

STAT

Declassified in Part - Sanitized Copy Approved for Release  
2013/02/13 : CIA-RDP80T00246A023400480001-9

**Page Denied**

Declassified in Part - Sanitized Copy Approved for Release  
2013/02/13 : CIA-RDP80T00246A023400480001-9

# Two-positional Functional Frequency Device for Automatic Regulation

I. A. MASLAROFF

*Bulgaria*

## Introduction

The complicated character of the technological processes has developed in parallel with other research methods of ascertaining ways of improving the qualities of the two-positional method for regulation. The simplicity of the device and the low price of the required elements have not detracted from its significance. From all published literature on this subject the extensive work of Campe Nemm<sup>1</sup> is particularly noted. The author analyses the existing methods of reducing the fluctuations of the unit to be regulated: increasing the extent of current: the use of cut-off two-positional regulation: and the introduction of inverse connections on the first and second derivative, etc.

This paper gives some results of the methods undertaken to improve the two-positional regulation by changing the frequency of the influenced impulses. The methods are mainly directed towards decreasing the fluctuations of the unit to be regulated.

## The Essence of Two-positional Functional Frequency Regulation

The present survey refers to the monotonous varying processes of a unit with a comparatively small changing rate of regulation and the form of the equation to be used:

$$C \frac{dA}{dt} = \sum Q \quad (1)$$

The principle of two-positional functional frequency regulation consists in the addition to the object of previously fixed identical portions of the utilized unit in the form of impulses. The frequency of these impulses depends on the difference  $\Delta A$  between the given and actual value of the unit to be regulated. Initially the influence of the net delay in the system is neglected in the survey.

Figure 1 shows the change of the unit to be regulated. During the time of impulses it is determined by:  $A = A_y (1 - e^{-t/T})$  and during the pauses, by:  $A = A_k e^{-t/T}$  ( $t = 0, A = A_k$ ). These two expressions are the integrals of (1) in the presence and absence of current. In such cases, at the end of the impulses and pauses, the unit to be regulated will be determined by:

$$\begin{aligned} A_1 &= A_y (1 - e^{-t_1/T}) \\ A_2 &= A_1 e^{-t_1/T} = A_y (1 - e^{-t_1/T}) e^{-t_1/T} \\ A_3 &= A_y (1 - e^{-t_1/T}) + A_2 e^{-t_1/T} = A_y (1 - e^{-t_1/T}) e^{-t_1/T} \\ A_4 &= \dots \\ &\vdots \end{aligned} \quad (2)$$

By using the method of full mathematical induction, we determine that the value of the unit to be regulated after  $n$  consecutive cycles (impulses and pauses) will be equal to:

$$A_{2n} = A_y (1 - e^{-t_1/T}) \sum_{a=1}^n e^{-\sum_{k=a}^n [t_k + (n-a)t_1]/T} \quad (3)$$

and after  $n + 1$  serial impulses:

$$A_{2n+1} = A_y (1 - e^{-t_1/T}) \left[ 1 + \sum_{a=1}^n e^{-\sum_{k=a}^n [t_k + [n-(a-1)t_1]/T} \right] \quad (4)$$

Eqns (3) and (4) show that by changing the duration of pauses one can effectively influence the unit to be regulated. In order to obtain the regulation we need the functional relation  $t = \varphi(\Delta A)$ , at which the time of the pause will increase with the decrease of the magnitude of the difference  $\Delta A$ . Such a dependence may be realized simply by introducing the exponential block in the scheme of the regulator (Figure 2).

The equation, characterizing the work of this scheme is:

$$k\Delta A (1 - e^{-t/T_1}) = B$$

The time constant of the exponential block of the scheme must be much smaller than the time constant of the object.

Then at  $\Delta A = \text{const}$ . the time of the pause is equal to:

$$t = T_1 \ln \frac{k\Delta A}{k\Delta A - B} \quad (5)$$

Eqn (5) shows large values of the difference when the percentage change in the pause time is insignificant. At an established regime when there are small values of the difference between the given and actual values of the unit to be regulated, the time of the pause is determined only by the parameters of the object ( $T \gg T_1$ ) where the delay due to the regulator is slightly neglected in comparison with the common time of the pause. In such a case the time of the pause is determined taking into consideration that the consecutive fluctuations of the unit to be regulated at a determined regime are also equal:

$$\delta A' = \delta A'' \quad (6)$$

where

$$\delta A' = A_{2n+1} - A_{2n+2}; \quad \delta A'' = A_{2n+3} - A_{2n+2}$$

Since

$$A_{2n+3} = A_y (1 - e^{-t_1/T}) + A_{2n+2} e^{-t_1/T}$$

$$A_{2n+2} = A_{2n+1} e^{-t_1/T}$$

031/2

the time of the pauses is equal to:

$$t_{n+1} = T \ln \frac{A_{2n+1}}{A_{2n+1} - A_y(1 - e^{-t_i/T})} \quad (7)$$

By exerting an influence on the coefficient of amplification and the internal limit of putting in motion  $B$  of the scheme it is always possible to receive an equalization for the maximal and given values for the unit to be regulated. Then eqn (7) is modified as:

$$t_{n+1} = T \ln \frac{A_g}{A_g - A_y(1 - e^{-t_i/T})} \quad (7a)$$

The maximum value of the fluctuations of the unit to be regulated is given by:

$$\delta A = \Delta A = \frac{B}{k} = A_g - A_{2n+2} = (A_y - A_g)(1 - e^{-t_i/T})e^{-t_i/T} \quad (8)$$

Eqn (8) shows that by decreasing the duration of the impulse  $t_i$  the fluctuations of the unit to be regulated may be most effectively reduced. The coefficient of amplification  $k$  may be determined at a previously chosen value  $B$  of the limit out of the duration of the impulse.

#### Influence of the Net Delay on the Two-positional Functional Frequency Method for Regulation

Usually, the effect of the delay which increases fluctuations of the unit to be regulated is shown in the systems of the type examined. In the following it is proved that the influence of the net delay upon the value of fluctuations may be substantially decreased using the functional frequency method for regulation. Actually *Figure 3* shows that the additional increase of fluctuations  $\delta A_{dt}$  which follows from the delay of the system, is equal to:

$$\delta A_{dt} = A_{2n+2}(1 - e^{-\Delta t/T}) \cong A_g(1 - e^{-\Delta t/T}) \quad (9)$$

With the usual two-positional regulation, the delay increases the fluctuations of the unit to be regulated in the direction of its decrease, as well as in the direction of its increase. These additional increases are of the same order.

It follows that with functional two-positional regulation the fluctuation of the unit to be regulated increases in the direction of its decrease and because of this the received additional fluctuation is about twice lower.

The total value of fluctuations is:

$$\delta A_x = \delta A + \delta A_{dt} = (A_y - A_g)(1 - e^{-t_i/T})e^{-t_i/T} + A_g(1 - e^{-\Delta t/T}) \quad (10)$$

If it is accepted that  $\delta A = \delta A_{dt}$ , then:

$$\frac{A_y}{A_g} = 1 + \frac{(1 - e^{-t_i/T})}{(1 - e^{-\Delta t/T})} e^{-t_i/T} \quad (11)$$

From eqn (11) some conclusions can be drawn for determining the parameters of the system to be regulated.

It is evident that at considerable values of the time of delay  $\Delta t$  it is apt to accept  $\Delta_y \gg \Delta_g$ , i.e. to use strong impulses. However, at small values of  $\Delta t$  it is apt to accept  $A_g \cong A_y$ , i.e. the impulses will be comparatively weaker.

From eqn (8) two fundamental parameters for the regulation may be determined—the internal limit for setting in motion  $B$  and the coefficient of the earlier amplification  $k$ . These parameters may be easily changed into parameters to be regulated in large limits, depending on the requirements of the object to be regulated.

#### Constructive Data of the Device for Functional Frequency Regulation

The device uses a vacuum-tube scheme (*Figure 4*) consisting of a measuring part 1, amplifier 2 and an integral group 3, two channels for constant current amplifiers 4 and 4' and an executive trigger 5. It differs from *Figure 2* by the use of a second channel for the constant current amplifier 4', which is included in a circulating chain of the integrating group and the base constant current amplifier 4. Its purpose is to accelerate the process for establishing the regime. When there are many large values of  $\Delta A$  the output voltage of 4' passes through the logical scheme 'IF'-6 and sets in motion the executive trigger. In this way the scheme works as an ordinary two-positional regulator. Placed in a regime, close to the one established, the output voltage of the second channel is not in position to set in motion the executive trigger, and the device works like a functional frequency regulator.

In parallel with the passing of each impulse from the trigger exit 5 to the object 7 the signal for clearing the integrating chain is simultaneously passed through an internal link.

#### Experimental Data

Initially the device was constructed and tested for regulating the concentration of solutions. Conductive transformers linked by a bridge scheme with temperature compensation were used as a measuring device\*.

The executive trigger exerts influence on an electromagnetic valve which adds a drop of concentrate to the solution at each impulse. The results obtained at the time of regulation were very good.

The device is used to regulate temperature, and for this purpose the executive trigger is replaced by a delay multivibrator. The time of the impulse may be regulated at will by changing the parameters of its device. *Figure 5* shows the diagrams of temperature change of one and the same object, recorded with the help of an electronic potentiometer. It is seen that the quality of regulation with the functional frequency method is much better than that of the ordinary two-positional method.

#### Conclusions

1. The two-positional functional frequency device for regulation allows the possibility of decreasing the fluctuations of the unit to be regulated, particularly those emerged out of the delay in the system.
2. By the character of its work, the device approaches the statistical regulators.
3. The devices for regulation can be realized by using practical simple means.
4. The test results prove the expedience of using this method for regulation in many cases.

\* Eng. D. Detcheva took part in the computing of the construction of the device.

**Nomenclature**

- $C$  Coefficient of the generalized capacity of the object to be regulated
- $A$  The unit to be regulated
- $A_y$  Fixed value of the unit to be regulated
- $A_g$  Given value of the unit to be regulated
- $\Delta A$  Difference between the given and actual value of the unit to be regulated
- $Q$  Generalized quantitative index of the process
- $\delta A$  Variation of the unit to be regulated in the period of one impulse or pause
- $t$  Time
- $t_i$  Time of the impulse

- $\Delta t$  Time of the net delay
- $n$  Number of the impulses
- $T$  Time constant of the object to be regulated
- $T_1$  Time constant of the exponential block of the scheme
- $B$  Internal limit for setting in motion the acting block of the scheme
- $k$  Coefficient of amplification

**References**

<sup>1</sup> CAMPE NEMM, A. A. Two-positional automatic regulation and methods of improving its characteristics. *Thermoenergetic and Chemicotechnological Devices and Regulators*. 1961. Moscow-Leningrad; Mashgiz

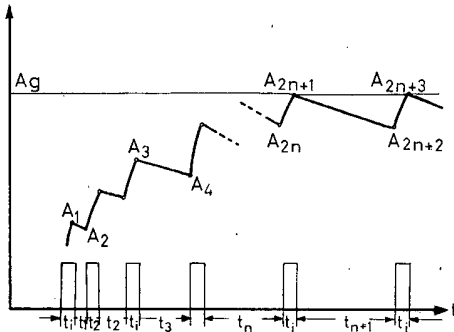


Figure 1

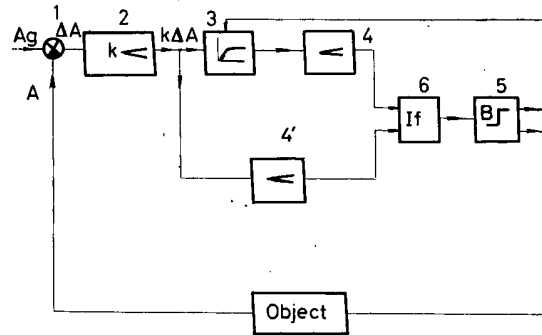


Figure 4

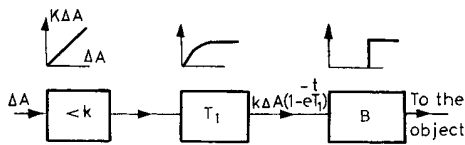


Figure 2

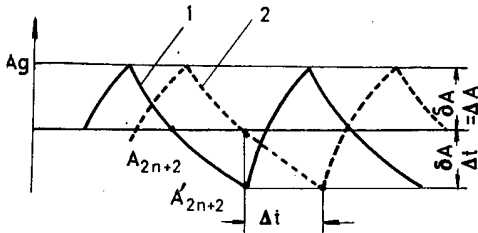


Figure 3

- Curve 1—Change of the regulated unit in close proximity to the source of the impulses
- Curve 2—Change of the regulated unit in the field of the sensitive element

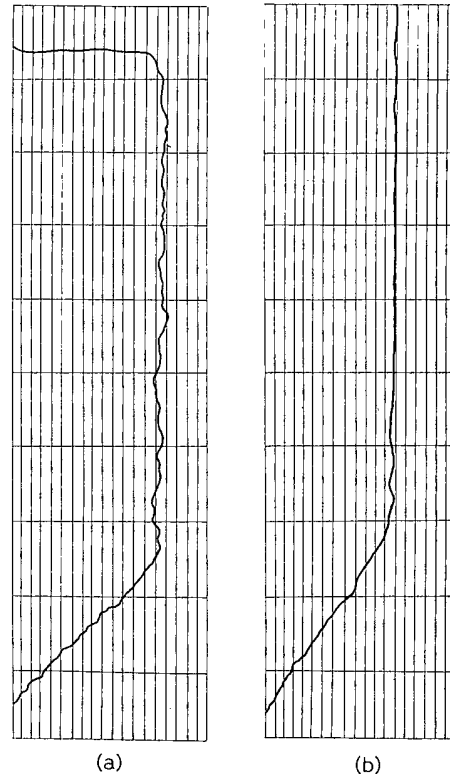


Figure 5

- (a) Change of temperature by using a contact thermometer for regulation
- (b) Change of temperature by using functional frequency regulation of the object

# On the Estimation of the Decaying Time

H. LING

CPR

The problem of stability is one of the basic problems in the proper operation of any dynamic system. After Liapunov's brilliant work<sup>1</sup>, very great effort has been expended on its theory in the last twenty years<sup>2-8</sup>. In the monograph by Zubov<sup>4</sup>, the stability of the invariant sets of the abstract dynamic systems in the metric space are treated in a very general sense.

Consider the differential system

$$\dot{x}_s = X_s(x_1, \dots, x_n, t) \quad (s=1, 2, \dots, n) \quad (1)$$

It is often necessary to study the stability problem for the given values  $\phi_1^0, \dots, \phi_k^0$  of the coordinate functions

$$\phi_i = \phi_i(x_1, \dots, x_n) \quad (i=1, 2, \dots, K) \quad (2)$$

Without loss of generality,  $\phi_1^0 = 0, \dots, \phi_k^0 = 0$  may be taken. Thus, the standard working state of the system is given by the equation

$$\phi_i(x_1, \dots, x_n) = 0 \quad (i=1, 2, \dots, K) \quad (3)$$

Let the set of points defined by (3) constitute an  $(n-k)$  manifold  $\mathcal{F}(n-k)$ . By means of (3), some particular and interesting motions of system (1) can generally be described (for example, the self-excited oscillations or the motions which demand their characteristic functions to take on given values). Here the generalized stability differs from the Liapunov stability in that the unperturbed motion is no longer a particular motion (e.g. trivial solution) but its  $K$  coordinate functions take on given values. In general, (3) represents a class of motions and constitutes a manifold in the phase-space. From the research point of view, the coordinate functions (2) are of more interest than the coordinates  $x_1, \dots, x_n$  themselves.

Obviously, when  $\phi_i = x_i, i=1, \dots, k, k < n$ , the stability of the partial coordinates is obtained, and when  $\phi_i = x_i, i=1, \dots, k, k = n$ , it agrees with the stability in Liapunov's sense.

In papers by Liapunov<sup>1</sup> and Rumyantsev<sup>8</sup> the stability of the partial coordinates and the stability for the given values of the functions are discussed. They require the absolute values of the initial perturbations  $x_1^0, \dots, x_n^0$  to be sufficiently small and thus essentially the unperturbed motion was supposed to be a point in the phase space. The approach in this paper differs from theirs, and it will be explained clearly below.

The set of points which satisfy the inequality

$$\sum_{i=1}^K \phi_i^2(x_1, \dots, x_n) \leq H^2 \quad (4)$$

is called the  $H$ -neighbourhood of  $\mathcal{F}(n-k)$  and is written as  $\mathcal{F}(n-k)(H)$ . It is assumed tacitly that through each point of  $\mathcal{F}(n-k)(H)$  there exists a unique solution of the system (1).

Of course, it is necessary that the standard working state (3) of the system has a certain upholding ability, which means that the functions

$$\Phi_i(x_1, \dots, x_n, t) = \text{grad } \phi_i X \quad (i=1, \dots, K) \quad (5)$$

should satisfy the conditions

$$\Phi_i(x_1, \dots, x_n, t) \equiv 0 \text{ as } \{x\}_n \in \mathcal{F}(n-k) \quad (6)$$

or equivalently,  $\mathcal{F}(n-k)$  is an invariant set of the system (1), where  $\{x\}_n$  represents the  $n$ -dimensional vector.

*Definition 1.* The system (1) is said to be stable with respect to the functions (2), taking on zeros (3) whenever, given any  $\varepsilon > 0$ , there is a  $\delta(t_0, \varepsilon) > 0$ , such that, for all trajectories  $x(t)$  with initial values satisfying

$$\{x(t_0)\}_n = \{x^0\}_n \in \mathcal{F}(n-k)(\delta), \quad t_0 > 0 \quad (7)$$

one has

$$\{x(t)\}_n \in \mathcal{F}(n-k)(\varepsilon) \quad (8)$$

for all  $t \geq t_0$ .

*Definition 2.* The system (1) is said to be asymptotically stable with respect to the functions (2) taking on zeros (3) if

(a) Definition (1) holds,

$$(b) \quad \lim_{t \rightarrow \infty} \Phi^2 = \lim_{t \rightarrow \infty} \sum_{i=1}^K \phi_i^2[x, (t), \dots, x_n(t)] = 0 \quad (9)$$

i.e. for any given  $\eta > 0$  there is positive number  $T = T(\eta, t_0, x^0)$  such that

$$\{x(t)\}_n \in \mathcal{F}(n-k)(\eta) \text{ as } t \geq t_0 + T \quad (10)$$

*Definition 3.* The system (1) is said to be equi-asymptotically stable with respect to the functions (2), taking on zeros (3) if

(a) Definition 2 holds,

(b) there exists  $H > 0$ , and  $T = T(\eta)$ , such that for all trajectories  $\{x(t)\}_n$  with initial values satisfying

$$\{x(t_0)\}_n \in \mathcal{F}(n-k)(H) \quad t_0 \geq 0 \quad (11)$$

then (10) holds.

