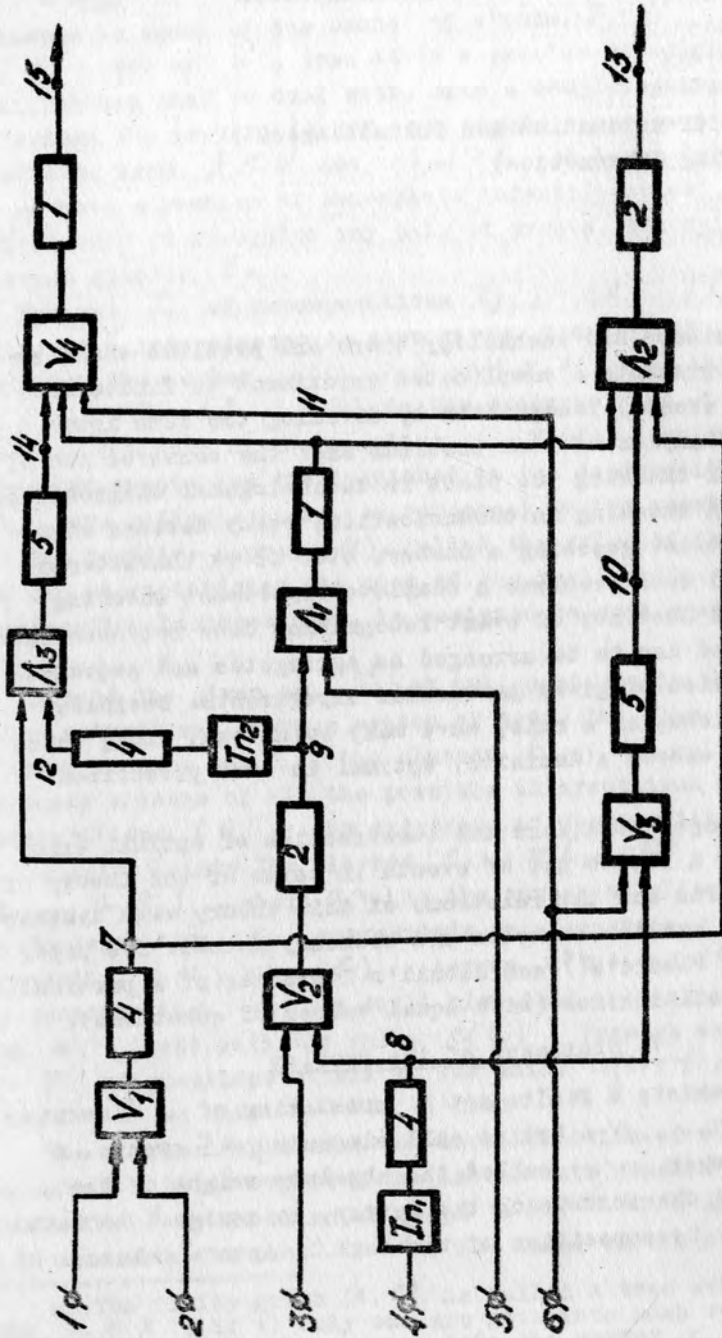
The above algorithm for construction of unitary tests for sequential digital devices yields tests which guarantee detection of faults that cause changes in logical or temporal characteristics of one element of the device at a time. An advantage of the algorithm is that construction of the test does not require a faults function table for the entire device. This results in a comparative simplicity and ease with which the algorithm is realized on computers.

REFERENCES

1. Черис И.А., Яблонский С.В. Логические способы контроля работы электрических схем. Труды Математ.ин-та им.Стеклова В.А., т.51, 1958, стр. 270-359.
2. Gill A., Introduction to the theory of finite-state Machines, McGraw-Hill Book Company, 1963.
3. Roth J.P., Diagnosis of automata failures: a calculus and a method. IBM J. of Research and Development, July 1966, vol. 10, No. 4, pp. 278-291.

$x/$ This may be the case for a sensitive path containing a loop with counterphase branches.

4. Ведешенков В.А. "О построении таблицы тестов для обнаружения логических неисправностей бесконтактных комбинационных устройств". Автоматика и телемеханика, том XXIX, № 3, 1968.
5. Ведешенков В.А., Волков А.Ф. "Блочный метод построения единичного проверяющего теста для бесконтактных комбинационных устройств без обратных связей". Тезисы докладов научно-технического семинара МДНТИ им.Ф.Э.Дзержинского "Проектирование радиоэлектронной аппаратуры с использованием вычислительной техники", июнь 1967, Москва.



Block-diagram of the sequential device analysed.

THE THEORY OF QUESTIONNAIRES AND PROBLEMS OF
TECHNOLOGICAL DIAGNOSTICS

P.P. Parkhomenko

Institute of Automation and Telemechanics
(Engineering Cybernetics)

Moscow

USSR

In science and technology there are problems whose solution represents a complicated experiment in finite sets elements (events) recognition by dividing the sets into classes. Examples of such problems are: the construction of programs of checking the plant in technological diagnostics, information encoding in communication, relay devices structure synthesis, guessing a number, etc. It is characteristic that in such problems a complete experiment ensuring the desired accuracy of event recognition does not necessarily exist and has to be arranged as aggregates and sequences of realization of given particular experiments. Besides, these problems, as a rule, have many solutions, among which one should choose a decision, optimal in some prescribed sense.

This paper considers the construction of optimal experiments for a finite set of events in terms of the theory of questionnaires and the relations of this theory with dynamic programming and the branches and boundary method. The paper is based on Picard's ¹ and Dubail's ² studies of a particular kind of questionnaires (with equal values of questions).

1. Statement of the Problem

There exists a finite set E consisting of elements y_i , $i = 1, 2, \dots, N$. Let us call elements $y \in E$ events. A positive number $w(y)$ called the absolute weight of the event y and characterizing the latter, is assigned to each event $y \in E$. Decomposition of the set E into λ classes of

E_μ , $1 \leq |E_\mu| < N$, $\mu = 1, 2, \dots, \lambda$ is recorded. If the number of classes is equal to the number of events, $\lambda = N$, and $|E_\mu| = 1$ for all μ , then it is a problem of complete identification that we deal with, when a complete experiment must ensure the recognition of each single event among all the others. When $\lambda < N$ and $|E_\mu| > 1$ at least for one μ , we have a problem of incomplete identification, which requires only to recognize any pair of events, belonging to different classes E_μ .

The set T_i of decompositions t_j , $j = 1, 2, \dots, |T_i|$ of the set E into classes is also given. Let us call $t \in T_i$ questions. The number $a(t)$, $2 \leq a(t) \leq N$ of the classes $E_{j(t)}$, $j(t) = 1, 2, \dots, a(t)$ in the decomposition $t \in T_i$ is called the base of the question. Features by which classes of events are distinguished in the decomposition $t \in T_i$ are called answers (or outcomes) to the question t .

The positive number $c(t)$, called the value of the question and characterizing the cost of the realization of the corresponding decomposition, is assigned to each question $t \in T_i$.

Expand the given set T_i of the questions in the following natural way. Form a system of sets, the elements \mathcal{E} of which are the set E , the classes $E_{j(t)}$ as well as non-empty classes of all the possible intersections of the decompositions $t \in T_i$. An aggregate of decomposition τ_j of the sets \mathcal{E} into the classes $\mathcal{E}_r(\tau_j) = E_{r(\tau_j)} \cap \mathcal{E}$, $r(\tau_j) = 1, 2, \dots, a(\tau_j)$, where $a(\tau_j)$ is the number of classes in the decomposition τ_j corresponds to a question $t_j \in T$. Apparently, $a(\tau_j) \leq a(t_j)$. Assume $c(\tau_j) = c(t_j)$. The decomposition τ_j we shall also call a question on the base $a(\tau_j)$ and with the value $c(\tau_j)$. Form an expanded set T of questions of all τ for which $a(\tau) \geq 1$. Since with $\mathcal{E} = E$ we have $\tau_j = t_j$, then $T_i \subseteq T$.

The complex experiment in identification of events of the set E by decomposing the latter into classes can be represented¹ by an oriented graph $G = (Q \cup E, \Gamma)$ of the type of a tree with the root x_0 .²⁾ The inner vertices and the

* The finite graph (X, Γ) is called a tree with the root $x_0 \in X$, if 1) only one arc goes into each vertex $\neq x_0$; 2) not a single arc goes into the vertex x_0 ; 3) the graph (X, Γ) does not contain any loops.

root of the graph form a set Q of questions; the arcs outgoing from the vertex $x \in Q$ are called the outcomes of the question x , while their number $a(x)$, $1 \leq a(x) \leq N$ is called the base of the question x ; the value $c(x) > 0$ is assigned to each question x . In the problems of the complete identification the events $y \in E$ are compared with and the weights $w(y)$ are assigned to the finite vertices of the graph G . In cases of incomplete identification of the events of the set E the finite vertices of the graph G are compared with v , $1 \leq v < N$, of the subsets

$E_k \subseteq E_\mu$, $k = 1, 2, \dots, v$, when $\bigcup_{k=1}^v E_k = E$, and the weights $w(E_k) = \sum_{y \in E_k} w(y)$ are assigned. The graph G is called a questionnaire for E (*).

Concerning the vertex $z \in Z = Q \cup E$ of the graph G we shall distinguish the set Γ_z of its followers, the set $\hat{\Gamma}_z \setminus Z = \Gamma_z \cup \Gamma(\Gamma_z) \cup \dots$ of its descendants the set Γ^{-1}_z of its forerunners and the set Γ^{-1}_z of its ancestors. Instead of the absolute weights $w(y)$ of the events it is more convenient to consider their relative weights $p(y) = \frac{w(y)}{w}$ where $w = \sum_{y \in E} w(y)$, further termed weights. The quantity $p(x) = \sum_{y \in \Gamma_x \cap E} p(y)$ is called the weight of the question $x \in Q$. The sum $c(x, z) = \sum_{x \in \Gamma^{-1}_z} c(x)$ is called the value of the path from x_0 to $z \in Z$. The value of the path from x_0 to $y \in E$ characterizes the costs of identification of the event y .

All the practically interesting characteristics of the experiments, described by questionnaires, can be tackled by these questionnaires. One of the general enough and at the same time practically useful features of a questionnaire is ³ the value of scanning determined by the expressions

$$C(x_0, E) = \sum_{y_i \in E} c(x_0, y_i) \cdot p(y_i) = \sum_{x_j \in Q} c(x_j) \cdot p(x_j) \quad (1)$$

corresponding to the mean weighed costs of identification of events in the entire questionnaire. For the questionnaires with equal values of questions one may obtain ¹ from (1) the length of the scanning, corresponding, for example, to the mean length of a coding combination in Shannon-Fano codes or to

*) Further, if the reverse is not specified, the questionnaires for complete identification problem will be considered.

the average number of operations in sorting problems. If the value of a question is the cost (time) of realization of a separate checking, and the weight of an event is the probability of the operable or inoperable state of the plant then eq. (1) gives us the mean cost (mean time) of determining the states in the conventional sequential checking program.

With the same sets E and T , one can, generally speaking, construct various questionnaires for E , differing in aggregate and in succession of question realization and having different values of scanning. The questionnaire for E with the minimum value of scanning will be called optimal. We shall discuss the features and the methods of constructing optimal questionnaires in this report.

2. Optimal Questionnaires

By transforming the given questionnaire G for E , invariantly relative to the number q_m of its questions with the base a_m , $m \in M$ and the number q_ℓ of its questions with the value c_ℓ , $\ell \in L$ (M, L being some numerical sets), it has been revealed³ that an optimal questionnaire, allowing to identify N events by means of $\sum_{m \in M} q_m = \sum_{\ell \in L} q_\ell$ questions, having bases a_m and costs c_ℓ is a tree with a root x_0 , where $N = \sum_{m \in M} q_m (a_m - 1) + 1$ and such that its vertices, arranged in a non-decreasing order relative to their ranks, have weights assigned to them in non-increasing order and the values of paths leading to them in a non-decreasing order, the questions being arranged in a non-increasing order relative to their costs; among all the vertices having the same value of the paths leading to them as the value of the followers of the question on the base a , there is not a single vertex, whose weight exceeds the sum of the weights of all the other a vertices with the same value of the paths leading to them.

Assume that to construct a questionnaire G for E there must be and there are q_m questions with the base a_m , $m \in M$, among which q_ℓ questions are of the value c_ℓ , $\ell \in L$. Then the following algorithm (let us call it A1) of constructing an optimal questionnaire results from the above definition.

Make list 1 of weights $p(y_i)$, $i = 1, 2, \dots, N$, arranged in a non-decreasing order, and list 2 of pairs (a_j, c_j) of question bases a_j , $j = 1, 2, \dots, \sum_{m \in H} q_m = |Q|$, arranged in a non-decreasing order and of question values c_j , $j = 1, 2, \dots, \sum_{c \in L} q_c = |Q|$ arranged in a non-increasing order. Eliminate the first pair (a_j, c_j) with the smallest base from list 2 and refer it to the question x_j . Eliminate a_j first (smallest) weights from list 1 and assign them to a_j followers of the question x_j . Determine the weight $p(x_j)$ of the question x_j and include it into list 1, and preserve the non-decreasing order of weights there. Return to list 2 and take its first not-eliminated pair (a_j, c_j) , etc. Repeat the described operations until the last pair $(a_{|a|}, c_{|a|})$ with the largest base and least value of a question is eliminated; the last $a_{|a|}$ weights will be eliminated out of list 1.

The described procedure is simple, almost does not require any search through all versions, ensures that the value of scanning of the questionnaire obtained will not exceed the value of scanning of any questionnaire for E , containing q_m questions with the base a_m and q_c questions of the value c_c .

A questionnaire, in which each question from Q is a question from the given set T of questions, i.e. for each question $x \in Q$ a question $t \in T$ can be found, such that $a(x) = a(t)$, $c(x) = c(t)$ and $\hat{I}_x \cap E = E_{x(t)}$ for each $x \in T$, or briefly, a questionnaire in which $Q \subseteq T$ will be called a realizable questionnaire. It is clear that the algorithm A1 does not always lead to an optimal realizable questionnaire. The condition $Q \subseteq T$ is always met for the problems, in which the set T contains all the possible decompositions of the set E into $a(t) = a = \text{const}$ classes, the values of all decompositions being the same, $c(t) = c = \text{const}$. The problems of encoding and sorting mentioned above are examples of such problems. It is not difficult to show that the algorithm of constructing optimal questionnaires with equal values of questions¹, representing a generalization of the well-known⁴ algorithm of constructing redundants codes, follows from the algorithm of constructing optimal ques-

tionnaires with unequal bases and values of questions.

The construction of optimal realizable questionnaires with restrictions existing for the given set T of questions is discussed below. We shall note here, that if from the given set of questions we choose $\sum_{m \in H} q_m$ questions with as large bases as possible and $\sum_{e \in L} q_e$ questions with as low values a_e as possible, then the value of scanning of the optimal (not necessarily realizable) questionnaire, constructed according to the algorithm described, represents the lower boundary of the value of scanning the questionnaires for E . This lower boundary is comparatively easy to calculate in the process of implementing the algorithm A1 operation and is achievable, when the condition $Q \subseteq T$ is met for the constructed optimal questionnaire. When the lower boundary is known it is possible to estimate the quality of questionnaires obtained by "approximated" methods and to construct optimal realizable questionnaires by the branches and boundaries method.

3. Recurrent Calculation of the Value of Scanning

Let us denote questionnaire vertices by the two indices - $[z, s]$ where z is the rank of a vertex and s determines its location among other vertices of the rank z . Let $G = (Z, \Gamma)$ be a questionnaire for E and $Z_{z,s}$ - its vertex. The subquestionnaire with the root $Z_{z,s}$ of the questionnaire G is a graph, $G_{z,s} = (Z_{z,s}, \Gamma_{z,s})$ where $Z_{z,s} = \hat{\Gamma}_{Z_{z,s}} \subseteq Z$ and the reflection $\Gamma_{z,s}$ is determined in the following way: $\Gamma_{z,s}z = \Gamma z \cap Z_{z,s}$. The subquestionnaire $G_{z,s}$ is a questionnaire for $E_{z,s} = \hat{\Gamma}_{Z_{z,s}} \cap E$.

If the vertex $z_{w,t} \in Z_{z,s}$, then its weight in the questionnaire G is $p(w, t)$ while its relative weight in the questionnaire $G_{z,s}$ is equal to $p_{z,s}(w, t) = \frac{p(w, t)}{p(z, s)}$. Assume that besides the subquestionnaire $G_{z,s}$, there is a questionnaire $G'_{z,s}$ for $E_{z,s}$. The operation of substituting subquestionnaires in the questionnaire G for E determines a new questionnaire, G'' , for E , obtained from by replacing the subquestionnaire $G_{z,s}$ by the questionnaire $G'_{z,s}$ with further re-calculation of the vertices weights of the latter by multiplying them by $p(z, s)$.

By eq. (1) for the value of the scanning $G_{z,s}$ we have:

$$C_{z,s} = C(Z_{z,s}, E_{z,s}) = \sum_{y_i \in E_{z,s}} c(Z_{z,s}, y_i) \cdot p_{z,s}(y_i) \quad (2)$$

Instead of (2) we may write the following:

$$C_{z,s} = \sum_{y_i \in E_{z,s}} [c(Z_{z,s}, y_i) - c(Z_{z,s})] p_{z,s}(y_i) + \sum_{y_i \in E_{z,s}} c(Z_{z,s}) \cdot p_{z,s}(y_i) \quad (3)$$

Noting that $E_{z,s} = \bigcup_{n=1}^{a(Z_{z,s})} E_{z+1,n}$;

$$p_{z,s}(y) = p_{z+1,n}(y) \cdot p_{z,s}(z+1, n); \quad c(Z_{z+1,n}, y) = c(Z_{z,s}, y) - c(Z_{z,s});$$

$$\sum_{y \in E_{z+1,n}} p(y) = p(z+1, n); \quad \sum_{z+1,n \in \Gamma Z_{z,s}} p(z+1, n) = p(z, s),$$

and taking into account (2), we obtain the following from eq. (3):

$$C_{z,s} = c(Z_{z,s}) + \sum_{z+1,n \in \Gamma Z_{z,s}} p_{z,s}(z+1, n) \cdot C_{z+1,n} \quad (4)$$

Therefore, the value of scanning $C_{z,s}$ of the subquestionnaire $G_{z,s}$ of the rank z is equal to the sum of its root's $Z_{z,s}$ value $c(Z_{z,s})$ and the sum of the values of the scanning $C_{z+1,n}$ of the subquestionnaires $G_{z+1,n}$ of the rank $z+1$, whose roots $Z_{z+1,n}$ are the followers of the root $Z_{z,s}$.

If $Z_{z,s} = y \in E$, then $\Gamma Z_{z,s} = \emptyset$ and it follows from eq. (4) that

$$C_{z,s} = c(Z_{z,s}) = c(y) = 0 \quad (5)$$

Formulae (5) and (4) allow to calculate recurrently the value of scanning every subquestionnaire $G_{z,s}$ of the questionnaire G , including the value of the scanning G itself:

$$C = C(x_0, E) = c(x_0) + \sum_{z,n \in \Gamma x_0} p(z, n) \cdot C_{z,n} \quad (6)$$

Let us prove, following Ref. 2, the following statement: an optimal questionnaire consists of optimal subquestionnaires.

Let G be an optimal questionnaire. Suppose that in G there is a non-optimal subquestionnaire $G_{z,s}$ with the value of scanning $C_{z,s}$. Then there exists an optimal questionnai-

re $G'_{z,s}$ for $E_{z,s}$ with the value of scanning $C'_{z,s}$ and

$$C'_{z,s} < C_{z,s} \quad (7)$$

Let z be the minimum rank of non-optimal subquestionnaires in G . If $z = 0$ then it is the questionnaire G that is non-optimal, which is impossible according to our condition. Assume that for all $w < z$, $0 < z \leq z_n$ (z_n being the maximum rank of a questionnaire) all subquestionnaires $G_{w,s}$ in G are optimal. Take a subquestionnaire, $G_{z-1,t}$, with the root $z_{z-1,t} \in \Gamma^{-1} z_{z,s}$. This subquestionnaire is optimal and, according to eq. (4), its value of scanning is

$$C_{z-1,t} = C(z_{z-1,t}) + \rho_{z-1,t}(z,s) \cdot C_{z,s} + \sum_{z_{z,n} \in \Gamma z_{z-1,t} \setminus z_{z,s}} \rho_{z-1,t}(z,n) \cdot C_{z,n} \quad (8)$$

In the questionnaire $G_{z-1,t}$ substitute the subquestionnaire $G'_{z,s}$ for $G_{z,s}$. For the new questionnaire $G'_{z-1,t}$ obtained we have:

$$C'_{z-1,t} = C(z_{z-1,t}) + \rho_{z-1,t}(z,s) \cdot C'_{z,s} + \sum_{z_{z,n} \in \Gamma z_{z-1,t} \setminus z_{z,s}} \rho_{z-1,t}(z,n) \cdot C_{z,n} \quad (9)$$

Comparing eq. (8) and eq. (9) we obtain by virtue of eq. (7):

$$C'_{z-1,t} < C_{z-1,t} \quad (10)$$

which is incorrect since $G_{z-1,t}$ is optimal by definition. The statement has been thus proved and it is also true for optimal realizable questionnaires.

4. Optimal Realizable Questionnaires

Let \bar{E} be a set of all possible different non-empty subsets E_j of the set E , containing $m_j \leq N$ events. When $m_j \geq 2$ we shall call the given decomposition of the subset E_j into $a(t)$, $2 \leq a(t) \leq m_j$ classes $E_{j,n}$, $n = 1, 2, \dots, a(t)$ a question $t \in \mathcal{T}$ feasible for E_j . Denote the set of questions feasible for E_j by \mathcal{T}_j . If $m_j = 1$, we assume $\mathcal{T}_j = \emptyset$. Then $\mathcal{T} = \bigcup_{E_j \in \bar{E}} \mathcal{T}_j$ is a set of feasible questions, which is called a compatible set if for it there is at least one realizable questionnaire for E . Let us call the pair (E_j, \mathcal{T}_j) a situation and m_j elements in E_j - the order of a situation.

It has been mentioned above that if there are constraints for the set T of the given questions the algorithm A1 may give an optimal realizable questionnaire only by chance, when all the questions of at least one of all versions possible according to the algorithm turn out to be feasible. Let us state the general algorithm A2 for construction of optimal realizable questionnaire, based on the procedure of recurrent calculation of the value of scanning a questionnaire.

Let \mathcal{G}_2 be a set of all realizable questionnaires for E and for the compatible set T of feasible questions and $S\mathcal{G}_2$ - a set of all (realizable) subquestionnaires of the above questionnaires from \mathcal{G}_2 .

Let us call the situation, to which the realizable subquestionnaire from $S\mathcal{G}_2$ corresponds, a possible situation. Let us call the questions of the realizable questionnaires from \mathcal{G}_2 possible feasible questions. Denote the set of feasible questions by T' .

Assume (\mathcal{E}_j, T'_j) to be a possible situation of the order m_j . If (\mathcal{E}_j, T'_j) is a possible situation of the first order, then the realizable subquestionnaire corresponding to it is a degenerated π) one and therefore optimal with the value of scanning

$$C_0(\mathcal{E}_j, T'_j) = C_0(y_i, \phi) = 0, \quad j = 1, 2, \dots, N \quad (11)$$

Now assume that for every possible situation (\mathcal{E}_k, T'_k) of the order m_k , $1 \leq m_k < m_j \leq N$ there is a corresponding optimal realizable subquestionnaire already constructed. Take a possible situation, (\mathcal{E}_j, T'_j) . The question t feasible for \mathcal{E}_j with the base $a(t)$ is a decomposition, \mathcal{E}_j , into $a(t)$ non-intersecting and non-empty subsets $\mathcal{E}_{j,n}$, each of which has the number of elements $m_{j,n} < m_j$. By our condition, an optimal realizable subquestionnaire has already been constructed for each of the possible situations $(\mathcal{E}_{j,n}, T'_{j,n})$. Then the subquestionnaire with the minimum value of scanning will be an optimal realizable subquestionnaire. This value is given by

π) We call questionnaires for one event degenerated ¹.

$$C_0(\varepsilon_j, T_j') = \min_{t \in T_j'} [c(t) + \sum_{n=1}^{a(t)} p_n C_0(\varepsilon_{j,n}, T_{j,n}')] \quad (12)$$

where

$$p_n = \frac{\sum_{y_i \in \varepsilon_{j,n}} p(y_i)}{\sum_{y_k \in \varepsilon_j} p(y_k)}$$

It is easy to observe that recurrent relation (12) is a functional equation of dynamic programming with the operator of the minimum⁵. Thus, the procedure of constructing

optimal realizable questionnaires by the algorithm A2 consists in examination of all possible situation of the order m_j , from $m_j = 1$ to $m_j = N$, and in constructing the corresponding optimal realizable subquestionnaires. The questionnaire corresponds to the possible situation (E, T') .

The algorithm A2 allows to construct realizable questionnaires, optimal not only in terms of the minimum value of optimization; the algorithm makes it possible to obtain either one or all the decisions, it can be easily modified for for the solution of incomplete identification problems. The drawback is the great number of calculating operations, determined in the worst case by $2^N \cdot |T'|^*$ references to eq. (12). Let us turn to the algorithm A3, representing a procedure of the branches and boundaries method⁶ and using the algorithm A1 for calculating the lower boundary of the value of scanning a questionnaire. In the version described the algorithm A2 constructs the desired questionnaire from its finite vertices towards the root x . We shall describe algorithm A3 in a version, where a questionnaire is being constructed, on the contrary, from its root x towards its finite vertices.

Assume that there is a tree with the root x , realizing the questionnaire G for E . In G let us single out a subtree $G_2 \subseteq G$ with the root x_2 , belonging to the set Q_2 of the inner vertices x_2 , and with a set E_2 , of finite vertices, $y_2 \in E_2$. Apart from the subtree G_2 the

*) In a number of cases it is possible to reduce the number of operations by determination of all possible situations in advance.

tree G contains subtree G_{y_2} with the roots $y_2 \in E_2$, realizing subquestionnaires G_{y_2} for E_{y_2} , corresponding to them (and including, perhaps, degenerated ones), while $\bigcup_{y_2 \in E_2} E_{y_2} = E$. The value of scanning the questionnaire G can be represented in the following way:

$$C(x_0, E) = \sum_{x_2 \in Q_1} p(x_2) \cdot C(x_2) + \sum_{y_2 \in E_2} p(y_2) C_{y_2} \quad (13)$$

where $p(y_2) = \sum_{y \in E_{y_2}} p(y)$ and C_{y_2} is the value of scanning the subquestionnaire G_{y_2} .

Let us determine the initial possible situation ($\varepsilon_j = E$, $\tau'_j = \tau'_1$) of the order N . Any question $t \in \tau'_1$ may turn out to be the first question x_0 of the optimal realizable questionnaire sought for, i.e. there are $|\tau'_1|$ versions (branches of the tree of solutions). Let us determine possible situations ($\varepsilon_{j,n}$, $\tau'_{j,n}$) for each version, by the algorithm A1 construct optimal (non-realizable in the general case) subquestionnaires, corresponding to each of these situations and find the value of scanning the latter $C_0(\varepsilon_{j,n}, \tau'_{j,n})$. Then by eq. 13, determine the value of scanning the entire questionnaire for each version and for the continuation of the solution choose the version that gives:

$$C(x_0, E) = \min_{t \in \tau'_1} \left[C(t) + \sum_{n=1}^{a(t)} p_n \cdot C_0(\varepsilon_{j,n}, \tau'_{j,n}) \right] \quad (14)$$

At the next steps the operations described above are repeated conformably to the possible situations, chosen at the previous steps. If the versions of the given step give a higher value than any of the versions of the preceding steps, then it is chosen for the continuation of the solution. The first realizable questionnaire for E obtained in the course of solution will be the optimal one.

In the majority of cases the branches and boundaries method allows to reduce the volume of calculating operations as compared with the dynamic programming procedure and is as general as the latter.

5. On "Approximated" Method of Constructing Questionnaires

The set E of events can practically always be represented as a complete system of events, whose probabilities are equal to their weights. Hence, the conventional entropy of the rank τ is rather often used in the questionnaire construction procedures as a function of preference in choosing the solution at every step. This corresponds to maximization at every step of the function

$$H = \sum_{n \in T_{z,s}} p_{z,s}(\tau+1, n) \cdot \log_{a_{z,s}} \frac{1}{p_{z,s}(\tau+1, n)} \quad (15)$$

i.e. to the choice of such a decomposition of the subset $E_{z,s}$, which makes it possible to obtain the closest to each other weights of the questions $[\tau+1, n]$, which are the followers of the question $[\tau, s]$.

Another "approximated" procedure of questionnaires construction, also often used in practice, is a procedure, known as Shannon-Fano method for codes with variable length of encoding combinations. By this procedure the best decomposition at every step is such that allows to obtain the weights of the questions $[\tau+1, n]$, closest to each other, provided that the weights of the events are arranged in a non-decreasing (or non-increasing) order.

The investigation of the features of the questionnaires, optimal in terms of the minimum value of their scanning, as well as of the questionnaires, optimal from the point of view of information reveals that neither of the preference functions mentioned guarantees that optimal questionnaires will be obtained even in case ^{1,2} of questionnaires with equal values of questions. In the general case the principle of the conventional entropy maximization gives decisions, which are farther from being optimal, than Shannon-Fano method. This is explained by the fact that in the first case there is a more pronounced violation of the property, characteristic of optimal questionnaires, which is that among all the vertices with the value of paths leading to them equal to that of the followers with the base a must be none whose weight exceeds the sum of weights of other a vertices with the same value of the paths leading to their forerunners.

In conclusion let us point out that the theory of questionnaires enables us to study the characteristics of questionnaires optimal in terms of the minimum or maximum of other objective functions of optimization differing from those mentioned in the paper.

References

1. Picard C. Theorie des questionnaires, Gauthiers-Villars, Paris, 1965.
2. Dubail F. Algorithmes de questionnaires realisables optimaux au sens de differents criteres. These presente a la faculte des sciences de l'Universit  de Lyon, 1967.
3. Пархоменко П.П. Оптимальные вопросники с неравными ценами вопросов. ДАН СССР, 184, № I, 1969 г.
4. Huffman K.E. A method for the construction of minimum-redundancy codes. Proc. IRE, vol. 9, 1952.
5. Беллман Р., Динамическое программирование, ИЛ, Москва, 1960 г.
6. Lawler E.L., Wood D.E., Branch-and-bound methods: a survey, Operat. Res., 14, No. 4, 1966.

EFFECT OF MONITORING PERIODICITY
ON RELIABILITY OF RESTORABLE DEVICES

Garcavi A.L., Gogolevsky V.B.,
Grabovezky V.P.
Reliability Problem Scientific
Council of the Academy of
Sciences of the USSR.
Moscow, USSR

Considered in this report are some methods of determining reliability characteristics of restorable devices under check. The proposed methods are used to solve two groups of problems.

I

The first group of problems is related to the reliability characteristics of devices, operating episodically (by special command), though they are under current for a long period of time. To be determined in this case is an effect of the periodicity and the character of the monitoring tests on the device efficiency. To be found further is the probability of the event when the device will be sound at a random moment t and will show trouble-free operation during time period $\tau - P(t, \tau)$.

It is supposed, that the efficiency testing of these devices is of intermittent type, since a continuous monitoring is either impossible or unsuitable.

The following assumptions are taken: the time intervals between the failures of the device and the time required for restoration are distributed by the exponential law with parameters λ and μ corresponding to the above values respectively. The device is monitored instantly; the monitoring equipment is fully reliable.

Two cases are considered for solving the problem.

a) The monitoring is accomplished at random time intervals $t_{\kappa i}$, distributed by the exponential law with parameter α (Fig. 1a). The device under check, switched on at time instant $t=0$, is in

the failure-free state during a random time interval, till the failure appears. The failure is detected only in a random time interval t_{or} at the instant, when the device is monitored. After detection of the failure the latter is eliminated for a random time t_{re} ; the device is switched on, and the process is repeated all over again.

Thus the device under check can be in three states:

State 1 - the device is sound;

State 2 - the device is faulty and fails to be restored;

State 3 - the device is faulty and restorable.

The random operation process of the device under check and the monitoring device, shown in the fig.1a, can be described by a following system of equations:

$$\begin{aligned} P_1'(t) &= -\lambda P_1(t) + \mu P_3(t) \\ P_2'(t) &= -\alpha P_2(t) + \lambda P_1(t) \\ P_3'(t) &= -\mu P_3(t) + \alpha P_2(t) \end{aligned} \quad (1)$$

where $\alpha = \frac{1}{T_K}$ - is monitoring intensity;
 T_K - is mean time between the monitoring testings.

$P_1(t), P_2(t), P_3(t)$ are probabilities of the fact that the device under check will be in state 1, 2 or 3 at time instant t .

The equation system (1) is solved under conditions that:

$$P_1(0) = 1, P_2(0) = P_3(0) = 0,$$

i.e. it is accepted that in the switching-on instant the device is always sound.

As result of the solution of the equation system (1) the expression for the probability $P_1(t)$, that is for the probability to find the device in the working order at any random time instant t will be represented as:

$$P_1(t) = \frac{\alpha \mu}{\alpha \mu + \alpha \lambda + \mu \lambda} +$$

$$+ \frac{\lambda[\lambda - \alpha - \mu - \sqrt{(\alpha - \mu - \lambda)^2 - 4\mu\lambda}]}{(\alpha - \mu - \lambda)^2 - 4\mu\lambda - (\alpha + \mu + \lambda)\sqrt{(\alpha - \mu - \lambda)^2 - 4\mu\lambda}} \cdot e^{-Ct} +$$

$$+ \frac{\lambda[\lambda - \alpha - \mu + \sqrt{(\alpha - \mu - \lambda)^2 - 4\mu\lambda}]}{(\alpha - \mu - \lambda)^2 - 4\mu\lambda + (\alpha + \mu + \lambda)\sqrt{(\alpha - \mu - \lambda)^2 - 4\mu\lambda}} \cdot e^{-Dt} \quad (2)$$

where

$$C = \frac{\alpha + \mu + \lambda - \sqrt{(\alpha - \mu - \lambda)^2 - 4\mu\lambda}}{2}$$

$$D = \frac{\alpha + \mu + \lambda + \sqrt{(\alpha - \mu - \lambda)^2 - 4\mu\lambda}}{2}$$

The probability to be found will be determined by the expression:

$$\mathcal{P}(t, \tau) = P_i(t) e^{-\lambda \tau} \quad (3)$$

The limiting cases for the probability $P_i(t)$ are:

1a. the permanent value ($t \rightarrow \infty$)

$$P_i(t) \rightarrow K'_{rk} = \frac{1}{1 + \varepsilon + \sigma},$$

where

$$\varepsilon = \frac{\lambda}{\mu}, \quad \sigma = \frac{\lambda}{\alpha}$$

2a. No monitoring ($t_k = \infty$, $\alpha = 0$)

$$P_i(t) = e^{-\lambda t}$$

3a. Continuous monitoring ($t_k = 0$, $\alpha = \infty$, $t = \infty$)

$$K'_{rk} = \frac{1}{1 + \varepsilon}$$

b) The device is monitored periodically in equal time intervals with duration h .

The random operation process of the device under check and the monitoring device for this case is shown in Fig.1b and can be described by the following system of the recurrent algebraic equations, based on the complete probability formula:

$$P(t) = P_n e^{-\lambda \tilde{t}} + \bar{P}_n z(\tilde{t})$$

$$P_n = P_{n-1} e^{-\lambda h} + \bar{P}_{n-1} z(h)$$

$$\bar{P}_n = P_{n-1} (1 - e^{-\lambda h}) + \bar{P}_{n-1} \bar{z}(h)$$

$$\dots$$

$$P_i = e^{-\lambda h}$$

$$\bar{P}_i = 1 - e^{-\lambda h}$$
(4)

The following symbols are used in these equations:

$\rho(t)$ - the probability of the fact that the device is in the working order at random time instant t ; it is accepted, that before the time instant t n monitoring tests have taken place;

\tilde{t} - the time interval elapsed from the instant of the last (n -th) monitoring test up to the present time instant t ;

ρ_i - the probability of the sound state of the device under check at the time instant of the i -th monitoring test;

$\bar{\rho}_i = 1 - \rho_i$ - the probability of the faulty state of the device under check at the time instant $i h$;

$z(\xi)$ - the probability of the fact that the faulty device under check will be restored and will not fail again during the time interval $(0, \xi)$;

$\tilde{z}(\xi) = 1 - z(\xi)$ - the probability of the fact that the faulty equipment will not be restored, or, when restored, will fail again during the time interval $(0, \xi)$.

The equation system is solved as follows:

$$\rho_n(t) = \frac{\mu[e^{-\lambda \tilde{t}}(1 - e^{-\mu h}) - e^{-\mu \tilde{t}}(1 - e^{-\lambda h})]}{\mu} +$$

$$+ \frac{(1 - e^{-\lambda h})(\mu e^{-\mu \tilde{t}} - \lambda e^{-\lambda \tilde{t}})}{\mu(1 - e^{-\mu h}) - \lambda(1 - e^{-\lambda h})} \cdot \left(\frac{\lambda e^{-\lambda h} - \mu e^{-\mu h}}{\mu - \lambda} \right)^n \quad (5)$$

The sign "+" is used at even n , the sign "-" is used at odd n .

The probability to be found will be expressed by the formula:

$$\mathcal{P}(t, \tau) = \rho_n(t) e^{-\lambda \tau} \quad (6)$$

The limiting case. At $n \rightarrow \infty$ and $\tilde{t} = 0$

$$\rho_n \rightarrow K_r' = \frac{\mu(e^{-\lambda h} - e^{-\mu h})}{\mu(1 - e^{-\mu h}) - \lambda(1 - e^{-\lambda h})}$$

The influence of the monitoring nature on the failure-free operation probability of the devices for the cases considered above, is illustrated by the dependences $P_n(t)$ and $P_i(t)$ shown in Fig.2.

II

In the following group of problems are considered the reliability characteristics of the devices, for which, under certain conditions, the operation with intervals is allowed. The intervals can be caused by necessity of performing monitoring tests and the restoration repairs. The principal requirement for such devices is the requirement of the necessity to process the preset volume of information for the preset time. The preset volume of information is divided into separate groups, which are processed in several stages. At the end of each stage the results of work at the given stage are checked by applying one of the known methods.

Then to be found are the following reliability characteristics of the device under check:

1. Probability $P_z(v, t)$ of the processing information volume v for the time $t \geq v$ (the time excess) at z stages;
2. The expectation $M[T]$ of time spent for processing information volume v at z stages;
3. The optimal number z_0 of stages, into which information processing cycle v is to be divided to get the minimum value of $M[T]$.

The determination of these characteristics makes it possible to connect the device reliability with their efficiency.

The problems are solved if we assume that:

- the failure flow^{is} of Poisson's form with parameter λ ;
- the checking time τ_n and the checking and restoration repair time τ_p is taken constant and equal to d and h correspondingly, with the exception of the cases specially mentioned, where these values are taken as random ones with expectations d and h .

Further on are considered three cases of the solution of

the problem.

a) The monitoring is accomplished at the end of each stage by using special fully reliable monitoring devices. In case of failure the results of all stages, preceding to the failure, get lost, because the intermediate storages are not present (Fig. 3a).

Thus, the preset volume of information will be processed, provided the device operates successfully during z stages continuously for the time interval t .

Let us take $\tau_n = \tau_z = d$. The stage number possible for time interval t is determined through the expression:

$$n = \left[\frac{t}{\frac{v}{z} + d} \right],$$

where $[x]$ is a whole part of x .

When we consider the particular stage as a test, and the trouble-free operation of the device during the stage as a event A with probability $p = e^{-\lambda v}$, a following problem arises: to find probability $P_z(n)$ of the event, that among n tests will be at least one series of z successive tests, where event A will occur. Let us consider the series of z of successive successful tests as series \mathcal{Y}_z .

Series \mathcal{Y}_z is among n tests only in two incompatible cases:

- 1) series \mathcal{Y}_z is among the first $(n-1)$ tests;
- 2) series \mathcal{Y}_z appeared for the first time at the n -th test.

Thus,

$$P_z(n) = P_z(n-1) + \pi_z(n), \quad (7)$$

where

$$\left. \begin{aligned} \pi_z(n) &= 0 & (n < z) \\ \pi_z(n) &= p^z & (n = z) \\ \pi_z(n) &= [1 - P_z(n-z-1)] q p^z & q = 1 - p \end{aligned} \right\} \quad (8)$$

Taking into account (8), with recurrent formula (7) successively get the following:

$$P_z(x, t) = P_z(z + x) = p^z (1 + xq), \quad (x = 0, 1, \dots, z)$$

$$P_z(x, t) = P_z(2z + i) = p^z \left\{ 1 + q \left[i + z - ip^z - \frac{i(i-1)}{2} q p^z \right] \right\} \quad (9)$$

etc.

$$(i = 1, 2, \dots, z)$$

The dependence of the information processing probability $P_e(\nu, t)$ on time t at various values z is shown in Fig. 4.

For determination of mathematical expectation of information processing time (τ_n and τ_p are taken here as random values) the following random value is introduced:

$$X_k = x_k - x_{k-1} \quad (k = 1, 2, \dots)$$

where x_k is the instant of the completion of the device restoring after the k -th failure, or the instant of the completion of information processing (in the latter case

$$x_{k+1} = x_{k+2} = \dots = 0)$$

Then the random time of the task execution is:

$$T = \sum_{k=1}^{\infty} X_k$$

and the mathematical expectation of the task execution time is:

$$M[T] = \sum_{k=1}^{\infty} M[X_k]$$

Using the known methods of finding the mathematical expectation, we easily get:

$$M[T] = \frac{(1 - e^{-\lambda\nu})}{(1 - e^{-\frac{\lambda\nu}{2}}) \cdot e^{-\lambda\nu}} \left[e^{-\frac{\lambda\nu}{2}} d + (1 - e^{-\frac{\lambda\nu}{2}}) h + \frac{\nu}{2} \right] \quad (10)$$

The nature of mathematical expectation $M[T]$ change depending on stage number z is shown in Fig. 5. Optimal stage number z_0 is derived from condition of finding a mathematical expectation minimum, which is determined through the formula (10):

$$\frac{\lambda\nu}{\sqrt{2\lambda d}} \leq z_0 \leq \frac{2\lambda\nu}{\sqrt{2\lambda d}} e^{\frac{1}{2}\sqrt{2\lambda d}} \quad (11)$$

Should λd be small, then with a ^{good} accuracy can be accepted, that

$$z_0 \approx \frac{\lambda\nu}{\sqrt{2\lambda d}} \quad (12)$$

b) The results of the information processing at each stage are checked by means of repetition of processing at the stage till two coincident results appear (not necessarily in succession).

The coincident results are considered as correct ones. Special monitoring devices are not used. In case of a failure at the stage the information processing results of only this stage depreciate (availability of information storages is supposed). The organization of each test - data output, comparison and restoration, if necessary - takes time, equal to d (Fig.3b).

As shown in fig.3b, the whole information processing cycle consists of 2 intervals. At each interval the calculations at one stage are made a random number of times (up to the appearance of two coincident results).

Successful processing of information volume v for time interval t with the volume of information divided into 2 stages and with the results to be checked several times, is equivalent to the event that among n checks not more than $n-22$ checks will be unsuccessful, because at each of 2 intervals it is necessary to have two successful checks, i.e. only 22 (or more) successful ~~checks~~ checks. Using the binomial law, we will get a probability

$P_2(v, t)$ of the following form:

$$P_2(v, t) = \sum_{\kappa=22}^n C_n^{\kappa} e^{-\frac{\kappa \lambda v}{2}} (1 - e^{-\frac{\lambda v}{2}})^{n-\kappa} \quad (13)$$

The information processing time expectation can be found with due regard to the fact that the mathematical expectation of test number before the appearance of a successful test with probability p is equal to $\frac{1}{p}$. The mathematical expectation to be found will be as follows:

$$M[T] = 2(v + 2d) e^{\frac{\lambda v}{2}} \quad (14)$$

The optimal number of stages is:

$$z_0 \approx v \sqrt{\frac{\lambda}{d}} \quad (15)$$

b) After each stage a fully reliable test with the aid of a special monitoring device is performed (Fig. 3b). The correct results of information processing at each stage are stored in storage devices. For processing information successfully for time interval t it is necessary that not more than $n-2$ stages from a possible number of the stages n are unsuccessful. Thus the probability to be found is determined in the following way:

$$P_2(v, t) = \sum_{\kappa=2}^n C_n^{\kappa} e^{-\frac{\kappa \lambda v}{2}} (1 - e^{-\frac{\lambda v}{2}})^{n-\kappa} \quad (16)$$

The information processing time expectation in this case is equal to:

$$M[\tau] = (v + 2d) e^{\frac{\lambda v}{2}} \quad (17)$$

The optimal stages number is:

$$z_0 \approx v \sqrt{\frac{\lambda}{d}} \quad (18)$$

CONCLUSIONS

1. In the present report is proposed a system which makes it possible to accomplish a well-founded evaluation of certain reliability parameters of restorable devices accordingly to the given conditions of use and maintenance of the devices.

2. The proposed system makes it possible to rationally organize a time diagram of device operation, which leads to a considerable rise of reliability. This advantage is obtained as a result of a time redundancy, and, secondly, as result of division of the whole information processing volume into stages, which permits rational use of the available time reserve (especially due to the fact that the successful stages results do not get lost and are used again when the device starts operating).

3. The approach described above and the results thus obtained make it possible to select some characteristics of the monitoring system - type of monitoring and its regularity which ensure pre-set reliability standards.

4. The results thus obtained make it possible to associate the devices reliability characteristics with such an important characteristic as efficiency (speed or throughput).

Indeed, it is quite possible to introduce such device characteristic as the real efficiency coefficient determined by a ratio of net time, required for processing information volume V (an absolute reliability of device is supposed) to the mathematical expectation of information processing time with due regard to device unreliability, stage-by-stage operation and time required for monitoring restoration. The real efficiency coefficient must be determined for a preset information processing probability value. Knowledge of these two characteristics permits, for example, settling the question of required operation speed choice when designing digital computers.

Reference

1. E.S.Kochetkov, Restorable systems reliability calculation with regard to the restoration time. "Automation and Telemechanics", N 5, 1965.
2. J.K.Beljaev, Efficiency at two types of failures. "Cybernetics in communism service", v.2, collection of papers, edited by A.I.Berg, N.G.Bruevich and B.V.Gnedenko. Publishing House "Energy", 1964.
3. G.N.Cherkessov, Efficiency of the restorable systems, Theses of reports made on the 3-rd scientific and technical conference on reliability, P.1, published by VSNT0, Leningrad, 1967.

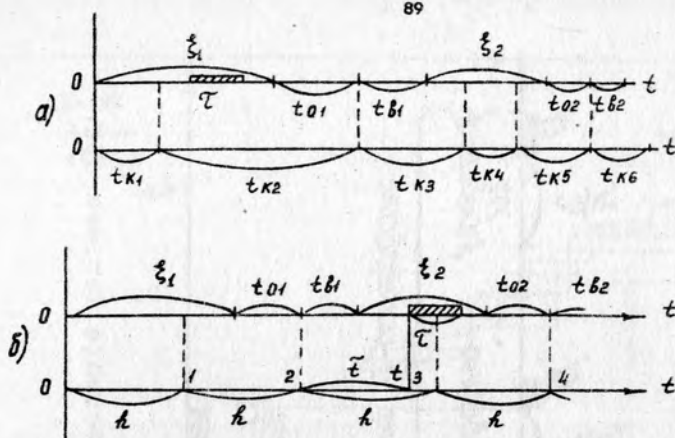


Fig.1. Random operation process of equipment under check and monitoring equipment.

a) random monitoring.

b) periodical monitoring.

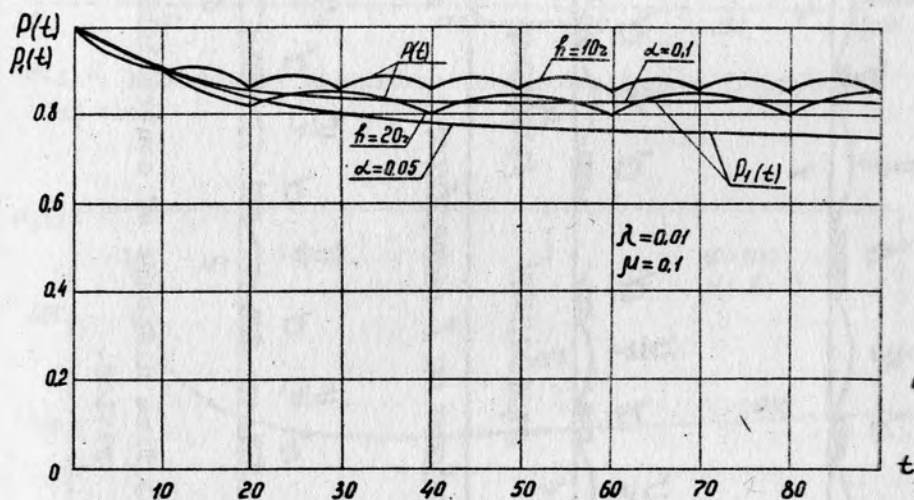


Fig.2. Dependence of probabilities $P(t)$ and $P_n(t)$ on time t .

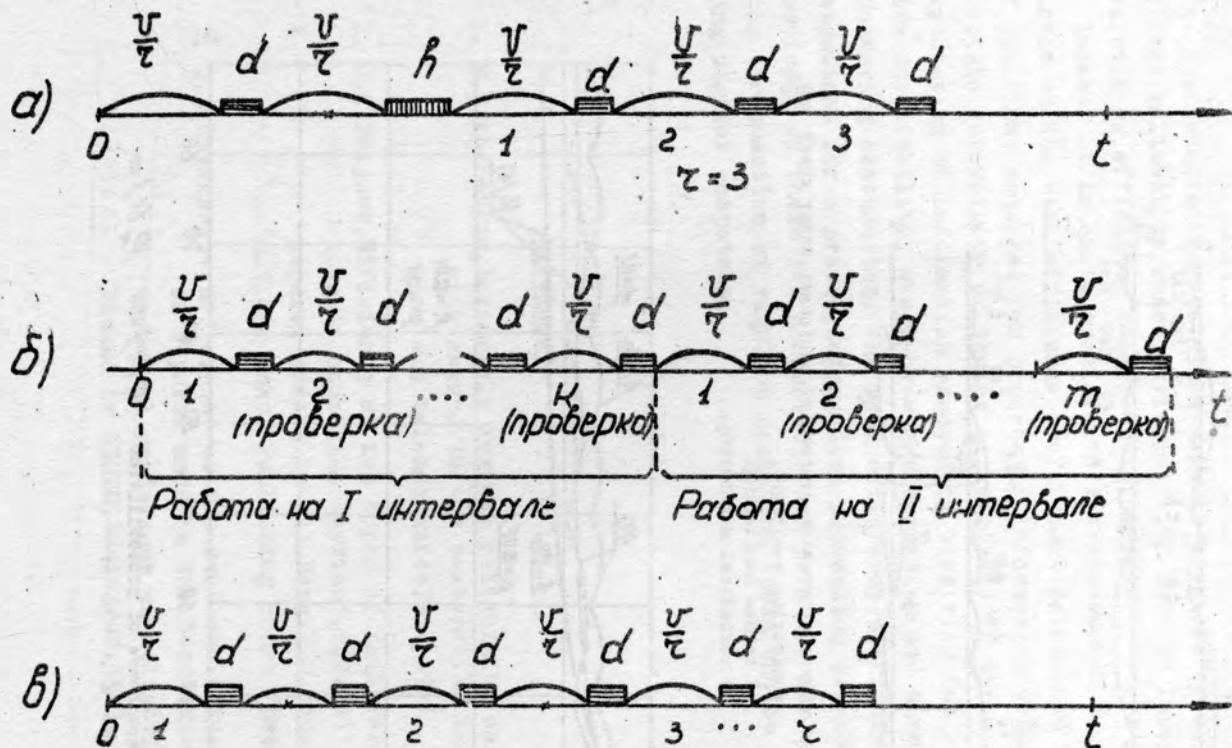


Fig.3. Time diagrams of controllable ~~and~~ restorable devices operation.

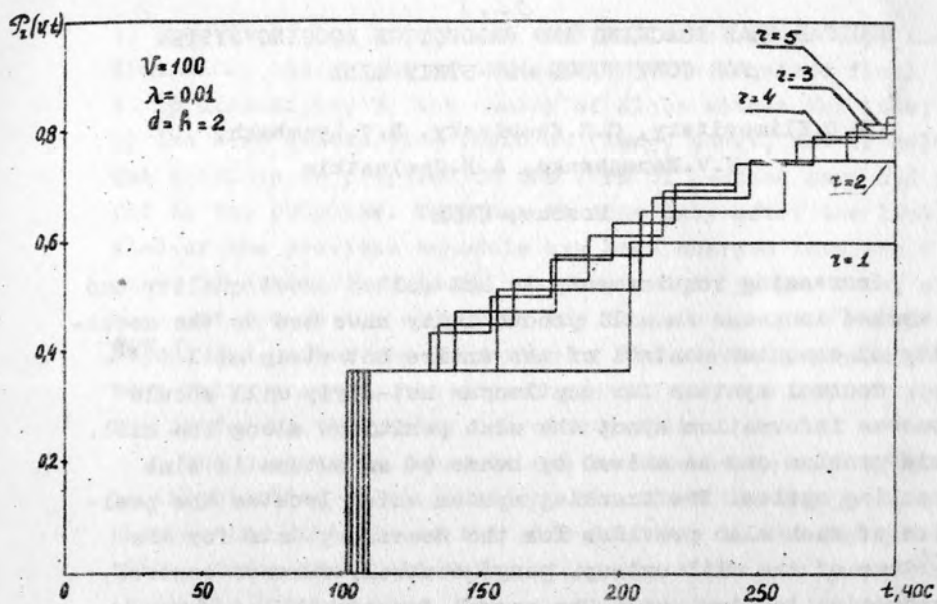


Fig.4. Probability of information processing as function of time.

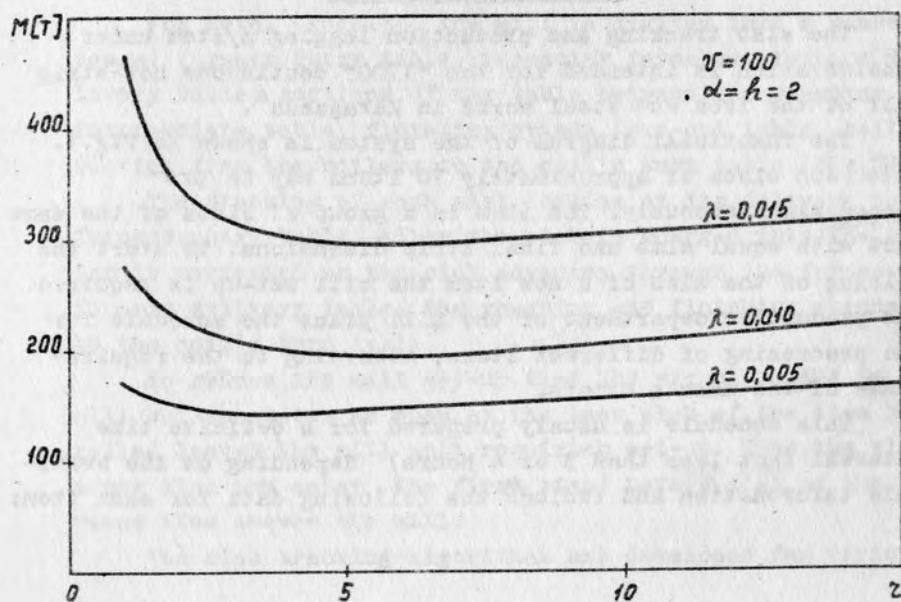


Fig.5. Information processing time expectation dependence on stages number .

DIGITAL SLAB TRACKING AND PRODUCTION LOGGING SYSTEM
FOR CONTINUOUS HOT-STRIP MILLM.D.Klimovitsky, O.S.Kozhinsky, R.V.Lyambakh,
V.V.Naumchenko, A.B.Chelustkin

Moscow, USSR

Increasing requirements to hot-rolled sheet quality and a marked increase in mill productivity have led to the necessity of computer control of the entire hot-strip mill^{1,2,3}.

Control systems for continuous hot-strip mill should possess information about the slab positions along the mill. This problem can be solved by means of an automatic slab tracking system. The tracking system which locates the position of each slab provides for the necessary data for the systems of the mill set-up, gauge control, furnace control, production logging, etc. The normal functioning of these systems requires data of the slab position with space and time.

1. Slab Tracking System

The slab tracking and production logging system under consideration is intended for the "1700" continuous hot-strip mill of the iron and steel works in Karaganda⁴.

The functional diagram of the system is shown in Fig.1. Up to 400 slabs of approximately 70 items may be processed simultaneously. The item is a group of slabs of the same cast with equal slab and final strip dimensions. To start the rolling of the slab of a new item the mill set-up is required. The production department of the mill plans the schedule for the processing of different items, according to the requirements of the user's orders.

This schedule is usually prepared for a definite time interval (not less than 3 or 4 hours) depending on the available information and includes the following data for each item:

1) numbers of the orders and cast; 2) steel grade ; 3) the dimensions and weight of the slab; 4) the required final strip dimensions; 5) the number of slabs within the item; 6) the slab distinctive features (long, short, hot or cold). The schedule is prepared in the form of punched card and then fed to the computer. This can be made only after the last slab of the previous schedule has been charged into the furnace. If a new schedule is not prepared in proper time, and thus not fed to the computer, the furnace operator receives a signal and, after the last item of the previous schedule had been fully charged, the signal prohibiting the further charging is fed to the pusher operator pulpits.

The production schedule can be altered in case of a change in the slabbing mill operation, etc. Then the charge operator changes the schedule of the slab delivery to the furnace entry table. Since the total information about the new cast is ^{not} yet available by this time, the charge operator may insert in the computer only the cast number and the number of slabs in it. According to this information, the computer assigns this cast an intrasystem item number. The rest of the cast data will be fed later by the computer operator.

For metal tracking, the mill is divided into a number of zones: furnace entry table, reheating furnaces, furnace delivery table, sections of the table between the roughing stands, intermediate table, finishing stands, run-out table, coilers, section from the coilers to the coil's turn table (Fig.2).

The tracking of each slab begins at its delivery to the furnace entry table. After the slab is charged into ^{the furnace,} the tracking is performed as the slab advances through the furnace, the furnace delivery table, the roughing and finishing stands up to the coil's turn table.

To reduce the mill set-up time, the signal is fed to the mill set-up system as soon as the last slab of the item being rolled leaves the mill unit requiring set-up. Thus, the slab of a new item may enter the first stand before ^{the last} slab of the previous item leaves the mill.

The slab tracking algorithms are developed for various

process areas. Fig. 3 shows the algorithm flow-chart of the slab tracking for the furnace entry table.

In developing the above algorithm, the following conditions were adopted:

1. The slab being on a furnace entry table is charged either into a furnace P_j or into a furnace row R_q depending on its length (long or short).

2. Only one long slab or two short slabs may be on the entry table of each furnace at a given time moment.

3. A slab can move forward and backward depending on the rotation direction of the table rollers.

4. A slab delivered to the furnace loading table can be charged into this furnace only.

The following designations are accepted for the flow-chart of Fig. 3:

$1MD_j(1MD_j')$ - metal detector signal;

$1TRS$ - signal of the roll-table reverse sensor;

$n_{V_i} HP_j(n_{V_i} HR_q)$ - the slabⁿ of the item V_i is on the entry table of the furnace P_j (or in front of the furnace row R_q);

$n_{V_i} OP_j(n_{V_i} OR_q)$ - the slabⁿ of the ~~xx~~ item V_i is on the loading table of the furnace P_j (or of the row R_q).

To locate the slabs moving from one mill zone to the other and within a zone, various detectors and sensors are provided, such as hot metal detectors, mill-motor load-current relays, etc. For increasing the reliability of the whole system, three sensors of one and the same or different types are installed at each of the appropriate points along the mill.

A slab removed from the process for any reason is deleted from the rolling schedule by the operator of the respective mill area.

The system operation is checked by manual feeding a signal that the last slab of the given item is supplied to the furnace entry table (or discharged to the furnace delivery table). In

case the signals of the operator and computer do not coincide, a corresponding signal to check the system operation is produced.

2. Production Logging System

Production logging is performed by quantitative and qualitative characteristics. All initial information is received from the tracking system and various process sensors and gauges such as thickness and width gauges, sensors of finishing and coiling temperatures, strip length gauge and coil weighbridge.

Some information (date, team number, numbers of the order and cast, steel chemical analysis, etc.) is printed without previous processing. The rest of information is printed after appropriate processing of the initial data. In particular, the coil theoretical weight G_t is determined as the product of steel specific gravity γ , required strip width S_t and thickness h_t , and actual strip length l_a :

$$G_t = \gamma S_t h_t l_a ; \quad (2.1)$$

ϵ is determined as the ratio of the actual coil weight G_a and theoretical weight:

$$\epsilon = \frac{G_a}{G_t} . \quad (2.2)$$

In addition, the total length and the location of the strip sections where the parameter U (thickness, width, finishing or coiling temperature) is outside the preset tolerances are determined. Readings of the parameters U are carried out at each half-meter length of the strip by the signals from the strip length gauge. For each reading the validity of the inequality

$$u_{\min} \leq u_i \leq u_{\max} \quad (2.3)$$

is checked; (U_{\max} and U_{\min} are the preset limits of the parameter U ; i is the reading ordinal number).

The strip length gauge is located at some distance from the process sensor or gauge, therefore the actual location of the strip section where the parameter U is outside the limits is determined by subtracting the number Δi constant for the given sensor from the number of the length gauge signals:

$$\Delta i = \frac{\Delta l}{0.5}.$$

Checking of the validity of the inequality (2.3) at arbitrary i and $i-1$ allows to distinguish the reading numbers where this inequality is violated and to indicate the violation sign. For example, if the parameter exceeds the upper limit at the point l_i (Fig.4), it is necessary to print " $0l_i+$ ". For Fig.4 we obtain:

$$0l_i+ + l_{i+1}0 \quad 0l_{i+2}- - l_{i+3}0.$$

The flow-chart of the strip classification algorithm is shown in Fig.5.

On the basis of this information, the system prints the coil report for each rolled slab.

3. Functional Part of the Tracking and Production Logging System

The slab tracking and production logging system uses the M2000 computers to cope both with the slab tracking and production logging problems and with the furnace and mill units control problems.

The tracking system discrete signals are fed by the discrete information input unit and delivered to the main frame of the computer through the communication unit. The process information from the gauges and sensors of the logging system is transmitted to the input by individual norming converters. To increase reliability, the information flows in parallel through two input units. The computer part of the system uses two processors. One, working in conjunction with the tracking

system with fixed logic, has tracking and logging functions, the other, in addition to these functions, is used to solve the control problems.

The back-up of the input and functional units of the tracking system is dictated by the increased requirements to the system faultless operation. Indeed, the tracking system failure results in loss of information about the rolling schedule and the operators do not know the final dimensions of the next slab. But failure of the processor or any other device performing the control functions requires only the manual input of control-point setting for local control systems, while the tracking and logging system provides for normal functioning of the mill.

Since the amount of data which must accompany each processed slab is very large (overall information for each slab is about 80 decimal digits), the direct transmission of all information through the tracking system (by file shift) is inexpedient because it involves large equipment expenditures. Due to this reason only the slab identity data are transmitted through the tracking system with fixed logic. The rest of information essential for the mill operation is stored in the computer memory. Slab tracking, input of information and its output (on digital display panels) are carried out according to conditional (intrasystem) item numbers.

The system block-diagram is shown in Fig.6. The tracking system operates as follows.

When a new slab item arrives to the furnace entry table, all necessary primary slab data are fed into the system by the automatic data input unit **B1** or by the furnace operator from its pulpit **P1** (when manual input). Simultaneously, this item is given automatically its own sequential intrasystem number. All this information together with this number (which plays the part of the address) enters the logging system computer and is stored in its memory. This item number is assigned to each slab of this item. When the slab comes to the furnace entry table this number is sent to the entry table tracking unit **B2**. The slab having come to the furnace loading

table, this unit transmits the appropriate signal to the receiving register^B of the furnace tracking unit^{B3}. Together with the item number, the tracking system transmits the slab indication comprising three binary digits and containing information about the slab proper (hot, cold, long, short, etc.).

The slab charged into the furnace, its item number and related information are fed into the furnace tracking unit. A new item appearing at the furnace exit, the enquire to the logging system memory is made according to the item number, and necessary information is sent to the digital display panels of the pusher operators (P2, P3), furnace operator (P4), roughing mill operator (P5), shears operator (P6), and finishing mill operator (P7). The number of the slab discharged from the furnace is sent to the display panels of the pusher operators and the furnace operator, and is transmitted to the tracking unit of the furnace delivery table (B4) together with the slab indication. The first slab of a new item discharged, the information on the quantity of the slabs of this item being in the furnaces (or on the quantity of the discharged slabs) is sent to the tracking unit of the furnace delivery table as well.

This unit tracks slabs on the furnace delivery table and counts up the slabs removed from the process and the number of slabs on the delivery table; the latter is sent to the furnace operator display panel. When the first slab of the new item reaches the edger entry table, the tracking unit enquires the logging system memory, and necessary information is transmitted to display panels of the roughing mill operator, shears operator, and finishing mill operator. Number of each slab reached the edger entry table is sent to the roughing mill tracking unit B5.

This unit B5 ensures slab tracking through the roughing mill and output of information about each slab position to the roughing mill operator display panel; moreover, after the last slab of the given item left a roughing stand, the unit sends signal for the set-up of this stand. The tracking unit enquires the logging system memory (according to ^{the} numbers of

the items being rolled in the roughing mill) and sends necessary data to the display panels of the roughing mill operator, shears operator, finishing mill operator, and coiler operator P8. The number of each slab having left the roughing mill is transferred to the intermediate table tracking unit B6.

This unit tracks the slab through the shears and intermediate table and deletes from the tracking system the slabs removed from the intermediate table. When the first slab of a new item reaches the roll table in front of the finishing scale-breaker, this unit enquires the logging system memory for output of necessary information to the display panels of the finishing mill operator and coiler operator.

The finishing mill tracking unit B7 tracks the slab through the finishing mill and enquires the memory for output of information to the display panels of the finishing mill and coiler operators, deletes from the tracking system the slabs removed from the finishing mill in case there were such ones, and assigns the ordinal numbers to the strips rolled during the day. When a strip leaves the finishing mill, the tracking unit B7 sends the item number, the slab number within the item and the strip ordinal number to the coiler tracking unit B8. These numbers are assigned to the data of the given slab as addresses.

The unit B8 tracks the strips according to the item numbers and slab numbers within the item up to the marker. When the first strip of a new item enters the coiler area, the memory is enquired according to the item number, and the data are delivered to the coiler operator display panel. The first coil having come to the turn table, the memory is enquired according to the item number and coil number, and the data are delivered to the marker display panel. The number of the item, of the slab within the item and of the coil are delivered by the coiler tracking unit B8 to the marker tracking unit B9.

The tracking unit B9 enquires the computer memory according to the numbers received and sends the necessary data to the display panel P9 and to the printer. At the same time all information about this coil is deleted from the computer memory. After information of the last coil of the given item is

printed, all data concerning the given item are deleted from the memory.

There are four printers in the system. Two of them (with printing speed of 10 characters per second) are located near the furnaces and print data about slabs being charged and discharged. The other two (printing speed of 400 lines per minute) are installed in the computer room and print coil reports containing the quantitative and qualitative data of the coil.

The table below shows the information about slabs being charged and discharged, the data which are brought out to the operator display panels, and the coil report data.

Table

Pusher operator display panel....	34 decimal	8 binary
Charge operator display panel...	53 decimal	7 binary
Furnace operator display panel...95	"	10 "
Roughing mill operator display panel.....	85 "	33 "
Shears operator display panel...	34 "	4 "
Finishing mill operator display panel.....	63 "	5 "
Coiler operator display panel...	63 "	5 "
Marker display panel	8 "	3 "
Coil report	From 168 to 600 decimal characters depending on the coil quality	
Two printers at the furnace operator pulpit	About 180 decimal characters per slab	

The information required for process control purposes is sent to the operators display panels installed in various areas of the mill.

The numbers of the next three items to be charged and the number of the slabs of the item being charged which are not yet delivered to the furnace entry table are displayed on the

charge operator display panel. When the last slab of the given item reaches the furnace entry table, the "check-up" signal appears on the display panel. If this signal coincides actually with the arrival of the last slab of the item, the charge operator confirms this coincidence by pressing the push-button.

To ensure that all slabs of the given item are discharged from the furnaces and to prevent the slabs of another item to be discharged not in time, the number of slabs of the given item being in each furnace at the present moment is shown on the display panels of the pusher operator and furnace operator. Besides, the "check-up" signal is displayed on the furnace operator panel after the last slab of the given item is discharged.

The data relevant to the item being rolled (required final strip dimensions, the number of slabs not rolled to the present moment, the steel grade, etc.) as well as the data required for the mill set-up before the next item, are displayed on the panels of the roughing and finishing mill operators, shears operator, and coiler operator. Besides, the slab positions on the tables between the roughing stands are displayed on the roughing mill operator panel.

The numbers of the coil and cast, the date, etc., are displayed on the marker panel.

The numbers being transmitted within the tracking system are checked for parity. When the item number is introduced into the tracking system, a check digit is assigned to it.

All data are stored in input/output devices together with the address (item number) and transmitted with it upon request. Comparison of the request address with the output data allows to check the operation of the input and output devices. In case an error in the information being transmitted between any tracking units is detected, the checking unit B10 sends an alarm signal to the display panels of each control pulpit (or to the control pulpits connected with this checking unit).

References

1. Roth J.F., A computer control hierarchy in a steel plant, Iron and Steel, 1964, No. 3,4.
2. Fapiano D.J., Technical and economic considerations in plate mill process control, Iron and Steel Eng., 1966, No.10.
3. Obelode G., Wisdika H., Automatisches Informationssystem zur Auftragsabwicklung in einem Warmbreitbandwalzwerk, Stahl und Eisen, 1967, No.22.
4. Лямбах Р.В., Добропразов Д.Н., Ромашкевич Л.Ф., Фельдман А.В. "Комплексная автоматизация широкополосного стана горячей прокатки с применением УВМ". Доклад на IV Всесоюзном совещании по автоматическому управлению (технической кибернетике), Тбилиси, 1968.

FUNDAMENTALS OF NONLINEAR CONTROL SYSTEMS WITH THE PULSE-FREQUENCY AND PULSE-WIDTH MODULATION

V.M. Kuntsevich, Yu.N. Chekhovoi

Institute of Cybernetics

Ukrainian Academy of Sciences

Kiev

USSR

During the last years the interest in investigating the pulse-frequency and pulse-width modulation (PFM and PWM) systems has grown considerably ¹⁻¹⁶. From the practical point of view this interest is justified by the simplicity of such systems technical realization in comparison with the pulse relay systems and by their much higher dynamical properties ^{1,2,6,9}. From the standpoint of theory the systems in question are interesting by their rather original nonlinear effects due to changing the pulse repetition rate and having no analogies in the theory of continuous and amplitude pulse systems. This fact permits to generalize some classical problems of the automatic control theory (for instance, the absolute stability problems) and enriches them with new content.