If the initial conditions were subjected to the following restrictions

$$\{x(t_0)\}_n \in G^0 \quad t_0 > 0 \quad (12)$$

then the above-mentioned stability, asymptotic stability and equi-asymptotic stability are said to be stable, asymptotically stable and equi-asymptotically stable under condition (12) respectively.

In the sequel, the function  $V(x_1, \dots, t)$  is called the Liapunov function with respect to functions  $\phi_1, \dots, \phi_k$  if

$$V(x_1, \dots, x_n, t) \equiv 0 \text{ as } \{x\}_n \in \mathcal{F}(n-k) \quad (13)$$

and  $V$  is assumed to have continuous partial derivations.

103/2

**Definition 4.** The function  $V(x_1, \dots, x_n, t)$  is said to be positive (negative) semi-definite with respect to (2) if

- (a) (13) holds,  
 (b)  $V \geq 0$  [ $V \leq 0$ ] in  $\mathcal{F}(n-k)(H)$ .

**Definition 5.** The function  $V(x_1, \dots, x_n, t)$  is said to be positive (negative) definite with respect to (2) if

- (a) Definition 4 holds,  
 (b) there is a positive function  $W_1(y_1, \dots, y_k)$  such that, in  $\mathcal{F}(n-k)(H)$ ,

$$V(x_1, \dots, x_n, t) \geq W_1[\phi_1(x_1, \dots, x_n), \dots, \phi_k(x_1, \dots, x_n)] \quad (14)$$

$$\{V(x_1, \dots, x_n, t) \leq -W_1[\phi_1(x_1, \dots, x_n), \dots, \phi_k(x_1, \dots, x_n)]\}$$

**Definition 6.** The function  $V$  is said to be uniformly small, if for any given  $\varepsilon > 0$ , there is  $\delta(t) > 0$  such that the conditions  $t > 0$  and  $\{x\}_n \in \mathcal{F}(n-k)(\delta)$  imply  $V \leq \varepsilon$ .

**Definition 7.** The function  $V(x_1, \dots, x_n, t)$  is said to have infinitely small upper bound with respect to (2) if there is a continuous function  $W_2(y_1, \dots, y_k)$  such that

- (a)  $W_2(0_1, \dots, 0) = 0$ ,  
 (b) in  $\mathcal{F}(n-k)(H)$ ,

$$W_2[\phi_1(x_1, \dots, x_n), \dots, \phi_k(x_1, \dots, x_n)] \geq V(x_1, \dots, x_n, t) \quad (15)$$

**Definition 8.** The function  $V(x_1, \dots, x_n, t)$  is said to have the property  $A$ , if there are two positive continuous functions  $W_1(S)$  and  $W_2(S)$  such that

$$(a) \quad W_1(0) = W_2(0) = 0, \quad W_1(\infty) = W_2(\infty) = +\infty \quad (16)$$

$$(b) \quad W_2(\|\phi\|_K) \geq V(x_1, \dots, x_n, t) \geq W_1(\|\phi\|_K) \quad (17)$$

Parallel to Definitions 1 and 3, one has the fundamental theorems shown in the following section.

### The Fundamental Theorems

(A) For the system (1), if there is a Liapunov function  $V(x_1, \dots, x_n, t)$  such that

- (a)  $V$  satisfies Definition 5 and it is positive definite with respect to (2),  
 (b)  $V$  satisfies Definition 7,  
 (c) the total derivative

$$\left. \frac{dV}{dt} \right|_{1,1} = \frac{\partial V}{\partial t} + \text{grad } V \cdot X \quad (18)$$

is negative definite with respect to (2), then the system (1) satisfies Definition 3.

(B) If the system (1) satisfies Definition 3, and the rank of matrix

$$\frac{D(\phi_1, \dots, \phi_K)}{D(x_1, \dots, x_n)} = \begin{pmatrix} \frac{\partial \phi_1}{\partial x_1}, \dots, \frac{\partial \phi_1}{\partial x_n} \\ \dots \\ \frac{\partial \phi_K}{\partial x_1}, \dots, \frac{\partial \phi_K}{\partial x_n} \end{pmatrix} \quad (19)$$

is  $K$ , and the functions  $\Phi_i$  ( $i = 1, \dots, K$ ) defined by (5) are uniformly bounded in  $\mathcal{F}(n-k)(H)$ , then there is a function  $V$

which satisfies all conditions in (A). The proof of this theorem is given in the Appendix. It is not difficult to prove the following corollaries.

**Corollary 1.** If  $V$  satisfies Definitions 5 and 6, and  $dV/dt|_{(1)}$  is negative semi-definite with respect to (2), then the system (1) satisfies Definition 1.

**Corollary 2.** If  $V$  satisfies Definition 8, and  $dV/dt|_{(1)}$  is negative definite with respect to (2) then (9) holds for any  $t_0 > 0$  and any  $x^0$  in the space.

Let the set of position points at time of the motions which take on the initial positions in  $G^0$  be written as  $G^{(t)}$ .

**Corollary 3.** If  $V$  satisfies (A), Corollary 1 or Corollary 2 in  $G^{(t)} \cap \mathcal{F}(n-k)(H)$ , then the system (1) is stable, equi-asymptotically stable, or asymptotically stable in the whole under condition (12) respectively.

**Corollary 4.** If the system (1) satisfies Definition 3 under condition (12) and the rank of matrix (19) is  $K$  in the neighbourhood  $\mathcal{F}(n-k) \cap G^{(t)}$  and  $\Phi_i$  ( $i = 1, \dots, K$ ) are uniformly bounded in  $G^{(t)} \cap \mathcal{F}(n-k)(H)$ , then there is a function  $V$  which satisfies the conditions in Corollary 3.

**Example—**Consider the system

$$\begin{cases} \dot{x} = ay - cx(bx^2 + ay^2) \sin \frac{1}{bx^2 + ay^2} \\ \dot{y} = -bx - cy(bx^2 + ay^2) \sin \frac{1}{bx^2 + ay^2} \end{cases} \quad (20)$$

$$(a \cdot b \cdot c > 0)$$

Obviously, if one takes  $\phi = bx^2 + ay^2$  then  $\phi = 1/k\pi$  ( $k = 1, 2, \dots$ ) are the invariant sets of (20), they are closed orbits. By means of the Liapunov functions  $V = \frac{1}{2}(\phi - 1/k\pi)_2$  with respect to  $\phi = 1/k\pi$ , the following statements can be proved:

- (a) In the exterior of the ellipse  $\phi = 1/\pi$ , there is no closed orbit;  
 (b) in the interior of the ellipse  $\phi = 1/\pi$ , there are infinitely many closed orbits;  
 (c) the closed orbit is asymptotically stable when  $K$  is even and it is unstable when  $K$  is odd;  
 (d) the origin  $x = y = 0$  is a singular point of (20) and it is stable. In any of its neighbourhood, there are infinitely many closed orbits, and hence the origin is not asymptotically stable.

In the regulating or the dynamic systems, it is often necessary to estimate the decaying time of perturbations for the standard working state. In this paper the problem of estimating the decaying time is considered. In the sequel it is assumed that system (1) is equi-asymptotically stable with respect to (2), and the following discussions are valid in certain attractive regions of  $\mathcal{F}(n-k)$ .

Let  $V$  be a Liapunov function of (1) which satisfies the conditions of the fundamental theorem (A). In the general case, there are two positive definite functions  $W_1(y_1, \dots, y_k)$  and  $W_2(y_1, \dots, y_k)$  such that

$$\begin{aligned} W_2[\phi_1(x), \dots, \phi_K(x)] &\geq V(x_1, \dots, x_n, t) \\ &\geq W_1[\phi_1(x), \dots, \phi_K(x)] \end{aligned} \quad (21)$$

Besides, it is assumed that there are two functions  $f_1(s)$  and  $f_2(s)$ , such that in  $\mathcal{F}(n-k)(H)$  the inequalities

$$f_1(v) \leq \frac{dv}{dt} \leq f_2(v) \quad (22)$$

hold. Furthermore, from (21) one has, in general,

$$\begin{cases} \{x_n\} \in \{W_2 \leq V_0\} \text{ implies } \{x\}_n \in \{V \leq V_0\} \\ \{x_n\} \in \{V \leq \varepsilon\} \text{ implies } \{x\}_n \in \{W_1 \leq \varepsilon\} \end{cases} \quad (23)$$

where  $V_0$  and  $\varepsilon$  are given position numbers. Denote

$$T_1 = - \int_{\varepsilon}^{V_0} \frac{d\lambda}{f_1(\lambda)}, \quad T_2 = - \int_{\varepsilon}^{V_0} \frac{d\lambda}{f_2(\lambda)} \quad (24)$$

Then, the following theorem estimates the decaying time.

**Theorem 1.** The decaying time  $T$  of the motion of the system (1) from an initial point in the region

$$W_2[\phi_1(x), \dots, \phi_k(x)] \leq V_0 \quad (25)$$

to a point in the region

$$W_1[\phi_1(x), \dots, \phi_k(x)] \leq \varepsilon \quad (26)$$

satisfies the inequality

$$T \leq T_2 \quad (27)$$

The decaying time  $T$  from an initial point in the region

$$W_2[\phi_1(x), \dots, \phi_k(x)] \geq V_0 \quad (28)$$

to a point in the region (26) satisfies

$$T \geq T_1 \quad (29)$$

Let  $M_2(R)$  be the maximum value of  $W_2$  on the boundary

$$\|\phi\|_k = R \text{ of } \mathcal{F}(n-k)(R)$$

and let  $m_1(\gamma)$  be the minimum value on the boundary  $\|\phi\|_k = \gamma$  of  $\mathcal{F}(n-k)(\gamma)$ . Again denoting

$$T_1 = \int_{m_1(\gamma)}^{M_2(R)} -\frac{d\lambda}{f_1(\lambda)}, \quad T_2 = \int_{m_1(\gamma)}^{M_2(R)} -\frac{d\lambda}{f_2(\lambda)} \quad (30)$$

the following theorem is obtained.

**Theorem 2.** The decaying time  $T$  of the motion of the system (1) from an initial point in the region  $\mathcal{F}(n-k)(R)$  to a point in the region  $\mathcal{F}(n-k)(\gamma)$  satisfies (27), and the decaying time  $T$  of the motion of the system (1) from an initial point in the region  $\|\phi\| \geq R$  to a point in the region  $\mathcal{F}(n-k)(\gamma)$  satisfies (29), where  $T_1, T_2$  are defined by (30).

By taking

$$f_1(v) = -\alpha v \quad f_2(v) = -\beta v \quad (\alpha > \beta) \quad (31)$$

one has

$$T_1 = \frac{1}{\alpha} \log \frac{M_2(R)}{m_1(\gamma)}, \quad T_2 = \frac{1}{\beta} \log \frac{M_2(R)}{m_1(\gamma)} \quad (32)$$

Particularly, when  $\phi_i = x_i, i = 1, \dots, k < n$  one obtains the formulae to estimate the decaying time for partial coordinates, and when  $\phi_i = x_i, i = 1, \dots, n$ , then one obtains the formulae to estimate the decaying time for total coordinates (9), (10) and (11).

The above method is used to solve the following example.

**Example**—Consider an autonomous system

$$\begin{aligned} \dot{x} &= x - a^2 x^3 - b^2 x y^2 \\ \dot{y} &= y - a^2 y x^2 - b^2 y^2 \end{aligned} \quad (a > b) \quad (33)$$

and its unique closed orbit

$$\phi = a^2 x^2 + b^2 y^2 - 1 = 0 \quad (34)$$

If one selects the Liapunov function with respect to  $\phi$  to be

$$V = (a^2 x^2 + b^2 y^2 - 1)^2 \quad (35)$$

then it may be asserted that:

(a) the system is asymptotically stable with respect to  $\phi = 0$ ;

(b) the decaying time  $T$  of the motion of the system from an initial point in the region  $|\phi| \leq \phi_0$  to a point in the region  $|\phi| \leq \varepsilon$  satisfies  $T \leq a^2 \log(1 + \varepsilon) \phi_0 / (1 + \phi_0) \varepsilon$ ;

(c) the decaying time  $T$  of the motion from an initial point in the region  $|\phi| \leq \phi_0$  to a point in the region  $|\phi| \leq \varepsilon$  satisfies  $T \geq b^2 \log(1 + \varepsilon) \phi_0 / (1 + \phi_0) \varepsilon$ .

### On the Estimation of Decaying Time for Linear System with Quasi-constant Coefficients

In the study of a practical dynamic system, one usually takes the linear system with constant coefficients as its first approximation. In general, the frequency method may be applied to estimate the time of transient process for the regulating system with constant coefficients. However, this method is only applicable to the case of single output under specific initial conditions. In addition, the method is not rigorous. This paper gives the formulae to estimate the decaying time in the general case, and the method is rigorous.

A large amount of work<sup>9-12</sup> is devoted to the estimation of decaying time for the asymptotically stable system

$$\dot{x}_S = p_{S1} x_1 + \dots + p_{SN} x_N, \quad S = 1, \dots, N \quad (36)$$

where the coefficients  $p_{ij}$  are constants. These results may be summarized as the following. For any given positive definite quadratic form

$$U = x' U x \quad (37)$$

there is a positive definite quadratic form

$$V = x' V x \quad (38)$$

such that

$$\left. \frac{dV}{dt} \right|_{(36)} = -U \quad (39)$$

If  $M_1$  and  $m_1$  are, respectively, the maximum and the minimum eigenvalues of the matrix  $V$ , and  $M$  and  $m$  are the maximum and the minimum eigenvalues of the matrix  $U$ , then the following results are obtained.

**Theorem 3.** The decaying time  $T$  of the motion of the system (36) from an initial point in

$$\sum_{S=1}^N x_S^2 = R^2$$

103/4

to a point in the region

$$\sum_{s=1}^N x_s^2 \leq r^2$$

satisfies the inequalities

$$\frac{m_1}{M_2} \log \frac{M_1 R^2}{m_1 r^2} \leq T \leq \frac{M_1}{m_2} \log \frac{M_1 R^2}{m_1 r^2} \quad (40)$$

In practice, it is of interest to select a suitable Liapunov function  $V$ , such that for the given system (36) the range defined by (40) is as accurate as it can be. It is very difficult to answer the above question in the general case. But if the system (36) is normal and the elementary divisors of the coefficient matrix  $P$  are all simple, it may be proved that when

$$V = \sum_{s=1}^N x_s^2$$

the equalities in (40) may be realized (i.e. the estimation is accurate).

Let the normal transformation be

$$y = Cx \quad (41)$$

where  $C$  is a matrix with real coefficients, and the system (36) is reduced to the normal system

$$\dot{y} = Jy \quad (42)$$

where

$$J = \begin{pmatrix} -\alpha_1 & & & & 0 \\ & \alpha_k & & & \\ & & -\beta_1 - \omega_1 & & \\ & & \omega_1 - \beta_1 & & \\ 0 & & & & -\omega_l \\ & & & & \omega_l - \beta_l \end{pmatrix} \quad (43)$$

$$V = x' C' c x \quad (44)$$

may be taken as a Liapunov function of the system (36). By means of (42) the following results may be proved.

**Theorem 4.** The decaying time  $T$  of the motion of the system (36), from an initial point in the  $(n-1)$ -dimensional ellipsoid  $V = V_0$  to a point in the ellipsoid  $V = \epsilon$ , satisfies

$$\frac{1}{2u} \log \frac{V_0}{\epsilon} \leq T \leq \frac{1}{2v} \log \frac{V_0}{\epsilon} \quad (45)$$

where  $u = \max(\alpha_i, \beta_j)$ ,  $v = \min(\alpha_i, \beta_j)$ . It is easy to select the initial points such that the equalities in (45) hold (i.e. this estimation is accurate).

In the following, the general formulae to estimate the decaying time is given.

All roots of the characteristic equation  $D(\lambda) = \det(R - \lambda I) = 0$  are assumed to have negative real parts. Let  $\lambda_1, \dots, \lambda_l$  be negative real roots, written as  $\lambda_i = -\alpha_i$  ( $i = 1, \dots, l$ ), and let  $\lambda_{l+1}, \dots, \lambda_n$  be the remaining roots, written as  $\beta_s \pm \omega_s i$  ( $S = 1, \dots, n-l/2 = k$ ), and the order of the corresponding elementary divisors be  $n_1, \dots, n_k$ .

It is known that there is a non-singular linear transformation

$y = Cx$  to reduce the system (36) to the normal form (42), in which

$$J = \begin{pmatrix} M_1 & & & 0 \\ & M_l & & \\ & & N_1 & \\ 0 & & & N_k \end{pmatrix} \quad (46)$$

$$M_i = \begin{pmatrix} -\alpha_i & 1 & 0 & \dots & 0 \\ 0 & -\alpha_i & 1 & \dots & 0 \\ & & & \ddots & \\ 0 & & & & 1 \\ & & & & -\alpha_i \end{pmatrix} \quad N_j = \begin{pmatrix} -\beta_i - \omega_i & 1 & 0 & \dots & 0 \\ \omega_i - \beta_i & 0 & 1 & \dots & 0 \\ & & \dots & & \\ 0 & 0 & \dots & -\beta_i - \omega_i & \\ 0 & 0 & \dots & \omega_i - \beta_i & \end{pmatrix} \quad (47)$$

If one writes the  $m_i \times m_i$  matrix

$$a^{(m_i)} = \begin{pmatrix} \frac{1}{\alpha_i} & \frac{1}{2\alpha_i^2} & \frac{1}{4\alpha_i^3} & \dots \\ \frac{1}{2\alpha_i^2} & \frac{1}{\alpha_i} \left[ 1 + \frac{1}{2\alpha_i^2} \right] & \dots \\ \frac{1}{4\alpha_i^3} & \dots & \dots \end{pmatrix} \quad (48)$$

as  $a_{s\sigma}^{(m_i)}$ , this is constructed according to the following rule:

(a) when  $s = \sigma$ ,  $a_{s\sigma}^{(m_i)}$  is equal to  $(1/\alpha)(1 + a_{s-1, \sigma}^{(m_i)})$ , and let

$$a_{11}^{(m_i)} = \frac{1}{\alpha_i};$$

(b) when

$$s > \sigma, a_{s\sigma}^{(m_i)} = \frac{1}{2\alpha_i} [a_{s, \sigma-1}^{(m_i)} + a_{s-1, \sigma}^{(m_i)}];$$

(c)

$$a_{s\sigma}^{(m_i)} = a_{\sigma s}^{(m_i)}$$

Thus the matrix is completely defined through the eigenvalues  $-\alpha_i$  and the order of its elementary divisor. The maximum eigenvalue of the matrix  $a^{(m_i)}$  is assumed to be  $v_i$

$$\begin{cases} \text{when } m_i = 1, v_i = \frac{1}{\alpha_i} \\ \text{when } m_i = 2, v_i = \frac{1}{\alpha_i} + \frac{1}{4\alpha_i^2} \left[ \frac{1}{\alpha_i} + \left( 4 + \frac{1}{\alpha_i^2} \right)^{\frac{1}{2}} \right] \end{cases} \quad (49)$$

Following the method of construction of the matrix  $a^{(m_i)}$ , the  $2n_i \times 2n_i$  matrix  $d^{2n_i}$  may be constructed in the following manner

(a)  $d_{2i-1, 2j}^{(2n_i)} = d_{2i, 2j-1}^{(2n_i)} = 0, \quad i, j = 1, \dots, n_i$

(b)  $d_{2i-1, 2j-1}^{(2n_i)} = d_{2j-1, 2i-1}^{(2n_i)} = a_{ij}^{(n_i)} \quad i, j = 1, \dots, n_i$

(c) to replace  $\alpha_i$  by  $\beta_i$  in the matrix  $a^{(n_i)}$ .

For example  $n_i = 2$ , one has

$$a^{(4)} = \begin{pmatrix} a_{11}^{(2)} & 0 & a_{12}^{(2)} & 0 \\ 0 & a_{11}^{(2)} & 0 & a_{12}^{(2)} \\ a_{21}^{(2)} & 0 & a_{22}^{(2)} & 0 \\ 0 & a_{21}^{(2)} & 0 & a_{22}^{(2)} \end{pmatrix}$$



Obviously, the formula for the maximum eigenvalue of  $d^{(2n_i)}$  is the same as that of  $a^{(m_i)}$  in which  $\alpha_i$  are replaced by  $\beta_i$ . Consider the Liapunov function for system (36) to be

$$V = x' C' A C x \quad (50)$$

where  $C$  is the normal transformation matrix, and

$$A = \begin{pmatrix} a^{(m_1)} & & & 0 \\ & a^{(m_1)} & & \\ & & d^{(2n_1)} & \\ 0 & & & d^{(2n_k)} \end{pmatrix} \quad (51)$$

It is not difficult to prove that  $V$  satisfies

$$\left. \frac{dV}{dt} \right|_{2.1} \leq -\frac{2}{\nu} V \quad (52)$$

where  $\nu$  is the maximum eigenvalue of the matrix  $A$ , and it can be calculated by the aforesaid method. When  $m_i = 1$  or  $m_i = 2$  it can be calculated through (49). If the maximum and the minimum eigenvalues of the symmetric matrix  $C'AC$  are assumed to be  $M$  and  $m$  respectively, then the following theorem is obtained.

**Theorem 5.** The decaying time  $T$  of the motion of the system (36) from an initial point in the  $(N-1)$ -dimensional ellipsoid  $V = V_0$  to a point in the  $(N-1)$ -dimensional ellipsoid  $V = \varepsilon$  satisfies

$$T \leq \frac{\nu}{2} \log \frac{V_0}{\varepsilon}$$

The decaying time  $T$  of the motion of the system (36) from an initial point in the sphere

$$\sum_{s=1}^N x_s^2 = R^2$$

to a point in the sphere

$$\sum_{s=1}^N x_s^2 = r^2$$

satisfies

$$T \leq \nu \log \frac{MR^2}{mr^2}$$

Moreover, the system

$$\dot{x} = px + X(x, t) \quad (53)$$

is considered, where  $X$  is a vector function which contains the non-linear terms and the unknown components. If one constructs a Liapunov's function (50) of its principal linear system

$$\dot{x} = px \quad (54)$$

and if one assumes  $X$  to satisfy the inequality

$$|\text{grad } V \cdot X| < bx' C' C x, \quad (b < 2) \quad (55)$$

then the following results are obtained.

**Theorem 6.** The decaying time  $T$  of the motion of the system (53) from an initial point in the sphere

$$\sum_{s=1}^N x_s^2 = R^2$$

to a point in the sphere

$$\sum_{s=1}^N x_s^2 = r^2$$

satisfies

$$T \leq \frac{\nu}{2(1-b/2)} \log \frac{MR^2}{mr^2} \quad (56)$$

As an application of this theorem, an example of a forced oscillation is considered.

*Example*—Consider the system

$$\dot{u} = pu + \varepsilon U(u, \dots, u_n) + F(t) \quad (57)$$

where  $p$  is assumed to have all its eigenvalues with negative real parts,  $\varepsilon$  is a small parameter,  $U$  is continuously differentiable and  $F(t)$  is the forcing term with period  $\tilde{T}$ .

Let the system (56) have a periodic solution

$$u_s = u_s^0(t), \quad u_s^0(t) = u_s^0(t + \tilde{T}) \quad (s = 1, \dots, N) \quad (58)$$

and let the linear transformation

$$y = Cu \quad (59)$$

transform the system  $\dot{u} = pu$  into its normal form

$$\dot{y} = Jy \quad (60)$$

By means of the transformation (59) the system (56) was reduced to a system

$$\dot{y} = Jy + \varepsilon Y(y) + \Phi(t) \quad (61)$$

where  $\Phi(t) = CF(t)$  has the same period as  $F(t)$ . Under this transformation, the periodic solution (58) is reduced to

$$y_s = y_s^0(t) = \sum_{\sigma=1}^N \cos u_\sigma^0(t) \quad (62)$$

Consider the perturbations  $x_s = y_s - y_s^0(t)$  then  $x$  satisfies

$$\dot{x} = Jx + \varepsilon q(t)x + \varepsilon X(x, t) \quad (63)$$

where  $q(t)$  is a periodic matrix with period  $T$ , and it may be evaluated through  $Y(t)$  and  $y_s^0(t)$ . If one takes  $\xi = C^{-1}x$ , then  $\xi = u - u^0(t)$  is the perturbation vector in  $u$  space.

By means of the above method the matrix  $A$  is constructed, with its maximum eigenvalue  $\nu$ , and the maximum and minimum eigenvalues of the matrix  $C'AC$  are  $M$  and  $m$  respectively. The following results are obtained.

**Theorem 7.** The decaying time  $T$  of the motion of the system (57) from an initial point in the  $R$  neighbourhood of the periodic solution (58) to a point in the  $r$  neighbourhood of the periodic solution (58) satisfies

$$T \leq \frac{\nu}{2 \left[ 1 - \frac{(b+\varepsilon)}{2} \varepsilon \right]} \log \frac{MR^2}{mr^2} \quad (64)$$

where the term  $X$  in (62) satisfies

$$|\text{grad } \bar{V} \cdot X| < bx'x \quad (b < 2\bar{V} = x'Ax) \quad (65)$$

and  $C$  is the maximum eigenvalue of the matrix  $q(t) + [q(t)]$  when  $t \in [0, \tilde{T}]$ . By the above-mentioned  $r$  neighbourhood of

103/6

the periodic solution (58), is meant the set of points which satisfies the inequality

$$\sum_{s=1}^N [u_s - u_s^0(t)]^2 \leq r^2$$

### On the Estimation of Decaying Time for Quasi reducible Linear System

In the study of the dynamic systems, one may sometimes fail to approximate it by a linear system with constant coefficients. In this case one may take a reducible system as its approximations, and construct the corresponding Liapunov function and estimate the decaying time.

Consider the non-linear system

$$\dot{x} = p(t)x + X(x, t) \quad (66)$$

Let the linear approximation system

$$\dot{x} = p(t)x \quad (67)$$

be a reducible system. Assuming the characteristic number to be all positive, there is a Liapunov transformation

$$y = C(t)x \quad (68)$$

which transforms the system (67) into its real normal form

$$\dot{y} = \tau y \quad (69)$$

By means of the method mentioned in the previous section,

$$V = x'c'ACx \quad (70)$$

is taken as a Liapunov's function for the system (67). Since (68) is the Liapunov's transformation when  $t \geq t_0$ , the maximum and minimum eigenvalues  $M$  and  $m$  cannot equal zero. Obviously, the maximum eigenvalue  $\nu$  of  $A$  can be calculated through the characteristic numbers of (67) by the same method. Parallel to Theorem 2 the following results may be obtained.

**Theorem 8.** The decaying time  $T$  of the motion of the system (67) from an initial point in the sphere

$$\sum_{s=1}^N x_s^2 = R^2$$

to a point in the sphere

$$\sum_{s=1}^N x_s^2 = r^2$$

satisfies

$$T \leq \frac{\nu}{2} \log \frac{MR^2}{mr^2} \quad (71)$$

Furthermore, if  $X$  satisfies the condition

$$|\text{grad } V \cdot X| < bx'c'Cx \quad (b < 2) \quad (72)$$

then one has the following results.

**Theorem 9.** The decaying time  $T$  of the motion of the system (66) from an initial point in the sphere

$$\sum_{s=1}^N x_s^2 = R^2$$

to a point in the sphere

$$\sum_{s=1}^N x_s^2 = r^2$$

satisfies

$$T \leq \frac{\nu}{2(1-b/2)} \log \frac{MR^2}{mr^2} \quad (b < 2) \quad (73)$$

where  $X$  satisfies (72).

Since the linear system with periodic coefficients is a reducible system, and its characteristic number can be represented through its characteristic exponentials, the results in this section can be applied to the general periodic systems.

### Appendix

#### Proof of the Fundamental Theorem

Obviously the system is stable with respect to (2) taking on zeros.

For any given  $\eta > 0$  the region  $H \geq \|\phi\|_K \geq \eta$  is considered. From the conditions mentioned in the theorem, the function  $V$  takes on maximum  $M > 0$  and minimum  $m > 0$  and the negative definite function  $dV/dt|_{1.1}$ , with respect to (2), takes on maximum  $-\alpha < 0$ . Let  $T = (M - m)/\alpha + 2\beta$ , where  $\beta$  is any arbitrary positive number. This is the required  $T$  and it is independent of the initial conditions. Thus, part (A) of the fundamental theorem holds.

Parallel to Theorem 3, from the conditions of the fundamental theorem, the following lemmas can be proved.

**Lemma 1.** If the system (1) is equi-asymptotically stable with respect to (2), taking zeros for the initial values in  $\mathcal{F}(n-k)(H)$ , then there is  $\psi(\tau)$  such that the motions, defined by the initial points of the above-mentioned region, satisfy

$$(a) \quad \|\phi[x(t_0 + \tau, x_1^0, \dots, x_n^0, t_0)]\|_K \leq \psi(\tau)$$

$$(b) \quad \lim_{\tau \rightarrow \infty} \psi(\tau) = 0 \quad \psi'(\tau) \leq 0 \quad \tau > 0$$

**Lemma 2.** For any given two positive functions  $M(\eta)$  and  $\psi(\eta)$ , where  $M(\eta)$  is an increasing function and  $\lim_{\eta \rightarrow 0} \psi(\eta) = 0$  there is a function  $G(\eta)$  such that

$$(a) \quad G(\eta) > 0, \quad G'(\eta) > 0 \quad \text{as } \eta > 0$$

$$(b) \quad G(0) = G'(0) = 0$$

$$(c) \quad \int_0^\infty G[\psi(\tau)] d\tau < \infty \quad \int_0^\infty G'[\psi(\tau)] M(\tau) d\tau < \infty$$

**Lemma 3.** If the system (1) is equi-asymptotically stable with respect to (2) taking on zeros, and the rank of the matrix (19) is  $K$ , then there are two positive constants  $A$  and  $\lambda$  independent of  $t_0$  and  $\phi^0$  such that

$$\left| \frac{\partial \phi^2}{\partial t_0} \right| < A e^{\lambda \tau} \quad \left| \frac{\partial \phi^2}{\partial \phi_s^0} \right| < A e^{\lambda \tau}$$

where

$$\phi^2 = \sum_{s=1}^K \phi_s^2(t, \phi_1^0, \dots, \phi_n^0, t_0)$$

$\phi_i^0$  ( $i = 1, \dots, n$ ) are the initial values and  $t$  is replaced by  $t_0 + \tau$ . Parallel to Theorem 3, the Liapunov function may be taken as

$$V = \int_0^{\infty} G \left[ \sum_{s=1}^K \phi_s^2(t+\tau, \phi_1, \dots, \phi_n, t) \right] d\tau$$

then one has

$$\left| \frac{\partial V}{\partial \phi_0} \right| = \left| \int_t^{\infty} G' \left[ \phi^2(t+\tau, \phi_1, \dots, \phi_n, t) \frac{\partial \phi^2}{\partial \phi_i} d\tau \right] \right| < \infty$$

(i = 1, ..., K)

It is convergent and uniformly bounded. This implies that  $V$  satisfies the Lipschitz condition and thus  $V$  has an infinitely small upper bound.

$$V > \frac{1}{2L} [\phi_1^2 + \dots + \phi_K^2] G \left[ \frac{1}{2} \sum_{i=1}^K \phi_i^2 \right]$$

where  $L$  is the Lipschitz constant. It implies  $V$  to be positive definite

$$\left. \frac{dV}{dt} \right|_{1.1} = -G[\phi_1^2 + \dots + \phi_K^2]$$

Hence

$$\left. \frac{dV}{dt} \right|_{1.1}$$

is negative definite with respect to (2).

The proof is complete.

#### References

<sup>1</sup> LIAPUNOV, A. M. *The General Problem of Stability of Motion*. 1950. Gostekhizdat

<sup>2</sup> MASSERA, J. L. Contributions to stability theory. *Ann. Math.* 64, No. 1 (1956)

<sup>3</sup> MALKIN, I. G. Contribution to the theory of the invertibility of Liapunov's theorem on asymptotic stability. *Prikl. Mat. i Mekh.* XVIII, B<sub>2</sub> (1954)

<sup>4</sup> ZUBOV, V. I. *The Methods of A. M. Liapunov and their Application*. 1957. Leningrad; Gos. Univ.

<sup>5</sup> KRASOVSKIY, N. N. *Certain Problems in the Theory of Stability of Motion*. 1959. Fizmatgiz

<sup>6</sup> CH'ING YÜAN-HSÜN. *A Lecture on the General Problem of Stability of Motion*. 1958. Peking

<sup>7</sup> KALMAN, R. E. and BERTRAM, J. E. Control system analysis and design via the 'Second Method' of Liapunov. *Trans. Amer. Soc. mech. Engrs Ser. D*, 82, No. 2 (1960)

<sup>8</sup> RUMYANTSEV, V. V. Stability of motion in relation to part of the variables. *Vestnik. Moskov. Univ.* B<sub>4</sub> (1957)

<sup>9</sup> HUANG LIN. Problems in estimation of damping time for multi-dimensional non-linear systems. *Acta Sci. Natur. Univ.* VI, No. 1 (1960)

<sup>10</sup> CHETAYEV, N. G. The choice of parameters for a stable mechanical system. *Prikl. Mat. i Mekh.* 15, B<sub>3</sub> (1951)

<sup>11</sup> LETOV, A. M. *The Stability of Non-linear Controlled Systems*. 1955. Gostekhizdat

<sup>12</sup> CHANG SSU-YING. Estimated solutions to systems of differential equations for accumulation, perturbation and stability of motion, over a finite time interval. *Prikl. Mat. i Mekh.* Vol. XXIII, B<sub>4</sub> (1959)

<sup>13</sup> TROYTSKIY, V. A. Canonical transformations of the equations of automatic control theory. *Prikl. Mat. i Mekh.* XXI, B<sub>4</sub> (1957)

<sup>14</sup> CHETAYEV, N. G. *Stability of Motion*. Gostekhizdat

# Quasi-invariant Hybrid Multi-parameter Control Loops

V. STREJC *Czech*

## Introduction

Previous papers by the author<sup>1-3</sup> contain the general theory of the synthesis of control systems and of the compensation of the effects of disturbances in hybrid, multi-parameter control loops, with due consideration of the conditions of autonomy, invariance and the finite number of control steps. A control loop is regarded as hybrid if the function of the controller is performed by a discrete filter (digital correcting member), the realization of which is assumed to be attainable by an automatic digital computer and a continuously-acting controller. In practice, hybrid control loops can be formed by the addition of an automatic computer to control loops containing continuously-acting controllers. This arrangement is made either in cases where it is necessary to improve the quality of control and to attain a higher stage of complex automation that would be difficult or too costly to realize by other means, or in newly designed control systems with the automatic computer as the principal technical means of realizing automation and in which the simple, continuously-acting controllers are used as a stand-by for sustaining the operation of the control system in the case of an outage of the automatic computer.

In practical applications the case of a multi-parameter control system may occur frequently where the desired values of the controlled variables remain constant (their relative deviations being zero), and the task of the control system is confined to the compensation of the effects of disturbances. If a control-system structure, according to *Figure 1*, is selected for a multi-parameter control loop of this kind, the conditions of invariance cannot be fulfilled. However, the existence of a solution will be presented according to which only the controlled variables  $x_i, i = k$ , are influenced by disturbances  $K_k$  with the possibility of determining the limits of this influencing, according to the selected criterion of the quality of control, or according to other suitable control conditions. Let control loops of this kind be designated as quasi-invariant control loops.

For a control loop according to *Figure 1*:

$$[K_u^*(z, 0)] = \{[1] + [\Omega^*(z, 0)][P^*(z, 0)]\}^{-1} [\Omega_u^*(z, 0)] \quad (1)$$

$$[K_u^*(z, \varepsilon)] = [\Omega_u^*(z, \varepsilon)] - [\Omega^*(z, \varepsilon)][P^*(z, 0)][K_u^*(z, 0)] \quad (2)$$

where

$$[\Omega(p)] = \{[1] + [S(p)][R(p)]\}^{-1} [G(p)]$$

$$[\Omega_u(p)] = \{[1] + [S(p)][R(p)]\}^{-1} [G_u(p)]$$

In the compensation of disturbance effects  $[K_u^*(z_1, \varepsilon)]$  and  $[K_u^*(z_1, 0)]$  are the matrices of the transfer functions of closed control loops with the elements of the matrices expressed as discrete Laplace transforms ( $Z$  transforms). In eqn (2)  $\varepsilon$  stands for the relative value of the independent time variable that

during one interval of sampling attains the value of  $\varepsilon < 0 \div 1 >$ . The sampling interval  $T$  is constant, and let the sampling be synchronous at all points of the control loop.

$[u(t)]$  is the  $(\xi; 1)$  type column matrix of the disturbances

$[x(t)]$  is the  $(\nu; 1)$  type column matrix of the controlled variables

$[S(p)]$  is the  $(\nu; \mu)$  type rectangular matrix,  $\mu \geq \nu$  of the transfer functions of the controlled system containing a servomotor and a final control element

$[G(p)]$  is the  $(\nu; \mu)$  type rectangular matrix,  $\mu \geq \nu$ , of the transfer functions of a controlled system containing a servomotor, final control element and a holding member

$[G_u(p)]$  is the  $(\nu; \xi)$  type rectangular matrix,  $\xi \leq \nu$ , of the transfer functions of the controlled system containing a holding member

$[P^*(p)]$  is the  $(\mu; \nu)$  type rectangular matrix,  $\mu \geq \nu$ , of the transfer functions of digital correcting members

$[R(p)]$  is the  $(\mu; \nu)$  type rectangular matrix,  $\mu \geq \nu$ , of the transfer functions of continuously-acting controllers

## Conditions of Stability

As it is desirable to express the quality of control by the requirements upon the transfer functions in matrix  $[K_u^*(z, 0)]$ , the matrix  $[P^*(z, 0)]$  is the function of matrix  $[K_u^*(z, 0)]$ . It can be calculated from eqn (1) that

$$[P^*(z, 0)] = [\Omega^*(z, 0)]^{-1} \{[\Omega_u^*(z, 0)] - [K_u^*(z, 0)]\} [\Omega_u^*(z, 0)]^{-1} \quad (3)$$

By substituting relation (3) for  $[P^*(z, 0)]$  into eqn (2)

$$[K_u^*(z, \varepsilon)] = [\Omega_u^*(z, \varepsilon)] - \frac{\Delta_{\Omega A}(z, 0)}{\Delta_{\Omega B}(z, 0)} [\Omega^*(z, \varepsilon)] [\omega^*(z, 0)] \{[\Omega_u^*(z, 0)] - [K_u^*(z, 0)]\} \quad (4)$$

where

$$\frac{\Delta_{\Omega A}(z, 0)}{\Delta_{\Omega B}(z, 0)} [\omega^*(z, 0)] = [\Omega^*(z, 0)]^{-1} \quad (5)$$

As the continuously-acting controllers, the transfer functions of which have the matrix  $[R(p)]$ , are determined by the condition of all loops of the control system remaining stable in the case of a computer outage, it may be stated that the elements of matrix  $[\Omega_u^*(z, \varepsilon)]$  will always be stable.