1. The Motion Equations

Let us consider the nonlinear sampled data control system (Fig.1) consisting of the continuous linear part (CLP), sequential correcting filter (CF) and nonlinear pulse modulator (PM) of the first-type.

CPL consists of the linear stationary units with lumped parameters and has the fractional rational response function

$$W(s) = \frac{B_\ell(s)}{A_m(s)} = \frac{\sum_{i=0}^{\ell} b_i s^i}{s^m + \sum_{i=0}^{m-1} a_i s^i} \quad (\ell < m), \quad (1.1)$$

CF - is the linear correcting filter defined by the equation

$$\dot{\sigma} = C^x (U - X), \quad (1.2)$$

where $X = (x, x', \dots, x^{(m-1)})$ is the system phase coordinates column-vector; $U = (u, u', \dots, u^{(m-1)})$ - the input signal; $C = (c_1, c_2, \dots, c_m)$ -

- the numerical column-vector; $C_i = 0$ when $i > K$ ($1 \leq K \leq m-l$);
 \mathcal{T}^* - is the operation of transformation.

Structural schemes of PM possible versions for different types of the pulse modulations are given on Fig.2. PE is the ideal sampler (pulse-amplitude modulator); the asterisk denotes the time quantizing operation made by PE: $\sigma^* = \sum_{n=0}^{\infty} \sigma_n \delta(t-t_n)$;

$\delta(t)$ - is a unit delta function; $t_n = \sum_{i=0}^{n-1} T_i$ - time of n-pulse appearance ($t_0 = 0$); $T_n = t_{n+1} - t_n$ - an interval between n- and (n+1)-pulses; $\sigma_n = \lim_{\epsilon \rightarrow 0} \sigma(t_n - \epsilon)$; Φ - a zero-order hold circuit (with constant or variable hold time); RE - relay element; f and F - elements for setting a pulse duration, which control the hold circuit and pulse elements PE. The ideal pulse-frequency modulator (Fig.2, a) modulates the frequency rate and sign of the $x^*(t)$ sequence of unit δ -pulses. Real pulse-frequency modulator (Fig. 2, b) acts on the frequency rate and sign of the $y(t)$ sequence of rectangular pulses having the constant duration τ and unit amplitude. The pulse-width modulator (Fig.2, c) acts on the sign and duration of the $y(t)$ rectangular pulse sequence coming with constant frequency $1/T$. Finally, the pulse-width-frequency modulator (Fig. 2, d) acts on the sign, frequency and duration of $y(t)$ sequence.

For all PM types the modulation by sign is determined by relay function (RE-characteristic)

$$z_n = z(\sigma_n) = \begin{cases} \text{sign } \sigma_n & \text{for } |\sigma_n| > \Delta; \\ 0 & \text{for } |\sigma_n| \leq \Delta; \end{cases} \quad (1.3)$$

the frequency modulation is determined by the FPM law

$$T_n = F(\sigma_n) \quad (1.4)$$

and the pulse-width modulation - by the PWM law

$$0 \leq \tau_n = f(\sigma_n) \begin{cases} < T_n & \text{for } |\sigma_n| < \Delta_0; \\ = T_n = \text{const} & \text{for } |\sigma_n| \geq \Delta_0. \end{cases} \quad (1.5)$$

Here, $F(\sigma)$ and $f(\sigma)$ are the even simple functions defined for all σ ; $F(\sigma) > 0$, $f(\sigma) \geq 0$, $f(\sigma)$ being zero only for $|\sigma| \leq \Delta$; $\Delta_0 > \Delta$ - PM saturation threshold.

Difference equations of the considered systems motion are reduced to the form [10, 11, 17, 18]:

$$X_{n+1} = H_n(X_n + K_n) \quad (1.6)$$

where $X_n = (x_n, x'_n, \dots, x_n^{(m-1)})$; $x_n^{(i)} = \lim_{\epsilon \rightarrow 0} x^{(i)}(t_n - \epsilon)$; $H_n = H[f(\sigma_n)] = \exp A f(\sigma_n)$ - is a transient matrix for the CLP; A - accompanying matrix for a CLP characteristic polynomial; $K_n = K[f(\sigma_n)] z(\sigma_n)$ - is the state changing vector of the CLP. Function $K(f)$ depends on the pulse modulation type. For the ideal PFM (Fig.2, a)

$$K(f) = G = (g, g', \dots, g^{(m-1)}); g^{(i)} = g^{(i)}(0); g(t) = L^{-1}[W(s)]; \quad (1.7)$$

For the real PFM (Fig.2, b)

$$K(f) = H(-\tau)R(\tau); R(\tau) = (z(\tau), z'(\tau), \dots, z^{(m-1)}(\tau)); z(t) = L^{-1}[\frac{1}{s}W(s)] \quad (1.8)$$

and, finally, for the PWM (Fig.2, c) and double frequency and width modulation (PFM and PWM, Fig.2, d)

$$K(f) = H(-\tau_n)R(\tau_n) = H[-f(\sigma_n)]R[f(\sigma_n)]. \quad (1.9)$$

The matrix equation (1.6) describes the system motion (Fig.1) in natural phase space $E^m = \{X_n\}$. Let us show that it is always possible to pass from (1.6) to an equation

$$\dot{X}_{n+1} = \dot{H}_n (\dot{X}_n + \dot{K}_n) \quad (1.10)$$

in difference phase space $D^m = \{\dot{X}_n\}$, $\dot{X}_n = (x_n, x_{n+1}, \dots, x_{n+m-1})$. Let us compose the following equation system [10, 11]:

$$\begin{vmatrix} -H_n & I & 0 & \dots & 0 & 0 \\ 0 & -H_{n+1} & I & \dots & 0 & 0 \\ 0 & 0 & -H_{n+2} & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & -H_{n+m-1} & I \end{vmatrix} \cdot \begin{vmatrix} X_n \\ X_{n+1} \\ X_{n+2} \\ \dots \\ X_{n+m-1} \\ X_{n+m} \end{vmatrix} = \begin{vmatrix} H_n K_n \\ H_{n+1} K_{n+1} \\ H_{n+2} K_{n+2} \\ \dots \\ H_{n+m-1} K_{n+m-1} \end{vmatrix} \quad (1.11)$$

where I - is a unit matrix. The system (1.11) has the rectangular matrix with $m^2 \times m(m+1)$ dimensions and of m^2 rank. Therefore (1.11) can be solved relatively to X_{n+m} assuming the variables x_{n+i} ($i=0, 1, \dots, m-1$) to be known. This result represented in matrix form provides the equation (1.10).

In future we shall use the motion equations in the form of (1.6) only, because its relation with system parameters is more simple and it is more convenient for investigations.

2. The Stability of the Equilibrium Points

Let us use discrete analogies of theorems of the Liapunov's direct method [11, 19, 20]. Consider the main case when CLP of the system is stable and the simplest critical case when CLP is neutral. In order the system (Fig.1) might have an equilibrium

point assume $u = \text{const.}$

The main case. According to (1.6) X_∞ column-vector of the equilibrium point coordinates must satisfy the equality

$$X_\infty = (H_\infty^{-1} - I)^{-1} K_\infty, \quad (2.1)$$

where H_∞ and K_∞ - matrix H_n and vector K_n , respectively, for $X_n = X_\infty$.

For $|u| \leq \frac{\Delta}{C_i}$ ($C_i > 0$ - the vector C element) equation (2.1) has zero solution $X_\infty = 0$, i.e. the system has the equilibrium point in the origin. This fact is easily checked by the simple substitution. For $|u| > \frac{\Delta}{C_i}$ the solution $X_\infty \neq 0$; in general case there may be several solutions²¹.

By substitution $X_n = E_n^\circ + X_\infty$ we shall display the phase origin E_n° in the system equilibrium point. In new coordinates instead of (1.6) we obtain

$$E_{n+1}^\circ = H_n (E_n^\circ + K_n^\circ), \quad (2.2)$$

where

$$K_n^\circ = (I - H_n^{-1}) X_\infty + K_n; \quad \theta_n = c^T (E_\infty - E_n^\circ); \quad E_\infty = V - X_\infty \quad (2.3)$$

The Liapunov's function we shall choose as a positively determined quadratic form

$$V_n = (E_n^\circ)^T P E_n^\circ; \quad P > 0. \quad (2.4)$$

The first difference of the function (2.4) owing to (2.2) is equal to:

$$\Delta V_n = -(E_n^\circ)^T (P - M_n) E_n^\circ + 2(E_n^\circ)^T M_n K_n^\circ + (K_n^\circ)^T M_n K_n^\circ, \quad M_n = H_n^T P H_n. \quad (2.5)$$

The system (2.2) is asymptotically stable on the whole if for any $E_n^\circ \neq 0$ function (2.5) is negative, i.e. if the inequalities are satisfied^[11, 19, 20]:

$$P - M_n > 0; \quad (2.6)$$

$$(E_n^\circ)^T (P - M_n) E_n^\circ - 2(E_n^\circ)^T M_n K_n^\circ > (K_n^\circ)^T M_n K_n^\circ. \quad (2.7)$$

It is shown that $\lim_{F \rightarrow \infty} M[F(\theta_n)] = a$, thus there is a class of sufficiently large functions (1.4), for which the condition (2.6) is satisfied^[22]. Assume that the function (1.4) belongs to this still unknown class and consider the condition (2.7).

Let us examine the equation of a surface on which the function (2.5) becomes zero:

$$(E_n^\circ)^T (P - M_n) E_n^\circ - 2(E_n^\circ)^T M_n K_n^\circ = (K_n^\circ)^T M_n K_n^\circ. \quad (2.8)$$

Condition (2.7) will be satisfied if for all $F_n^0 \neq 0$ this surface does not exist. We substitute in (2.8) σ_n defined from (2.3), by expression $\sigma_n - \sigma$ where $\sigma_n = C^T E_n^0$ and σ - is the arbitrary real parameter not depending on E_n^0 :

$$(E_n^0)^T (P-M) E_n^0 - 2(E_n^0)^T M K^0 = (K^0)^T M K^0; M = M_n \Big|_{\sigma_n = \sigma_n - \sigma}; K^0 = K_n^0 \Big|_{\sigma_n = \sigma_n - \sigma} \quad (2.9)$$

Equation (2.9) describes a family of n -dimensional ellipsoids depending on σ parameters, the surface (2.8) not existing for $E_n^0 \neq 0$, if for all $\sigma \neq 0$ the ellipsoid (2.9) is not contiguous to $C^T E_n^0 = \sigma$ plane (Fig.3, a) ^{11, 20}. Let us construct the plane

$$C^T E_n^0 = \rho(\sigma), \quad (2.10)$$

tangent to ellipsoid (2.9). Condition (2.7) is transformed into an inequality

$$|\sigma| > |\rho(\sigma)|, \quad \sigma \neq 0, \quad |\sigma| \leq \Delta. \quad (2.11)$$

To determine $\rho(\sigma)$ function we write down the general equation of a plane tangent to ellipsoid (2.9) in some point ²¹⁻²³

$$(E_n^0)^T (P-M) E^0(A) - [E_n^0 + E^0(A)]^T M K^0 = (K^0)^T M K^0, \quad (2.12)$$

where $E^0(A)$ - radius-vector of the point of tangency (Fig.3, a). Planes (2.10) and (2.12) coincide if for some $\alpha > 0$

$$(P-M) E^0(A) - M K^0 = \alpha C; \quad (2.13)$$

$$(K^0)^T M K^0 + [E^0(A)]^T M K^0 = \alpha \rho(\sigma). \quad (2.14)$$

The combined solution of (2.9), (2.13) and (2.14) gives ^{21, 22}

$$\rho(\sigma) = \left\{ (K^0)^T [M + M(P-M)^{-1} M] K^0 C^T (P-M)^{-1} C \right\}^{\frac{1}{2}} \text{sign } \sigma + C^T (P-M) M K^0 \quad (2.15)$$

Note that when $|u| \leq \frac{\Delta}{C}$ function (2.15) becomes zero for $|\sigma| \leq \Delta$ and it is quite sufficiently to check the inequality (2.11) for $|\sigma| > \Delta$.

Condition (2.11) and equation (2.15) are obtained on the base of assumption that there exists (2.6) for all $\sigma \neq 0$, but it is unnecessary to check the satisfaction of this condition along all σ axis. Really let us assume that inequality (2.6) is satisfied for $\sigma = \sigma_1$ and violated for $\sigma = \sigma_2 \neq \sigma_1$; then there is such $\sigma_3 \in (\sigma_1, \sigma_2)$ for which matrix $(P-M)$ is degenerated and function (2.15) has the gap of the continuity. It means that it is sufficiently to check (2.6) only in one arbitrary point of every interval of the function continuity. Now the final re-

sult may be formulated as follows: in the main case the equilibrium point X_∞ of system (1.6) is asymptotically stable on the whole if condition (2.11) is satisfied and inside any continuity interval of function (2.15) it is possible to point out at least one value of σ for which the inequality (2.6) is satisfied.

The simplest critical case. In this case one root of the characteristic equation $A_m(s)=0$ is zero ($a_0=0$), matrix A is degenerated and has $m-1$ rank. Transform (1.6) multiplying it by the matrix^[22]

$$R = \begin{vmatrix} 1 & \frac{a_1}{a_1} & \dots & \frac{a_{m-1}}{a_1} & \frac{1}{a_1} \\ 0 & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & \dots & 0 & 1 \end{vmatrix}; \quad (2.16)$$

after simple transformations we obtain:

$$\tilde{E}_{n+1}^\circ = \tilde{H}_n (\tilde{E}_n^\circ + \tilde{K}_n), \quad \sigma_n = \tilde{C}^T (U - \tilde{X}_n) = -\tilde{C}^T \tilde{E}_n^\circ \quad (2.17)$$

where $\tilde{E}_n^\circ = \tilde{X}_n - U$; $\tilde{X}_n = R X_n$; $\tilde{K}_n = R K_n$; $\tilde{C} = (R^T)^{-1} C$; $\tilde{H}_n = R H_n R^{-1} = \text{diag} \{1, (H_n)_{i,i}\}$; $(H_n)_{i,i}$ - the matrix of $m-1$ order, obtained from H_n by crossing out the first row and first column. To equation (2.17) a set \mathfrak{E} of equilibrium points corresponds over which

$$(\tilde{X}_n)_1 = (x_n)_1 = 0; \quad |u - \tilde{x}_{i,n}| = |u - x_{i,n}| \leq \frac{\Delta}{\tilde{c}_i} = \frac{\Delta}{c_i}. \quad (2.18)$$

Here $(X)_1$ is a column-vector obtained from X by crossing the first element out; $\tilde{x}_{i,n}$ and \tilde{c}_i are the elements of vectors \tilde{X}_n and \tilde{C} .

We choose the Liapunov's function in a form

$$V_n = (\tilde{E}_n^\circ)^T P \tilde{E}_n^\circ, \quad (2.19)$$

in addition demanding from matrix P that $P = \text{diag} \{1, (P)_{i,i}\}$; $(P)_{i,i} > 0$. Then

$$\Delta V_n = -(\tilde{E}_n^\circ)_1^T (P - \tilde{M}_n)_{i,i} (\tilde{E}_n^\circ)_1 + 2(\tilde{E}_n^\circ)^T \tilde{M}_n K_n + \tilde{K}_n^T \tilde{M}_n \tilde{K}_n; M_n = \tilde{H}_n^T P \tilde{H}_n \quad (2.20)$$

In accordance with discrete analog of La-Sall theorem the set \mathfrak{E} of equilibrium points of system (2.17) is asymptotically stable on the whole^[11, 20] if

$$\Delta V_n = 0, \quad \tilde{E}_n^\circ \in \mathfrak{E} \quad \Delta V_n < 0, \quad \tilde{E}_n^\circ \notin \mathfrak{E} \quad (2.21)$$

From (2.20), (2.17) and definition of \mathcal{E} it follows that the first of two conditions (2.21) is always satisfied. The check of the second condition (2.21) is analog to the previous one and leads us to the geometrical problem of finding conditions under which m -dimensional paraboloid

$$(\tilde{E}_n^0)^T (P - \tilde{M})_{ii} (\tilde{E}_n^0)_i - 2(\tilde{E}_n^0)^T \tilde{M} \tilde{K} = \tilde{K}^T \tilde{M} \tilde{K}; \tilde{M} = \tilde{M}_n \Big|_{\tilde{\sigma}_n = \tilde{\sigma}}, \tilde{K} = \tilde{K}_n \Big|_{\tilde{\sigma}_n = \tilde{\sigma}} \quad (2.22)$$

is not tangent to the plane $\tilde{C}^T \tilde{E}_n^0 = \tilde{\sigma}$ (Fig. 3, b). As before, the problem solution leads to inequality (2.11). Omitting the intermediate calculations (analog to those realized in the main case) let us write the expression for function $\rho(\sigma)$ ²²:

$$\rho(\sigma) = -\frac{C_i}{2\tilde{K}_i} \left\{ \tilde{K}_i^2 + (\tilde{K})_i^T \left[(\tilde{M})_{ii} + (\tilde{M})_{ii} (P - \tilde{M})_{ii}^{-1} (\tilde{M})_{ii} \right] (\tilde{K})_i + \right. \\ \left. + \frac{\tilde{K}_i^2}{\tilde{C}_i^2} (\tilde{C})_i^T (P - \tilde{M})_{ii}^{-1} (\tilde{C})_i \right\} + (\tilde{C})_i^T (P - \tilde{M})_{ii}^{-1} (\tilde{M})_{ii} (\tilde{K})_i,$$

where \tilde{K}_i - element of \tilde{K} vector.

Function (2.23) becomes zero when $|\sigma| \leq \Delta$, so in the simplest critical case (irrespective of u value) it is sufficient to check the inequality (2.11) only for $|\sigma| > \Delta$.

In the simplest critical case the additional condition (2.6) takes on the form:

$$(P - \tilde{M})_{ii} > 0. \quad (2.24)$$

Taking into account the last observations the final result is formulated as follows: in the simplest critical case the set \mathcal{E} of system (1.6) equilibrium points is asymptotically stable on the whole, provided that for $|\sigma| > \Delta$ condition (2.11) is satisfied and that it is possible to indicate at least one $\tilde{\sigma}$ value in every interval of function (2.23) continuity for which inequality (2.24) is satisfied.

CLP Critical Transfer Coefficient. Let us represent (1.1) as $W(s) = k W_0(s)$, where constant k (we call it CLP transfer coefficient) is $\lim_{s \rightarrow 0} W(s)$ mainly, and $\lim_{s \rightarrow 0} s W(s)$ in the simplest critical cases. Then inequality (2.11) can be reduced to the form

$$\frac{|\sigma|}{k} > |\rho_0(\sigma)|, \quad (2.26)$$

where $\rho_0(\sigma)$ is determined by (2.15) or (2.23) if instead of (1.1) we use $W_0(s)$. From (2.26) it follows that the critical transfer coefficient (or the smallest k for which the stability condition is not satisfied) corresponds to the largest k for which the function $\rho_0(\sigma)$ is in sector $[0, \frac{1}{k}]$ (Fig. 4).

The Absolute Stability of the Equilibrium Points. The system consisting of a nonlinear element (NE) and a linear part (LP) is called as absolute stable if it is asymptotically stable on the whole for all NE characteristics of some class^{24, 25}. In the system under investigation (Fig. 1) the nonlinear element is the PM whose properties are completely defined by (1.3), (1.4) and (1.5). Thus the absolute stability problem can be considered here in three different statements: 1) (1.4) and (1.5) are known, a class of feasible functions (1.3) is searched for; 2) (1.3) and (1.5) are known, a class of suitable PFM laws (1.4) is searched for; 3) (1.3) and (1.4) are known, a suitable class of PWM laws (1.5) is searched for²⁶.

To simplify the problem let us assume $u=0$ in the main case. Then the stability condition on the whole will be

$$|\sigma| > |\rho(\sigma)|, \quad \Delta < |\sigma| \leq \Delta_0. \quad (2.27)$$

The first statement of the absolute stability problem is trivial and is solved completely by condition (2.27). Really if for some $\Delta = \Delta_*$ condition (2.27) is satisfied it will be satisfied for any $\Delta \geq \Delta_*$. The last inequality solves the problem completely as it defines the class of functions (1.3) searched for.

For solving the problem in its 2nd and 3rd statements let us represent the function $\rho(\sigma)$ for $\sigma > 0$ in the following form:

$$\rho(\sigma) = \rho_\Delta[F(\sigma), f(\sigma)] = \rho_\Delta(F, f), \quad \sigma > 0 \quad (2.28)$$

For $F \rightarrow \infty$ and $f \rightarrow 0$ the system under investigation will be broken out. In accordance with the condition the CLP is stable (or limiting stable²⁴) so the feasible class of functions $F(\sigma)$ is bounded below and $f(\sigma)$ - is bounded above. Let us consider the equation

$$\sigma = \rho_\Delta(F, f), \quad \sigma > 0, \quad (2.29)$$

corresponding to the boundary of the stability domain given by (2.27). If function (1.5) is set the equation (2.29) gives inexplicitly the function

$$\sigma = \rho_{\Delta, f}(F), \quad \sigma > 0. \quad (2.30)$$

Function (2.30) is defined and positive over the interval

$F \in (F_0, \infty)$, where $F_0 \geq 0$ - is the highest value of F for which conditions (2.6) or (2.24) are violated. Therefore there is the positive inverse function $F_*(\sigma) = \rho_{\sigma, f}^{-1}(\sigma)$ which use makes it possible to reduce (2.27) condition to the form:

$$F(\sigma) > F_*(\sigma) = \rho_{\sigma, f}^{-1}(\sigma); \quad \sigma \in (\Delta, \Delta_0]. \quad (2.31)$$

The obtained inequality defines the class of functions (1.4) searched for. Similarly the class of feasible functions (1.5) is defined by the inequality

$$f(\sigma) < f_*(\sigma) = \rho_{\sigma, f}^{-1}(\sigma), \quad \sigma \in (\Delta, \Delta_0). \quad (2.32)$$

In general case one fails to find inverse functions $\rho_{\sigma, f}^{-1}(\sigma)$ and $\rho_{\sigma, F}^{-1}(\sigma)$ in analytic form but it is simple to find them graphically²⁶.

The Asymptotic Stability Region. If condition (2.27) is violated for $\sigma = \Delta, \in [\Delta, \Delta_0]$, the system under investigation does not possess the asymptotic stability on the whole but it has the region (in E^m space) of the asymptotic stability. The largest of the open regions is the estimation of this region, bounded by the surfaces $\psi_n = \beta = \text{const}$ and satisfying inequality $|\sigma| < \Delta$. It is easy to show that in the main and the simplest critical cases, respectively, β value is equal to²⁶:

$$\beta = \frac{\Delta_i^2}{C^T P^{-1} C}; \quad \beta = \frac{\Delta_i^2}{\tilde{C}^T P^{-1} \tilde{C}}. \quad (2.33)$$

Remark about Choosing the Liapunov Function. From the stated above it follows that inequalities (2.31) and (2.32) ensure the stability region for any $P > 0$ and $(P)_{ii} > 0$, respectively. Satisfactory results and material simplifications of calculations can be obtained when choosing

$$P = S^* S, \quad P = S^* D S, \quad (P)_{ii} = S^* S, \quad (P)_{ii} = S^* D S, \quad (2.34)$$

where S - matrix transforming the matrix A or $(A)_{ii}$, respectively into the normal Jordan's form; S^* - matrix, Hermitianly conjugate to S ; D - the positive elements diagonal matrix^{11, 17, 18, 20-22, 26}.

3. The Forced Stationar Regime Stability

Let us consider the (1.6) system in the simplest critical case for $u(t) = \omega t$ (the regime of following the linearly increasing input signal). Using the substitution $X_n = U_n - E_n$ we reduce (1.6) to the following form:

$$E_{n+1} = H_n(E_n + L_n), \quad L_n = H_n^{-1} U_{n+1} - U_n - K_n. \quad (3.1)$$

It may be shown that for $u(t) = \omega t$ this equation does not depend explicitly on t_n and there exists for it a limiting equality $\lim_{n \rightarrow \infty} E_n = E_\infty$, where E_∞ is a numerical vector²⁷. According to (3.1) vector E_∞ must satisfy the equality

$$(H_\infty^{-1} - I) \cdot E_\infty = L_\infty, \quad L_\infty = \lim_{n \rightarrow \infty} L_n. \quad (3.2)$$

It may be shown that there is always interval $\omega \in (\omega_1, \omega_2)$ ($\omega_2 > \omega_1 > 0$) for which equation (3.2) has at least one solution²⁸. By substituting $E_n = E_\infty - E_n^\circ$ we replace the origin of phase space coordinates to the point corresponding to the investigated stationary regime. In new coordinates instead of (3.1) we obtain

$$E_{n+1}^\circ = H_n(E_n^\circ + L_n^\circ), \quad L_n^\circ = -(I - H_n^{-1})E_\infty - L_n, \quad \sigma_n = C^T(E_\infty - E_n^\circ). \quad (3.3)$$

Similarly to the previous let us transform equation (3.3) by multiplying it by the matrix (2.16)

$$\tilde{E}_{n+1}^\circ = \tilde{H}_n(\tilde{E}_n^\circ + \tilde{L}_n^\circ), \quad \sigma_n = \tilde{C}^T(\tilde{E}_\infty - E_n^\circ) = \tilde{\sigma}_\infty - \tilde{C}^T \tilde{E}_n^\circ, \quad (3.4)$$

where $\tilde{L}_n^\circ = R L_n^\circ$, $\tilde{E}_\infty = R E_\infty$ (the rest of indices correspond to those taken above). Equation (3.4) is similar to (2.17) but instead of \mathcal{E} set of equilibrium points the system (3.4) has one equilibrium point $\tilde{E}_n^\circ = 0$. When Liapunov function is chosen as (2.19) the stability condition (similarly to the previous one) will be obtained in the form of (2.11), where

$$\rho(\sigma) = -\frac{C_1}{2\tilde{C}_1} \left\{ (\tilde{E}_1^\circ)^2 - (\tilde{L}_1^\circ)^T [(\tilde{M})_{11} + (\tilde{M})_{11}^T (P - \tilde{M})_{11}^{-1} (\tilde{M})_{11}] (\tilde{L}_1^\circ) + \right. \\ \left. + \frac{(\tilde{E}_1^\circ)^2}{C_1} (\tilde{C})_1^T (P - \tilde{M})_{11}^{-1} (\tilde{C})_1 \right\} + (\tilde{C})_1^T (P - \tilde{M})_{11}^{-1} (\tilde{M})_{11} (\tilde{L}_1^\circ). \quad (3.5)$$

Here $\tilde{L}_1^\circ = \tilde{L}_1^\circ|_{\sigma_n = \tilde{\sigma}_\infty - \sigma}$; $\tilde{M} = \tilde{M}_n|_{\sigma_n = \tilde{\sigma}_\infty - \sigma}$; \tilde{E}_1° is vector element.

4. Limited Boundedness (Dissipativeness)

In those cases when the condition of asymptotical stability is not satisfied on the whole in (2.2) system existence of periodical regimes is possible. The exact and approximate methods for analysis of such regimes have small practical efficiency being cumbersome and requiring the apriori knowledge about the form of the periodical process (number of pulses on

on a half-period, order of their alternation, etc.) or searching all possible variants^{10, 11, 28}. The method of investigating a limited boundedness (dissipativeness) of automatic systems gives the best results.

System (2.2) is called a limited bounded (dissipative) on the whole if there is such a compact set \mathcal{G} (asymptotically stable set) that for any initial conditions $X_n \in \mathcal{G}$ when $n \rightarrow \infty$. On the ground of discrete analog of the T. Yoshisava theorem system (2.2) is limited bounded, if the set \mathcal{B} over which function (2.5) is positive is limited. The estimation of the asymptotically stable set \mathcal{G} is closed domain limited by the surface $V_n = \mu$ ($0 < \mu = \text{const}$) which describes \mathcal{B} ^{11, 29, 30}. Let us express matrix P of the quadratic form (2.4) as $P = Q^* Q$ and transform equation (2.2) by multiplying it by Q from the left:

$$Y_{n+1} = Q H_n Q^{-1} (Y_n + Q K_n^0), \quad Y_n = Q E_n^0. \quad (4.1)$$

In the space $E_q^m = \{Y_n\}$ the surface $V_n = \mu$ is a sphere and, consequently, the sphere circumscribed the set \mathcal{G} is a boundary of asymptotically stable set \mathcal{B} ^{11, 29, 30}. For the second order systems this conclusion leads to the simple graphical procedure (Fig. 5) that sufficiently simplifies the investigation. The examples of investigating the stability and limited boundedness of the concrete systems can be found in the works^{11, 17, 18, 20-22, 26-30}.

REFERENCES

1. Цыпкин Я.З. Дискретные автоматические системы, их теория и перспективы развития. В сб. "Теория и применение дискретных автоматических систем". Изд-во АН СССР, 1960.
2. Кунцевич В.М. Исследование переходных и установившихся процессов в релейно-импульсных и экстремальных системах с постоянным и переменным периодом регулирования. Изв. АН СССР, ОТН, Энергетика и автоматика, № 5, 1961.
3. Dorf R.C., Farren M.C., Phillips C.A. Adaptive Sampling Frequency for Sampled-Data Control Systems. IRE Transactions, V. AC-7, No. 1, 1962.
4. Murphy G., West R.L. The Use of Pulse-Frequency Modulation for Adaptive Control. Proceedings of the National

Electronics Conference, V. XVIII, Pt I, Chicago, 1962.

5. Ли С.С., Джонс Р.В. Интегральные системы управления с частотно-импульсной модуляцией. Докл. на 2-м межд. конгр. ИФАК, М., 1963.

6. Кунцевич В.М. Исследование импульсных самонастраивающихся систем и способов повышения их качества. Дисс. на соискание уч. ст. доктора техн. наук, Ин-т кибернетики АН УССР, Киев, 1964.

7. Murphy G.J., Wu S.H. A Stability Criterion for Pulse-Width Modulated Feedback Control Systems. IEEE Transactions, V. AC-9, No. 4, 1964.

8. Pavlidis T., Jury E.I. Analysis of a New Class of Pulse-Frequency Modulated Feedback Systems. IEEE Transactions, V. AC-10, No. 1, 1965.

9. Кунцевич В.М. Импульсные экстремальные и самонастраивающиеся системы. Изд-во "Техника", Киев, 1966.

10. Чеховой Ю.Н. Динамика импульсных систем с частотно-импульсной модуляцией. Точные методы анализа. Ч. I и II. В сб. "Теория автоматического управления", вып. I, изд. "Наукова думка", Киев, 1966.

11. Чеховой Ю.Н. Динамика нелинейных импульсных систем управления с частотно-импульсной модуляцией. Дисс. на соискание уч. ст. канд. техн. наук, Ин-т кибернетики АН УССР, Киев, 1966.

12. Pavlidis T. Stability of a Class of Discontinuous Dynamical Systems. Information and Control, V. 9, No. 3, 1966.

13. Clark J.P.C., Noges E. The Stability of Pulse-Frequency Modulated Closed Loop Control Systems. IEEE International Conventional Record, V. 14, No. 6, 1966.

14. Чеховой Ю.Н. Нелинейный частотно-импульсный регулятор. Энергетика и электрификация, № I, Киев, 1967.

15. Гелиг А.Х. Стабилизация нелинейных систем с частотно-импульсной модуляцией. Автоматика и телемеханика, № 6, 1967.

16. Дымков В.И. Об абсолютной устойчивости частотно-импульсных систем. Автоматика и телемеханика, № 10, 1967.

17. Кунцевич В.М., Чеховой Ю.Н. Устойчивость систем с двойной (частотной и широтной) импульсной модуляцией. Автоматика и телемеханика, № 7, 1967.

18. Кунцевич В.М., Чеховой Ю.Н. Исследование устойчивости систем управления с двойной (частотной и широтной) импульсной модуляцией прямым методом Ляпунова. В сб. "Сложные системы управления", вып. 3, изд-во "Наукова думка", Киев, 1967.

19. Бромберг П.В. Матричные методы в теории релейного и

импульсного регулирования. Изд-во "Наука", М., 1967.

20. Кунцевич В.М., Чеховой Ю.Н. Исследование устойчивости импульсных систем управления с частотно-импульсной модуляцией прямым методом Ляпунова. Автоматика и телемеханика, № 2, 1967.

21. Чеховой Ю.Н. Устойчивость систем управления с частотно-импульсной модуляцией при постоянных внешних воздействиях. В сб. "Дискретные системы управления", вып. I. Изд-во "Наукова думка", Киев, 1969, в печати.

22. Чеховой Ю.Н. К задаче об устойчивости в целом импульсных систем управления с частотно-импульсной модуляцией. Автоматика и телемеханика, № 6, 1968.

23. Розенфельд Б.А. Многомерные пространства. Изд-во "Наука", М., 1966.

24. Айзерман М.А., Гантмахер Ф.Р. Абсолютная устойчивость регулируемых систем. Изд-во АН СССР, М., 1963.

25. Гантмахер Ф.Р., Якубович В.А. Абсолютная устойчивость нелинейных регулируемых систем. В трудах 2-го всесоюзного съезда по теор. и прикл. механике, вып. I, изд-во "Наука", М., 1965.

26. Чеховой Ю.Н. Абсолютная устойчивость импульсных систем управления с частотной и широтной импульсной модуляцией. В сб. "Теория автоматического управления", вып. I, Киев, 1968.

27. Кунцевич В.М. Исследование устойчивости стационарного вынужденного режима в системах управления с частотно-импульсной модуляцией. В сб. "Дискретные системы управления", вып. I, изд-во "Наукова думка", Киев, 1969, в печати.

28. Чеховой Ю.Н. Динамика импульсных систем с частотно-импульсной модуляцией. Ч. II, В сб. "Теория автоматического управления", вып. 3, Киев, 1967.

29. Чеховой Ю.Н. Об одной форме устойчивости нелинейных импульсных систем. Изв. АН СССР. Техническая кибернетика, № 4, 1967.

30. Чеховой Ю.Н. Предельная ограниченность нелинейных импульсных систем. В сб. "Теория автоматического управления", вып. 2, Киев, 1967.

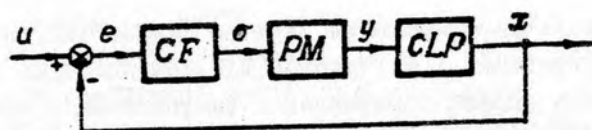


Fig.1. Block diagram of nonlinear pulse system of automatic control.

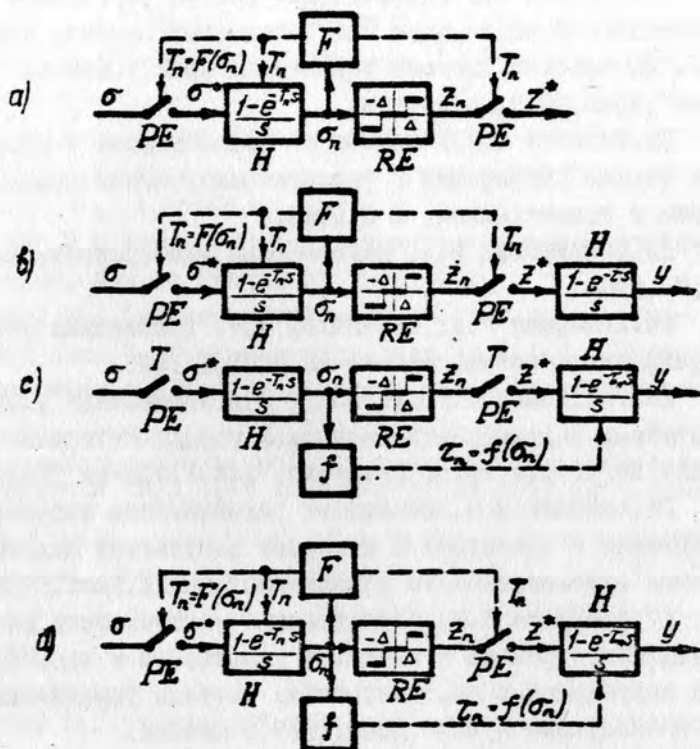


Fig.2. Block diagrams of nonlinear pulse modulators:
 a - ideal frequency-pulse modulator; b - real frequency-pulse modulator; c - width-pulse modulator; d - frequency-width pulse modulator.

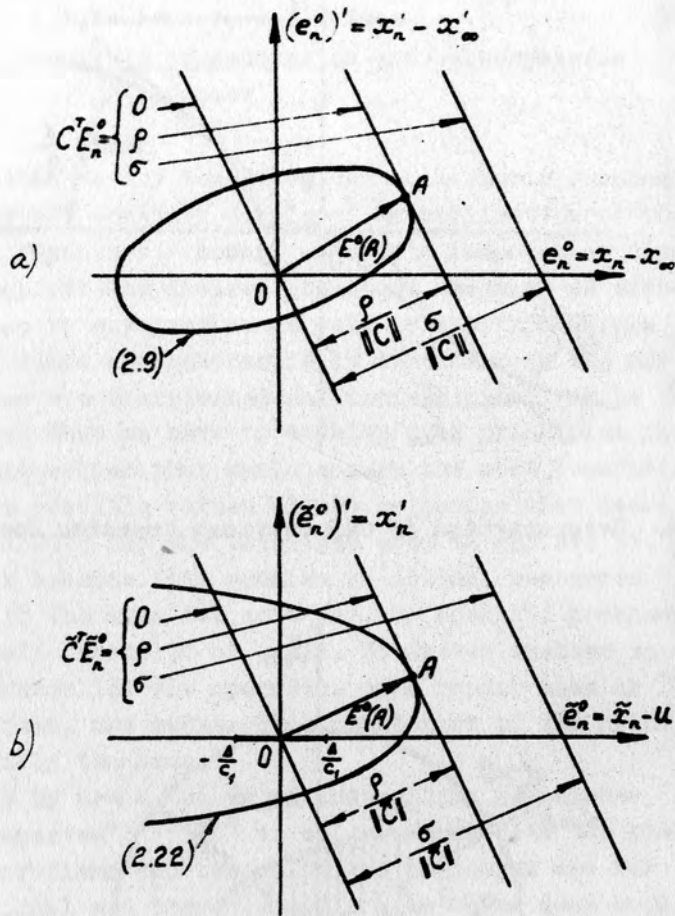


Fig.3. Graphical interpretation of stability condition (2.11) in the main (a) and simplest critical (b) cases.

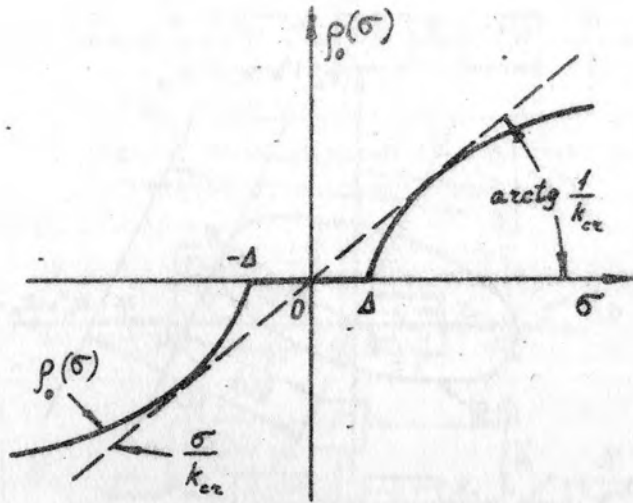


Fig.4. Determination of CLP critical transfer coefficient.

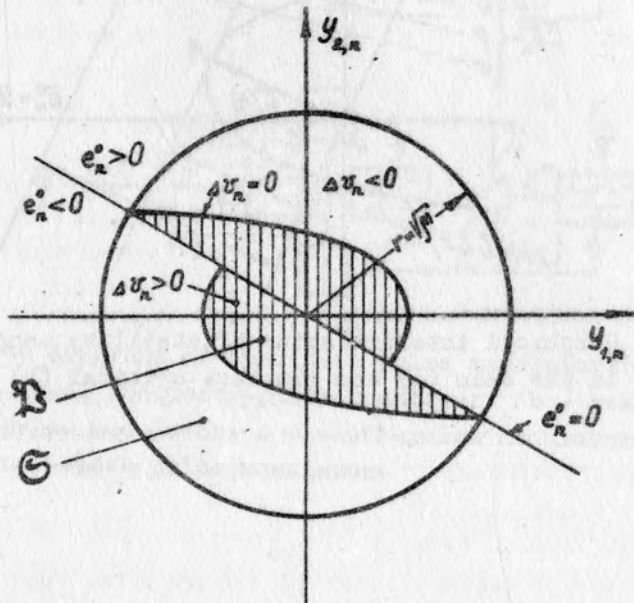


Fig.5, Determination of asymptotically stable set of limited bounded system of the 2nd order.

ON OPTIMAL RESOURCES ALLOCATION

A.Ya.Lerner, A.I.Teiman

Institute of Automation and Telemechanics

Moscow

USSR

A solution to any technological or economic problem requires certain activity displayed as utilization of the appropriate resources - human, material, temporal or financial - to achieve the desired objective. Because an effective solution to any problem as well as the expenditure of time and funds are substantially dependent on the way the resources are distributed and redistributed, it is quite natural that we have to develop such principles for resources allocation that would ensure the most beneficial value of the possible values of the criterion that describes the result and the technique used to achieve it.

In its essence this problem of optimal resources allocation is the same for most various specific problems whatever their variables or scale. No matter whether we deal with scheduling the operation of a repair team or a major project, the mathematical statement of the problem is actually the same.

Though by now a number of interesting researches have been reported^{1,2,3} etc we believe that both the statements of problems and the solutions suggested are far from being final and there is still much to be done both in theoretical research and new applications.

This paper is an attempt to give strict definitions to basic concepts and suggest possible solutions to certain theoretical problems. The research was undertaken in the Large-Scale Systems Division of the Institute of Automation and Telemechanics^{x/}.

x/ Among the collaborators of the Division who took part in the research we would note V.N.Burkov, A.A.Voronov, S.Ye.Lovetsky and others.

A number of definitions as well as classification of problems and techniques are not conventional. The writers hope that a discussion will result in a uniform point of view on these matters which is now very much desirable.

Basic Concepts and Definitions

1.1. An activity is a process described by equation of the form

$$\omega_i(t) = \frac{dx_i(t)}{dt} = \varphi_i[U_{i1}(t), \dots, U_{im}(t)] \quad (1)$$

where $x_i(t)$ is the status of the i -th activity at time t ,

$U_{ij}(t)$ is the amount of resources of the j -th form in the i -th activity at time t ,

$\omega_i(t)$ is the rate at which the i -th activity is being accomplished at time t . Since resources are involved in an activity in certain ratios the concept of resources set will be useful.

A resources set in the i -th activity is a set of resources allocations $\{U_{ij}(t)\}$ such that $U_{ij}(t) = \alpha_{ij} \vartheta_i(t)$ where $\{\alpha_{ij}\}$ are the parameters of the set and $\vartheta_i(t)$ is the capacity of the set at time t . Different sets are generally feasible in an activity, but for simplicity we will discuss only the case of one feasible set per activity. With the given parameters of a set the rate of each activity depends only on the capacity of the set and time elapsed (the rate will be assumed not depending on time explicitly). Thus

$$\omega_i(t) = f_i[\vartheta_i(t)] \quad (2)$$

where f_i is a non-decreasing continuous function of ϑ_i and

$f_i(0) = 0$. Completion of an activity implies the change of its state from the initial one $x_i(0) = 0$ to the final one

$x_i(T) = W_i$ where W_i is the volume of the activity and T is the time when all activities have been completed.

1.2. A final set of activities makes a complex of activities.

This is normally shown as a network. The state of the complex is the vector $X(t) = (x_1(t), \dots, x_n(t))$ whose components are the states of activities in the complex. The volume of the complex is the vector $W = (W_1, \dots, W_n)$. However, the concept of an equivalent volume is sometimes more con-

venient as a scalar quantity depending on the volumes of activities (e.g. $W_e = (\sum_i w_i^{\alpha})^{1/\alpha}$, $\alpha > 1$).

1.3. Let us introduce two types of constraints on resources, constraints on capacity (instantaneous constraints)

$$\sum_i \alpha_{ij} \varphi_i(t) \leq N_j(t), \quad j = 1, 2, \dots, m \quad (3)$$

and constraints on expenditure (integral constraints)

$$\sum_i \alpha_{ij} S_i \leq S_j, \quad j = 1, 2, \dots, m. \quad (4)$$

where $S_i = \int_0^T \varphi_i(t) dt$ are expenditures of capacity.

2. Description of the model

2.1. The statement of and solution to the problems of optimal resources allocation may be based on a model of a complex of activities including the following constituents

a) A network or a matrix which represents the composition, volumes and feasible sequences of activities in the complex.

b) A graph or a matrix representing the feasible distributions of resources among the activities of the complex, permissible expenditures of resources to complete these activities, and expenditures of time and funds for redistribution of resources.

c) A set of equations representing the activities of the complex: $\omega_i(t) = f_i[\varphi_i(t)]$

d) A functional J which represents the total efficiency in completion of the activity complex.

2.2. The problem of optimal resources allocation consists in selecting programs $\varphi_{ij}^k(t)$ (the amount of resources of the k -th form that are reallocated from the i -th activity to the j -th activity at time t) such that with the constraints on feasible sequences of operations (2.1a) and on feasible values of resources fluxes (2.1b) would ensure completion of the complex described by eqs of (2.1b) so as the functional J takes the minimal value.

2.3. Depending on specific conditions the functional may depend on different arguments or their combinations. In practice the most important are cases where the objectives of optimization include the following factors:

- T - the time of completing the complex;
 R - resources allocated for completion of a complex;
 P - the probability that the complex will be completed within the time interval
 K - the quality of the result.

Generally optimization of the schedule for completion of a complex should include minimization of a certain functional of the above factors

$$J(T, R, P, K).$$

In specific cases some arguments of the functional J can be neglected. Depending on the parameters that are included either in the functional J or in the constraints we can distinguish the classes of problems represented in Fig. 1 (the time parameter is assumed to be always incorporated in the problem). The problems of the class T (time) are not optimization problems. Such are those of the complex analysis (finding the time when the complex has been completed, finding the time slack, etc.). The class TP (time-probability) is made chiefly by the problems of allocation the dependent time slacks between operations with the view to maximizing the probability of meeting the schedule. The classes which incorporate the quality are relatively little known. The remaining discussion will cover the problems of the classes time-resources and time-probability (TR and TP).

3. Statement of the problems

3.1. Two basic classes of problems can be distinguished. In the problems of the first type the time T for completion of the complex is given, in the problems of the second type it is not.

Problem 1. Find $U_i(t)$, $i=1, 2, \dots, n$ that satisfy eqs (3), (4) so that the complex be completed within the time T and the optimality criterion take the minimal value.

Problem 2. Find $U_i(t)$, $i=1, 2, \dots, n$ so that

the complex be completed and the optimality criterion take the minimal value (time is not given).

3.2. The following optimality criteria are most widely used

a) $J_1 = T$ (minimization of the complex time)

b) $J_2 = \sum_j C_j \max_i \sum_j \alpha_{ij} \vartheta_i(t)$ (minimization of the resources levels)

c) $J_3 = \sum_j C_j \int_0^T (\sum_i \alpha_{ij} \vartheta_i(t))^2 dt$ (steady of resources)
(minimization of costs)

d) $J_4 = \sum_j C_j \sum_i \alpha_{ij} \alpha_{ij}$

e) $J_5 = \alpha T + J_2$

f) $J_6 = \alpha T + J_4$

In these criteria $C_j \geq 0$ represents the costs of resources of the j -th form,

α is the losses when the complex completion is delayed by a time unit. If there are intermediate objectives, e.g. the earliest possible occurrence of certain events, the criterion

$$J_8 = \sum_i \sigma_i(t_i - \Delta_i)$$

is used where $\sigma_i(t_i - \Delta_i) = 0$ at $t_i \leq \Delta_i$ and is a non-decreasing function t_i at $t_i > \Delta_i$, being the time t_i when the i -th activity is completed.

4. Complex aggregation

4.1. Complex aggregation is replacement of a complex or its part by one activity⁴. Aggregation is used when a complex consists of a large number of operations. Direct solution to the optimization problem is difficult due to its high dimensionality. Aggregation is made in three stages:

1. Ordering the states of the complex.

2. Finding the parameters of the set of the aggregated operation.

3. Finding the relation between the rate of the aggregated operation and the capacity of the set.

4.2. Let us consider an example of aggregating a complex of P sequential activities completed by resources of the same kind. The first and the second stages are unnecessary since the states of the complex have already been ordered and the single parameter of the set of the aggrega-

ted activity can be assumed to be equal to one. Let $\tau_i(\vartheta)$ denote the time for completion of the i -th activity at capacity ϑ of its set, i.e.

$$\tau_i(\vartheta) = \frac{w_i}{f_i(\vartheta)}$$

Then by definition the rate of the aggregated activity

$$f(\vartheta) = \frac{W}{\sum \tau_i(\vartheta)} \quad (5)$$

Assume the volume of the aggregated activity

$$W = \sum \beta_i w_i.$$

Estimate the aggregation error. Let the capacity of the set of the i -th activity be ϑ_i . Then the time for completion of an aggregated activity

$$\tau(\vartheta_1, \dots, \vartheta_p) = \sum_i \frac{\beta_i w_i}{f(\vartheta_i)} = \frac{1}{W} \sum_j \sum_i \beta_j w_j \tau_i(\vartheta_j).$$

The actual time of the complex

$$\tau'(\vartheta_1, \dots, \vartheta_p) = \sum_i \tau_i(\vartheta_i) = \frac{1}{W} \sum_j \sum_i \beta_j w_j \tau_i(\vartheta_i).$$

The aggregation error

$$\varepsilon = \tau'(\vartheta_1, \dots, \vartheta_p) - \tau(\vartheta_1, \dots, \vartheta_p) = \frac{1}{W} \sum_j \sum_i \beta_j w_j [\tau_i(\vartheta_i) - \tau_i(\vartheta_j)]$$

In this case the problem of optimal aggregation is in finding β_1, \dots, β_p so as to make the aggregation error minimal.

Theorem 1. Let $\tau_i(\vartheta_j) = \alpha_i + \beta_i \tau(\vartheta_j)$, $i = 1, \dots, p$ then the aggregation error is zero at $\beta_j = \frac{C \beta_j^0}{w_j}$ ($C > 0$ is an arbitrary constant). The proof is found by direct verification. This theorem makes it possible to solve the problem of optimal aggregation by representing the relations $\tau_i(\vartheta)$ in the approximate form $\alpha_i + \beta_i \tau(\vartheta)$ and assuming

$$\beta_i^0 = C \beta_i / w_i.$$

4.3. The problem of optimal aggregation in which the maximal relative downtime of various types of resources was minimized in Ref.⁵ for the case of linear relations between the rates of activities and the capacity of resources set. That ideal aggregation is possible (where the aggregation error is zero) in a network of an arbitrary form was shown in Ref.⁶ for the case of exponential relation between the rates of activities, on the one hand and the set capacity

w_i and the resources of the same type, on the other.

Let us proceed to the precise methods of solution developed for various specific statements of the problem, such as the case of independent activities, the case of ordered events, the problem of expenditure allocation and a number of problems of the combinatorial type.

5. Independent activities

5.1. Let $f_i(v_i)$ be convex upwards functions of v_i . Then it can be shown⁷ that in an optimal solution all operations are performed at constant intensity and complete simultaneously. Hence

$$f_i(v_i) = w_i/T, \quad v_i = \varphi_i(w_i/T)$$

where φ_i is a function, reverse to f_i .

The minimal time of completing a complex is defined as a minimal T for which following set of inequalities is true

$$\sum_i \alpha_{ij} \varphi_i(w_i/T) \leq N_j, \quad j = 1, 2, \dots, m. \quad (7)$$

Example 1 $f_i(v_i) = v_i, \quad v_i \leq \beta_i, \quad i = 1, 2, \dots, n$.

We have

$$\sum_i \alpha_{ij} w_i/T \leq N_j, \quad j = 1, 2, \dots, m.$$

$$w_i/T \leq \beta_i, \quad i = 1, 2, \dots, n$$

$$T_{min} = \max \left[\max w_i/\beta_i; \max_j \frac{1}{N_j} \sum_i \alpha_{ij} w_i \right] \quad (8)$$

5.2. If $f_i(v_i)$ are non-convex functions then the optimal solution generally consists of intervals of n constants. For what follows of importance is

Theorem 11. The minimal time for completion of a complex $T_{min}(w_1, \dots, w_n)$ is a convex (downwards) function of activity volumes (even if $f_i(v_i)$ are non-convex functions).

6. Ordered events

6.1. The events of a complex are ordered if the time when the i -th activity occurs is no longer than the time of the j -th event, provided that $i < j$.

Denote by Δ_S the duration of the interval between $(S-1)$ and S -th event $S = 1, 2, \dots, 2, R_S$ is a set of activities that can be completed within the S -th inter-

val Q is the set of intervals where the i -th activity can be completed, $x_{i,j}$ is the volume of the i -th activity completed in the j -th interval. Let $z_j = \{x_{i,j} : i \in R_j\}$

be given. Then for each activity we have a case of independent activities and can find the minimal duration of the j -th interval of $\Delta_j(z_j)$. The time for completing a complex

$$T = \sum_j \Delta_j(z_j) \quad (9)$$

is by Theorem 11 a convex (below) function of $x_{i,j}$. Thus we have the problem of minimizing a convex function at linear constraints

$$\sum_{j \in Q_i} x_{i,j} = W_i, \quad i = 1, 2, \dots, n. \quad (10)$$

solvable by any techniques of convex programming.

6.2. Let us take now the problem of steady use of resources. Assume that each activity requires resources of just one kind. Denote as $R_{j,j} \subset R_j$ set of activities which require resources of the j -th kind. We have to minimize

$$J_j = \sum_j c_j \sum_s \frac{1}{\Delta_s} \left(\sum_{i \in R_{j,j}} x_{i,s} \right)^2 \quad (11)$$

at constraints (10) and

$$\sum_j \Delta_j = T. \quad (12)$$

Since conditions (10) and (12) are independent we can improve successively a certain feasible solution by varying first $x_{i,j}$ at fixed Δ_j and then Δ_j at fixed $x_{i,j}$.

1st stage. Assume that feasible values of $\Delta_j \geq 0$ are given such that $\sum_j \Delta_j = T$. Then the problem of minimization decomposes into m independent problems (by the number of resources kinds). Each of these is in minimizing

$$J_j = \sum_s \frac{1}{\Delta_s} \left(\sum_{i \in R_{j,j}} x_{i,s} \right)^2 \quad (13)$$

at constraints (10). This is a problem of quadratic programming.

2nd stage. Denote by $x_{i,s}^0$ the optimal values obtained at the 1st stage.

$$B_s^2 = \sum_j c_j \left(\sum_{i \in R_{sj}} x_{i,s}^0 \right)^2 \quad (14)$$

We have the minimization problem

$$J_3 = \sum_s \frac{B_s^2}{\Delta_s} \quad (15)$$

at constraint (12). By using the technique of Lagrangian multipliers we obtain directly

$$\Delta_s^0 = \frac{B_s T}{\sum_s B_s} \quad (16)$$

Then the 1st and the 2nd stage are made in the reverse order.

6.3. The most difficult is the 1st stage. To perform this a successive improvement technique is suggested which is based on hydrodynamic analogy^{8,9} which is a modification of the quadratic programming technique¹⁰.

An effective algorithm for the 1st stage can be obtained from the following rule for resources allocation: the activities with the minimal number of a final event are performed in the first turn. Its application is illustrated by this example.

Example 3. Let us consider a network of 12 activities shown in Fig.2 ($A_i(w_i)$ denotes the i -th activity of the volume w_i). Let $\Delta_1 = \Delta_2 = \dots = \Delta_6 = 7$. We find^{x/} that $\lambda_0 = \frac{7}{6} \sum w_i = 12$. Allocate the resources and denote the events with signs. The sign (-) marks only event 3. The event closest to it on the left-hand side, event 1, is denoted as (+). Since in intervals 2 and 3 no activity with a final event above 3 is completed, we make a subnetwork of events (1,2,3) and in the remaining part of the network these events are united into one. The remainder of the network is a serial connection of two networks, one of which includes event and the united event (1,2,3), and the other the united event (1,2,3) and

x/ The algorithm was suggested by V.N.Burkov.

and events 4,5,6. Therefore the problem can be solved separately for each network. We obtain three subnetworks shown in Fig.3.

For the network $G_1, \lambda_1 = 8$ and the feasible solution is $x_{11} = 3, x_{21} = 5$.

For the network $G_2, \lambda_2 = 17$ and the feasible solution $x_{32} = 14, x_{42} = 3, x_{43} = 2, x_{63} = 15$.

For the network $G_3, \lambda_3 = 10$ and the feasible solutions $x_{74} = 5, x_{84} = 2, x_{94} = 2, x_{54} = 7, x_{55} = 7, x_{10,5} = 1, x_{11,5} = 2, x_{11,6} = 2, x_{12,6} = 8$.

The values of $\{x_{is}\}$ obtained make the optimal solution to the 1st stage.

The value of the criterion $J = 7\lambda_1^2 + 2\lambda_2^2 + 3\lambda_3^2 = 942$.

7. Problems of the class TP

7.1. Among the problems of the TP class one can distinguish the problems of maximal complex reliability, i.e. of scheduling the realization of the complex so as to ensure a minimal possibility of not meeting the schedule. Two types of problems can be distinguished.

A. The probability distribution $P(t)$ of completing the activity is known. In this case there are two problems.

Problem 1. For the given time T at which the complex is to be realized find a distribution of times for completion of activities $\{t_i\}$ such that $P\{t_n \leq T\} \rightarrow \max$.

Problem 2. For the desired level of reliability in complex P_0 find a distribution of $\{t_i\}$ such that

$$P\{t_n \leq T\} \geq P_0$$

B. Probability distribution $P(t)$ is unknown.

In this case a certain criterion $J(t, \Delta t)$ is formed that depends on the time t of completing the activities and on the amount of reserved time slack. The value of the criterion J is an estimate of the complex reliability; problems similar to Problems 1 and 2 can be formulated for it. A number of statements for problems of the A and B type as well as the algorithms for their solution were discussed in^{11,12}.

7.2. It is a specific feature of the TP time probability

and time resources classes, types A and B problems that the complex is assumed given. In practice this is preceded by the process of making the complex; the problem of scheduling this process optimally arises. Below we will present the results obtained by one of the writers¹³.

7.3. A model of complex formation

Let the complex S_0 consist of one activity at the initial time. Then in the process of scheduling it decomposes into parts, detailed and in the long run the complex S_0 contains n activities. The activities of the complex are decomposed, specified and aggregated.

Let us define these basic concepts. It is assumed that ξ is a random duration of an activity and $0 \leq a \leq \xi \leq b < \infty$, $c = b - a > 0$, $F(x) = P\{\xi \leq x\}$, $m = M\xi$, $\sigma^2 = D\xi$. Assuming that $\alpha = \sigma/(b-a)$; $\beta = (m-a)/(b-a)$ it easily follows that $0 \leq \alpha \leq 1/2$, $0 \leq \beta \leq 1$.

Specification of activities

The activity ξ_2 is specification of ξ_1 , if $b_2 - a_2 \leq b_1 - a_1$ and $\alpha_2 \leq \alpha_1$.

Activities decomposition

1. Let ξ decompose into K successive activities $\xi_1, \xi_2, \dots, \xi_K$. This decomposition is regular if
 - a. random values ξ_i are independent
 - b. $\sum_{i=1}^K a_i = a$; $\sum_{i=1}^K b_i = b$
 - c. $\alpha_i \leq \alpha$
 - d. $\sum_{i=1}^K \beta_i c_i / c = \beta$.
11. Let ξ decomposes into k parallel activities $\xi_1, \xi_2, \dots, \xi_K$. This decomposition is regular if
 - a. random values of ξ_i are independent
 - b. $\max a_i = a$; $\max b_i = b$; $c_i \leq c$
 - c. $\alpha_i \leq \alpha$
 - d. $\max [\beta_i c_i + a_i] = \beta c + a$
111. Let the networks S_i have the corresponding decomposition Ω_i while the networks S_{i+1} the decomposition Ω_{i+1} . The decomposition $\Omega_{i+1} \supset \Omega_i$ and regular, if obtained by regular breakdown of activities in the network S_i .

Activities aggregation

Activities aggregation is unification of several activities into one. Like decomposition it can be made in series or in parallel. Aggregation will be defined as operation reverse to decomposition.

7.4. Analysis of complex formation scheduling. Basic theorems.

Theorem 1. Let t_{i_s} be the time for occurrence of the i_s event, $t_{i_s}^*$ is a fixed scheduled deadline. Then for each event i_s there is a number τ_{i_s} such that whatever probability P_0 is given, there always is a sequence of regular complex network network decomposition $\Omega_0 \leftarrow \Omega_1 \leftarrow \dots$ for which, starting from Ω_0 ,

$$P\{t_{i_s} > t_{i_s}^*\} \geq P_0 \quad \text{at } t_{i_s}^* < \tau_{i_s}$$

and

$$P\{t_{i_s} < t_{i_s}^*\} \geq P_0 \quad \text{at } t_{i_s}^* > \tau_{i_s}.$$

If separate classes of decompositions are discussed then stricter statements can be proved. E.g. there is

Theorem 11. Let $\{\Omega_i\}, i=1,2,\dots$ be regular successive decompositions. Theorem 1 is true then for Ω_0 and any $\Omega_i > \Omega_0$. To prove these Theorems¹³ those networks are taken that consist only of serial or parallel activities and a network of a general form.

The results obtained show that there are certain critical deadlines for a complex and if deadlines are chosen below below critical values the probability of failures will be very high.

A number of scheduling procedures, e.g. those that use averaged indices can be shown to be incorrect in this sense.

Therefore two classes of problems appear, formulation of upper-bounds for $\bar{t} \geq \tau$ and optimal choice of t_i^* .

8. Conclusion

1. Problems of optimal resources allocation found in practice of recent years are of great economic importance. However, neither basic concepts nor strict and sufficiently general statements of problems have not been defined so far.

2. The paper is an attempt to satisfy the urgent need in formulating the problems of optimal resources allocation by using models that represent sufficiently the actual conditions for implementing a complex of activities.

It has been found that no solution to this class of problems can be based on a single mathematical tool; this requires the entire range of various tools such as mathematical programming, theory of graphs, optimal control theory, etc.

3. The efforts of scientists have undoubtedly be directed at development of techniques for solution to a number of urgent resources allocation problems such as finding paths for transportation of resources optimal in terms of maximal reliability of achieving the objective, development of techniques for solution to a general problem of optimal resources allocation, both deterministically and stochastically.

4. Some problems raised in the paper are reported in more detail in References.

References

1. Altaev V.Ya., Burkov V.N., Teiman A.I. Automation and Remote Control, vol. 27, No 5, May 1966.
2. McGee A.A. and Markarian M.D. Optimal allocation of research (engineering manpower within a multi-project organisational structures). IRE Trans. on Eng.Manag., 1962, v.9, N 3.
3. Lambourn S. Resource allocation and multi - project scheduling (RAMPS). A new tool in planning and control. Computer J., 1963, v.5.
4. Бурков В.Н., Лернер А.Я. Новые задачи теории сетевого планирования и управления. Сборник "Вопросы управления большими системами". Изд-во "Оптиприбор", 1967.
5. Чеботарев О.Г. распределение ресурсов в многотемных разработках на основе агрегирования комплекса операций. Доклад на IV Всесоюзном совещании по автоматическому управлению, Тбилиси, 1968.

6. Бурков В.Н. Оптимальное управление комплексами операций. Доклад на IV Всесоюзном совещании по автоматическому управлению, Тбилиси, 1968.
7. Burkov V.N. Automation and Remote Control, vol 27, No 7, July 1966.
8. Razumikhin B.S. Automation and Remote Control, vol 26, No 7, July 1965.
9. Razumikhin B.S. Automation and Remote Control, vol.28, No.1, January 1967.
10. Voronov A.A., Petrushinin Ye.P. Automation and Remote Control, vol. 27, No 11, 1966.
11. Тейман А.И. Календарное планирование комплексов операций. Материалы Всесоюзного семинара по теоретическим проблемам управления большими системами и исследованию операций. Изд-во "Знание", Москва, 1967.
12. Тейман А.И. Оптимальное планирование комплексов операций. Сб. "Вопросы управления большими системами". Онти-прибор", 1967.
13. Тейман А.И. Некоторые задачи управления комплексами операций в условиях неопределенности. Доклад на IV Всесоюзном совещании по автоматическому управлению, Тбилиси, 1968.

ON STOCK CONTROL THEORY

V. Avdiysky, A. Voronov, S. Lovetsky

The problem considered is to define the optimum level of raw materials to be maintained by a large plant when the demand is stochastic and an outside source of the raw materials is available.

§1. Statement of Deterministic Stock Control Problem

We shall assume that the warehouse of a large plant (say in the steel industry) holds stocks of different items of row materials P_i ($i = 1, 2, \dots, m$) (e.g. metal bars of different sizes); from which n types of product π_j ($j = 1, 2, \dots, n$) (e.g. different types of rolled metal) must be produced. Generally speaking each type of product can be manufactured from one row item not necessarily always the same one. However we shall restrict ourselves to the case when having chosen the row material P_i for the production π_j we shall now only use this row item for manufacturing this product.

Furthermore we shall assume that the demand b_j for product π_j is known and not affected by type of row material used.

Let s_{ij} denote the cost of manufacturing the product π_j from row material P_i . If for some reason it is impossible to produce π_j from the row material P_i (e.g. the length of the requested rolled bar is greater than the length of the row bar) then we shall give s_{0j} the value ∞ .

Apart from the manufacturing costs we have to take into account the cost of holding stock in the warehouse. This type of expence is usually difficult to estimate. However, this is usually small in comparison with production costs, so we can arrive at an approximate figure if we take the total cost of maintaining the warehouse during the period of production and divide it by the quantity of row materials held during this period. Then total cost of a unit π_j of P_i will be

$$c'_{ij} = s_{ij} + \sigma$$

where σ is the estimated cost of holding a unit of row material in the warehouse.

In cases where it is necessary one can take into account the expence d_i of ordering a fresh supply of row material P_i (cost of clerical work etc.).

Under some circumstances it may be more profitable to purchase raw materials from an outside source by passing warehouse. Let h_j' denote unit cost of raw source, which includes cost of raw item cost of delivery and another possible articles and η_j quantity bused. Then the following problem can be formulated:

Find x_{ij}' and η_j' minimizing total expencies expressed in functional:

$$[\sum c_{ij} x_{ij}' + \sum_i d_i \operatorname{sgn} \sum_j x_{ij}' + \sum h_j' \eta_j'] \rightarrow \min \quad (1)$$

subject to

$$\sum x_{ij}' + \eta_j' = b_j \quad (2)$$

$$x_{ij}' \geq 0, \quad \eta_j' \geq 0, \quad i = 1, \dots, m; \quad j = 1, \dots, n; \quad (3)$$

i.e. demand must be satisfied and all variables are nonnegative.

For convenience let's make the following transformation of the variables:

$$\text{let } x_{ij} = \frac{x_{ij}'}{b_j} \quad \text{and} \quad \eta_j = \frac{\eta_j'}{b_j}$$

and include a new variable

$$y_i = \operatorname{sgn} \sum_j x_{ij} = \begin{cases} 0 & \text{if } \sum_j x_{ij} = 0; \\ 1 & \text{if } \sum_j x_{ij} > 0; \end{cases}$$

Then the problem (1) - (3) will be

$$Z = [\sum_{ij} c_{ij} x_{ij} + \sum_i d_i y_i + \sum h_j \eta_j] \rightarrow \min \quad (4)$$

subject to

$$\sum_i x_{ij} + \eta_j = 1, \quad j = 1, \dots, n; \quad (5)$$

$$y_i = \{0, 1\}, \quad i = 1, \dots, m; \quad (6)$$

$$0 \leq x_{ij} \leq 1; \quad 0 \leq \eta_j \leq 1; \quad i = 1, \dots, m; \quad j = 1, \dots, n; \quad (7)$$

where

$$c_{ij} = c'_{ij} b_j, \quad \eta_j = \eta'_j b_j$$

Problem (4) - (7) is a programming problem with mixed variables. It may be effectively solved by the algorithms based on branch and bound techniques.

We shall consider in the next section a method for solving (4) - (7) which in fact is a generalisation of the algorithm given in (2) for the plant location problem.

§2. The algorithm for solving deterministic stock control problem.

The branch and bound algorithm with modifications for the problem (4) - (7) which will be described later, is an effective and finite method for solving combinatorial optimisation problems arising in the integer programming (3), scheduling theory (4), etc. which might considerably decrease the number of trials. During the last few years many papers concerned with this method and its applications were issued (1).

The branching processes will be done as usual by integer variable y_i putting it in turn equal 0 and 1. Let's assume that it is possible to find the lower bound (estimate) of the functional (4) (way of calculating this estimate will be described later) in each vertex of the obtained decision tree, i.e. for each subset of the feasible solutions set, neglecting the integer nature of the variables y_i . Then starting from the initial vertex of the decision tree U_0 one can calculate the value of the lower bound $Z_0 \leq Z$ in this vertex. If all y_i are found to be integer then the problem is solved. In the case when some y_k is fractional putting $y_k = 0$ then $y_k = 1$ one generates two new branches (U_0, U_1) and (U_0, U_2) from the vertex U_0 which terminates correspondingly in the new vertices U_1 and U_2 . (If there are several fractional y_k , then an interesting problem concerned

with the optimal strategy of choosing the next variable for branching arises). For each vertex obtained one calculates the lower bounds, which we denote \bar{x}_1 and \bar{x}_2 correspondingly. In the set of vertices $Q = \{v_1, v_2\}$ one can find the vertex with the minimal lower bound. Let it be v_1 for example, i.e.

$$\bar{x}_1 = \min(\bar{x}_1, \bar{x}_2), \quad \bar{x}_1 \geq \bar{x}_2.$$

Then branching goes on from vertex $v_1 \in Q$ when the next fractional variable, e.g. y_k ($k \neq 1$) is set to 0 and to 1. As a result two new vertices v_3 and v_4 are generated for which the lower bounds \bar{x}_3 and \bar{x}_4 are calculated as earlier. Again in the vertices set $Q = \{v_2, v_3, v_4\}$ the vertex with the minimal lower bound is found, say v_3 , i.e.

$$\bar{x}_3 = \min\{\bar{x}_2, \bar{x}_3, \bar{x}_4\} \quad \bar{x}_3 \geq \bar{x}_2$$

From this vertex v_3 branching procedure goes on.

Process described continues until coming to the "terminal" vertex of the decision tree, or alternatively speaking until working out a solution with all integer y_i .

If value \bar{x} obtained by functional (4) on this solution satisfies inequality

$$\bar{x} \leq \bar{x}_i \quad \text{for all } v_i \in Q$$

then this solution is optimal.

Before describing the procedure for calculating the estimates \bar{x}_i notice that the optimal solution of the problem (4) - (7) will contain exactly n variables x_{ij} which will be equal to 1, and all the rest will be 0. It becomes obvious if one take into account that under the conditions of the problem (4) - (7) there are no restrictions on the quantity of products shipped. This means that the optimal solution for (4) - (7) must be found among x_{ij} equal to 0 or 1 ($i = 1, \dots, m; j = 1, \dots, n$).