On the other hand, the elements of the second term on the right-hand side of eqn (4) can be unstable if the polynomial  $\Delta_{\Omega B}(z, 0)$ , which is the numerator of the determinant  $\dagger \Delta_{\Omega}(z, 0)$

<sup>†</sup> The 'determinant' of the  $(m; n)$  type rectangular matrix  $A$ , with  $m - n$ , is to be considered as being identical with the determinant of matrix  $A^T A$  where  $A^T$  is the matrix transposed towards matrix  $A$ .

120/2

of matrix  $[\Omega^*(z, 0)]$ , has its zero outside the zone of stability. The unstable zeros of polynomial  $\Delta_{\Omega B}(z, 0)$  must be assumed, however, to be compensated by the numerator of the elements of matrix  $\{[\Omega_u^*(z, 0)] - [K_u^*(z, 0)]\}$  which, as can be seen from eqn (3), is a cofactor of the matrix  $[P^*(z, 0)]$ . In accordance with the assumptions stated previously, the elements of matrices  $[\Omega^*(z, \varepsilon)]$  and  $[\omega^*(z, 0)]$  in eqn (4) are stable, while the stability of the elements of matrix  $[K_u^*(z, 0)]$  must be presupposed. On the basis of the above findings, it is possible to state the condition of the stability of a hybrid, multi-parameter control system for the compensation of disturbance effects, as follows:

$$[\Omega_u^*(z, 0)] - [K_u^*(z, 0)] = \Delta_{\Omega B}^-(z, 0) [D_u^*(z, 0)] \quad (6)$$

where  $[D_u^*(z, 0)]$  is the matrix of auxiliary functions that must be determined in more detail, while  $\Delta_{\Omega B}^-(z, 0)$  follows from equation

$$\Delta_{\Omega B}(z, 0) = \Delta_{\Omega B}^+(z, 0) \Delta_{\Omega B}^-(z, 0) \quad (7)$$

where  $\Delta_{\Omega B}^+(z, 0)$  signifies the product of the stable, and  $\Delta_{\Omega B}^-(z, 0)$  the product of the unstable root factors of the numerator of the determinant of matrix  $[\Omega^*(z, 0)]$ .

Introduce

$$[K_u^*(z, 0)] = [\bar{Q}_u^*(z, 0)] [(1 - z^{-1})^m C^*(z, 0)] \quad (8)$$

where  $[\bar{Q}_u^*(z, 0)]$  is the matrix of auxiliary functions, and  $[(1 - z^{-1})^m C^*(z, 0)]$  is a diagonal matrix of the  $(\xi; \xi)$  type, the elements of which should be polynomials independent of the properties of control loop members. Let these elements be the denominators of the Z transforms of the general form of the disturbances

$$[U^*(z, 0)] = \left[ \frac{F^*(z, 0)}{(1 - z^{-1})^m C^*(z, 0)} \right] \quad (9)$$

Now, eqn (6) can be rewritten in the form

$$[\Omega_u^*(z, 0)] - [\bar{Q}_u^*(z, 0)] [(1 - z^{-1})^m C^*(z, 0)] = \Delta_{\Omega B}^-(z, 0) [D_u^*(z, 0)] \quad (10)$$

After substituting relations (8) and (10) into eqn (3), matrix  $[P^*(z, 0)]$  expressed by this equation will acquire the form

$$[P^*(z, 0)] = \frac{\Delta_{\Omega A}(z, 0)}{\Delta_{\Omega B}^+(z, 0)} [\omega^*(z, 0)] [D_u^*(z, 0)] \{[\bar{Q}_u^*(z, 0)] [(1 - z^{-1})^m C^*(z, 0)]\}^{-1} \quad (11)$$

Now, let the real functions in eqn (2) be marked with the subscript  $s$ , and the imaginary functions in eqn (11) with the subscript  $p$ . After substituting relation (11) into eqn (2), it follows

$$[K_{us}^*(z, \varepsilon)] = [\Omega_{us}^*(z, \varepsilon)] - \frac{\Delta_{\Omega A}(z, 0)}{\Delta_{\Omega B}^+(z, 0)} [\Omega_s^*(z, \varepsilon)] [\omega_p^*(z, 0)] [D_{up}^*(z, 0)] \quad (12)$$

provided that

$$[\bar{Q}_{up}^*(z, 0)]^{-1} [\bar{Q}_{us}^*(z, 0)] = [1] \quad (13)$$

Assumption (13) can be fulfilled only if the zeros and the poles of the determinant of matrix  $[\bar{Q}_u^*(z, 0)]$  are inside the zone of stability.

The following holds for the elements of matrices in eqn (10):

$$\Omega_{u, ik}^*(z, 0) - \bar{Q}_{u, ik}^*(z, 0) (1 - z^{-1})_{kk}^m C_{kk}^*(z, 0) = \Delta_{\Omega B}^-(z, 0) D_{u, ik}^*(z, 0) \quad (14)$$

Introduce

$$\left. \begin{aligned} K_{u, ik}^*(z, 0) &= \frac{K_{uB, ik}^*(z, 0)}{K_{uA, ik}^*(z, 0)} & \bar{Q}_{u, ik}^*(z, 0) &= \frac{\bar{Q}_{uB, ik}^*(z, 0)}{\bar{Q}_{uA, ik}^*(z, 0)} \\ \Omega_{u, ik}^*(z, 0) &= \frac{M_{uB, ik}^*(z, 0)}{M_{uA, ik}^*(z, 0)} & D_{u, ik}^*(z, 0) &= \frac{D_{uB, ik}^*(z, 0)}{D_{uA, ik}^*(z, 0)} \end{aligned} \right\} \quad (15)$$

where fractions (15) represent the ratios of polynomials in  $z^{-1}$  with a finite number of terms. By using relations (15), eqn (14) can be rewritten in the form

$$\frac{M_{uB, ik}^*(z, 0)}{M_{uA, ik}^*(z, 0)} \frac{\bar{Q}_{uB, ik}^*(z, 0)}{\bar{Q}_{uA, ik}^*(z, 0)} (1 - z^{-1})_{kk}^m C_{kk}^*(z, 0) = \Delta_{\Omega B}^-(z, 0) \frac{D_{uB, ik}^*(z, 0)}{D_{uA, ik}^*(z, 0)} \quad (16)$$

Similarly, the following holds for the elements of the matrices in eqn (8)

$$\frac{K_{uB, ik}^*(z, 0)}{K_{uA, ik}^*(z, 0)} = \frac{\bar{Q}_{uB, ik}^*(z, 0)}{\bar{Q}_{uA, ik}^*(z, 0)} (1 - z^{-1})_{kk}^m C_{kk}^*(z, 0) \quad (17)$$

In the case of

$$D_{uA, ik}^*(z, 0) = M_{uA, ik}^*(z, 0) \bar{Q}_{uA, ik}^*(z, 0) \quad (18)$$

eqn (16) will acquire the form

$$\bar{Q}_{uA, ik}^*(z, 0) M_{uB, ik}^*(z, 0) - \bar{Q}_{uB, ik}^*(z, 0) M_{uA, ik}^*(z, 0) (1 - z^{-1})_{kk}^m C_{kk}^*(z, 0) = \Delta_{\Omega B}^-(z, 0) D_{uB, ik}^*(z, 0) \quad (19)$$

Denote

$$\left. \begin{aligned} \Delta_{\Omega B}^-(z, 0) &= 1 + \sum_{v=1}^{L_1} b_v z^{-v} & \bar{Q}_{uA, ik}^*(z, 0) &= 1 + \sum_{v=1}^{Q_A} p_v z^{-v} \\ C_{kk}^*(z, 0) &= 1 + \sum_{v=1}^C c_v z^{-v} & \bar{Q}_{uB, ik}^*(z, 0) &= \sum_{v=1}^{Q_B} q_v z^{-v} \\ M_{uA, ik}^*(z, 0) &= 1 + \sum_{v=1}^{M_A} \alpha_v z^{-v} & D_{uB, ik}^*(z, 0) &= \sum_{v=1}^{D_B} d_v z^{-v} \\ M_{uB, ik}^*(z, 0) &= \sum_{v=1}^{M_B} \beta_v z^{-v} \end{aligned} \right\} \quad (20)$$

Let the degree of polynomial  $\bar{Q}_{uB, ik}^*(z, 0)$  be assumed as

$$Q_B = Q + N \quad (21)$$

where  $Q$  is the lowest possible degree of the polynomial  $\bar{Q}_{uB, ik}^*(z, 0)$  that follows from eqn (19), and  $N$  the number of degrees of freedom.

Assuming that

$$Q_A + M_B \leq Q_B + M_A + m + C \quad (22)$$

the degree of the resultant polynomial on the left-hand side of eqn (19) will be

$$Q_B + M_A + m + C = L_1 + D + N \quad (23)$$

From eqn (23) follows the degree of polynomial  $\bar{Q}^*_{uB, ik}(z, 0)$

$$Q_B = L_1 + N = Q + N$$

and the degree of polynomial  $D^*_{uB, ik}(z, 0)$

$$D = M_A + m + C \quad (24)$$

$$D_B = D + N \quad (25)$$

By comparing the coefficients of the equal powers  $z^{-1}$  in the resultant polynomials on both sides of eqn (19), the system of  $Q_B + D$  linear algebraic equations is obtained where

$$Q_B + D = L_1 + D + N \quad (26)$$

To this system of equations it is necessary to add further  $N + Q_A$  equations of conditions that follow from the selected conditions of control. The system of  $Q_B + D_B + Q_A$  equations obtained in this way determines the coefficients of polynomials  $\bar{Q}^*_{uB, ik}(z, 0)$ ,  $\bar{Q}^*_{uA, ik}(z, 0)$  and  $D^*_{uB, ik}(z, 0)$  of the auxiliary functions, provided that the determinant of the equation system does not equal zero. The number of such coefficients is

$$Q_A + Q_B + D_B = L_1 + D + 2N + Q_1 \quad (27)$$

A more detailed analysis would prove that  $\beta_1 = q_1$  and  $d_1 = 0$  holds generally, and consequently the number of conditions necessary for the determination of the coefficients of auxiliary functions may be reduced by two.

The solution is somewhat simplified if it can be stated that

$$K^*_{uA, ik}(z, 0) = \bar{Q}^*_{uA, ik}(z, 0) = M^*_{uA, ik}(z, 0) \quad (28)$$

It follows

$$D^*_{uA, ik}(z, 0) = M^*_{uA, ik}(z, 0) \quad (29)$$

and eqn (16) will assume the form

$$\begin{aligned} M^*_{uB, ik}(z, 0) - \bar{Q}^*_{uB, ik}(z, 0)(1 - z^{-1})^m C^*_{kk}(z, 0) \\ = \Delta_{\bar{Q}_B}(z, 0) D^*_{uB, ik}(z, 0) \end{aligned} \quad (30)$$

The above simplification does not allow the inclusion, in the characteristic equation of transfer functions  $K^*_{u, ik}(z, 0)$ , of additional requirements above those asserted in the characteristic equation of the terms  $\bar{Q}^*_{u, ik}(z, 0)$ .

After the determination of all elements of matrix  $[\bar{Q}^*_u(z, 0)]$  it is necessary to check the zeros in the numerator of the determinant of this matrix.

Now, the conditions of stability can be summarized as:

**Theorem 1**—In the defined hybrid control loop where  $\Delta_{\bar{Q}_B}(z, 0)$  is the product of the unstable root factors of the numerator of the determinant of matrix  $[\Omega^*(z, 0)]$ , with the poles of this determinant lying within the stable zone of plane  $z$ , the transfer functions of the control loops, i.e. the elements of matrix  $[K^*_u(z, \varepsilon)]$ , are stable, provided that: (a) the poles and zeros of the determinant of matrix  $[\bar{Q}^*_u(z, 0)]$  lie within the stable zone of plane  $z$ , (b) the matrix  $[\bar{Q}^*_u(z, 0)]$  is in accordance with the equation of conditions (10) and none of its elements is equal to zero, and (c) the poles of the elements of matrix  $[D^*_u(z, 0)]$  in eqn (10) also lie in the stable zone of plane  $z$ . These conditions are necessary and sufficient.

### The Conditions of Zero Offset

Provided that the conditions of stability are fulfilled, it is possible to state the condition of zero offset according to the theorem of finite values by the following equation:

$$\lim_{z \rightarrow 1} K^*_{u, ik}(z, 0) = 0 \quad (31)$$

The above condition can be fulfilled if the value of  $m$  in the general relation (8) is at least  $m = 1$ . In other words, the product of the numerator root factors of transfer functions  $K^*_{u, ik}(z, 0)$  must necessarily contain the factor  $(1 - z^{-1})$ .

### Quasi-invariant Control Loops

Provided that all zeros of the determinant of matrix  $[\Omega^*(z, 0)]$  lie within the stable zone of plane  $z$ , the following substitution can be made in eqn (14):

$$\Delta_{\bar{Q}_B}(z, 0) = 1 \quad (32)$$

If the selectable functions are stated as

$$\bar{Q}^*_{u, ik}(z, 0) = 0 \quad \text{for } i \neq k \quad (33)$$

it follows

$$D^*_{u, ik}(z, 0) = \Omega^*_{u, ik}(z, 0) \quad \text{for } i \neq k \quad (34)$$

and the remaining functions  $D^*_{u, ik}(z, 0)$ ,  $i = k$ , can be determined by the same method as shown earlier in this paper. In this case the matrix  $[\bar{Q}^*_u(z, 0)]$  will be a diagonal matrix and consequently, with regard to eqn (8),  $[K^*_u(z, 0)]$  will also be a diagonal matrix. This solution permits a situation to be reached where disturbances  $U^*_k(z, 0)$ , ( $k = 1, 2, \dots, \xi$ ), (where  $\xi \geq \nu$  and  $\nu$  is the number of controlled variables) will influence only the controlled variables  $X^*_i(z, \varepsilon)$ , for which  $i = k$ , and will have no influence upon the controlled variables  $X^*_i(z, \varepsilon)$ , for which  $i \neq k$ . If  $\xi < \nu$ , the effect of disturbances  $U^*_k(z, 0)$  will be confined to the controlled variables  $X^*_i(z, \varepsilon)$ , for which  $i = k$  and  $i = 1, 2, \dots, \xi$  and with no effect upon the controlled variables  $X^*_i(z, \varepsilon)$ , for which  $i \neq k$  and also those for which  $i = k$  but  $i = \xi + 1, \xi + 2, \dots, \nu$ .

Due to this solution the transfer functions in diagonal matrix  $[K^*_u(z, 0)]$  can have an arbitrary number of degrees of freedom that can be utilized for the fulfilment of further conditions of control, or for the compliance with a suitable criterion of the quality of control. In this way it is possible to reach a solution at which the effect of disturbances, that cannot be eliminated by the introduction of condition (33), is kept within admissible limits.

In principle, this method of the compensation of disturbance effects can also be applied to continuously-acting control systems. However, up to the time of writing this paper, this possibility has not been mentioned in any technical literature accessible to the author.

### Finite Number of Control Steps

In the compensation of disturbance effects a multi-parameter control loop complies with the requirement of a finite number of control steps, if the same requirement is complied with by all components of output signals  $X^*_i(z, \varepsilon)$  of the controlled system. In this case

$$X^*_{ik}(z, \varepsilon) = K^*_{u, ik}(z, \varepsilon) U^*_k(z, 0) \quad (35)$$

120/4

The finite number of control steps is understood as the number of sampling intervals, at the attainment of which the offset is permanently zero or constant at any one instant of sampling. In the intervals between the instants of sampling this condition need not be fulfilled. If it is possible to express the  $Z$  transforms of the general forms of the disturbances by eqn (9), and the matrix of transfer functions  $[K_u^*(z, 0)]$  by eqn (8), it follows

$$X^*(z, 0) = [\bar{Q}_u^*(z, 0)] F^*(z, 0) \quad (36)$$

It follows from eqn (36) that the requirement of the finite number of control steps can be complied with only if the elements of matrix  $[\bar{Q}_u^*(z, 0)]$  are polynomials having a finite number of terms.

In this case it is necessary to substitute in eqns (18) and (19)

$$\bar{Q}_{uA, ik}^*(z, 0) = 1 \quad (37)$$

Then for individual components

$$X_{ik}^*(z, 0) = \bar{Q}_{u, ik}^*(z, 0) F_k^*(z, 0) \quad (38)$$

and the degree of polynomial  $X_{ik}^*(z, 0)$  is

$$X = L_1 + N + F \quad (39)$$

The transform of the controlled variable is

$$X_i^*(z, 0) = \sum_{k=1}^{\xi} X_{ik}^*(z, 0) \quad (40)$$

with the degree of polynomial  $X_i^*(z, 0)$  being

$$X_i = L_1 + (N + F)_i \quad (41)$$

where  $(N + F)_i$  is the highest value of the sum  $N + F$  in the polynomials  $X_{ik}^*(z, 0)$ ,  $(k = 1, 2, \dots, \xi)$ .

The number of the control steps is thus

$$n_{ki} = L_1 + (N + F)_i + 1 \quad (42)$$

The highest value of  $n_{ki}$  ( $i = 1, 2, \dots, v$ ), is regarded to be the finite number of the control steps of the whole multi-parameter system.

If disturbances  $u_k(t)$  can be regarded as the linear combination of the function  $t^{m-1}/(m-1)!$ , it follows  $F = m - 1$  and the number of control steps is

$$n_{ki} = L_1 + (N + m)_i \quad (43)$$

It can equally be proved that it is also possible to obtain a zero deviation of the controlled variables  $X_i^*(z, \varepsilon)$  for  $\varepsilon < 0 \div 1 >$  beginning with the instant  $n = n_{ki}$  provided that: disturbance  $u_k(t)$  varied from instant  $n = 0$  according to the function  $u_k(t) = t^{m-1}/(m-1)!$ , the elements of matrices  $[G(p)]$  and  $[G_u(p)]$  have at least one  $m$ -fold zero pole (in hybrid loops the elements of matrices  $[G(p)]$  and  $[G_u(p)]$  must have a holding member at least of the order  $(m - 1)$ , in the equation of conditions (6) and in the equations derived from it  $\Delta_{\Omega B}(z, 0)$  is substituted for  $\Delta_{\Omega B}(z, 0)$  and the auxiliary functions in matrices  $[D_u^*(z, 0)]$  and  $[\bar{Q}_u^*(z, 0)]$  are only polynomials in  $z^{-1}$ .

Then it follows from eqn (16)

$$\begin{aligned} M_{uB, ik}^*(z, 0) &= \\ &= \bar{Q}_{uB, ik}^*(z, 0) (1 - z^{-1})_{kk}^{m-1} C_{kk}^*(z, 0) M_{uA, ik}^*(z, 0) \\ &= \Delta_{\Omega B}(z, 0) D_{uB, ik}^*(z, 0) M_{uA, ik}^*(z, 0) \end{aligned} \quad (44)$$

The degree of eqn (44) is

$$Q_B + m - 1 + C + M_A = l_1 + D + N + M_A \quad (45)$$

$$D_B = D + N \quad (46)$$

$$D = m - 1 + C + M_A \quad (47)$$

$$Q_B = l_1 + N + M_A \quad (48)$$

The number of control steps is then

$$n_{ki} = (Q_B + F)_i \quad (49)$$

$$n_{ki} = l_1 + (M_A + N + C + m)_i - 1 \quad (50)$$

where  $l_1$  is the degree of the polynomial  $\Delta_{\Omega B}(z, 0)$ .

Let it be noted further that in eqn (44)

$$D_{uB, ik}^*(z, 0) = \sum_{v=0}^D d_v z^{-v} \quad (51)$$

differently from the polynomial in eqn (20). Owing to  $d_0 \neq 0$  it has been possible to reduce exponent  $m$  in eqn (44). It should also be mentioned that in this case the number of control steps cannot generally be lowered by the value of  $M_A$  by setting

$$D_{u, ik}^*(z, 0) = \frac{D_{uB, ik}^*(z, 0)}{M_{uA, ik}^*(z, 0)}$$

because the output signal of the digital correction members is

$$E_2^*(z, 0) = -\Delta_{\Omega A}(z, 0) [\omega^*(z, 0)] [D_u^*(z, 0)] U^*(z, 0) \quad (52)$$

and it cannot be assumed that in a general case  $M_{uA, ik}^*(z, 0)$  is contained in  $\Delta_{\Omega A}(z, 0)$ . The denominators of the elements of matrix  $\omega^*[(z, 0)]$  are contained in  $\Delta_{\Omega A}(z, 0)$ .

#### The Optimum Compensation of Disturbance Effects in Wiener's Sense

A method has been shown how to limit the effect of disturbances in quasi-variant control loops by the criterion of the finite number of control steps being considered as the criterion of the quality of control. Another method of solution will be shown where the least square of the deviations of the controlled variables is taken as the criterion of the quality of control. Let the problem be stated by the application of the conventional diagram shown in Figure 2 with the following meaning of denotations:

- $[u(t)]$  = stationary random disturbances
- $[m(t)]$  = parasitic noise
- $[K_u^*(z, \varepsilon)]$  = the  $(v; \xi)$  type rectangular matrix of the transfer functions of the control system that are to be determined
- $[I^*(z, \varepsilon)]$  = the  $(v; \xi)$  type rectangular matrix of the ideal transfer functions of the control system
- $[\Delta(t)]$  = the deviations of the controlled variables  $x[(t)]$  from the ideal output signals  $y[(t)]$ .

$$\Delta_i(n, \varepsilon) = x_i(n, \varepsilon) - y_i(n, \varepsilon) \quad (53)$$

For the sake of brevity the analysis that follows deals only with the case of non-correlated input signals. By using the

results published in an earlier paper by the author, the transfer functions sought for, i.e. the elements of matrix  $K_u^*(z, \epsilon)$ , can be determined by the solution of equation

$$K_{u, ik}^*(j\bar{\omega}, \epsilon) {}^1S_{kk}^*(\bar{\omega}, 0) - \Gamma_{ik}^*(j\bar{\omega}, \epsilon) {}^2S_{kk}^*(\bar{\omega}, 0) = 0 \quad (54)$$

where  $K_{u, ik}^*(j\bar{\omega}, \epsilon)$  and  $\Gamma_{ik}^*(j\bar{\omega}, \epsilon)$  are the  $z$  transforms of the above-mentioned transfer functions,  $z = e^{j\bar{\omega}}$ ,  $\bar{\omega} = \omega T$  while  ${}^1S_{kk}^*(\bar{\omega}, 0)$  and  ${}^2S_{kk}^*(\bar{\omega}, 0)$  are discrete forms of the performance spectral densities:

$${}^1S_{kk}^*(\bar{\omega}, 0) = S_{u_k u_k}^*(\bar{\omega}, 0) + S_{m_k m_k}^*(\bar{\omega}, 0) \quad (55)$$

$${}^2S_{kk}^*(\bar{\omega}, 0) = S_{u_k u_k}^*(\bar{\omega}, 0) \quad (56)$$

Eqn (54) is representing the discrete Laplace transform of the Wiener-Hopf integral equation, the solution of which

$$K_{u, ik}^*(j\bar{\omega}, \epsilon) = \frac{\left[ \frac{\Gamma_{ik}^*(j\bar{\omega}, \epsilon) {}^2S_{kk}^*(\bar{\omega}, 0)}{{}^1S_{kk}^-(\bar{\omega}, 0)} \right]_+}{{}^1S_{kk}^+(\bar{\omega}, 0)} \quad (57)$$

by the known method fulfils the condition

$$k_{u, ik}(n, \epsilon) = 0 \quad \text{for } n < 0 \quad (58)$$

where  $k_{u, ik}(n, \epsilon)$  is the original of the transform  $K_{u, ik}^*(j\bar{\omega}, \epsilon)$ . In eqn (57)

$${}^1S_{kk}^+(\bar{\omega}, 0) {}^1S_{kk}^-(\bar{\omega}, 0) = {}^1S_{kk}^*(\bar{\omega}, 0) \quad (59)$$

where all the poles of  ${}^1S_{kk}^+(\bar{\omega}, 0)$  are inside, and all the poles of  ${}^1S_{kk}^-(\bar{\omega}, 0)$  are outside the zone of stability of plane  $z$ . The + sign in the place of a subscript of the brackets in the numerator of eqn (57) signifies that the function in the brackets has all its poles inside the stability zone of plane  $z$ .

For  $\epsilon = 0$  the elements of matrix  $[K_u^*(z, 0)]$  can be determined from eqn (57), and subsequently the matrix of the transfer functions of the digital correcting members is determined from eqn (3).

Equation (2) represents the unequivocal relationship that exists between transfer functions  $K_{u, ik}^*(z, 0)$  and  $K_{u, ik}^*(z, \epsilon)$ , with the former necessarily fulfilling eqn (6) and also the conditions of stability attached to this equation. Eqns (2) and (6) must equally be complied with by the ideal transfer functions  $\Gamma_{ik}^*(z, 0)$  and  $\Gamma_{ik}^*(z, \epsilon)$ . Consequently transfer functions  $\Gamma_{ik}^*(z, 0)$  cannot be selected arbitrarily. They must fulfil eqn (6) and the conditions of stability that follow from it. If transfer functions  $\Gamma_{ik}^*(z, 0)$  are determined in this way, the stable transfer functions  $\Gamma_{ik}^*(z, \epsilon)$  are obtained unequivocally from eqn (2).

If, in the opposite way, transfer functions  $\Gamma_{ik}^*(z, 0)$  are selected arbitrarily, the required course of the controlling actions can be ensured only at the instants of sampling by the digital correcting members calculated from eqn (3) and with the aid of transfer functions  $K_{u, ik}^*(z, 0)$  determined by eqn (57). However, during the periods between the sampling instants, the course in time of the controlled variables cannot be guaranteed, and it may even be labile.

The new concepts can be summarized in the following theorem.

**Theorem 2**—For the determination of the digital correcting members in Wiener's sense, i.e. in a control loop containing a continuously-acting controlling system and exposed to the effects

of stationary random disturbances the mean square of the deviations of the real output signals from the ideal output signals should attain its minimum value, the command transfer functions  $K_{u, ik}^*(z, 0)$  must fulfil the conditions of stability issuing from the solution of the Wiener-Hopf integral equation, and, in order to ensure the stability of transfer functions  $K_{u, ik}^*(z, \epsilon)$ , the ideal transfer functions  $\Gamma_{ik}^*(z, 0)$  of this solution must comply with the conditions of stability pertaining to eqn (6).

These conditions are necessary and sufficient. This is the fundamental difference between the described control loops and control loops containing only discretely-acting or only continuously-acting members.

The conditions stipulated in Theorem 2 can be fulfilled, if the root factors in the numerators and denominators of transfer functions  $K_{u, ik}^*(z, 0)$  derived from the relations  ${}^1S_{kk}^*(\bar{\omega}, 0)$  and  ${}^2S_{kk}^*(\bar{\omega}, 0)$  are introduced as a condition into auxiliary functions  $\bar{Q}_{u, ik}^*(z, 0)$  and  $D_{u, ik}^*(z, 0)$  calculated from the equation of conditions (6). From this point of view, the solution according to the least square of deviations in Wiener's sense represents only the utilization of a possible application of the required criterion of quality of control within the determinative synthesis theory, and the possibility of extending the auxiliary functions  $\bar{Q}_{u, ik}^*(z, 0)$  and  $D_{u, ik}^*(z, 0)$  by the required number of degrees of freedom.

**Reference**

- 1 STREJC, V. Ensuring reliability in complex automation by automatic digital computers. *Automatizace* 5 (1962) 123-125
- 2 STREJC, V., and RIKA, J. Pt 1: *Izv. Akad. Nauk SSSR, OTN. Energetika i Avtomatika* 5 (1961. 57-71. Pt 2: *Izv. Akad. Nauk SSSR, OTN. Energetika i Avtomatika* (1962). In the press
- 3 STREJC, V. The theory of the synthesis of multi-parameter control systems containing automatic computers and acted upon by random signals. *Sbornik 9. Stroje na zpracovani informaci*. 1963. Prague; NCSAV. In the press

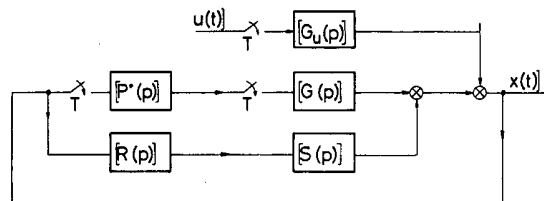


Figure 1

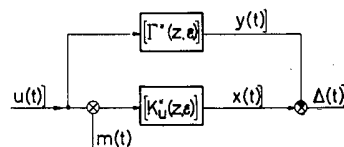


Figure 2



# On Systems with Automatic Control of Configuration

J. BENEŠ

*Ozech*

## Introduction

A further development of the theory of automatic control results from its application to complexes of many elements subject to automatic ordering action. Owing to the impossibility of following the dynamics of the very numerous elements of the complex, statistical characteristics, accessible to macroscopic measurement, may be used for the description of the evolution of the ensemble of the elements. These characteristics are to be compared with the corresponding theoretical ones, derived from the mathematical model of the process of configuration based upon the properties of the elements and upon the conditions of the process, including the influence of the control action. The problem of controlling the development of the ensemble of elements leads to the introduction of a deterministic control of certain frequency functions of events pertaining to the formation of configurations. The corresponding mathematical models use probabilities instead of frequency functions.

## Definition and Basic Scheme

A system with automatic control of configuration is a system with a complex of elements which develops by automatic control towards an assigned state or set of states, characterized by the configuration of these elements.

During this development one (or more) of the following basic operations of configuration, pertaining to the elements of the complex, is realized: the aggregation; the orientation; the liaison; the arrangement; the connection.

The general scheme of a system with automatic control of configuration is shown in *Figure 1*. Here  $K$  = the complex;  $F$  = the formator;  $S$  = the measured state variables of the complex;  $A$  = the acting variables of the formator;  $P$  = the perturbing signals acting upon the complex;  $V$  = the output variables of the complex;  $R$  = the command variables.

The function of the formator is to elaborate the acting signals for the influencing of the configuration of the elements of the complex. The output variables  $V$  of the complex are, in general, different from the measured state variables  $S$  which are chosen so as to inform, by their ensemble, about the configuration of the elements of the complex.

## Description of the State of the Complex

An approach to the description of a complex with a great number of elements consists in its *division into equal zones* and in considering the ensemble of elements contained in each of these zones. The interrelation of the ensembles contained in the zones, especially in the neighbouring ones, may be of interest. In two-dimensional representation we draw the meshwork of zones as, for example, in *Figures 7* and *8*. As it may not be possible to measure the state variable in all these zones, we

introduce *sample zones*. The measured state of the complex may be expressed by the measured state variables  $S$  in the form: (a) of a column vector with elements  $s_1(t), s_2(t), \dots, s_n(t)$ ; (b) or, in the case of a two-dimensional arrangement of measuring points, which may be advantageous for the expression of the configuration, in the form of a quadratic matrix  $\|s_{ij}(t)\|$  ( $i, j = 1, 2, \dots, n$ ); (c) or, in the case of a three-dimensional arrangement of measuring points, in the form of a cubic matrix

$$\|s_{ijk}(t)\| \quad (i, j, k = 1, 2, \dots, n) \quad (1)$$

Measurements of state variables in three-dimensional or in two-dimensional arrays of zones of the complex can be reduced by scanning to a sequence of measurements. Similarly, one can express the required state of the complex using the command variables  $R$ .

A theoretical measure of the state of the complex with many elements is the configurational redundancy

$$R_{fm} = 1 - \frac{\Delta S_{fm}}{\Delta S_{fv}} \quad (2)$$

where  $\Delta S_{fm} < \Delta S_{fv}$ .

The index  $m$  applies to the intermediate state between the initial state (index  $v$ ) and the final state (index  $c$ ) and where the differences of configurational entropy are

$$\Delta S_{fm} = k \log(Z_m - Z_c) \quad (3)$$

$$\Delta S_{fv} = k \log(Z_v - Z_c) \quad (4)$$

where  $k$  is a scale factor.  $Z_v$  is the number of possible ways of having the elements ordered at the initial state,  $Z_m$  is the number of different ways of ordering of the elements suiting the definition of the intermediate state at a certain phase of development and  $Z_c$  is the number of different ways of ordering of the elements suiting the requirements upon the final state.

The configurational entropy is a concept used in statistical physics. In crystallography, the entropy change for a transition in the crystalline phase is divided into: the change of the configurational entropy and the change of thermal entropy<sup>8</sup>. The configurational entropy of the arrangement of atoms in a lattice is determined by the number of different ways in which the atoms may be arranged over the available number of lattice sites. In chemistry, the information content  $I_t$  of a protein is divided into:  $I_s$  depending upon the amino acid sequence, and  $I_c$  depending upon the configuration of the polypeptide chain in the native molecule<sup>9</sup>.

The state of the complex may be described by different concepts and measures depending upon the basic operation of configuration of the elements<sup>6, 12</sup>.

121/2

*Quantitative Expression of State During Aggregation*

The simplest expression of the state is in terms of the number of elements or of their concentration in the different zones. The information connected with the concentration into a single zone  $s$  of elements of a certain type  $i$ , which previously have been distributed over the whole complex, is

$$I_{ci} = \log_2 \left( \frac{c_{is}}{c_{ik}} \right) \quad [\text{bit}] \quad (5)$$

where  $c_{is}$  is the concentration of the elements of type  $i$  in the zone  $s$ ,

$c_{ik}$  is the concentration of the elements of type  $i$  in the whole complex.

When several types  $i$  ( $i = 1, 2, \dots$ ) of elements are involved, it is

$$I_c = \sum_i \log_2 \left( \frac{c_{is}}{c_{ik}} \right) \quad [\text{bit}] \quad (6)$$

*Quantitative Expression of State During Orientation*

The orientation of the elements of a zone of the complex may be expressed by angular measure. The information connected with the orientation of the  $i$ th element of the complex may be expressed by its orientation information  $I_{or [i]}$ . If the configuration of the complex requires that the orientation of the  $i$ th element be fixed within  $\Delta\Theta$ ,  $\Delta\phi$ ,  $\Delta\psi$ , where  $\Theta$ ,  $\phi$ ,  $\psi$  are Euler angles, it is

$$I_{or [i]} = \log_2 \left\{ \frac{8\pi^3}{\Delta\Theta_i \Delta\phi_i \Delta\psi_i} \right\} \quad [\text{bit}] \quad (7)$$

*Quantitative Expression of State During Liaison*

Consider the operation of liaison of the elements in a complex of constant volume, with elements of different types  $i$ , where  $i = 1, 2, \dots, k$  and denote  $n_1, n_2, \dots, n_k$  the numbers of these elements. They move and combine at random to form new types of elements by liaison. To characterize the development of the state of the complex use the probability  $P(\mathbf{n}, t)$ , which is the probability, that in time  $t$  the complex has the composition  $\mathbf{n}$ , where  $\mathbf{n}$  is a vector, whose components are the numbers of the elements of the different types. By the action of the formator one wishes to influence the probability  $P(\mathbf{n}, t)$ .

*Quantitative Expression of State During Arrangement*

The number of correctly occupied sites or lattice sites by the elements of the complex is a simple measure of arrangement. The information connected with the position of the  $i$ th element in the complex of volume  $V$  may be expressed by its placement information  $I_p [i]$ . If the configuration of the complex requires that the  $i$ th element remain within a space  $\Delta x_i$ ,  $\Delta y_i$ ,  $\Delta z_i$ , it is

$$I_p [i] = \log_2 \left\{ \frac{V}{\Delta x_i \Delta y_i \Delta z_i} \right\} \quad [\text{bit}] \quad (8)$$

*Quantitative Expression of State During Connection*

As a characteristic quantity of a random net, Clark and Farley have used the connectivity. An element,  $i$ , is connected

to an element  $j$  with a probability  $P_{ij}$  which may depend upon both  $i$  and  $j$  and on other characteristic quantities of the net as a whole. Uttley has considered the probability of the connection of an element in a given position to an input point. This probability is a function of the position. Beurle has used a probability of connection of an element to all elements which are in a distance  $r$  of it, this probability being a function of the co-ordinates  $x, y, z$  of the element and of the distance  $r$ .

As we are interested in characterizing the state of a developing random net, the use of test impulses, applied in points of sample zones and the measurement of the number of elements activated at a given distance in a given direction, or of the speed of the signal spreading, can be suggested, according to the particular case. These quantities are related to the statistical characteristics of the random net.