Now let's consider the procedure for calculating lower bounds \bar{x}_i in each vertex of the decision tree v_i . One should realise that this procedure will be repeated

many times during the processes of solving the problem. That means that it is desirable to have this procedure as simple as possible.

Let N_j denote the set of different types of row materials from which the product π_j may be manufactured. Let M_i denote the set of products which may be manufactured from the i -th type of row material ρ_i in stock and n_i number of entries M_i .

Now, for some vertex of the decision tree, let S_1 and S_0 denote the sets of y_i which, during the process of solution, were given the values 1 and 0 correspondingly and let S denote the set consisting of the remainder of y_i . Then the value of the lower bound can be calculated if one substitute in (4) the values of x_{ij} , η_j and y_i following:

$$x_{ij} = 1$$

$$\eta_j = 0, \quad \text{if } \min_{i \in S_1, US} \{ [c_{ij} + \frac{g_i}{n_i}], \eta_j \} = [c_{kj} + \frac{g_k}{n_k}]; \quad (8)$$

$$\begin{aligned} x_{ij} &= 0 \\ \eta_j &= 1, \quad \text{if } \min_{i \in S_1, US} \{ [c_{ij} + \frac{g_i}{n_i}], \eta_j \} = \eta_j \end{aligned}$$

in all other cases x_{ij} and η_j are equal 0.

$$y_i = \frac{1}{n_i} \sum_{j \in M_i} x_{ij}, \quad i \in S, \quad (9)$$

where

$$g_i = \begin{cases} d_i, & \text{if } i \in S_0 \\ 0, & \text{if } i \in S_1 \end{cases}$$

It is obvious that expressions (8) - (9) define the optimal solution for the problem (4) - (7) without integer restriction on y_i , i.e.

$$0 \leq y_i \leq 1, \quad i = 1, \dots, m \quad (10)$$

and therefore the value of the functional (4) which obtained from it is the lower bound for the problem (4) - (7).

To prove this notice that for

$$\sum_{j \in M_i} x_{ij} \leq n_i y_i$$

and for the optimal solution of problem (4) - (7), (10) we shall have equality, i.e.

$$\sum_{j \in M_i} x_{ij} = n_i y_i \quad \text{or} \quad \sum_{j \in M_i} \frac{x_{ij}}{n_i} = y_i$$

By substituting this value of y_i for $i \in S$ in (4) we shall have

$$\min Z = \sum_{i \in S_1} d_i + \min \left\{ \sum_{i \in S_1} c_{ij} x_{ij} + \sum_{i \in S} \left[c_{ij} + \frac{d_i}{n_i} \right] x_{ij} + \sum_j h_j \eta_j \right\}.$$

providing that

$$\sum_i x_{ij} + \eta_j = 1, \quad j = 1, \dots, n$$

§3. Formulation of Stock Control problem under probabilistic demand.

In the majority of practical applications the vector of demand β_i is not known in advance. We may know the probability distribution $F(\beta_1, \dots, \beta_n)$ for the vector β or sometimes only limits in which demand may change. In this case the problem (1) - (3) becomes meaningless and one must define what the solution of the problem is.

There are a number of possible approaches to the statement of the stochastic programming problems given in (5). One of the possibilities consists in finding the x_{ij} which gives minimum to the functional (1) subject to q :

$$\sum x_{ij} + \eta_j \geq b_j(q) \quad \text{for all } q \in Q \quad (2)$$

$$x_{ij} \geq 0, \eta_j \geq 0 \quad (3)$$

where q is a random variable or set of random parameters the sample of which specifies implementation of random entries in the problem, and Q is set of q values q appeared with non-zero probability.

This type of problems usually called "rigid" or one-stage stochastic programming problems because this problems assumed to be solved once and any amendments of the accepted solution not allowed even if one have got some additional data about the environment.

It seems more reasonable to take two-stage stochastic programming problem. In this kind of the problems one might imagine the process of decision making in two stages.

Firstly one can choose some solution X , not necessary satisfying all constraints for all possible samples of q . Then some sample of q and therefore implementation of the vector $b(q)$ are fixed up and the vector η which should amend the accepted decision X is included.

Assuming that we undergo by additional expense for amending the decision made on the first stage it is appropriate to state the following problem: minimize the expected value of the total expense, i.e. cost of the solution taken on the first stage plus additional cost of amendments on the second stage. Let's consider the process of decision making in the problem (1) - (3) as two-stages. The plan X is defined during the first stage. Then when the vector b becomes known and there is not enough row materials the penalty proportional to outstanding quantity

$\eta_j = b_j - \sum_i x_{ij}$ is paid (one might consider this penalty as a cost of urgent purchase of row

^{*}) From now we shall consider x_{ij} and η_j as actual (not relative) values which in §1 was denoted as x'_{ij} and η'_j .

materials from the more expensive outside source). Then one can formulate the problem of minimizing the expected values of the total expanse

$$M[\sum_{ij} c_{ij} x_{ij} + \sum_i d_i y_i + \sum_j h_j \eta_j] \rightarrow \min \quad (11)$$

subject to

$$\sum x_{ij} + \eta_j = b_j \quad (12)$$

$$x_{ij} \geq 0, \eta_j \geq 0 \quad (13)$$

where η_j is unsatisfied demand, h_j penalty for unit. Notice, that if h_j grows to infinitude the two-stage problem (11) - (13) becomes one-stage problem (1), (2'), (3).

In order to find the deterministic equivalent (i.e. the problem of mathematical programming which solution is also solution of the corresponding stochastic programming problem) for the problem (11) - (13) let's fix up some feasible solution X and consider (11) - (13) as the minimizing problem with variable η . This one may be written as follows:

$$\sum h_j \eta_j \rightarrow \min$$

subject to

$$\sum_i x_{ij} + \eta_j \geq b_j \quad \text{or} \quad \eta_j \geq b_j - \sum_i x_{ij}$$

because the rest entries of (11) not affected by η . One can easily find that the dual problem of maximizing is follows

$$\sum_j (b_j - \sum_i x_{ij}) \xi_j \rightarrow \max$$

subject to

$$\xi_j \leq h_j$$

$$\xi_j \geq 0$$

It is obvious that the solution of the last problem is given by the equalities

$$\xi_j = h_j$$

Therefore by wellknown duality theorem

$$\min_{\eta} \sum_j h_j \eta_j = \max_{\xi} \sum_j (b_j - \sum_i x_{ij}) \xi_j = \sum_j (b_j - \sum_i x_{ij}) h_j$$

Substituting the right part of the last equality in the functional (11) one get the following deterministic equivalent for the problem (11) - (13):

minimize functional (14)

$$\left[\sum_{i,j} c_{ij} x_{ij} + \sum_i d_i y_i + \sum_j h_j M[b_j - \sum_i x_{ij}] \right]_{\min} \quad (14)$$

subject to

$$x_{ij} \geq 0 \quad (15)$$

$$y_i = \{0, 1\} \quad (16)$$

As some of the variables in the problem (14) - (16) are integer this problem is the mixed programming problem again which might be solved by the algorithm based on branch and bound techniques described in the next section.

§4. Algorithm for solving problem (14)-(16).

Method for solving problem (14) - (16) is based on the branch and bound techniques. The branching procedure done by the integer variable y_i putting it in turn equal to 0 and 1 as described earlier.

Let's now consider the problems of estimating the lower bounds in this problem.

First of all let's pay attention to some features of the "truncated" mathematical expectation $M[b_j - \sum_i x_{ij}]$ for $\sum_i x_{ij} \leq b_j$, which will be usefull for

calculating the estimates.

Consider the behavior of the function $F(x) = M[b - x]$ where b is a random variable with known probability density $\varphi(b)$ for $x \leq b$.

One can easily see that

$$F(x) = \int_x^{\infty} (b - x) \varphi(b) db$$

Using this expression it is easy to calculate the first derivative $F'(x)$

$$F'(x) = \int_0^x \varphi(b) db - 1$$

Similary, calculating the second derivative

$$F''(x) = \varphi(x) \geq 0$$

one can easily see that the function $F(x)$ is convex. Hence the function

$$F(\sum_i x_{ij}) = M[b_j - \sum_i x_{ij}]$$

is a convex function of all variables.

Now consider the problem of minimizing the functional

$$\Phi(x) = \sum_i k_i x_i + F(\sum_i x_i) \rightarrow \min$$

subject to $x_i \geq 0$, and assuming that $F(\sum_i x_i)$ is convex and all $k_i > 0$. Also we assume the existence of the continuous derivative of $F(\sum_i x_i)$. The assumptions of a positive k_i and convex $F(\sum_i x_i)$ provides that the minimum exists and attained inside the cone $x_i \geq 0$ or on the bound.

Let's assume that minimum attained inside the cone.

Then because the derivative is continuous the coordinates of the minimum point must satisfy the system of the equalities

$$\frac{\partial \Phi}{\partial x_i} = 0$$

or write it in more details

$$k_i + F'_z(\sum_i x_i) = 0, \quad z = \sum_i x_i$$

It is obvious that if $k_i \neq k_j$ for $i \neq j$ this system is inconsistent, moreover, there are no two equations from this system which are compatible except case $k_i = k_j$ for $i \neq j$.

This contradiction proves that the minimum can be attained only on the bound of the cone $x_i \geq 0$. Then one can easily find out that the minimum is attained on a vector all components of which are zero, except one. The index of this coordinate may be defined by choosing minimal k_{i_0} , i.e.

$$k_{i_0} = \min_i k_i$$

and its value is a root of the equation

$$k_{i_0} + F'(x, 0) = 0 \quad (*)$$

In the case when there are several minimal k_i (which are all equal of course) it makes no difference what quantity of each is taken, provided their sum is equal to the quantity defined by equation (*).

Another point is that, obviously if one allows y_i to be a continuous variable between 0 and 1, then the problem of minimizing (14) might be split up on the n subproblems for each index j , and each such subproblem is a problem of minimizing the functional similar to $\Phi(x)$.

All points mentioned above allows us to construct an algorithm for calculating the lower bounds of functional (14) in some vertex V_i of the decision tree.

Let S_0 and S_1 denote the sets of y_i to which during the process of solution were attached values 0 and 1 correspondingly, and let S be the set of the remaining y_i . Let n_i denote number kinds of products P_j which may be produced from i -th kind of row material D_i .

The value of the lower bound may be found if one substitute in (14) following values of x_{ij} and y_i . Let

$$c_{i_0 j} = \min_i c_{ij}, \text{ then}$$

$$\begin{aligned} x_{ij} &= 0 & \text{if } i \in S_0 \\ x_{ij} &= 0 & \text{if } i \neq i_0 \\ x_{ij} &= x_{ij}^{(0)} & \text{if } i = i_0 \end{aligned} \quad (17)$$

where $x_{ij}^{(0)}$ is a root of the equation

$$\frac{\partial}{\partial x_{ij}} [c_{ij} x_{ij} + h_j M [b_j - x_{ij}]] = 0$$

If one remembers the expression for the derivative of the "truncated" mathematical expectation, it is obvious that

$x_{ij}^{(0)}$ is a root of the equation

$$\int_0^{x_{ij}} \varphi(b_j) db_j = 1 - \frac{c_{ij}}{h_j}$$

Now

$$y_i = 0 \quad \text{if } i \in S_0$$

$$y_i = 1 \quad \text{if } i \in S_1$$

(18)

$$y_i = \frac{\sum_j \text{sgn } x_{ij}}{n_i} \quad \text{if } i \in S$$

The values (17) - (18) of x_{ij} and y_i gives the solution for the minimizing problem (14) - (15) if

$0 \leq y_i \leq 1$, and therefore the value of the functional attained by it is a lower bound for the problem (14) - (16) in some vertex v_i of the decision tree.

References

1. Lowler E.L., Wood D.E., "Branch-and-Bound Methods. A Survey". Opns. Res. v. 14, No.4, 1966.
2. Efroymsen M.A., Ray T.L., "A Branch and Bound Algorithm for Plant Location". Opns. Res. v. 14, No.3, 1966.
3. Balas E., "An Additive Algorithm for Solving Linear Programs with Zero-One Variable". Opns. Res. v. 13, No. 14, 1965.
4. Ignall Edward, Schrage Linis. "Applications of the Branch-and-Bound Technique to some Flow-Shop Scheduling Problems". Opns. Res. v.13, No.3, 1965.
5. Golshtain E.G., Yudin D.B., "New Directions in The Linear Programming". "Sovetskoe Radio", ch. 6, 1966.

RESOURCE ALLOCATION IN MULTI-PROJECT
BASED ON AGGREGATION OF THE PROJECT NETWORKS

Oleg G. Tchebotarev
Central Economic Mathematical Institute
(Moscow, USSR)

The problem of PERT/CPM system design and the problem of limited resource allocation among several projects activities is very complicated but it is of the great practical importance at the same time.

The paper is concerned with a technique for the solution of resource allocation problem in multi-project when the functions of the activity rate is a linear one.

1. Basic Definitions and Concepts

The activity is a process described by the following equation:

$$v_i(t) = \frac{dx_i(t)}{dt} = f_i[r_j(t), t], \quad (1)$$

where $x_i(t)$ is the state of the i -th activity at the moment t ; $x_i(t_{in})=0$, $x_i(t_{fin})=W_i$ (t_{in} and t_{fin} is the initial and the final moment of an activity correspondingly, W_i is the volume of the i -th activity); $v_i(t)$ is the rate of the i -th activity at the moment t ; $r_j(t)$ are the resource group parameters of the i -th activity, $j=1,2,\dots,m$ (If the j -th type of resources doesn't take part in the i -th activity the parameter $r_j(t)=0$).

Realization of the activity is the variation of its state from some initial value $x_i(t_{in})=0$ up to the final one $x_i(t_{fin})=W_i$.

$\vec{r}_i(t) = \{r_j(t)\}$ is the m -vector of the i -th activity resource group.

Let the modulus of the vector $\vec{r}_i(t)$ depends on the

value of t only. So vector $\vec{r}(t)$ can be represented as $\vec{r}_i(t) = \rho_i(t) \vec{\alpha}_i$, where $\rho_i(t)$ is the power of the i -th activity resource group at the moment t , $\vec{\alpha}_i = \{\alpha_{ij}\}$ is the m -vector of resource group parameters for the i -th activity and resource group parameter α_{ij} is the value of the j -th resource when $\rho_i(t) = 1$.

Project is a partially ordered set which consists of the final number of the activities. We can represent the project as a network. Realization of the project is the variation of its state $x(t) = [x_1(t), \dots, x_n(t)]$ from some initial value $x(0) = 0$ up to the final one $x(T) = W$, $W = (W_1, W_2, \dots, W_n)$, where W is the volume of the project and T is the final moment of the project. The project is finished if all its activities are finished.

Multi-project is a set independent projects, which must be realized by common resources.

Resources constraints. There are two types of resources constraints

$$\sum_{p=1}^{\ell} \sum_{i=1}^{n_p} \alpha_{ij} \rho_i(t) \leq N_j(t), \quad j = 1, 2, \dots, m, \quad (2)$$

$$\text{or } \sum_{p=1}^{\ell} \sum_{i=1}^{n_p} \alpha_{ij} q_i \leq S_j, \quad j = 1, 2, \dots, m, \quad (3)$$

where $q_i = \int_0^T \rho_i(t) dt$ is the power consumption of the i -th activity ($i = 1, 2, \dots, n_p$), S_j is the permissible consumption of the j -th resource in the multi-project, N_j is the given value of the j -th resource in the multi-project at the moment t , i is the number of the project activity ($i = 1, 2, \dots, n_p$), p is the number of the project ($p = 1, 2, \dots, \ell$).

If there are constraints of the type (2) only, a resources are considered, as a power type resources.

If there are constraints of the type (3) only, a resources are considered as a cost type resources.

The project aggregation is the representation of a project network by one activity. The values W , $\vec{\alpha}$ and function

$f[p(t), t]$ are defined for this activity using the given values of $W_i, \bar{\alpha}_i$ and functions $u_i(t) = f_i[q_i(t), t]$ for each project activity.

The definitions and concepts introduced are basically analogous to those ones used by V. Burkov and A. Lerner¹.

2. Resource Allocation Problem in Multi-Project

Let the multi-project consisting of ℓ projects each having the volume W_1, W_2, \dots, W_ℓ must be realized under definite constraints in resources. The problem is to allocate resources to projects activities so that to minimize a criterion under constraints (2) and/or (3).

Several functions may be used as the criterion, for example $\max_p T_p$ and $\sum_{p=1}^{\ell} \varphi_p(T_p)$, where T_p is the final moment of the p -th project and $\varphi_p(T_p)$ is a non-decreasing function of the T_p . $\varphi_p(T_p)$ is a penalty function for the delay of the p -th project (for example, $\varphi_p(T_p) = 0$, if $T_p \leq T'_p$ and $\varphi_p(T_p) = \alpha_p(T_p - T'_p)$, if $T_p > T'_p$, where T'_p is the given final moment the p -th project and α_p is a constant).

The aggregation of the project networks permits to obtain the solution of the problem as a sequence of the following actions.

1. Aggregate the project networks that is given the values of $W_i, \bar{\alpha}_i$ and functions $u_i(t) = f_i[q_i(t), t]$ for each activity define the values of $W_p, \bar{\alpha}_p$ and functions $v_p(t) = f_p[q_p(t), t]$.

2. When aggregation is completed solves the resource allocation problem with ℓ independent activities. The result of this step is the values of $N_{pi}(t)$ and S_{pi} for each project.

3. Using the values found at the previous step solves the allocation problem for each separate project.

Thus the problem of resource allocation with $n = \sum_{p=1}^{\ell} n_p$ activities is transformed to ℓ resource allocation problems with n_p activities. As an example of such the approach consider the solution of the resource allocation problem in the multi-project with resources of cost type and with the criterion

$\max_p T_p \rightarrow \min$. It is assumed that each activity in the network is subject to a continuous upward - concave time - cost relationship.

In this case the functions $S_p(T_p)$ may be found by means of Berman's algorithm². Because each $S_p(T_p)$ is the non-decreasing function of T_p , all projects have the same final moment that is $T_1 = T_2 = \dots = T_\ell = T$, if the final moment of the multi-project is minimized³. Thus the value of T may be found from the equation

$$S_1(T) + S_2(T) + \dots + S_\ell(T) = S.$$

ℓ separate the resource allocation problems may be solved by means of Berman's algorithm too².

3. Aggregation of the project network

The project aggregation process consist of the following steps:

1. Ordering of the project states.
2. Determination of the permissible vector of resource group parameters for the project.

3. Determination of the function $f_p[\rho_p(t), t]$, ($p=1,2,\dots,\ell$).

I step. Each activity is the ordered sequence of its states that is $x_i^1 < x_i^2$ or $x_i^1 > x_i^2$ is valide for any two states: x_i^1 and x_i^2 . The same condition must be valide for the project states if the project is transformed to single activity. One may easily show that to make ordering of the project states it is nesessory and sufficient:

- 1) to order the events of the project network
- 2) to find the values of y_{is} corresponding to those parts of the activities volumes, which must take place within the s -th interval (within the interval from the $(s-1)$ -th to s -th event, $s=1,2,\dots,g$). The values of y_{is} must satisfy to condition $\sum_{s \in Q_i} y_{is} = W_i$, where Q_i is the set of the intervals, within which the i -th activity takes place.

3) to define the policy of the activity realization within each interval. This may be done by means of parametric equations $x_{is} = \Psi_{is}(\beta_s)$, where $i \in R_s$ (R_s is the set of activities, which may take place in the s -th interval), β_s is a parameter, which has the same value for each activity of the set R_s , ($0 \leq \beta_s \leq 1$), Ψ_{is} are continuous non-decreasing functions of the β_s ; $\Psi_{is}(0) = 0$ and $\Psi_{is}(1) = y_{is}$.

II step. Suppose that the values of y_{is} , $\vec{\alpha}_i$ and functions $v_i(t) = f_i[q_i(t), t]$ are given and the values of α_j should be found.

Note. We shall suppose hereafter that the functions $f_i[q_i(t), t]$ are

$$f_i(p_i) = p_i. \quad (4)$$

Consider the case of the simultaneous realization of the activity parts y_{is} within each interval $s \in Q_i$. The condition of simultaneity may be written as

$$p_{is} = \beta_s y_{is}, \quad 0 \leq \beta_s \leq 1, \quad (5)$$

where p_{is} is the power of the i -th activity resource group in the s -th interval. From (4) and (5) one has: $p_{is} = \beta_s y_{is}$. Let $\sum_{i \in R_s} \beta_s y_{is} \leq \rho$ if the power ρ of the project is given. Then the minimum value of the Δ_s (where Δ_s is the value of the s -th interval) is equal to $1/\beta_s$, that is

$$\Delta_s(\rho) = \frac{1}{\rho} \sum_{i \in R_s} y_{is}.$$

The value of the j -th resources required in the s -th interval is

$$N_{js}(\rho) = \sum_{i \in R_s} \alpha_{ij} p_{is} = \rho \frac{\sum_{i \in R_s} \alpha_{ij} y_{is}}{\sum_{i \in R_s} y_{is}}.$$

If the project power of the resource group has the given value ρ , the value of the j -th resources required for realization of the project is equal to the maximum value of $N_{js}(\rho)$, that is

$$N_j = \max_s \rho \frac{\sum_{i \in R_s} \alpha_{ij} y_{is}}{\sum_{i \in R_s} y_{is}}.$$

Because it is defined that $N_j(\rho) = \alpha_j \rho$ one has

$$\alpha_j = \max_s \frac{\sum_{i \in R_s} \alpha_{ij} y_{is}}{\sum_{i \in R_s} y_{is}}.$$

III step. If $\Delta_s(\rho)$ and α_j are known, the rate of the project is

$$f(\rho) = \frac{W}{T(\rho)} = \frac{W}{\sum_{s=1}^j \Delta_s(\rho)}.$$

4. The aggregation of the project network with linear functions of the activity rates

Let each activity has a linear power-rate relationship, that is $f_i(\rho_i) = \rho_i$ and let the network events are ordered.

We shall consider the case of serial activities realization within each interval. In that case the s -th interval is splitted into several subintervals (the number of the subinterval is equal to the number of the activities $i \in R_s$) under the condition $\sum_{i \in R_s} \Delta_{is} = \Delta_s$. One has $\Delta_{is} = y_{is}/\rho$ and $\Delta_s = \frac{1}{\rho} \sum_{i \in R_s} y_{is}$. The project time is

$$T = \sum_{s=1}^j \Delta_s = \frac{1}{\rho} \sum_{s=1}^j \sum_{i \in R_s} y_{is} = \frac{1}{\rho} \sum_{i=1}^n \sum_{s \in Q_i} y_{is} = \frac{1}{\rho} \sum_{i=1}^n W_i.$$

The maximum value of the j -th resource taking part in the project is

$$N_j = \max_s \max_{i \in R_s} \alpha_{ij} \rho = \rho \max_i \alpha_{ij}.$$

Taking into account that $N_j(\rho) = \alpha_j \rho$, one has $\alpha_j = \max_i \alpha_{ij}$.

As each component of the vector $\vec{\alpha}$ is the maximum of the corresponding components of vectors $\vec{\alpha}_i$, the use of resources is irregular (if corresponding components of $\vec{\alpha}_i$ is not equal). If the j -th resource takes part in the i -th

activity the value of the resource used within each subinterval $\Delta_{ij} \in \Theta_i$ is $\alpha_{ij} \rho$. However the value of the resource used may be $\alpha_j \rho$ if α_j and ρ are given. Thus we can define the values of the j -th resource standing idle during the time W_i/ρ of the i -th activity realization. This value is equal to $(\alpha_j - \alpha_{ij})\rho$. The value $(\alpha_j - \alpha_{ij})W_i$ will be called by the loss of the j -th resource for the i -th activity. The relative loss of the j -th resource for the project is

$$\delta_j = \frac{\sum_{i=1}^n (\alpha_j - \alpha_{ij})W_i}{\sum_{i=1}^n \alpha_j W_i} = 1 - \frac{\sum_{i=1}^n \alpha_{ij} W_i}{\alpha_j W}. \quad (8)$$

Because the parameters α_{ij} defines the correlation of different resources only one may take $\alpha'_{ij} = \alpha_{ij} \mu_i$ as parameters. (In this case the function of an activity rate is $f_i(\rho_i) = \mu_i \rho_i$. After substitution the equation (8) takes the form:

$$\delta_j = 1 - \min_i \frac{\sum_{i=1}^n \alpha_{ij} W_i}{\alpha_j \mu_i \sum_{i=1}^n W_i / \mu_i}.$$

The the problem arises to define μ_i so that

$$\max_j \delta_j \rightarrow \min. \quad (9)$$

The optimal solution of this problem is obviously the same as the optimal solution of the problem

$$\min_j \min_i \frac{\sum_{i=1}^n \alpha_{ij} W_i}{\alpha_j \mu_i \sum_{i=1}^n W_i / \mu_i} \rightarrow \max. \quad (10)$$

It is evident that the multiplying μ_i by a constant doesn't change the (10). So one may suppose that $\sum_{i=1}^n W_i / \mu_i = 1$. Denote $\eta_i = W_i / \mu_i$ and $\epsilon_{ij} = \frac{1}{\alpha_j W_i} \sum_{i=1}^n \alpha_{ij} W_i$.

Thus we have the problem:

$$\min_j \min_i \epsilon_{ij} \eta_i \rightarrow \max, \\ \text{under constrain } \sum_{i=1}^n \eta_i = 1.$$

Having changed the order of minimizing we come to the problem $\min \lambda_i \eta_i \rightarrow \max$ under constrain $\sum_{i=1}^n \eta_i = 1$, where $\lambda_i = \min \sigma_{ij}$.

The value of $\lambda_i \eta_i$ remains constant for each activity at the optimal solution.

Denote $z = \lambda_i \eta_i$. We have $\eta_i = z_i / \lambda_i$, $\sum_{i=1}^n \eta_i = z_i \sum_{i=1}^n 1/\lambda_i = 1$, $z = \left(\sum_{i=1}^n 1/\lambda_i \right)^{-1}$, $\eta_i = \left(\lambda_i \sum_{i=1}^n 1/\lambda_i \right)^{-1}$, and $\mu_i = W_i / \eta_i = W_i \lambda_i \sum_{i=1}^n 1/\lambda_i$ that is

$$\mu_i = W_i \left(\min_j \frac{\sum_{i=1}^n \alpha_{ij} W_i}{\alpha_{ij} W_i} \right) \sum_{i=1}^n \max_j \frac{\alpha_{ij} W_i}{\sum_{i=1}^n \alpha_{ij} W_i}.$$

The relative loss of a resource doesn't exceed the value

$$\delta = 1 - z = 1 - \left(\sum_{i=1}^n \max_j \frac{\alpha_{ij} W_i}{\sum_{i=1}^n \alpha_{ij} W_i} \right)^{-1}. \quad (11)$$

The parameters of the project resource group are

$$\alpha_j = \max_i \alpha_{ij} W_i \left(\min_j \frac{\sum_{i=1}^n \alpha_{ij} W_i}{\alpha_{ij} W_i} \right) \sum_{i=1}^n \max_j \frac{\alpha_{ij} W_i}{\sum_{i=1}^n \alpha_{ij} W_i}.$$

The duration of the project is

$$T(\rho) = \sum_{i=1}^n W_i / \rho \mu_i = \frac{1}{\rho} \sum_{i=1}^n \eta_i = \frac{1}{\rho}.$$

Consider the case of simultaneous realization of activities $i \in R_s$ within each interval s .

The condition of simultaneity and values of Δ_s and α_j are defined by relationships (5), (6) and (7). In this case the relative loss of the j -th resource of the project is

$$\delta_j = 1 - \frac{\sum_{s=1}^S \sum_{i \in R_s} \alpha_{ij} y_{is}}{\left(\max_s \frac{\sum_{i \in R_s} \alpha_{ij} y_{is}}{\sum_{i \in R_s} y_{is}} \right) \sum_{s=1}^S \sum_{i \in R_s} y_{is}}. \quad (12)$$

Like in the previous case, we take as a parameter of the resource group $\alpha'_j = \alpha_j \mu_i$ substituting ρ_i by $\rho'_i = \mu_i \rho_i$.

at the equation (4). After that equations (4), (6), (7) and (12) are replaced by

$$f_i(\rho_i) = \rho_i \mu_i, \quad (4a)$$

$$\Delta_s(\rho) = \frac{1}{\rho} \sum_{i \in R_s} \tilde{y}_{is}, \quad (6a)$$

$$\alpha_j = \max_s \frac{\tilde{x}_{js}}{\sum_{i \in R_s} \tilde{y}_{is}}, \quad (7a)$$

$$\delta_j = 1 - \frac{A_j}{\left(\max_s \frac{\tilde{x}_{js}}{\sum_{i \in R_s} \tilde{y}_{is}} \right) \sum_{s=1}^g \sum_{i \in R_s} \tilde{y}_{is}}, \quad (12a)$$

where $\tilde{y}_{is} = y_{is}/\mu_i$, $\tilde{x}_{js} = \sum_{i \in R_s} \alpha_{ij} y_{is}$ and $A_j = \sum_{s=1}^g \tilde{x}_{js}$.

Consider the problem: to find μ_i so that $\max_j \delta_j \rightarrow \min$ or

$$\min_j \min_s \frac{A_j \sum_{i \in R_s} \tilde{y}_{is}}{\tilde{x}_{js} \sum_{i=1}^n W_i / \mu_i} \rightarrow \max. \quad (13)$$

Because the value of (13) doesn't change when each μ_i is multiplied by the same constant one may assume that $\sum_{i=1}^n \eta_i = 1$, where $\eta_i = W_i / \mu_i$.

The problem (13) is equivalent to the problem

$$\min_s \left(\min_j \frac{A_j}{\tilde{x}_{js}} \right) \sum_{i \in R_s} \frac{y_{is} \eta_i}{W_i} \rightarrow \max \text{ under } \sum_{i=1}^n \eta_i = 1. \quad (14)$$

Denote $B_s = \min_j \frac{A_j}{\tilde{x}_{js}}$, $\epsilon_{is} = \frac{y_{is}}{W_i}$ and write (14) in the form

$$\min_s B_s \sum_{i \in R_s} \epsilon_{is} \eta_i \rightarrow \max \text{ under } \sum_{i=1}^n \eta_i = 1.$$

Let $C_{is} = B_s \epsilon_{is}$. So the problem is

$$\min_s \sum_{i \in R_s} C_{is} \eta_i \rightarrow \max \text{ under } \sum_{i=1}^n \eta_i = 1.$$

Denote $\min_s \sum_{i \in R_s} C_{is} \eta_i = \gamma$. We have the next problem: to maximize $\Phi = \gamma$ under constraints

$$\begin{cases} \sum_{i \in R_s} C_{is} \eta_i = \gamma, \\ \sum_{i=1}^n \eta_i = 1. \end{cases} \quad (15)$$

Constrains (15) are equivalent to

$$\begin{cases} v - \sum_{i \in R_s} C_{is} \eta_i \leq 0, \\ \sum_{i=1}^n \eta_i = 1. \end{cases} \quad (16)$$

Transfer the problem to a dual one. Let $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_g$ and ξ are new variables (ξ correspond to the second constrain). We have the next problem: to minimize ξ under constrains

$$-\sum_{s \in Q_i} C_{is} \varepsilon_s + \xi \geq 0 \quad (17)$$

$$\text{and } \sum_{s=1}^g \varepsilon_s \geq 1. \quad (18)$$

The constrain (17) may be replaced by equation

$$\xi = \max_i \sum_{s \in Q_i} C_{is} \varepsilon_s$$

and constrain (18) may be replaced by equation

$$\sum_{s=1}^g \varepsilon_s = 1. \quad (18a)$$

So the dual problem is

$$\max_i \sum_{s \in Q_i} C_{is} \varepsilon_s \rightarrow \min \quad \text{under } \sum_{s=1}^g \varepsilon_s = 1.$$

Take v, η_i as basic variables of the activities that may realized within a single interval only. Denote the set of those activities by H (remark that $|H| = g$).

The basic solution is the one of the next problem:

$$\min_{i \in H} C_{is} \eta_i \rightarrow \max \quad \text{under } \sum_{i \in H} \eta_i = 1.$$

The all values $C_{is} \eta_i$ is readily seen to be the same in the optimal solution of the problem. So the basic solution is

$$\eta_i = \begin{cases} \left(C_{is} \sum_{i \in H} \frac{1}{C_{is}} \right)^{-1} & \text{for } i \in H, \\ 0 & \text{for } i \notin H. \end{cases} \quad (19)$$

In this case $v = \left(\sum_{i \in H} \frac{1}{C_{is}} \right)^{-1}$.

Prove that the solution (19) is optimal.

For this purpose define the solution of the dual problem.

It follows from the relation of duality that the strict equality in equation (17) is true for the activities belonged to set H . Taking into account that each activity from set H may be realized within a single interval only, one has

$$C_{is} \varepsilon_s = \xi \quad \text{for } i \in H. \text{ So} \\ \varepsilon_s = \frac{\xi}{C_{is}}. \quad (20)$$

After substituting (20) in (18a) we have

$$\xi = \left(\sum_{s=1}^j \frac{1}{C_{is}} \right)^{-1} = \left(\sum_{i \in H} \frac{1}{C_{is}} \right)^{-1} = \gamma. \quad (21)$$

Since the solutions of the direct and dual problems are the same, the solution (19) is optimal.

Let η_s is the η_i of the activity which realizing within the s -th interval. Correspondingly C_s is the C_{is} of the same activity. We have

$$C_s = \min_j \frac{\sum_{i \in R_s} \alpha_{ij} y_{is}}{\sum_{i \in R_s} \alpha_{ij} y_{is}} = \min_j \frac{\sum_{i=1}^n \alpha_{ij} W_i}{\sum_{i \in R_s} \alpha_{ij} y_{is}};$$

$$\eta_s = \left(\sum_{s=1}^j \max_j \frac{\sum_{i \in R_s} \alpha_{ij} y_{is}}{\sum_{i=1}^n \alpha_{ij} W_i} \right)^{-1} \max_j \frac{\sum_{i \in R_s} \alpha_{ij} y_{is}}{\sum_{i=1}^n \alpha_{ij} W_i};$$

$$\alpha_j = \left(\max_s \frac{\sum_{i \in R_s} \alpha_{ij} y_{is}}{\max_j \frac{\sum_{i \in R_s} \alpha_{ij} y_{is}}{\sum_{i=1}^n \alpha_{ij} W_i}} \right) \sum_{s=1}^j \max_j \frac{\sum_{i \in R_s} \alpha_{ij} y_{is}}{\sum_{i=1}^n \alpha_{ij} W_i};$$

$$T(\rho) = \frac{1}{\rho} \sum_{s=1}^j \sum_{i \in R_s} \frac{y_{is} \eta_i}{W_i} = \frac{1}{\rho} \sum_{s=1}^j \eta_s = \frac{1}{\rho}.$$

Suppose that $W=1$ for the project. So we have $f(\rho) = \rho$

for the project. The relative loss for a resource doesn't exceed

$$\delta = 1 - \nu = 1 - \left(\sum_{s=1}^S \max_j \frac{\sum_{i \in R_s} \alpha_{ij} y_{is}}{\sum_{i=1}^n \alpha_{ij} W_i} \right)^{-1}$$

The results to be found here permit us to state the next problem of the optimal aggregation: define y_{is} so that $\sum_{s=1}^S \max_j \frac{1}{D_j} \sum_{i \in R_s} \alpha_{ij} y_{is} \rightarrow \max$ under constrain $\sum_{s \in G_i} y_{is} = W_i$, where $D_j = \sum_{i=1}^n \alpha_{ij} W_i$.

In conclusion let's compare the value of the relative losses under the serial and simultaneous realization within each interval. For this purpose consider permissible solution in case of the simultaneous realization of the activities: $y_{is} = W_i$, if the activity is realized within the s -th interval and $y_{is} = 0$ other wise. In such a case the equation of relative losses is

$$\delta_{sim} = 1 - \left(\sum_{s=1}^S \max_j \frac{\sum_{i \in R_s^*} \alpha_{ij} W_i}{\sum_{i=1}^n \alpha_{ij} W_i} \right)^{-1}, \quad \text{where } R_s^* = \{i: y_{is} > 0\}.$$

The equation of relative losses in the case of series realization of activities may be represented in the next form

$$\delta_{ser} = 1 - \left(\sum_{s=1}^S \sum_{i \in R_s^*} \max_j \frac{\alpha_{ij} W_i}{\sum_{i=1}^n \alpha_{ij} W_i} \right)^{-1}.$$

So far as the maximum of the sum is larger than the sum of maxima, $\delta_{sim} \leq \delta_{ser}$. Consequently, the case of simultaneous realization is not worse than the case of serial realization in the sense of the value of the relative resources loss. (If the values of y_{is} are optimal).

References

1. Бурков В.Н., Лернер А.Я. Новые задачи теории сетевого планирования. В сб. "Вопросы управления большими системами", 1968.
2. Berman E.B. Resource allocation in a PERT network under continuous activity time-cost functions. Management Science, Vol. 10, No 4, (July 1964), 734-745.
3. Бурков В.Н. Распределение ресурсов как задача оптимального быстрогодействия. Автоматика и телемеханика, № 7, 1966.

V.N.Burkov

Institute of Automation and Remote Control

Moscow

USSR

1. Statement of the problem.

The problems of Resource allocation in the project were set and classified Comparatively not long ago¹. At the moment resources allocation problem is one of the main part of the network Planning and Control Theory (NPC-Theory).

A minimum-time resource allocation problem is considered in the paper. Denote

$w_i(t)$ the speed of the i -th activity at the moment t ,

$v_i(t)$ the value of the i -th activity resources at the moment, t ,

W_i the volume of the i -th activity,

$N(t)$ the general level of the resources at the moment t (given function),

$s_i(t)$ the i -th activity cost at the moment t ,

$S(t)$ the project cost at the moment t (given function)/

The following relations take place:

$$w_i(t) = f_i[v_i(t)], \quad (1)$$

$$W_i = \int_0^T f_i[v_i(t)] dt, \quad (2)$$

$$\sum_{i=1}^n v_i(t) \leq N(t) \quad (3)$$

$$s_i(t) = \int_0^t v_i(\tau) d\tau, \quad (4)$$

$$\sum_{i=1}^n s_i(t) \leq S(t), \quad (5)$$

$$v_i(t) \geq 0, W_i \geq 0, N(t) \geq 0, S(t) \geq 0, i=1, \dots, n, t \in (0, T), \quad (6)$$

where f_i - nondecreasing function of v_i ($v_i \geq 0$), $f_i(0) = 0$,
 T - the project time.

The problem 1. Find $v_i(t)$ satisfying (3) so that the project time is minimal.

The problem 2. Find $v_i(t)$ satisfying (5) so that the project time is minimal.

The problem 3. Find $S_i(T)$ under the Condition

$\sum_1^n S_i(\tau) \leq S(\tau)$, so that the project time is minimal.

2. Solution of problem 1.

Let $N(t) = N = \text{const.}$ Suppose the events of the network are ordered; that is the moment of the s -th event t_s is less than the moment of the k -th event t_k if $s < k$. Denote

R_s - the set of activities which can be done at the s -th interval that is at the interval (t_{s-1}, t_s) , $s = 1, 2, \dots, m$, where $(m+1)$ is the number of the events.

Q_i - the set of the intervals where the i -th activity can be done $i = 1, 2, \dots, n$,

x_{is} - the volume of the i -th activity done at the s -th interval, $i \in R_s$, $s = 1, 2, \dots, m$.

$\Delta_s(z_s)$ - the minimal duration of the s -th interval
 $(z_s = \{x_{is} : i \in R_s\}, x_{is} \geq 0, s = 1, 2, \dots, m)$.

Note that

$$\sum_{s \in Q_i} x_{is} = W_i, \quad i = 1, 2, \dots, n. \quad (7)$$

Given $\{x_{is}\}$ the minimal time problem for each interval can be stated³. It can be shown that $\Delta_s(z_s)$ is a convex function of x_{is} .

Then the project time

$$T(z_1, z_2, \dots, z_m) = \sum_{s=1}^m \Delta_s(z_s) \quad (8)$$

is the convex function of x_{is} as well.

So we have the problem of convex function minimization under the linear constraints (7). This problem can be solved by any method of the convex programming. If some $\Delta_s(z_s) = 0$ the solution can be improved by changing the order of the events $(s-1)$ and s and Solving the problem under new ordering⁴. However the solution obtained by this procedure is not necessarily optimal.

3. The case of power functions.

Let $f_i(v_i) = v_i^{\alpha}$, $\alpha > 1$, $i = 1, 2, \dots, n$. Then³

$$\Delta_s(z_s) = N^{-1/\alpha} \left[\sum_{i \in R_s} x_{is}^{\alpha} \right]^{1/\alpha} \quad (9)$$

Applying the Lagrange Multipliers Method one has

$$\frac{\partial \Delta_s(z_s)}{\partial x_{is}} = N^{-1/\alpha} \left(\frac{x_{is}}{\Delta_s} \right)^{\alpha-1} = \bar{\pi}_i \quad (10)$$

Denote v_{is} - the i -th activity resource value at the s -th interval.

From (10) it follows that $\tilde{v}_{is} = (N \tilde{x}_i^\alpha)^{1/\alpha-1}$.
 Consequently doesn't \tilde{v}_{is} depend on S . So the important property of the optimal solution can be formulated: each activity has a constant level of resources. (property 1).

Suppose that $\Delta_s(\bar{z}_s) > 0$, $s = 1, 2, \dots, m$ under the optimal solution. Then taking into account property 1 property 2 can be stated as follows: the resource values $\{\tilde{v}_i\}$ are the flow N in the network (property 2).

It can be proved that property 2 takes place ~~when~~ when some $\Delta_s = 0$.

Properties 1 and 2 are valid for any ordering of the events. Consequently they are valid without any assumption of ordering. Therefore the ordering of the events isn't supposed here after.

To get other properties consider q -dimensional region Y corresponding to the project network². Define the distance ρ between any two points y^1 and y^2 of the region as

$$\rho(y^1, y^2) = \left(\sum_i |y_i^1 - y_i^2|^\alpha \right)^{1/\alpha} \quad (11)$$

Comparing (11) and (9) one can notice that some trajectory corresponds to any set $\{x_{is}\}$. From this fact it follows that:

the shortest trajectory connecting initial and final points corresponds to the optimal solution (property 3).

Denote W_e - the length of the shortest trajectory and call it the equivalent project volume.

Then

$$T_{\min} = W_e N^{-1/\alpha} \quad (12)$$

As it follows from (12) the phase point velocity is equal to $N^{1/\alpha}$. As the shortest trajectory is the only one under any velocity it follows that values $\{x_{is}\}$ don't depend on the resources level $N(t)$ under the optimal solution (property 4)

The minimal project time is defined by the equation

$$\int_0^T N^{1/\alpha}(t) dt = W_e \quad (13)$$

Describe an algorithm of foundation W_e . First of all represent the network in the new coordinat system where axis are parallel. The number of axis is equal to the project

dimension q .

Each point of the ordinary (Decart) system is represented by a curve (front) in the parallel coordinat system. Such representation is intermediate one between network and phase-space representation. Fig:1 shows all three types of the project representation. Here $F_0 = (0,0,0)$ - the initial state of the project, $F_k = (y_1^k, y_2^k, y_3^k)$ - the final state of the project. The main idea of the algorithm is based on shortest route foundation.

In order to find the shortest route between points 0 and A define the line OA in the ordinary coordinat system. This line will be represented by ~~system~~ the set of fronts in the parallel coordinat System

$$y_j(t) = t y_j^k, \quad 0 \leq t \leq 1, \quad j=1, 2, \dots, q. \quad (14)$$

If this line is completely inside Y corresponding solution will be optimal. In the other case the shortest trajectory will consist of two segments OD and DA. The following relation takes place

$$\frac{\partial C}{\partial B} = \frac{\partial D}{\partial A} \quad (15)$$

Point D defines the front $F_1 = (y_1^1, y_2^1, y_3^1)$ where $y_1^1 = W_1$,

$$y_2^1 = W_2, \quad y_3^1 = W_5 \frac{(W_1^\alpha + W_2^\alpha)^{1/\alpha}}{(W_1^\alpha + W_2^\alpha)^{1/\alpha} + (W_3^\alpha + W_4^\alpha)^{1/\alpha}} \quad (16)$$

Notice that (16) is equivalent to (15). Similary one can define a number of basic fronts such that . The trajectory between them is a line. The solution obtained can be improved by correcting the position of the each ^{front} basic with respect to the position of two neighboring ones:

4. The relation between problem 1 and 3.

Define the activity time-cost function

$$S_i(\tau_i) = v_i(\tau_i) \cdot \tau_i = \tau_i \int^{-1} \left(\frac{W_i}{\tau_i} \right) = \frac{W_i^\alpha}{\tau_i^{\alpha-1}} \quad (17)$$

taking into consideration property 1,

So the activity time-cost function is the power convex function of τ_i .

Consider minimal cost problem under the given project time T . The optimality conditions can be written as follow⁵.

$$\sum_{i \in V_p^-} \frac{ds_i(\tau_i)}{d\tau_i} = \sum_{i \in V_p^+} \frac{ds_i(\tau_i)}{d\tau_i}, \quad p=2, \dots, m-1, \quad (18)$$

where V_p^- is the set of the activities for which the p -th event is the final one, V_p^+ is the set of the activities for which the p -th event is the initial one.

From (17) and (18) we obtain

$$\sum_{i \in V_p^+} v_i = \sum_{i \in V_p^-} v_i, \quad p=2, \dots, m-1 \quad (19)$$

In other words (19) is the property 2 for the problem 1. Let S_0 be the minimal project cost under the project time T_0 , τ_i^0 and s_i^0 - the time and cost of the i -th activity correspondingly. Then the equivalent volume W_e ; minimal project time T , under the given resource level N the values τ_i and s_i can be found from the following formulas

$$\begin{aligned} W_e &= \left(\frac{S_0}{T_0} \right)^{1/2} T_0, \quad T = T_0 \left(\frac{S_0}{N T_0} \right)^{1/2} \\ v_i &= \frac{s_i^0}{\tau_i^0} \frac{N T_0}{S_0}, \quad \tau_i = \tau_i^0 \left(\frac{S_0}{N T_0} \right)^{1/2} \end{aligned} \quad (20)$$

So knowing the solution of problem 3 it is possible to find the optimal solution of problem 1. The relation between the minimal cost problem and minimal resources level problem can be extended to the case of the arbitrary functions $f_i(v_i)$ if the solution satisfying properties 1 and 2 is obtained. In this case the optimal activity times τ_i in problem 1 are the same as in problem 3 under the same project time T if

$$s_i(\tau_i) = \int_{\tau_i}^{\infty} f_i \left[\frac{W}{\tau} \right] d\tau \quad (21)$$

The proof follows immediately from the fact that property 2 is equivalent to (18) if (21) takes place.

Example. Let $s_i(\tau_i) = a_i - b_i \tau_i$,
 $d_i \leq \tau_i \leq \bar{d}_i$, $i=1, 2, \dots, n$.

From (21) we have

$$\tau_i(v_i) = \begin{cases} \bar{d}_i & \text{if } v_i < b_i \\ d_i & \text{if } v_i > b_i \end{cases}$$

and $\tau_i(v_i)$ has any value from d_i to \bar{d}_i if $v_i = b_i$. Consequently the minimal flow problem under the given project time T is obtained.

5. Solution of problem 2.

Consider problem 2 for the case of the power functions $f_i(v_i)$. Let W_e be "the equivalent project volume. Then $w(t) = N^{1/2}(t)$ is the project speed and the constrain (5) takes form

$$\int_0^t N(\tau) d\tau \leq S(t) \quad (22)$$

The problem is to find $N(\tau)$ satisfying (22) so that the project time is minimal.

The main point of the algorithm is to build the convex $\hat{S}_T(t)$, $t \in [0, T]$, so that $\hat{S}_T(T) = S(T)$ (Fig.4).

Step 1. Find T_0 using the equation

$$T_0 [S(T_0) T_0^{-1}]^{1/2} = W_e$$

If $S(T_0) T_0^{-1} t \leq S(t)$, $t \in [0, T_0]$ then T_0 is the minimal time. If there are intervals $\varepsilon \in [0, T_0]$ such that $S'(T_0) T_0^{-1} t > S'(t)$ for $t \in \varepsilon$ find minimal Δ

such that $S(T_0) T_0^{-1} (t - \Delta) \leq S(t)$ for $t \in [0, T_0]$. To provide this find $t(\Delta)$ minimizing $S(t) - S(T_0) T_0^{-1} (t - \Delta)$

The value Δ is determined by the equation $t'(\Delta) = 0$. When the value Δ is calculated,

find $T_1 = T_0 + \Delta$

Step 2. Construct $\hat{S}_{T_1}(t)$, $t \in [0, T_1]$

Calculate

$$W(T_1) = \int_0^{T_1} \left[\frac{d\hat{S}_{T_1}(t)}{dt} \right]^{1/2} dt$$

If $W(T_1) = W_e$ then T_1 is the minimal time.

If $W(T_1) > W_e$ then find

$$T_2 = \frac{1}{2} (T_0 + T_1)$$

and go to the 3-rd step.

Step i.. Construct $\hat{S}_{T_{i-1}}(t)$, $t \in [0, T_{i-1}]$

Calculate

$$W(T_{i-1}) = \int_0^{T_{i-1}} \left[\frac{d\hat{S}_{T_{i-1}}(t)}{dt} \right]^{1/2} dt$$

If $W(T_{i-1}) = W_e$ then T_{i-1} is the minimal time

If $W(T_{i-1}) > W_e$ then determine $T_i = T_{i-1} - \frac{1}{2} |T_i - T_{i-1}|$,

If $W(T_{i-1}) < W_e$ then determine $T_i = T_{i-1} + \frac{1}{2} |T_i - T_{i-1}|$

and go to the next step.

If $S'(t)$ is the convex function the problem can be solved more simply.

If $S'(t)$ is convex (upwards) then T_{min} is determined by the equation

$$\left[\frac{S(T_{min})}{T_{min}} \right]^{1/2} T_{min} = W_e$$

If $S'(t)$ is convex (downwards) then T_{min} is determined by the equation

$$\int_0^{T_{min}} \left[\frac{dS(t)}{dt} \right]^{1/2} dt = W_e$$

6. Analysis of the property 1.

The question of interest is what types of functions $f_i(v_i)$ give the optimal solution of problem 1, satisfying properties. The optimality conditions for the arbitrary convex functions $f_i(v_i)$ are following

$$\frac{1}{\Delta_i} \frac{d\varphi_i(w_{is})}{dw_{is}} = \sum_{j \in R_s} w_{js} \frac{d\varphi_j(w_{js})}{dw_{js}}$$

where $\varphi_i = f_i^{-1}$ and $w_{js} = x_{js}/\Delta_s$.

If $w_{is} = w_i$ i.e. w_{is} doesn't depend on S then

$$\sum_{i \in U_s} w_i \frac{d\varphi_i(w_i)}{dw_i} = \sum_{i \in U_s} \frac{d\varphi_i(w_i)}{dw_i} \quad (23)$$

Denote $F_i(v_i) = w_i \frac{d\varphi_i(w_i)}{dw_i} = f_i(v_i) \left[\frac{df_i(v_i)}{dv_i} \right]^{-1}$. As it follows from (23) $\{F_i(v_i)\}$ form the flow in the network.

Theorem. To provide that $\{F_i(v_i)\}$ is a flow for any flow $\{v_i\}$, the following necessary and sufficient Conditions must be satisfied

$$F_i(v_i) = a_i + \alpha v_i, \quad i = 1, 2, \dots, n,$$

where $\{a_i\}$ is a flow (It is supposed the network stays connected after input and output are removed).

Proof. Making variations of v_i along any path in the network we obtain

$$\frac{dF_i(v_i)}{dv_i} = \frac{dF_j(v_j)}{dv_j} \quad \text{for any } i, j \quad \text{and} \quad \frac{d^2 F_i(v_i)}{dv_i^2} = 0$$

for any i . It follows from the written above that

$$F_i(v_i) = a_i + \alpha v_i, \quad i = 1, 2, \dots, n,$$

where $\{a_i\}$ is a flow because $F_i = a_i$ when $v_i = 0$.

It follows from the theorem that $f_i(v_i) = C_i(a_i + \alpha v_i)^{1/\alpha}$, where $\alpha > 1$ because $f_i(v_i)$ must be a convex function when $a_i + \alpha v_i \geq 0$. Let W_e - the equivalent project volume and \tilde{x}_{is} - optimal values of x_{is} . Then if

$$\tilde{x}_{is}^{\alpha} \left(\sum_{i \in R_s} \tilde{x}_{is}^{\alpha} \right)^{-1} \geq a_i (\alpha N + A)^{-1}$$

where A is the value of the flow $\{a_i\}$ then the minimal

project time

$$T_{min} = w_e (\alpha N + A)^{-1/\alpha}$$

REFERENCES

1. В.Н.Бурков, А.Я. Лернер. Новые задачи теории сетевого планирования и управления. В сб. "Вопросы управления в больших системах". Изд. "Онтиприбор", 1968.
2. В.Н.Бурков. Применение теории оптимального управления к задачам распределения ресурсов. Труды III Всесоюзного совещания по автоматическому управлению /Одесса, 1965 г./ . т.Управление производством. м., Наука, 1967.
3. В.Н.Бурков. Распределение ресурсов как задача оптимального быстрогодействия. Автоматика и телемеханика, т. XXV, №7, 1966.
4. Б.С.Разумихин. Задачи об оптимальном распределении ресурсов. Автоматика и телемеханика, № I, 1967.
5. E.B.Berman. Resource allocation in a PERT network under continuous activity time-cost functions. Manag.Sci.,v.10, No. 4, 1964.

SOME QUESTIONS OF THE TESTING AND CONSTRUCTION PRINCIPLES
OF AN OPTIMUM MULTILEVEL CONTROL STRUCTURE IN SYSTEMS WITH
A SPECIFIC OBJECTIVE FUNCTION

M.K.Badunachvili, D.I.Golenko,
S.S.Naumov

Electronic Control Computers Institute, Moscow U S S R

Let us consider a class of controlling subsystems with some objective function and let us assume that there exists an integral measure of the target vicinity in such subsystems and this integral measure is some nondecreasing time function $V(t)$. It is evident that for each of the subsystems under discussion there exists some value of the integral measure corresponding to the target attainment V_{sch} . Depending on a variety of factors this measure can be obtained at different time instants. For the given subsystem an earliest early target attainment time t_e is assumed to exist, objectively, also (there exists the maximum movement speed to the target) which corresponds to the maximum utilization of the subsystem internal resources and the minimum external obstacles; generally speaking, this time t_p differs from the optimistic time t_o used in PERT Systems. At the same time there exists a deadline for the target attainment t_{cr}^x . Of the target is reached later than t_{kp} , the subsystem pays a fine. Hence, with the above constraints the subsystem must reach the target within the time interval (t_e, t_{cr}) .

It is necessary to note that during the movement to the given target the subsystem tries to minimize the efforts directed to the attainment of this target because repeated and lengthy operation in critical conditions (corresponding to the maximum speed of the movement to the target) can an early wear the subsystem too early. Therefore in this subsystem the scheduled time of the target attainment $t_p \leq t_{sch} < t_{kp}$ is chosen such that it is

x) Systems for which $t_{cr} = \infty$ obviously do not belong to the class of systems discussed.

ahead of t_{cr} by some time interval and provides certain confidence boundaries simultaneously both for the reliability of the target attainment not later than t_{cr} and for the minimization of the efforts spent in this movement.

Let us assume the curve $V_{sch}(t)$ /fig.1/ is the scheduled curve of the movement to the target and the integral measure corresponding to the target attainment is

$$V_{sch} = \int_{t_i}^{t_{sch}} V_{sch}(t) dt, \quad / 1 /$$

where t_i - the initial instant of the movement to the aim. Then in the process of the system movement, for the testing purposes it is necessary to compare the true measure

$$V(t_i) = \int_{t_i}^{t_j} V(t) dt, \quad / 2 /$$

which is a random value with the value of the target vicinity $V_{sch}(t_i)$ at certain time instants t_i , which must be defined by the choice of the quantization unit $V(t)$ is the true curve of the movement to the target in the above formula.

The control function of this system is elimination of the divergence between the actual and the scheduled measures of the target vicinity appearing due to the effect disturbances with the view to ensuring the target attainment by the subsystem by the time t_{sch} .

Before we discuss the algorithm let us shift the curve $V_p(t)$, corresponding to the movement to the target at the time t_e , in the direction of the X axis so that the end of the curve $V_p(t)$ would coincide with the point (V_{sch}, t_{sch}) /fig.1/. In this case we shall receive the point t_I as the intersection of the X axis with $V_e(t)$. It is easy to see that if the subsystem does not move to the target $V(t_I) = 0$ before the time t_I , obviously there is a non zero probability that we shall reach the target by the time moment t_{sch} starting from the time instant t_I , and using the possibilities of the subsystem to the utmost. We can consider t_I as the critical permissible time of the first interrogation. If the interrogation will be carried out after the time instant t_I then the deadline t_{sch} cannot be met because the highest velocity to the target is defined by the curve $V_e(t)$.

Thus, to meet the deadline t_{sch} , the first interrogation of the subsystem must be performed within the period specified by the interval $t_i < t \leq t_I$. Assuming the critical value $t = t_I$ and performing the interrogation, the subsystem receives the information on the target vicinity by comparing $V(t_I)$ and $V_{sch}(t_I)$. With this information effective controlling actions aimed at suppression of the divergence must be worked out. Drawing now through the point $[t_I, V(t_I)]$ a straight line, parallel to the X-axis, to the intersection with the shifted curve $V_e(t)$, we shall receive a point with the abscissa (t_2) which, as indicated above, determines the critical value of the second interrogatory period.

The next times of the interrogation are found in a similar way.

We have noted above that the actual curve the movement to the target $V(t)$ is a random curve. Thus, depending on the nature of its change, the same volume of V_{sch} can be reached at a certain actual time t_{act} . Here we have three cases:

1. The case $t_{sch} > t_{act}$ can occur a) when the subsystem is uncontrollable; it is evident that for this subsystem the interrogation is generally useless; b) when, in spite of the fact that the subsystem is controlled, the condition t_{sch} is obtained. In this case, at some interrogation step we shall be on the shifted curve $V_e(t)$. As the possibility of the movement strictly along the curve $V_e(t)$ is low, the unit controlled in this case is to obtain the least delay as far as the time t_{sch} is concerned,

2. In the case where $t_i < t_{act} < t_{sch}$ a finite number of the interrogatory steps for the control of the movement to the target are evidently required.

3. In the case $t_{act} = t_{sch}$ the sequence of all points of the interrogation has the limit of convergence - t_{sch} and the process of approximation occurs in an number of steps. But since in practice the volume V_{sch} is to be reached by the time t_{sch} with some specified accuracy $t_{sch} \pm \Delta t_{sch}$, in this case what is needed is to get into the region $[t_{sch} - \Delta t_{sch}, t_{sch} + \Delta t_{sch}]$

that, in contrast to the convergence to t_{sch} can be done in a finite number of steps. The condition $\Delta t < \tau$ must evidently be met.

An analytical record of the geometrical interpretation of the above algorithm leads us to the following relation which describes the critical value of the $(i+1)$ -th interrogation depending on the information on the measure of the target vicinity at the i -th control step of the subsystem.

$$t_{j+1} = t_1 + \frac{V(t_j)}{V_{sch}} (t_{sch} - t_1) \quad / 3 /$$

here t_1 is the time of the first interrogation V_{sch} and $V(t_i)$ are respectively of the target vicinity expressed by formulas (1) and (2) the planned and the actual integral measures.

Let us suppose now that $V_{sch}(t)$ and $V_e(t)$ are given in the form of straight lines and the actual curve of the movement to the target coincides with $V_{sch}(t)$. In this case $t_{act} = t_{sch}$ and the interrogation of the subsystem is needed only to prevent possible troubles

By using the evident relation
$$\frac{V(t_j)}{V_{sch}} = \frac{t_j}{t_{sch}} >$$

formula / 3 / can be written as

$$t_{j+1} = t_1 + \left(1 - \frac{t_1}{t_{sch}}\right) t_j$$

Using this expression the easy to receive formula for the value of the $(i+1)$ step is $t_{j+1} - t_j = \left(1 - \frac{t_1}{t_{sch}}\right) (t_j - t_{j-1})$.

This relation may be recorded also as

$$t_{j+1} - t_j = \left(1 - \frac{t_1}{t_{sch}}\right) t_1,$$

and as

$$t_{sch} = \sum_{j=1}^{\infty} (t_{j+1} - t_j),$$

it is easy to see that

$$t_{sch} = t_1 \sum_{j=0}^{\infty} \left(1 - \frac{t_1}{t_{sch}}\right)^j + t_1 \sum_{j=n+1}^{\infty} \left(1 - \frac{t_1}{t_{sch}}\right)^j \quad / 4 /$$

Here

$$R_n = t_1 \sum_{i=n+1}^{\infty} \left(1 - \frac{t_1}{t_{sch}}\right)^i,$$

is a remanence

term for which the following constraint will be reasonable

$$R_n \leq \Delta t_{sch}.$$

Now since the terms under the sums sign in the right-hand part of ed / 4 / are the sums of the decreasing geometrical progression, we can write that

$$t_{sch} = t_1 + t_{sch} \left[\left(1 - \frac{t_1}{t_{sch}}\right) - \left(1 - \frac{t_1}{t_{sch}}\right)^{n+1} \right] + \Delta t_{sch}.$$

After simple transformations we can obtain an expression for number of the preventive interrogations required to test the subsystem with the parameters t_{sch} , t_e , t_{sch}

$$n = \frac{\ln \frac{\Delta t_{sch}}{t_e}}{\ln \frac{t_e}{t_{sch}}} \quad / 5 /$$

Let us consider now a complex system consisting of subsystems of the above form. Thus, we consider two fixed levels of some multilevel structure. With this we carry out the numbering of the ranks in the following way: the number K is given to the rank of the subsystems and the number $K-1$ to the rank of the system.

As the functioning of every subsystem is directed to of a definite target, this system evidently will generally possess with some network of the targets which are parts of the overall target. The system has the integral curves of the aim nearness $V_{sch}^{(c)}(t)$, $V_e^{(c)}(t)$ and $V^{(c)}(t)$

showing the different conditions of the system movement to the overall target as well as the appropriate parameters.

The existence of the target network also means that there is some critical path with a certain number of targets on it and this means the existence of a certain number m of the critical subsystems M . Let us choose the i subsystem from all these subsystems for which the number of the preventive interrogations given by formula (5) will be minimal, i.e.

$$\min_{1 \leq j \leq m} N_j = \min_{1 \leq j \leq m} \frac{\ln \frac{\Delta t_{schj}}{t_{ej}}}{\ln \frac{t_{ej}}{t_{schj}}}$$

Here we have the obvious inequality

$$\Delta t_{schj} < t_{ej} < t_{schj} \quad / 6 /$$

If every subsystem the critical path is interrogated a specified number of times, the appropriate minimal number of the interrogation required for the interrogation of all critical path will be equal to

$$\min_{1 \leq j \leq m} N_K = m \cdot \min_{1 \leq j \leq m} N_j \quad / 7 /$$

where N_K is the number of the interrogations needed for the critical path with the given detailing Δt_{schj}

Let us consider now the system with the parameters t_{sch} , t_e and Δt_{sch} of its movement to the overall target

$$\begin{cases} T_{sch} = \gamma_1 \cdot t_{sch}, \\ T_e = \gamma_2 \cdot t_{ej}, \\ \Delta T_{sch} = \gamma_3 \cdot \Delta t_{schj}, \end{cases}$$

where $\gamma_1, \gamma_2, \gamma_3$ are arbitrary nonnegative number above 1 and let us assume that

$$\gamma_1 > \gamma_2 > \gamma_3. \quad / 9 /$$

Let us quite naturally require that in the discussed multilevel structure of the number of interrogations needed for control decrease with the increase of the rank (with the decrease of the rank number) of the hierarchy. With this we shall give more rigid form to our demand

$$N_{K-1} \leq \min_{1 \leq j \leq m} N_K. \quad / 10 /$$

Where N_{K-1} is the number of interrogations needed for the system interrogations in terms of detailing the $(K-1)$ th rank, i.e. by parameters (8). In fact, taking into consideration expressions (6), (7), (8) and (9) we can write the following inequalities

$$\begin{aligned} N_{K-1} &= \frac{\ln \frac{\gamma_3 \Delta t_{schj}}{\gamma_2 t_{ej}}}{\ln \frac{\gamma_2 t_{ej}}{\gamma_1 t_{schj}}} \leq \frac{\ln \frac{\gamma_3}{\gamma_2}}{\ln \frac{\gamma_2}{\gamma_1}} + \frac{\ln \frac{\Delta t_{schj}}{t_{ej}}}{\ln \frac{t_{ej}}{t_{schj}}} = \frac{\ln \frac{\gamma_3}{\gamma_2}}{\ln \frac{\gamma_2}{\gamma_1}} + \min_{1 \leq j \leq m} N_j \leq \\ &\leq \frac{\ln \frac{\gamma_3}{\gamma_2}}{\ln \frac{\gamma_2}{\gamma_1}} + m \cdot \min_{1 \leq j \leq m} N_j = \frac{\ln \frac{\gamma_3}{\gamma_2}}{\ln \frac{\gamma_2}{\gamma_1}} + (m-1) \min_{1 \leq j \leq m} N_j + \min_{1 \leq j \leq m} N_j. \end{aligned}$$

Hence it is easy to see that for the realization of the condition

$$N_{K-1} \leq \min_{1 \leq j \leq m} N_j$$

automatically meeting condition (10) it is necessary that

$$\frac{\ln \frac{\gamma_3}{\gamma_2}}{\ln \frac{\gamma_2}{\gamma_1}} + (m-1) \min_{1 \leq j \leq m} N_j \leq 0. \quad (II)$$

Now, taking into consideration the evident inequality

$$\gamma_1 t_{schj} \gg \gamma_2 t_{ej}, \quad (I2)$$

condition (II) may be written as

$$\frac{\gamma_3}{\gamma_2} \gg \left(\frac{\Delta t_{schj}}{t_{ej}} \right)^{m-1} \quad (I3)$$

Let us prove that condition (I3) is also a sufficient condition.

In fact, assuming that $\frac{\gamma_3}{\gamma_2} < \left(\frac{\Delta t_{schj}}{t_{ej}}\right)^{m-1}$,

then at $m \rightarrow \infty$ $\frac{\gamma_3}{\gamma_2} < \lim_{m \rightarrow \infty} \left(\frac{\Delta t_{schj}}{t_{ej}}\right)^{m-1} = 0$.

This means that either γ_3 or γ_2 is negative that contradicts the statement of the problem. This contradiction proves the assumption. However, as in fact there exists only a certain finite value of the number of subsystems along the critical path $m = m_{ch}$ evidently there exists an interval for the positive values $\frac{\gamma_3}{\gamma_2}$

$0 < \frac{\gamma_3}{\gamma_2} \leq \left(\frac{\Delta t_{schj}}{t_{ej}}\right)^{m_{ch}-1}$
within which it is not clear whether the sufficiency for condition (I3).

The rate of convergence for the value $y = \left(\frac{\Delta t_{schj}}{t_{ej}}\right)^{m_{ch}-1}$ to the zero when m increases may be characterized by the length of the interval $J = \int_1^{\infty} \left(\frac{\Delta t_{schj}}{t_{ej}}\right)^{m-1} dm = \frac{1}{\ln \frac{t_{ej}}{\Delta t_{schj}}}$. Hence follows that the value J changes in proportion to the detailing index Δt_{schj} . The higher the number of secondary targets, i.e. the Δt_{schj} the higher the rate of convergence for to the zero the indicated function and thus the smaller the zone where the property of the sufficiency is doubted. Obviously, by selecting the parameters Δt_{schj} , t_{ej} and m this zone may be done smaller than any number given in advance.

Now let us solve inequality (II) for $\frac{\gamma_2}{\gamma_1}$ and with analogous reasonings; it is easy to see that

$$\frac{\gamma_2}{\gamma_1} \leq \left(\frac{t_{schj}}{t_{ej}}\right)^{\frac{1}{m-1}}$$

with relation (I2) as well as inequalities (II) and $\gamma_2 t_{ej} > \gamma_3 \Delta t_{schj}$ the entire system of the constraints that are put in this case on the relation $\frac{\gamma_2}{\gamma_1}$ and $\frac{\gamma_2}{\gamma_3}$ will have this form

$$\left\{ \left(\frac{\Delta t_{schj}}{t_{ej}}\right)^{m-1} \leq \frac{\gamma_3}{\gamma_2} \leq \frac{t_{ej}}{\Delta t_{schj}}, \quad \frac{t_{ej}}{t_{schj}} \leq \frac{\gamma_2}{\gamma_1} \leq \left(\frac{t_{schj}}{t_{ej}}\right)^{\frac{1}{m-1}} \right\}$$

Taking into consideration relation (8), the necessary condition for the optimality of the given multilevel structure may be finally given by.

$$\begin{cases} \left(\frac{\Delta t_{schj}}{t_{ej}}\right)^m \leq \frac{T_{sch}}{T_e} \leq 1 \\ \left(\frac{t_{ej}}{t_{schj}}\right)^2 \leq \frac{T_e}{T_{sch}} \leq \left(\frac{t_{schj}}{t_{ej}}\right)^{\frac{2-m}{m-1}} \end{cases}$$

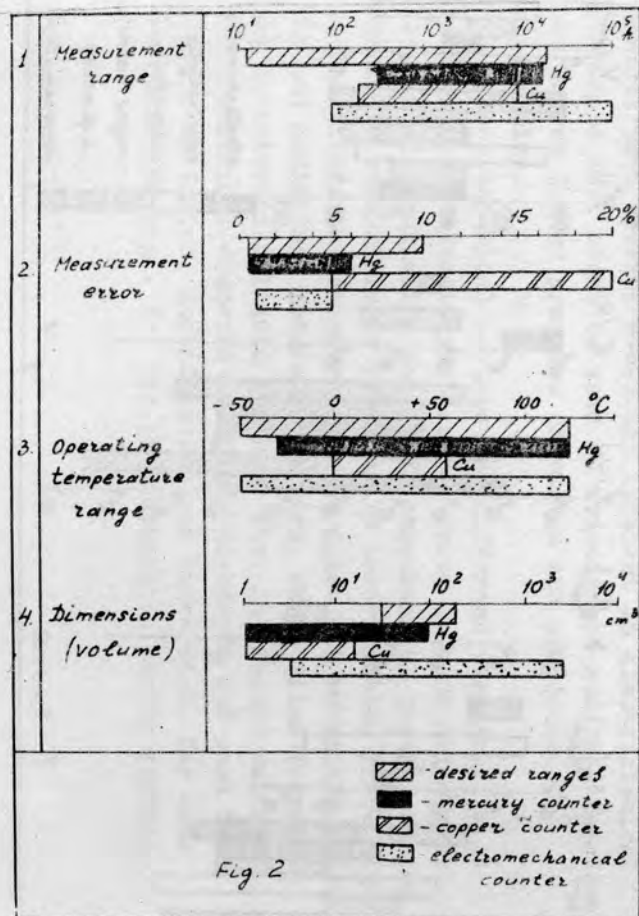
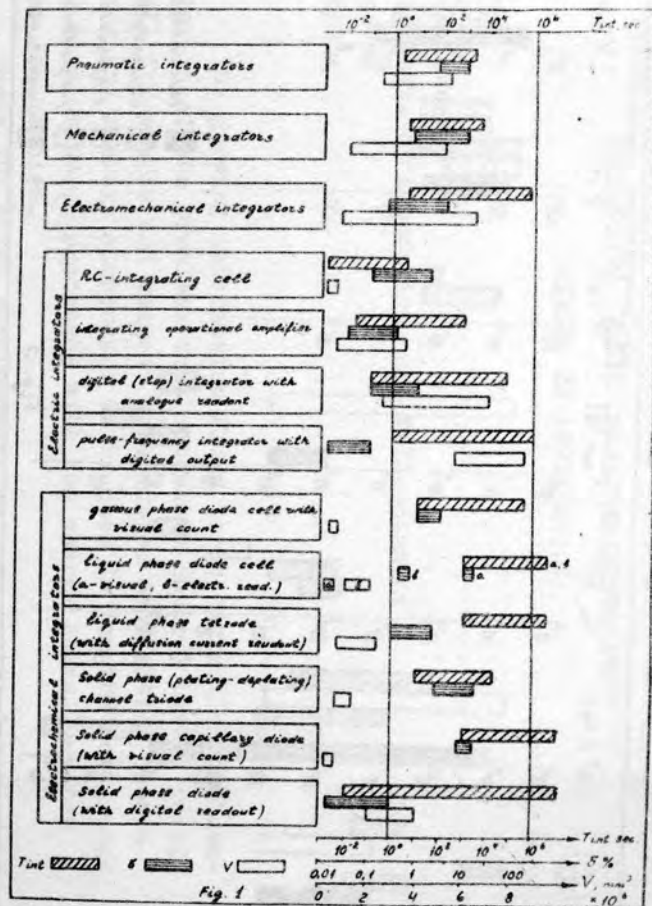
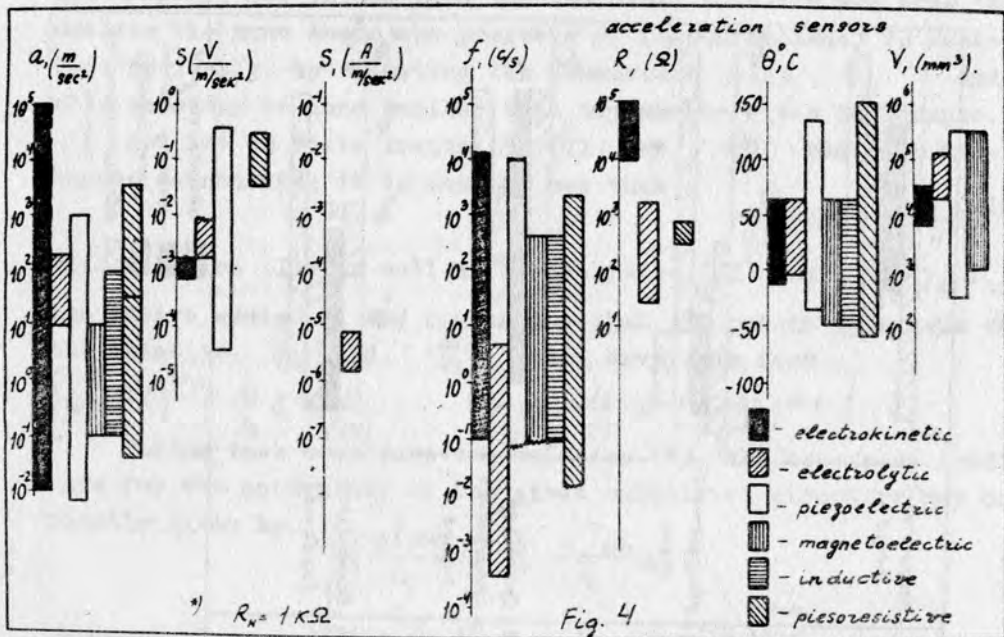
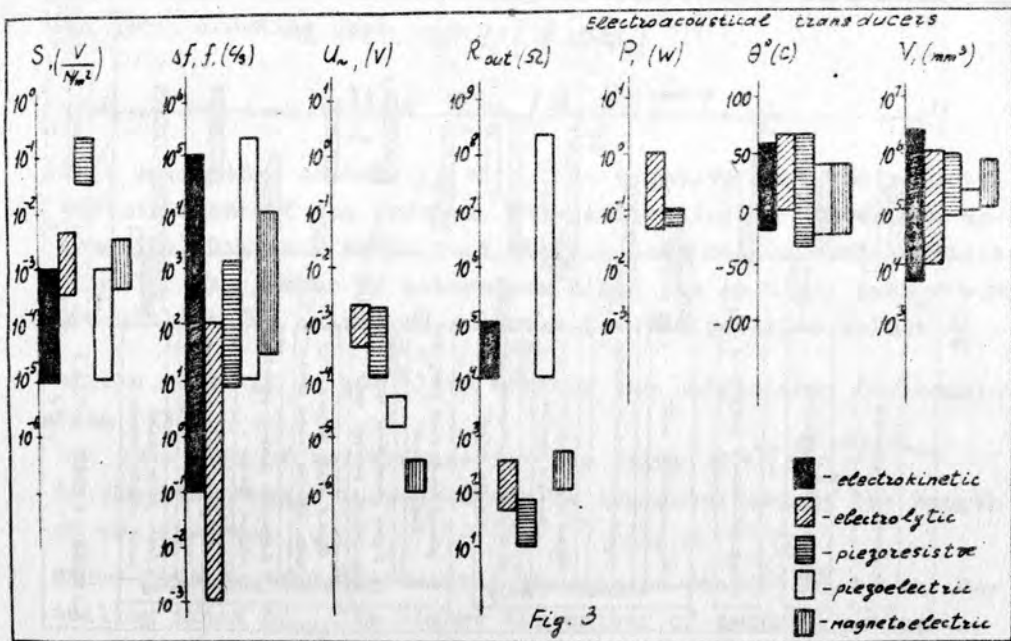


Fig. 2





ELECTROCHEMICAL TRANSDUCERS, COMPARATIVE
PROPERTIES, BASIC CHARACTERISTICS AND
FIELDS OF APPLICATION

A.P.Shorygin

Institute of Automation and Telemechanics

Moscow

USSR

The growing information flows to be controlled and complexity of the problems that automatic control has to tackle call for simple and compact information acquisition, processing and storage components to be used in lengthy continuous or discontinuous processes with low energy stimuli and/or low frequencies.

These and other requirements to components, devices and instruments made by new problems of automatic monitoring and control have aroused a certain interest in electrochemical phenomena that might be used in transducers and this despite the unprecedented progress of vacuum and semiconductor electronics.

Electrochemical devices use different reactions in liquid or solid electrolytes and/or polar liquids (electrolytic and electrokinetic cells). Electrochemical systems may be partially reversible or practically irreversible. There might be variations in concentrations of dissolved materials, plating and deplating, metals deposition, oxide films formation and deterioration, liquid transfer etc. At present various electrochemical elements are being developed and used for signals detection, amplification and conversion, for data processing, storage, for use in adaptation, remote control,

automata etc. Investigations are underway on electrochemical matrices; attempts are made to use electrochemical devices for simulation of biological processes; electrochemical systems and processes that could be used in transducers are investigated.

The existing electrochemical devices can be divided into two major groups. The first includes those that already are or can be used in industry due to their clear-cut advantages or specific features.

These units are:

- 1) integrating units with visual, electric or photoelectric read-out;
- 2) electric rectifying diodes for very low currents and low frequency;
- 3) sensors of variable or pulsed pressures and acoustic pressure gradients of low and infralow frequency;
- 4) conductive concentration (salination) meters, conductometric and galvanic gas analysers, composition sensors, etc.

Of substantial interest for a number of fields are first of all the electrolytic integrating devices.

$$P_{\text{output}} = \varphi \left[\int_0^{\tau} i_{\text{inp}}(t) dt \right] \quad (1)$$

where i_{inp} is the input current, τ - the duration of its flow. Electrochemical integrating devices vary widely in their characteristics and practical possibilities due to the type of redox system, electrochemical effects used and the read-put procedure (Fig. 2).

The output of "gas phase" integrating diodes is the displacement Δx of electrolyte in the capillary which connects the hydrogen-filled electrode chambers

$$\Delta x \approx \frac{RT}{K_2 p S_K} \int_0^{\tau} i_{\text{inp}}(t) dt \quad (2)$$

where p is the hydrogen pressure, S_K the crosssection area of the capillary, K_2 - the instrument constant.

The redox reaction in liquid phase cells (solion) follows the equation



where A_1 is the oxidized and A_2 the reduced form of the solution components. For typical conditions where the rate of the process practically depends only on diffusion towards a cathode of the form A_1 , the output e.m.f. of an integrating diode is found by the equation

$$C_{01} n F S \ell \cdot \text{th} \frac{E_c}{E_T} = \int_0^T i_{\text{inp}}(t) dt \quad (4)$$

where $E_T = RT/nF$, F the Faraday number, S and ℓ the cross-section area and the length of integral chamber separated from the remainder of the cell by a semi-permeable membrane,

C_{01} the initial concentration of the form A_1 (non-basic carrier). With $E_c \ll E_T$ the relation becomes linear and the sensitivity of the integrator diode will be

$$S_E = \frac{E_c}{q_{\text{inp}}} = \frac{2RT}{C_{01} (nF)^2 S \ell} \quad (5)$$

The maximal value of the input charge will depend on the capacity of integral chamber and the initial concentration of form A_1 :

$$(q_{\text{inp}})_{\text{max}} = n F C_{01} V \quad (6)$$

To ensure the required integration accuracy the upper bound of the input signal frequency spectrum should meet the condition

$$f_2 \ll \frac{D}{2\pi \ell^2} \quad (7)$$

where D is the diffusion factor.

The storage time for the integral value depends on the self-discharge time constant

$$\tau_{s.d} = \frac{V \ell_K}{D S_K} \quad (8)$$

Here S_K is the total equivalent cross-section area of the capillaries in the semi-permeable membrane,

l_k is the equivalent thickness of the membrane.

The output of liquid state triodes and tetrodes is the diffusion current between the anode of the integrating device and the read-out electrode inserted in the anode chamber which receives the negative bias. It is described as

$$(I_o)_{outp} = \frac{2q_{inp} D}{l_{zA}^2} \cdot \frac{e^{U_z/E_T} - e^{U_A/E_T}}{e^{U_z/E_T} + e^{U_A/E_T}} \quad (9)$$

where U_z and U_A are the electrode-electrolyte potential for the read-out electrode and the anode, while l_{zA} is the distance between these electrodes. Since the current in the output circuit is limited by the load resistance R_H , if the input charge reaches the value

$$q_{inp} \geq \frac{E}{R_H} \cdot \frac{l_{zA}^2}{2D} \quad (10)$$

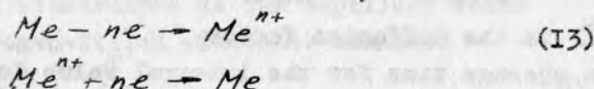
saturation occurs and the saturation diffusion current will be

$$I_{sd} = \frac{nFD C_{A1} S_{zA}}{l_{zA}} \quad (11)$$

Under these conditions which are used in practical integrators

$$I_{out} = S_I \int_0^{\tau} i_{inp}(t) dt \quad (12)$$

($S_I = \frac{k_T n D}{l_{zA}}$ and k_T is the solion constant). The typical triodes and tetrodes characteristics are shown in Table 1. Solid phase plate-deplate integrating devices, (Faraday coulombmeters) are now represented by a great variety of types (Table I). The redox reaction follows the scheme



The mass of the substance plated and/or deplated during the interval τ (e.g. during one cycle of action of integrating diode with digital read-out)

$$m_{Me} = A\eta/nF \cdot \int_0^{\tau} i_{inp}(t) dt \quad (14)$$

where A is the atomic weight, η the input factor ($\eta \approx 1$). In the case of a capillary cell the output is the displacement of the electrode-solution boundary

$$\Delta x_{out} = S_e \int_0^t i_{inp}(t) dt \quad (15)$$

where $S_e = A\eta/nF\rho S_k$ the sensitivity (ρ is the density of the material of the electrode column, S_k - its cross-section area). When such a device is employed as the elapsed time indicator

$$t \approx \frac{\Delta x_{out}}{S_e \frac{U}{R_0} \frac{1}{1+R_c/R_0}} \quad (16)$$

where U - is the voltage, R_0 is the limiting resistor, R_c - the cell resistance.

Of special interest are solid phase channel triodes (memistors) with monotonous changes of linear output resistance R_{out} against an electric charge that has flown in the input circuit [2,3]. At sufficiently low R_{out} and uniform plating of the metal sheet on the storage electrode

$$\sigma_{out} = \frac{R_{out}(t)}{R_{out}(0)} = \frac{1}{1 + R_{out}(0) / \psi \int_0^t i_{inp}(t) dt} \quad (17)$$

where $R_{out}(0) = (R_{out})_{max}$ is the initial value, $R_{out}(t)$ the same at the instant t and $\psi = nF/\rho A \cdot \rho \gamma \ell_s^2$ (γ is the mean density of the deposited metal, ρ its specific resistance, ℓ_s the length of storage electrode). At sufficiently large q_{inp} value the resistance variation velocity [3] is

$$\frac{d\sigma_{out}}{dt} = \frac{\psi}{R_{out}(0) \cdot I_{inp} \cdot t^2} \quad (18)$$

and the shortest time of complete resistance change is

$$t_{min} = \frac{\psi [R_{out}(0) / ((R_{out})_{min} - 1)]}{R_{out}(0) \cdot (I_{inp})_{max}} \quad (19)$$

($(I_{inp})_{max}$ is the largest permissible input current).

Under typical conditions the resistance of the reading current for a/c will be [4]:

$$Z_2 \approx 2 \frac{\sqrt{R_{s1} A_{F1}}}{f^{0.25}} \cdot \tanh\left(\frac{\ell_s}{2} \sqrt{\frac{R_{s1}}{A_{F1}}} \cdot f^{0.25}\right), \quad (20)$$

where ℓ_s is the length of the storage electrode (resistive

electrode), R_s , the resistance per unit length, A_{F_i} is the factor which describes the electrode-electrolyte Faraday impedance per unit length of the same electrode.

By their integration time range, magnitudes of currents, accuracy and dimensions in a number of cases the electrolytic integrating devices are advantageous (see Fig. 1). They have linear characteristics, very low consumption (tens or hundreds of microwatts), they can integrate fractions of microcoulombs, operate with gating photocalls, thermocouples,

Hall sensors etc. Their dimensions do not exceed those of miniature electronic tubes, while their weight is 2 to 20 gr. The highest accuracy is achieved in solid phase electrochemical integrating devices with digital read-out.

The processes in electrochemical elements develop in very thin layers measured in fractions of a micron; more miniature ("planar") elements are expected.

The near future will see a wide application of elapsed time indicators used on the above principles (since it is necessary to monitor the operability and longevity of instruments, devices and machinery as stipulated by state standards) (Fig. 2)[6], as well as integrating units especially with large integration time, the adjusting units of adaptive systems and optimisers etc. Of substantial interest are also electrochemical devices that process automatically the measurements data and realize an algorithm of the form

$$y_i = \frac{K_i \int_{t_i}^{t_{i+1}} [x(t) - x(t_0)] dt}{\sum_1^n K_j \int_{t_j}^{t_{j+1}} [x(t) - x(t_0)] dt} \quad (21)$$

In liquid phase electrochemical rectifying diodes a redox reaction represented by eq.(2). If the rate of the reaction depends practically on diffusion alone, the integral equation for current through a diode will be [5]:

$$\begin{aligned} \left[1 - \frac{1}{U_1} \cdot \frac{d}{dt} \int_0^t f_1(t-\tau) \cdot i(\tau) d\tau \right] \exp n F U_{inp} / RT = (22) \\ = \left[1 - \frac{1}{U_2} \cdot \frac{d}{dt} \int_0^t f_2(t-\tau) \cdot i(\tau) d\tau \right]^3, \end{aligned}$$

where

$$U_{1,2} = n F S C_{1,2}$$

$$f_{1,2} = \frac{a}{D_{1,2}} \left\{ 1 - \exp \frac{t D_{1,2}}{a^2} \operatorname{erfc} \sqrt{\frac{t D_{1,2}}{a^2}} \right\} \quad (23)$$

while a - is the radius of a spherical microelectrode; indices 1 and 2 denote the quantities that relate to forms A_1 and A_2 respectively.

$$\text{At } t \rightarrow \infty \quad \lim i(t) = i_\infty$$

$$\lim \frac{1}{U_{1,2}} \cdot \frac{d}{dt} \int_0^t f_{1,2}(t-\tau) \cdot i(\tau) d\tau = \frac{a \cdot i_\infty}{n F C_{1,2} D_{1,2}} \quad (24)$$

so that the static current-voltage characteristics of a diode will be

$$E_\infty = E_T \left[3 \ln \left(1 - \frac{a i_\infty}{n F S C_{ox} D_2} \right) - \ln \left(1 - \frac{a i_\infty}{n F S C_{ox} D_1} \right) \right] \quad (25)$$

Electrolytic rectifying diodes handle very low currents, from hundredths or tenths of a microampere to hundreds microamperes which gives an advantage (at low frequencies) over diodes using other techniques. Electrolytic diodes of this type underlie elements which realize operations $\frac{1}{\sqrt{p}} x$ and $\sqrt{p} x$ [5] in a wide range of infralow frequencies as well as negative resistors.

Electrolytic and electrokinetic non-resonance sensors of low variable and discontinuous pressures, sensors of pressure gradients (in particular electroacoustic transducers of infralow and acoustic frequency bands that can operate at very high static pressures) and acceleration sensors have, as follows from Fig. 3, 4, 5, advantages in terms of possible frequency bands and coverage of low frequencies, although they are not as sensitive as piezoresistive transducers. The input/output characteristics of electrokinetic transducers are linear in a wide range of amplitudes. For free ran

$$E_{out} = \frac{K_9 \xi E \Delta p_{inp}}{4 \pi \mu \lambda_5} \quad (26)$$

where ξ is electrokinetic potential, $E, \mu, \lambda_5 \dots$ are the dielectric factor, dynamic viscosity and the resulting electroconductivity of the polar liquid, $K_9 \dots$ is the

factor which incorporates the hydrodynamic properties of the semi-permeable membrane. The amplitude-and-frequency response of that transducers is

$$E_m = \frac{\xi \varepsilon}{4\pi\mu\lambda_0} \cdot \frac{\omega \tau_m \rho_m}{\sqrt{1 + (\omega \tau_m)^2}} \quad (27)$$

where τ_m - is the mechanical time constant of the transducer.

The input/output characteristic of liquid phase solion electrolytic transducers (which operate by mode (2)) for one of the configurations of the tiny cathode orifice

$$I_{out} = \frac{2FC_0\delta^3\ell}{3\mu a\pi^3} \cdot \Delta\rho_{inp} \left[1 - \exp\left(-\frac{12\pi^2 D\mu a^2}{\delta^4} \cdot \frac{1}{\Delta\rho_{inp}}\right) \right] \quad (28)$$

where C_0 is the concentration of non-basic carriers in free volume, D - is the diffusion factor, δ , a , ℓ are geometric dimensions of the slot orifice, μ - is the dynamic viscosity of the solution. By changing the shape of the cathode orifice one can obtain various amplitude responses (linear, logarithmic, etc.).

Versions of these elements are employed as seismic sensors, sensors of biological parameters, acoustic receivers of infralow frequency, etc. The electrokinetic devices have highly linear amplitude responses and the widest dynamic range, while electrolytic ones have very low frequency band. However, at present both are inferior to new piezoelectric and piezoresistive sensors in terms of the operational temperature range.

Various types of electrochemical (a/c, contact and contactless, galvanic) concentration meters for solutions and gases are now manufactured on a large scale and used in laboratories and industry.

Another group is made by electrochemical devices that need further study and improvement. We can expect that scientific research will reveal new fields of application where such transducers will be advantageous. This group may incorporate, e.g.:

- 1) electrolytic sensors of vibrations;
- 2) electrokinetic sensors of vibrations;
- 3) low voltages indicators;

- 4) infralow frequency liquid phase electronic amplifiers;
- 5) electrolytic and electrokinetic data processing transducers;
- 6) solid phase electrolytic static switch and power amplifiers etc.

In the writers opinion this field of control needs research along these main lines:

- a) new basic principles that would expand technological prospects of electrochemical devices;
- b) characteristics of sensors, transducers and other units and most practical field of their application;
- c) theory and calculation procedures for electrochemical devices of various types;
- d) circuits with electrochemical elements to be incorporated in automatic control and monitoring;
- e) most advanced technology of electrochemical elements manufacturing with a view to securing their high reliability.

Table I

Types of integrating elements	Effect used	Output physical quantity	Read-out
<u>Redox systems with inert electrodes</u>			
1. Gas phase	a) hydrogen evolution absorption in electrochemical reactions;	a) change of volume (or pressure) difference over passive elec- trodes of opposite polarities;	Visual
		b) an electrolytic contact swit- ching;	Electric
	b) change of properties of an electrode due to absorption of hydrogen formed in elec- trochemical reactions;	a) change in electrode electric conductivity;	Electric
		b) change in concentration e.m.f.	
2. Liquid phase (solution)	Change of the components concen- tration distribution in the elec- trolyte cell.	a) change in concentration e.m.f.	Electric
		b) change in electrolyte optical density	Visual or electric
		c) change of diffusion saturation current	Electric

1	2	3	4
<u>Redox systems with electrochemically soluble electrodes</u>			
3. Solid phase	a) change of electrode mass due to electrolytic plating or deplating;	a) change in electrode electric conductivity;	Electric
		b) change in optical density of a thin film electrode;	Photoelectric
		c) change in electrode linear dimensions;	
		d) change in electrode weight;	Electric
		e) change in electric conductivity of a cell (due to change in electrode geometric dimensions)	Electric
	b) change in equilibrium potential when the anodeactive metal dissolves over an inert support	Change in voltage drop on electrochemical cell	Electric

Conclusion

Comparison of characteristics leads one to believe that electrochemical devices based on different redox reactions can be used in sensing and data processing, in control of processes characterized with very low frequencies below tenths or hundredths of a cycle/sec, very low currents (to fractions of μA), very long time intervals (up to tens or hundreds of days and more). Very small sizes and low power consumption can be secured.

Such principles can underlie adaptive elements where microsecond pulses are used for non-destructive read-out.

References

1. А.П. Шорыгин Электрохимические элементы /основные особенности; классификация/ ЭИКА вып. 8, изд "Энергия", 1967.
2. A.P.Shorygin Automation and Remote Control, No 8, 1966.
3. В.С. Боровков, В.В.Трейер Электрохимические аналоговые запоминающие устройства /редактор А.П.Шорыгин/, изд.ВЗЭИ, М., 1967.
4. А.П.Шорыгин Исследование входных, выходных и передаточных характеристик твердофазных канальных электролитических триодов. Электрохимия, т.IV
5. Р.Ш.Нигметуллин Теория электрохимического диода, ДАН т. 150, 600, 1963, №3
6. А.П.Шорыгин, Э.В.Казарян Счётчики машинного времени, Труды АН Арм.ССР 1968 Вып. 4.
7. К.Ф.Гуссейнов, М.С.Касимзаде Переходные процессы в электрокинетических сейсмоприемниках, Автометрия, стр.19, 1966, №4

RATIONAL ALGORITHM OF CONTROLLING THE THERMAL CONDITION OF
BLAST FURNACE USING COMPUTERS

E.L.Suchanov, V.S.Shvidki, B.I.Kiteev, Ju.G.Yaroschenko,
Ju.N.Ovchinnikov, V.G.Lisienko

(The Urals Polytechnical Institute, Sverdlovsk, USSR.)

The blast furnace is a complicated multidimensional controlled plant having distributed parameters the estimation of which is a very difficult problem. The most difficult thing here is to find out the criteria and ways of estimation of the blast furnace heat regime which is the main controlling factor for reduction processes of the blast furnace operation. The very notion of heat regime or thermal condition of the blast furnace has no uniform interpretation and needs more accurate definition.

The modern conception of the fundamental laws governing the heat exchange in the blast furnace run at combined blast is presented in recently published works^{1,2}. These two and earlier studies³ have brought to light S-shaped nature of changing temperatures over the height of the blast furnace (fig.1), the temperature of gas and burden at its middle levels varying only slightly and approximating to each other. Such an experimentally determined dependence leads to conclusion that heat exchange processes are concentrated in the upper and lower steps of intensive heat exchange. There is an intermediate region between these steps with mean temperature being practically constant under particular conditions of melting.

Under a certain stability of this temperature with respect to time and in a certain height range the intermediate region is a kind of a damper that removes the direct influence of heat exchange steps on each other. This accounts for revealed by us property of non-interaction of the thermal operation in the upper and lower sections of the blast furnace.

The property of non-interaction becomes apparent in different and sometimes opposite response to the same controlling factors displayed by the upper and lower sections of the furnace (we mean here both statics and dynamics of transient heat processes).

Relative independence and great difference in the nature and character of heat exchange in various furnace zones is reflected in the mathematical description of these phenomena¹.

As a result of the abovementioned the authors came to a conclusion that because of the property of non-interaction it would be wrong in the main to evaluate the thermal condition of the furnace as a whole. Also it would be wrong to judge of the temperature regime of the furnace according to the thermal condition of the hearth without taking into account the temperature profile of the furnace shaft, as sometimes the case may be. It goes without saying, one can't do vice versa, too.

Therefore, when dealing with a blast furnace there must be only differentiated approach to the estimation of the thermal condition in the upper and lower sections of the furnace. In carrying out such independent inspection over the thermal regime of the furnace it would be helpful to divide the entire operating volume of the furnace into two independent zones - the upper and the lower, the relative boundary between them lying on quite definite isothermic surface corresponding to the mean temperature or to some other most stable with respect to time temperature of the intermediate region.

It was assumed that the relative boundary between the upper and lower furnace zones is temperature t_0 (fig.1) which characterizes the beginning of the intensive evolution of the endothermic reaction of carbon dioxide reduction in the ore ridge area. This temperature usually ranges from 850°C to 920°C depending upon the type of iron produced, the blast composition and some other quite

definite operating conditions of the furnace.

The mean temperature over the mass \bar{t}_w of the burden volume contained between the charging level and the assumed division line may serve as a general quantity estimation of the thermal condition of the upper zone of the furnace (in fig. 1, a the volume in question is shaded). The relative value of the temperature $i_s = \bar{t}_w / t_o$ was termed by the authors the index of the temperature profile of the upper section of the furnace. The revealed laws of heat exchange in the upper zone of the furnace¹ make it possible to determine index i_s from the current information about the temperature t_k and consumption of top gas G_k and also the consumption of burden materials charged G_w . Besides, it is necessary to have the information about the heat exchange coefficient α , specific heat of the burden and top gases, the average shaft cross-sectional area S and the average height of the furnace upper zone assumed for the calculation H . The calculation are made according to the following formula⁴:

$$i_s = 1 - \left[\frac{t_k - t_{wk}}{t_k - m(t_{wk} + \Delta t_o)} \cdot \frac{1 - \exp(-A)}{A} - \frac{m \cdot \Delta t_o}{t_k - m(t_{wk} + \Delta t_o)} \right], \quad (1)$$

where t_k and t_{wk} = temperature of gas and burden on the top of the furnace, °C;

$m = 0,5 \left(\frac{W_{wk}}{W_k} + 1 \right)$ = average burden-gas thermal capacity flow ratio for the upper zone;

W_{wk} - W_k = gas and burden flow thermal capacities on the top level, w/°C;

$\Delta t_o = t_o - t_{wo}$ = temperature difference between gas and burden on the division line between the zones, °C;

$A = \frac{\alpha H S}{W_w} (1 - m)$ = auxiliary coefficient.

Index i_s is a generalized parameter of the thermal condition of the upper blast furnace zone. Recently this parameter has become more essential for the heat regime

inspection and control of the blast furnace due to the wide use of oxygen for the intensification of the blast furnace production. With the oxygen enrichment of the blast the volume of hearth gases is decreased, which reduces warming up the burden, thereby slowing down the course of the main reducing reactions. Later on you will see that it is impossible to solve the problem of optimization of the blast furnace without using index i_g .

As a generalized parameter of the thermal condition of the lower blast furnace zone the authors suggest index i_H equal to the ratio of actual quantity of heat $Q_{\text{факт}}$ in the tapped products to the heat consumption Q_0 theoretically required for physical and chemical heating of the smelting products of the predetermined composition:

$$i_H = \frac{Q_{\text{факт}}}{Q_0} = \frac{Q_{\text{жж}} + Q_{\text{шл}} + Q_{\text{Si, Mn, P...}}}{Q_0}, \quad (2)$$

where $Q_{\text{жж}}$ = iron enthalpy allowing for the melting heat;
 $Q_{\text{шл}}$ = slag enthalpy without the slag formation heat;
 $Q_{\text{Si, Mn, P...}}$ = the heat consumed for reduction of silicon manganese, phosphorus and other elements.

Index i_H which is determined according to the amount, temperature and chemical analysis of the tapped products is quite a reliable criterion of the heat condition of the lower part of the furnace and may be taken as a reference input for the automatic control system. The efficiency of the automatic control system greatly increases if we consider not only the past but also the current and even predicted thermal condition of the lower part of the furnace. Thus to determine index i_H during the interval between the tappings of iron we may use the balance inspection method applying the following equation^I:

$$i_H = \frac{1}{Q_c} \left\{ \frac{V_d \cdot \delta}{(V_{O_2})_K - (V_{O_2})_D} \left[(c_{H_2} \cdot t_d - 1312) H_2 + (c_{O_2} \cdot t_d + 7840) O_2 + \right. \right. \\ \left. \left. + 1,244 \cdot 10^{-3} \cdot c_{H_2O} \cdot t_d \cdot \varphi_D - 11,53 (1 + 0,154 \cdot \eta_{H_2}) \cdot \varphi_D \right] + \right. \\ \left. + [c_A \cdot G_A + c_K \cdot K \cdot (1 - 0,01 \cdot \Pi)] \cdot t_0 - 6912 (1 - 0,411 \cdot \eta_{H_2}) \cdot \Gamma - \right. \\ \left. - 31750 \cdot Fe_o \cdot z_d - 104500 \frac{d}{P_{cyr}} - 24600 \right\}. \quad (3)$$

For the calculation according to the above equation it is necessary to have some current information about the blast and top gas parameters (altogether 14 variables) and besides some recurrent data about the changes in the burden composition and some reference data. The main coefficient τ_d characterizing the degree of direct reduction can be periodically defined more precisely by comparing the results of the calculation according to the formulae (2) and (3). Index i_H obtained from equation (3) may be considered a predicted value since the changes in the top gas composition to a certain degree leave behind the evolution of transient processes in the thermal condition of the lower part of the furnace.

The works of C.Staib and J.Michard (IRCID) repeated recently by P.Jourde and C.Remont⁵ have shown the practical fitness of the parameter W_u similar to index i_H for the inspection and control of the thermal condition in the lower part of the furnace.

The peculiarities of the blast furnace operation, namely the discrete and cyclic charging of raw materials into furnace impose certain conditions on the methods of calculation of indexes i_g and i_H from equations (2) and (3). The analysis of all the input information led to conclusion that charging cycle time, namely the period covering 5 - 7

chargings (depending upon the agreed charging program) should be considered an optimal interval between the successive calculations of indexes i_g and i_n . It requires averaging the information during the charging cycle. The generalized parameters i_g and i_n determined under such conditions characterize the thermal condition of the upper and lower parts of the furnace during the previous charging cycle, that is, during the last half an hour of the furnace operation.

The revealed non-interaction in thermal operation of the upper and lower zones, of the furnace as well as the above shown possibility of independent estimation of the thermal condition of the two zones enable us to consider the blast furnace to consist of two interconnected but self-contained controlled plants with their own static and dynamic properties. Such quite new treatment of the blast furnace as a controlled plant reflects the nature of the blast furnace melt more accurately thereby increasing the ways of controlling this process.

Differentiated approach to the upper and lower thermal zones of the furnace as independent controlled plants has enabled us to make an analysis of their static and dynamic properties when using different controlling factors. Fig. 2 and 3 show the static and dynamic characteristics determined by calculation, the static characteristics being calculated according to known dependences for the settled thermal condition of the furnace while the dynamic characteristics are determined by the method of approximate calculation of transient processes.

A slight initial cooling of the upper and part of the furnace after a stopped rise of specific coke consumption is accounted for by gradual increase (fig.3) of the thermal capacity of burden materials with changing the relative proportion of coke. Warming up the shaft of the furnace begins only after the burden of a new composition reaches the tuyers and as a result of combustion of additional coke the yield of hearth gas increases.

The initial drop of index i_n during the natural gas

supply (fig.3) is due to the heat consumed for converting the injected fuel. Later on due to increasing hydrogen content in reduction gases the amount of coke carbon consumed for direct reduction decreases and a slight increase of the proportion of coke burning out at the tuyers results (the predetermined fuel consumption and other controlling factors except the one under review are assumed to be unchanged). Similar results were obtained experimentally⁷, the injected fuel being masont.

The revealed difference in static and dynamic properties of the controlled plants we have examined enables us to combine the controlling factors into the complex actions which have necessary selective (local) effects on the heat regime of only one zone of the blast furnace. Such quick-response complex actions which are to correct the detected deviations in the heat regime of individual blast furnace zones are called corrective.

As an example of such simplest corrective action which is efficient for the upper zone of the furnace and practically neutral for its lower zone we may give changing the oxygen content in the blast (in case $t_p > 1000^\circ\text{C}$). Quite opposite is the response observed when applying another simplest corrective action, namely the change in the blast temperature (Fig.3).

The corrective actions are composed in such a way that while applying them the thermal condition of the main (limiting) zone of the furnace should return within the limits of the zone of permissible deviations from the optimum during one or two charging cycles, and afterwards should not come beyond the above limits. It should be born in mind that each carefully selected corrective action represents not only a set of necessary discrete changes of different controlling factors but also a definite programme of action (sequence and speed of change in blast parameters must be indicated). It makes possible to take into account the dynamics of transient processes and the imposed limitations when applying actions "from below".

There may be several corrective actions exerting similar

influence upon the thermal condition of the furnace. The selection of the optimum variant is due to the available possibilities of using various controlling factors and depends also on such limiting conditions as the necessity to retain gas dynamics of charge column and reduction potential of gases. For particular conditions of melting in the blast furnace it is helpful to calculate in advance all the possible variants of solving problem, and to present the obtained results as ready-made tables with logical schema of selecting necessary corrective action.

Thus the application of generalized parameters i_g and i_n and the possibility of using quick-response corrective actions proved above makes it possible to raise the problem of stabilization and optimization of the heat regime of the blast furnace.

As it follows from the very notion of the index

$(i_n)_{opt} = 1$, the value of the index i_n corresponds to the optimum thermal condition of the lower section of the furnace.

To find the optimum value of index $(i_g)_{opt}$ special investigations⁶ were made. These investigations have revealed the extreme character of the dependence between the thermal condition of the upper section of the furnace and its output (Fig.4), the configuration of the curves being determined by the particular conditions of furnace operation. For experimental data shown in Fig.4 the optimum values of index $(i_g)_{opt}$ corresponding to maximum output of the furnace range from 0.65 to 0.77. For other furnaces and conditions of operation the extreme character of the dependence under investigation is retained but the values of index $(i_g)_{opt}$ may be different.⁶

The results obtained affirm convincingly enough considerable influence of the thermal condition of the shafts of modern furnaces using oxygen injected fuel on the main indexes of the blast furnace process. It proves that the optimum variant of melting technique in the blast furnace (forced smooth operation of the furnace when melting iron

of pre-determined composition with least amount of coke consumed) is possible only if the heat regime is stabilized at the optimum level not only in the lower section of the furnace but at the same time in the upper one as well.

The suggested algorithm of controlling thermal condition of the blast furnace is as follows: The optimum value of index $(i_g)_{opt}$ is determined statistically or by some other method and the value Q_o is specified in the same way. The current values of indexes i_g and i_H are recurrently calculated once during the charging cycle. Then the deviations $\Delta i_g = (i_g)_{opt} - i_g$ and $\Delta i_H = (i_H)_{opt} - i_H$ are defined, and according to the sign and the value of the latter the necessary corrective actions are calculated (or selected from ready solutions) by a computer allowing for additional conditions.

The agreed decision about the value and sequence of changing various controlling factors holds true until the steady deviations from the heat regime rates of the upper or lower sections of the furnace are fixed again. In this case the pre-selected controlling action is repeated or replaced by the new one according to the particular conditions of furnace operation.

In order to increase reliability of the taken decisions about the selection of controlling actions it was decided:

1. Definite zones of insensitiveness whose range somewhat exceeds permissible miscalculations in measuring input values must be established just the same as in the case of the discrete control systems. If the permissible deviations of measurements from the range of controlled values of these indexes are $\pm 2,5\%$ for index i_g and $\pm 6\%$ for index i_H (for the index i_g in the range of 0.65-0.85 and for index i_H in the range of 0.75-1.25), the zones of insensitiveness may be assumed to be equal $\pm 0,005 \cdot i_g$ and $\pm 0,03 \cdot i_H$.

2. To eliminate the influence of occasional variations in the heat regime of the furnace and to secure the necessary holding for the estimation of the results of

steps taken, the authors suggest to take into account only those deviations $\pm \Delta i_g > 0.005$ and $\pm \Delta i_H > 0.03$, which are retained at two successive calculations of these values. Only such steady deviations may be considered discrete controlling signals on the basis of which one may estimate the changes in the thermal condition of the furnace and re-consider the decisions taken before.

Therefore, the rate of the work of a computer can't be pre-determined but is defined by technological process itself. For example, for the blast furnace of 1513 m^3 the charging cycle which was assumed to be basic period for all main calculating operations is 15-30 minutes if the output of the furnace is high and 40-50 minutes if it is low. The possible replacement of controlling commands by others in 15-50 minutes is quite acceptable in controlling blast furnace process and requires no quick-acting computers.

It is known that the main controlling factor which determines the heat regime of the melting in the blast furnace is the specific consumption of coke or the amount of ore loads. Considerable persistence and some other properties of this action, however, make it difficult to use it for stabilization of thermal condition in different zones of the furnace. In connection with this the authors suggest that the actual influence of the fixed coke consumption (or ore loads) on the heat regime of different zones of the furnace be corrected through small changes of the blast temperature Δt_a , its humidity $\Delta \varphi_a$, the content of oxygen in the blast ΔO_2 and the consumption of natural gas. $\Delta \Gamma$. It is not the suggestion of advisability of applying these controlling factors, which has long been proved in practice, that is new but the possibility to arrange them into purposeful controlling actions exerting necessary influence on the thermal condition of different zones of the furnace determined by us.

Below we give the example explaining the selection of corrective actions:

Assume that under particular operating conditions of the blast furnace of 1242 m^3 ($K = 550 \text{ kg/ton iron}$), $t_2 = 1000^\circ\text{C}$, $\varphi_2 = 15 \text{ g/m}^3$, $O_2 = 22\%$, $\Gamma = 30 \text{ m}^3/\text{ton iron}$, $(i_g)_{opt} = 0.78$) the control system detected in succession deviations $+\Delta i_g > 0.005$ and $-\Delta i_H > 0.03$, which shows steady indications of warming-up upper section and cooling lower section of the furnace. For correcting the revealed deviation of heat regime from optimum one can apply the corrective action providing for increase in oxygen content of the blast and at the same time decrease in blast humidity: $\Delta O_2 = +0.9\%$, $\Delta \varphi_2 = -4 \text{ g/m}^3$ of the blast.

The rate of changing these controlling factors is assumed to be the same, since their dynamic properties coincide. Such corrective action, in fact, will not affect the aerodynamic and reducing processes of melting but under the influence of this action the thermal condition of upper and lower sections of the furnace will return to its optimum values.

The block diagram of the UPI (the Urals Polytechnical Institute) system (Fig.5) is suggested as one of the variants of carrying out the conception of independent inspection and local control of the thermal condition of the upper and lower sections of the furnace. Besides conventional systems of control and stabilization of different input parameters we suggest that a simple computer be used which calculates values of indexes i_g and i_H , determines the deviations of these indexes from their optimum values $\pm \Delta i_g$ and $\pm \Delta i_H$ once during the charging cycle, and according to the sign and value of the latter selects the necessary corrective action allowing for imposed limitations. The determined decision is carried out by means of changes according to a definite programme of targets for controllers and stabilizers of natural gas consumption and blast parameters or by means of delivery of recommendations for changing specific coke consumption.

REFERENCES

1. Б.И.Китаев, Ю.Г.Ярошенко, Б.Л.Лазарев. Теплообмен в доменной печи. Изд. "Металлургия", 1966.
2. B.I.Kitaev, Yu.Yaroschenko, V.D.Suchkov. Heat Exchange in Shaft Furnaces. Pergamon Press, Oxford, 1967.
3. B.I.Kitaev, Y.C.Yarochenko, B.L.Lasarev. Etat actuel de la theorie des echanges thermiques dans le haut fourneau. Troisiemes Journees Internationales de Siderurgie, Luxemburg, 1962.
4. Ю.Г.Ярошенко, Б.И.Китаев, В.С.Швыдкий, Б.Л.Лазарев, Б.Л.Суханов. Контроль теплового состояния шахты доменной печи. Сб. "Автоматизация доменного производства", Изд. "Металлургия", 1966.
5. P.Jourde, C.Remont, C.Staib, N.Jusseau, A.Schaller. Control et reglage automatique d'un fourneau par un calculateur industriel. Journees Internationales de Siderurgie, Amsterdam, 1965.
6. Ю.Г.Ярошенко, Б.И.Китаев, Б.Л.Лазарев, В.С.Швыдкий, И.А.Тациенко, Б.М.Герман. Тепловое состояние зоны непрямого восстановления доменных печей и его контроль. Черметинформация, Серия 4, инф. 15, 1966.
7. P.Guillemain, N.Jusseau, P.Bece. Conduite automatique et etude des regimes transitoires de deux hauts fourneaux alimentes et agglomere. Revue de Metallurgie, 1967, Nr 11.

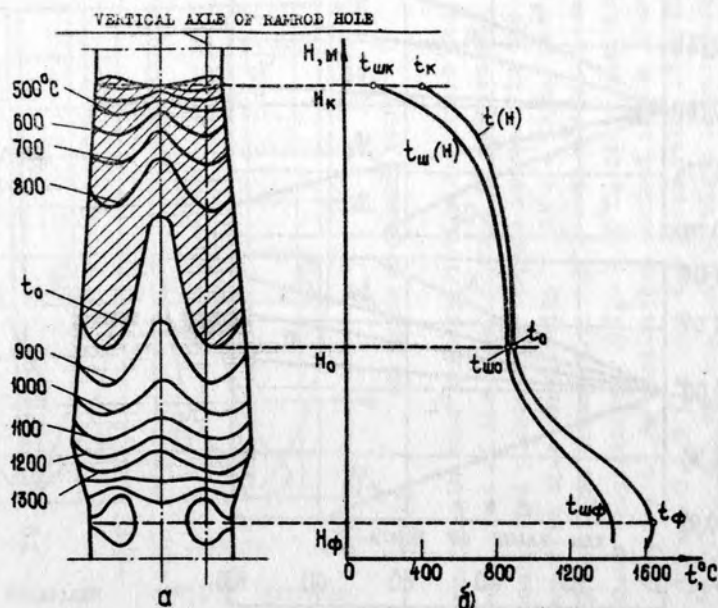


Fig.I

The nature of temperature profile of the blast furnace of 1242 m³ when melting pig-iron and under normal operating conditions: a - the working space of the furnace is divided into upper (shaded) and lower heat zones; b - the curves of changing temperatures of countercurrents of charge and gas along the axle of ramrod hole.

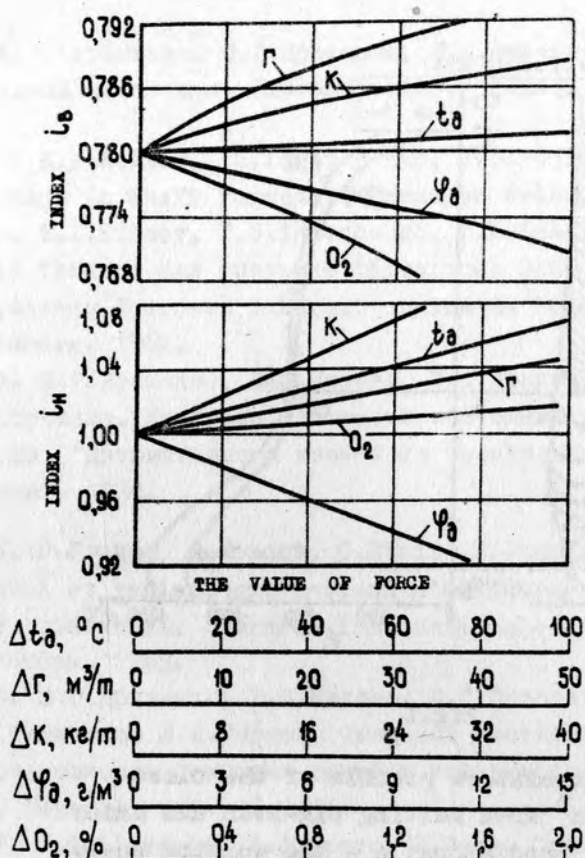


Fig. 2

Calculated statistical characteristics of individual (different) controlling actions for particular operating conditions of the furnace of 1242 m^3 when changing specific coke consumption K from 550 to 590 kg/t of iron, blast temperature from 1000 to 1100°C , blast humidity φ_a from 7 to 22 g/m³, oxygen content in the blast O_2 from 21 to 23 % and natural gas consumption r from 10 to 60 m³/ton iron.

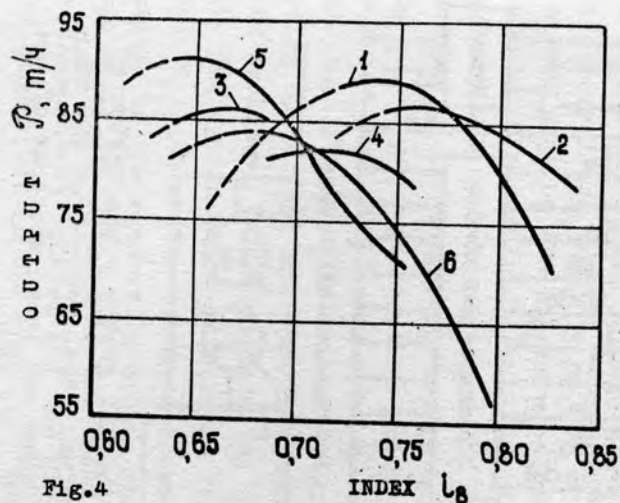


Fig. 4

The relationship between index l_g and the output of the blast furnace of 1242 m³ defined, according to average data, for 24 hours for 6 months' operation: 1-, 2- and 3- October, November and December, 1965; 4-, 5-, 6- January, February and March, 1966.

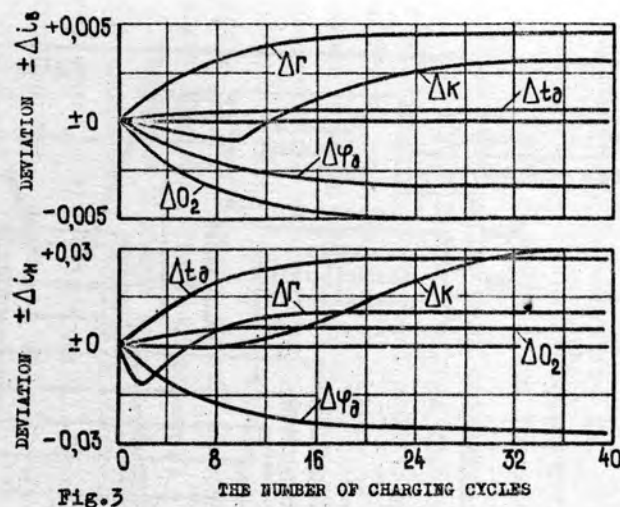


Fig. 3

Approximate time characteristics for upper and lower heat zones of the blast furnace as independent controlled plants with stepped increase in the specific coke consumption $\Delta\kappa = +10$ kg/ton iron, natural gas consumption $\Delta r = +10$ m³/ton iron, blast temperature $\Delta t_z = +40^\circ\text{C}$, blast humidity $\Delta\varphi_z = +4$ /m³ and oxygen content in the blast $\Delta O_2 = +0.6\%$ (in case the index value is optimum $(l_g)_{opt} = 0.78$, the initial blast temperature $t_z = 1000^\circ\text{C}$ and other particular conditions).

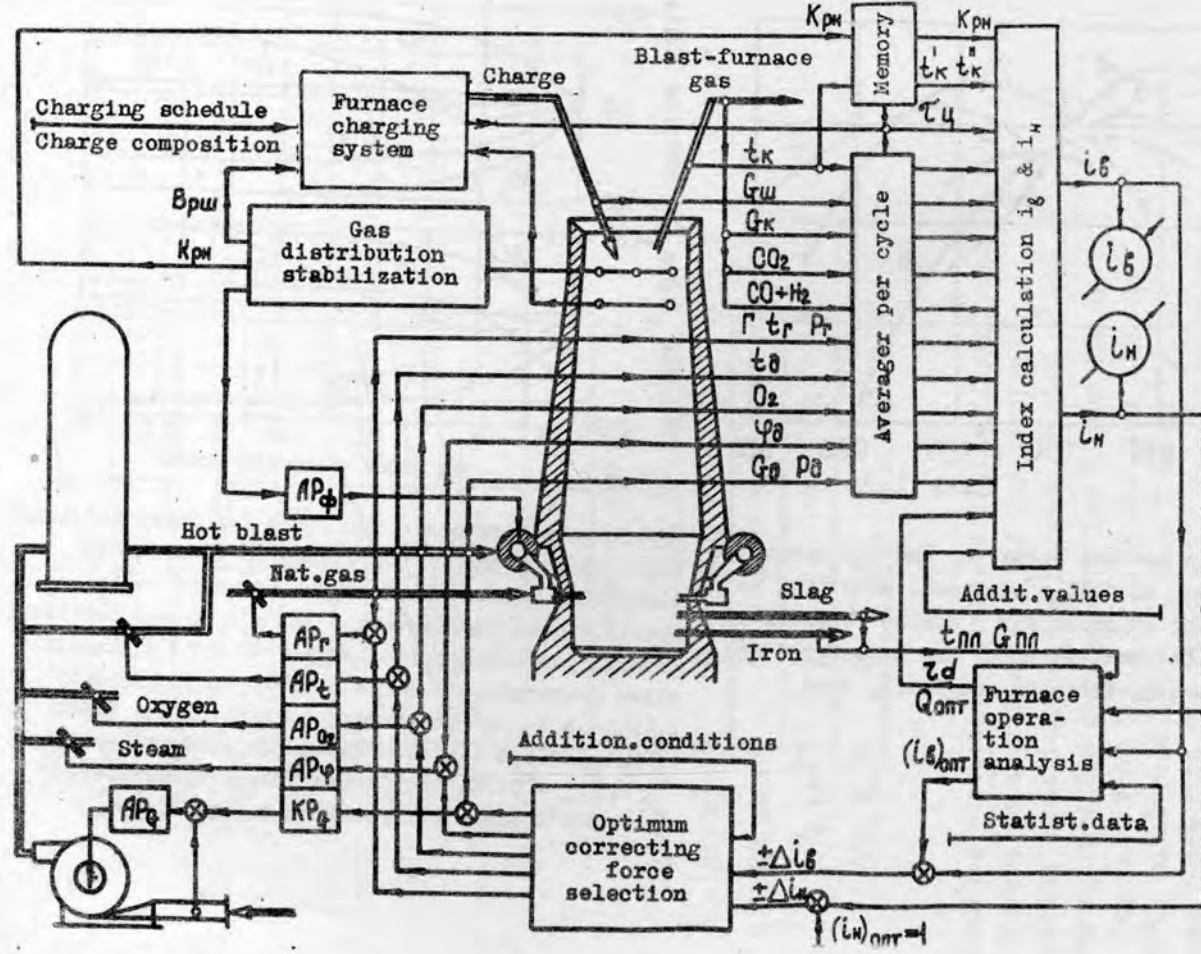


Fig.5 The block diagram of the UPI system - the system of independent inspection and local control of the thermal condition of the upper and lower sections of the blast furnace using combined blast.

Rational algorithm of controlling the thermal condition of blast furnace using computers

E.L.Suchanov, V.S.Shvidki, B.I.Kitaev, Ju.G.Yaroschenko,
Ju.N.Ovchinnikov, V.G.Lisienko

(The Urals Polytechnical Institute, Sverdlovsk, USSR).

The autonomy in heat operation of upper and lower stages in a blast furnace called for the necessity of different evaluation of their heat conditions. „The index i_B ” is offered as a general variable of the heat condition in the upper part of the blast furnace and „the index i_H ” is proposed for controlling heat conditions in the lower part of the furnace. These variables are calculated only once per cycle of the furnace charge according to usual information about the technological process.

The blast furnace should be considered as two connected but independent objects of controlling with their own static and dynamic characteristics. When analysing those characteristics the conclusion about real combination of such controlling factors is obtained. The general influence of these factors on the furnace heat conditions has its necessary local action.

Both statics and dynamics of transient processes in objects of controlling are considered. Each controlling effect is a certain programme of the necessary digital changes in the temperature and humidity of the blast, in the higher oxygen concentration and in the expense of the injected fuel. The lower specific expense of the coke can be assumed because of economic consideration. All the calculating operations are performed by the information controlling machines. In the article the possibility of stabilization and optimum heat conditions in the modern blast furnace which is worked at a combined blast is proved. The block-diagram of the automation system which is resulted from the idea about the optimum in the blast furnace process under the independent control and the local stabilization of the heat conditions in the upper and lower stages in the furnace is given.

TIME SUB-OPTIMUM CONTROL OF THE WORK OF CRANES
WITH SPECIAL REGARD TO ITS REALIZATION IN PRACTICE

M. Sc. Roman Gorecki
Academy of Mining and Metallurgy
Department of Automation and Industrial
Electronics
Kraków, Poland

1. Introduction

There is already ample literature available in the field of optimum control and new items are constantly being published. As a rule these are theoretical research works and seldom concerned with a practical solution. The methods suggested in literature usually lead to very complex solutions requiring the use of computers or necessitating carrying out of numerous measurements.

The present work is an attempt to solve time sub-optimum control by means of a specific process with a special regard to practical realization.

2. Description of the process

We shall concern ourselves with the design of a system for an optimum, meaning "minimum time", control of the work of a crane. We shall limit ourselves to the investigation on the optimum control of cranes of the travelling bridge type, i.e. those in which we have at our disposal three independent drives for transporting the load. Each of them shifts the load along a different

axis of the Cartesian system. A situation of this type is presented in Fig.1. The load may be shifted in the vertical plane /movement C/ and transferred to an optional position in a horizontal plane by means of drives A and B. The load may be shifted simultaneously in three directions. This is often used by the operator. In order not to set the load in excessive swinging motion the operators apply low travelling rates. In spite of this a-fter arriving above the destination the load will as a rule oscillate enough to make it impossible for the operator to lower it to the foreseen place required without damping of the oscillations.

When the crane is in constant use, e.g. in the case of harbours concerned with reloading of goods, it is worthwhile considering the possibility of modernizing the control, i.e. changing the control in order to reduce the time required for transporting loads as well as making greater use of the installed power.

To begin with a somewhat simpler case than the one described above will be considered. It is assumed that the load hangs below the trolley on a weightless rope of a constant length and has the whole mass concentrated in one point at a distance "l" from the point of suspension, otherwise it constitutes a mathematical pendulum. Controlling of the trolley motion is limited to the direction "Z" only, /Fig.2/. The trolley drive is located beyond the trolley. The drive consists of a motor with rotations controlled within certain limits, transmission and coupling through with the rope hauling the trolley is coupled with the operating drive or brake. In addition, it is assumed that the trolley mass is so small in relation to the inertia of mas-

ses in motion as to be negligible. The shifting of the point of the engagement from a certain initial position is marked by the letter "y" and the shifting of the load along the axis "x" by the symbol "x" /Fig.3./.

Several forces will be acting on the load being in motion, i.e. the force of gravity "G", force of inertia "B", resistance of dynamic friction "T", and the force of the rope reaction "R". The equilibrium conditions will be written by projecting the forces in the direction normal to the track of motion, i.e. along the line

$$G_n - R = 0 \quad 2.1$$

where G_n is the projection of the force of gravity in a direction perpendicular to the track of motion

R - reaction in the rope.

The next condition of the equilibrium of the forces projection in the "x" direction will be given:

$$-B_x - T_x + G_s \cdot \cos \alpha - G_n \sin \alpha + R \sin \alpha = 0 \quad 2.2$$

B_x is the component of the force of inertia in the direction of the "x" axis. It is proportional to the load mass "m" and its acceleration in the direction of this axis.

$$B_x = mx'' \quad 2.3$$

T_x is the component of resistances to motion proportional to the rate of travel of the load in the direction of the "x" axis.

$$T_x = r x' \quad 2.4$$

where "r" is the coefficient of the resistances to motion.

G_s is the component of the forces of gravity "G" tangent to the track of motion.

$$G_s = G \sin \alpha \quad 2.5$$

Taking into account the relationship 2.1 to 2.5 in the formula 2.2. one may write:

$$-mx'' - rx' + G \sin \alpha \cdot \cos \alpha = 0 \quad 2.6$$

For small angles " α " one may assume

$$\cos \alpha = 1 \quad 2.7$$

and since from Fig.3 one sees that:

$$\sin \alpha = \frac{y - x}{l} \quad 2.8$$

therefore for small deviation the relation 2.6 will assume

$$\text{the form: } -mx'' - rx' + G \cdot \frac{y - x}{l} = 0 \quad 2.9$$

After ordering one obtains:

$$x'' + \frac{r}{m} x' + \frac{G}{m \cdot l} (y - x) = 0 \quad 2.10$$

and since

$$G = m \cdot g \quad 2.11$$

therefore:

$$x'' + \frac{r}{m} x' + \frac{g}{l} x = \frac{g}{l} y \quad 2.12$$

By designating:

$$\frac{r}{m} = a$$

$$\frac{g}{l} = b \quad 2.13$$

one obtains:

$$x'' + ax' + bx = by \quad 2.14$$

Since in fact the rate of the trolley travel is controlled therefore its shifting equals;

$$y = \int v \cdot dt \quad 2.15$$

The rate of the trolley travel can assume values contained within the interval from 0 to $v_{\max} = u > 0$, taking into account 2.15 in 2.14:

$$x'' + ax' + bx = \int_0^t bv \, dt \quad 2.16$$

Differentiating both sides of 2.16 one obtains:

$$x'''' + ax''' + bx'' = bv$$

$$0 \leq v \leq v_{\max} \quad 2.17$$

The plant described by the above equation should be transported from the initial state:

$$x'' / 0 = x' / 0 = x / 0 = 0 \quad 2.18$$

at the instant $t_0 = 0$ to the terminal state:

$$x'' / t_k = x' / t_k = 0 ; x / t_k = x_k \quad 2.19$$

at the instant " t_k " in the possibly shortest time.

3. Application of the "maximum principle"

Since equation 2.17 is a linear one the control function "v" for obtaining the time optimum course should according to the "maximum principle" be constant in the intervals and assume extreme values alternately. The maximum principle gives us in this case the necessary condition of optimality in the sense of minimum time.

The third order equation 2.17 is converted into a first order system of equation which one obtains by substituting:

$$x_1 = x ; x_2 = x' ; x_3 = x'' \quad 3.1$$

and now we obtain:

$$x_1' = x_2$$

$$x_2' = x_3 \quad 3.2$$

$$x_3' = bx_2 - ax_3 + bv$$

The Hamiltonian function:

$$H = \sum_{k=1}^3 \psi_k \frac{dx_k}{dt} \quad 3.3$$

taking into account 3.2 we have:

$$H = \psi_1 \cdot x_2 + \psi_2 \cdot x_3 + \psi_3 / -bx_2 - ax_3 + bv / \quad 3.4$$

This function will attain a maximum by a control in which

$$v = \text{sign } /b \psi_3/ \quad 3.5$$

which in the case discussed is reduced to assuming the value:

$$v = \begin{cases} u & \text{for } \psi_3 > 0 \\ 0 & \text{for } \psi_3 \leq 0 \end{cases} \quad 3.6$$

For determining ψ_3 as a function of time we make use of the dependence:

$$\frac{d\psi_i}{dt} = - \frac{\partial H}{\partial x_i} \quad / i = 1, 2, 3 / \quad 3.7$$

The following system of equations is to be solved:

$$\begin{aligned} \psi_1' &= 0 \\ \psi_2' &= -\psi_1 + b \psi_3 \\ \psi_3' &= -\psi_2 + a \psi_3 \end{aligned} \quad 3.8$$

wherefrom

$$\begin{aligned} \psi_1 &= c_1 \\ \psi_3'' &= -\psi_2 + a \psi_3' \\ \psi_3'' &= \psi_1 - b \psi_3' + a \psi_3' \end{aligned} \quad 3.9$$

Finally the following equation is to be solved:

$$\psi_3'' - a \psi_3' + b \psi_3 = c_1 \quad 3.10$$

After solving equation 3.10 the function ψ_3 is expressed by the formula:

$$\begin{aligned} \psi_3 / t / &= \frac{s_1^2 c_3 + s_1 c_2 - s_1 a c_3 + c_1}{3s_1^2 - 2as_1 + b} \cdot e^{s_1 t} + \\ &+ \frac{s_2^2 c_3 + s_2 c_2 - s_2 a c_3 + c_1}{3s_2^2 - 2as_2 + b} \cdot e^{s_2 t} + \frac{c_1}{b} \end{aligned} \quad 3.11$$

where:

$$\begin{aligned} c_2 &= \psi_3' / 0 / \\ c_3 &= \psi_3 / 0 / \end{aligned} \quad 3.12$$

If we assume $a = 0$, which in the case under discussion we may

safely do, since the damping of the motion by the air resistance is minimum then the function ψ_3 assumes a simpler form:

$$\psi_3 / t = \frac{c_1}{b} / 1 - \cos \sqrt{b} t + c_3 \cos \sqrt{b} t + c_2 \frac{1}{\sqrt{b}} \sin \sqrt{b} t \quad 3.13$$

designating:

$$\sqrt{b} = \omega \quad 3.14$$

$$\psi_3 / t = \frac{c_1}{\omega^2} / 1 - \cos \omega t + c_3 \cos \omega t + c_2 \frac{1}{\omega} \sin \omega t \quad 3.15$$

On the basis of equation 3.8 and 3.15 we may compute the remaining auxiliary functions:

$$\psi_1 = c_1$$

$$\psi_2 = /c_3 \omega - \frac{c_1}{\omega} / \sin \omega t + c_2 \cos \omega t \quad 3.16$$

$$\psi_3 = \frac{c_1}{\omega^2} + /c_3 - \frac{c_1}{\omega^2} \cos \omega t + c_2 \frac{1}{\omega} \sin \omega t$$

The system of equations 3.16 may be presented as:

$$\psi_1 = c_1$$

$$\psi_2 = R \sin / \omega t + \alpha /$$

$$\psi_3 = \frac{c_1}{\omega^2} + R \cos / \omega t + \alpha /$$

where:

$$R = \frac{1}{\omega} \sqrt{/c_3 \omega - \frac{c_1}{\omega} /^2 + c_2^2}$$

$$\alpha = \arctg \frac{c_2}{c_3 \omega - \frac{c_1}{\omega}}$$

3.17

In the Cartesian system of axes ψ_1, ψ_2, ψ_3 , the parametric equation 3.17 presents an ellipse laying in a plane perpendicular to axis ψ_1 which intersects this axis at point

$$\psi_1 = c_1.$$

In the system of axes $\frac{\psi_2}{\omega}$, ψ_3 we will have a circle with a radius "R" and centre laying on axis ψ_3 at point $\psi_3 = \frac{c_1}{\omega^2}$.

Fig.4. From the computations carried out the course ψ_3 as the function of time is known from which it follows that the character of the switching function "v" is known. Unfortunately, the Pontryagin method does not furnish us with information about the value of initial conditions c_1, c_2, c_3 , therefore on the basis of this function we know only that function "v" will be reproducible together with period "T", but the length of the switchings is not known neither is the time after which the first and last switching will occur. From the character of the control function it results that it may be considered as a sum of shifted in time constant inputs of "u" values.

With the assumption of zero initial conditions the output signal $x / t /$ will be equal to the sum of responses to each input signal. From formula 2.17 one can compute the response of the system to the step input signal

$$v = \begin{cases} 0 & \text{for } t < 0 \\ u & \text{for } t \geq 0 \end{cases} \quad 3.18$$

will be:

$$\begin{aligned} h / t / &= u / t - \frac{1}{\omega} \sin \omega t / \\ h' / t / &= u / 1 - \cos \omega t / \\ h'' / t / &= u \cdot \omega \sin \omega t \end{aligned} \quad 3.19$$

The response of the system to the control signal composed of constant functions alternately equalling "u" and "0" will be:

$$x / t / = u \sum_{k=0}^n (-1)^k \cdot h / t - t_k / \quad 3.20$$

4. Optimum control

To determine the optimum control additional information is necessary. In the case of a second-order equation of state it is very convenient to draw the phase trajectories in the x, x' system of coordinates. In the considered case of a third-order equation the trajectories should be per analogiam considered in a special system x, x', x'' which we are unable to illustrate clearly on a plane.

The equation of the system considered is degenerated. This is due to the fact that the value of the coefficient at " x " is equal to zero hence when observing the performance of the system in the coordinates x', x'' we shall have full information about all the derivatives and the function itself.

In the $x', \frac{x''}{\omega}$ system the trajectories of the system considered will be circles with centres located on the axis " x'' " at point $x' = v$.

Only the trajectories of the system around the extreme positions $v = 0$ and $v = u$ constitute the optimum control. From the initial and terminal conditions it results that the process must begin and end with a trajectory crossing through the origin of the system $\frac{x''}{\omega}$, x' shifting simultaneously x from x_0 to x_k , but this is not to be seen in Fig.4.

If the input signal $v = u$ acts on the system one will see that the leading point after time $t = T$ will return to the origin. The load will rest as at the moment of starting perpendicularly below the trolley, hence the shifting of the trolley and load will be the same and equal to:

$$x_{\Gamma} = x_{\Gamma} = u \cdot T$$

Therefore we see that for shifting the load over a distance

$$x = n \cdot x_T \quad n = 1, 2, 3, \quad 4.2$$

the necessary control will be limited to one switching period, This is in agreement with the result obtained by Pontryagin's method. This is a particular case in which the switching periods are equal to the period of function and the times of the trolley standstills are reduced to zero.

The number of switchings in an optimum control time are equal to:

$$i = 2 \left[E^* / \frac{x_k}{x_T} / + 2 \right] \quad 4.3$$

In Fig. 5 an example of a phase trajectory for $x_T \geq x_k \geq 0$ is shown. Owing to the relation

$$\alpha = \omega t \quad 4.4$$

and the observation that, according to the initial and terminal conditions, the track of shifting the load is equal to the track of the trolley travel and the time of travel corresponds to the time of moving along the circular trajectory with a centre at $x^* = u$ one may write:

$$x_k = u / T_1 + T_3 - T_2 / \quad 4.5$$

and

$$\alpha = \tilde{\eta} - \gamma \quad 4.6$$

and after taking into account the relation 4.4:

$$\begin{aligned} T_2 - T_1 &= \frac{T}{2} - T_1 \\ T_2 &= \frac{T}{2} \end{aligned} \quad 4.7$$

Reasoning thus:

$$T_k - T_1 = \frac{T}{2} \quad 4.8$$

$E^*/z/$ designates a complete part from number z

Finally we obtain:

$$\begin{aligned} T_1 &= \frac{x}{2u} \\ T_2 &= \frac{T}{2} \\ T_k &= \frac{x}{2u} + \frac{T}{2} \end{aligned} \quad 4.9$$

Determination of the optimum process for $2x_T > x_k > x_T$ can be carried out on the basis of the phase characteristics from the information already acquired which is presented in Fig.6.

The main indication when determining optimum trajectory is the fact that except the first and last intervals of time the sum of two consecutive switchings must be equal to the period T , according to Pontryagin's method.

From Fig.6 one may derive the relationship between angles α and δ :

$$\alpha = \frac{1}{2} \arccos \frac{2 - \cos \delta}{\sin \delta} \quad 4.10$$

and for x from an arbitrary interval $(n-1) \cdot x_T \leq x \leq n \cdot x_T$:

$$\alpha = \frac{1}{2} \arccos \frac{2 - \cos \delta}{\sin \delta} \quad 4.11$$

Knowing the magnitude of angle α corresponding to the intervals of the stillstand time and δ corresponding to the first and last time intervals of travel, the times of consecutive switchings may easily be completed. The operation total time will be:

$$T_{kn} = \frac{T}{360^\circ} [2\delta + \alpha + (n-1) \cdot 360^\circ] \quad 4.12$$

The load will be shifted to the point of destination:

$$x_k = \frac{x_T}{360^\circ} [2\delta + (n-1)(360^\circ - \alpha)] \quad 4.13$$

With the increment of the assumed travel distance, i.e. with increment of "n", α aims at "0" and the operation time will tend towards a sum time of travel T_j which will be:

$$T_j = \frac{x_k}{u} \quad 4.14$$

Since this is the shortest time during which the trolley is able to cover the assumed distance of shifting regardless of the oscillation of the load, the possible symmetric limitations of the control function $-u \leq v \leq u$ at a greater distance x_k would not shorten the time of duration of the process.

The increasingly shorter standstill periods with increment of x_k the accuracy of their realization conditioning the achievement of the final state is the reason that for greater x_k the process of time-optimum control becomes un-realizable. Moreover, considerable difficulties would be encountered if it proved necessary to introduce corrections taking into account the non-fulfillment of the assumptions made.

5. Time sub-optimum control

On practical grounds it may be profitable not to aim at obtaining optimum control, but to be satisfied with a control somewhat worse than optimum as regards the operation speed, but easier for realization.

Moreover, it is also worthwhile to consider the possibility of realizing a control which would limit the swinging of the load during travel, in particular when covering greater distances.

A control which would ensure non-oscillating travel over longer distances may easily be read from the phase plane

$$\frac{x''}{\omega}, x', \text{ Fig. 7.}$$

To this purpose the load rate should be made equal with the trolley rate equalling "u". The point on the phase trajectory should be carried from the origin of the system to the point "u" on the axis "x". This is carried out most quickly making use of a circular trajectory crossing through the system origin with a centre at point "u", and then from the trajectory circumscribed around the origin of the system shown in the figure and going through the point "u".

From the relationships seen in the figure, i.e. from the equilaterality of the triangle it results that angle:

$$\alpha_1 = \alpha_2 = \frac{1}{6} 2\pi \quad 5.1$$

which corresponds to:

$$T_1 = T_2 - T_1 = \frac{1}{6} T \quad 5.2$$

Therefore the control equivalizing the rate of the load shifting with the trolley travel rate will consist in switching the drive for a period $T_1 = \frac{1}{6} T$, then stopping the trolley for the period $T_2 - T_1 = T_1 = \frac{1}{6} T$ after which travel will be smooth during any desired length of a time interval $T_3 - T_2$. In order to stop the load an analogous operation must be carried out, i.e. the stopping of the trolley for the period $T_4 - T_3 = \frac{1}{6} T$, then again shifting it with a maximum travel rate "u" during the time $T_k - T_4 = \frac{1}{6} T$.