*The Interaction of the Formator and of the Complex*

The acting variables of the formator are from the point of view of automatic control the output variables of a multi-dimensional controller with many inputs, some of which are the measured state variables of the complex. The information about the behaviour of the complex only from the measurement of external input and output variables of the complex would be insufficient. This is also in compliance with the principle of uncertainty in the structural behaviour of multivariable systems, formulated by Mesarović<sup>1</sup>.

Therefore direct measurement of the state variables, amended by theoretical relations yielded from the mathematical model of the configuration process, is required.

A methodical approach towards the identification of the process occurring in the complex consists in the following stages:

(1) The formation of a mathematical model of the process of configuration of the elements.

(2) Computation of the relevant mean values  $S^{(T)}(t)$  on a mathematical machine on the basis of the mathematical model and using complementary information about the physical conditions of the process.

(3) Comparison of the computed mean values  $S^{(T)}$  and their development in time with the measured state variables  $S$ .

(4) The appropriate adaptive correction of the model sub 2 in order to minimize the difference of the comparison sub 3.

The objective is to obtain an approximate model of the process in the complex based on theoretical results about the statistical dynamics of the elements of the complex. Such a model is intended to bring a better insight into the mechanism of configurational changes and to help in the choosing the method of control.

Two basic deviations which have to be considered in the theory of these systems are:

$$(1) \quad \varphi(t) = \mathbf{R}(t) - \mathbf{S}(t) \quad (9)$$

where  $\mathbf{R}(t)$  are the command variables and  $\mathbf{S}(t)$  the measured state variables,

$$(2) \quad \varphi^{(T)}(t) = \mathbf{S}^{(T)}(t) - \mathbf{S}(t) \quad (10)$$

where  $\mathbf{S}^{(T)}(t)$  are theoretical state variables computed from the mathematical model. The minimization of  $\varphi^{(T)}(t)$  is a form of the identification problem of the process in the complex.

When using the cubic matrix expressions the basic deviations are:

$$(a) \quad \|\phi_{ijk}(t)\| = \|r_{ijk}(t)\| - \|s_{ijk}(t)\| \quad (11)$$

$$(b) \quad \|\phi_{ijk}^{(T)}(t)\| = \|s_{ijk}^{(T)}(t)\| - \|s_{ijk}(t)\| \quad (i, j, k = 1, 2, \dots, n) \quad (12)$$

Consider a set of discrete simultaneous values of the time functions. A typical operation with the cubic matrix consists in applying to the trilinear form

$$F = \sum_{i, j, k=1}^n \phi_{ijk} x_i y_j z_k \quad (13)$$

the linear transformation

$$x_\gamma = \sum_{i=1}^n g_{\gamma i} X_i \quad (\gamma = 1, 2, \dots, n) \quad (14)$$

with the quadratic matrix

$$g = \| \|g_{\gamma i}\| \quad (\gamma, i = 1, 2, \dots, n) \quad (15)$$

There is then obtained the product cubic matrix in the index  $i$ :

$$\phi \{i\} g = \| \| \sum_{\gamma=1}^n \phi_{\gamma j k} g_{\gamma i} \| \| \quad (16)$$

Similarly,

$$\phi \{j\} g = \| \| \sum_{\gamma=1}^n \phi_{i \gamma k} g_{\gamma j} \| \| \quad (17)$$

$$\phi \{k\} g = \| \| \sum_{\gamma=1}^n \phi_{i j \gamma} g_{\gamma k} \| \| \quad (18)$$

Applying certain bilinear transformations with a cubic matrix of  $n$ th order to the trilinear form we would get product tetric matrices of the  $n$ th order, which are already more complex to handle<sup>3</sup>.

Owing to the number of input variables of the formator and to the complexity of its function, the formator should be a digital computer.

A special role in the introduction of this point of view of the influencing of probabilities of events by the formator is to be assigned to the concept of controlled probabilistic transducer. As an example of a type of controlled probabilistic transducer let us derive from the transducer type, mentioned by Kochen<sup>11</sup>, the transducer of Figure 2, where the conditional probabilities  $c_1$  and  $c_2$  are, respectively, functions of the acting variables  $a_1, a_2$  of the formator. The conditional probabilities are

$$c_1 = P[y(t) = 1 | x(t-1) = 1] \quad (19)$$

$$c_2 = P[y(t) = 1 | x(t-1) = 0] \quad (20)$$

where  $x(t)$  is a two-valued variable (value 1 or zero), and  $t$  denotes the number of the time interval.

Another type of controlled probabilistic transducer (Figure 3) can be derived from the concept of binomial probabilistic transform, studied by Sugimori<sup>4</sup>, by making the parameters  $p$  and  $\delta$  functions of the acting variables  $a_1$  and  $a_2$  of the formator. Let  $x, y$  be two-valued variables (value 1 or zero) and  $\delta$  a time lag. For the special case, that  $p$  and  $\delta$  are constants, and when  $x$  has a Poisson distribution with parameter  $\lambda$

$$p(x) = \frac{\lambda^x}{x!} e^{-\lambda} \quad (21)$$

then in time  $\delta$  and only in this time

$$p(y) = \sum_{x=0}^{\infty} p(x) b(y; x, p) = \frac{(\lambda p)^y}{y!} e^{-\lambda p} \quad (22)$$

where  $b(y; x, p)$  is the binomial distribution of the random variable  $S_x$  denoting the number of successes in  $x$  Bernoulli trials.

Further, a simple scheme of a controlled probabilistic transducer without time-lag for a sequence of equidistant pulses can be represented by an element for logical product, steered by a generator of random pulses with controlled statistical parameter, under the assumption of the synchronization of the pulse sequences.

### The Optimization of the Process of Configuration

The attainment of a certain state of configuration (or set of states) can be considered as a result of two opposite actions: the one of configurational ordering, the other of disordering.

An important problem in connection with the systems with automatic control of configuration can be raised: the stability of the controlled complex in remaining in a definite set of states. For the description of the development of some complexes with many elements, mathematical models based on Markov processes seem suitable.

The number of states which the complex with many elements can take is infinite and countable. The theory of the Markov processes with an infinite and countable number of states is to be applied.

In order to formulate certain basic relations concerning the optimization of the process of automatic configuration we shall willingly limit ourselves to consider it as a Markov process with a finite number  $N$  of states. Let this process be defined by the matrix of transition probabilities  $P = [p_{jk}]$  and by the matrix of rewards  $W = [w_{jk}]$ , where the indexes  $j$  and  $k$  apply to the transition from the state  $j$  to the state  $k$ .

The result  $w_{jk}$  is the increment of configurational ordering associated with this transition. It is to be expressed in the units of configurational measure. A physical interpretation may be, for example, as the increase of the number of a new type of element formed by the liaison of two types of elements, or of the number of correctly occupied sites in a lattice *a. s. o.* During the development of the complex some  $w_{jk}$  can be negative.

Applying the method of Howard<sup>2</sup>, one expresses the mean reward from a transition

$$g = \sum_{j=1}^N \pi_j q_j \quad (23)$$

where  $\pi_j$  = the probability of the complex being in state  $j$  after a large number of steps,

$q_j$  = the immediate reward expected at state  $j$ , i.e. the expected reward connected with the transition of the complex from the state  $j$  to the next one.

The immediate result  $q_j$  is

$$q_j = \sum_{k=1}^N p_{jk} w_{jk} \quad (j = 1, 2, \dots, N) \quad (24)$$

## 121/4

The aim is to maximize the mean reward  $g$  of the Markov process of configuration.

Supposing that at each state  $j$  one of the alternatives of the action of the formator upon the complex can be chosen. To each alternative  $a$  corresponds a transition probability  $p_{jk}^a$  and a reward  $w_{jk}^a$ . When chosen, the alternative becomes a decision  $d$ . One denotes by  $d_j(n)$  the decision taken at the state  $j$ , which means  $n$  states before the attainment of the final state. A column vector  $d$ , with elements  $d_j(n)$  expresses the chosen policy. The total expected reward during the development of the complex in  $n$  steps starting from the state  $j$  and applying a specific policy is  $v_j(n)$ . Under the assumption of the Markov process being completely ergodic, it is for large  $n$

$$v_j(n) = ng + v_j \quad (j=1, 2, \dots, N) \quad (25)$$

Between the introduced quantities there is the relation

$$g + v_j = q_j + \sum_{k=1}^N p_{jk} v_k \quad (j=1, 2, \dots, N) \quad (26)$$

Following further the method of Howard, let  $v_N = 0$  and call  $v_j$  the relative values of the policy. By a judicious choice of the  $p_{jk}$  and  $q_j$  for each state  $j$  the reward  $g$  is to be maximized.

(1) For each state  $j$  the alternative  $a'$  which maximalizes the value [see the relation (24)]

$$q_j^a + \sum_{k=1}^N p_{jk}^a v_k \quad (27)$$

is to be determined. Here the index  $a$  denotes the values belonging to the alternative  $a$ . Then by putting  $p_{jk}^a = p_{jk}$ ;  $q_j^a = q_j$ , the resulting values are used below.

(2) Using these  $p_{jk}$  and  $q_j$  in the system of linear simultaneous equations (26) and by solving this system one gets the  $v_j$  and the  $g$  which will be again used in (1). By the iterative computation process involving (1) and (2) one finally gets the  $g$ , and the  $p_{jk}$  and  $q_j$ . The speed of computation would be too high. On the other hand, assuming that the complex evolves relatively slowly, the change of the  $p_{jk}$  would be done by the action variables of the formator. Without the action of the formator, the isolated complex would develop 'spontaneously' with transition probabilities  $p_{jk}^{(s)}$ .

Theoretical investigations require the application of the theory of Markov processes with an infinite and countable number of states. Some notions are common with the theory of Markov processes with a finite number of states, as for example, the notion of undecomposable groups, of transition groups and of final groups.

### Suggested Examples of Systems

As an example of a system with automatic control of configuration, consider a system with controlled operation of liaison of three different types  $A$ ,  $B$ ,  $C$  of the very numerous elements of its complex. The elements move with random Brownian movement and have the following properties: (i) when  $A$  and  $B$  collide, a new element  $D$  results, (ii) when  $D$  and  $C$  collide, a new element  $E$  results, (iii) the collisions of the elements are at random, (iv) the liaisons are irreversible, (v) the direct liaison of  $C$  with  $A$  or  $B$  is impossible, (vi) the liaison rate parameters are  $k_1$  and  $k_2$  (Figure 4).

The state variables of the complex at time  $t = 0$  are  $n_{a0}$ ,  $n_{b0}$ ,  $n_{c0}$ ,  $0$ ,  $0$ —the numbers of elements of the different types. The command variables  $r_1$ ,  $r_2$  are the required numbers of the elements  $D$  and  $E$  respectively (Figure 4). The acting variables of the formator are  $a_1$ ,  $a_2$ , influencing the liaison rate parameters  $k_1$  and  $k_2$  respectively. Let  $x_1$ ,  $x_2$  be the number of elements  $D$  and  $E$  respectively at time  $t$ .

The control of the operation of liaison is based on making the liaison rate parameters appropriate functions of time:  $k_1(t)$ ,  $k_2(t)$ .

In order to simplify the expressions, first consider  $k_1$  and  $k_2$  as constants in the mathematical model. The differential equation describing the development of the complex is then:

$$\begin{aligned} \frac{dP(x_1, x_2, t)}{dt} = & k_1(n_{a0} - x_1 + 1)(n_{b0} - x_1 + 1)P(x_1 - 1, x_2, t) \\ & - k_1(n_{a0} - x_1)(n_{b0} - x_1)P(x_1, x_2, t) \\ & + k_2(x_1 + 1)(n_{c0} - x_2 + 1)P(x_1 + 1, x_2 - 1, t) \\ & - k_2(x_1)(n_{c0} - x_2)P(x_1, x_2, t) \end{aligned} \quad (28)$$

Using the method of the generating function one takes

$$F(s_1, s_2; t) = \sum_{x_1, x_2=0}^{\infty} P(x_1, x_2, t) s_1^{x_1} s_2^{x_2} \quad (29)$$

and finds

$$\begin{aligned} \frac{\partial F}{\partial t} = & F[-k_1 n_{a0} n_{b0} (1 - s_1)] \\ & + \frac{\partial F}{\partial s_1} \{k_1 s_1 (1 - s_1)(n_{a0} + n_{b0} - 1) \\ & + k_2 [(s_1 - s_2) - n_{c0}(s_1 - s_2)]\} \end{aligned} \quad (30)$$

The boundary conditions are

$$F(0, 0; 0) = 1 \quad \text{and} \quad F(1, 1; t) = 1$$

The mean values of the numbers  $n_d$  and  $n_e$  of the elements  $D$  and  $E$  respectively are then as functions of time

$$m_{n_d}(t) = \left[ \frac{\partial}{\partial s_1} \log F(s_1, s_2; t) \right]_{s_1=s_2=1} \quad (31)$$

$$m_{n_e}(t) = \left[ \frac{\partial}{\partial s_2} \log F(s_1, s_2; t) \right]_{s_1=s_2=1} \quad (32)$$

The dispersions are

$$\sigma_{n_d}^2(t) = \left[ \frac{\partial^2}{\partial s_1^2} \log F(s_1, s_2; t) + \frac{\partial}{\partial s_1} \log F(s_1, s_2; t) \right]_{s_1=s_2=1} \quad (33)$$

$$\sigma_{n_e}^2(t) = \left[ \frac{\partial^2}{\partial s_2^2} \log F(s_1, s_2; t) + \frac{\partial}{\partial s_2} \log F(s_1, s_2; t) \right]_{s_1=s_2=1} \quad (34)$$

Thus the model of the complex may be represented as a black box with 2 inputs:  $k_1$  and  $k_2$ , with an initial state, characterized by the vector components  $n_{a0}$ ,  $n_{b0}$ ,  $n_{c0}$ ,  $0$ ,  $0$  at time  $t = 0$ , and with 2 outputs:  $m_{n_d}(t)$ ,  $m_{n_e}(t)$  or  $\sigma_{n_d}^2(t)$ ,  $\sigma_{n_e}^2(t)$ .

More generally, if  $k_1$  and  $k_2$  are not constants, set in advance, but change in time under the control action, the corresponding Markov process is non-homogeneous.

On *Figure 6* there is a closed oriented chain of reservoirs  $B_1, B_2, B_3$ , containing many elements of the same type. This chain forms the complex. The acting variables of the formator  $a_1, a_2, a_3$  control the probability transducers represented by full points. The probabilities of the random transitions of elements from one reservoir to another are thus controlled. The aim of the function of the system is to reach a repartition of the elements over the reservoirs, prescribed by the command variables  $r_1, r_2, r_3$ . The total number of elements is  $N$ . The number of elements contained in the reservoir  $B_2$  at time  $t$  is  $x_2$ . The probability of the transition  $x_2 \rightarrow x_2 + 1$  in the time interval  $(t, t + \Delta t)$  is

$$\lambda_{x_2} \Delta t + \sigma(\Delta t)$$

where, at first, let the  $\lambda$  in the relation

$$\lambda_{x_2} = \lambda_{x_2} \quad (35)$$

be constant in time.

If at time  $t$  the reservoir  $B_2$  contains  $x_2$  elements, the probability of the transition  $x_2 \rightarrow x_2 - 1$  in the time interval  $(t, t + \Delta t)$  is

$$\mu_{x_2} \Delta t + \sigma(\Delta t)$$

where, at first, let the  $\mu$  in the relation

$$\mu_{x_2} = \mu \cdot x_2 \quad (36)$$

be constant in time.

The probability of the transition to a number of elements other than  $x_2 + 1$  or  $x_2 - 1$  is  $\sigma(\Delta t)$ .

The probability of no change in the time interval  $(t, t + \Delta t)$  is

$$1 - (\lambda_{x_2} + \mu_{x_2}) \Delta t + \sigma(\Delta t)$$

By making, in addition, similar assumptions for the reservoirs  $B_3$  and  $B_1$ , one gets the system of differential equations:

$$\frac{dP_{x_i}(t)}{dt} = \lambda_{x_i-1} P_{x_i-1}(t) - (\lambda_{x_i} + \mu_{x_i}) P_{x_i}(t) + \mu_{x_i+1} P_{x_i+1}(t),$$

$$\text{for } x_i = 1, 2, \dots \text{ and } i = 1, 2, 3 \quad (37)$$

When the parameters  $\lambda$  and  $\mu$  in the relations as (35) and (36) change in time, as under the control action of the formator, one has to deal with Markov processes of the birth and death type non-homogeneous in time, describing the development of each of the reservoirs.

Another example of system may be suggested with complexes whose schematic representation is in the form of a two-dimensional array of zones, which may have, for example, rectangular (*Figure 7*) or triangular (*Figure 8*) form. The zones contain many elements. The transition of the elements from one zone to other zones is controlled by probabilistic transducers steered by the acting variables of the formator and represented as full dots. On *Figure 7* the state of the selected zone  $R_{22}$  is a function of the states of the neighbouring zones. Considering at first a Markov process homogeneous in time as a model of the development of the zone  $R_{22}$ , which can be represented as a rectangle with two inputs and two outputs, one makes the following assumptions:

The number of elements in the zone  $R_{22}$  at time  $t$  is  $x_{22}$ .

The probability of the transition  $x_{22} \rightarrow x_{22} + 1$  in the time interval  $(t, t + \Delta t)$  is

$$\lambda_{x_{22s}} \Delta t + \lambda_{x_{22z}} \Delta t + \sigma(\Delta t)$$

The probability of the transition  $x_{22} \rightarrow x_{22} - 1$  in the time interval  $(t, t + \Delta t)$ , if at time  $t$  the zone is in state  $x_{22}$  ( $x_{22} = 1, 2, \dots$ ), is

$$\mu_{x_{22j}} \Delta t + \mu_{x_{22v}} \Delta t + \sigma(\Delta t)$$

The probability of the transition to a state other than  $x_{22} + 1$  or  $x_{22} - 1$  is  $\sigma(t)$ .

The probability of no change of state is

$$1 - (\lambda_{x_{22s}} + \lambda_{x_{22z}} + \mu_{x_{22j}} + \mu_{x_{22v}}) \Delta t + \sigma(\Delta t)$$

The corresponding Markov process pertaining to the zone  $R_{22}$  is of the birth and death type.

Because of the interrelation of the zones there is an interdependence between the parameters  $\lambda$  and  $\mu$  relative to neighbouring zones. It may be e.g.

$$\left. \begin{aligned} \lambda_{x_{22s}} &= \mu_{x_{12j}} \cdot x_{12}; & \mu_{x_{22j}} &= \mu_{x_{22j}} \cdot x_{22} \\ \lambda_{x_{22z}} &= \mu_{x_{21v}} \cdot x_{21}; & \mu_{x_{22v}} &= \mu_{x_{22v}} \cdot x_{22} \end{aligned} \right\} \quad (38)$$

The quantities  $\mu_{ikj}$  and  $\mu_{ikv}$ , where  $i, k = 1, 2, \dots, n$ , may be arranged into a quadratic matrix. Owing to the action of the probabilistic transducers, the  $\mu_{ikj}$  and  $\mu_{ikv}$  change in time.