The load will remain motionless and the track covered will be:

$$x_k = [(T_k - T_4) + (T_3 - T_2) + T_1] u \quad 5.3$$

or

$$x_k = (T - \frac{1}{6} T) u \quad 5.4$$

This control may be used only for:

$$x_k = \frac{1}{3} T \cdot u \quad 5.5$$

Moreover, it may be estimated that the duration of the process in the most disadvantageous case differs from the optimum T_0 by the standstill time.

$$T_k - T_0 = \frac{1}{3} T \quad 5.6$$

and in the case of travel over long distances and small period of oscillations the time loss will constitute a negligible percent. In Fig.8 the time of the duration of the process is illustrated in proportional magnitudes at a time sub-optimum control and an optimum control with symmetric and asymmetric limitation of the control function for small x_k .

It is obvious that where four of the five time intervals are equal and are functions of only one magnitude T , and the fifth time interval is also presented by a rather uncomplicated formula

$$T_3 - T_2 = \frac{x_k}{u} + \frac{1}{3} T \quad 5.7$$

the realization of a control is comparatively easy.

In the Laboratory of the Department of Automation and Industrial Electronics of the Academy of Mining and Metallurgy in Cracow a system/shown in Fig.2/ controlled by the above described programme has been realized. A load of 4 kg weight suspended on a 1.75 m coupling bar is shifted for a distance of up to 3.5 m with a maximum trolley travel rate 0.25 m/sec. After having reached the destination point the oscillation amplitude does not exceed 0.5 cm. When controlling the travel by hand it is very difficult to attain the damping of oscillations simultaneously with the determination of the process at a determined

point. The loss of time as compared with the control according to an optimum time control depends on the operator's skill, but on the average it considerably exceeds the double time of the programmed control.

It should be stressed that this programme may remain unchanged in spite of the necessity of omitting the vertical obstacle if only the vertical motions are located in the period of the smooth travel, i.e. during the time interval from T_2 to T_3 and one returns to the previous level.

If the point aimed at is on another level it is enough to compute the two terminal time intervals as being equal to one third of the new period of proper oscillation.

If for the transporting of the load two drives, A and B, Fig.1, have to be used then the programme will continue to retain its simple form as opposed to the time optimum programme where the control with for example drive A must depend on the character of motions B and C.

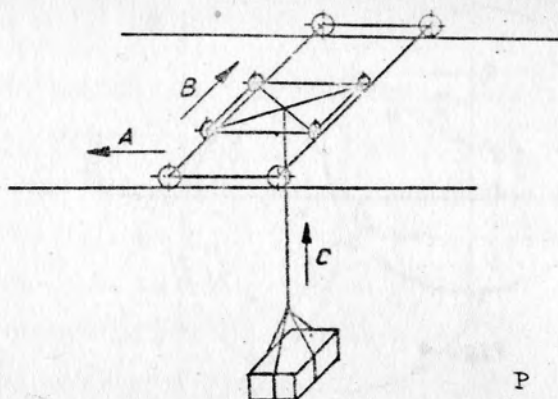


Fig. 1

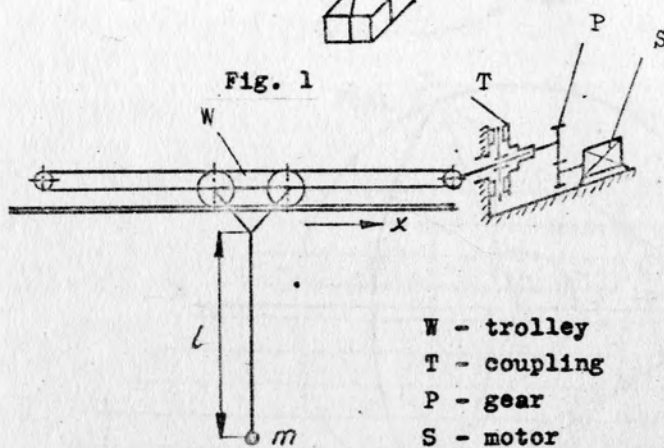


Fig. 2

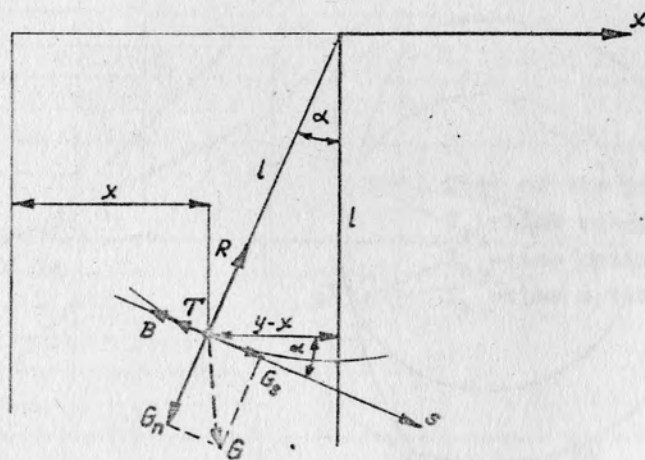


Fig. 3

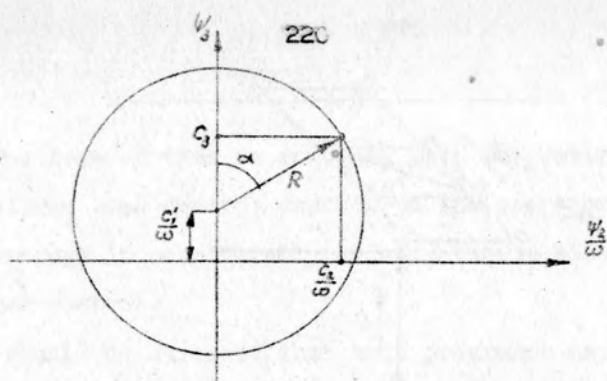


Fig. 4

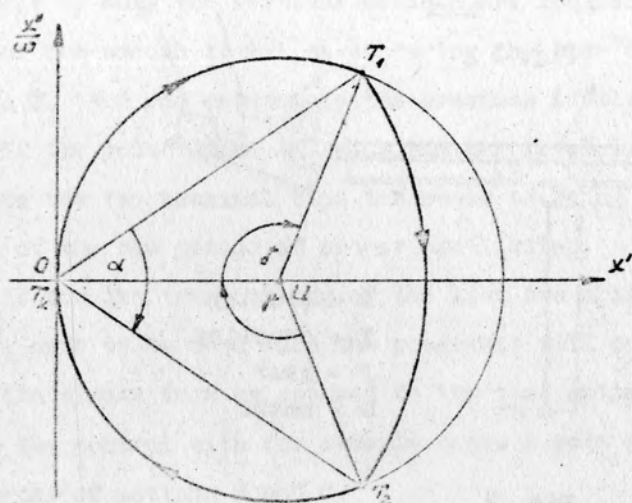


Fig. 5

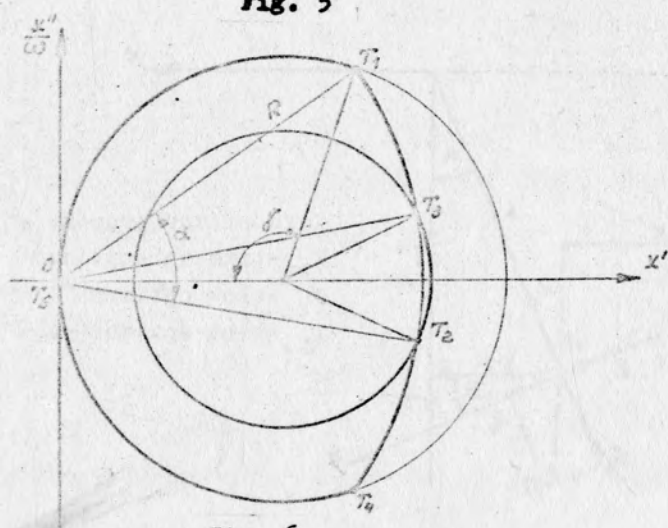


Fig. 6

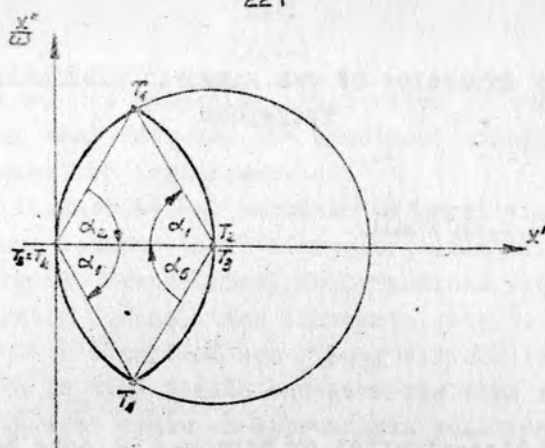


Fig. 7

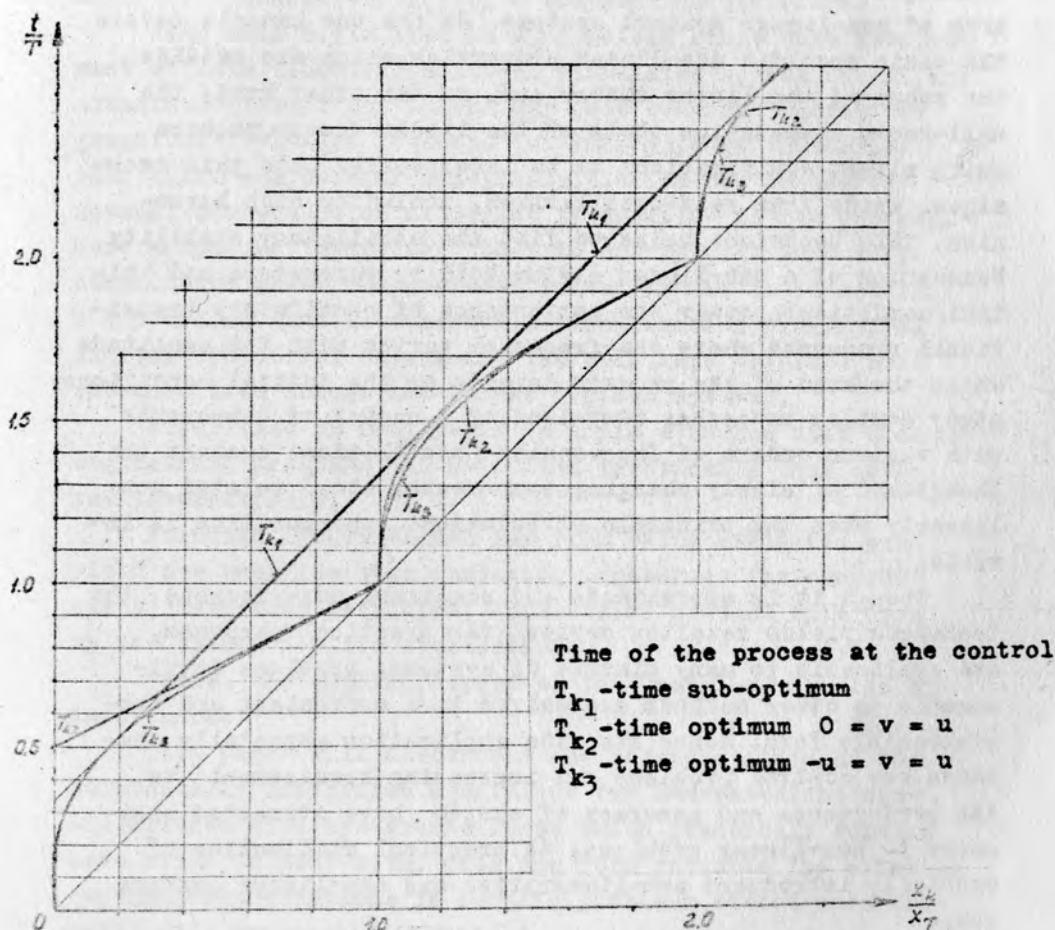


Fig. 8

AN EXTENSION OF THE HARMONIC LINEARIZATION
TECHNIQUE

Ye.P.Popov, Ye.I.Khlypalo

Moscow

USSR

Harmonic linearization or harmonic balance has become a technique used widely in study and design of a wide spectrum of non-linear control systems. On the one hand it covers the basic specific non-linear properties which are outside the scope of the linear theory and, on the other hand, the well-known computation tools of the linear theory require quite slight modifications to be incorporated into this technique. Aside from self-oscillations, including high harmonics, this technique helps to find the oscillatory stability boundaries of a non-linear system both by parameters and initial conditions; study the performance of oscillatory transitional processes where the frequency varies with the amplitude while the kind of the process depends on the initial conditions; study complex processes comprised of a number of components with various orders of frequencies (oscillations against the background of slowly changing components, etc.) related non-linearly when the principle of solutions superposition is invalid^{[1],[2]}.

Though it is approximate and sometimes non-rigorous, the technique yields results, correct for practical purposes, and applicable to many classes of systems; problems invulnerable to other methods are solved in a convenient and comprehensible form. Hence its wide application especially because new control problems and increasing requirement to the performance and accuracy of system have attracted engineers to non-linear problems, to practical utilization of specially introduced non-linearities and non-linear control laws.

The advances and expanding application of computers increase rather than decrease the practical significance of such approximation techniques.

Harmonic linearization (harmonic balance) yields readily the relations between basic characteristics of processes such as amplitudes, frequencies, decay indices and parameters of the system (gains, time constants, etc.). No accurate solution to the initial non-linear differential and other equations in time domain can give the same results. Furthermore, in most cases an approximate solution by the harmonic linearization technique cannot be replaced by more accurate - and more involved - computation procedures.

With this basic fact in mind we can state that development of this technique has been stimulated by the advent and advances of computing technology just because "manual" or graphical treatment of these problems appears cumbersome when there are several non-linearities in the system or several components of different frequencies in the process. Here computing technology comes to our rescue. Many very complicated problems have been solved so far in industrial R & D bodies that used harmonic linearization programs on digital computers, e.g. M-20; the solutions obtained were embedded into actual non-linear control systems.

Experience of many years is ample evidence that much more engineering problems can be solved practically than covered theoretically.

Harmonic linearization far from being exhausted can yield new benefits when computing procedures for complex cases have been developed; theoretical studies will open up now fields for its application.

Hence, the applied aspect of the technique is also of importance.

This paper will describe a new representation of the harmonically linearized equivalent for non-univalent non-linearities with hysteresis loops which practically workers have often to deal with. This new form corrects the existing non-rigorous notations in certain problems stated when formally this approach is utilized.

A harmonically linearized expression for non-univalent non-linearities

$$y = F(x) \quad (1)$$

with hysteresis loops is normally given by

$$y = \left[\alpha(A) + \frac{\beta(A)}{\omega} p \right] x, \quad \beta(A) < 0, \quad (2)$$

where A, ω are the oscillations amplitude and frequency, $\alpha(A)$ and $\beta(A)$ are coefficients

$$\alpha = \frac{1}{\pi A} \int_0^{2\pi} F(A \sin \psi) \sin \psi d\psi, \quad \beta = \frac{1}{\pi A} \int_0^{2\pi} F(A \sin \psi) \cos \psi d\psi.$$

As a result the transfer function of an open-loop system has a pole in the right-hand semi-plane (which changes its position when the amplitude of the oscillations changes). In terms of the linear control theory an equivalent linearized non-minimal phase element appears here.

Sometimes when the loop is closed no drastic changes occur and all operations proceed normally. However, there is a class of systems where this makes certain coefficients of the characteristic equation negative which would seem to prove that the oscillations in question are unstable, but actually these remain stable and the solution for the oscillations amplitude and frequency is valid.

An example of such a system is shown in Fig. 1. The dynamics of the system is described by

$$\begin{aligned} \varepsilon &= g_1 - g_2, & u_3 &= \frac{k_3}{1 + T_2 p} (u_2 - u_{ps}), & \Omega &= k_4 u_4, \\ u_1 &= k_1 \varepsilon, & u_{ps} &= k_{ps} p \Omega, & g_2 &= \frac{k_5}{p} \Omega, \\ u_2 &= \frac{k_2 u_1}{1 + T_1 p}, & u_4 &= F(u_3) = \alpha(A) u_3 + \frac{\beta(A)}{\omega} p u_3, \\ k &= k_1 k_2 k_3 k_4 k_5, & k_6 &= k_3 k_4 k_{ps}. \end{aligned}$$

The transfer function of that part of the system which has a feedback loop is

$$W_1(p) = \frac{\Omega}{u_2} = \frac{k_3 k_4 \left[\alpha(A) + \frac{\beta(A)}{\omega} p \right]}{1 + T_2 p + k_6 \left[\alpha(A) + \frac{\beta(A)}{\omega} p \right] p}.$$

The transfer function of an open-loop system is

$$W(p) = \frac{k[a(A) + \frac{b(A)}{\omega} p]}{p(1+T_1 p) \{1 + T_2 p + k_6 p [a(A) + \frac{b(A)}{\omega} p]\}}.$$

The characteristic equation of the system under consideration is

$$W(p) + 1 = 0$$

or

$$A_0 p^4 + A_1 p^3 + A_2 p^2 + A_3 p + A_4 = 0.$$

The factors of the characteristic equation will be

$$A_0 = T_1 k_6 \frac{b(A)}{\omega},$$

$$A_1 = T_1 T_2 + T_1 k_6 a(A) + k_6 \frac{b(A)}{\omega},$$

$$A_2 = T_1 + T_2 + k_6 a(A),$$

$$A_3 = 1 + k \frac{b(A)}{\omega}, \quad A_4 = k a(A).$$

By eq. (2) the coefficient A_0 is negative. The coefficient A_3 can also prove negative whereas all other coefficients are necessarily positive.

This is the result of erroneous form for equivalent harmonically linearized expression (2) which is normally employed in harmonic balance. We can preserve the same validity the solution for the amplitude and frequency and avoid the above non-rigorous intermediate expression (characteristic equation) if non-linear component (1) with an ambiguous loop non-linearity of a hysteresis type is assumed to be in the form of an inertial component equivalent to harmonically linearized transfer function given by

$$y = \frac{k_* A}{1 + T_* p} x. \quad (4)$$

New harmonic linearization coefficients $k_*(A)$ and $T_*(A)$, that is the gain and the time constant of an inertial component equivalent to a non-linear component with hysteresis loops can be expressed identically through the former harmonic linearization coefficients $a(A)$ and $b(A)$.

In the case of periodic oscillations when $p = j\omega$ from the desired identity

$$\frac{k_*}{1 + T_* j\omega} = a + \frac{b}{\omega} j\omega \quad (b < 0)$$

follows

$$\frac{k_*}{1 + T_*^2 \omega^2} = a(A), \quad \frac{k_* T_* \omega}{1 + T_*^2 \omega^2} = -b(A),$$

hence

$$T_* = \frac{-b(A)}{\omega a(A)}, \quad k_* = \frac{a^2(A) + b^2(A)}{a(A)}, \quad (5)$$

where $b(A) < 0$; $a(A)$ and $b(A)$ are found by eq. (3).

In the case of decaying and diverging oscillations the equivalent harmonic linearized expression for non-univalent loop non-linearity is given, instead of by eq. (2), by [1]

$$y = \left[a(A) - \frac{\xi}{\omega} b(A) + \frac{b(A)}{\omega} p \right] x, \quad (6)$$

where $\xi < 0$ for decaying oscillations and $\xi > 0$ for diverging oscillations.

In this case when $p = \xi + j\omega$ from the desired identity

$$\frac{k_*}{1 + T_* (\xi + j\omega)} = a - \frac{\xi}{\omega} b + \frac{b}{\omega} (\xi + j\omega)$$

follows

$$\frac{k_* (1 + T_* \xi)}{(1 + T_* \xi)^2 + T_*^2 \omega^2} = a(A), \quad \frac{k_* T_* \omega}{(1 + T_* \xi)^2 + T_*^2 \omega^2} = -b(A),$$

Hence

$$\omega T_* = \frac{-b(A)}{a(A) + \frac{\xi}{\omega} b(A)}, \quad k_* = \frac{a^2(A) + b^2(A)}{a(A) + \frac{\xi}{\omega} b(A)}, \quad (7)$$

where $b(A) < 0$; $a(A)$ and $b(A)$ are also found eq. (3).

In the case of complex processes when the oscillations are superimposed on a slowly changing component, an approximate solution instead of $x = A \sin \omega t$ is found in

the form $x = x^0 + x^*$, $x^* = A \sin \omega t$, and the usual form of harmonic linearization instead of (2) is given by [2]

$$y = F^0(A, x^0) + \left[\alpha(A, x^0) + \frac{\beta(A, x^0)}{\omega} p \right] x^*, \quad (8)$$

where

$$\left. \begin{aligned} F^0 &= \frac{1}{2\pi} \int_0^{2\pi} F(x^0 + A \sin \psi) d\psi, \\ \alpha &= \frac{1}{\pi A} \int_0^{2\pi} F(x^0 + A \sin \psi) \sin \psi d\psi, \\ \beta &= \frac{1}{\pi A} \int_0^{2\pi} F(x^0 + A \sin \psi) \cos \psi d\psi. \end{aligned} \right\} \quad (9)$$

Then the non-linear system equation will be

$$Q(p)x + R(p)F(x) = N(p)f(t), \quad (10)$$

where $f(t)$, a external action changing slower than the ω oscillations frequency decomposes, after harmonic linearization, into two non-linearly related equations

$$Q(p)x^0 + R(p)F^0(A, x^0) = N(p)f(t), \quad (11)$$

$$Q(p)x^* + R(p) \left[\alpha(A, x^0) + \frac{\beta(A, x^0)}{\omega} p \right] x^* = 0. \quad (12)$$

In this case the new form of an equivalent transfer function can be applied in eq. (12), that is the expression in brackets can be replaced by

$$\frac{k_*(A, x^0)}{1 + T_*(A, x^0)p},$$

where

$$T_* = \frac{-\beta(A, x^0)}{\omega \alpha(A, x^0)}, \quad k_* = \frac{\alpha^2(A, x^0) + \beta^2(A, x^0)}{\alpha(A, x^0)},$$

while α and β are found by eq. (9), $\beta < 0$.

We can treat similarly the forced oscillations where in the right-hand part of the non-linear system of eq. (10) we have

$$f(t) = B \sin \omega t \quad \text{or} \quad f(t) = f^0(t) + B \sin \omega t.$$

We will describe the computations with the new form of harmonic realization for the above example of a non-linear system (fig. 1).

For the kind of non-univalent non-linearity under consideration with hysteresis loops the expressions for coefficients of harmonic linearization are given by

$$\alpha = \frac{2h}{\pi A} \left(\sqrt{1 - \frac{c^2}{A^2}} + \sqrt{1 - \frac{m^2 c^2}{A^2}} \right),$$

$$\beta = -\frac{2ch}{\pi A^2} (1-m), \quad \text{at } A \geq c,$$

where A is the oscillation amplitude of the voltage u_3 .

E.g. for numerical values of a non-linear components parameters

$$h = 110 \text{ v}, \quad c = 24 \text{ v}, \quad m = 0.2,$$

the graphs of harmonic linearization coefficients are given in Fig. 2.

These graphs with eq. (5) can easily give relations of equivalent parameters k_* and ωT_* against oscillations amplitudes represented in Fig. 3.

It should be remembered that the formulae of harmonic linearization and therefore the graphs of Fig. 3 are meaningful only at $A > c$. Therefore the magnitude of T_* is bounded. Fig. 3 shows that with increase of the amplitude A the time constant T_* decreases; with the given non-linearity this is equivalent to a decreased effect of the hysteresis loops at large oscillation amplitudes.

With eq. (4) for the equivalent transfer function the dynamics of the system (Fig. 1) is given by the equations

$$\varepsilon = g_1 - g_2, \quad u_3 = \frac{k_3}{1 + \frac{1}{T_2} p} (u_2 - u_{ps}), \quad u_{ps} = k_{ps} p \Omega,$$

$$u_1 = k_1 \varepsilon, \quad u_4 = \frac{k_4}{1 + T_* p} u_3, \quad g_2 = \frac{k_5 \Omega}{p},$$

$$u_2 = \frac{k_2}{1 + T_1 p} u_1, \quad \Omega = k_4 u_4.$$

The transfer function of the part of the circuit with a feedback is given by

$$W_{c.f.} = \frac{k_3 k_4 k_*}{(1+T_2 p)(1+T_* p) + k_6 k_*}, \quad k_6 = k_3 k_4 k_{pB}.$$

The transfer function of an open-loop system will be

$$W = \frac{k k_*}{p(1+T_1 p)[(1+T_2 p)(1+T_* p) + k_6 k_* p]}, \quad k = k_1 k_2 k_3 k_4 k_5.$$

The characteristic equation of a non-linear harmonically linearized system will be given by

$$A_0 p^4 + A_1 p^3 + A_2 p^2 + A_3 p + A_4 = 0,$$

where

$$\begin{aligned} A_0 &= T_1 T_2 T_*, \\ A_1 &= T_1 T_2 + T_1 T_* + T_2 T_* + T_1 k_6 k_*, \\ A_2 &= T_1 + T_2 + T_* + k_6 k_*, \\ A_3 &= 1, \quad A_4 = k k_*. \end{aligned}$$

We can see that all coefficients of a characteristic equation are, as distinct from the previous one, positive.

Assume that the objective of further computation of the system is to obtain a non-linear transition processes performance diagram^[1] in order to obtain the transition process decay index ξ and the oscillation frequency ω for each value of the selected parameter of the system and the amplitude.

The gain of the system linear part $k = k_1 k_2 k_3 k_4 k_5$ will be taken as the parameter selected.

To solve our problem we will introduce into the characteristic equation of a harmonically linearized system the value $p = \xi + j\omega$ instead of p . As a result we will obtain an equation of the form

$$X(A, \omega, \xi, k) + jY(A, \omega, \xi) = 0.$$

Because the selected parameter k is included only into the coefficient A_4 , the parameter k will be only

in the real part of X after the real and imaginary parts have been separated. Therefore, from the equation

$$Y(A, \omega, \xi) = 0$$

we can find the relations $\xi(A)$ for different values of $\omega = \text{const.}$ and by substituting these into the expression

$$X(A, \omega, \xi, k) = 0,$$

find the relations $k(A)$ for different $\omega = \text{const.}$ This will immediately yield the lines of $\omega = \text{const.}$ on the desired performance diagram of non-linear transition processes with the system of coordinates (k, A) and the relations $\xi(A)$ at different $\omega = \text{const.}$ found before will make it possible to trace the major part of the diagram as lines of $\xi = \text{const.}$ in the same system of coordinates.

Similarly to this case we can easily handle other problems solved in practice by harmonic linearization at various non-univalent non-linearities with hysteresis loops by using the new, more rigorous notation suggested here for equivalent transfer function. To facilitate practical calculations for all specific kinds of non-linearities, formulae and graphs can be prepared in advance for new harmonic linearization coefficients $k_*(A)$ and $\omega T_*(A)$ similar to those that exist now for the coefficients $\alpha(A)$ and $\beta(A)$. This can also be done in more complex cases which include an additional relation with the ratio $\frac{\xi}{\omega}$ or the constant constituent x^0 .

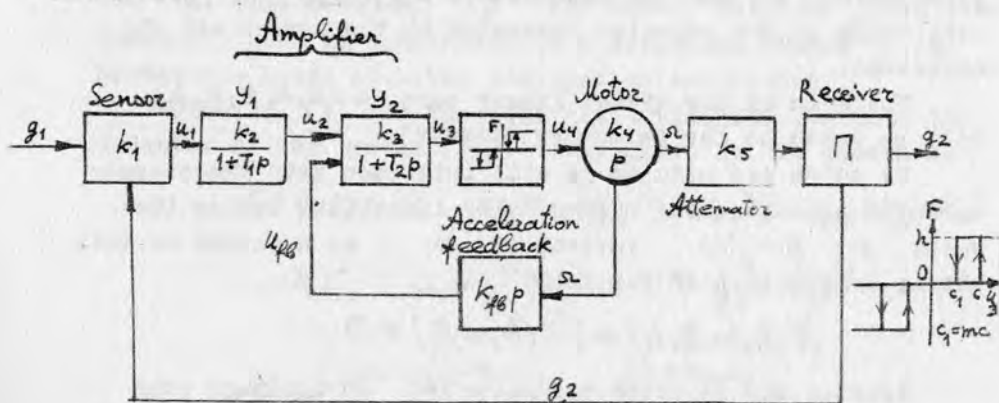


Fig. 1 System block-diagram

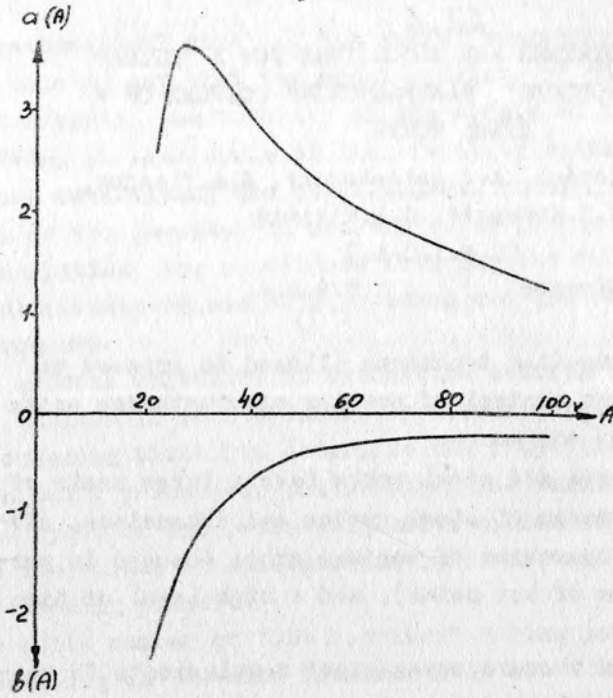


Fig. 2 Harmonic linearization coefficients

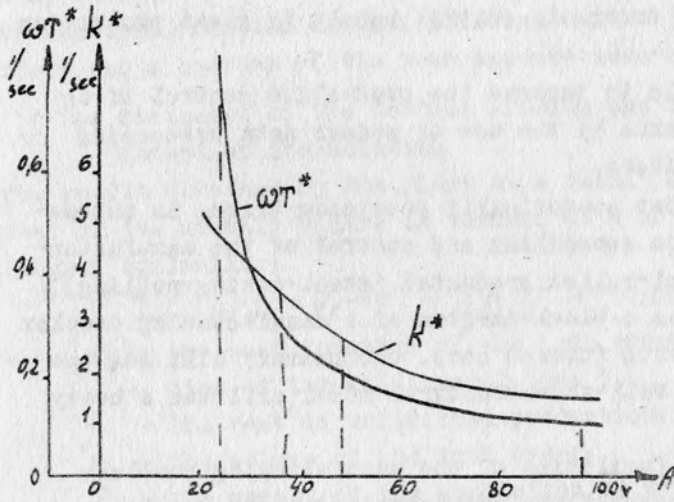


Fig. 3 Equivalent values of gain and time constant

CONTROL SYSTEMS AND ALGORITHMS FOR A "STEEL-ROLLED PRODUCTS" MANUFACTURING COMPLEX OF A STEEL WORKS

A.P.Kopelovich, A.A.Belostotsky, B.A.Vlasjuk,
V.M.Khrupkin, G.I.Nikitin
(C.R.I.C.A.)
Moscow U.S.S.R.

Progress in computing technique allowed to proceed to automatic production control of complex manufacturing units such as large shops and works.

Large modern iron and steel works have a large scale of production a wide range of steel grades and dimensions, close interaction of operation of various shops (caused in particular by the flow of hot metal), and a high level of disturbances.

These and other factors cause great requirements to control and planning systems for an iron and steel works: complex problems and large size problems must be solved in short periods of time, errors in control result in great production losses.

It is possible to improve the production control of an iron and steel works by the use of modern data processing means, i.e. computers.

One of the most economically promising trends is automation of production scheduling and control of the manufacturing complex "steel-rolled products" (steel-making-rolling)¹.

Fig.1 contains a block diagram of a manufacturing complex with one open hearth furnace shop, one primary mill and two jobbing mills (a rail-and-structural steel mill and a heavy section mill).

All the main facilities of the manufacturing complex "steel-rolled products" are characterized by the cyclic mode of production process. Every cycle of a facility operation corresponds to the processing of a certain production unit - a discrete portion of metal (a cast, an ingot, a bloom).

Lack of space prevents the authors from describing a de-

tailed mathematical model of the works' operation. It will be only pointed out that the model consists of three main interrelated parts: the totality of the models of the processes occurring at each cycle of each facility operation; the conditions determining the metallurgical routing and the constraints on the sequence of product units processing on the works facilities; the conditions ensuring the fulfillment of the predetermined volume of production and the variety of steel products.

The general objective of production control is to choose control actions in each process, to schedule the operation of each processing plant (to determine the properties of the products being processed, the starting and termination points for each cycle of the processing units operation) and to determine the volume of discrete portions for each cycle of each processing plant.

The whole number of "the machines" taking part in the metal processing in the area considered is about 200; the number of product units (the casts processed in the area during a month) reaches 2000. The plant operates according to the orders received from the central planning office; the number of orders per a quarter of the year amounts several thousands.

The Statement of the Control Problem and the Method of its Solution

The profit obtained by the plant as a result of the fulfillment of its monthly orders is assumed as a criterion of the control optimality;

$$\sum_i C_i(t_i) \cdot y_i + y_T - K T - \sum_j b_j x_j \rightarrow \max \quad (1)$$

where $C_i(t_i)$ - the cost function of the i -th order as a function of its fulfillment time t_i ;

y_T - the cost of unfinished production;

y_i - the volume of the i -th order;

T - the moment of the completion of the total number of orders;

K - the cost of the time unit of the plant operation;

x_j - the amount of the j -resource used;

b_j - the j -resource cost.

The first term of the criterion (1) is the cost of the whole set of orders taking into account its decrease when the orders are not fulfilled by the proper time. The second criterion term contributes to the production process in the preceding areas of the metallurgical routing with the predetermined volume of the finished product units for the plant. The third and the fourth terms require minimization of the idle time of the production facilities and of the resource used.

The analysis of the control problem for the plant shows that it does not belong to any known class of the extremum problems; it involves the theory of combinations, elements of variational problems and the problems of mathematical programming. To solve accurately a problem of such complexity and of such an enormous size and to correct it according to its disturbances would require nonexistent computing devices.

The main ways of overcoming the "damnation of size" are: making the most use of the specific nature of the particular problem, i.e. the simplifying features following from the controlled plant structure and the optimality criterion; and the employment of approximate solutions assuming as a basis the previous management experience (ranking of priorities, reasonable restriction of a great number of variants etc.)

The simplifying features of the problem arise due to the fact that the ties between different areas of great manufacturing complexes give certain freedom of choice in performing separate operations. Besides there is some freedom of control within a shift, twenty four hours and so on, i.e. within separate time intervals of the total period of scheduling.

The specific features of the production control system are represented mathematically by a specific structure of the system of equations and inequalities describing the controlled plant, e.g. in a block character of this system. This allows to apply special computational algorithms which are more efficient than general methods of solution; decomposition methods of solutions of the linear programming block problems^{2,3}.

local algorithms of integer problems solutions⁴, dynamic programming of Markovian processes⁵, sequential analysis scheme of variants⁶, etc.

To make use of such computational algorithms it is necessary to represent the controlled plant model structure and to select an efficient structure of the problem solution for the given model.

It is reasonable to define the structures of the model and of the solution in several steps⁷. The first step is to decompose the controlled plant into comparatively few parts or blocks. The totality of the blocks and their interdependences (the conditions combining the variables of separate blocks) represents the controlled plant structure at the first step.

Analysis of this structure allows to obtain the efficient structure of the problem solution, i.e. a set of the control subproblems and their interdependence in the process of solution.

At the second and subsequent stages, the control subproblems obtained at the preceding stages are examined; the structure of the controlled plant of each subproblem is detailed and the solution structure selected.

The works operation during a month is represented as a complex multistep process. The general model of the works operation is decomposed into blocks in accordance with the separate manufacturing areas and the time subintervals of the total planning period. A separate block of overall structure describes the work of one manufacturing area within one time subinterval.

It is reasonable to select the size of one manufacturing area covered by one block and the duration of the planning subintervals so that, firstly, the total number of the blocks were not too great allowing to derive the most efficient succession of the subproblems solution, and, secondly, the number of the parameters connecting the separate blocks were as few as possible.

The general structure of the controlled plant analysed

at the first decomposition step is given in Fig. 2, where a block is designated with two indices: K - the number of the time subinterval, and P - the number of the manufacturing area. Here $P = 1$ - the open-hearth furnaces area; $P = 2$ - the blooming mill area; $P = 3$ - the first N 2 jobbing mill area; $P = 4$ - the second jobbing mill area.

Each block has ties of two types: the "static" ties with the other block (area) within one planning time subinterval, and the "dynamic" ties - with the blocks describing the same manufacturing area but in other subintervals.

Decomposition of the control problem comprises the formulation of the local control subproblems of the separate blocks and determination of the method of coordination of the separate subproblems solutions with the aim to optimize solution of the general control problem.

The local control subproblems of the blocks can be derived using the optimality principle of dynamic programming: "any process occurring between two fixed finite points must be controlled optimally", i.e., at the fixed input and output coordinates of the given block, only those control variants might be optimal which provide the maximum block efficiency.

The local control subproblems are the problems of the production control; they involve the determination of the mode and succession of operations on the plants of the given area for the present initial states of all the plants, predetermined delivery schedules of the product units in the area, and at the predetermined time of processing completion of the plants.

The coordinating control subproblems are the problems of production scheduling; they involve determination of the above variables of the local control subproblems. A great size of the general coordination problem requires its decomposition into separate subproblems in each of which only a part of the ties among the blocks is solved.

Each succession of the block ties resolution may be defined by respective expenditures for the problem solution us-

ing the computer (required memory capacity, duration of solution etc.). For the problems in question, the solution expenditures depend on the number of parameters. It is reasonable to select such succession of the blocks ties resolution that the total number of parameters in all parametric coordinating subproblems would be minimal.

The structure of the solution of the steel works control problem based upon the above methods is given in fig.3. Here blocks 1-9 are the production scheduling system; blocks 10-13 - the open-hearth furnace shop control system; blocks 14-16 - the control system of conveying metal to the soaking pits; blocks 17-20 - the primary mill control system; and 21-24 are the jobbing mills control system.

The main subsystems of the "steel-rolled products" manufacturing complex for different hierarchial levels being designed at The Central Research Institute of Complex Automation are given below.

The Production Planning System

One of the most promising (so far as economy is concerned) trends of the modern high speed computers use is automation of the production planning of an iron and steel works operations.

The general iron and steel production planning problem is the following: to make optimal schedules of the main processing units operations based on the totality of orders from consumers, the current state of the processing units and unfinished processing⁸.

The optimal schedule of processing units operation is such a schedule which is firstly made in accordance with the operation schedules of other processing units and areas and secondly it minimizes the plant losses through processing units idle time, undue delivery dates, hot metal cooling, etc.

Various disturbances require rather frequent corrections of the schedules formed earlier. That is why the planning automation was considered as developing the planning system, including a system of production data collection and trans-

mission, the systems of algorithms with which the optimal schedules calculations and corrections are made and finally the systems of schedules transmission to the personnel at the production areas.

The plant operates according to monthly orders, which the central planning department sends to the plant once a quarter; the number of orders per a quarter totals several thousands. The orders are mainly for structural shape but the plant can also deliver unfinished products - steel ingots, blooms.

The quality of all the schedules was estimated according to the optimality criterion for the whole plant. The profit obtained by the plant as a result of fulfillment of the monthly book of orders was assumed as an optimality criterion.

The general problem of finding optimal schedules of all the production units operation on the basis of the monthly book of orders has a great number of variables which are connected by complex and numerous relations. That is why the centralized solution of such a problem is impossible even using modern computers and decomposition of the general problem into a set of interrelated subproblems of smaller size is necessary.

In performing this decomposition the choice of planning algorithms system is made on the basis of the following principles; the algorithms structure should be hierarchical; the optimal algorithms structure should minimize the total expenditures on the problem solution; the input data which a higher level of planning works out for a lower level should be either certain production units schedules or the prices of unfinished products.

The structure obtained envisages the use of the following basic algorithms:

- 1) Making up a set of lots, i.e. the totality of orders rolled on jobbing mills having one set of rolls. The problem is formulated as a combinatorial one solved by means of a heuristic algorithm based on the priority rule.
- 2) The determination of lots sequence rolled on jobbing mills. An algorithm combining dynamic programming and direct-

ed selection is developed to solve it.

3) Scheduling of the manufacturing complex "steel-rolled products" operations. Daily orders of open-hearth furnace steel production of rolling on the blooming mill and jobbing mills, of metal dispatching are determined. A block method of linear programming is used for solving this problem.

4) Daily scheduling of an open-hearth furnace shop: the distribution of steel grades according to furnaces and cast numbers on each furnace within a day (twenty four hours). The problem is solved as an integer linear programming problem.

5) Daily scheduling of cast processing on a blooming mill. The algorithm is formed as a result of formalization of the rules employed by the plant personnel.

6) Daily scheduling for jobbing mills. The algorithm is based on heuristic rules.

The efficiency estimation of the foregoing algorithms of scheduling showed that the system allowed to increase the productivity of jobbing mills, of the open-hearth furnace shop and the blooming mill by 2.5-3%.

To solve the problems the planning system uses computers of the iron and steel plant computing center. The production progress data are fed by means of teleprinters, placed in the open-hearth furnace shop, on the blooming mill and jobbing mills. The data are printed on the printers in the computing center and also directly for the controllers of the plant.

The Open-Hearth Furnace Shop Control System

The automatic control system is designed for the shop having a stockyard, a mixer, several open-hearth furnaces, a furnace bay and a teeming bay¹. The shop control problem is stated as a problem of the shop relative profit maximization:

$$\Pi = \sum_j c_j g_j - \sum_i \kappa_i \tau_i - a \sum_i g_i \rightarrow \max \quad (2)$$

where g_j - the production volume of the j -th steel grade;
 τ_i - the time spent by the i -th furnace for production;
 g_i - the consumption of the technological oxygen by the i -th furnace;

- C_j - the relative cost of the j -th steel grade;
 K_i - the time unit relative cost of the i -th furnace operation;
 α - the relative cost of technological oxygen.

The adopted open-hearth furnace shop criterion form is flexible enough and a proper choice of factors may reduce it to various particular forms.

The linear nature of the criterion is very convenient to solve certain optimization problems.

The control system uses the following algorithms:

1) Current scheduling of open-hearth furnaces: for each cast on the scheduling interval a certain oxygen consumption is calculated; the totality of the scheduled casts in all the furnaces forms a sequence of demands for the maintenance of the furnaces with auxiliary equipment. The schedule of meeting these demands is made for each kind of equipment; in the equipment assignment checking up the fulfillment of the schedule for the furnaces maintenance with the earlier assigned equipment is performed through a designated for each cast numerical parameter - "activity slack" (the difference between the moment of service demand arrival and the moment of releasing the equipment designated for servicing this demand).

After the assignment of all kinds of equipment reassignment of oxygen is made in such a way that the total oxygen consumption per shop at any moment would not exceed the predetermined value. The oxygen consumption of the cast period with negative activity slack increases and it decreases for periods with positive activity slack. After the oxygen reassignment the auxiliary equipment assignment is made again. In accordance with the furnaces operation ^{the schedule} schedule of charge conveying to furnaces and trains arriving for teeming is calculated, the schedule of other technological operations in furnaces is calculated as well.

2) The optimal strategy determination for lack of coincidence of steel teeming while chargings in adjacent furnaces coincide; the strategy consists in determining the sequence of chargings and cast weights for every possible interval between

the chargings starting points on adjacent furnaces. The problem is formulated as a problem of optimal control search for the controlled Markovian chain and it is solved through the modification of Howard's method.

3) The auxiliary equipment control; the problem is in assignment of the available equipment (a charger, equipment of the teeming bay, etc.) to fulfil certain operations; the problem is also in determining the next charging of cast iron to furnaces. The equipment distribution is accomplished according to priorities of furnaces based on "activity slack".

4) The production progress control; it is performed by means of comparing the schedules of operations based on the algorithms with their actual fulfillment and the production recording.

The automatic system is developed for the centralized control shop. All the production progress information is concentrated at the shop control point. The production progress information is fed automatically and manually. The traffic control (trains and locomotives) is performed by means of relay centralization of the route and train information devices providing information about the progress and N of the trains and locomotives.

The heart of the system is a computer. The computer is to operate in real time providing multiprogram operation with interruptions for information input - output. The computer speed is about 30 000 short operations/sec; the memory capacity is about 20 000 of 24-bit words.

The benefit obtained from the implementation of the system is achieved as a result of increasing the shop productivity. Besides, the gain is expected from reduced fuel consumption and reduced ingots heat losses thanks to more uniform cast yield from the shop.

The Control System of Conveying Metal to the Soaking Pits

The control system of conveying metal to the soaking pits is intended for that area of the works which covers the teeming bay, stripper bay, mould preparation shop and cold ingots

stock. The peculiarity of this area is the closest interaction of technological and transport operations⁹.

The functions of the control system under consideration include:

1) The organization of conveying the trains with casts to soaking pits according to the requirements of the central planning system.

The basic criterion in organizing the traffic is the minimization of deviations from the order of casts arrival; the necessity of making up trains of moulds in the mould preparation shop in due time must be taken into account. A real traffic situation is proceeded from: the availability of empty transit rails, the possibility of the locomotive approach to the train or the necessity of shunting, the availability of vacant locomotives, etc.

The criteria of each particular control problem solution (the choice of locomotives, optimal routing, etc.) are determined by the specific features of rails in the controlled area; however, the system always tries to ensure the train movement for the shortest time.

2) Forming a sequence of trains directed to the mould preparation shop for making up. The criterion for this problem solution is the minimization of the deviation of the train arrival time for teeming from the scheduled starting point of teeming in the open-hearth furnace shop.

3) The optimal distribution of the locomotives' work which may be considered as a transport linear programming problem. The criterion for solving this problem is the minimization of the time when the trains wait for locomotives and idle running of locomotives.

The solution of these problems results in making up a schedule for the area for a certain period of time. Partial or full corrections of the schedule are made for every change of a plant state.

According to the schedule the traffic automatic devices at the controlled area receive at certain moments instructions about the automatic choice of the route (switching on

proper railway points and lights, switching off the route which was completed); these traffic instructions are transmitted to locomotives.

The instructions about the order of making up and stripping are transmitted to locomotive drivers and operators in the stripper bay and mould preparation shop. The instructions about making up trains with cold ingots of certain kind are transmitted to the cold ingot stock.

The automatic traffic devices are the source of information about the traffic situation in the controlled area.

To get information about the change of situation at an area sensing devices are installed in certain points of the rails. They automatically transmit a passing train or locomotive number and also their direction into automatic traffic devices. It is possible to follow the trains and locomotive movement through the area with the aid of such devices. The numbers of the trains and locomotives are transmitted from one sensing device to the next one following the movement of the locomotives.

The data about the progress of certain technological operations (the starting and finishing points of making up trains, stripping, etc.) are fed by operators through keyboards.

The whole system is expected to increase the temperature of the metal arriving at the soaking pits and to reduce the blooming mill idle time.

REFERENCES

1. М.Д.Климовицкий, А.П.Копелович, Автоматический контроль и регулирование в черной металлургии. Изд-во "Металлургия", 1967 г.
2. G.B.Dantzig, Ph.Wolfe, "Decomposition principle for linear programs", Oper.Res., 1960, vol.8, N 1.
3. Е.Г.Гольштейн, Д.Б.Юдин, Новые направления в линейном программировании. М., 1966 г.
4. Ю.М.Журавлев, Ю.Ю.Финкельштейн, Локальные алгоритмы для задач целочисленного программирования, "Проблемы кибернетики", Вып.14, 1965.
5. Р.Беллман, С.Дрейфус, Прикладные задачи динамического программирования, Изд-во "Наука", М., 1965.
6. В.Михалевич, Б.Шор, Метод последовательного анализа вариантов для численного решения задач оптимизации, "Труды конференции по вопросам применения ЭВМ в народном хозяйстве", г.Горький, 1964.
7. Г.И.Никитин, Выбор структуры автоматизированной системы управления металлургическим заводом, "Труды ЦНИИКА", вып.19, 1968.
8. Б.А.Власюк, А.П. Копелович, Г.Р.Кюсснер, Применение ЭВМ для оперативного планирования прокатного производства, "Приборы и системы управления", № 6, 1967.
9. А.А.Белостоцкий, Ю.С. Вальденберг, Система оперативно-диспетчерского управления. "Импульс" для участка металлургического комбината, в сб. "Управление производством", Изд-во "Наука", 1967.

FIGURES

Fig.1. The works structure:

- 1) Area 1; 2) Area 2; 3) Area 4; 4) Area 3; 5) Open-hearth furnace shop; 6) Dispatch; 7) Purchased ingots; 8) Purchased blooms; 9) Dispatch; 10) Stockyard; 11) Furnace bay; 12) Teeming bay; 13) Stripper bay; 14) Soaking pits; 15) Blooming mill; 16) Blooming yard; 17) Reheating furnaces; 18) Mill N 2; 19) Mill N 1; 20) Finishing; 21) Cold ingots stock.

Fig.2. The works model structure.

Fig.3. The steel works control structure:

- 1) Scheduling of lots sequence at the jobbing mills;
- 2) Making up lots at the jobbing mills; 3) Monthly production scheduling; 4) Scheduling for the open-hearth furnace shop; 5) Scheduling of casts processing; 6) Scheduling of ordered lots processing sequence; 7) Scheduling of lots processing sequence; 8) Scheduling of ingots processing sequence; 9) Scheduling for N1 jobbing mill;
- 10) Scheduling of open-hearth furnace operating time;
- 11) Control of the stockyard equipment operation;
- 12) Control of the furnace bay equipment operation;
- 13) Control of the teeming bay equipment operation;
- 14) Traffic control; 15) Cold ingots stock control;
- 15) Cold ingots stock control; 16) Stripper bay control;
- 17) Soaking pits control; 18) Ingot cranes control;
- 19) Roughing mill control; 20) Shears control; 21) Blooms stock control of N 1 mill; 22) Reheat furnaces control;
- 23) Compartment furnaces control; 24) Jobbing mill control.

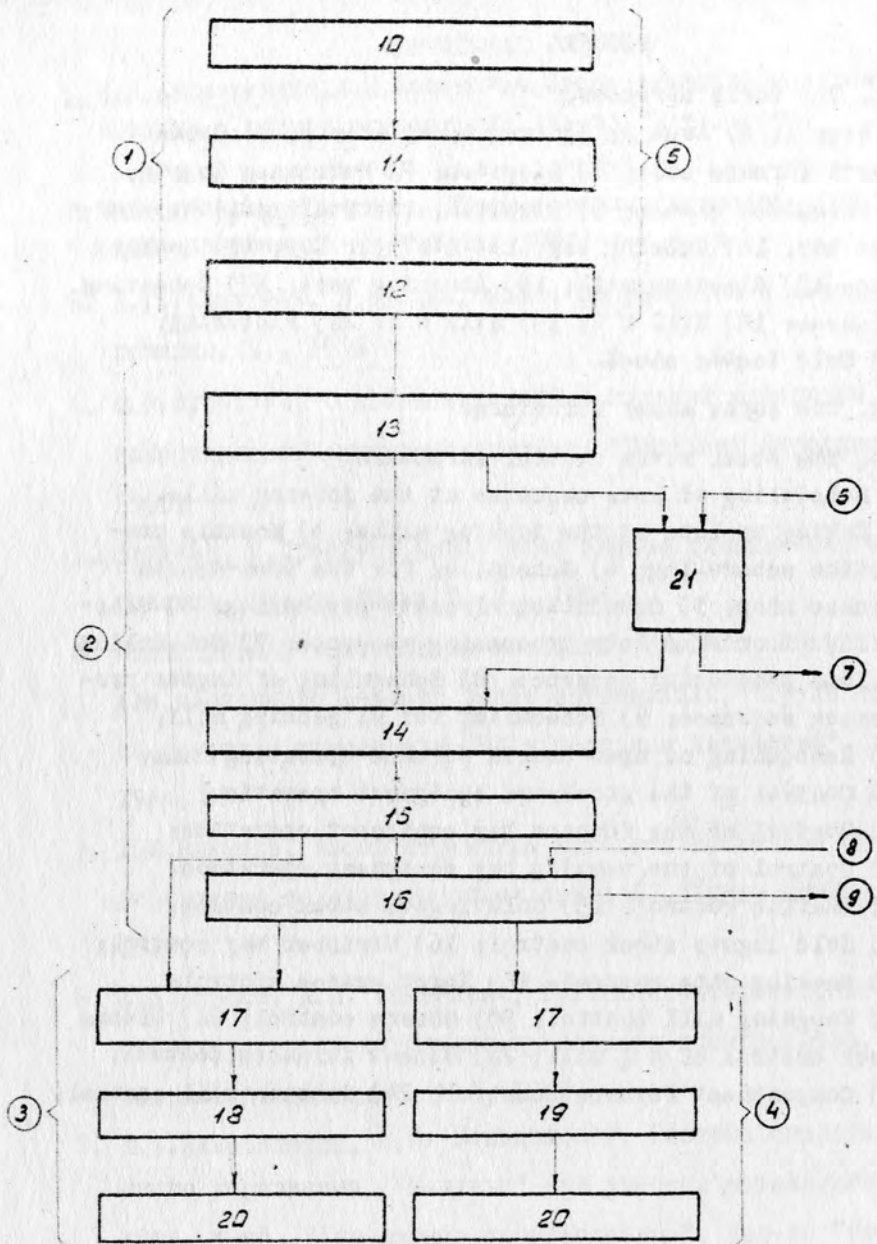


Fig. 1

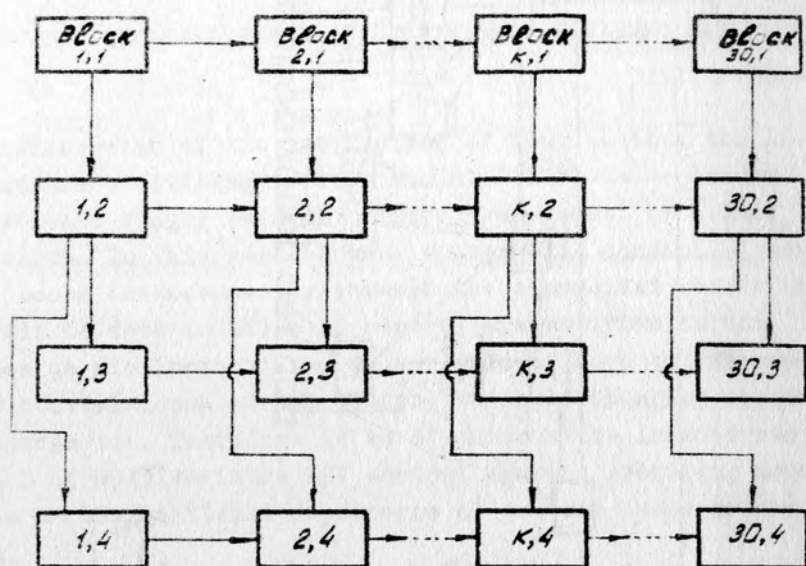


Fig. 2

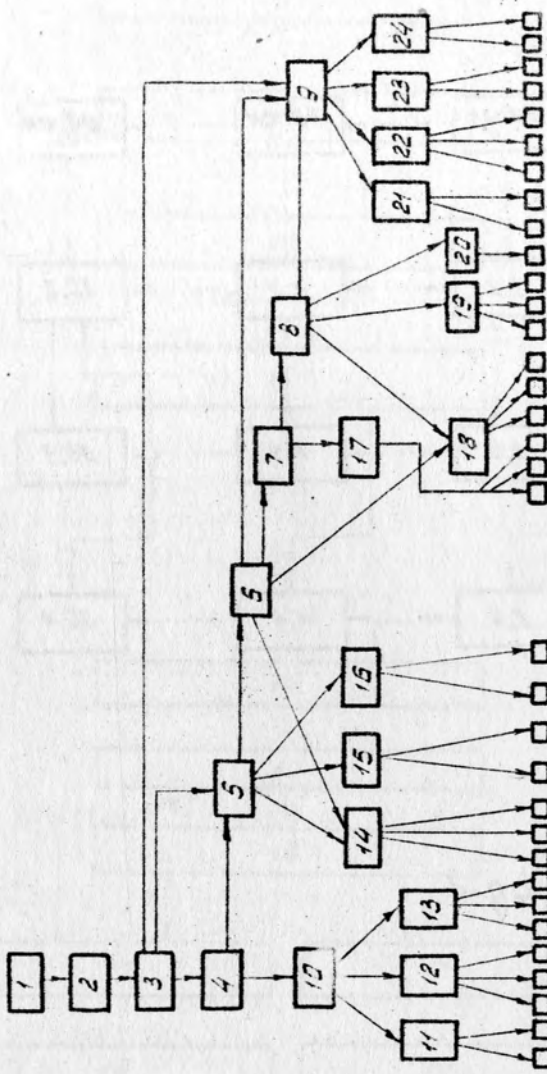


Fig. 3

CONSIDERING SYNTHESIS OF LIFTING REENTRY VEHICLE
CONTROL SYSTEM STRUCTURES IN ATMOSPHERIC MANEUVERPETROV B.N. KOLPAKOVA N.P. VASILYEV V.A. PAVLENKO A.I.
(MOSCOW)

Consideration of the interaction of longitudinal and lateral motions of lifting reentry vehicle (LRV) is necessary for hypersonic flight even at comparatively small ^{angles} of attack and sideslip. In this case control system will consist of many control loops interconnected through the controlled member LRV that leads to deterioration of control system dynamics and sometimes to its instability. As a result of such an interaction of control loops at one control variable changing the others would change too. Therefore it is of interest to investigate the set G of multivariable LRV control systems providing independence or insignificant dependence of control loops or groups of them.

In this paper the linear multivariable control system where the control loops interact through LRV have been considered and the problem of obtaining the set $G \{M, S, U\}$ of structures providing selective invariance has been considered to select the best structure in the sense of control quality and simplicity of implementation.

I. EQUATIONS OF THREE-DIMENSIONAL LRV MOTION

LRV is a vehicle with comparatively high hypersonic L/D . There LRV can fulfil an effective aerodynamic maneuver [1] owing to lift force arising in its atmospheric flight. Such flight regimes at great bank angle (as high as 90°) are typical for LRV maneuver; then new problem for control system have arisen.

In the ~~of~~ LRV atmospheric flight with almost circular velocity it is impossible to neglect of Coriolis and centrifugal force as in the case of aircraft. Therefore to compute LRV trajectory we can use a set of LRV centre of mass motion equations given in velocity semi-state coordinates considering

curvature and rotation of the Earth [2].

Adding the equations describing LRV rotating motion, cinematic equations of LRV centre of mass motion converted at geographic coordinates, cinematic equations describing LRV rotating motion about Earth axes and a number of geometric correlations connecting angles at velocity semi-state coordinates with angles at vehicle state coordinates we get a set of 15 nonlinear equations.

Owing to the nonlinearity of these equations exact parameters of LRV motion can be calculated only with methods of numerical integrating or digital computers. At the same time at the first stages of the control system synthesis we are mainly interested in analytical estimations which allow to investigate LRV dynamics.

For that purpose let's simplify the received set of equations by linearizing them about some nominal regime. But it should be noted that it is impossible to use the known techniques of separating of LRV longitudinal and lateral motions to investigate its dynamics as LRV fulfils aerodynamic maneuver with large angles of attack and usually large bank angles, which leads large angles of sideslip.

After linearizing the set of equations is resolved as to the most important controlled variables.

The set of linearized equations of LRV three-dimensional motion is as follows:

$$\left\{ \begin{aligned} V &= a_{11}\alpha + a_{12}\beta + a_{13}\gamma + a_{14}\psi + a_{15}\varphi + a_{17}H + a_{18}\delta_{sp} + a_{19}\delta_B + a_{110}\delta_H, \\ \psi &= a_{21}\alpha + a_{22}\beta + a_{23}\gamma + a_{25}\varphi + a_{26}V + a_{27}H + a_{28}\delta_{sp} + a_{29}\delta_B + a_{210}\delta_H, \\ \varphi &= a_{31}\alpha + a_{32}\beta + a_{33}\gamma + a_{34}\psi + a_{36}V + a_{37}H + a_{38}\delta_{sp} + a_{39}\delta_B + a_{310}\delta_H, \\ \gamma &= a_{41}\alpha + a_{42}\beta + a_{44}\psi + a_{45}\varphi + a_{46}V + a_{47}H + a_{49}\delta_B + a_{410}\delta_H + a_{411}\delta_T, \\ \beta &= a_{53}\gamma + a_{54}\psi + a_{55}\varphi + a_{56}V + a_{57}H + a_{510}\delta_H + a_{511}\delta_T, \\ \alpha &= a_{62}\beta + a_{63}\gamma + a_{64}\psi + a_{65}\varphi + a_{66}V + a_{67}H + a_{69}\delta_B, \\ H &= a_{71}\alpha + a_{72}\beta + a_{73}\gamma + a_{74}\psi + a_{76}V \end{aligned} \right. \quad \dots (1)$$

where: Q_{ij} - operator factors; V - velocity deviation;

ϑ - angle of pitch deviation; ψ - yaw deviation; γ - bank angle deviation; β - angle of sideslip deviation; α - angle of attack deviation; H - altitude deviation; δ_p - changing of thrust controller position; δ_e - deflection of an altitude control; δ_h - deflection of an yaw rudder; δ_a - deflection of ailerons.

We can see the interaction between the variables $V, \vartheta, \psi, \gamma, \beta, \alpha, H$ from the equations (I). The control system with four loops which would be connected with others through the controlled member is composed of four controllers and four controlled variables. Obviously it would be difficult to develop the control system with no additional connections providing independence or insignificant dependence of control loops.

Therefore a set of structures in the form of connected systems will be considered using the algorithms for composing a set of selectively invariant systems. The set of combined systems are not considered here.

2. GRAPH-A MATHEMATICAL MODEL OF MOTION EQUATIONS. SYNTHESIS OF LRV CONTROL SYSTEM.

The graphs with no loops are chosen as the basis of structural representation of LRV control systems are chosen to make easy the investigation of interconnections of LRV controlled variables and to allow choosing the controllers under the conditions of their most efficiency in control and also determining the set $G\{M, S, U\}$ of selectively invariant systems.

Then using equations (I) compose a graph-a mathematical model of LRV three-dimensional motion in earth atmosphere as of a controlled member (Fig.1).

Take V, γ, α as LRV controlled variables. So far as angle of sideslip β can be large at some LRV flight regimes and would greatly influence LRV three-dimensional trajectories it is necessary to consider β as a controlled variable too.

Then separate the graph $G(x, r)$ (Fig.2) of controlled variables from the graph of LRV three-dimensional motion (Fig. 1) and consider a set of possible selectively invariant

systems structures.

All the kinds of connections providing selective invariance of multivariable systems have been considered in detail [4]. It has been shown that autonomy in the sense of Voznesensky I.N. or invariency as to no own control commands or disturbances and both of them at the same time are possible depending upon the types of connection.

Let us begin consideration of the set of LRV control system structures with the set M in the form of great-trees. A number of the great-tree variances for a graph whose root has been chosen as bank angle γ can be determined with the theorem III (Appendix).

$$\Delta_{\gamma} = \begin{vmatrix} 2 & -1 & -1 \\ -1 & 3 & -1 \\ 0 & -1 & 3 \end{vmatrix} = 12$$

All the structures of the set M with bank angle γ as a root are shown in Fig.3. In a similar way we can determine a number of the great-trees for V ($\Delta_V = 12$) and α ($\Delta_{\alpha} = 16$), β ($\Delta_{\beta} = 8$). Thus all the set M of control system structures in the form of great-trees would be 48. set S of the structural graphs. As well as before

Then let us begin consideration of set S with root γ will be $\Delta_{\gamma} = 3$.

The structural graphs for γ have been presented in Fig.4. A number of connections necessary for constructing the control system in the form of the structural graph is equal to $K=6$.

The investigation of the set S of structures gives large possibilities for choosing LRV control system satisfying to the requirements presented.

And at last let us consider the set L of autonomous groups of loops. Four versions of construction of autonomous groups are possible. ~~In the form of the symmetric~~. In this case a number of connections K will be:

- a) for groups in the form of the symmetric subgraphs $K=6$,
- b) for groups in the form of the great-trees $K=9$,

c) at triangular (or quasi-triangular) system matrix and in the presence of groups $K=3$.

Some structures in form of the autonomous groups or the symmetric subgraphs are presented in Fig.5.

It should be noted that in more complex cases when a number of controlled variables > 4 it is possible to construct an algorithm of looking all the set $G\{M, S, L\}$ of structures over and to obtain with digital computer all structures satisfying the requirement of implementation and given quality of control.

As an example of the above suggested method of synthesis of LRV control system structures the synthesis of control system of hypothetical LRV has been made.

For the purpose of illustration the particular regime of LRV flight (namely thrust flight at constant altitude in the plane of minor circle) has been considered (Fig.6.)

Analysis of values of operator-factors Q_j shows that for this regime of flight some of them can be neglected because of their insignificance and then the controlled member is simplified to the graph presented in Fig.7. The reason for such a simplification have been proved with the transients presented in Fig.8. and computed by digital computer as a result of solving LRV linearised equations (I) at impulse changing of controllers $\delta_{\varphi}, \delta_H, \delta_\beta, \delta_\gamma$.

Comparison of obtained control system structures allows to conclude that the structure (Fig.9.) in which additional connection (dotted line) is required seems to be preferable.

Peculiarity of the given structure is that the angle β is autonomous and invariant as to angles γ, α and flight velocity V . In this case the influence of the angle of sideslip β on the angles γ and α is simply a disturbance, whose calculation is very easy. The known disturbances influencing the velocity control loop will be those with the changing of the angles α, β, γ .

For example Laplace transformation of changing angle with disturbance of β loop will be as follows:

$$\Delta(s) = a_{\alpha\beta}(s) a_{\beta\gamma}(s) \cdot f_{\beta}(s),$$

where: $f_{\beta}(s)$ -disturbance in β control loop;

$a_{\beta\gamma}(s)$ -transfer function of β control loop at disturbance $f_{\beta}(s)$.

In such a way Laplace transformations of controlled variables can be written corresponding to the graph presented in Fig. 9.

Let us take the following sequence of controllers synthesis assuming that α, γ and flight velocity controllers can be chosen without taking into consideration the connections through the controlled member. Obtain the transfer functions of the controllers $W_{\alpha}(s), W_{\beta}(s), W_{\gamma}(s)$ for V, α, γ control and β stabilization on the basis of desired dynamics as follows:

$$\Phi_{\kappa}^{y_i}(s) = \Phi_i^{y_i}(s), \quad (3)$$

where: $\Phi_{\kappa}^{y_i}(s)$ -transfer function of closed loop, corresponding the desired dynamics;

$\Phi_i^{y_i}(s)$ -transfer function of synthesised control loop.

From the equation (3) we can obtain formula of the controller transfer function for control regime:

$$W_i(s) = \frac{1}{W_{ii}(s)} \cdot \frac{A_i(s)}{[B_i(s) - A_i(s)]}, \quad (4)$$

where: $\Phi_{\kappa}^{y_i} = \frac{A_i(s)}{B_i(s)}$; $W_{ii}(s)$ -transfer function of control loop for control command y_i .

Synthesis of transfer function $W_{\beta}(s)$ will be made taking into account interconnections determined with the transfer functions of controlled variables α, γ connections.

From the graph presented in Fig. 9 the controller

transfer function $W_y(s)$ according to ⁽⁴⁾ will be:

$$W_y(s) = \frac{1}{W_{yy}(s)} \cdot \frac{A_y(s)}{[B_y(s) - A_y(s)]} [1 - W_e(s)], \quad (5)$$

where: $W_e(s) = \frac{W_{ay}(s) W_{ya}(s)}{1 + W_{ad}(s) W_{\alpha}(s)}$

The coordinate transients can be computed from the graph presented in Fig.9 at control loops interaction.

For example

$$y_{\beta}(t) = \frac{1}{2\pi j} \int_{-j\infty}^{+j\infty} \bar{\Phi}_{\beta\beta}(s) \beta_f(s) e^{st} ds,$$

where:

$$\bar{\Phi}_{\beta\beta}(s) = \frac{A_{y\beta}(s) + A_{y\alpha}(s) A_{y\alpha}(s)}{1 - A_{y\alpha}(s) A_{\alpha y}(s)}; \quad A_{y\beta}(s) = \frac{W_{y\beta}(s)}{1 + W_{y\beta}(s) W_{\beta}(s)}; \quad A_{\alpha y}(s) = \frac{W_{\alpha y}(s)}{1 + W_{\alpha y}(s) W_{\alpha}(s)};$$

$$A_{y\alpha}(s) = \frac{W_{y\alpha}(s)}{1 + W_{y\alpha}(s) W_{\alpha}(s)}; \quad A_{\alpha y}(s) = \frac{W_{\alpha y}(s)}{1 + W_{\alpha y}(s) W_{\alpha}(s)}; \quad \beta_f(s) = \frac{W_{\beta f}(s)}{1 + W_{\beta f}(s) W_{\beta}(s)} f_{\beta}(s)$$

f_{β} - disturbance;

$$v_{\beta}(t) = \frac{1}{2\pi j} \int_{-j\infty}^{+j\infty} \bar{\Phi}_{\beta\beta}(s) \beta_f(s) e^{st} ds,$$

where: $\bar{\Phi}_{v\beta}(s) = \frac{W_{v\beta}(s)}{1 + W_{v\beta}(s) W_{v}(s)}$

In a similar way other coordinate transients at control commands and disturbances $\alpha_{\beta}(t), \alpha_y(t), \alpha_y(t)$ can be computed.

For example the coordinate transients have been computed with digital computer at control command y_v . As can be seen from Fig.10 the flight velocity transient satisfies to control quality requirements presented.

CONCLUSIONS

In this paper it has been shown that:

1. LRV in aerodynamic maneuver is a complex controlled member with connections of all coordinates. Therefore it is reasonable to consider a set of control system structures as multivariable systems providing independence or insignificant dependence of control.

2. Using the algorithms for composing the set of selectively

invariant system structures and decomposing a control system with the graph theory allows to choose the best structure in the sense of the quality of control and simplicity of implementation.

A P P E N D I X .

Some theorems used in this paper and proved in [5] are presented:

Theorem I. If a graph contains arcs of a loop the determinant $|S_n^a|$ of the incidence matrix for the arcs would be equal to zero.

Autonomous systems.

Let us consider the autonomous systems in sense of Voznesensky I.N.; first we shall transform the linear composition matrix of graph arcs and apexes as follows:

$$\left| S_i^j \right| = \begin{cases} -1 & \text{at } i \neq j & (x_i, x_j) \in u \\ 0 & \text{at } i \neq j & (x_i, x_j) \notin u \\ 1 & \text{at } i = j \end{cases} \quad (I)$$

Theorem II. A graph is autonomous when the square matrix determinant $|S_n^a|$ is equal to I at $[x = \phi(\text{directly})]$. A number of connections necessary for the formation of this graph is

$$K = n(n-1)$$

Structural graph.

Define the structural graph with a root x_i [4] :

- 1) The structural graph is one $G(x, \Gamma)$ with a root x_i
- 2) The graph $G(x, \Gamma)$ has no loops .
- 3) No apex (root x_i) is included in the arc.

Let's begin with the consideration of the structural graphs which often appear to be the great-trees. Define a great-tree:

- 1) a great-tree is a graph with a root $x_i \in u$,
- 2) the graph (x, u) does not include loops,
- 3) each graph apex includes only one arc,

- 4) no graph root x_i includes an arc,
 5) the graph (x, u) includes $(n-1)$ arcs. If n is a number of apex, m is a number of loops, then $m = n-1$.

Using the known Trent's theorem^{*} about the number of trees and remembering that the great-tree is a tree with oriented arcs let us present the matrix as follows:

$$\|a_{ij}\| = \begin{cases} |\Gamma_{x_i}| & \text{at } i=j \\ -1 & \text{at } i \neq j \quad (x_i, x_j) \in u \\ 0 & \text{at } i \neq j \quad (x_i, x_j) \notin u \end{cases}$$

Theorem III. The number of great-trees in a graph is equal to an element minor Δ_i of main diagonal in the square matrix with order n .

For symmetric graph a minor with a root x_i would be:

$$\Delta_i = \begin{vmatrix} (n-1) & -1 & \dots & -1 \\ -1 & (n-1) & \dots & -1 \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ -1 & -1 & \dots & (n-1) \end{vmatrix} \quad (2)$$

As follows from (2) for symmetric graph $\Delta_i = n^{(n-2)}$. Then:

- 1) the number of the connections K providing composition of a great-tree is equal to $K = n(n-1) - (n-1) = (n-1)^2$;
- 2) all the set M of the great-trees in symmetric graph is equal to $B = n^{(n-1)}$
- 3) in general the number B of possible versions of the great-trees is equal to $B = \Delta_i n$.

Theorem IV. The graph G is a structural one with a root x_i when and only when $\Delta = I$; in other cases $\Delta = 0$.

Theorem V. The number of the structural graphs is equal to a minor Δ_i of an element in the main diagonal of the square incidence matrix $\|a_{ij}\|$. In symmetric graph $G(x, \Gamma)$ the number of structural graphs with a root x_i is $(n-1)!$.

^{*}) Trent H.M. A note of the enumeration and listing of all possible trees in a connected linear graph. Proc. Nat. Ac. Sciences, 40, 1954, 1004.

From this theorem the formula for determining number of connections providing composition of the structural graph has been obtained

$$K = \frac{n(n-1)}{2}$$

All the set of graph structures in symmetric graph is $B = n!$. In general case $b = \Delta_n$.

Autonomous Groups decomposition of system.

We can define the autonomous system as some control loops in the system having actual interconnections through the controlled member.

Let's present the matrix $\|S_n\|$ in the form of:

$$\|S_n\| = \begin{vmatrix} S_K & 0 \\ 0 & S_{n-K}^{n-K} \end{vmatrix} \quad \text{or} \quad \|S_n\| = \begin{vmatrix} S_K & S_{n-K}^K \\ 0 & S_{n-K}^{n-K} \end{vmatrix} \quad (3)$$

then as before:

$$\|S_n\| = |S_K| \times |S_{n-K}^{n-K}|$$

Four versions of the system are possible:

- 1) if $|S_K| = 0$, $|S_{n-K}^{n-K}| = 0$, then $|S_n| = 0$
- 2) if $|S_K| = I$, $|S_{n-K}^{n-K}| = 0$, then $|S_n| = 0$
- 3) if $|S_K| = 0$, $|S_{n-K}^{n-K}| = I$, then $|S_n| = 0$
- 4) if $|S_K| = I$, $|S_{n-K}^{n-K}| = I$, then $|S_n| = I$.

The first version of multivariable systems corresponds to the case when all the autonomous groups are constructed in the form of the symmetric subgraphs. Some part of the groups, defined by $|S_K|$ or $|S_{n-K}^{n-K}|$ is constructed in form of structural graphs in the second and third versions; the fourth version is a graph composed of the autonomous groups in the form of structural graphs.

In these cases the number of necessary connections (K) depends on the number of groups (ν) and control ^(ed)variables in group m :

- 1) for groups in the form of the symmetric subgraphs:

$$K = n(n-1) - \frac{1}{2} \sum_{z=1}^y m_z(m_z-1) \quad (5)$$

- 2) for groups in the form of the structural graphs:

$$K = n(n-1) - \frac{1}{2} \sum_{z=1}^y (m_z-1) \quad (6)$$

- 3) for groups in the form of the great-trees:

$$K = n(n-1) - \sum_{z=1}^y (m_z-1) \quad (7)$$

- 4) at the triangular (or quasi-triangular) system matrix and the presense of groups:

$$K = \frac{1}{2} \left[n(n-1) - \sum_{z=1}^y m_z(m_z-1) \right] \quad (8)$$

REFERENCES.

1. Bruce R.W. The combined aerodynamic propulsive orbit plane change maneuver. AIAA Paper 65-20.
2. Остославский И.В., Стражева И.В. Динамика полета, М., Оборонгиз, 1963.
3. Петров Б.Н. О построении и преобразовании структурных схем, Изв.ОТН АН СССР, 1945, № 12.
4. Колпакова Н.П. К теории систем с несколькими регулируемы-ми координатами. Труды II Всесоюзного совещания по теории инвариантности и ее применению. М., Наука, 1964.
5. Колпакова Н.П. К вопросу синтеза селективно инвариантных систем. Труды III Всесоюзного совещания по теории инвариантности и ее применению (в печати).
6. Техническая кибернетика. Теория автоматического регулирования. Книга I, М., Машиностроение, 1967.

OPTIMAL PARAMETRIC CONTROL FOR

THE RE-ENTRY SPACE VEHICLE

Ponomarev V.M., Gorodezky V.I.

Leningrad State University, Leningrad, USSR.

Stochastic disturbances, acting on a space vehicle (SV) on the atmospheric part of the descent and the initial conditions scattering of the SV re-entry cause the great landing coordinate scatter, especially when the re-entry velocity is higher the first cosmic one. In this case the use of standard linear control equation of SV mass centre stabilization with isochronous parameter variations of the program motion may not give high accuracy. This tends to use more effective control methods.

So far as with the board space vehicle navigation unit and the computer: one may continuously define motion parameters and form control signals by enough complex algorithms it makes possible to use parametric control programs when the program trajectory is given against some kinematic or dynamic trajectory parameter.

I. The problem formulation.

Motion equations.

Let the mass center motion of the S.V. descending in the atmosphere, is described by the following differential equations:

$$\begin{aligned} \frac{dU_x}{dt} = & -K C_x \frac{P}{P_0} U U_x - K C_y \frac{P}{P_0} U U_y - \frac{g_z}{c} (x - x_c) - \frac{g_\omega}{\omega_3} \omega_{3x} + \\ & + a_{11} (x - x_c) + a_{12} (y - y_c) + a_{13} (z - z_c) + b_{12} U_y + b_{13} U_z ; \\ \frac{dU_y}{dt} = & -K C_x \frac{P}{P_0} U U_y + K C_y \frac{P}{P_0} U U_x - \frac{g_z}{c} (y - y_c) - \frac{g_\omega}{\omega_3} \omega_{3y} + \end{aligned} \quad (1)$$

$$+ a_{21}(x-x_c) + a_{22}(y-y_c) + a_{23}(z-z_c) + \beta_{21}U_x + \beta_{23}U_z ;$$

$$\frac{dU_z}{dt} = -K C_x \frac{\rho}{\rho_0} U U_z + K C_z \frac{\rho}{\rho_0} U^2 \beta - \frac{g_z}{2}(z-z_c) - \frac{g_\omega}{\omega_3} \omega_{3z} + \quad (1)$$

$$+ a_{31}(x-x_c) + a_{32}(y-y_c) + a_{33}(z-z_c) + \beta_{31}U_x + \beta_{32}U_y ;$$

$$\frac{dx}{dt} = U_x ; \quad \frac{dy}{dt} = U_y ; \quad \frac{dz}{dt} = U_z .$$

These equations are given in the start coordinate system, related to the rotating Earth. The origin of the system is taken from the vertical cross-point passing through the SV mass centre at the time reference with the Earth surface. Y- axis coincides the vertical, X- axis is directed to the termination point of landing and Z- axis form the right coordinate system.

In equations (1) the following notation is given:
 V_x, V_y, V_z, x, y, z - vector velocity components and the SV mass centre coordinates in the given coordinate system,

C_x, C_y, C_z - aerodynamic factors,

ρ/ρ_0 - the relative air density,

$U = \sqrt{U_x^2 + U_y^2 + U_z^2}$ - the velocity vector modulus,

x_c, y_c, z_c - coordinates of the Earth centre,

\vec{r} - position vector of SV mass centre with the origin in the Earth centre,

$K = \frac{S \rho_0}{2m}$ - ballistic coefficient depending on the cross section area of the space vehicle S, its mass M and atmospheric density ρ_0 on Earth surface,

g_z, g_ω - components of gravitational acceleration,

α, β - vehicle angles of attack and sideslip accordingly,

$\omega_3, \omega_{3x}, \omega_{3y}, \omega_{3z}$ - angular velocity of the Earth and components of angular velocity in the start coordinate system related to the rotating Earth,

α_{ij}, β_{kl} - constants, which depend on $\omega_3, \omega_{3x}, \omega_{3y}, \omega_{3z}$ entrance latitude, initial azimuth trajectory of re-entry 1,

As it was shown above programmer trajectory should be given depending on any parameter with following stabilization of the SV mass centre in programmer trajectory according with a some control equation. From physical point of view such a control is equal feedback programmer control, that indirectly permit to take into account the stochastic disturbances, influenced over SV. In such case "the tube" of disturbed trajectories will be more narrow and therefore the linear control quite sufficient.

The choice of the programmer argument is sufficiently complicated and now is not solved. In the case under consider programmer argument is the down range L , distance to destination along great circle route.

Suppose that control in bank motion is performed with the angle of attack, and in side motion - with the angle of sideslip. α and β programmes are given in such a manner, that the accelerations along the nominal trajectory were less, than 3 g.

Those programmes are choosed in result of preliminary calculations

$$\alpha_{nom} = \begin{cases} 0.2 & ; \text{ if } L \leq 2730 \text{ km} \\ 0.1 & , \text{ if } 2730 \text{ km} < L \leq 3350 \text{ km} \\ 0.05 & , \text{ if } 3350 \text{ km} < L \leq 3428 \text{ km} \\ 0 & , \text{ if } L > 3428 \text{ km} \end{cases} \quad (2)$$

$$\beta_{nom} = 0$$

Stability control equation of SV on the programmer trajectory should be search in such form:

$$\Delta L = K_1(U_x - U_{xnom}) + K_2(U_y - U_{ynom}) + K_3(h - h_{nom}) + K_4(Z - Z_{nom}) \quad (3)$$

$$\Delta \beta = K_5(U_z - U_{znom}) + K_6(Z - Z_{nom}) \quad (4)$$

where $U_{xnom}(L)$, $U_{ynom}(L)$, $U_{znom}(L)$, $h_{nom}(L)$, $Z_{nom}(L)$ — are the nominal values of velocity components.

altitude and side deflection accordingly,

K_1, \dots, K_6 — the unknown (searched) coefficients of linear control equation,

$\Delta L, \Delta \beta$ — controlling variations of angles of attack and sideslip.

Then

$$L = L_{nom} + \Delta L \quad (5)$$

$$\beta = \Delta \beta \quad (6)$$

There are control boundaries of values L and β in this problem which are firstly because of limited control effectiveness, and secondly for limiting of maximum accelerations acting on SV:

$$\begin{aligned} |L| &\leq \bar{L} \\ |\beta| &\leq \bar{\beta} \end{aligned} \quad (7)$$

The basic aim of the re-entry SV control is to reduce the terminal scattering. Therefore as minimization functional was taken the such one.

$$I = D[L] + D[Z] + (M[L] - L_{nom})^2 + (M[Z] - Z_{nom})^2, \quad (8)$$

where $D[L], D[Z]$, $M[L]$, $M[Z]$ — dispersions and means L and Z in point of the SV landing ($h = 0$),

L^p and Z^p — nominal values of the corresponding parameters in point of the landing.

Minimization of the functional permits to solve the problem of SV leading to given point with the minimum scatter.

Disturbances, acting on the re-entry SV are stochastic.

The basic scattering of the SV ground-points arise because of the atmospheric density, c_y/c_x variations, and deflexion of coordinates and re-entry angle initial conditions. Suppose that initial conditions are given in form

$$\begin{aligned} x_0 &= 0, \quad y_0 = (100\,000 + \omega_1) [m], \quad z_0 = \omega_3 [m] \\ v_0 &= 7800 [m/sec], \quad \theta_0 = (0.055 + \omega_2) [rad], \quad u_{z0} = 0 \end{aligned} \quad (9)$$

Here w_1, w_2, w_3 are uncorrelated and normal distributed stochastic coefficients. Their mean are zero and root-mean-square variations are known

$$\begin{aligned} \sigma(\omega_1) &= 2000 \text{ m} \\ \sigma(\omega_2) &= 0.0005 \text{ rad} \\ \sigma(\omega_3) &= 3300 \text{ m} \end{aligned}$$

The disturbances of atmospheric density as altitude function are given by with such canonic decomposition

$$\rho = \rho_{nom} \left(1 + \sum_{i=1}^m \omega_{i+3} f_i(h) \right) \quad (10)$$

Suppose, the dispersions of the stochastic coefficients to be equal 1. Coordinate functions can be get by calculation stochastic data about atmosphere density.

Suppose, that variations of parameter c_y/c_x may be run up to 15%, and its distribution is normal.

Let us set the following problem of terminal control.

It is necessary to find such a control in form (3), which minimizes the functional (8) with given initial conditions (9) and boundaries (7).

2. Algorithm of problem solution.

When solving this problem we shall use the Consecutive Optimization Method³. This Method is founded on replacement of initial problem by consequences of convex quadratic programming problem with linear boundaries.

There are stochastic functional in our problem, therefore additional difficulties arise here.

Algorithm of one step problem solving in such case consist of two independent stages: the achievement of quadratic approximation of functional and linearisation of

boundaries. 2. The solving of quadratic programming problem.

Let us consider the peculiarities of each stages. The general difficulty in getting of quadratic approximation of minimized functional is connected with its stochastic character. It means, that functional's value may be calculated only approximately (by Monte-Carlo Method, Thernetsky's Interpolating Method⁴, Dostupov's Method⁵). Moreover, it is not to state that the functional (8) is convex as for as differential equation (1) are nonlinear. But even it is convex it may be turn out after quadratic interpolation unconvex, that prevent using such well known methods of convexive programming⁶.

Therefore when solving this problem one ought to carry out the quadratic interpolation of given initial functional with additional boundaries.

$$\lambda_i \geq 0, \quad i = 1, \dots, m \quad (11)$$

where λ_i - eigenvalues of quadratic matrix form got in result of interpolation the initial functional. It is evidently, that condition (12) is equal to requirement of convexity of the initial functional (8) quadratic approximation. Then interpolation problem is solved as some smoothing problem.

Let us give the short description of the best in root-mean-square sense convex approximation, which was used for solving synthesis parametric control equation to re-entry SV problem.

Let I^m - is the value of functional (8) in point K^m ($m = 1, \dots, l$) of control parameters space from permitting region. Let be some initial convex quadratic approximation of functional (8). It may be rather far from given functional, therefore its getting is not difficult.

Suppose, that its form is following

$$\bar{I}(K) = C^0 + \sum_{i=1}^n b_i^0 K_i + \sum_{i,j=1}^n a_{ij}^0 K_i K_j \quad (12)$$

All values $K^{(m)}$ are given, and the values of functional I in points K , where are given the accurate values of functional I are given may be written in form.

$$\bar{I} = C^0 + \sum_{i=1}^n \beta_i^0 K_i^{(m)} + \sum_{i,j=1}^n \alpha_{ij}^0 K_i^{(m)} K_j^{(m)} \quad (13)$$

Matrix (α_{ij}^0) is symmetrical and positive. Let us set the following problem of quadratic programming. The problem is to find such values C , β_i , α_i , which realize the minimum of the auxiliary functional

$$\Phi = \sum_{m=1}^{\ell} (\bar{I}^m - I^m)^2 \quad (14)$$

and satisfy to boundaries

$$\lambda_i \geq 0 \quad (15)$$

$$|\bar{I}^m - I^m| \leq \delta, \quad m = 1, \dots, \ell \quad (16)$$

Boundaries (16) are made for excluding of the great difference between the functional (8) I and its convex approximation \bar{I} . However, it is not essential in the most number of problems.

Let us introduce the vector

$$Y^m = (y_0^m, y_1^m, \dots, y_p^m) \quad (17)$$

by formulas

$$\begin{aligned} y_0^m &= 1, \\ y_1^m &= K_1^m, \quad y_2^m = K_2^m, \quad \dots, \quad y_n^m = K_n^m, \\ y_{n+1}^m &= K_1^m K_1^m, \quad y_{n+2}^m = K_1^m K_2^m, \quad \dots, \quad y_{2n}^m = K_1^m K_n^m, \\ y_{2n+1}^m &= K_2^m K_2^m, \quad y_{2n+2}^m = K_2^m K_3^m, \quad \dots, \quad y_{3n-1}^m = K_2^m K_n^m, \\ &\dots \\ y_p^m &= K_n^m K_n^m, \\ p &= \frac{n(n+3)}{2}. \end{aligned} \quad (18)$$

It is evidently, that the components of vector Y^m are the ordered coefficients at c, b_i and d_{ij} in the equality (13).

Let us introduce also the vector to make a short note

$$\Omega = (\omega_0, \omega_1, \dots, \omega_p) \quad (19)$$

and its components ω_j are coefficients c, b_i and d_{ij} in the equality (13) at y_i^m in according to designations (18). Then the minimized functional (18) is the quadratic function of components of vector Ω

$$\Phi(\Omega) = \sum_{m=1}^L (I^m - \sum_{i=0}^p y_i^m \omega_i)^2 \quad (20)$$

It is obviously, that the functional Φ may be made convex by choice of the magnitude y_i^m .

Boundaries (16) are linear. Boundaries (15) may be made linear:

$$\lambda_i(\Omega) \cong \lambda_i(\Omega^0) + \sum_{j=n+1}^p \frac{\partial \lambda_i(\Omega^0)}{\partial \omega_j} (\omega_j - \omega_j^0) \geq 0 \quad (21)$$

that may be written in such a way

$$-\sum_{j=n+1}^p \frac{\partial \lambda_i(\Omega^0)}{\partial \omega_j} \omega_j \leq \lambda_i(\Omega^0) - \sum_{j=n+1}^p \frac{\partial \lambda_i(\Omega^0)}{\partial \omega_j} \omega_j^0 \quad (22)$$

Thus, the problem of the best in the root-mean-square convex approximation of the initial functional is reduced to the problem of minimization quadratic functional by the linear boundaries (16), (22).

The convex programming problem was solved by Hildreth and D'Esopo's Method⁶.

Let us consider some peculiarities of the second solving stage.

It is known that the several variables convex functional minimizing is quite difficult problem and for its solving it is required many of iterations. Below it is offered some methods for decreasing the number of the iterations and value of calculations.

In our practice it was used the quadratic functional optimization method widely in which on every step it was executed the optimization on the subset U , which dimension less than dimension of vector K . In that case every iteration consist of the solving the such problem:

$$\min_{\tilde{U} \subset U} I(K_1, K_2, \dots, K_m) \quad (23)$$

where U is the region of the permitted control. In such method the general difficulty is in selection of subset V sequence which determines the number simultaneously optimized coefficients and the subset \tilde{V} structure.

Sometimes the essential simplification may be get it as the quadratic functional matrix has the quasidiagonal form and in the such case it is possible the independent optimization on the coefficients of every block.

While practic calculating the quadratic form matrix elements are seldom equal zero. It is connected with the approximative character derivatives calculus as well as with correlation which take place.

However it values of some second mixed partial derivatives are far smaller as other it may consider that they are equal zero. That permit to reduce number of iterations and calculation volume as initial task is divided on several simpler problems.

However it may not be used always.

Another way may be indicated which may improve process of task solution both in regard to reduce volume of calculations and in regard to increase improving of the process convergence.

While numerical solution of optimization problem solving it may introduce some scale coefficients on different parameters of control in such a way that in the presence of limitations on parameter

$$\bar{K}_i \leq K_i \leq \bar{\bar{K}}_i, \quad i=1, \dots, m \quad (24)$$

differences $(\bar{\bar{K}}_i - \bar{K}_i)$ would be quantities of identical

order. It was pointed out in the book⁸ about it, and this fact was confirmed by the numerical optimization experience.

For convex functional it may be shown, then absolute value of a partial derivatives of functional on parameters of control $\frac{\partial I}{\partial K}$ are decreased while approaching to optimum. Therefore it may be stated, that on every optimization step value of gradient

$$\frac{\partial I}{\partial K} = \left(\frac{\partial I}{\partial K_1}, \frac{\partial I}{\partial K_2}, \dots, \frac{\partial I}{\partial K_m} \right) \quad (25)$$

permits to estimate crudely possibility of functional decreasing

$$|\Delta I|_{\max}^j \leq \left| \frac{\partial I}{\partial K_j} \right| |K_{j\max}| \quad (26)$$

where $\Delta K_{j\max}$ - as much as possible change of parameter K_j to the side, which corresponds to functional decrease with accordance to the normalized limitations.

Often the some parameter derivative are smaller than others, therefore it is carried out:

$$\left| \frac{\partial I}{\partial K_j} \right| |\Delta K_{j\max}| \gg \left| \frac{\partial I}{\partial K_i} \right| |\Delta K_{i\max}| \quad (27)$$

so in that case

$$\left| \frac{\partial I}{\partial K_j} \right| \gg \left| \frac{\partial I}{\partial K_i} \right| \quad (28)$$

and values $\Delta K_{i\max}$ and $\Delta K_{j\max}$ have identical order in accordance with the normalization. It gives some grounds on r -th optimization step to carry out minimum search in the permissible subset of those parameters which have comparable derivative value and all they are more than others parameter derivatives.

It is required to repeat from step to step this method.

This algorithm simplification method leads to good results.

For the large dimension controlling vector K problem it may be given the good results with the combination of this two methods.

3. Results of calculations

The problem of optimal parametric control for the re-entry SV synthesis was solved on the computer M-20 in accordance with the described in section 2 algorithm.

There was obtained the nominal trajectory with the α - and β - program (2) and with zero disturbances ($\omega_1 = \omega_2 = \dots = \omega_7 = 0$):

While solving this problem it was identified that the $K_1 \div K_4$ and K_5, K_6 optimization may be done independently because of quasidiagonal form of quadratic functional. Then the sensitive analysis showed, that the sensitive of the functional (8) coefficients K_4 and K_5 variations in environment of initial approximation

$$K_1 = K_2 = \dots = K_6 = 0$$

are enough little and the essential decrease of functional may be arise by means of coefficients $K_1 \div K_3$ and K_6 optimization at first.

While solving it was used the consecutive Optimization Method³.

The stochastic characteristics L and Z were calculated by Dostupov's Method⁵. The initial functional value

($K_1 = K_2 = \dots = K_6 = 0$) was as follows

$$I = 0.162 \cdot 10^{12} [m^2]$$

and the maximum variation of SV landing point was

$$\Delta R_{max} = 1298 [km]$$

Such solution was got after the first optimization step.

$$\begin{aligned} K_1 &= -0.3 \cdot 10^{-3}, & K_2 &= -0.35 \cdot 10^{-3}, & K_3 &= -0.1 \cdot 10^{-4}, \\ K_4 &= 0, & K_5 &= 0, & K_6 &= 0.4 \cdot 10^{-4}, \\ I^{(1)} &= 0.552 \cdot 10^8 [m^2], \end{aligned}$$

$$\Delta R_{max} = 14.02 [km]$$

Sensitive analysis on second and third optimization steps was permitted on every step correctly the subset of optimizing coefficients.

The final solution of this problem is as follows

$$K_1 = -0.49 \cdot 10^{-3}, K_2 = -0.26 \cdot 10^{-3}, K_3 = -0.8 \cdot 10^{-5}, \\ K_4 = -0.3 \cdot 10^{-5}, K_5 = 0.6 \cdot 10^{-4}, K_6 = 0.46 \cdot 10^{-4},$$

$$I^{(3)} = 0.407 \cdot 10^6 [m^2]$$

$$\Delta R_{max} = 1.2 [km]$$

How the analysis of disturbed SV centre of mass trajectory shows the transient responses on velocity vector components and coordinates are good, and maximal value of manipulated variables $\Delta \alpha$ and $\Delta \beta$ should take place at initial variations x, z and θ are great enough.

This solution was got to SV re-entry azimuth $A = 0$, but as the checking showed, this solution gave the good results when SV re-entry azimuth had values from the -90° to 90° range.

Small sensitivity of the criterion (8) to variations of coefficients $K_1 \div K_6$ was made clear in result of the solution character in the neighbourhood of the optimal control analysis. It has been permitted to reduce the requirements of their tasking accuracy.

Literature

1. Аппазов Р.Ф., Лавров С.С., Мишин В.П. "Баллистика управляемых ракет дальнего действия". "Наука", 1966.
2. "Управление космическими летательными аппаратами" Под ред. К.Т.Леондеса. "Машиностроение", 1967.
3. Пономарев В.М. "Метод последовательной оптимизации в задачах управления". Известия АН СССР "Техническая кибернетика", № 2, 1967.
4. Чернецкий В.И. "Анализ точности нелинейных систем управления". "Машиностроение", 1968.
5. Доступов Б.Г. "Приближенное определение вероятностных характеристик выходных координат нелинейных систем автоматического регулирования". "Автоматика и телемеханика", т.18, 1957, № II.
6. Кюнц Р.П., Крелле В. "Нелинейное программирование". "Советское Радио", 1965.
7. Д.Дж.Уайльд "Методы поиска экстремума". "Наука", 1967.

STOCHASTIC OPTIMIZATION OF SPACESHIP REENTRY CONTROL
IN ATMOSPHEREA.G.Vlasov, E.I.Mitroshin, I.S.Ukolov
Moscow, USSRAbstract

The paper considers the spaceship atmospheric reentry control problem which is reduced to a problem of optimal control of the terminal state of a certain stochastic system. The control algorithm is produced by using the nominal trajectory, the acceleration measurements being considered as the information source. The problems of programming the nominal trajectory of motion and providing the conditions of its actual realization are investigated simultaneously.

- x -

Among the problems connected with the manned spaceship flight a prominent place belongs to the problem of safe reentry into the Earth atmosphere. In solving this problem there arise great difficulties, particularly in case of reentry into the atmosphere with superorbital velocities. A reentry control system must secure the spaceship landing into the given region with due regard of limitations of accelerations, aerodynamic heating etc.

Since disturbing factors in the reentry process (initial scatter of reentry parameters, fluctuations of atmospheric density etc.) are stochastic with given probability characteristics the reentry control problem should be considered as

a stochastic one. Without loss of generality one may assume the latter to be reduced to the optimal control problem of the terminal state of a stochastic dynamic system described by the set of nonlinear differential equations:

$$\dot{X} = X(x, u, h, t) \quad (1)$$

where X - is a state vector of dimension $(n \times 1)$ of the system; U is a control vector of dimension $(z \times 1)$, usually belonging to a closed region \mathcal{U} ; h - is a stochastic disturbance vector; X is a known vector function; t is an independent variable (time or one of state coordinates; for simplicity later on t denotes time; $t \in [0, T]$); $\dot{}$ - is a symbol of differentiating with respect to t .

Information about the current state of the system as a result of observations usually made by using autonomous devices on board the spaceship are represented in the form

$$y = Y(x, e, t) \quad (2)$$

where y - is an observed vector of dimensions $(e \times 1; e \leq n)$ (e.g. acceleration vector); e - is a stochastic error vector; Y is a known vector function.

It is necessary to obtain the extremum of a terminal state function (e.g. the minimum of down range scatter or the minimum of heat supplied to the spaceship during the reentry etc.).

$$J = M \tilde{\omega} [x(T)] \quad (3)$$

where M - is a symbol of mathematical expectation; $\tilde{\omega}$ - is a scalar nonnegative function, with due regard of the following

boundary conditions: $X(0)$ is a vector of stochastic quantities with definite probability characteristics; at the moment T :

$$T \in \{T: \Omega[X(T), T] = 0\} \quad (4)$$

where Ω is a nonlinear function, the relation

$$X(T) \in \{X: g_K(X) = 0\} \quad (5)$$

(where g_K is a nonlinear vector function) must be satisfied with a definite degree of probability.

Spaceship reentry problems are characterized by a particular requirement of meeting the current phase restriction of inequality type (e.g. maximum acceleration limit); in other words, with a definite degree of probability there must be satisfied the following condition:

$$g(t, x) \leq 0 \quad (6)$$

where g is a nonlinear vector function of dimension $(m \times 1)$.

Stochastic disturbance vectors and the measured errors in a general case include the stochastic parameters ξ (e.g. scatter in spaceship parameters and initial conditions of reentry) and the stochastic processes q (e.g. atmospheric density fluctuations) as well.

By using the shaping filters the stochastic processes q can be approximated in the form of solutions of differential equations of the form

$$\dot{q} = fq + \varepsilon \quad (7)$$

here f is a known matrix; ε is a vector of stochastic

δ -correlated processes ("white noise").

If for stochastic parameters η we write formally the equation of the shaping filter in the form

$$\dot{\eta} = 0 \quad (8)$$

and if we assume now X to be the expanded state vector $\begin{pmatrix} x \\ \eta \end{pmatrix}$ and X to be the expanded vector function, respectively, then without loss of generality equations (1) and (2) can be written as follows:

$$\dot{x} = X(x, u, \varepsilon, t) \quad (9)$$

$$y = Y(x, \xi, t) \quad (10)$$

where ε, ξ - denote "white noise".

In such a general formulation the problem investigation is greatly complicated due to nonlinearity of the set of equations (9)-(10).

However assuming the disturbing factors are negligible and, consequently, the disturbed trajectory is close to the theoretical (nominal) one it is possible to use a linearization method for describing the disturbed motion thus significantly simplifying the investigation of the problem.

Supposing that disturbances acting on the nominal trajectory are equal to zero one may assume the investigated dynamic system be described by the following set of differential equations:

$$\dot{x} = X(x, u, t) \quad (11)$$

$$\dot{x}_1 = A(x, u, t)x_1 + B(x, u, t)u + \varepsilon(x, u, t) \quad (12)$$

$$\dot{Y}_1 = H(x, u, t)X_1 + \xi(x, u, t) \quad (13)$$

where X - is a system state vector in case of motion along the nominal trajectory; X_1 - is a generalized system state vector for the case of disturbed motion; Y_1 - is an observed vector; ξ, ξ - denote "white noise"; U - is a control vector for the case of motion along the nominal trajectory; V is a control vector for the case of motion along the disturbed trajectory; X - is a known vector function corresponding to the motion along the nominal trajectory; A, B, H are matrices of corresponding dimensions.

The boundary conditions at the ends are written in the form:

$$X(0) \in q_0; \quad X(T) \in \{x: g_k(x) = 0\}; \quad T \in \{T: \Omega[X(T), T] = 0\} \quad (14)$$

where $X_1(0)$ is a vector of stochastic quantities with definite probability characteristics.

For the problems of spaceship descent in the atmosphere it is typical the following relation imposed on the control vector in case of motion along the nominal and disturbed trajectories:

$$U + V \in V \quad (15)$$

In accordance with the above mentioned assumptions, the current restriction of the inequality type (6) and optimized functional (3) (for simplicity later on it will be assumed as minimized) may be written as follows:

$$g(t, x) + g_x(t, x)X_1 \leq 0 \quad (16)$$

where g_x is a matrix of partial derivatives of g with

respect to X ,

$$J = M \tilde{\omega} [x(T) + x_1(T)] \quad (17)$$

If $\tilde{\omega}$ is considered to be the function

$$\tilde{\omega} [x(T) + x_1(T)]^T \Lambda [x(T) + x_1(T)] \quad (18)$$

where $[]^T$ is a symbol of transposition; Λ is a weight matrix, then the expression for J can be written in the following full form:

$$J = \omega [x(T)] + M \omega_1 [x_1(T)] + \omega_2 [x(T), M x_1(T)] \quad (19)$$

where $\omega, \omega_1, \omega_2$ are scalar functions, the first two of which being non-negative.

The problems of programming the nominal motion and securing its actual realization are usually investigated separately. The former is considered as a dereminate optimal guidance problem from the minimization conditions

$$J_0 = \omega [x(T)] \quad (20)$$

The latter problem being stochastic reduces to the optimal regulator synthesis from the minimization conditions

$$J_1 = M \omega_1 [x_1(T)] \quad (21)$$

It is natural to solve both problems simultaneously. The importance of such an approach to solving the optimization problem was pointed out in Reference ¹.

Consider at first the disturbed motion, since the optimal regulator, as it will be shown later on, is structurally

invariant relative to the nominal motion parameters. It will be supposed that the nominal motion is prescribed and matrices A, B, H and "white noise" in equations (I2)-(I3) are only time functions.

For the optimal regulator synthesis we shall use the dynamic programming method and the concept of sufficient coordinates³. The sufficient coordinates are the coordinates of a space in which Bellman functional equation is being considered. Their introduction permits to formally separate the problems of information processing and optimal synthesis. In solving the first problem sufficient coordinates are defined by using linear and nonlinear optimum filtration techniques. The defined sufficient coordinates are used in the second problem for optimum synthesis by solving Bellman equation.

It is assumed that $x_1(0), \varepsilon(t), \xi(t)$ are independent and normally distributed:

$$\begin{aligned} M[x_1(0)] &= 0 & M[\varepsilon(t)] &= M[\xi(t)] = 0 \\ M[x_1(0)x_1^T(0)] &= K_0; & M[\varepsilon(t)\varepsilon^T(\tau)] &= Q(t)\delta(t-\tau); \\ & & M[\xi(t)\xi^T(\tau)] &= G(t)\delta(t-\tau); \end{aligned} \quad (22)$$

where K_0 is an a priori covariance matrix of the vector $x_1(0)$; and $Q(t), G(t)$ - are known ^{intensity} matrices of gaussian "white noise" $\varepsilon(t), \xi(t)$.

Then for the investigated dynamic system (I2)-(I3) the vector of sufficient coordinates coincides with the vector of a posteriori mathematical expectation \hat{x} of the generalized state vector x_1 that can also be defined as the vector of the optimal estimation according to the method of the a posteriori probability density maximum by solving the following set

of differential equations (Kalman filter):

$$\dot{Z} = AZ + BU + PH^T G^{-1} (y - HZ) \quad (23)$$

where P - is a covariance matrix of the estimation errors and is determined from the differential equations (24) and

$[]^{-1}$ is a matrix inversion symbol

$$\dot{P} = AP + PA^T - PH^T G^{-1} HP + Q \quad (24)$$

Equations (23)-(24) are solved under the initial conditions

$$Z(0) = M[x_0(0)] = 0; \quad P(0) = K_0 \quad (25)$$

Accordingly Bellman equation is written as follows:

$$-\frac{\partial \varphi(z, \dot{z})}{\partial z} = \min_{v(t) \in U(t)} [\varphi_z(Az + Bv) + \frac{1}{2} \text{tr} \{ \dot{P} \varphi_{zz} \dot{z} \}] ; \quad \varphi(z, t) = \min_{v(\tau), \tau \in [t, T]} M \{ \omega_1[x_1(\tau)] / Z(t) \} ; \quad (26)$$

where φ is a loss function; φ_z is a vector of gradient φ with respect to z ; φ_{zz} is a matrix of the second partial derivatives φ with respect to z ; $\text{tr} \{ \dot{P} \}$ is a spur of matrix $\dot{P} = PH^T G^{-1} HP$; $M[]$ is a symbol of conditional mathematical expectation. This equation must be solved with the boundary condition

$$\varphi[z, T] = M \{ \omega_1[x_1(\tau)] / z(\tau) \} \quad (27)$$

The synthesis of an optimal regulator on the basis of the solution of Bellman equation is extremely difficult in a general case. Consider a simplified problem formulation assuming that

ω is a scalar quantity and the function $\omega_1[x_1(\tau)]$ is

$$\omega_1[x_1(\tau)] = \{ \chi[x_1(\tau)] \}^2 \quad (28)$$

where χ is a linear form.

Then, adding an additional component $\chi[x_r(\tau)]$ to the vector $x_r(t)$, multidimensional synthesis of an optimal regulator can be reduced to a one dimensional form by transition to a new variable according to the formula²:

$$\rho(t) = \Phi(\tau, t) x_r(t) \quad (29)$$

where $\Phi(\tau, t)$ is a fundamental matrix of a homogeneous equation corresponding to equation (12) and $\Phi(\tau, t)$ satisfies the following relations:

$$\dot{\Phi}(\tau, t) = -\Phi(\tau, t) A; \quad \Phi(\tau, \tau) = E \quad (30)$$

where E is a unity matrix.

$$\begin{aligned} \dot{\rho} &= \bar{B}v + \bar{\varepsilon}; \quad y_r = \bar{H}\rho + \bar{f}; \quad \rho(\tau) = x_r(\tau); \\ \bar{B} &= \Phi(\tau, t) B; \quad \bar{H} = H\Phi^{-1}(\tau, t); \quad \bar{\varepsilon} = \Phi(\tau, t) \varepsilon; \end{aligned} \quad (31)$$

Equations for the a posteriori mathematical expectation of ρ and the a posteriori covariance matrix of estimate errors \bar{P} are written respectively

$$\dot{\bar{x}} = \bar{B}v + \bar{P}\bar{H}^T G^{-1}(y_r - \bar{H}\bar{x}); \quad \dot{\bar{P}} = \bar{Q} - \bar{P}\bar{H}^T G^{-1} \bar{H}\bar{P};$$

where: \bar{x} is a vector of the optimal estimate ρ and \bar{P} is a covariance matrix of the estimate errors ρ ;

$$\bar{Q} = \Phi(\tau, t) Q \Phi^T(\tau, t); \quad \bar{x}(0) = 0; \quad \bar{P}(0) = \Phi(\tau, 0) K_0 \Phi^T(\tau, 0)$$

Assuming that $\chi[x_r(t)]$ is a first coordinate of the state expanded vector $x_r(t)$ and denoting by the symbol " ' ", here and later on, the first coordinates and the first lines

of vectors and matrices respectively, we may write Bellman equation as:

$$-\frac{\partial \varphi(\bar{x}', t)}{\partial t} = \min_{u(t) \in U(t)} [\varphi_{\bar{x}'}' \bar{B}' u - \frac{1}{2} \varphi_{\bar{x}'}' \bar{x}' \bar{x}'] \quad (32)$$

where $\bar{x}' = (\bar{p} \bar{H}'' \bar{G}^{-1}) \bar{G} (\bar{p} \bar{H}'' \bar{G}^{-1})'$

This equation must be solved under the boundary conditions

$$\varphi(\bar{x}', T) = M \left\{ [\rho'(T)] \frac{1}{\bar{x}'}(\pi) \right\} \quad (33)$$

However, for the optimal synthesis, in this case, it is not necessary to solve Bellman equation.

Consider the case of discrete information input that is typical for atmospheric reentry problems, provided the information is processed by the digital computer on board the spaceship. Assume the information enters at the discrete time $0 \leq t_N < t_{N-1} < \dots < t_1 < t_0 (t_0 < T)$ with the discrete interval Δt in which \bar{v} is supposed to be constant. Then the differential equation for $\rho'(t)$ and Bellman equation are substituted by recurrent equations:

$$\rho_{K-1}' = \rho_K' + \bar{v}_K \int_{t_K}^{t_{K-1}} \bar{B}' dt + \int_{t_K}^{t_{K-1}} \bar{G}' dt \quad (34)$$

$$\varphi_K(\bar{x}_K') = \min_{\bar{v}_K \in \bar{V}_K} M [\varphi_{K-1}(\bar{x}_K' + \bar{v}_K \int_{t_K}^{t_{K-1}} \bar{B}' dt + \bar{v}_K') / \bar{x}_K'] \quad (35)$$

$$\bar{V}_K' = \int_{t_K}^{t_{K-1}} \bar{G}' dt - \bar{\delta}_K' + \bar{\delta}_{K-1}', \quad \bar{\delta}_K' = \bar{x}_K' - \rho_K'$$

where index "K" corresponds to the moment t_K . Supposing that the region of permissible values of \bar{v} is asymmetrical,

$$\bar{v}_{Kmin} \leq \bar{v}_K \leq \bar{v}_{Kmax} \quad (36)$$

what is natural to expect taking into account that the controls

are interconnected in nominal and disturbed motions following from the condition (15), it can be shown that the optimal control $v_{K \text{ opt}}$ in the interval $[t_K, t_{K-1}]$ is

$$\begin{aligned}
 & \left. \begin{aligned}
 & v_{K \text{ min}}, \text{ if } \text{sign} \left\{ \left[\bar{x}'_K + \sum_{i=K-1}^0 \frac{v_{i \text{ max}} + v_{i \text{ min}}}{2} \int_{t_i}^{t_{i-1}} \bar{B}' dt \right] \cdot \int_{t_K}^{t_{K-1}} \bar{B}' dt \right\} = +1; \\
 & v_{K \text{ max}}, \text{ if } \text{sign} \left\{ \left[\bar{x}'_K + \sum_{i=K-1}^0 \frac{v_{i \text{ max}} + v_{i \text{ min}}}{2} \int_{t_i}^{t_{i-1}} \bar{B}' dt \right] \cdot \int_{t_K}^{t_{K-1}} \bar{B}' dt \right\} = -1; \\
 & v_{K \text{ max}} \left| \int_{t_K}^{t_{K-1}} \bar{B}' dt \right| \leq - \left[\bar{x}'_K + \sum_{i=K-1}^0 \frac{v_{i \text{ max}} + v_{i \text{ min}}}{2} \int_{t_i}^{t_{i-1}} \bar{B}' dt \right] \text{sign} \int_{t_K}^{t_{K-1}} \bar{B}' dt \neq \\
 & \leq v_{K \text{ min}} \left| \int_{t_K}^{t_{K-1}} \bar{B}' dt \right|; \\
 & \bar{x}'_K + v_K \int_{t_K}^{t_{K-1}} \bar{B}' dt + \sum_{i=K-1}^0 \frac{v_{i \text{ max}} + v_{i \text{ min}}}{2} \int_{t_i}^{t_{i-1}} \bar{B}' dt = 0; \\
 & v_{K \text{ max}} \left| \int_{t_K}^{t_{K-1}} \bar{B}' dt \right| > - \left[\bar{x}'_K + \sum_{i=K-1}^0 \frac{v_{i \text{ max}} + v_{i \text{ min}}}{2} \int_{t_i}^{t_{i-1}} \bar{B}' dt \right] \text{sign} \int_{t_K}^{t_{K-1}} \bar{B}' dt > \\
 & > v_{K \text{ min}} \left| \int_{t_K}^{t_{K-1}} \bar{B}' dt \right|
 \end{aligned} \right\} \quad (37)
 \end{aligned}$$

Hence it appears that on every segment Δt the optimal regulator with maximum speed tends to combine the optimal estimate \bar{x}' with a nonzero quantity in

$$\left[- \sum_{i=K-1}^0 \frac{v_{i \text{ max}} + v_{i \text{ min}}}{2} \int_{t_i}^{t_{i-1}} \bar{B}' dt \right]$$

conditioned by the asymmetry of the region of permissible values of the control except a certain region of non-unique values of the optimal control

$$\begin{aligned}
 & v_{K \text{ max}} \left| \int_{t_K}^{t_{K-1}} \bar{B}' dt \right| > - \left[\bar{x}'_K + \sum_{i=K-1}^0 \frac{v_{i \text{ max}} + v_{i \text{ min}}}{2} \int_{t_i}^{t_{i-1}} \bar{B}' dt \right] \cdot \\
 & \cdot \text{sign} \int_{t_K}^{t_{K-1}} \bar{B}' dt > v_{K \text{ min}} \left| \int_{t_K}^{t_{K-1}} \bar{B}' dt \right|; \quad (38)
 \end{aligned}$$

which depends on the discrete interval Δz and the permissible values of the control.

Coming to the limit $\Delta z \rightarrow 0$ in (37) we have

$$v_{opt}(t) = \begin{cases} v_{min}(t), & \text{if } \dot{x}(t) \operatorname{sign} \bar{B}'(t) > - \int_t^T \frac{v_{max}(\tau) + v_{min}(\tau)}{2} \cdot \bar{B}'(\tau) d\tau \operatorname{sign} \bar{B}'(t); \\ 0, & \text{if } \dot{x}(t) = - \int_t^T \frac{v_{max}(\tau) + v_{min}(\tau)}{2} \bar{B}'(\tau) d\tau; \\ v_{max}(t), & \text{if } \dot{x}(t) \operatorname{sign} \bar{B}'(t) < - \int_t^T \frac{v_{max}(\tau) + v_{min}(\tau)}{2} \bar{B}'(\tau) d\tau \operatorname{sign} \bar{B}'(t). \end{cases} \quad (39)$$

Thus, the algorithm of the optimal control during the continuous information input is completely determined without direct solving Bellman equation. In case of continuous information input, in contrast to the determinate case $[E(t) = \dot{x}(t) = 0]$, where the optimal control has not unique values in the so called "complete controllability zone", in the stochastic problem the optimal control is defined uniquely on the basis of (39); at every time the optimal control with maximum speed tends to combine the optimal estimate $\dot{x}(t)$ with the quantity

$$- \int_t^T \frac{v_{max}(\tau) + v_{min}(\tau)}{2} \bar{B}'(\tau) d\tau$$

conditioned by the asymmetry of permissible value region of the control.

Substituting the expression for the optimal control v_{opt} (39) in Bellman equation (32), from the solution of the latter we can define that part of the minimized functional J which is conditioned by the disturbed motion - J_d .

$$\min J_d = \min M[p'(\tau)] \stackrel{L}{=} p[\dot{x}'(0), t=0]; \quad (40)$$

Since it is very difficult to obtain the exact solution of Bellman equation we use an approximate method of solution, that is, parametric method.

The parametric method consists in approximating the function $\varphi(\bar{x}', t)$ as a function of a finite number of parameters

$$\varphi(\bar{x}', t) = \psi^T(\bar{x}') a(t); \quad (41)$$

where ψ is a vector function of the known expansion functions; a is a $(S \times 1)$ -dimensional vector of unknown parameters.

Imposing on the parameters a certain natural conditions for the exactness of approximating function $\varphi(\bar{x}', t)$ and using Bellman equation, we can get a set of ordinary differential equations with boundary condition at the moment T which describes the evolution of the parameters in reverse time.

These equations, respectively, for different parametric methods described in Refs. 3, 4, are written as follows:

$$1) \dot{a} = - \left[\int_{\bar{x}_n^*}^{\bar{x}_0^*} \psi \psi^T d\bar{x}' \right]^{-1} \int_{\bar{x}_n^*}^{\bar{x}_0^*} (\bar{B}' v_{0n} \psi_{\bar{x}'}^T + \frac{1}{2} \psi_{\bar{x}' \bar{x}'}^T \psi') a \cdot \psi d\bar{x}' \quad (42)$$

$$2) [\psi^T / \bar{x}' = \bar{x}'_n(t) - \psi^T / \bar{x}' = \bar{x}'_n(t)] \dot{a} = [(v_{0n} \bar{B}' \psi_{\bar{x}'}^T + \frac{1}{2} \psi_{\bar{x}' \bar{x}'}^T \psi') / \bar{x}' = \bar{x}'_n(t)] - (v_{0n} \bar{B}' \psi_{\bar{x}'}^T + \frac{1}{2} \psi_{\bar{x}' \bar{x}'}^T \psi') / \bar{x}' = \bar{x}'_n(t)] \cdot a; \quad K=2, 3, \dots, (S+1) \quad (43)$$

where $[\bar{x}'_n, \bar{x}'_0]$ denotes the region of permissible values \bar{x}' ; $\bar{x}'_n(t)$ are certain prescribed functions \bar{x}' ; $\psi_{\bar{x}'}$ is a vector of gradient ψ with respect to \bar{x}' ; $\psi_{\bar{x}' \bar{x}'}$ is a matrix of the second partial derivatives ψ with respect to \bar{x}' .

The equations are solved under the boundary conditions

$$a(T) = a_T \quad (44)$$

where a_{γ} follows from the conditions of securing (33).

The quantity $\min J_{\gamma}$ is written as follows:

$$\min J_{\gamma} = P[\bar{x}'(0), t=0] = \gamma^T [\bar{x}'(0)] \cdot a(0) \quad (45)$$

As mentioned above, because of the stochastic nature of the component of the vector x_{γ} we may speak about solving the inequality (16) only with some a priori level of probability; therefore it is necessary to know the a priori probability density of vector x_{γ} or its substitute parameters at every time t . In this case the estimation of the exactness of the stochastic vector means obtaining the vector of an a priori mathematical expectation and an a priori covariance matrix.

As we are interested in the estimation of the exactness of the stochastic vector $g_x x_{\gamma}$ (of dimension $m \times 1$) then by introducing the expanded vector of state $x_2 = \begin{pmatrix} g_x x_{\gamma} \\ \frac{f}{2} \end{pmatrix}$ we can state that for the system of the type

$$\dot{x}_2 = S(t)x_2 + F(t, x_2) + G(t) \quad (46)$$

where the matrix $S(t)$, ~~is~~ nonlinear vector-function $F(t, x_2)$ and "white noise" $G(t)$ correspond to the expanded state vector

x_2 . It is necessary to define the first m components of the vector

$$\bar{x}_2(t) = M x_2(t) \quad (47)$$

and the matrix consisting of the first m lines and columns of the matrix

$$\bar{\bar{x}}_2(t) = M [x_2(t) - \bar{x}_2(t)] [x_2(t) - \bar{x}_2(t)]^T \quad (48)$$

Supposing that $x_2(0)$ and $G(t)$ are independent and normally distributed

$$\bar{x}_2(0) = x_{20}; \quad \bar{\bar{x}}_2(0) = K_{20}; \quad MG(t) = 0; \quad (49)$$

$$M[G(t)G^T(\tau)] = K(t) \delta(t - \tau)$$

and using the method of the statistic linearization, we can show that $\bar{x}_2(t)$ and $\bar{\bar{x}}_2(t)$ approximately satisfy the following set of differential equations:

$$\dot{\bar{x}}_2 = S(t)\bar{x}_2 + \Psi_1(t, \bar{x}_2, \bar{\bar{x}}_2); \quad \bar{x}_2(0) = x_{20} \quad (50)$$

$$\dot{\bar{\bar{x}}}_2 = [S(t) + \Psi_2(t, \bar{x}_2, \bar{\bar{x}}_2)]\bar{\bar{x}}_2 + \bar{x}_2[S(t) + \Psi_2(t, \bar{x}_2, \bar{\bar{x}}_2)] + K \quad (51)$$

$$\bar{\bar{x}}_2(0) = K_{20}$$

where Ψ_1, Ψ_2 are functional relations whose definition ultimately reduces to the solution of the following relation

$$\lambda(t, M_{\bar{x}}', \bar{\sigma}_{\bar{x}}') = \frac{1}{2} [\bar{v}(t) + \bar{\bar{v}}(t)] + [\bar{v}(t) - \bar{\bar{v}}(t)] \cdot \Phi \left[\frac{\bar{x}' - M_{\bar{x}}'}{\bar{\sigma}_{\bar{x}}'} \right] \quad (52)$$

where $M_{\bar{x}}'$ and $\bar{\sigma}_{\bar{x}}'^2$ are components of \bar{x}_2 and $\bar{\bar{x}}_2$ corresponding to the coordinate \bar{x}' :

$$\bar{v}(t) = \begin{cases} v_{\max}(t), & \text{if } \text{sign } \bar{B}' = +1; \\ v_{\min}(t), & \text{if } \text{sign } \bar{B}' = -1; \end{cases} \quad \bar{\bar{v}}(t) = \begin{cases} v_{\max}(t), & \text{if } \text{sign } \bar{B}' = -1; \\ v_{\min}(t), & \text{if } \text{sign } \bar{B}' = +1; \end{cases}$$

$$\bar{x}'^* = - \int_t^{\infty} \frac{v_{\max}(\tau) + v_{\min}(\tau)}{2} \bar{B}'(\tau) d\tau;$$

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-\frac{t^2}{2}} dt \quad (53)$$

where Φ is a probability integral.

Neglecting, for simplicity, the correlation between components of the vector $(x_2 - \bar{x}_2)$ it is possible to assume, in accordance with the rule "36", that the requirement for solving the stochastic inequality (16) reduces to the necessity of solving the following determinate inequality

$$g(t, x) + \bar{x}_m(t, x) + 3\bar{b}_m(t, x) \leq 0 \quad (54)$$

where \bar{x}_m is a vector including the first m components of the vector x_2 and \bar{b}_m is a vector including the square roots of the first m diagonal elements of the matrix \bar{x}_2 .

Now, having determined the structure of the optimal regulator, estimated that part of the functional which is conditioned by the disturbed motion and estimated the exactness of the generalized state vector in case of disturbed motion we may come to the conclusion that the joint consideration of the determinate problem of programming the nominal motion and stochastic problem of securing its actual realization reduces to the following determinate extremal problem.

Given a set of ordinary differential equations (notations are given above) prescribed on the segment $[0, T]$

$$\begin{aligned} \dot{x} &= X(x, u, t); \quad \dot{\phi}(T, x, t, u) = -\phi(T, t, x, u) A(x, u, t); \\ \dot{\bar{P}}(x, u, t) &= \bar{Q}(x, u, t) - \bar{P}(x, u, t) \bar{H}^T(x, u, t) \bar{G}^{-1}(x, u, t) \bar{H}(x, u, t) \bar{P}(x, u, t); \\ \dot{\alpha}(x, u, t) &= - \int_{\bar{x}_n}^{\bar{x}_n^*} \bar{\psi}(\bar{x}') \bar{\psi}^T(\bar{x}') d\bar{x}' - \int_{\bar{x}_n}^{\bar{x}_n^*} [\bar{v}_{on}^T(x, u, t; \bar{x}') \bar{B}'(x, u, t) \bar{\psi}(\bar{x}') + \\ &\quad + \frac{1}{2} \bar{\psi}^T(\bar{x}') (\bar{x}' \bar{A}'(x, u, t)) \bar{\psi}(\bar{x}') d\bar{x}']; \\ \dot{\bar{x}}_2(x, u, t) &= S(x, u, t) \bar{x}_2(x, u, t) + \bar{\psi}_2(x, u, t); \\ \dot{\bar{\bar{x}}}_2(x, u, t) &= [S(x, u, t) + \bar{\psi}_2(x, u, t)] \bar{\bar{x}}_2(x, u, t) + \bar{\bar{x}}_2(x, u, t) [S(x, u, t) + \\ &\quad + \bar{\psi}_2(x, u, t)]^T + K(x, u, t); \end{aligned}$$

with boundary conditions

$$\begin{aligned} x(0) \in g_0; \quad x(\tau) \in \{x: g_x(x) = 0\}; \quad \tau \in \{T: \Omega[x(\tau), \tau] = 0\}; \\ \Phi(\tau, \tau) = E; \quad \bar{p}(0) = \Phi(\tau, 0) K_0 \Phi^T(\tau, 0); \quad a(\tau) = a_\tau; \\ \bar{x}_2(0) = x_{20}; \quad \bar{x}_2(0) = K_2 0; \end{aligned} \quad (56)$$

It is necessary to select the program $u(t) (u \in U)$ and the initial conditions $x(0)$ according to the conditions

$$\min J = \min \left\{ \omega[x(\tau)] + \psi[\bar{x}'(0)] a(0) + \omega_x[x(\tau), \bar{x}_2(\tau)] \right\} \quad (57)$$

in meeting the current restrictions of the type of inequality

$$g(t, x) + \bar{x}_m(x, u, t) \pm 35_m(x, u, t) \leq 0 \quad (58)$$

This is a determinate problem and it pertains to the problem of optimal control

having the limits for control and phase coordinates and can be solved by using well known approximate calculation methods. Thus the algorithm for defining the optimal nominal motion is completely determined.

The realization of the optimal algorithms of controlling the disturbed motion in reentry problems obtained as stated above requires having a high speed digital computer on board the spaceship.

The specific conditions of the spaceship motion in the atmosphere with superorbital velocity completely define the self-sufficiency of the spacecraft control system. Accordingly as the source of information it is possible to use time and

accelerations measured in inertial or fixed axes of the spacecraft.

A high speed digital computer on board the spacecraft by statistical processing acceleration measurement results permits to obtain complete information about the spacecraft motion parameters. The necessary information about the nominal motion that is selected beforehand by the above described optimal method is stored on board the spacecraft. The deviations of the actual values of motion parameters from the nominal ones are used for obtaining the necessary input control signal.

R E F E R E N C E S

1. LETOV A.M., "Optimal Control Theory," Proceedings of the 2nd Congress of IFAC, Vol. "Optimal Systems. Statistical Methods". "Nauka" Publishing House, Moscow, 1965 (In Russian).
2. BOGUSLAVSKY I.A., On Statistical Optimal Control of Terminal State. Journal "Automatika i Telemekhanika", No.5, 1966.
3. STRATONOVICH R.L., Recent Development of Dynamic Programming Methods and their Application for Synthesis of Optimal Systems". Proceedings of the 2nd Congress of IFAC, Vol. "Optimal Systems" (Statistical Methods.) "Nauka" Publishing House, Moscow, 1965 (In Russian).
4. KROTOV V.F., Approximate Synthesis of Optimal Control". Journal "Automatika i Telemekhanika", No.II, 1964.

ATMOSPHERE RE-ENTRY CONTROL PROBLEM

Okhotsimski D.E., Bukharkina A.P., Golubiev Yu.F.

Institute of Applied Mathematics

Moscow

USSR

This paper deals with the problem of re-entry control at parabolic velocity. The multistep algorithm of the downrange control is described. Some results of the digital computer simulation are given. It is shown that the algorithm provides a sufficiently smooth and accurate regulation process despite large-scale variations in the atmospheric density distribution. This paper develops further the results of previous work.¹⁻⁴

I. Statement of the problem. This paper deals with the problem of re-entry control algorithm at parabolic velocity. It is supposed that the space vehicle has three accelerometers installed on the stabilized platform. By changing the roll angle the vehicle total lift direction changes. This makes it possible to control the motion of the vehicle. The onboard digital computer calculates the needed roll angle values taking into account the measuring data processing results.

The papers^{1, 2} contain an example of the downrange control algorithm. The algorithm built in the paper³ provides lateral-range control as well. The paper⁴ contains the algorithm of the re-entry initial conditions calculation and considers the effect of the systematic, instrumental and executive errors.

The above-mentioned control algorithm of the first-dip portion of the re-entry motion (I) provided a rather narrow range of kinematic parameters at the end of the skip-out portion (II) and thus gave the possibility to compensate the deviations using the restricted control potentialities of the second-dip portion (III) (fig. 1). The algorithm displayed fortitude to the errors and density distribution uncertainties, and provided a sufficient accuracy in the re-entry corridor of about ± 15 km.

These results cleared the way to the more perfect algorithms. In particular, it was desirable that the roll angle should be a more usable continuous function of time instead of a step-function employed in.^{I-4} It was also desirable to take into consideration the actual constraints on the roll control torques and also the calculation time for decision taking. This paper presents a step in this direction.

2. Decision logics. Let us at first neglect the time for the calculation and divide the first-dip portion into equal time intervals. Let us suppose that the control decisions are taken when passing from one time interval to another, and that the instant roll angle jumps are impossible, and that the roll rate $\dot{\gamma}$ is constant during each time interval. The function $\gamma(t)$ will be a piecewise linear function the angular points of which coincide with the interval ends.

Let us assume that the roll rate jump is constrained by the condition

$$|\dot{\gamma}_{i+1} - \dot{\gamma}_i| \leq \Gamma \quad (2.1)$$

where Γ is the increment of the roll angle caused by the control torque during a one time interval. The fulfilment of the condition (2.1) provides obtaining such a piecewise linear function $\gamma(t)$ which may be approximately realized.

During each decision-making the values of the roll rate and, consequently, the roll angle for the nearest time interval are chosen in such a way that, with the roll angles of the rest of the first-dip portion, the desired downrange is provided. Like in^{I-4} two integrations of the motion equations forward up to the end of the first-dip portion are carried out. The data thus obtained give the grounds for decision taking.

It appeared reasonable to use different decision logics during the initial part of the motion up to the velocity head maximum and during the rest of the portion. Let us describe some of the variants that have been investigated and give a number of arguments for selection.

In one of simplest variants we choose $\dot{\gamma}$ for the ne-

arest time interval (t_i, t_{i+1}) to provide the prescribed down-range. Discontinuity at t_{i+1} is admissible, but at t_i the roll angle must be continuous. The repeating of the process gives the roll angle as a piecewise linear function of time. But the permission of a jump at the end of the time interval may result in a saw-tooth function.

It is possible to avoid this unpleasant event by appropriate choice of γ not for one, but for two nearest time intervals (fig. 3). Let the value of γ at t_i join continuously the previous values and at t_{i+2} the beforehand nominated function $\tilde{\gamma}(t)$. Let us choose the value γ at t_{i+1} or the value of γ for the interval (t_i, t_{i+1}) (which is actually the same) to provide the required down-range. The repeating of the process gives a continuous piecewise linear function without saw-teeth.

Fig. 3 shows that if all the conditions for decision-making at point t_i would also be true for point t_{i+1} and if the two-links construction built at t_i would satisfy condition (2.1), in this case the first link of the new two-links construction during the next step of decision-making at t_{i+1} would coincide with the second link of the previous step and the new second link would coincide with the portion of the function $\tilde{\gamma}(t)$. Thus, there would be a two-links passage to the function $\tilde{\gamma}(t)$ by means of the construction obtained at t_i . As a consequence, if we deviated from the function $\tilde{\gamma}(t)$ we would return to it after two time intervals. Also, if $\tilde{\gamma}(t)$ were any admissible polygon the roll angle γ would arrive at $\tilde{\gamma}(t)$ and then coincide with it.

Actually, the atmospheric density distribution is partly unknown and conditions for decision-making vary from one point to another. It is reasonable to modify the function $\tilde{\gamma}(t)$ in accordance with the density variations which are revealed during the flight. But, nevertheless, the algorithm with two-links logics appeared to be effective enough and was used at first on the second part of the first-dip portion.

Within the bounds of the first part of the portion the effectiveness of the roll angle change is small at first and

increases with time. This causes large-scale variations of the roll angle and insufficient smoothness of the regulation process. In view of this it appears to be more reasonable on the initial part of the first-dip portion to vary a constant value of the roll angle within a more lengthy period of time which includes several standard time intervals, and to add two links (at its beginning and at its end) for smooth connection with the previous and consequent course of the function $\gamma(t)$ (fig. 4). The varied position of the function $\gamma(t)$ is shown by a dashed line. The closer to the region of the maximum velocity head the shorter is the interval of the constant γ value variation, and the described logics turns smoothly enough into two-links logics.

To successfully overcome the atmospheric density uncertainties it appears reasonable to take the standard time intervals shorter, and to take decisions more often. But the decrease of the time interval length causes the decrease of the regulation smoothness. It was suitable to modify two-links logics in such a way so that each link would include two standard time intervals. One time interval later, at t_{i+1} we have to build a new two-links construction with the same length of each link, but biased by one standard interval further (fig. 5).

It appears also reasonable to modify a little the logics for the initial part of the portion and to change from one to two standard intervals the length of the link connecting the previous roll angle course with the interval of the constant roll angle value, but to take decisions after each standard time interval.

When choosing $\dot{\gamma}$ for the nearest time interval it may occur that the $\dot{\gamma}$ jump will be greater than the permissible one, according to (2.1). In this case the maximum possible value must be chosen.

For the downrange control it is sufficient to vary γ within the range of 0 to π . The whole lift force is directed upwards, if $\gamma = 0$ and downwards, if $\gamma = \pi$. If γ has an intermediate value, the vertical component of the total lift force is also intermediate. In particular, if $\gamma = \frac{\pi}{2}$

the vertical lift force component equals zero and the longitudinal motion of the space vehicle is like the motion without a lift force.

During the regulation the roll angle can reach the upper or the lower stop and remain there for some time. To avoid violation of (2.1) when reaching the stop it was assumed a limitation for $\dot{\gamma}$ to provide the stop reaching smooth enough.

The algorithm is able to carry out an analysis of the real density distribution. Knowing the vehicle position and velocity and using the measuring data given by the accelerometers we can calculate for each time the real density ρ and the quantity

$$\xi = \frac{\rho}{\rho_{st}}$$

where ρ_{st} is the density value at the same point according to the standard density distribution. Along the motion, ξ appear to be a function of time the previous course of which is well-known. Let us assume, as in I-4, that ξ is a rather smooth function of coordinates. In this case $\xi(t)$ along the motion is a rather smooth function of time, and it is possible to extrapolate this function forward for a short period of time. Such an extrapolation appeared to be useful during the integrating of the motion equations forward for decision taking.

The information about the possible $\xi(t)$ course in the nearest future was also used for the appropriate choice of the function $\tilde{\gamma}(t)$ outside the interval of time where the choice was carried out to provide the desired down-range. The choice of $\tilde{\gamma}(t)$ ensured adaptation and created the controllability reserve to compensate the predicted density variations.

The original function $\tilde{\gamma}(t)$ is shown in fig. 6. The portion of the constant value $\gamma = \ell$ at t_1 turns into an inclined line and then at t_2 into a second portion of the constant γ value which was originally taken equal $\frac{\gamma}{2}$. The values t_1 and t_2 are chosen beforehand and remain fixed. When varying ℓ to provide down-range the inclined portion of $\tilde{\gamma}(t)$ varies too. It is shown by a dashed line (fig. 6). In the course of time the length of the varied por-

tion becomes shorter, and near t_1 the logics turns into a two-links one or into its modified variant.

The adaptation is carried out by changing the position of the function $\tilde{y}(t)$ after the time t_2 . The quantity m (fig. 6) may be calculated according to the formula

$$m = \frac{T}{2} + \Delta \tilde{y} \quad (2.2)$$

where $\Delta \tilde{y}$, as in $I-4$, is

$$\Delta \tilde{y} = f(t) [A \dot{\xi} + B \ddot{\xi}] \quad (2.3)$$

The coefficients A and B are constants found empirically; $\dot{\xi}$ and $\ddot{\xi}$ are the first and the second derivatives of the function $\xi(t)$ calculated for the decision time according to the polynomial which approximately describes the previous course of the function $\xi(t)$; $f(t)$ is the function of time which equals zero before the region of the maximum velocity head, then it increases linearly and, after reaching the prescribed value, it remains constant. The parameters of the function $f(t)$ were defined empirically. A number of reasons associated with the introduction of such an adaptation type was discussed in papers^{1, 2}.

3. Taking account of the calculation time. In the previous section of this paper we neglected the calculation time. Let us take it into account. Let us assume that during one standard time interval all the calculations connected with the decision-making were completed. Let us also assume when calculating that only that information may be used which was obtained prior to the beginning of the time interval. In the course of this calculation we obtain information about the current position and velocity and about the course of the function $\xi(t)$, but these data may be employed only during the next decision-making in the next time interval. Thus we have a lag between the end of the last information inflow and the beginning of the decision execution. The lag is one standard time interval long.

Let us change a little the previous decision logics. Let us assume that everything remains ^{true}, but when choosing \tilde{y} for the nearest time interval the information of the previous time interval cannot be used. To enter the previous algorithm

logics we have to calculate all the quantities forward for the onetime interval. We obtain the values of the function $\xi(t)$ by extrapolation. Let us calculate coordinates and velocity components by integrating the motion equations for a one time interval using extrapolated values of the function $\xi(t)$. Let us obtain the quantities in the formula for the function

$\tilde{\gamma}(t)$ by extrapolation. After all this additional calculations we can enter the previous algorithm.

This approach appears to be very suitable for it gave a simple method to investigate the effect of the time lag using the algorithm without the time lag.

In fact, the measurement information obtained in the course of the calculation may be partly used to diminish the time lag effect. This would be useful for the fortitude of the algorithm.

4. Simulation results. The algorithm performance test was simulated on a digital computer. The components of the aerodynamic acceleration were calculated by integration of the motion equations which imitated the moving space vehicle. The variations of the standard density distribution and the roll angle produced by the control algorithm were inserted into this equation system. The control algorithm itself was realized as a program block for the computer. It received the imitated measurement information, integrated the navigational equations, predicted the $\xi(t)$ course and carried out the adaptation and other calculations which were necessary to take decisions for the roll angle for the nearest time interval.

The density variations were taken as in^{2, 3}. These variations imitated the density deviations depending both on the altitude and on the longitudinal range. Their magnitude is apparently larger than really existing. Therefore the algorithm which is able to struggle successfully with them will have a certain reliability reserve.

Some simulation results are presented in fig. 7 and in others. The thick polygon represents the roll angle versus time during the first-dip portion. The initial point of the time is located at the altitude of 150 km. The algorithm

function begins when the integral of the acceleration reaches the prescribed value. At first, when the efficiency is small, the flight is uncontrolled and the roll angle is maintained as a constant which depends on the perigee altitude. When $t = 45$ sec the algorithm begins to act. The thin line represents the function $\xi(t)$ built up during the flight. The dashed line represents the quantity η (fig. 6) which depends on the $\xi(t)$ course according to the (2.2) and (2.3).

It is seen that for the variant in fig. 7 the $\gamma(t)$ course is at first rather quiet. The tendency towards the density increase is compensated by the decrease of the roll angle and causes an increase of the vertical lift force component which controls the longitudinal movement. The increase of the $\xi(t)$ changes to decrease. Performing the $\xi(t)$ forecast the algorithm increases the roll angle in advance so as to avert the space vehicle from skipping out of the atmosphere earlier than its velocity would be sufficiently braked.

The fast change of the situation regarding the $\xi(t)$ function course requires a fast response. The control torques restrictions and the calculation time lag diminish the response of the system. Therefore the fast change in the $\xi(t)$ course leads the roll angle to the upper stop and leaves it there up to the end of the first-dip portion. Despite the measures taken the time lag of the system causes the skip-out with an excessive velocity and the initial point of the second dip appears to be biased 210 km forwards. Such a deviation may be easily compensated by the second-dip control, for it is deeply within the permissible deviation range.

It should be noted that the variant in fig. 7 was the most difficult of all the atmospheric density variations investigated. This variant was very hard on the system, for the $\xi(t)$ function began a fast decrease. Therefore the controllability reserve appears to be insufficient. The investigation shows that a similar fast change from decrease to increase of the atmospheric density appears to be much easier for the control algorithm.

All other variants investigated have considerably less second-dip initial point deviations than the variant shown in fig. 7. In the cases where the regulation process ended within the stops the deviations were usually no more than some scores of kilometres.

Fig. 8 shows a smoother course of the $\xi(t)$ function when increase changes to decrease. The roll angle reaches the upper stop but leaves it after some time. The accuracy of the second-dip initial point is rather high.

Fig. 9 illustrates a longer stay on the upper stop.