#### Perspectives of Development of the Theory and of its Applications

There is a large field for the development of the theory of systems with automatic control of configuration. The methods and results of the statistical mechanics form the basis for the dynamics of complexes with many elements. The representativeness of mathematical models is to be checked against physical measurements. The solution of problems, related to the Markov processes of configurational development non-homogeneous in time, could be aided by modelling the process on simulators using random signal generators.

The field of application of the theory may be seen, for example, in these directions: the influencing of the formation of strips of molecules; the automatic control of the cultivation of microorganisms, such as algae; the formation of random nets.

*Fruitful suggestions from Professor Robert Fortet of Paris and from Professor Jaroslav Kožesník of Prague are gratefully acknowledged.*

#### References

- MESAROVIC, M. D. *The Control of Multivariable Systems*. 1960. New York; The Technology Press of MIT and Wiley,
- HOWARD, R. *Dynamic Programming and Markov Processes*. 1960. New York; The Technology Press of MIT and Wiley,
- СОЛОЛОВ, Н. Т. *Пространственные матрицы и их приложения*. Физматгиз, Москов, 1960
- SUGIMORI, MAKOTO. Binomial probabilistic sequential circuit. *Rev. Electr. Commun. Lab., Tokyo*, 9, Nos. 9-10 (1961) 627-654
- KOZESNIK, J. A simple stochastic model of continuous cultivation of microorganisms in several basins. *Second Internat. Symp. Continuous Cultivation of Microorganisms*, Prague, 18-23 (1962) 6
- JACOBSON, H. The Informational content of mechanisms and circuits. *Inform. Control*, 2 (1959) 285-296
- BHARUCHA-REID, A. T. *Elements of the Theory of Markov Processes and their Applications*. 1960. New York; McGraw-Hill
- MØLLER, CHR. KN. Electrochemical investigation fo the transition from tetragonal to cubic caesium plumbo chloride. *Mat. Fys. Medd. Dan. Vid. Selsk., Copenhagen*, 32, No. 15 (1960)

<sup>9</sup> AUGENSTINE, L. G. Protein structure and information content, *Symposium on Information Theory in Biology*, pp. 102-123. 1958. London; Pergamon Press  
<sup>10</sup> FOERSTER, H. V. On self-organizing systems and their environments. *Self-organizing Systems*. 1960. New York. Pergamon Press

<sup>11</sup> KOCHEN, M. Circle networks of probabilistic transducers. *Inform. Control*. 2 (1959) 168-182  
<sup>12</sup> SINGER, K. Application of the theory of stochastic processes to the study of irreproducible chemical reactions and nucleation processes. *J. roy. Statist. Soc.*, ser. B. 15 (1953) 92-106

Summary

The definition, the scheme and five basic operations of systems with automatic control of configuration are given. For the description of the complex with many elements its division into zones and statistical characteristics as state variables are introduced. The configurational redundancy and simpler measurable quantities are measures of ordering. The probabilities involved in the configuration are influenced by the formator through controlled probabilistic transducers. The methodical approach to the solution of the formator and complex interaction is in the formation of a mathematical model, checked against physical measurements and in using it in the choice of the

control algorithm. The two basic deviations and cubic matrices of variables of this type of multivariate systems are introduced. The optimization of the process of configuration is described using the method of Howard in terms of a Markov process with decisions and rewards. Three examples of systems are suggested, pertaining to formation of strips of elements and to the migration of elements in one-dimensional and two-dimensional arrays of zones. Further progress is expected from physical modelling of non-homogeneous Markov processes. The field of application of the theory can be seen in chemistry, in automatic cultivation of algae and in the random net formation.

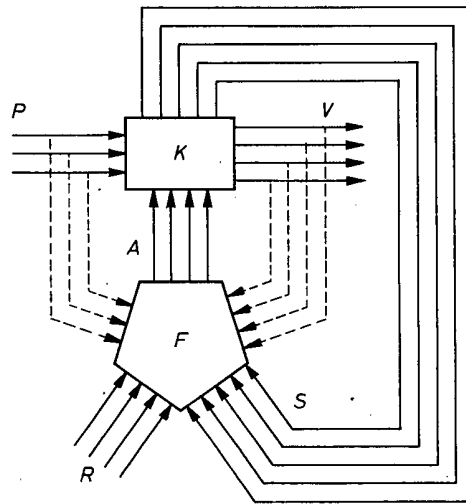


Figure 1

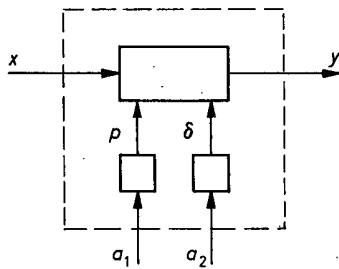


Figure 2

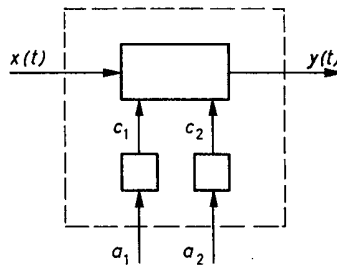


Figure 3

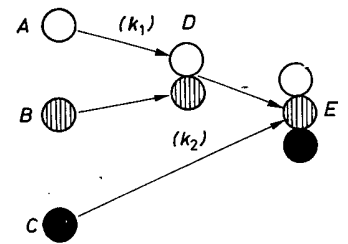


Figure 4

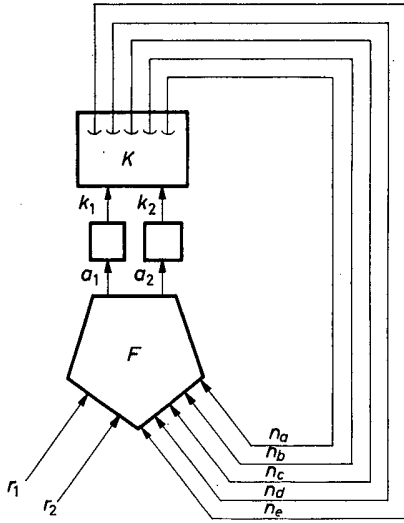


Figure 5

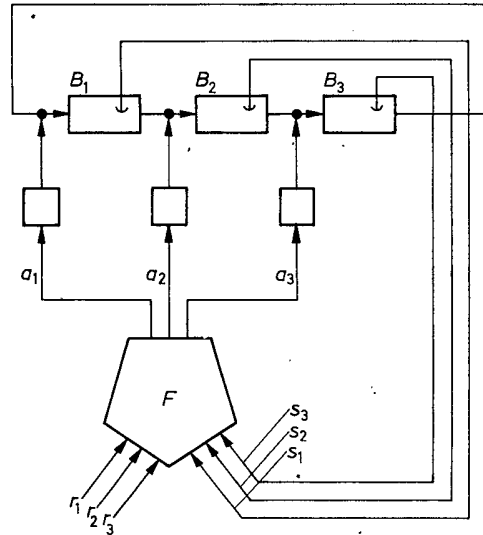


Figure 6

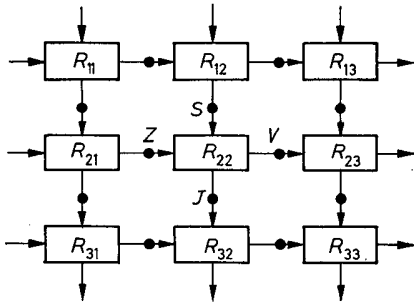


Figure 7

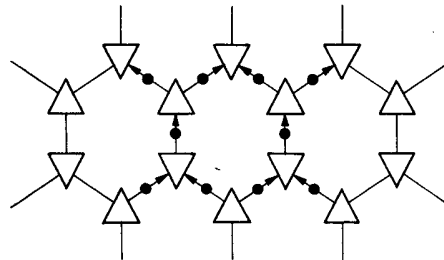


Figure 8

# Combination of Finite Settling Time and Minimum Integral of Squared Error in Digital Control Systems

V. PETERKA *Brno*

## Introduction

The aim frequently followed in the design of digital control systems is to eliminate the system error, caused by an input signal of a typical form (step, ramp, constant acceleration), within a minimum time<sup>1-3</sup>. Cases are often encountered in such systems with the fastest response where the transient error has a very short duration, but an inadmissible magnitude. This shortcoming can be removed by extending the response by one, two, or more sampling periods, as required, and by the application of a further criterion suppressing the system errors<sup>4</sup>. This article deals with a simple numerical method of digital controller design where the criterion of finite settling time has been combined with the minimum integral of squared error.

## The Statement of the Problem

Consider a control system compensated by a digital controller according to *Figure 1*. It is assumed that the transfer function of the plant is a rational fraction

$$S(p) = \frac{\sum_{v=0}^m b_v p^v}{\sum_{v=0}^n a_v p^v} = \frac{b(p)}{a(p)} = \frac{b(p)}{a_n \prod_{v=1}^n (p - p_v)} \quad (1)$$

with all its poles  $p_v$  in the left-hand side semi-plane  $p$ , or with maximum one pole equalling zero. For simplicity, let the problem be confined to an input signal having the form of a unit step

$$W(z) = \frac{1}{1 - z^{-1}}$$

and the holding device being of zero order

$$H(p) = \frac{1 - e^{-Tp}}{p}$$

Cases with another type of input signal, and with a holding device of a higher order, can be investigated in a similar way. It will also be assumed that the time required for the computing operation can be neglected, and the system has no dead time; however, it is possible to show that the consideration of both these lags is possible without any fundamental difficulties.

Let the pulse-transfer function of the continuously acting member be denoted

$$G(z) = \frac{B_0 + B_1 z^{-1} + \dots + B_n z^{-n}}{A_0 + A_1 z^{-1} + \dots + A_n z^{-n}} = \frac{B(z)}{A(z)} \quad (2)$$

The conditions of a finite settling time have been discussed in detail<sup>1-4</sup> and here they are stated only briefly in a form suited for the case.

For attaining a zero steady-state error at the sampling instants, after a finite number of sampling periods and under the conditions stated above, it is necessary that the overall pulse-transfer function

$$F(z) = \frac{P(z)G(z)}{1 + P(z)G(z)} \quad (3)$$

should be a polynomial in  $z^{-1}$ , and

$$F(1) = 1 \quad (4)$$

If the intersampling ripples also are to be eliminated, it is necessary to attain the settling of the manipulated variable  $y(t)$ . This will happen, if the pulse-transfer function of  $E_2(z)/W(z)$  is also a polynomial in  $z^{-1}$ . The equation for this transfer function can be modified by the relations

$$E_2(z) = \frac{X(z)}{G(z)}, \quad F(z) = \frac{X(z)}{W(z)}$$

into the form

$$\frac{E_2(z)}{W(z)} = \frac{F(z)}{G(z)} = \frac{F(z)A(z)}{B(z)} \quad (5)$$

It follows from eqn (5) that all conditions stated will be fulfilled, if the overall pulse-transfer function has the form

$$F(z) = \frac{1}{B(1)} B(z) D(z) \quad (6)$$

where

$$D(z) = D_0 + D_1 z^{-1} + \dots + D_L z^{-L} \quad (7)$$

is a selectable polynomial for which

$$D(1) = \sum_{i=0}^L D_i = 1 \quad (8)$$

From relations (3) and (6) the necessary pulse-transfer function of the digital computer follows

$$P(z) = \frac{D(z)A(z)}{B(1) - D(z)B(z)} \quad (9)$$

If  $D(z) = 1$  is selected the system will have the fastest response, nevertheless the transient error can reach an inadmissible magnitude as shown in the example that follows. Therefore let polynomial  $D(z)$  be of the general order of  $L$ , state the problem



122/2

as the determination of the coefficients  $D_0, D_1, \dots, D_L$  of the polynomial with the integral of the squared error

$$J = \int_0^\infty q(t) e_1^2(t) dt \quad (10)$$

having a minimum value.

Now it remains to select a suitable weighting function  $q(t)$ , so the following cases will be investigated.

In the first the same importance is allotted to all errors during the control process and  $q(t) = 1$  is selected, Figure 2(a). However, this selection need not be necessarily the most advantageous, namely the largest share in integral (10) belongs to errors at the beginning of the control process that cannot be physically eliminated in plants with a step function response starting from the origin. The minimalization of integral (10) can produce rather large overshoots that are not always desirable.

For this reason it is necessary to investigate the second case where no errors in the first sampling period are contained in integral (10), i.e. the weighting function is selected in the form of a unit-step function in time  $T$ ,  $q(t) = 1(t - T)$  Figure 2(b). The physical meaning of this condition is the requirement of the computer liquidating the error, as far as possible, during one step, and not instantaneously as demanded in the former case. That is to say, in this second case the requirement put forward is less severe, and technically easier to realize.

The method of calculation is arranged in such a way that both cases can be investigated simultaneously, and thus it is possible to reach a decision in favour of the case that is more beneficial at the given concrete application.

**The Survey of Results**

Coefficients  $D_1, D_2, \dots, D_L$  of the selectable polynomial (7), fulfilling the condition of the minimum integral (10) of the squared error, can be found by the solution of the system of linear equations

$$[K_{rs}][D_s] = -[R_{r0}] \quad (11)$$

where the square matrix  $[K_{rs}]$  is symmetrical the elements of which, and also the elements of column matrix  $[R_{r0}]$ , are independent of the selected degree  $L$  of polynomial  $D(z)$ . Two different cases are considered in the calculation of the elements of matrices  $[K_{rs}]$  and  $[R_{r0}]$ : (a) the transfer function  $S(p)$  of the plant has no zero pole, and (b)  $S(p)$  has one zero pole.

*Case (a)*

In the case of the transfer function  $S(p)$  having all its poles different from zero, the step function response of the system is given by the equation

$$s(t) = \mathcal{L}^{-1} \left\{ \frac{S(p)}{p} \right\} = C_0 + \sum_{v=1}^n C_v e^{p_v t} \quad (12)$$

In this, and in all other equations that follow, the assumption is made that all the poles differ from each other. The case of multiple poles can be introduced by means of limits. The elements of matrices  $[K_{rs}]$  and  $[R_{r0}]$  are calculated by the following procedure. First of all the following expressions are solved numerically

$$\theta(k) = \sum_{v=1}^n \rho_v z_v^k, \quad k=0, 1, 2, \dots, n+L \quad (13)$$

where

$$z_v = e^{p_v T}$$

$$\rho_v = C_v \left( \frac{C_0}{P_v} + \delta_v \right), \quad \delta_v = \sum_{\mu=1}^n \frac{C_\mu}{P_v + P_\mu} \quad (14)$$

The calculation is made for  $k = 0, 1, 2, \dots, n + L$ , where  $L$  is the selected degree of polynomial  $D(z)$ .

The procedure is continued in such a way that all elements of the same row of matrix  $[K_{rs}]$ , and also of matrix  $[R_{r0}]$ , are calculated simultaneously for the weighting function  $q(t) = 1$ , and also for  $q(t) = 1(t - T)$ . As  $[K_{rs}]$  is a symmetrical matrix, it is sufficient to calculate the numerical values only of the elements lying below and on the main diagonal.

In order to calculate the elements of the  $r$ th row, the following equations have to be solved numerically

$${}^r \Gamma_{-k} = \theta(k) - \theta(r+k), \quad k=0, 1, 2, \dots, n$$

$${}^r \Gamma_k = c \cdot \min(k, r) + \theta(k) - \theta(|r-k|), \quad k=1, 2, \dots, n+r \quad (15)$$

where

$$c = C_0^2 T \quad (16)$$

and  $\min(k, r)$  denotes the lower of the numbers  $k, r$ .

The figures  ${}^r \Gamma_{-n}$  and  ${}^r \Gamma_{n+r}$  obtained in this way are entered into a column as shown in Table 1(a).

Table 1

(a)		(b)	
	${}^r \Gamma_{-n}$		${}^r U_0$ $R_{r0}$
	${}^r \Gamma_{-n+1}$		${}^r U_1$ $R_{r1}$
	$\vdots$		$\vdots$
$A_n$	${}^r \Gamma_{-n+k}$	$A_0$	${}^r U_s \rightarrow R_{rs}$
$A_{n-1}$	${}^r \Gamma_{-n+k+1}$	$A_1$	${}^r U_{s+1}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$A_1$	${}^r \Gamma_{k-1}$ ${}^r U_{k-1}$	$A_{n-1}$	${}^r U_{s+n-1}$
$A_0$	${}^r \Gamma_k \rightarrow {}^r U_k$	$A_n$	${}^r U_{s+n}$
	$\vdots$		$\vdots$
	${}^r \Gamma_{n+r}$		${}^r U_{n+r}$

A slip of paper is laid beside the column with the coefficients  $A_n, A_{n-1}, \dots, A_0$  of the denominator of pulse-transfer function  $G(z)$  written on it one below the other. The product of figures lying beside each other (see column (a) of Table 1) supplies the value of

$${}^r U_k = \sum_{i=0}^n A_i {}^r \Gamma_{k-i} \quad (17)$$

which is then entered in the next column into the row containing the coefficient  $A_0$ . All the required values of  ${}^r U_k$  ( $k = 1, 2, \dots, n+r$ ) are then obtained by a gradual shifting of the paper slip with the coefficients  $A_i$  written on it.

The next operation represented by

$$R_{rs} = \sum_{i=0}^n A_i {}^r U_{i+s} \quad (18)$$

is performed again by using the paper slip with coefficients  $A_i$ . However, this time the coefficients are written in the opposite order as can be seen from column (b) of Table 1 where the paper slip is drawn in a position at which the numerical value of  $R_{rs}$  is being determined. The value  $R_{r0}$  is already the requested element of column matrix  $[R_{r0}]$  for the weighting function  $q(t) = 1$ . The elements of matrix  $[K_{rs}]$  are obtained simply as the difference

$$K_{rs} = R_{rs} - R_{r0} \quad (19)$$

The correctness of the calculations made so far can be checked by the relation

$$K_{rr} + 2R_{r0} = rc \left( \sum_{i=0}^n A_i \right)^2 \quad (20)$$

The elements of matrices  $[\hat{K}_{rs}]$  and  $[\hat{R}_{r0}]$  pertaining to the weighting function  $q(t) = 1(t - T)$  are obtained by adding the same figure of  $\Delta K$  respectively  $\Delta R$  to all elements of matrices  $[K_{rs}]$  and  $[R_{r0}]$  respectively.

$$\begin{aligned} \hat{K}_{rs} &= K_{rs} + \Delta K \\ \hat{R}_{r0} &= R_{r0} + \Delta R \end{aligned} \quad (21)$$

The figures to be added are obtained from

$$\begin{aligned} \Delta K &= -\lambda A_0 \\ \Delta R &= \lambda A_0^2 - \kappa A_0 \sum_{i=1}^n A_i \end{aligned} \quad (22)$$

where

$$\begin{aligned} \lambda &= c - \theta(0) - C_0 \sum_{v=1}^n \frac{C_v}{P_v} + \sum_{v=1}^n C_v z_v \left( \frac{2C_0}{P_v} + \delta_v \right) \\ \delta_v &= \sum_{\mu=1}^n \frac{C_\mu z_\mu}{P_v + P_\mu} \\ \kappa &= \sum_{v=1}^n \frac{C_v}{P_v} (1 - z_v) - C \end{aligned} \quad (23)$$

Case (b)

In the case where one pole of transfer function  $S(p)$  lies in the origin, the step function responses of the system is given by the equation

$$s(t) = \mathcal{L}^{-1} \left\{ \frac{S(p)}{p} \right\} = C_1 t + C_0 + \sum_{v=1}^N C_v e^{p_v t}, \quad N = n - 1 \quad (24)$$

The elements of matrices  $[K_{rs}]$  and  $[R_{r0}]$  are calculated by the same method, only some values are calculated according to changed formulas.