Fig. 10 demonstrates the density variation which causes the increase of the density mainly. The rather fast change of the $\xi(t)$ course from decrease to increase leads the roll angle to the lower stop in order to leave the atmosphere as fast as possible. The deviation is negligible.

Fig. 11 shows the regulation process which takes place within the stops. The deviation is small. It is possible to observe the influence of the $\xi(t)$ course on the adaptation parameter γ and on the roll angle.

Fig. 12 presents the variant in which the fast decrease of the density causes the upper stop of the roll angle for some time. The downrange deviation is very small.

The simulation results demonstrate that a reasonable choice of the parameters of the algorithm makes it sufficiently resistant to the density distribution variations. The algorithm provides for the second-dip initial point parameters to be within the domain where the second-dip control ensures precise space vehicle landing in the prescribed region.

References

1. Д.Е.Охочимский, Г.И. Белъчанский, А.П. Бухаркина, Ю.Ф. Голубев, Н.И. Золотухина, Ю.Н. Иванов. Оптимальное управление при входе в атмосферу. Космические исследования, т. VI, вып. I, 1968 г.
2. Д.Е. Охочимский, А.П. Бухаркина, Ю.Ф. Голубев, Ю.Н. Иванов. Управление продольным движением при входе космического аппарата в атмосферу. Отчет ИПМ АН СССР, 1967 г.

3. А.П. Бухаркина, Ю.Ф. Голубев, Д.Е. Охотимский. Управление пространственным движением при входе космического аппарата в атмосферу. Отчет ИПМ АН СССР, 1968 г.
4. Ю.Ф. Голубев. Определение условий входа в атмосферу и анализ влияния погрешностей. Отчет ИПМ АН СССР, 1968 г.



Fig. 1 Re-entry motion scheme

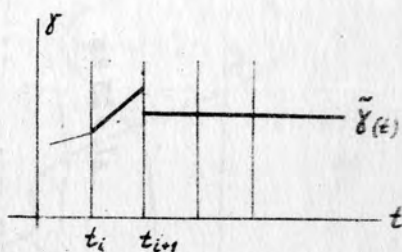


Fig. 2

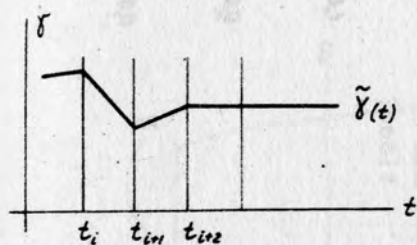


Fig. 3

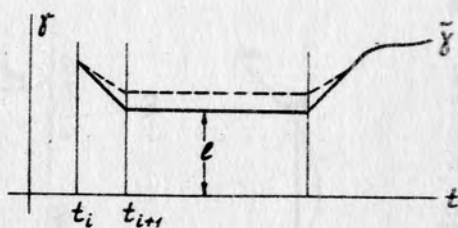


Fig. 4

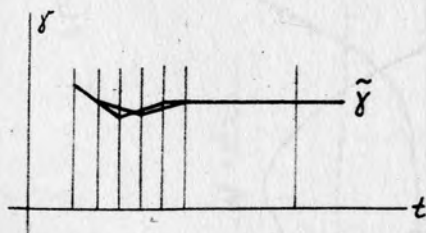


Fig. 5

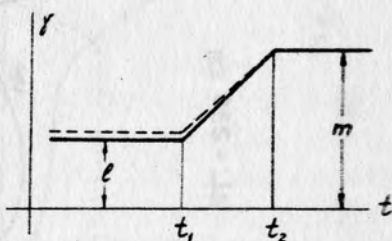


Fig. 6

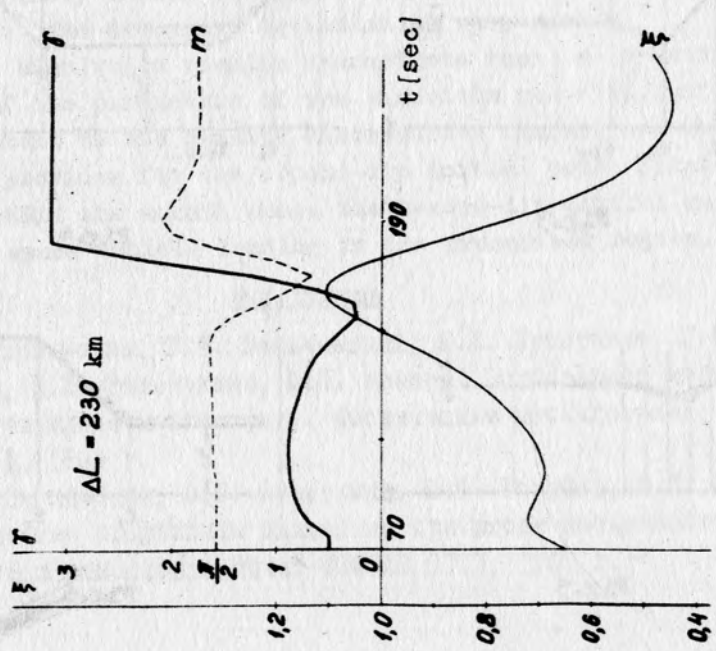


Fig. 7

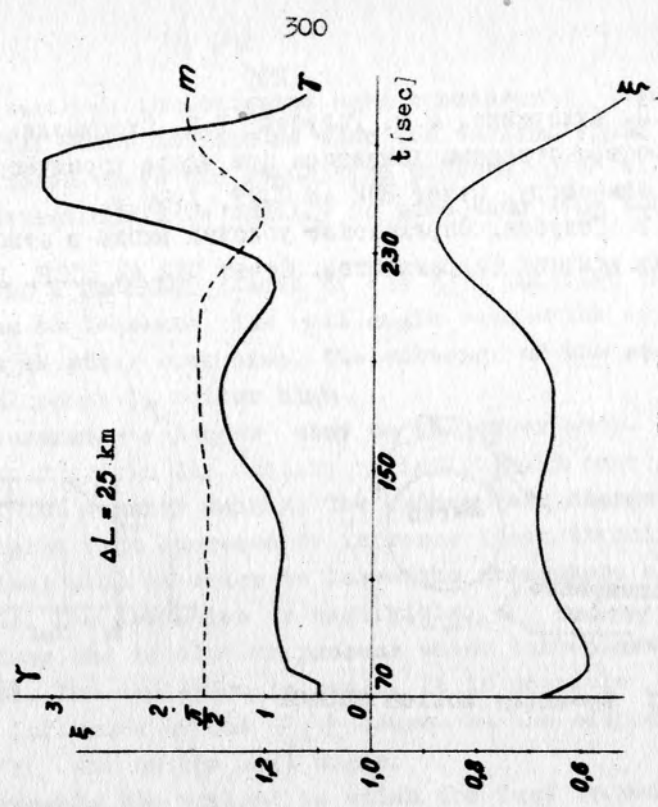


Fig. 8

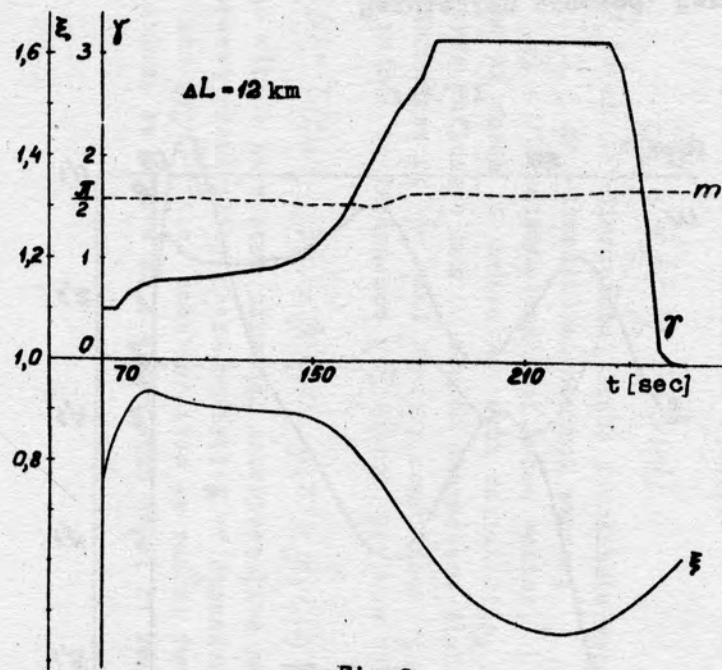


Fig. 9

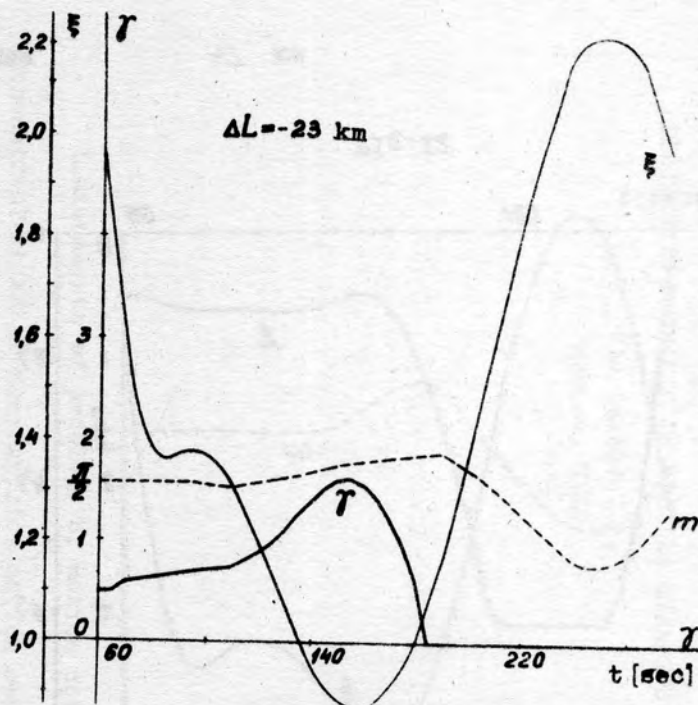


Fig. 10

Regulation process. Perigee

47 km

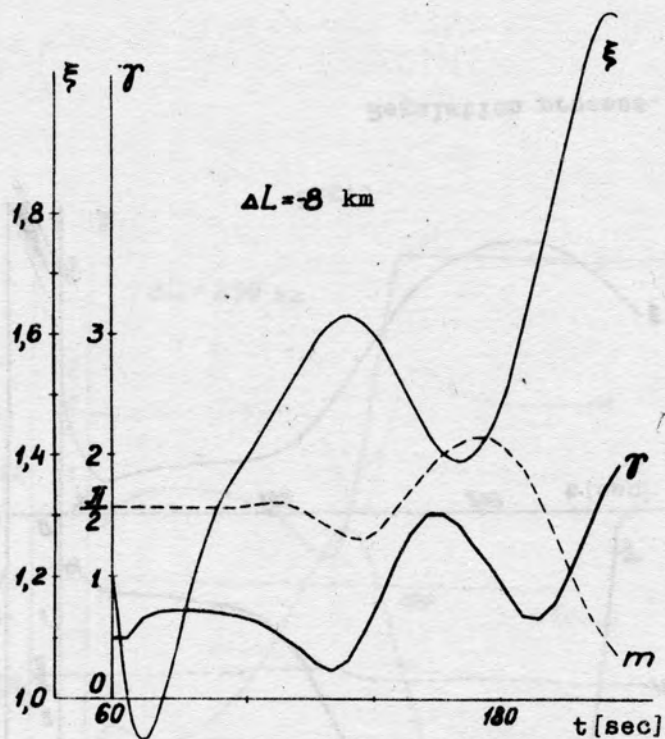


Fig. II

Regulation process. Perigee

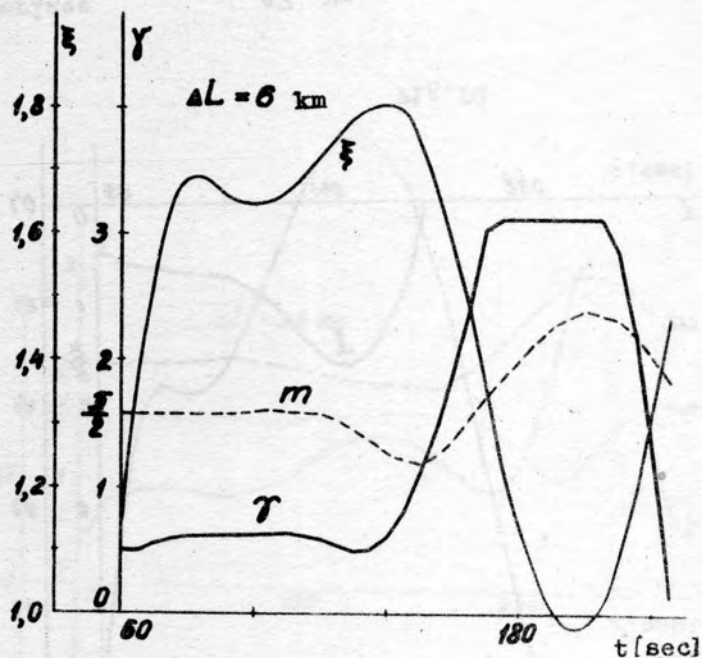


Fig. I2

47 km

STOCHASTIC PROBLEMS OF MISSILE DYNAMICS

Plo tnikov Y.P.

Moscow, USSR

I. Statements of problems of motion control
to be met in practice

Mathematical model of many applied motion control problems is preset by a system of equation

$$\frac{dy}{dt} = f(t, y, u, \xi(t))$$

here y is an n -dimensional state coordinate vector;

u - m -dimensional control vector;

ξ - r -dimensional influence vector,

the value of which is unknown with reliability, i.e. in other words influence has random character. Statement

about randomness of ξ and boundary conditions $y(0), y(\tau)$

formally reduce to dependence $\xi, y(0)$ and $y(\tau)$ from "chance" ω :

$$\xi(t) = \xi(t, \omega), \quad y(0) = y(0, \omega), \quad y(\tau) = y(\tau, \omega).$$

ω can either have some physical interpretation here or not.

As we preset probability characteristics of dependence of $\xi(t)$

$y(0)$ and $y(\tau)$ on ω , i.e. consider them as random functions or variables, we arrive at a stochastic model of real

control process,

$$\frac{dy}{dt} = f(t, y(t, \omega), v(t, \omega), \xi(t, \omega)) \quad (1)$$

Vector function f is supposed to have the properties, ensuring the existence of a solution of the system (1) for the class of controls $v(t, \omega)$ and perturbations $\xi(t, \omega)$ under discussion.

Not all trajectories $y(t, \omega)$ and controls $v(t, \omega)$ possible for (1) are considered on practice, but only those which with a given probability

1. belong to the region

$$g(t, y, v) \leq 0 \quad (2)$$

where g is a given vector function of t, y and v .

2. begin under condition

$$y(q, \omega) \in \{y: g_q(y, \xi(q)) = 0\} \quad (3)$$

3. and end, when

$$y(\tau, \omega) \in \{y: g_\kappa(y, \xi(\tau)) = 0\} \quad (3')$$

Performance characteristic of motion to be minimized by admissible $v(t, \omega)$ choice usually represents the expectation of some functional, including integral and terminal parts.

By now the optimization theory of stochastic control systems has not yet reached such a completeness, which optimal deterministic system theory already has. This is the particular reason for the fact that the routine practice has worked out its own way of the solution of the above raised problem. This way assumes stage by stage solution of the problem, with its being divided into two basic ones.

The first of them is the problem of optimal motions programming; to be more exact, the problem of individual dynamics of a controlled object, where forces and motion

are partially given and missing forces and motion are to be found so that it would offer stipulated optimal properties². In that instance, physical model is described here by a deterministic system of differential equations

$$\frac{dx}{dt} = f(t, x, u, \bar{\xi}(t)) \quad (4)$$

where x means state vector, u - control vector sought for, being function of time and initial conditions.

At any instant functions x and u are such that

$$\bar{g}(t, x, u) \leq 0$$

when $t = 0$, $x(0) \in \{x: \bar{g}_\mu(x) = 0\}$

and when $t = T$, $x(T) \in \{x: \bar{g}_\mu(x) = 0\}$,

one of coordinates of vector $x(t)$ being optimized.

Mathematical model of program motion (deterministic as a rule) is obtained from (1), replacing $\bar{\xi}(t, \omega)$, $y(0)$ and $y(T, \omega)$ by average values of $\bar{\xi}(t)$, $x(0)$ and $x(T)$.

The consideration of the first problem belongs to the sphere of interest of scientists and engineers engaged in a particular field of dynamics, which is concerned with the given object.

By its construction the process, treated in the first problem, differs from original one. The reason for this is other initial conditions, as well as difference in values of the right parts of motion equations in any instant of time. The idea of considering the first problem is that, knowing its solution $x(t)$, $u(t)$, we can represent the real solution in the following form:

$$y = x + z, \quad v = u + w, \quad \bar{\xi}(t, \omega) = \bar{\xi}(t) + \bar{\xi}^0(t, \omega)$$

The perturbed motion z, w is defined as

$$\begin{aligned} \frac{dz}{dt} &\equiv \dot{z} = f(t, y, v, \bar{\xi}(t, \omega)) - f(t, x(t), u(t), \bar{\xi}(t)) = \\ &= A(t)z + B(t)w + C(t)\bar{\xi}^0(t, \omega) + R(t, z, w, \bar{\xi}^0) \end{aligned} \quad (5)$$

Also

$$g(t, x(t) + z, u(t) + w) =$$

$$Z(0, \omega) \in \{z: g_{\mu}(x(0) + z, \xi) = 0\} = \{z: g_{\mu}^0(z, \xi^0) = 0\} \quad (6)$$

and at the end of the motion ^{with} some probability

$$Z(T, \omega) \in \{z: g_{\mu}(x(T, \omega) + z, \xi_T) = 0\} = \{z: g_{\mu}^0(z, \xi^0) = 0\}$$

The right part of (5) is an expansion in Taylor series, so matrices A , B and C have the next form

$$A(t) = \frac{\partial f(t, y, v, \xi)}{\partial y} \Big|_{y=x(t), v=u(t), \xi=\bar{\xi}}$$

$$B(t) = \frac{\partial f(t, y, v, \xi)}{\partial v} \Big|_{y=x(t), v=u(t), \xi=\bar{\xi}}$$

$$C(t) = \frac{\partial f(t, y, v, \xi)}{\partial \xi} \Big|_{y=x(t), v=u(t), \xi=\bar{\xi}}$$

A usual linear model of perturbed motion is obtained, if it is possible to neglect the remainder R in (5).

The second basic problem is formulated for system (5). It consists in determination of law W for transformation of information about perturbed motion so that the chosen control $w(t, z)$ should put out this motion to the best advantage, with $w(t, z)$ being able to depend evidently only upon coordinates of process to be measured.

2. Shortcomings of generally used statements of the problems

A Specialist on control who has to point out the method and means for actual realization of program motion, regards it as the prescribed one, which in particular is reflected in the way of presentation of the perturbed motion in the form of (5), where matrices A , B and C are functions of time^I. They obtain this form only after the program motion construction as $x = x(t)$, $u = u(t)$, with the system of equation (4), which determines these $x(t)$ and $u(t)$, containing no information about future perturbed motion.

This approach to original problem solution that has already become canonical in literature, is connected with the necessity to simplify equation of motion (I) in order to obtain an answer with our restricted computing and algorithmical means. Along with that the separate consideration of the first and the second problems and original system (I) coarsening introduced with it have as a consequence the inadequacy of its properties to those of our differentiated model (4)-(5). The reason for inadequacy is in particular our neglecting the relation between the problem of program $x(t)$, $u(t)$ choice and the one perturbed motion z, w control. But the relation - which is a bilateral one - takes place irrespective of whether $x(t)$, $u(t)$ is a solution of some variational problem or they have been chosen out of some other consideration.

On the one hand, the interrelation proves in the fact that system (5) takes its specific shape only upon the solutions of system (4), when matrices $\partial^2 H / \partial y^2$, $\partial^2 H / \partial v^2$ and $\partial^2 H / \partial z^2$ become the known functions of time, corresponding to specific $y = x(t)$, $v = u(t)$, $z = \bar{z}(t)$. Therefore all perturbed motion properties, such as controllability, stabilizability and their quantitative equivalents depend evidently on the choice of x and u ³.

The program motion acceptability is determined to a great extent by the accuracy, with which a prescribed motion can be realized for given perturbed motion control structure⁹. To select at once the program motion extremal in some sense and satisfying the accuracy requirements, the joint consideration of (4) and (5) system is needed. The same thing is necessitated by the well founded separation of subregions of coordinate changing of perturbed and program motion from region defined by inequality (2). In the first place this refers to the region V of u and w control vectors changes¹. It is necessary to choose the vector so that the set $V-u(t)$ might be substantial enough if we want the problem of control of perturbed motion and its optimization in particular to be have sense.

3. Some new statements of motion control problems

The above mentioned permits to consider it advisable not to fix the program motion control while describing the perturbed motion springing up in its neighbourhood. For these purposes it is necessary to consider the program motion equation system together with perturbed motion equations, thus reformulating the problems, stated for original system (I) to the following one

$$\frac{dx}{dt} = f(t, x, u, \xi) \quad (7')$$

$$\frac{dz}{dt} = \frac{\partial f(t, y, v, \xi)}{\partial y} z + \frac{\partial f}{\partial v} w + \frac{\partial f}{\partial \xi} \xi^0(t, w) + R \quad (7'')$$

($\partial f / \partial y$, $\partial f / \partial v$ and $\partial f / \partial \xi$ have been calculated for $y = x$, $v = u$ and $\xi = \xi(t)$) with

$$g(t, x + z, u + w) \leq 0 \quad (8)$$

$$x(0) \in \{x: \bar{g}_\mu(x) = 0\}, \quad x(t_{np}) \in \{x: \bar{g}_\kappa(x) = 0\} \quad (9)$$

and

$$z(0) \in \{z: g_\mu(x(0) + z, \xi) = 0\}. \quad (10)$$

$$z(\tau) \in \{z: g_\kappa(x(t_{np}) + z, \xi) = 0\} \quad (11)$$

$T - T_{np}$ is a difference of time intervals which perturbed and program motion is exist on, caused by nonsimultaneity to end conditions fulfilment.

Let us first of all pay attention to those problems, where perturbed motion control $w(t, z)$ have^{been} calculated and such program control $u(t)$ is required to be found which would give to system (I) minimal (or maximal) value of one coordinate of vector $x(t_{np})$ under conditions that (8), (9) and (11) is kept with the prescribed probability.

A relative problem to that is the following one: to choose the control in the program movement (7') so that maximize the probability of the event (11) with the

conditions (8) - (10) kept.

In this wording the problem of program motion control defining is a new one which is not be met either when constructing the program motion or stabilizing the perturbed one. At the first turn an practical prototype of this problem occurs where perturbed motion control $w(t, z)$ owing to the restricted composition of magnitudes measured during the motion, admits the existence side by side with perturbations, compensated by this control, noncompensated ones also. The practical value of such statement of the problem consists in the reduction of influence of the latter.

Purposes point out above are achieved not only by the choice of program control $u(t)$. It is permittable to consider that $w(t, z)$ is not give completely, but as a structure only, for example,

$$w(t, z) = K(t) \hat{z}$$

where \hat{z} is part of z , that can be measured. To define the control w it is necessary to set a matrix $K(t)$, choice of it and control $u(t)$ permitting to achieve the purposes in this statement of original problem.

If we don't confine to give structure of perturbed motion control the solution of raised problem (to maximize the probability of event (II) with the rest of restrictions on system (7) being fulfilled). means joint choice of $u(t)$ - control on program trajectory, as well as control synthesis $w(t, z)$ in perturbed motion, that is finding it in terms of time and coordinates measured.

In missile dynamics problems those of them, where necessity of joint program and perturbed systems investigation is obvious (the latter being stochastic one), constitute the bulk of problem which may be called stochastic missile dynamics problems.

The way on which we look for the solution of such stochastic problems is rested on the basic lemma stated in⁴.

4. Sufficient conditions of absolute minimum for stochastic systems

The starting-point of the given direction can be found in⁴. Results obtained in this paper^x are true for a

^x) These results constitute a part of investigation, carried out by the author and V.F.Krotov

vast class of stochastic processes $\xi(t, \omega)$. Limits to that class are given below.

Let the right parts of (I) or (7) are

$$\frac{dy(t, \omega)}{dt} = f(t, y(t, \omega), v(t, \omega), \xi(t, \omega)), \quad t \in [0, T]$$

where random vector function $\xi(t, \omega)$, given on probability space (Ω, \mathcal{B}, P) , with the prescribed control $v(t, \omega)$ satisfies the conditions, which define $y(t, \omega)$ as n -dimensional random process^{5,6}. In a moment T we know the value of functional, calculated on solutions of system

$$\begin{aligned} J &= M \left[\int_0^T f^0(t, y(t, \omega), v(t, \omega), \xi(t, \omega)) dt + F(y(0, \omega), y(T, \omega)) \right] = \\ &= \int_0^T M f^0(t, y(t, \omega), v(t, \omega), \xi(t, \omega)) dt + M F(y(0, \omega), y(T, \omega)) \end{aligned} \quad (I2)$$

Functions $F(y(0, \omega), y(T, \omega))$ and $f^0(t, y(t, \omega), v(t, \omega), \xi(t, \omega))$ are P -integrable here on Ω and on $\Omega \times [0, T]$ respectively for functions y, v and ξ to be met below.

In any moment of time control $v(t, \omega)$ and vector function $y(t, \omega)$ belong to sets $Q(t, y)$ and $B(t)$ of spaces R_n^{Ω} and R_n^{Ω} [see (2), (3) or (8)-(II)]. Let \mathcal{D} be a set of "pairs" $y(t, \omega)$ and v , satisfying to differential equations and restrictions pointed out^x.

Let's state the problem: "from a set "of pairs" $y(t, \omega), v$ we have to find such one for which the functional J would have the least value (if such a pair is absent in the class \mathcal{D} , it is necessary to look for a minimizing sequence $(y_n(t, \omega), v_n) \in \mathcal{D}$ upon which $J_n \rightarrow \inf J > C$).

Lemma mentioned given the opportunity to replace the problem of functional minimization on set \mathcal{D} by the same^x) I don't indicate here functions of what "parameters" the control v is because in each case this defines its own class \mathcal{D} . At this general case it is important only to point out that v belong to given set for $t \in [0, T]$.

problem on more vast set E of independent pair of vector functions $(y(t, \omega), v)$, satisfying to all conditions raised above, besides equation (I).

Let's put to consideration a functional

$$\varphi[t, y(\omega)], \quad y(\omega) \in R_n^{\Omega}, \quad t \in [0, T],$$

differentiated with respect to t and also having restricted continuous Gateau - derivalive $D_y \varphi$ for arbitrary random variable $y(\omega)$ from R_n^{Ω} .⁸

Let

$$R[t, y(\omega), v] = \frac{\partial \varphi[t, y(\omega)]}{\partial t} + D_y \varphi[t, y(\omega)] \times f(t, y(\omega), v; \xi(t, \omega)) - M f^0(t, y(\omega), v; \xi(t, \omega)) \quad (I3)$$

and

$$\Phi[y(0, \omega), y(T, \omega)] = M F(y(0, \omega), y(T, \omega)) + \varphi[T, y(T, \omega)] - \varphi[0, y(0, \omega)] \quad (I4)$$

Theorem A. Let us have the "Pair" of functions $\bar{y}(t, \omega), \bar{v}$, then in order that this "pair" would minimize the functional J on D the existence of such functional $\varphi[t, y(\omega)]$ is sufficient with properties mentioned that

$$1. R[t, \bar{y}(t, \omega), \bar{v}] = \sup_{\substack{y(\omega) \in B(t) \\ v \in Q(t, y)}} R[t, y(\omega), v] = \mu(t) \quad (I5)$$

$$2. \Phi[\bar{y}(0, \omega), \bar{y}(T, \omega)] = \inf_{\substack{y(0, \omega), y(T, \omega) \\ y(0, \omega), y(T, \omega)}} \Phi[y(0, \omega), y(T, \omega)], \quad (I6)$$

$y(0, \omega)$ and $y(T, \omega)$ belonging to the sets (3)^x. When absolute minimal does not exist on D theorem terms 1 and 2, defining a minimizing sequence $y_n(t, \omega), v_n$, coincide with cited ones, if the sign of equality is replaced by the symbol: " \rightarrow as $n \rightarrow \infty$ ".

To prove the theorem let us define on set E functional

$$L = M F(y(0, \omega), y(T, \omega)) + \varphi[T, y(T, \omega)] - \varphi[0, y(0, \omega)] - \int_0^T \left\{ \frac{\partial \varphi[t, y(\omega)]}{\partial t} + D_y \varphi[t, y(\omega)] \times f(t, y(\omega), v; \xi(t, \omega)) - M f^0(t, y(\omega), v; \xi(t, \omega)) \right\} dt \quad (I7)$$

^x) The theorem is also true in that case, when and are elements of some Banach spaces.

This functional is a continuation on set E of functional J , defined on D . Indeed, on D owing to (I) and that for (I),

$$d\varphi[t, y(t, \omega)] = \left[\frac{d\varphi[t, y(t, \omega)]}{dt} + D_y \varphi[t, y(t, \omega)] \times f(t, y(t, \omega), v, \xi(t, \omega)) \right] dt$$

$$L = \Phi[y(0, \omega), y(T, \omega)] - \int_0^T [d\varphi - M^0(t, y(t, \omega), v, \xi(t, \omega))] dt = J$$

If functional $\varphi[t, y(t, \omega)]$ and "pair" of $\bar{y}(t, \omega)$ and \bar{v} , satisfying to theorem terms 1 and 2, exists then from (I5)-(I7) it results that this ^{point} minimizes L on E and by lemma the functional J on D .⁴

Let us pay attention to a case, when $\bar{y}(t, \omega)$ and $v \in L_2(\Omega, \mathcal{B}, P)$ that is when admissible control and corresponding to it (in D) vector $y(t, \omega)$ are random vectors (for any $t \in [0, T]$ the coordinate squares of which are P -integrable ~~with respect to measure~~⁵). Let $\xi(t, \omega) \in L_2(\Omega, \mathcal{B}, P)$

for $t \in [0, T]$ also. Then under the same conditions for the right parts of system (I), which provided the existence of solution $y(t, \omega)$ of (I), $f(t, y(t, \omega), v(t, \omega), \xi(t, \omega)) \in L_2(\Omega, \mathcal{B}, P)$ for $t \in [0, T]$ ^{6, 7}. But for a linear functional in $L_2(\Omega, \mathcal{B}, P)$ is an integral of measure $P(\cdot)$, so

$$D_y \varphi[t, y(\omega)] \times f(t, y(\omega), v, \xi(t, \omega)) =$$

$$= \int_{\Omega} \varphi_y(t, y(\omega), \omega) f(t, y(\omega), v, \xi(t, \omega)) P(d\omega)$$

Here row vector $\varphi_y(t, y(\omega), \omega) \in L_2(\Omega, \mathcal{B}, P)$.

Theorem A terms, putting the constraints on φ and f at minimal $(\bar{y}(t, \omega), \bar{v})$ give an arbitrariness broad enough in setting of functional $\varphi[t, y(\omega)]$ out of it. This allows to choose the most fitted algorithm for solution of problem.

Let us dwell on two of them. We arrive at the first algorithm it we demand from $R[t, y(\omega), v(\omega)]$ to satisfy identically to some conditions in the region $[0, T] \times B$ this is an analog of Hamilton-Iacobi-Bellman formalism for deterministic case⁴.

If we confine ourselves by carrying out the theorem terms on the minimal sought for only, we shall arrive at

an analog of Lagrange formalism. Similar to that case we reduce the problem to a boundary problem for ordinary differential (but stochastic already) equations.

5. An analog of Hamilton-Iacobi-Bellman formalism. Systems, linear in state coordinate

Let here and below the region of change for $y(t, \omega)$ when $t \in (0, T]$ be an open one and

$$\begin{aligned} S[t, y(\omega)] &= \sup_{v \in Q(t, y)} R[t, y(\omega), v] = \sup_{v \in Q} \left[\frac{\partial \varphi[t, y(\omega)]}{\partial t} + \right. \\ &+ \left. \int_{\Omega} [\varphi^v(t, y(\omega), \omega) f(t, y(\omega), v, \xi(t, \omega)) - f^0(t, y(\omega), v, \xi(t, \omega))] P(d\omega) \right] = \\ &= \frac{\partial \varphi[t, y(\omega)]}{\partial t} + \sup_{v \in Q} \int_{\Omega} [\varphi^v f - f^0] P(d\omega) \end{aligned} \quad (18)$$

We shall take $\varphi[t, y(\omega)]$ so that S does not depend on $y(\omega)$, that is

$$S[t, y(\omega)] = \frac{\partial \varphi}{\partial t} + \sup_{v \in Q} \int_{\Omega} [\varphi^v f - f^0] P(d\omega) = c(t).$$

$c(t)$ is a function of time. Then $S[t, y(\omega)] = \mu(t)$ for any $y(\omega)$.

If for $v = \bar{v}[t, y(\omega)]$ $R[t, y(\omega), \bar{v}]$ has a supremum in point $(t, y(\omega))$ that is

$$R[t, y(\omega), \bar{v}] = S[t, y(\omega)], \quad (19)$$

then solution of system (I), $y(t, \omega)$ together with \bar{v} belong to \mathcal{D} , satisfying term I of theorem A. For fulfilment of the second term it is sufficient to demand, that for $t = T$ does not depend on $y(T, \omega)$, that is

$$\Phi[T, y(T, \omega)] = \text{const} \quad (20)$$

As an example for application of above mentioned we can take the systems, linear in state coordinates.

Let right parts of (I) have the form

$$\dot{y} = A(t, \omega)y + h(t, v, \omega) \quad (21)$$

and

$$f^0 = a^0(t, \omega)y + h^0(t, v, \omega), \quad F=0 \text{ here} \quad (22)$$

A is an matrix here, a^0 and h are vectors, h^0 is a scalar function.

We shall find the minimum of functional, putting the -boundaries of control change, depending on t only.

For our functional J R have the form $R[t, y(\omega), v] =$
 $= \frac{dy}{dt} + \int_{\Omega} [\psi^T A - a^0] y + \psi^T h(t, v, \omega) - h^0(t, v, \omega)] P(d\omega)$

and $\mathcal{P}[t, y(\omega)] = \frac{dy}{dt} + \int_{\Omega} [\psi^T A - a^0] y(\omega) P(d\omega) + \mathcal{R}(t)$,
 where

$$\mathcal{R}(t) = \sup_{v \in Q(t)} \int_{\Omega} [\psi^T h(t, v, \omega) - h^0(t, v, \omega)] P(d\omega)$$

To satisfy the system (I) it is necessary to choose functional $\mathcal{P}[t, y(\omega)]$ so that functional would ^{not} depend on $y(\omega)$.

$$\text{Let } \mathcal{Q}[t, y(\omega)] = \int_{\Omega} \Psi(t, \omega) y(\omega) P(d\omega) \quad (24)$$

Then

$$\mathcal{P}[t, y(\omega)] = \int_{\Omega} \left[\frac{d\Psi(t, \omega)}{dt} + \Psi A - a^0 \right] y(\omega) P(d\omega) + \mathcal{R}(t)$$

$$\mathcal{R}(t) = \sup_{v \in Q(t)} \int_{\Omega} [\Psi(t, \omega) h(t, v, \omega) - h^0(t, v, \omega)] P(d\omega)$$

Given row vector function $\Psi(t, \omega)$ by the system

$$\frac{d\Psi(t, \omega)}{dt} + \Psi A = a^0 \quad \text{a.e.}, \quad (25)$$

$$\Psi(T, \omega) = 0 \quad \text{a.e.}, \quad (26)$$

We shall see, that $\mathcal{P}[t, y(\omega)] = \mathcal{R}(t)$ does not depend on $y(\omega)$, and $\mathcal{P}'[t, y(\omega)] = 0$ ($\mathcal{Q}[t, y(\omega)] = \text{const}$ here, for $y(t, \omega)$ has been taken as fixed one). Thus the two terms, put upon the functional in Hamilton algorithm are fulfilled. In this case the choice of functional \mathcal{J} was reduced to Cauchy problem for system of linear differential equations.

6. An analog of Lagrange formalism

Let us suppose that functional $\mathcal{Q}[t, y(\omega)]$ has the second Gateau-derivative and the second mixed derivative $D_y \mathcal{Q}[t, y(\omega)]$ is continuous. We shall look for this functional together with the minimal $\bar{y}(t, \omega), \bar{v}$ from condition of maximum of R over y and v on this minimal. The R and Φ over $y(\omega)$ stationarity condition in $(\bar{y}(t, \omega), \bar{v})$ point can be written as equality to zero of linear part of R 's increment for some $y(t, \omega) = \bar{y}(t, \omega) + h(t, \omega)$, with $h(t, \omega) \in L_2$ and sufficiently small.

Namely,

$$D_y \Psi_t[t, \bar{y}(t, \omega)] \times h(t, \omega) + \int_{\Omega} \frac{\partial}{\partial y} \Psi^y(t, \bar{y}(t, \omega), \omega) h(t, \omega) f(t, \bar{y}(t, \omega), v, \xi) P(d\omega) + \\ + \int_{\Omega} \left[-\frac{\partial f^0}{\partial y} h(t, \omega) + \Psi^y(t, \bar{y}(t, \omega), \omega) \frac{\partial}{\partial y} f(t, \bar{y}(t, \omega), v, \xi(t, \omega)) h(t, \omega) \right] P(d\omega) = 0$$

If we designate $\Psi^y[t, \bar{y}(t, \omega)] = \Psi(t, \omega)$ the first two terms of above equality will be $\int_{\Omega} \frac{d\Psi(t, \omega)}{dt} h(t, \omega) P(d\omega)$ because, due to the second mixed derivative continuouness,

$$D_y \Psi_t = \frac{\partial}{\partial t} D_y \Psi. \text{ Thus this equality can be rewritten as } \int_{\Omega} \left[\frac{d\Psi(t, \omega)}{dt} + H_y(t, \bar{y}(t, \omega), v, \xi(t, \omega)) h(t, \omega) \right] P(d\omega) = 0$$

where $H(t, \bar{y}(\omega), v, \xi) = \Psi(t, \omega) f(t, \bar{y}(\omega), v, \xi) - f^0$

Owing to arbitrariness of $h(t, \omega)$ the condition of R over $\bar{y}(\omega)$ stationarity is equivalent to

$$\frac{d\Psi(t, \omega)}{dt} + H_y(t, \bar{y}(t, \omega), v, \xi(t, \omega)) = 0 \quad \text{a.e.} \quad (27)$$

with such end condition (condition of Φ over $\bar{y}(\tau, \omega)$ stationarity for now $\bar{y}(\tau, \omega)$ is fixed)

$$\Psi(\tau, \omega) = -\frac{\partial}{\partial y} F(\bar{y}(\omega))|_{t=\tau} \quad \text{a.e.} \quad (28)$$

For $\Psi_t[t, \bar{y}(t, \omega)]$ is a number, when $t \in [0, T]$, R over v supremum, the condition defining optimal v , takes the form of

$$\mathcal{H}[t, \bar{y}(t, \omega), \bar{v}] = \int_{\Omega} [\Psi^y(t, \bar{y}(t, \omega), \omega) f(t, \bar{y}, \bar{v}, \xi) - f(t, \bar{y}, \bar{v}, \xi)] P(d\omega) = \\ = \int_{\Omega} [\Psi(t, \omega) f(t, \bar{y}, \bar{v}, \xi) - f^0] P(d\omega) = \quad (29)$$

$$= \sup_{v \in Q} \int_{\Omega} [\Psi(t, \omega) f(t, \bar{y}(t, \omega), v, \xi) - f^0(t, \bar{y}, v, \xi)] P(d\omega) = \sup_{v \in Q} \mathcal{H}[t, \bar{y}, v]$$

In such a case when for $\bar{y}(\tau, \omega)$ end conditions are given, say,

$$\int_{\Omega} G(\bar{y}(\tau, \omega)) P(d\omega) = C \quad (30)$$

the condition of stationarity of $MF(\bar{y}(\tau, \omega)) + \Psi[\tau, \bar{y}(\tau, \omega)]$ takes the form of

$$\Psi(\tau, \omega) = - \left[\frac{\partial}{\partial y} F(\bar{y}(\omega)) + \lambda \frac{\partial}{\partial y} G(\bar{y}(\omega)) \right] |_{t=\tau} \quad (28')$$

λ is an isoperimetric constant to be defined by equality (30).

If the functional $\Psi[t, \bar{y}(\omega)]$ can be preset in neigh-

bourhood of extremal $(\bar{y}(t, \omega), \bar{v})$ so that the sufficient condition of $R[t, y(\omega), v]$'s maximum for $t \in [0, T]$ occurs on extremal itself, the given extremal owing to theorem A is the absolute minimal. Now theorem A may be reformulated so:

Theorem B. Let the aggregate of vector functions $\bar{y}(t, \omega)$, \bar{v} , $\Psi(t, \omega)$ be the result of systems (I), (27) and (29) solution. In order that extremal $\bar{y}(t, \omega), \bar{v}$ would be the absolute minimal the existence of such functional $\Phi[t, y(\omega)]$ is sufficient with mentioned properties that

1. $\Phi(t, \bar{y}(t, \omega), \omega) = \Psi(t, \omega)$
2. $R[t, \bar{y}(t, \omega), \bar{v}] = \sup_{y, v} R[t, y(\omega), v] = \mu(t), t \in [0, T]$
3. $\Phi[\bar{y}(t, \omega)] = \inf_{y(t, \omega)} \Phi[y(t, \omega)]$

7. An linear problem

As an example let us consider for system (I) with the right parts (2I) the problem of functional (22) minimization with the initial $y(0, \omega)$ and end vector $y(T, \omega)$ fixed.

Let us use the theorem B.

Solving with given end conditions the system (27) together with (I) and equation

$$\int_{\Omega} [\Psi(t, \omega) h(t, \bar{v}, \omega) - h^0(t, \bar{v}, \omega)] P(d\omega) = \sup_{v \in Q} \int_{\Omega} [\Psi h(t, v, \omega) - h^0(t, v, \omega)] P(d\omega) \quad (3I)$$

we shall obtain the extremal $\bar{y}(t, \omega), \bar{v}, \Psi(t, \omega)$.

Let us see that this extremal is the absolute minimal of functional J.

Let $\Phi[t, y(\omega)] = \int_{\Omega} \Psi(t, \omega) y(\omega) P(d\omega)$ then term I of theorem B is fulfilled. In this case owing to (25) and (3I)

$$\begin{aligned} R[t, \bar{y}(t, \omega), \bar{v}] &= \int_{\Omega} \left[\left(\frac{d\Phi}{dt} + \Psi A - Q^0 \right) \bar{y}(t, \omega) - \Psi h(t, \bar{v}, \omega) - h^0(t, \bar{v}, \omega) \right] P(d\omega) = \\ &= \int_{\Omega} [\Psi h(t, \bar{v}, \omega) - h^0(t, \bar{v}, \omega)] P(d\omega) = \sup_{y, v} R[t, y(\omega), v] \end{aligned}$$

i.e. the second term of theorem B is also fulfilled.

8. About possibility to solve of first and second statements of original problem

We shall use the results of section 7 for solution of problem which can be considered as equivalent of the first or second statement of original problem (see section 3)

Let us preset the process over $[0, T]$ by the system

$$\begin{aligned}\dot{x} &= f(t, x, u) \\ \dot{z} &= Q(t, x, u)z + R(t, x, u)\xi(t, \omega)\end{aligned}\quad (32)$$

Let $x(0) = x_0$ and some coordinates of vector $x(\tau)$ be fixed.

The random influence $z(t, \omega)$ and $\xi(t, \omega)$ is given by its probability distributions. The performance characteristic of the process is

$$P\{ |Q(x(\tau))z(\tau, \omega)| < C \} \quad (33)$$

We shall maximize this probability by the choice of $u(t)$. Let us notice, that $P\{ |Qz(\tau, \omega)| < C \} = - \int_{\Omega} F(\alpha z(\tau, \omega)) P(d\omega)$, where $-F(p)$ is the characteristic functions of $[-C, C]$ interval. In usual sense this function is not differentiable, so we shall consider its derivative as a limit of approximate functions derivatives $\frac{dF_k(p)}{dp}$, as $k \rightarrow \infty$. Therefore it is necessary to consider all relations including the derivative of this function as limitary as $k \rightarrow \infty$.

Theorem B allow to look for optimal among those which minimize the integral

$$\int_{\Omega} \{ \psi(t, \omega) + \mu(t, \omega) [Qz(t, \omega) + R\xi(t, \omega)] \} P(d\omega)$$

with condition that row-vector functions ψ and μ , satisfy the equations

$$\dot{\psi}_j = - [\psi f_{x_j} + \mu (Q_{x_j} z + R_{x_j} \xi)] \quad j=1, \dots, n \quad (34)$$

$\psi_j(t, \omega) = - \frac{dF}{dp} \left(\frac{\partial a}{\partial x_j} z(t, \omega) \right)$, if coordinate $x_j(\tau)$ is not fixed,

$$\dot{\mu} = - \mu Q(t, x, u), \quad \mu(T, \omega) = - \frac{dF}{dp} a \quad (35)$$

It follows from (35) that $\mu(t, \omega) = \frac{dF}{dp} \bar{\mu}(t)$, where

$$\dot{\bar{\mu}} = - \bar{\mu} Q(t, x, u), \quad \bar{\mu}(T) = -a \quad (35')$$

Therefore the last integral is transformed into

$$\int_{\Omega} \left\{ \Psi f + \frac{dF}{d\mu} \bar{\mu} (QZ + R\bar{\xi}(t, \omega)) \right\} P(d\omega) = \int_{\Omega} \Psi(t, \omega) P(d\omega) \cdot f + \\ + \bar{\mu} \left\{ Q \int_{\Omega} \frac{dF}{d\mu} Z(t, \omega) P(d\omega) + R \int_{\Omega} \frac{dF}{d\mu} \bar{\xi}(t, \omega) P(d\omega) \right\} = \\ = \bar{\Psi}(t) f + \bar{\mu} (Q \bar{Z}_c(t) + R \bar{\xi}_c(t)) \quad (36)$$

Here $\bar{\Psi}(t) = \int_{\Omega} \Psi(t, \omega) P(d\omega)$; $\bar{Z}_c(t) = \int_{\Omega} \frac{dF}{d\mu} Z(t, \omega) P(d\omega) =$

$$= \gamma(c) M Z / \frac{dZ(t, \omega)}{dZ(t, \omega)=c} - \gamma(-c) M Z / \frac{dZ(t, \omega)}{dZ(t, \omega)=-c}, \quad \bar{\xi}_c(t) = \gamma(c) M \bar{\xi} / \frac{dZ(t, \omega)}{dZ(t, \omega)=c} - \gamma(-c) M \bar{\xi} / \frac{dZ(t, \omega)}{dZ(t, \omega)=-c}$$

where $\gamma(c)$ and $\gamma(-c)$ is the values of probability density of random variable $\chi = QZ(t, \omega)$, when $\chi = QZ(t, \omega) = c$ and $\chi = -c$. It follows from (32) that

$$\dot{\bar{\Psi}}_j = - [\bar{\Psi} f_{x_j} + \bar{\mu} (Q_{x_j} \bar{Z}_c + R_{x_j} \bar{\xi}_c(t))] \quad (37)$$

$$j = 1, 2, \dots, A$$

$\bar{\Psi}_j(\tau) = -\frac{\partial Q}{\partial x_j} \bar{Z}_c(\tau)$, if $x_j(\tau)$ is not fixed, and from equations of the process

$$\dot{\bar{Z}}_c = Q \bar{Z}_c + R \bar{\xi}_c(t) \quad (38)$$

Thus we have had the possibility to calculate the components of Hamiltonian by integrating of deterministic system of equations (35'), (37) and (38) (boundary conditions for them are not stochastic also). The difficulty of its integration is that we don't know the conditional expectation $\bar{Z}_c(0)$ and $\bar{\xi}_c(t)$ from right parts of (37) and (38) due to their values depend on the choice of control $u(t)$ over all interval of motion.

Now we shall consider one of those cases when this difficulty can be overcome. Let $\xi(t, \omega) = \xi(\omega)$ and $Z(\omega)$ are centred random vectors (this case included such a class of perturbations met on practice, when it is possible to consider them depending on some finite set of parameters, i.e. random variables) with the contour lines of joint distribution of $Z(0)$ and ξ being hyperspheres. Then due to linearity of system for \bar{Z} , $\bar{\xi}$ and $\bar{\mu}$ we have $\bar{\mu}(t) = -a A(t)$, $\bar{Z}_c(t) = B(t) \bar{Z}_c(0) + S(t) \bar{\xi}_c$

Here $A(t)$ and $B(t)$ are matrices of fundamental systems of solutions for (35') and (36) respectively, with $A(T) = B(0) = E$ (identity matrix). $S(t)$ is a matrix of particular solutions for (38) when initial conditions are zero ones, S^k , k -th column of S , is the solution that correspond the influence of k -th component of vector ξ . Now owing to distribution properties mentioned

$$\bar{Z}_c(0) = 2\gamma(c) M Z(0) \Big|_{x=c} = \frac{2\gamma(c) a B(T)}{\sqrt{a B(T) B'(T) a' + a S(T) S'(T) a'}}$$

$$\bar{\xi}_c = \frac{2\gamma(c) a S(T)}{\sqrt{a B(T) B'(T) a' + a S(T) S'(T) a'}}$$

Due to it the Hamiltonian (35) is equal

$$\begin{aligned} \bar{\Psi} f - \sum_{i=1}^n \frac{a A(t) Q B^i(t) 2\gamma(c) a B^i(T)}{\sqrt{a B(T) B'(T) a' + a S(T) S'(T) a'}} - \sum_{k=1}^m a A(t) \frac{Q S^k(t) 2\gamma a S^k(T) + R^k 2\gamma a S^k(T)}{\sqrt{a B(T) B'(T) a' + a S(T) S'(T) a'}} \\ = \bar{\Psi} f - \sum_{i=1}^n \frac{2\gamma(c) a B^i(T)}{\sqrt{\dots}} a A(t) Q B^i(t) - \sum_{k=1}^m \frac{2\gamma(c) a S^k(T)}{\sqrt{\dots}} [Q S^k(t) + R^k] \end{aligned}$$

If we set

$$\begin{aligned} \dot{m}^j = m^j Q(t, x, u), \quad j=1, \dots, n \\ m^j(T) = \frac{-2\gamma(c) B^j(T) a}{\sqrt{a B(T) B'(T) a' + a S(T) S'(T) a'}} \end{aligned} \quad (39)$$

$$B^j = Q(t, x, u) B^j, \quad j=1, \dots, n \quad B^j(0) = E^j \quad (40)$$

$$\dot{S}^k = Q(t, x, u) S^k + R^k(t, x, u), \quad S^k(0) = 0 \quad (41)$$

$$\begin{aligned} \dot{\gamma}^k = -\gamma^k Q(t, x, u), \quad \gamma^k(T) = \frac{-2\gamma(c) a S^k(T)}{\sqrt{a B(T) B'(T) a' + a S(T) S'(T) a'}} \\ k=1, 2, \dots, m \end{aligned} \quad (42)$$

$$\begin{aligned} \dot{\Psi}_j = -[\bar{\Psi} f_{x_j} + \sum_{i=1}^n m^j Q_{x_j} B^i + \sum_{k=1}^m \gamma^k [Q_{x_j} S^k + R^k_{x_j}]] \quad j=1, \dots, n \\ \bar{\Psi}_j(T) = \frac{-2\gamma(c) \frac{\partial a}{\partial x_j} [B(T) B'(T) a' + S(T) S'(T) a']}{\sqrt{a B(T) B'(T) a' + a S(T) S'(T) a'}} \quad c \neq x_j(T) \end{aligned} \quad (43)$$

is not fixed, it follows the final expression for (35)

$$\mathcal{H} = \bar{\Psi} f + \sum_{i=1}^n m^j Q B^j + \sum_{k=1}^m \gamma^k [Q S^k + R^k] \quad (35')$$

with optimal control being maximized this Hamiltonian for each fixed t .

Now we restrict ourselves to the solution of such applied problems mathematical model of which correspond the two first formulations of original one.

We reduce them to the solution of some boundary problem for the systems (39) - (43), (35').

References

1. Летов А.М. (Letov A.M.) Теория оптимальных систем. Труды II Конгресса ИФАС . Изд. Наука, 1965.
2. Летов А.М. (Letov A.M.) Аналитическое конструирование регуляторов. Автоматика и телемеханика, т. XXI, № 4, 5, 6; т. XXII, № 4; т. XXIII, № II
3. Красовский Н.Н. (Krasovskiy N.N.). Теория управления движением. Наука, 1968.
4. Кротов В.Ф. (Krotov V.F.). Докторская диссертация, 1963.
5. Колмогоров А.Н. (Kolmogorov A.N). Основные понятия теории вероятностей. ОНТИ, 1936.
6. Гихман И.И. (Gichman I.I.), Скороход А.В. (Skorokhod A.V.) Введение в теорию случайных процессов. Изд. "Наука", 1965.
7. Зубов В.И. (Zoobov V.I.). Теория оптимального управления. "Судостроение", 1966.
8. Вайнберг М.М. (Vienbergh M.M.). Вариационные методы исследования нелинейных операторов, ГИТТИ, 1956.
9. Пономарёв В.М. (Ponomarev V.M.). Теория управления движением космических аппаратов. Изд. "Наука", 1965.
10. Кожевников Ю.В. (Kogevnikov J.V.). ПММ т. 30 , № 4, 1966

STATISTICAL SYNTHESIS OF OPTIMAL PULSE
CONTROL SYSTEMS WITH REGARD TO SYSTEM'S
STRUCTURE CONSTRAINTS

A.Ya. Andrienko

Institute of Automation and Telemechanics
(Engineering Cybernetics)

Moscow

U S S R

Introduction

In development of automatic closed-loop systems with a digital control computer the designers try sometimes to reproduce control laws chosen earlier for continuous systems. These laws usually foresee the formation of signals action over instantaneous derivative of the controlled variable. As a result the requirements to the speed of a digital computer become more strict and the efficiency of use of computing technique is less, especially in cases when the controlled variables are measured with random errors. At the same time digital computers enable one to realize specifically discrete algorithms of control formation on the basis of analysis of case history of the control process. Apart from this it becomes possible to additionally improve quality of control at the expense of aptimization of time sequence of quantization intervals of a pulse system.

The requirements to the characteristics of a digital control computer as well as reliability of the whole control system are specified, to a considerable extent, by the system's structure complexity realized with the aid of the digital computer. Therefore, it is expedient to synthesize optimal systems with taking into account the constraints imposed on to the structure.

The paper discusses, in connection with the terminal control systems, the statistical methods of synthesis of pulse systems with taking into account the constraints of some basic kinds.

1. Synthesis of a Pulse Control System with Memory Capacity Constraint of Control Device

A controlled plant is considered whose output variable is measured at discrete time intervals Δ . The equations describing the plant are considered to be known.

$$x_{\nu(i+1)} = F_{\nu(i+1)}(\bar{x}_i, \bar{V}, u_i) \quad (1)$$

$$(\nu = 1, 2, \dots, N; i = 0, 1, \dots, I),$$

where $\bar{x}_i = (x_{1i}, x_{2i}, \dots, x_{Ni})$ - vector of plant variables at the i -th time moment;

$\bar{V} = (V_1, V_2, \dots, V_R)$ - random vector of disturbances affecting the plant;

u_i - value of controlled variable.

The value of y_i , of the plant output variable $x_i = x_{1i}$, measured with random error f_i , is fed to the input of the control device. In this case

$$y_i = x_i + f_i \quad (i = 1, 2, \dots, I). \quad (2)$$

The control should satisfy the inequality

$$|u_i| \leq U_i \quad (i = 0, 1, \dots, I), \quad (3)$$

where U_i - limit permissible value of the control u_i .

There are assumed to be set up the a priori densities $P(\bar{V})$, $P(\bar{x}_0)$, $P(f_i)$ of random vectors \bar{V} , \bar{x}_0 and of errors values of the measurement f_i .

The following risk function is assumed to be a criterion of optimality

$$W = M[w(x_{I+1})]. \quad (4)$$

The problem consists in defining the operation algorithm of the system's control device when a minimal value of risk is achieved; here the formation of the next control should be carried out with the use of limited volume of information on the state of control process so that control is the function of the form

$$u_i = u_i(\bar{y}_{ji}, \bar{u}_{j(i-1)}) \quad (i = 0, 1, \dots, I; j = i - n), \quad (5)$$

where

$$y_{ji} = (y_j, y_{j+1}, \dots, y_i),$$

$$u_{ji} = (u_j, u_{j+1}, \dots, u_i).$$

The memory capacity limitation is realized by setting up number n which characterizes the value of observation time interval.

The solution of the given problem is made on the basis of the theory of stochastic solutions and dynamic programming.

In formation of control at the i -th time moments there are regarded to be known, firstly, the control functions $u_s = u_s(\bar{y}_{ts}, \bar{u}_{t(s-1)})$ ($s = i+1, i+2, \dots, I; t = s-n$) at all subsequent time moments; secondly, control values u_s ($s = j, j+1, \dots, i-1$) and measured variables y_s ($s = j, j+1, \dots, i$) at preceding time moments; and, thirdly, control functions $u_s = u_s(\bar{y}_{ts}, \bar{u}_{t(s-1)})$ ($s = 1, 2, \dots, j; j = i-n; t = s-n$) at time moments preceding to the observation time interval. The latter permits to determine¹ from the a priori densities $P(\bar{V}), P(\bar{x}_0), P(f_s)$ ($s = 1, 2, \dots, j$) the distribution density of the so-called vector of the reduced disturbances $\bar{V}_s = (\bar{V}, \bar{x}_s; u_s = 0)$.

The set of equations (1) solved with respect to the output variable x_{i+1} , the vector \bar{V}_s being fixed, has the form

$$x_{i+1} = \psi_{i+1}(\bar{V}_s, \bar{u}_{si}) \quad (i = 0, 1, \dots, I; s = 0, 1, \dots, j; j = i-n). \quad (6)$$

The totality ($I-n+1$) of the auxiliary functions is introduced

$$\int_{\Omega(\bar{V}_s)} \bar{u}_{s(i-1)} d\Omega(\bar{V}_s) = \int_{\Omega(\bar{V}_s)} w[\psi_{i+1}(\bar{V}_s, \bar{u}_{si})] P(\bar{V}_s) \prod_{l=s}^I P(y_l/\bar{V}_s, \bar{u}_{li}) d\Omega(\bar{V}_s) \quad (s = 0, 1, \dots, I-n).$$

The optimal control $u_i = u_i^*$ is computed at sequential minimization with respect to u_I, u_{I-1}, \dots, u_i of the functions $\gamma_{jI}, \gamma_{(j-1)(I-1)}, \dots, \gamma_{ji}$ ($j = I - n$), so that

$$\gamma_{ji}^* = \min_{u_i \in \omega(u_i)} \gamma_{ji}(\bar{y}_{ji}, \bar{u}_{j(i-1)}; u_i) \quad (i = 0, 1, \dots, I),$$

where

$$\gamma_{ji} = \int \gamma_{jI}(\bar{y}_{ji}, \bar{u}_{ji}; u_{i+1}^*, u_{i+2}^*, \dots, u_I^*) d\Omega(\bar{y}_{(i+1)I}) \Omega(\bar{y}_{(i+1)I})$$

$\omega(u_i)$ - region of permissible controls with respect to (3).

For $n = I$ the totality of auxiliary functions degenerates into one function and the discussed methods coincide with the procedure of defining dual control.²

Similar methods can be presented for the case³ when the optimal increment $\Delta u_i = u_i - u_{i-1}$ of control is determined in the class of functions,

$$\Delta u_i = \Delta u_i(\bar{y}_{ji}, \bar{\Delta u}_{j(i-1)}) \quad (i = 0, 1, \dots, I; j = i - n), \quad (7)$$

where

$$\bar{\Delta u}_{ji} = (\Delta u_j, \Delta u_{j+1}, \dots, \Delta u_i)$$

The above given relations do not permit, in a general case, to obtain in the explicit form the algorithm of operation of the system's control device. Therefore we shall discuss the method of construction of a suboptimal control system, the system with independent identification of a plant.

Let us divide the coordinates of the vector \bar{V}_j into two groups.

$$\bar{V}_j = (\bar{V}_j^{(1)}, \bar{V}_j^{(2)})$$

The random vector $\bar{V}_j^{(1)}$ consists of such variables of the vector \bar{V}_j , whose action on the plant leads to that the plant loses the property of neutrality².

Generally, the vector $\bar{V}_j^{(1)}$ involves random deviations of controlled plant parameters. The vector $\bar{V}_j^{(2)}$ consists of the rest variables of the vector of reduced disturbances.

Replace the controlled plant, described by (6), by its model variable in discrete time

$$x_{i+s} = \psi_{i+s}(\bar{V}_{ji}^*, \bar{V}_j^{(2)}, \bar{u}_{j(i+s-1)}) \quad (8)$$

$$(i = 0, 1, \dots, I; s = 1, 2, \dots, I-i+1; j = i-n).$$

Here, vector \bar{V}_{ji}^* , the estimate of the vector $\bar{V}_j^{(1)}$, is assumed to be known but successively corrected at each interval i in the result of minimization of a certain adopted risk function.

$$W_i^{(1)} = M\{w^{(1)}[\bar{V}_j^{(1)}, \bar{V}_{ji}^*(\bar{y}_{ji}, \bar{u}_{j(i-1)})]\}$$

$$(i = 0, 1, \dots, I).$$

Correspondingly, the optimality criterion (4) is replaced by the sequence of risk functions.

$$W_i = M[w(x_{i+s})] \quad (i = 0, 1, \dots, I; S = I-i+1).$$

The control plant model (8) refers to neutral plants. Determination of optimal, with respect to risk W_i , control of this model is made sufficiently simply according to the discussed relationships. Here the obtained control function does not depend on the function of studying the vector $\bar{V}_j^{(1)}$ at intervals $i+1, i+2, \dots, I$.

2. Approximate synthesis of Systems Limited in Number of Devices Reproducing Coefficients of Control Algorithm

The system's accuracy estimated by the risk function (4) in optimal control can be improved either at the expense of increasing the number quantization intervals I of a pulse system or at the expense of optimization of program of variation in time of quantization intervals. In certain cases the efficiency of intervals optimization turns out to be the same as that achieved at increase of the number I equal intervals by several orders. Therefore a problem may be formulated of

program optimization of quantization intervals variation, which is formulated as the problem of determining the optimal "control" coordinates (i.e. intervals), independent of current values of coordinates of the pulse system. This problem has already been solved on the basis of a specific method of statistical optimization.⁴

While synthesizing a system with limited number of coefficients of the control algorithm it is necessary, in addition to the original data (1-5), to set the number L of permissible variations of algorithm's coefficients and numbers of intervals i_k ($i_1 = 0$; $k = 1, 2, \dots, L$), at which these variations are permitted.

Introduce into consideration L controlled plants with smaller numbers of quantization intervals which are described by equations

$$x_{i+1}^{(K)} = \psi_{i+1}(\bar{v}_{j_k}, u_{j_k}^{(K)}) \quad (i = j_k, j_k + 1, \dots, I; j_k = i_k - n), \quad (9)$$

obtained from (6). The plants are affected by the vectors \bar{v}_{j_k} which coincide with the vectors of reduced disturbances in the desired system. Define for each plant (9) the optimal function $u_i^{(K)} = u_i^{(K)}(\bar{y}_{j_k}^{(K)}, \bar{u}_{j_k}^{(K)}(i-1))$ without taking into account limitations with respect to the number L of permissible variations of coefficients. In a case, when quantization intervals of the original system vary according to optimal program the desired control function at intervals $(i_{k-1} \div i_k)$ approximately coincides with optimal function of control of the K -th plant (9) at the interval i_k :

$$u_{i_k}(\bar{y}_{j_k i_k}, \bar{u}_{j_k}(i_k-1)) \equiv u_{i_k}^{(K)}(\bar{y}_{j_k}^{(K)}, \bar{u}_{j_k}^{(K)}(i_k-1)).$$

This control can be practically calculated by iterative methods with the use of the previously given relationships.

3. Control System Synthesis with Taking
into Account the Constraints with Respect
to Operations Which Can be Realized in a
Control Device

A control plant is considered which consists of serially connected nonlinear inertia-free part, described by the function

$$v_i = v_i(u_i) \quad (i=0, 1, \dots, I), \quad (10)$$

and that under the influence of random disturbances of the linear part described by the equations

$$x_{\nu}(i+1) = \sum_{\mu=1}^N a_{\nu\mu} x_{\mu} + \sum_{\nu=1}^R b_{\nu i} v_{\nu} + c_{\nu i} (1 + \delta c) v_i \quad (\nu=1, 2, \dots, N), \quad (11)$$

where v_i - value of the output variable of the nonlinear part of the system,

δc - random deviation of the coefficient $c_{\nu i}$

Assume that the a priori densities of the vectors \bar{v} , \bar{x}_0 , errors of measurements f_i and the value of δc are described by the normal distribution laws.

As a criterion of optimality the risk function is adopted

$$W = M(x_{I+1}^2).$$

Regarding the peculiarities of digital computers it is reasonable to adopt algebraic operations as operations which can be realized in the control device. Considering at first the case when dispersion D_c of the value δc is equal to zero we shall define optimal control in the class of linear functions of the type

$$u_i = \sum_{j=1}^n A_{je} y_e + \sum_{j=1}^n B_{je} v_e \quad (i=1, 2, \dots, I; j=i-n). \quad (13)$$

For the linear part of the plant (11) the equation (6) has the form

$$x_{i+1} = \sum_{q=1}^q \alpha_{siq} V_{sq} + \sum_{\ell=1}^i \beta_{i\ell} v_{\ell} \\ (i=1, 2, \dots, I; s=1, 2, \dots, j; j=i-n).$$

The conditional mathematical expectation $M(x_{I+1} / \bar{y}_{ji}, \bar{v}_{j(i-1)}; \bar{v}_{iI} = 0)$ is the known linear function of the measured variables of the system⁵

$$M(x_{I+1} / \bar{y}_{ji}, \bar{v}_{j(i-1)}; \bar{v}_{iI} = 0) = \sum_{j=1}^n A_{je}^* y_e + \sum_{j=1}^{i-1} B_{je}^* v_e.$$

It can be shown that the desired coefficients of the control algorithm (12) must satisfy the following conditions

$$\begin{aligned}
 2A_{je}^* M(\dot{y}_e^2) + \sum_{\substack{k=j \\ k \neq e}}^i A_{jk}^* M(y_e y_k) + \sum_{\substack{k=j \\ k \neq e}}^{i-1} B_{jk}^* M(y_e v_k) = \\
 = -2\beta_{iI} M[v_i(u_i) y_e] \quad (l=j, j+1, \dots, i) \\
 2B_{je}^* M(v_e^2) + \sum_{\substack{k=j \\ k \neq e}}^i A_{jk}^* M(v_e y_k) + \sum_{\substack{k=j \\ k \neq e}}^{i-1} B_{jk}^* M(v_e v_k) = \\
 = -2\beta_{iI} M[v_i(u_i) v_e] \quad (l=j, j+1, \dots, i-1).
 \end{aligned} \tag{14}$$

The correlation moments of these relations should be calculated with the use of specific methods.⁶ The values satisfying the conditions (14) can be determined by iterative methods.

In case of $D_c \neq 0$ there can be set the problem of approximate determination of optimal control in the class of rational functions. This control can be calculated from the relation

$$u_i = K_i M(x_{I+1}/\bar{y}_{ji}, \bar{v}_{j(i-1)}; \bar{v}_{iI}=0; \delta c = \delta c^*),$$

where the coefficients K_i are determined from the condition

$$\begin{aligned}
 -\beta_{iI} M[M(x_{I+1}/\bar{y}_{ji}, \bar{v}_{j(i-1)}; \bar{v}_{iI}=0; \delta c = \delta c^*) v_i(u_i)] = \\
 = M\{[M(x_{I+1}/\bar{y}_{ji}, \bar{v}_{j(i-1)}; \bar{v}_{iI}=0; \delta c = \delta c^*)]^2\},
 \end{aligned}$$

and the value δc^* is calculated from the condition of minimization of the risk function

$$R_{ci} = M\{[\delta c - \delta c^*(\bar{y}_{ji}, \bar{v}_{j(i-1)})]^2\}.$$

4. Statistical Synthesis of Invariant Pulse Systems with Taking into Account Structure Constraints

In a number of cases the accuracy of a system obtained in statistical synthesis with taking into account structure constraints turns out to be insufficiently high comparing with ultimate accuracy possibilities of the system designed without taking into account the constraints.

Improvement of control accuracy can be achieved either at the expense of weakening of requirements with respect to structure constraint or at the expense of use of the principle of invariancy in synthesis of the system.

Let us assume that a control plant is described by the equations (10), (11). Set the problem of finding the control function minimizing the risk function (12) under the condition of minimal influence on the risk of variation of parameters of the distribution law of separate disturbances V_1, V_2, \dots, V_r . The control should be found in the class of functions linear with respect to variables of the vectors $\bar{y}_{ji}, \bar{v}_{j(i-1)}$.

Such a control at the first r intervals can be found from the condition

$$M_i^* = \min_{u_i \in \omega(u_i)} [M_{i(I+1)}(x_{I+1}/\bar{y}_{ii}, \bar{v}_{i(i-1)}; \bar{v}_{iI}=0) + \\ + v_i(u_i) \sum_{k=i}^I c_{ik}] \quad (i=1, 2, \dots, r), \quad (15)$$

where in calculation of the value $M_{iI}(x_I/\bar{y}_{ii}, \bar{v}_{i(i-1)}; \bar{v}_{i(I-1)}=0)$ it is assumed that the measurement errors at these intervals are equal to zero.

Introduce at the control interval $(r+1) \div (I+1)$ a new controlled variable

$$y_i^* = y_i - x_i^*,$$

where the values x_i^* are calculated from the relation

$$x_i^* = M_{ri}(x_i/\bar{y}_{ri}, \bar{v}_{r(i-1)}; \bar{v}_{ri}=0) + v_r^* \sum_{k=r}^I c_{ik}$$

and the value v_r^* is determined so that $x_{I+1}^* = 0$.

The desired control at the intervals $i = \tau+1, \tau+2, \dots, I$ can be defined from the condition

$$M_i^* = \min_{u_i \in \omega(u_i)} \{ M(x_{I+1} / \bar{y}_{j,i}, \Delta \bar{v}_{j,i(i-1)}; \Delta \bar{v}_{iI} = 0) + \quad (16) \\ + [\bar{v}_i(u_i) - \bar{v}_\tau] \sum_{k=1}^I c_{1k} \quad (i = \tau+1, \tau+2, \dots, I; \\ j_1 = i - n - \tau), \}$$

where

$$\Delta \bar{v}_{j,i} = (\Delta \bar{v}_{\tau j_1}, \Delta \bar{v}_{\tau(j_1+1)}, \dots, \Delta \bar{v}_{\tau i}), \quad \Delta \bar{v}_{\tau k} = \bar{v}_k - \bar{v}_\tau.$$

The conditions (15) and (16) lead to the relationships of the type (14), which can be solved by iterative methods.

R e f e r e n c e s

1. V.S. Pugachev. Teoriya sluchainykh funktsii i yego primeneniye k zadacham avtomaticheskogo upravleniya. Fizmatgiz, 1962.
2. A.A. Fel'dbaum. Osnovy teorii optimal'nykh avtomaticheskikh sistem. Fizmatgiz, 1963.
3. A.Ya. Andrienko. Statisticheskii sintez optimal'nykh impul'snykh sistem upravleniya s ogranicheniem po yomkosti pamyati upravlyaushchego ustroystva. "Avtomatika i Telemekhanika", N.7, 1968.
4. A.Ya. Andrienko. Metod statisticheskoi optimizatsii nelineynykh impul'snykh sistem avtomaticheskogo upravleniya. Technicheskaya kibernetika, N.4, 1967.
5. Yu.V. Linnik. Metod naimen'shikh kvadratov i osnovy matematicheskoi teorii obrabotki nablyudenii, Fizmatgiz, 1958.
6. A.Ya. Andrienko, B.N. Petrov, Yu.P. Portnov-Sokolov. Otsenka tochnosti sistem upravleniya kosmicheskikh ob'ektov.

Trudy II simpoziuma IFAC po upravleniyu v kosmose, 1967.

OPTIMAL CONTROL SYSTEM FOR STATIONARY
ARTIFICIAL CIRCUMTERRESTRIAL SATELLITE
ORBIT

A.A. Lebedev, M.N. Krasilshchikov, V.V. Malishev
Moscow, USSR

- Solution of technical problem of optimization the stationary artificial circumterrestrial satellite (SACS) process of transfer to a given position with necessary accuracy at minimum energetical loss is being considered in this report.

By SACS we mean 24-hr equatorial satellite. For an observer on Earth such satellite will seem motionless. Thanks to this particularity the satellite can be used at global communications' system creation. A number of technical reasons do not permit to transfer the satellite directly to the longitude required. To transfer SACS to a required longitude method of transfer ellipse is used¹ (Figure 1); the essence of this method is in the fact that satellite moves at elliptical orbit with apogee (perigee) at altitude corresponding to the radius of stationary orbit, the period of revolution being less (more) than 24-hr. Because of it satellite will drift to East (West). Very high final accuracy required is particular for this system, which leads to the necessity to create control algorithm on feedback principle. This feedback is realized by definition of SACS's angle (azimuth) α_u and SACS's distance r_u in coordinates of ground tracking station situated at any point on Earth. Since final accuracy and total energy loss depend not only on control algorithm but on satellite initial insertion velocity as well, it is interesting to solve also the problem of optimal calculated velocity of transfer the satellite to the altitude of stationary orbit. The given problem is the one of control based on incomplete data and its exact solution on this data is difficult. Then the approximate method is used², the synthesis problem being divided into two independent problems. The first one

consists of transfer process optimization in assumption that all the necessary for control data are available. The second one is to get this information by optimal processing the data of measurement. Principal assumptions at solution of formulated problems are given below: 1) flat problem is being considered, that is satellite deviations at latitude are being neglected; 2) initial transfer errors at realization of corrective impulses is supposed further on that control is realized by velocity impulses) and errors of phase coordinates determination are introduced into discussion as irritations.

I. Let us take into account vector of SACS's phase coordinates x , the components of which being x_1 - satellite deviation from given position at longitude at the moment of passing through apogee (perigee); x_2 - satellite's drift velocity per one revolution at transfer ellipse (Figure 1). Equations of satellite motions in terms of above given coordinates are the following:

$$x_{i+1} = \Phi x_i + V(1 + \mu_i)u_i + V\epsilon_i, \quad (1)$$

$$(i = 0, 1, \dots, N)$$

where x_i is state vector x at the moment of apogee i ,

$$\Phi = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}, \quad V = \begin{bmatrix} 1 \\ 1 \end{bmatrix}. \quad (2)$$

Here u_i is calculated value of corrective velocity impulse i , μ_i , u_i , ϵ_i are multiplex and additive components of control influence error, N - number of revolutions at transfer ellipse orbit. It is supposed that μ_i and ϵ_i are accidental centered independent Gaussian numbers and

$$E[\mu_i^2] = \sigma_\mu^2, \quad E[\epsilon_i^2] = \sigma_\epsilon^2. \quad (3)$$

where E - denote mathematical expectation.

So μ_i and ε_i may be considered as discret white noise.

Mathematical formula of problem to search the optimal correction algorithm comes to the following. It is necessary to find such a sequence $\{u_i\}$ ($i = 1, \dots, N$), which assures the final precision required, characterized by given risk.

$$R_1 = E[x_{N+1}^T \lambda x_{N+1}] \quad (4)$$

at minimum energy loss measured by risk

$$R_2 = E\left[\sum_{i=1}^N u_i^2\right]. \quad (5)$$

solution of problem given has been received in³ and comes to the following. Algorithm of optimal correction looks like

$$u_i = -\lambda_i x_i, \quad (6)$$

where

$$\lambda_i = \Gamma_i^{-1} V^T \lambda_{i+1} \Phi, \quad \Gamma_i = V^T \lambda_{i+1} V (1 + \sigma_{\mu}^2) + \alpha \quad (7)$$

Matrix λ_i is determined according to the recurrent relation

$$\lambda_i = \Phi_i^T [\lambda_{i+1} - \lambda_{i+1} V \Gamma_i^{-1} V^T \lambda_{i+1}] \Phi_i \quad (8)$$

with an initial condition $\lambda_{N+1} = \lambda$.

Multiplier α should be determined by method of successive approximations or by graphic techniques on condition

$$\bar{x}_0^T \lambda_{10} \bar{x}_0 + C_{10} + Sp(\lambda_{10} K_{x0}) = R_1^* \quad (9)$$

where \bar{x}_0 is mathematical expectation of vector x_0 , $K_{x0} = E[x_0 x_0^T]$ being its covariance matrix, R_1^* , being the required value of risk R_1 .

Matrix λ_{10} and coefficient C_{10} are found by recurrent correlations

$$\lambda_{ii} = \Phi^T \lambda_{ii+1} \Phi - \Phi \lambda_{ii+1} V h_i - h_i^T V^T \lambda_{ii+1} \Phi + h_i^T \Gamma_{ii}^{-1} h_i \quad (10)$$

$$C_{ii} = C_{ii+1} + \sigma_e^2 Sp(\lambda_{ii+1} V V^T) \quad (11)$$

where $\Gamma_{ii} = V^T \lambda_{ii+1} V (1 + \sigma_e^2)$

with initial conditions $\lambda_{1N+1} = \lambda$, $C_{1N+1} = 0$.

Satellite drift optimal calculated velocity at its transfer to the altitude of stationary orbit according to ³ is equal

$$\bar{x}_{20opt} = -(\lambda_{22})_0^{-1} (\lambda_{21})_0 \bar{x}_{10} \quad (12)$$

where $(\lambda_{22})_0$, $(\lambda_{21})_0$ are corresponding elements of matrix λ_0 . As follows from the mentioned above, for optimal control one should know state vector x and time τ_d of SACS's passing through apogee. Since these components can't be measured, problem of their optimal estimation arises. System

(I) is not convenient for determination of optimal estimations of vector on base of ground measurements, because vector of measurements itself is not included in this system. That is why equations of satellite motion in rotating geocentric coordinate system, linearized relatively to reference circular orbit are used. On condition that measurement are discret, these equations is:

$$\mathbf{z}_k = A(t_k - t_{k-1}) \mathbf{z}_{k-1} + B(1 + \mu_{k-1}) \mathbf{u}_{k-1} + B \mathbf{e}_{k-1} \quad (I3)$$

where \mathbf{z}_k - state vector at the time t_k , its components being $z_1 = \Delta z$ - deviations of SACS altitude from the one corresponding to designed circular orbit, $z_2 = \Delta \lambda$ - SACS deviations at longitude from the position at calculated circular orbit, $z_3 = \Delta V_r$ - radial component of satellite velocity deviation, $z_4 = \Delta V_s$ - deviation of SACS's tangential velocity component, $A(t_k - t_{k-1})$ - state transition matrix for discrete equation of movement, corresponding to the reference circular orbit

$$A(\tau) = \begin{vmatrix} 4 + 3\cos\omega_0\tau & 0 & \frac{\sin\omega_0\tau}{\omega_0} & \frac{2}{\omega_0}(1 - \cos\omega_0\tau) \\ -\frac{6}{z_0}(\omega_0\tau - \sin\omega_0\tau) & 1 & -\frac{2}{\omega_0 z_0}(1 - \cos\omega_0\tau) & -\frac{3}{z_0} + \frac{4\sin\omega_0\tau}{z_0\omega_0} \\ 3\omega_0\sin\omega_0\tau & 0 & \cos\omega_0\tau & 2\sin\omega_0\tau \\ -6\omega_0(1 - \cos\omega_0\tau) & 0 & -2\sin\omega_0\tau & -3 + 4\cos\omega_0\tau \end{vmatrix} \quad (I4)$$

(I5)

$$B = \begin{vmatrix} 0 \\ 0 \\ 0 \\ 0 \end{vmatrix}, \quad \mathbf{0} = \frac{\omega_0 z_0}{2g}, \quad \tau = t_k - t_{k-1}$$

Measured coordinates λ_{ux} and z_{ux} are related with components of state vector \tilde{x} by equations:

$$\lambda_{ux} = azctg \left[tg \left(-\alpha_{10} - \Delta\lambda_x - \lambda_u + azctg \frac{R_3 \cos \vartheta \sin(-\alpha_{10} - \Delta\lambda_x - \lambda_u)}{z_0 + \Delta z_x - R_3 \cos \vartheta \cos(-\alpha_{10} - \Delta\lambda_x - \lambda_u)} \right) \right] \quad (I6)$$

$$\times \cos azctg \frac{R_3 \sin \vartheta}{z_0 + \Delta z_x - R_3 \cos \vartheta} \Big];$$

$$z_{ux} = \frac{\sqrt{R_3^2 \cos^2 \vartheta \sin^2(-\alpha_{10} - \Delta\lambda_x - \lambda_u) + (z_0 + \Delta z_x - R_3 \cos \vartheta \cos(-\alpha_{10} - \Delta\lambda_x - \lambda_u))^2}}{\cos azctg \frac{R_3 \sin \vartheta}{z_0 + \Delta z_x - R_3 \cos \vartheta}} \quad (I7)$$

Here R_3 , λ_u , ϑ_x - are respectively radius of Earth, longitude and latitude of ground tracking site. To receive a more simple estimatic algorithm equations (I6) and (I7) are linearized. To provide a better convergence of estimations at it, one linearization about previous optimal predicted estimation at a previous step of measurement is being used⁴. In this case it corresponds to linearization of equations (I6), (I7) in vicinity of transfer ellipse predicted on the base of optimal estimations after the next correction. Then (I6) and (I7) may be rewritten: can be expressed by

$$y_x = H_x \tilde{x}_x + v_x \quad (I8)$$

where y_x is vector of measurements; H_x is 2×4 matrix of partial derivatives, elements of which depend on time and on previous estimations; v_x is vector of measurement errors,

which are Gaussian independent random numbers with zero mean and $E[v_k v_k^T] = Q_v$.

Taking into account the above said, optimal estimations of state vector \hat{x}_k , are optimal estimations according to Bayess rule \hat{x}_k^* , \hat{z}_k and are determined by the following equations⁵

$$\hat{x}_k^* = \hat{P}_k (P_k')^{-1} \hat{z}_k + \hat{P}_k H_k^T Q_v^{-1} y_k. \quad (19)$$

$$\hat{P}_k = [(P_k')^{-1} + H_k^T Q_v^{-1} H_k]^{-1}.$$

$$P_k' = A_{k,k-1} \hat{P}_{k-1} A_{k,k-1}^T + B_{k-1} G_\delta^2 B_{k-1}^T + B_{k-1} G_\mu^2 u_k^2 B_{k-1}^T. \quad (20)$$

$$\hat{z}_k = A_{k,k-1} \hat{z}_{k-1} \quad (21)$$

$$(22)$$

If after cycle of measurements completed the estimation error of components \hat{x}_k not exceed the additive error of corrective impulse realization than the determination of moment of satellite passing through apogee takes place from the condition $\Delta V_z = 0$. In that case optimal estimation of τ_d is equal accuracy of linear approximation

$$\tau_d^* = \frac{1}{\omega_0} \arctg \left(- \frac{\Delta V_z^*}{\omega_0 \Delta \tau^* + 2 \Delta V_s^*} \right). \quad (23)$$

Estimation of components of vector x is being determined

for the time τ_d^* .

2. In the numerical treatment 24-hr satellite has been considered^{1,6} for which $\alpha_{10} = -35^\circ$. Parameters of stationary orbit were:

$$V_0 = \omega_0 r_0 = 3075 \text{ m/sec}; \quad r_0 = 42165 \text{ km}, \quad T_0 = 86164 \text{ sec.}$$

The transfer process is considered completed if the following conditions are fulfilled: 1) satellite is within the region of admissible longitudinal deviations $|\alpha_1| \leq \Delta \lambda_m$; 2) residual drift velocity α is such that staying of satellite within the region $|\alpha_1| \leq \Delta \lambda_m$ during the time t_0 is guaranteed. $\Delta \lambda_m = 2,5^\circ$; $t_0 = 60$ revolutions 60 days were assumed. Region of admissible final errors on phase plane (α_1, α_2) is shown at figure 2.

Let us approximate the boundary of region Γ by ellipse Γ'

$$\frac{1}{t_0^2} \alpha_1^2 + \frac{1}{t_0} \alpha_1 \alpha_2 + \alpha_2^2 = \frac{3}{4} \frac{\Delta \lambda_m^2}{t_0^2} \quad (24)$$

To provide the final admissible accuracy with probability 0,997 as the required value of risk R_1 , the value R_1^* should be taken

$$R_1^* = \frac{1}{9} \cdot \frac{3}{4} \frac{\Delta \lambda_m^2}{t_0^2} \quad (25)$$

As follows from (24), matrix λ in risk R_1 (4) is equal

$$\lambda = \begin{vmatrix} \frac{1}{t_0^2} & \frac{1}{2t_0} \\ \frac{1}{2t_0} & 1 \end{vmatrix} \quad (26)$$

Taking into consideration (25), (26) and using equations (7) - (10), calculations were carried out for different values of λ and N and $\sigma_{\mu}^2 = 25 \cdot 10^{-4}$; $\sigma_{\epsilon}^2 = 0.003 \text{ grad}^2 / \text{rev}^2$ results are given in figures 3, 4, 5.

The definition of the necessary interval of measurements and of their number per each revolution is given, considering the estimated error of the satellite velocity components should not exceed the additive error resulting from the application of the corrective velocity impulses.

On the grounds of the covariance matrix numerical calculations in accordance with equations (20), (21) for $\sigma_{\mu} = 0$, the interval between measurements was taken equal to 1 hr, their number being 5. Curves characterizing the convergence of state vector Z coordinate are given in fig. 6. It was assumed $\sigma_{\Delta u} = 0.003^\circ$; $\sigma_{\Delta v} = 25 \text{ km.}$, the initial launching errors being $\sigma_{\Delta z} = 50 \text{ km.}$; $\sigma_{\Delta \lambda} = 1^\circ$; $\sigma_{\Delta v_1} = \sigma_{\Delta v_2} = 1.7 \text{ m/sec} / 1/2 / 6/$.

3. In order to examine the actual possibilities of SACS's satisfactory transfer to a given position using the suggested control algorithm, the simulation of the closed-loop control system by Monte Carlo method was applied. When simulating, attitude and stabilization errors during the SACS correction were taken into account, assuming that they do not exceed the measurement error and the corrective impulse application error by more than $\pm 2^\circ$. The satellite dynamics was defined by the equations of the elliptical theory [7]. Furthermore it was assumed that Δu and Δv relations to vector Z components were linear.

When simulating the closed-loop system the following principal versions were examined:

1) supposing that the exact values of x_1 , x_2 and τ_1 were known, it was possible to check the correction algorithm. The dispersion pattern for this case is given in fig. 7. As may be seen from fig. 7, the dispersion near to normal and is in good agreement with the calculated dispersion;

2) if measurement errors are present, taking into consideration all the above mentioned factors. The dispersion pattern for this case is given in fig. 8. As it may be seen from fig. 8 the dispersion pattern is rather different from the previous one, but even in this case it is within the admissible limits.

References

1. Кент Г., Кенехен М.Е. Выведение на заданную долготу и управление орбитой 24-часового экваториального спутника Земли. "Ракетная техника и космонавтика", русский пер., 1964, № 6.
2. Ли Р. Оптимальные оценки, определение характеристик и управление. Изд. "Наука", 1966.
3. Малышев В.В. Задача об оптимальном дискретном управлении конечным состоянием линейной стохастической системы. "Автоматика и телемеханика", 1967, № 5.
4. Красильщиков М.Н. Некоторые вопросы определения достаточных координат при управлении динамическими объектами. Доклад на I Всесоюзном симпозиуме по статистическим проблемам в технической кибернетике. Москва, 1967.
5. Богуславский И.А. Об уравнениях стохастического управления. "Автоматика и телемеханика", 1966, № II.
6. Боучер Р.А. Использование электрических ракетных двигателей в системе управления стационарными спутниками. "Вопросы ракетной техники", 1965, № 4.
7. Балк М.Б. Элементы динамики космического полета. Изд. "Наука" 1965 .

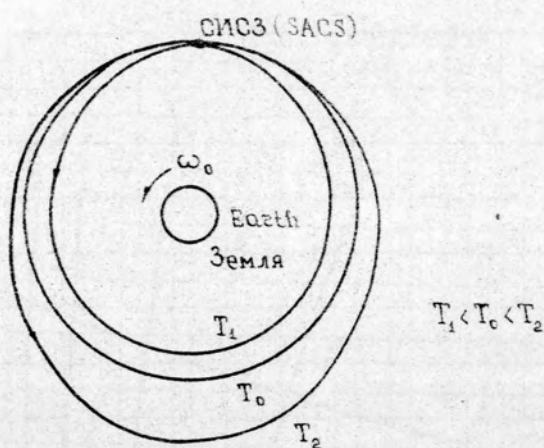


Figure 1

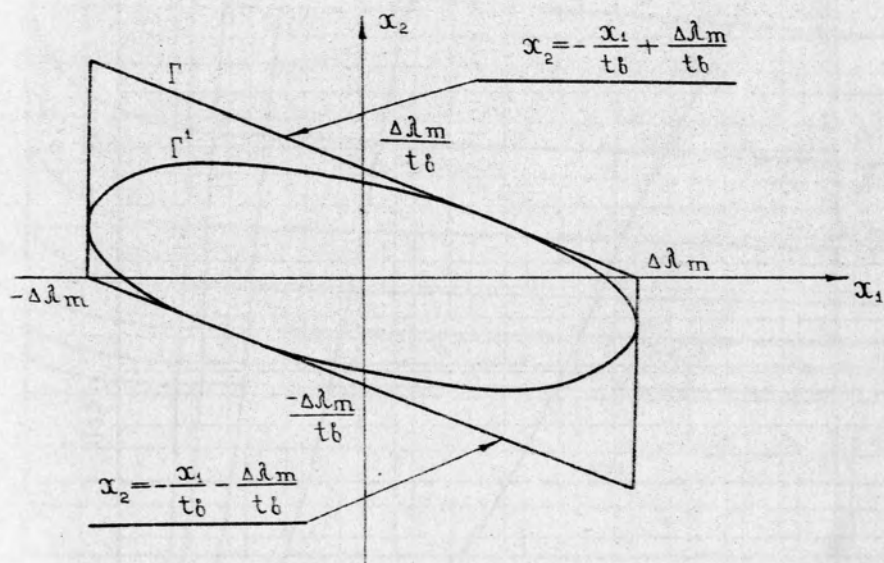


Figure 2

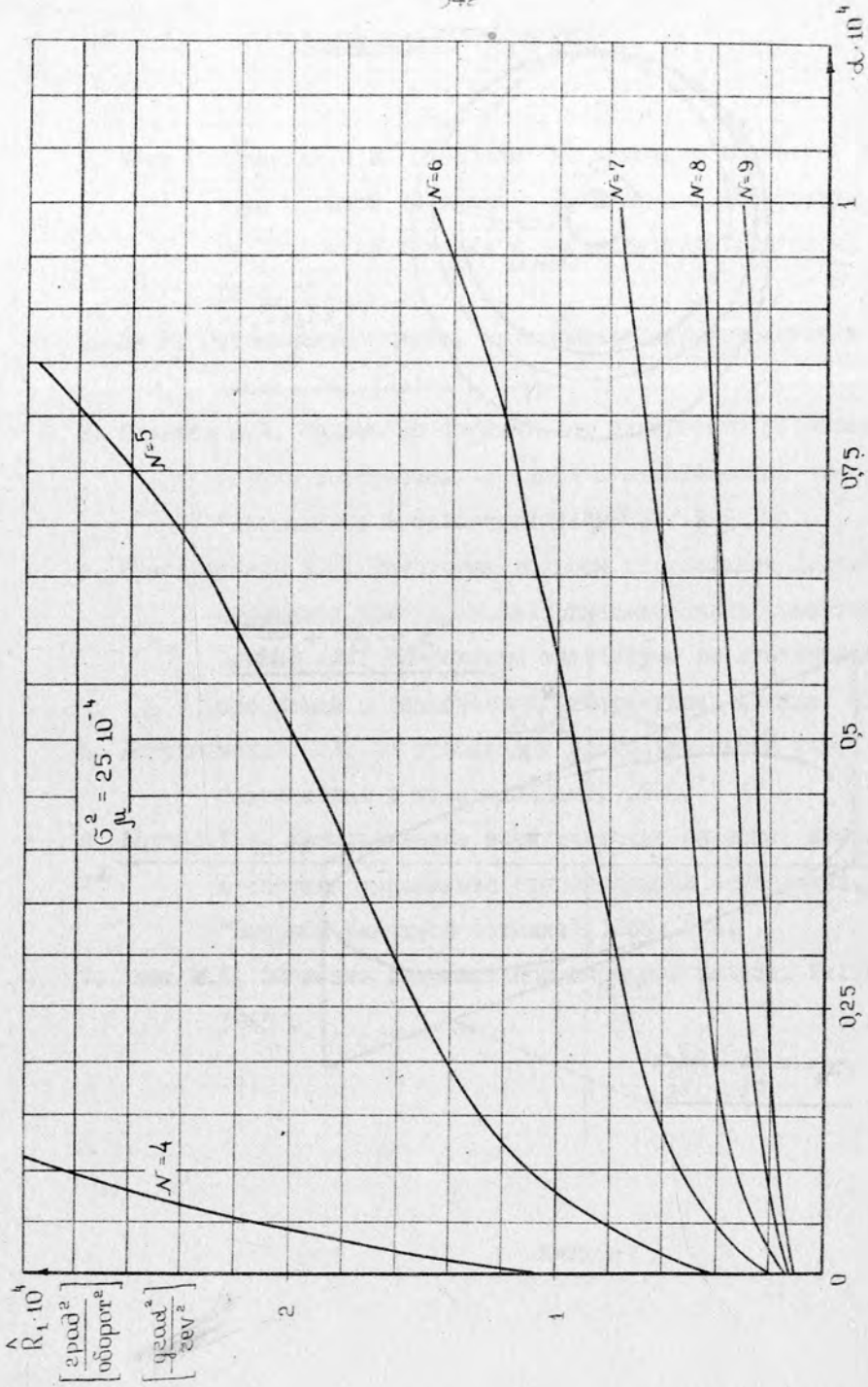


Figure 3

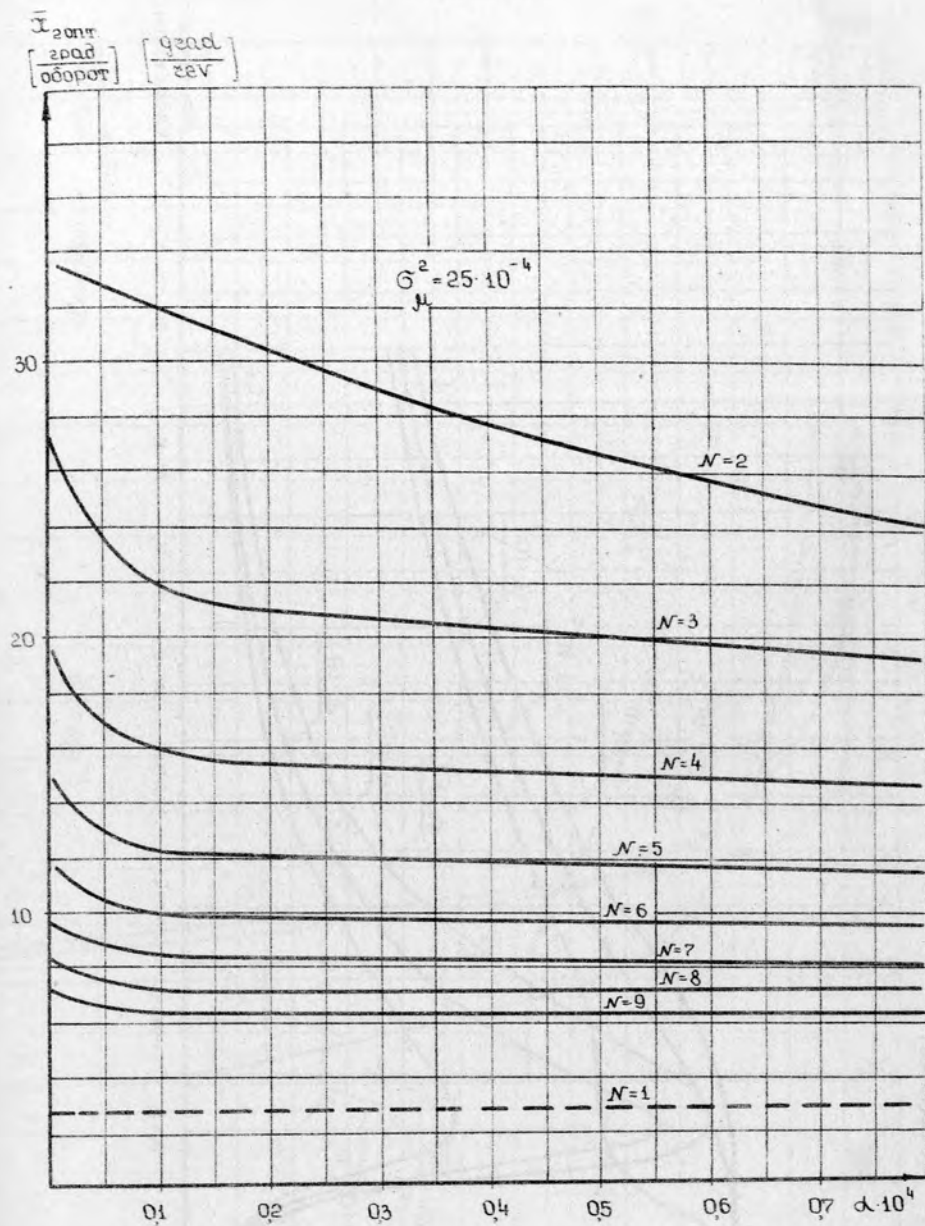


Figure 4

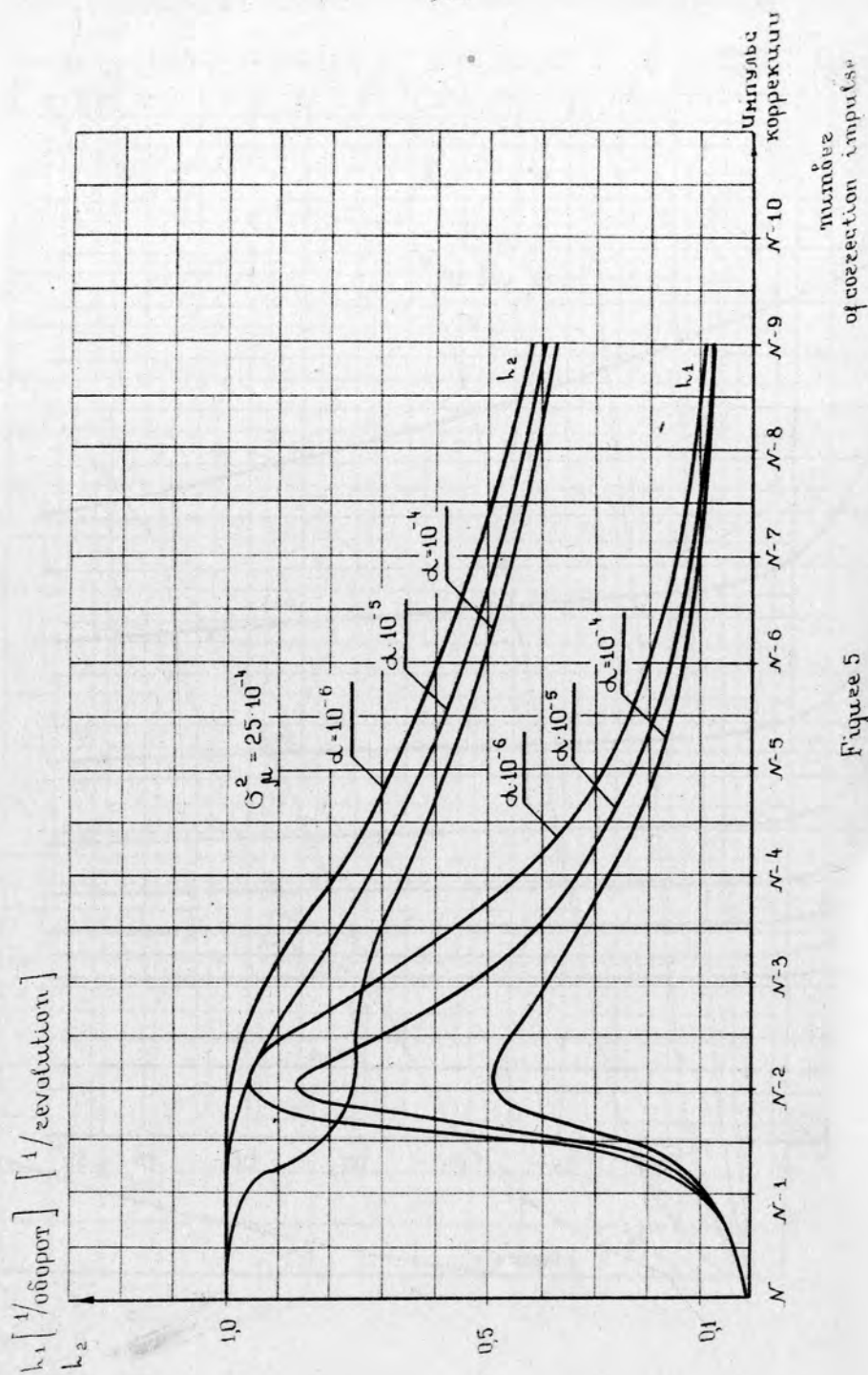
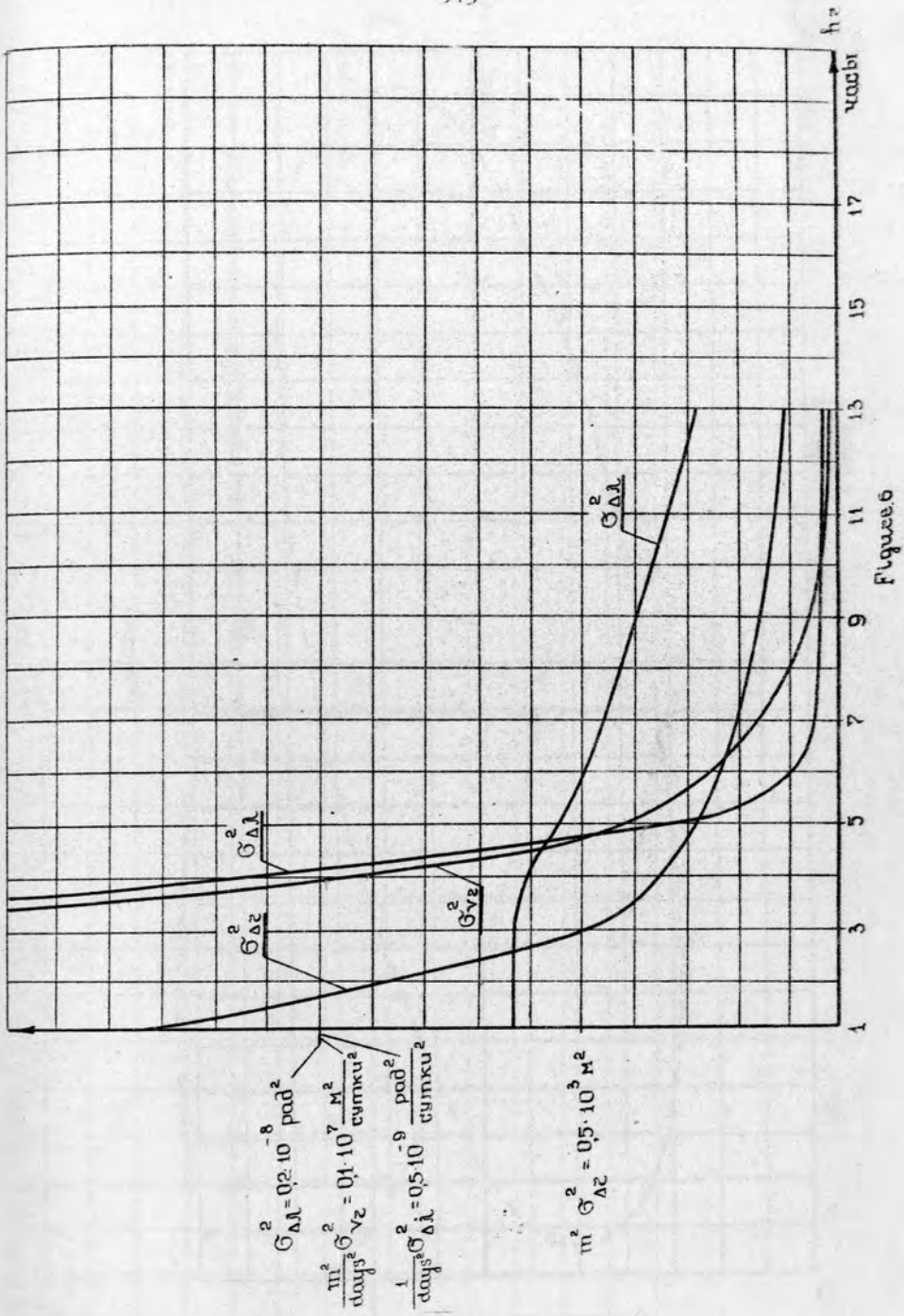
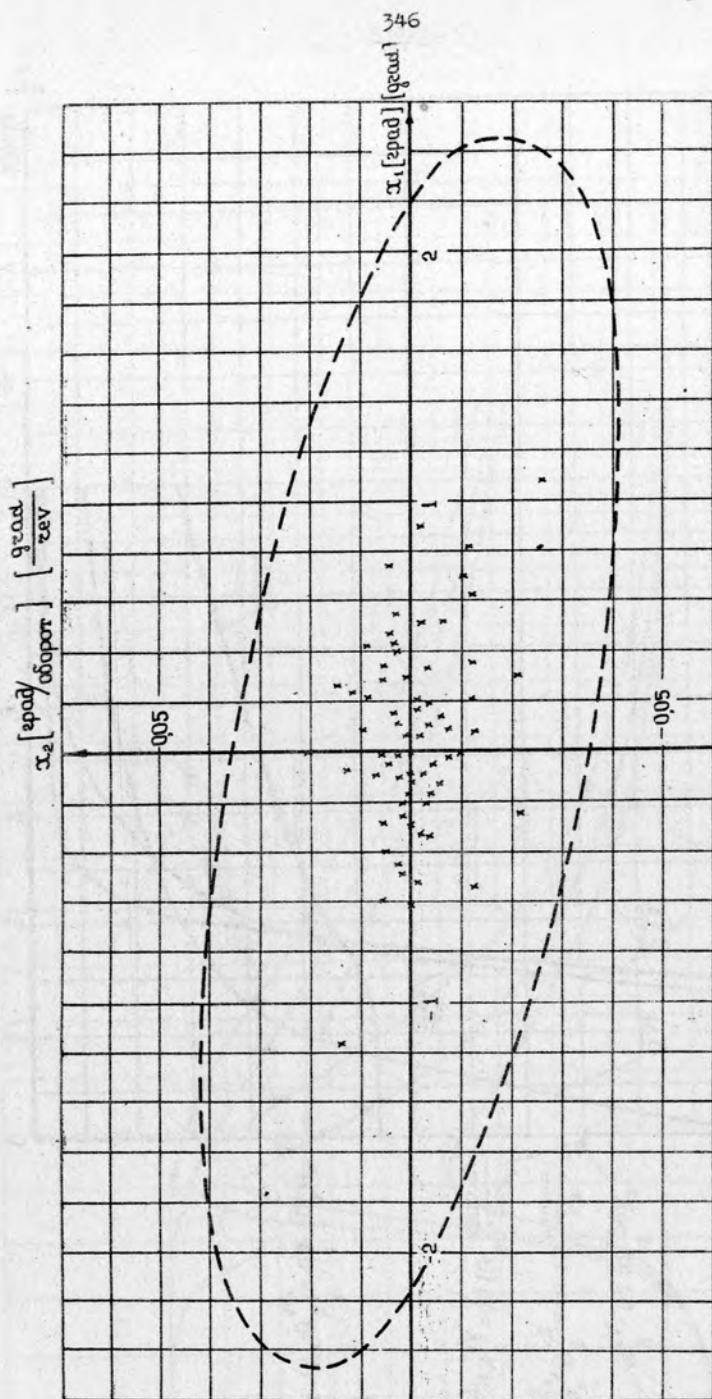


Figure 5





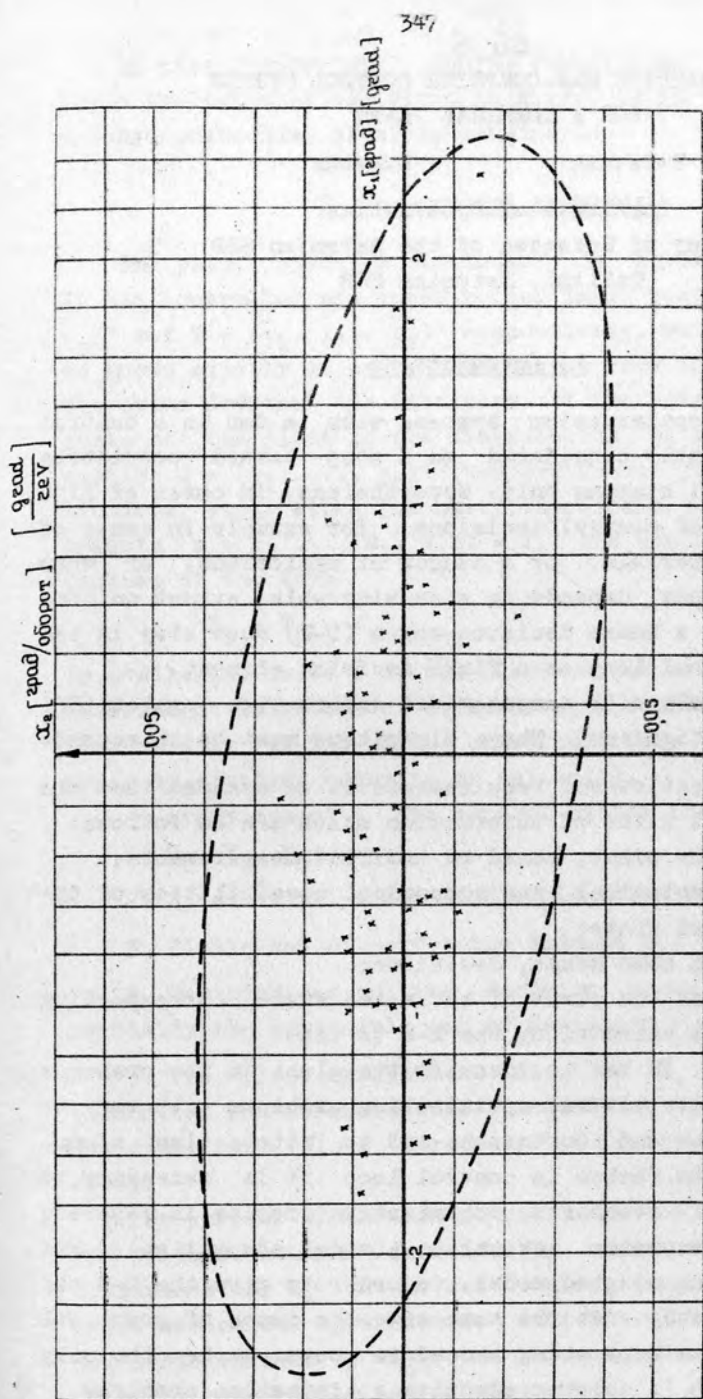


Figure 8

AN ADAPTIVE MAN-COMPUTER CONTROL SYSTEM
FOR A CHEMICAL PLANT

R.Tavast

L.Mytus

Institute of Cybernetics
Academy of Sciences of the Estonian SSR
Tallinn, Estonian SSR1. Introduction

Supervisory optimization systems with a man in a control loop are frequently considered as a step toward completely automatic control systems only. Nevertheless, in cases of high responsibility of control decisions, for example in cases of high costs of materials, or a danger of explosions, or when a plant performance depends to a considerable extent on non-formal factors, a human decision-maker (D-M) must stay in the supervisory control loop as a final decision element.

The paper deals with computerized information system (IS) algorithms investigation. These algorithms must be in accordance with D-M practice and requirements. It is assumed that the D-M needs several kinds of information which are as follows:

- state of the plant, based on indirect measurements;
- limit technological and economical possibilities of the plant of estimated state;
- optimal, in some sense, decisions;
- behaviour of the plant of estimated state, corresponding to various inputs selected by the D-M to test.

Consequently, IS has to identify the plant in the presence of noises, to solve several optimization problems with various performance indexes and constraints and to imitate plant operating. Due to human factor in control loop it is necessary to solve the overall stochastic optimization problem in separate steps, plant parameter estimation (model adaptation) and optimization using adapted model, in order to give the D-M the results of each step. At the same time, in cases of practical importance such a separating procedure seems to be the only possible approach in solving adaptive optimization problems.

In this report the problem formulation and its solution for a typical chemical plant is given. As an example the formaldehyde production plant is considered.

2. The control problem.

The plant operates continuously on a time interval $[0, T]$. It has controlled and uncontrolled input vectors $U = (u_1, \dots, u_r)'$ and $V = (v_1, \dots, v_s)'$ respectively, which are assumed to be known exactly at every time moment from interval considered. The prime denotes the transpose of the matrix. The physical state of the plant is the distribution of some physical properties $Q = (q_1, \dots, q_t)'$ (temperatures, pressures, concentrations, etc.) along the space coordinate η . At discrete time moments $n = 1, \dots, N_1$, $N_1 \Delta t = T$, the stochastic m -vector process of the form

$$Z_n = \bar{Y}_n + W_n \quad (1)$$

is available, where \bar{Y} is the plant output vector and W_n is the measurement noise vector sequence with bounded variances. The time varying of the plant characteristics is caused for example by catalyst decaying or poisoning.

At every moment n the control system knows:

1. U_n, V_n, Z_n ;
2. Hypothetical model of the plant $Y = F(U, V, \theta, n)$;
3. Finite set of performance indexes $w_i, i \in A = \{1, \dots, e\}$;
4. Restrictions on the decision space that determines partially the admissible set of decisions Γ_U .
5. The plant environment situation S_n characterized by economical and technological time-varying, partially ill-defined conditions (plans, arrangements, technological conditions, etc.).

The purpose of the man-computer control system is to make and adjust control decisions U_n at moment n to satisfy conditions in S_n in the best manner, in some sense. In most cases economical optimization criteria (performance indexes) are predetermined by S_n , but there is always some freedom in choosing admissible set of decisions Γ_U . Sometimes S_n may

change in a manner that demands the optimization criterion to be altered. The optimization problem formulation is the first task of the D-M. (It must be emphasized that the problem of mapping S_n into Γ_U and $w_i, i \in \Lambda$ is the psychological one and it is not discussed here). The task of IS algorithms is to solve several optimization problems formulated by the D-M for the time-varying plant. The final decision U_n at S_n must be made by the D-M again on the bases of the solutions, given by IS.

For getting additional information he can access to the computer.

Optimization problem. Find point U_i^* from admissible set $\Gamma_U \subset E^r$ to minimize performance index expectation

$$J_i = \mathcal{E} w_i(U_n, \bar{Y}_n), i \in \Lambda, \quad (2)$$

where w_i are known functions (profit, production yield, product quality, etc.). Admissible set Γ_U is of the form of joint probabilistic conditions

$$\Gamma_U = \left\{ U: P[G(\bar{U}_n, \bar{Y}_n, V_n, \underline{G}, \bar{G}) \leq 0] \geq \gamma_G \right\} \quad (3)$$

where $P[A]$ denotes probability of A , $1 - \gamma_G$ is the subjectively allowable probability of violating condition $G(\cdot) \leq 0$, the latter being an inequality system

$$\begin{aligned} g_1 &\leq g_1(U_n, \bar{Y}_n) \leq \bar{g}_1 \\ &\dots \dots \dots \dots \dots \dots \\ g_c &\leq g_c(U_n, \bar{Y}_n) \leq \bar{g}_c. \end{aligned} \quad (4)$$

Some of the elements of real positive vectors $\underline{G} = (g_1, \dots, g_c)$ $\bar{G} = (\bar{g}_1, \dots, \bar{g}_c)$ are fixed by S_n , some can be chosen by the D-M.

Optimization problem solving is the task of computer algorithms, which will be considered in the following section.

Final decision-making. The D-M can accept one of the $U^*, i \in \Lambda$ and make decision $U_n = U_{in}^*$, or interpolate between several solutions, or formulate a new problem. Before adjusting on the plant the final decision candidate will be tested on the model.

3. The plant model.

The stationary behaviour of a distributed plant is described by lumped parameter model of the following structure. The model has input vectors U, V that coincide with those of the plant. Output vector $Y = (y_1, \dots, y_m)'$ elements are known functions of the physical state Q at points η_1, \dots, η_v

$$y_j = h_j[Q(\eta_1), \dots, Q(\eta_v)], \quad j = 1, \dots, m. \quad (5)$$

Differential equation which has been constructed on the hypotheses concerning the plant

$$D[Q(\eta), U, V, \Delta] = 0, \quad (6)$$

involves space derivatives $\partial Q / \partial \eta$ and unknown parameter vector Δ , but does not involve time derivatives (heat and mass transfer equations in stationary conditions together with equations of chemical kinetics). Equation (6) subject to boundary conditions

$$B[Q(\eta_b), U, V] = 0, \quad (7)$$

where η_b are coordinates of plant boundary, determines the solution

$$q_k = \varphi_k(\eta, U, V, \Delta), \quad k = 1, \dots, t, \quad (8)$$

with U, V, Δ fixed. Substituting (8) into known functions h_j , $j = 1, \dots, m$ we have static operator as the plant model

$$Y = F(U, V, \Delta), \quad F = (f_1, \dots, f_m)'. \quad (9)$$

In all cases of practical importance we can obtain numerical solutions of (6), (7) only. The elements of parameter p_Δ -vector Δ are [from equation (6)] heat and mass transfer coefficients, chemical reaction rate coefficients, etc. Several elements of Δ may slowly vary in time. There is extensive literature concerning chemical plant modelling in the form (6), (7) (as, for instance in³) but theoretical models with time dependent parameters are seldom available. Hence we shall take some formal model approach

$$\Delta_n = C(n, \theta), \quad (10)$$

where θ is an p -vector of unknown real parameters and n is a time index as before. Some simple examples of (10) are

$$\Delta_n = \theta, \quad p_\Delta = p, \quad (11)$$

$$\Delta_n = H \begin{bmatrix} 1 \\ n \end{bmatrix}, \quad p = 2p_\Delta, \quad (12)$$

here H is an $p_\Delta \times 2$ matrix of unknown parameters, θ is known to lie in a domain Γ_θ :

$$\Gamma_\theta = \{\theta: \underline{\theta} \leq \theta \leq \bar{\theta}\}. \quad (13)$$

Taking (10) into consideration we have time dependent memoryless plant model

$$Y = F(U, V, \theta, n). \quad (14)$$

4. Multidimensional parameter estimation.

Consider the m -vectors residuals

$$\varepsilon_n = Z_n - F(U_n, V_n, \theta, n), \quad n = 1, \dots, N, \quad (15)$$

which are random independent vectors with unknown multidimensional distributions. It is assumed that at the beginning of the interval $[0, T]$ N sets of vectors $U_n, V_n, Z_n, n = 1, \dots, N$ as the basis of initial estimation, are available. Thus we are led to the problem of nonlinear simultaneous estimation of unknown parameters¹. We introduce following notations

$$\begin{aligned} X_n &= (U_n, V_n, n)' \\ Z_{jN} &= (z_{j1}, \dots, z_{jN})' \\ \varepsilon_{jN} &= (\varepsilon_{j1}, \dots, \varepsilon_{jN})' \\ F_{jN} &= [f_j(X_1, \theta), \dots, f_j(X_N, \theta)]' \\ Z_N &= (Z'_{1N}, \dots, Z'_{mN})' \\ \varepsilon_N &= (\varepsilon'_{1N}, \dots, \varepsilon'_{mN})' \\ F_N &= (F'_{1N}, \dots, F'_{mN})' \end{aligned} \quad (16)$$

So the problem is to find estimate $\hat{\theta}$ that minimizes quadratic form

$$\Phi(\theta, \Sigma) = \varepsilon_N' \Omega^{-1} \varepsilon_N, \quad \varepsilon_N = Z_N - F_N, \quad (17)$$

with $\Omega = \Sigma \otimes I_{N \times N}$, where Σ is an unknown residual covariance matrix, \otimes is the symbol of direct product, $I_{N \times N}$ is the unit matrix.

The estimate of Σ can be computed in several ways. For example, the following iteration process can be used. Denote by $\theta_N(\Sigma_1)$ the parameter estimate corresponding to minimization of $\Phi(\theta, \Sigma_1)$:

1. Take as the initial estimate $\hat{\Omega}_0 = \Sigma_0 \otimes I$ based on known measurement properties.

2. Minimize $\Phi_N(\theta, \Sigma_1)$ with $\Omega^{-1} = \hat{\Sigma}_1^{-1} \otimes I$ to determine $\hat{\theta}_N(\Sigma_1)$.

3. Evaluate estimates of $\hat{\sigma}_{jk}$, $j, k = 1, \dots, m$ elements of $\hat{\Sigma}_{i+1}$:

$$\hat{\sigma}_{jk} = \frac{\hat{\varepsilon}_j^* \hat{\varepsilon}_k}{N}, \quad \hat{\varepsilon}_j = (\hat{\varepsilon}_{jn}, \dots, \hat{\varepsilon}_{jn})',$$

$$\hat{\varepsilon}_{jn} = z_{jn} - f_j[X_n, \hat{\theta}_N(\hat{\Sigma}_1)]. \quad (18)$$

4. Compute $|\det \Sigma_{i+1} - \det \Sigma_i| \leq \xi$, if yes - go to 5., or else take $\Sigma_i = \Sigma_{i+1}$ go to 2.

5. Take $\hat{\theta}_N = \hat{\theta}_N(\hat{\Sigma}_{i+1})$.

The convergence properties of this procedure are not known.

Another method¹ uses each of the residual ε_j , $j = 1, \dots, m$ least-square estimate $\hat{\theta}^{(j)}$ approach:

1. Minimize $\Phi^{(j)}(\theta) = \varepsilon_j^* \varepsilon_j$ to determine $\hat{\theta}^{(j)}$, $j = 1, \dots, m$.

2. Compute elements $\hat{\sigma}_{jk}$ of the matrix $\hat{\Sigma}$ as in point 3. of previous method, where

$$\hat{\varepsilon}_{jn} = z_{jn} - f_j(X_n, \hat{\theta}^{(j)}), \quad j = 1, \dots, m. \quad (19)$$

3. Minimize $\Phi(\theta, \hat{\Sigma})$ to determine initial estimate $\hat{\theta}_N$.

If we assume, due to complex interlacing of independent random effects, the normality of ε_j , the least-square estimate will be consistent and the estimate $\hat{\Sigma}$ converges in probability to Σ^1 .

It is necessary to apply here minimization procedures with global minima seeking property. In⁷ such an algorithm has been developed. Its brief description will be given in section 5. After the initial estimation of $\hat{\theta}_N$ the corresponding recursive estimate correction $\hat{\theta}_{N+1}$, $\hat{\theta}_{N+2}$ should be done since X_n , Z_n , $n = N+1, N+2, \dots$ arrive successively. Consider a slight modification of stochastic approximation scheme^{4, 5, 6} for the case of vector output which is closely related to the "batch processing" scheme of reference⁵.

Using the found matrix estimate $\hat{\Sigma}$ the unit step loss function equals to

$$\zeta = \epsilon_n' \hat{\Sigma}^{-1} \epsilon_n \quad (20)$$

Its gradient is the p-vector

$$\zeta = \left[\frac{\partial \zeta}{\partial \theta} \right]' = -2 \left[\frac{\partial F}{\partial \theta} \right]' \hat{\Sigma}^{-1} \epsilon_n, \quad (21)$$

where $\frac{\partial F}{\partial \theta}$ is an $m \times p$ matrix. Introduce an $p \times m$ matrix

$$K_n = \frac{1}{\sum_{k=1}^n \left\| \left[\frac{\partial F}{\partial \theta} \right]'_{\theta_k} \right\|^2} \left[\frac{\partial F}{\partial \theta} \right]' \hat{\Sigma}^{-1}, \quad (22)$$

where $\|B\|^2 = \lambda_{\max}$ with λ_{\max} equals to max eigenvalue of $B'B$. Then the truncated recursive scheme of the form

$$\hat{\theta}_{n+1} = [\hat{\theta}_n + K_n \epsilon_n]_{\Gamma_\theta}, \quad (23)$$

with

$$[\hat{\theta}]_{\Gamma_\theta} = \begin{cases} \bar{\theta}, & \text{if } \hat{\theta} \geq \bar{\theta} \\ \hat{\theta}, & \text{if } \underline{\theta} < \hat{\theta} < \bar{\theta} \\ \underline{\theta}, & \text{if } \hat{\theta} \leq \underline{\theta} \end{cases}, \quad (24)$$

satisfies the conditions of theorem 6.4 ref⁵, and hence $(\hat{\theta}_n - \theta) \rightarrow 0$, in m.s.q. and w.p. one while $n \rightarrow \infty$.

5. The approximate solution of the optimization problem.

In this section, instead of plant output vector \bar{Y}_n , in the inequalities (3) and performance indexes w_i the model output function $F(U_n, \hat{V}_n, \hat{\theta}_n, n)$ at stage n would be used. So the original optimization problem is substituted by its deterministic approximation: find U_{in}^* to minimize

$$J_{id} = w_i[U_n, F(U_n, \hat{V}_n, \hat{\theta}_n, n)], i \in A \quad (25)$$

subject

$$\Gamma_{Ud} = \{U: \underline{G}_d \leq G[U_n, F(U_n, \hat{V}_n, \hat{\theta}_n, n)] \leq \bar{G}_d\} \quad (26)$$

Here is assumed that the D-M is able to choose vectors \underline{G}_d , \bar{G}_d to guarantee satisfaction of (3) with probability at least γ_G .

There are numerous possible (25), (26) solution methods. We shall reduce this to the seeking of unconstrained minima of the function

$$\tilde{J} = J_{id} + \chi \Pi(G, \underline{G}_d, \bar{G}_d), \quad (27)$$

where

$$\Pi = \sum_{j=1}^c E_{jd} - \bar{E}_{jd} + |E_j - E_{jd}| + |E_j - \bar{E}_{jd}| \quad (28)$$

and χ is a positive constant. Here again algorithm Ru-237 is useful. It turns to account the values of \tilde{J} only and does not use the derivatives, scanning on some successively decreasing domain in Γ_U with successively decreasing steps.

The IS computer algorithm is schematically on fig. 3. The initial estimation of $\hat{\Theta}_N$, $\hat{\Sigma}$ along the second method is performed using Ru-237. Further $\hat{\Sigma}$ is used in the recursive model parameter correction. Parameter estimates and the smoothed values \hat{V}_n of V_n are used in the model for finding the optimal decision U_1^* with w_i , Γ_{Ud} predetermined by the D-M, and for computing $Y_k = F(U_k, \hat{V}_n, \hat{\Theta}_n, n)$ with U_k set by the D-M.

6. Formaldehyd plant control system.

The plant (fig. 2). Methanol (CH_3OH) and water are mixed in the boiler (position 1) and are converted into steam (2). The steam mixed with air is preheated (3) and flows into the reactor (5). Reactor output gas, after cooling, goes into absorbing system where formaldehyd (CH_2O) and methanol are absorbed, gas consisting of CO_2 , CO , H_2 , O_2 , N_2 is thrown off. The plant is controlled by air flow rate (u_1), methanol-water correlation (u_2), contact zone temperature (u_3), and inlet air temperature (u_4). The controller outputs (u_i , $i = 1, 2, 3, 4$) are filtrated. Filtration error is assumed to be negligible. At discrete time moments, measured values of reaction gas and the product analysis, methanol and product flow rates are available.

Decision-maker. The D-M is the plant technologist who chooses $U = (u_1, \dots, u_4)^T$ in complex varying situations. This decision depends on catalyst preparing mode, on nominal (planned) values of product yield and quality, on season, conditions of absorption, in a word it depends on S_n . The main reasons for making a new decision are catalyst decaying and nominal values changing. An experienced technologist can choose for each S_n decisions that are admissible but may be far from optimality. D-M does not make a new decision if S_n

changes although decision stated previously appeared to be admissible. IS improves the D-M's decisions by complementing his conceptual model.

The control goals. The admissible decisions are such as prevent explosion and guarantee accomplishing of nominal values of product cost, quality and yield. In each state of the plant (catalyst activity, absorption coefficients) the D-M can choose several control goals as situation S_n changes:

1. Maximum average product yield;
2. Minimum average production cost;
3. Maximum average profit;
4. Minimum average methanol concentration in the product.

The plant model. Conversion of methanol into formaldehyde is a heterogeneous catalytic process in a fixed isothermal bed of silver catalyst. The main reactions are methanol oxydation and dehydrogenation together with side and consecutive reactions. Reaction rates are diffusion controlled, hence rate coefficients depend on gas flow rate. Model equations arose from heat and mass balance and chemical kinetics. The conversion process is an autothermal one: reaction zone temperature depends on the concentration of reagents, steam temperature and rate coefficients which in turn depend on zone temperature. Therefore the model output can be computed iteratively. The behavior of model and plant is similar (existence of stationary points, directions of gradients of Y). By identification procedure the minimum weighed mean-square model error can be obtained.

Identification. Parameters to be identified are reaction rate coefficients, heat losses, formaldehyde and methanol absorption coefficients.

Optimization problems. The admissible set Γ_U is constructed on bases of following set of inequalities:

1. $\underline{g}_U \leq U_n \leq \bar{g}_U$, \underline{g}_U , \bar{g}_U - 4-vectors of technological restrictions;

2. $g_2 \leq g_2[U_n, F(U_n, \hat{v}_n, \hat{\theta}_n, n)]$ - safety conditions;

3. $g_3[U_n, F(U_n, \hat{v}_n, \hat{\theta}_n, n)] \leq \bar{g}_3$ - normal absorption conditions
 4. $g_4 \leq g_4[U_n, F(U_n, \hat{v}_n, \hat{\theta}_n, n)] \leq \bar{g}_4 - \text{CH}_2\text{O}$
 5. $g_5 \leq g_5[U_n, F(U_n, \hat{v}_n, \hat{\theta}_n, n)] \leq \bar{g}_5 - \text{CH}_3\text{OH}$ } concentrations in the product: quality conditions;

6. $g_6 \leq g_6[U_n, F(U_n, \hat{v}_n, \hat{\theta}_n, n)]$ - nominal yield requirement;

7. $g_7[U_n, F(U_n, \hat{v}_n, \hat{\theta}_n, n)] \leq g_7$ - nominal production cost requirement, where g_2 is the CH_3OH concentration before the reactor, g_3 is gas volume flow rate, g_4 and g_5 are CH_2O and CH_3OH concentrations in the product, g_6 is yield, g_7 is production cost.

The goal number	Performance index	Γ_{Ud}
1	$-g_6[U_n, F(U_n, \hat{v}_n, \hat{\theta}_n, n)] = \min_{\Gamma_{Ud}}$	conditions 1.-5., 7.
2	$g_7[U_n, F(U_n, \hat{v}_n, \hat{\theta}_n, n)] = \min_{\Gamma_{Ud}}$	conditions 1.-6.
3	$-g_8 = -g_6(g_7 - \lambda) = \min_{\Gamma_{Ud}}$	conditions 1.-7.
4	$g_5[U_n, F(U_n, \hat{v}_n, \hat{\theta}_n, n)] = \min_{\Gamma_{Ud}}$	conditions 1.-4., 6.-7.

where g_8 is profit, λ is product price.

Control system realization. Data accumulation rate is low - 16 averaged quantities per shift (8 hours), therefore the IS computation may be performed in an off-line mode.

The present IS algorithm enables to control several independent formaldehyd plants. Plants are connected with the computing center by telegraph lines. In a plant data are booked into special forms, punched and transmitted to the computer. They contain three different data sets (B1, B2, B3) and plant number label together with computation code. The latter shows what kind of information the D-M wants to acquire from computer. B1 consists of X_n, Z_n for the identification purposes, B2 presents constraint vectors and constructive parameters, B3 consist of vector U_k to be tested on the model. Information from B2 is stored and used until a new arrives. The results are punched by computer in the following forms (Fig. 4-6):

1. The means of the plant variables, Fig. 4a;
2. Estimated parameters, Fig. 4b;
3. Model outputs \hat{Y}_k at the state $\hat{\theta}_n$ and \hat{V}_n for decision U_k , Fig. 5;
4. Optimal decisions corresponding to goals 1.-4. and the plant output for those decisions, Fig. 6.

The algorithm is realized on a Minsk-22 computer. The program consists of 6700₍₁₀₎ words in core memory and magnetic tape of $8192 + \langle \text{number of plants} \rangle \cdot 1024$ words.

R E F E R E N C E

1. Beauchamp, J.J., Cornell, R.G. Simultaneous nonlinear estimation. *Technometrics*, V. 8, nr. 2, May, 1966.
2. Андерсон, Т. Введение в многомерный статистический анализ. М. 1963.
3. Арис, Р. Анализ процессов в химических реакторах. Л. 1967.
4. Цыпкин, Я.Э. Адаптация, обучение и самообучение в автоматических системах. *Автоматика и телемеханика*, т. XXVII, № I, 1966.
5. Albert, A.E., Gardner, L.A. Stochastic approximation and nonlinear regression. MIT Press, Cambridge, Mass., 1967.
6. Zhivoglyadov, V.P., Kaipov, V.Kh. Identification of distributed parameter plants in the presence of noises. Prepr. of the IFAC symp. Identification on automatic control systems, Prague, 1967.
7. Руубель, Х.В. Поиск точки минимума функции В со. программы для ЦВМ "Минск-22", вып. 7 (to be appear).

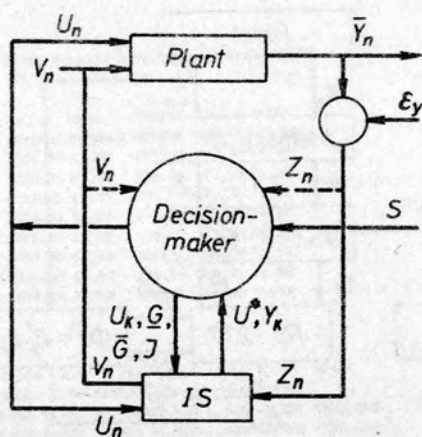


Fig. 1.

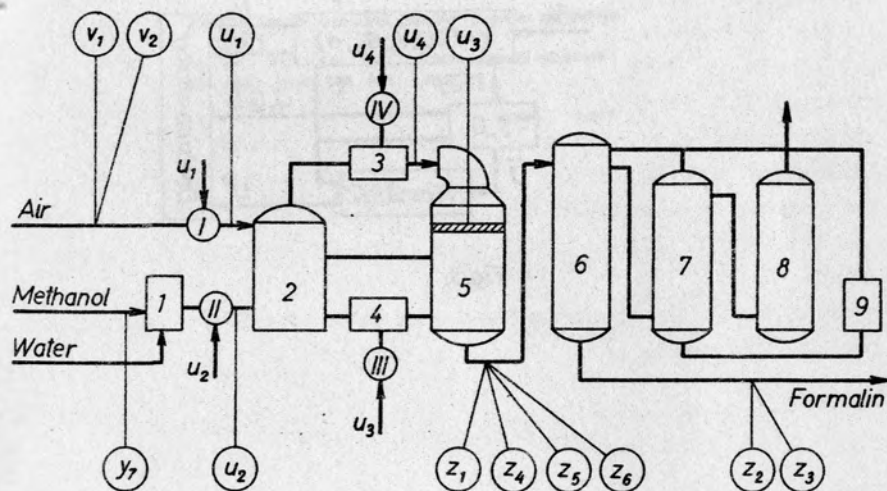


Fig. 2.

СВОДКА РЕЖИМОВ И ВЫХОДОВ АГРЕГАТА КОМПЕР

2

ЧИСЛО КОМПЕР	РАСХОД	КОМ	ТЕМПЕР	ТЕМПЕР	УСЛОВ
МЕСЯЦ СМЕН	ВОЗДУХА	МЕТАНОЛ	КОНТАК	ПЕРЕГР	МЕТАНОЛ
4,09	82,00	2050,00	63,40	695,00	112,00 17,42
5,09	85,00	2050,00	63,70	695,00	113,00 17,30
6,09	88,00	2050,00	63,30	695,00	112,00 17,66
7,09	91,00	2050,00	63,60	694,00	113,00 17,32
8,09	94,00	2050,00	63,50	695,00	111,00 16,59
9,09	97,00	2050,00	63,50	694,00	110,00 17,04
10,09	100,00	2050,00	63,50	695,00	110,00 16,00
11,09	103,00	2050,00	63,90	694,00	110,00 14,92

ВЫХОД	СРЕДН	ОСНОВ	РАСХ	УСЛОВ	ПРОЦ	ПРОЦ	ЭКОН	ДОХОД
КОМ	ТРИН	КОМ	КОД	КОМ-НА	КОМ-ДА	МЕТАН	МЕТАН	
0,756	0,883	0,856	0,523	30,307	37,400	6,600	-0,091	2313,018
0,755	0,879	0,860	0,523	30,293	37,400	6,400	-0,105	2303,836
0,755	0,885	0,853	0,525	30,674	37,500	6,700	-0,159	2329,996
0,748	0,883	0,846	0,529	29,864	37,400	7,100	-0,269	2431,801
0,744	0,874	0,852	0,531	28,477	37,400	6,800	-0,323	2331,244
0,747	0,873	0,855	0,529	29,384	37,400	7,300	-0,266	2412,255
0,745	0,870	0,856	0,526	27,888	37,100	7,200	-0,173	2297,420
0,740	0,864	0,856	0,533	25,554	37,300	7,100	-0,326	2088,409

СОСТОЯНИЕ ПРОЦЕССА

ПРЕДЪЯСНОВЕНИЯ СКОРОСТЕЙ ПОТЕРИ СТЕПЕНИ ПОГЛОЩЕНИЯ

КОД	К18	К20	К30	К40	ТЕПЛА	КОМ-ДА	МЕТАНОЛ
3037	4309	2476	185	593	55000		

1,000 0,917

Figure 4

ИСПЫТАНИЕ ЗАДАНИЙ РЕШИМО

1, 2050,0	63,0	695,0	113,0					
2, 2050,0	63,5	695,0	114,0					
3, 2100,0	64,0	695,0	110,0					
	1	2	3	0	0	0	0	0
УСЛОВ								
	31,0	31,0	32,7	0,0	0,0	0,0	0,0	0,0
РАСХ КОС								
	0,512	0,517	0,609	0,000	0,000	0,000	0,000	0,000
ДОХОД								
	2597	2582	2422	0	0	0	0	0
СОБЕРИ МЕТАЛ								
	6,4	6,8	16,6	0,0	0,0	0,0	0,0	0,0
СОБЕРИ ОРОН								
	49,2	49,1	40,6	0,0	0,0	0,0	0,0	0,0
ЭКОНОМИЯ								
	0,24	0,09	-2,90	0,00	0,00	0,00	0,00	0,00
ВЫХОД ОРОН								
	0,767	0,760	0,645	0,000	0,000	0,000	0,000	0,000
СЕЛЕКТИВ								
	0,007	0,007	0,914	0,000	0,000	0,000	0,000	0,000
ГАЗО О2								
	0,20	0,20	0,24	0,00	0,00	0,00	0,00	0,00
ВУН СО2								
	3,92	3,92	4,02	0,00	0,00	0,00	0,00	0,00
АНА СО								
	1,60	1,60	0,30	0,00	0,00	0,00	0,00	0,00
ВУН Н2								
	10,0	10,0	17,1	0,0	0,0	0,0	0,0	0,0
КОМЛ О2								
	0,096	0,096	0,007	0,000	0,000	0,000	0,000	0,000
КОМЛ СН3ОН								
	0,266	0,269	0,294	0,000	0,000	0,000	0,000	0,000

Figure 5

РЕКОМЕНДУЕМЫЕ РЕЖИМЫ

МАКС ВЫРАБОТКА	2050,0	61,0	692,0	129,0
2 МИН РАСХ КОЗ**	2050,0	63,5	712,0	125,0
3 МАКС ДОХОД	2050,0	61,5	696,0	120,0
4 МИН СОДЕРЖ МЕТАН	2050,0	64,0	718,0	115,0

РЕЗУЛЬТАТ ОПТИМАЛЬНЫХ РЕЖИМОВ

	1	2	3	4	0	0	0	0
УСЛОВ	31,1	30,4	31,0	30,1	0,0	0,0	0,0	0,0
РАСХ КОЗ**	0,519	0,506	0,508	0,491	0,000	0,000	0,000	0,000
ДОХОД	2587	2568	2612	2587	0	0	0	0
СОДЕРЖ МЕТАН	7,2	5,0	6,0	3,2	0,0	0,0	0,0	0,0
СОДЕРЖ ФОРН	46,6	51,2	48,4	53,9	0,0	0,0	0,0	0,0
ЭКОНОМИЯ	0,02	0,43	0,39	0,88	0,00	0,00	0,00	0,00
ВЫХОД ФОРН	0,757	0,777	0,774	0,801	0,000	0,000	0,000	0,000
СЕЛЕКТИВН	0,891	0,870	0,887	0,861	0,000	0,000	0,000	0,000
ГАЗО O2	0,21	0,19	0,20	0,19	0,00	0,00	0,00	0,00
ВЫР CO2	3,93	3,86	3,92	3,82	0,00	0,00	0,00	0,00
АНА CO	1,41	2,37	1,60	2,80	0,00	0,00	0,00	0,00
АНЗ H2	17,8	18,5	18,0	18,8	0,0	0,0	0,0	0,0
КОНЦ O2	0,093	0,098	0,095	0,101	0,000	0,000	0,000	0,000
КОНЦ СН3ОН	0,262	0,264	0,260	0,260	0,000	0,000	0,000	0,000

Figure 6

NOT APPROVED BY THE BOARD OF DIRECTORS