Now, the numerical values of  $\theta(k)$  for  $k = 1, 2, \dots, N+L$  (where  $N = n - 1$  is the number of non-zero poles) are obtained from the relations

$$k > 0, \quad \theta(k) = \sum_{v=1}^N \rho_v z_v^{k-1} \quad (25)$$

where

$$\begin{aligned} \rho_v &= C_v (1 - z_v)^2 \left( \frac{C_{-1}}{P_v^2} - \frac{C_0}{P_v} - \delta_v \right) \\ z_v &= e^{p_v T}, \quad \delta_v = \sum_{\mu=1}^N \frac{C_\mu}{P_v + P_\mu} \end{aligned} \quad (26)$$

The value of  $\theta(0)$  is calculated separately from

$$\theta(0) = \frac{C_{-1}^2 T^3}{6} - C_0^2 T - 2C_{-1} T \sum_{v=1}^N \frac{C_v}{P_v} - 2 \sum_{v=1}^N \frac{\rho_v}{1 - z_v} \quad (27)$$

The further procedure of calculation remains the same, except that for 0 we substitute everywhere

$$c = C_{-1}^2 T^3 \quad (28)$$

and instead of coefficients  $A_i$  ( $i = 0, 1, \dots, n$ ) we use everywhere the coefficients  $\bar{A}_i$  ( $i = 0, 1, 2, \dots, N$ ). Their relationship can be seen from the arrangement of the denominator of the pulse-transfer function  $G(z)$

$$\begin{aligned} A(z) &= A_0 + A_1 z^{-1} + \dots + A_n z^{-n} \\ &= (1 - z^{-1})(\bar{A}_0 + \bar{A}_1 z^{-1} + \dots + \bar{A}_N z^{-N}), \quad N = n - 1 \end{aligned} \quad (29)$$

This arrangement is made possible just because one pole of the transfer function  $S(p)$  equals zero.

The last difference in comparison with case (a) lies in the determination of the numerical values of  $\lambda$  and  $\kappa$  which are used in the determination of the matrices pertaining to the weighting function  $qct) = 1(t - T)$ . They are calculated from the formulas

$$\begin{aligned} \lambda &= \frac{C}{2} + C_0 C_{-1} T^2 - \theta(0) - \sum_{v=1}^N C_v (1 - z_v) \left( \frac{2C_{-1} T}{P_v} - \delta_v \right) \\ \delta_v &= \sum_{\mu=1}^N \frac{C_\mu (1 - z_\mu)}{P_v + P_\mu} \\ \kappa &= \frac{C}{2} + C_0 C_{-1} T^2 - C_{-1} T \sum_{v=1}^N \frac{C_v}{P_v} (1 - z_v) \end{aligned} \quad (30)$$

Example

In order to illustrate the method of calculation described generally in the preceding section, the calculation of a concrete case is given below. The transfer function of the plant is

$$S(p) = \frac{6p + 4.5}{(p+2)(p+1)(p+0.5)}$$

All poles of this transfer function are different from zero, the problem discussed is thus of the type of Case (a). The unit-step function response of the system is

$$s(t) = \mathcal{L}^{-1} \left\{ \frac{S(p)}{p} \right\} = C_0 + C_1 e^{p_1 t} + C_2 e^{p_2 t} + C_3 e^{p_3 t}$$

$$p_1 = -2; \quad p_2 = -1; \quad p_3 = -0.5; \quad C_0 = 4.5; \quad C_1 = 2.5; \quad C_2 = -3; \quad C_3 = -4.$$

The continuously acting member of the system has a pulse-transfer function

$$G(z) = \frac{B(z)}{A(z)} = \frac{1.309 z^{-1} - 0.092 z^{-2} + 0.248 z^{-3}}{1 - 1.110 z^{-1} + 0.355 z^{-2} - 0.030 z^{-3}}$$

122/4

Let the calculation of the coefficients of polynomial  $D(z)$  for its selected degree  $L = 1, 2, 3$  be presented. By solving eqns (14) and (13) we obtain

$$\begin{aligned} \delta_1 &= 4.9375; & \delta_2 &= -10; & \delta_3 &= -20 \\ \rho_1 &= -0.6875; & \rho_2 &= 3.5; & \rho_3 &= 16 \end{aligned}$$

$$k = 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6$$

$$\theta(k) = 18.813 \quad 10.899 \quad 6.347 \quad 3.743 \quad 2.229 \quad 1.337 \quad 0.805$$

Table 2 contains the calculation of the elements of the second row ( $r = 2$ ) of matrices  $K_{rs}$  and  $R_{r0}$ .

Table 2

$k$	${}^2\Gamma_k$	${}^2U_k$	$R_{2k}$	$K_{2k}$
-3	2.406			
-2	4.118			
-1	7.156			
0	12.465	5.913	-0.694	
1	20.250	8.833	0.937	1.630
2	28.035	9.771	2.567	3.260
3	33.344	9.045		
4	36.382	8.720		
5	38.094	8.710		

The first column in Table 2 has been compiled according to eqns (15), the second and third have been calculated schematically according to Table 1. The fourth column containing elements  $K_{2k}$  has been obtained by means of relation (19).

The elements of the remaining two rows of matrices  $[K_{rs}]$  and  $[R_{r0}]$  are calculated in a similar way. As  $[K_{rs}]$  is a symmetrical matrix, it is sufficient to calculate only its elements lying to the left of the main diagonal and those on the diagonal itself. The correctness of the calculation is checked by substituting into relation (20) which is the means of checking almost all numerical operations represented in Table 2 including the compilation of the first column  ${}^2\Gamma_k$ .

By this method it has been possible to obtain a system of linear equations for the sought after coefficients pertaining to the weighting function  $q(t) = 1$ :

$$\begin{bmatrix} 1.697 & 1.630 & 13.27 \\ 1.630 & 3.260 & 2.890 \\ 1.327 & 2.890 & 4.217 \end{bmatrix} \begin{bmatrix} D_1 \\ D_2 \\ D_3 \end{bmatrix} = \begin{bmatrix} 0.380 \\ 0.694 \\ 0.704 \end{bmatrix}$$

As the elements of matrices  $[K_{rs}]$  and  $[R_{r0}]$  are independent of the chosen degree  $L$  of polynomial  $D(z)$  the mere reduction of the respective matrices will suffice to meet the case of  $L = 1, 2$ . By the solution of the above system of equations coefficients  $D_i$  are obtained for  $i \neq 0$ , while the coefficient  $D_0$  follows from condition (8)

$$D_0 = 1 - \sum_{i=1}^L D_i$$

In this way the following results have been obtained

$L$	$D_0$	$D_1$	$D_2$	$D_3$
3	0.758	0.046	0.140	0.057
2	0.768	0.038	0.194	
1	0.776	0.224		

In order to obtain the system of equations for the coefficients  $\hat{D}_i$  pertaining to the weighting function  $q(t) = 1(t - T)$  it will suffice, in accordance with relation (21), to add to each element of matrices  $[K_{rs}]$  and  $[R_{r0}]$  respectively the following quantities

$$\Delta K = -0.4548 \quad \text{and} \quad \Delta R = -0.0646$$

obtained by the numerical solution of eqns (22) and (23). In this way one obtains

$L$	$\hat{D}_0$	$\hat{D}_1$	$\hat{D}_2$	$\hat{D}_3$
3	0.611	0.186	0.120	0.084
2	0.631	0.170	0.199	
1	0.642	0.358		

The pulse-transfer function of the continuously acting member of the system  $G(z) = B(z)/A(z)$  and the polynomial  $D(z)$ , the coefficients of which have just been calculated, determine completely the necessary transfer function (9) of the computer. The respective curves of the controlled variable  $x$  following the unit-step change of input signal  $w$  are represented in Figure 3 for the weighting function  $q(t) = 1$ , and in Figure 4 for the function  $q(t) = 1(t - T)$ .

It can be seen from Figures 3 and 4 that, compared with the minimum number of steps ( $L = 0$ ), a considerable improvement has been attained, especially in the case where in the minimalization of the integral of squared error the errors have been considered as occurring only after the first sampling period.

**Derivations and Proofs**

First of all it will be proved that the above stated results hold for the case where all the poles of transfer function  $S(p)$  are different from zero.

The sequence of the increments of the variable  $e_2^*(t)$

$$\Delta e_2 [i] = e_2^*(iT) - e_2^*(iT - T)$$

has, according to eqns (5) and (6), the z-transform of

$$\mathcal{L} \{ \Delta e_2 [i] \} = (1 - z^{-1}) E_2(z) = \frac{1}{B(1)} D(z) A(z) \quad (31)$$

From this z-transform it follows obviously

$$\Delta e_2 [i] = \frac{1}{B(1)} \sum_{s=0}^L D_s A_{i-s} \quad (32)$$

where  $A_k = 0$  for  $k < 0$  and  $k > n$ ;  $\Delta e_2 [i] = 0$  for  $i > n + L$ . Eqn (32) contains all the  $L + 1$  coefficients of polynomial  $D(z)$ ; however, only  $L$  of them can be selected, as it is necessary to fulfil condition (8) that is  $D(1) = 1$ . For the purpose of fulfilling this condition let coefficient  $D_0$  be detached

$$D_0 = 1 - \sum_{s=1}^L D_s \quad (33)$$

and eliminated from eqn (32)

$$\Delta e_2 [i] = \frac{1}{B(1)} \left[ \sum_{s=1}^L D_s (A_{i-s} - A_i) + A_i \right] \quad (34)$$

Now, the time curve of the manipulated variable  $y(t)$  is dissolved into the sum of unit-step functions

$$y(t) = \sum_{i=0}^{n+L} \Delta e_2 [i] 1(t-iT)$$

and the curve of the controlled variable  $x(t)$  can then be represented by the superposition of the unit-step responses

$$x(t) = \sum_{i=0}^{n+L} \Delta e_2 [i] s(t-iT) \quad (35)$$

Then for error  $e_1(t)$  it holds that

$$\begin{aligned} e_1(t) &= 1 - x(t) = 1 - \sum_{i=0}^{n+L} \Delta e_2 [i] s(t-iT) \\ &= \sum_{i=0}^{n+L} \Delta e_2 [i] \bar{s}(t-iT) \end{aligned} \quad (36)$$

where

$$\begin{aligned} \bar{s}(t-iT) &= s(\infty) - s(t-iT) \\ t > iT, \bar{s}(t-iT) &= - \sum_{v=1}^n C_v e^{P_v(t-iT)} \\ t < iT, \bar{s}(t-iT) &= C_0 \end{aligned} \quad (37)$$

If the integral of squared error (10) has a minimum value, the coefficients of polynomial  $D(z)$  must fulfil the equations

$$\frac{\partial J}{\partial D_r} = 0, \quad r = 1, 2, \dots, L \quad (38)$$

After the above indicated derivative of the integral it follows

$$2 \int_0^{\infty} q(t) e_1(t) \frac{\partial e_1(t)}{\partial D_r} dt = 0 \quad (39)$$

The necessary partial derivative is determined from relations (36) and (34)

$$\frac{\partial e_1(t)}{\partial D_r} = \frac{1}{B(1)} \sum_{j=0}^{n+L} (A_{j-r} - A_j) \bar{s}(t-jT) \quad (40)$$

By substituting (40) and (36) together with (34) into condition (39), and by altering the sequence of addition, it follows

$$\begin{aligned} \int_0^{\infty} q(t) \left\{ \sum_{s=1}^L D_s \sum_{i=0}^{n+L} \sum_{j=0}^{n+L} (A_{i-s} - A_i) (A_{j-r} - A_j) \bar{s}(t-iT) \right. \\ \left. \bar{s}(t-jT) + \sum_{i=0}^{n+L} \sum_{j=0}^{n+L} A_i (A_{j-r} - A_j) \bar{s}(t-iT) \bar{s}(t-jT) \right\} dt = 0 \end{aligned} \quad (41)$$

Under the accepted pre-condition of transfer function  $S(p)$  having all its poles in the left semi-plane the integrals

$$\sigma_{ij} = \int_0^{\infty} q(t) \bar{s}(t-iT) \bar{s}(t-jT) dt \quad (42)$$

converge and in eqn (41) the integral of the sum can be expressed as the sum of the integrals

$$\begin{aligned} \sum_{s=1}^L D_s \sum_{i=0}^{n+L} \sum_{j=0}^{n+L} (A_{i-s} - A_i) (A_{j-r} - A_j) \sigma_{ij} \\ + \sum_{i=0}^n \sum_{j=0}^{n+L} A_i (A_{j-r} - A_j) \sigma_{ij} = 0 \end{aligned}$$

and in the abbreviated form

$$\sum_{s=1}^L D_s K_{rs} + R_{r0} = 0, \quad r = 1, 2, \dots, L \quad (43)$$

with the following denotations

$$K_{rs} = \sum_{i=0}^{n+L} \sum_{j=0}^{n+L} (A_{i-s} - A_i) (A_{j-r} - A_j) \sigma_{ij} \quad (44)$$

$$R_{r0} = \sum_{i=0}^n \sum_{j=0}^{n+L} A_i (A_{j-r} - A_j) \sigma_{ij} \quad (45)$$

For the degree  $L$  of the selectable coefficients  $D_s$  the system of linear equations is thus obtained that can be written in the matrix form (11) as

$$[K_{rs}] [D_s] = -[R_{r0}]$$

By interchanging the subscripts in relation (44) it can be easily proved that  $[K_{rs}]$  is a symmetrical matrix.

By the solution of integral (42) it follows for the case of weighting function  $q(t) = 1$

$$\sigma_{ij} = c \min(i, j) + C_0 \sum_{v=1}^n \frac{C_v}{P_v} - \theta(|j-i|) \quad (46)$$

where the function  $\theta(k)$  is determined by relation (13), and  $c$  according to relation (16). Integral (42) for the weighting function  $q(t) = 1(t-T)$  is to be denoted by  $\hat{\sigma}_{ij}$ . It holds

$$i \neq 0, j \neq 0, \hat{\sigma}_{ij} = \sigma_{ij} - c$$

$$\hat{\sigma}_{0k} = \hat{\sigma}_{k0} = C_0 \sum_{v=1}^n \frac{C_v}{P_v} z_v - \theta(k) \quad (47)$$

$$\hat{\sigma}_{00} = - \sum_{v=1}^n C_v z_v \delta_v$$

where  $\delta_v$  is determined by the second of relations (23), and  $z_v = e^{P_v T}$ .

The calculation of the elements of matrices  $[K_{rs}]$  and  $[R_{r0}]$  according to relations (44) and (45) would be very laborious. For this reason let some arrangements be introduced that will simplify this calculation considerably.

First, let us divide relation (44) into two terms

$$\begin{aligned} K_{rs} &= \sum_{i=0}^{n+L} \sum_{j=0}^{n+L} \bar{A}_{i-s} (A_{j-r} - A_j) \\ &\quad - \sum_{i=0}^n \sum_{j=0}^{n+L} A_i (A_{j-r} - A_j) \sigma_{ij} \end{aligned}$$

Now, if the summing subscript  $i$  in the first member is shifted by  $s$ , i.e.  $i-s \rightarrow i$ , and considering that  $A_i = 0$  for  $i > n$  and  $i < 0$ , it follows

122/6

$$K_{rs} = \sum_{i=0}^n \sum_{j=0}^n A_i (A_{j-r} - A_j) \sigma_{i+s, j} - \sum_{i=0}^n \sum_{j=0}^n A_i (A_{j-r} - A_j) \sigma_{ij}$$

According to (45) the second term equals  $R_{r0}$ , and let the first be denoted

$$R_{rs} = \sum_{i=0}^n A_i \sum_{j=0}^n (A_{j-r} - A_j) \sigma_{i+s, j} \quad (48)$$

In this way relation (19) has been obtained

$$K_{rs} = R_{rs} - R_{r0} \quad (49)$$

As the term  $R_{r0}$  represents a special case of  $R_{rs}$  with  $s = 0$ , it will suffice further to seek only the numerical solution of (48) for  $R_{rs}$ . Let it be written in the following form

$$R_{rs} = \sum_{i=0}^n A_i \sum_{j=0}^n A_j (\sigma_{i+s, j+r} - \sigma_{i+s, j})$$

If we denote

$${}^r U_{i+s} = \sum_{j=0}^n A_j (\sigma_{i+s, j+r} - \sigma_{i+s, j}) \quad (50)$$

we obtain relation (18)

$$R_{rs} = \sum_{i=0}^n A_i {}^r U_{i+s} \quad (51)$$

All values of  ${}^r U$  required for the calculation of the  $r$ th row of matrices  $[K_{rs}]$  and  $[R_{r0}]$  can be obtained as the product of the rectangular matrix

$$\boxed{\text{Eqn (52)}}^*$$

and of the column matrix

$$\begin{aligned} [A_i] &= [A_0, A_1, \dots, A_n] \\ [{}^r U] &= [\sigma] [A_i] \end{aligned} \quad (53)$$

It will be proved that in the case of weighting function  $q(t) = 1$  matrix  $[\sigma]$  has all its elements lying on the lines parallel to the main diagonal of the same value. For the  $m$ th element of the  $k$ th parallel above main diagonal it holds

$$\begin{aligned} \sigma_{m, k+r+m} - \sigma_{m, k+m} &= cm + C_0 \sum_{v=1}^n \frac{C_v}{P_v} - \theta(k+r) \\ - \left[ cm + C_0 \sum_{v=1}^n \frac{C_v}{P_v} - \theta(k) \right] &= \theta(k) - \theta(k+r) \end{aligned}$$

As this relation is independent of  $m$ , all elements lying on this parallel are equal, and they may be denoted by the same symbol

$${}^r \Gamma_k = \theta(k) - \theta(k+r)$$

Similarly it holds for the elements on the  $k$ th parallel below the main diagonal

$${}^r \Gamma_k = \sigma_{k+m, r+m} - \sigma_{k+m, m} = c \min(k, r) + \theta(k) - \theta(|r-k|)$$

For the main diagonal  $k = 0$ .

Due to this property of matrix  $[\sigma]$  it is possible to arrange the numerical solution of matrix product (51) into a scheme shown in Table 1 (a) which can be easily found by comparing both methods of calculation.

It remains yet to prove the validity of formulas (21), (22), and (23) by which the former results are to be corrected, if errors are being considered only after the first sampling period. By substituting into matrix (50) for  $\sigma_{ij}$  (46) the terms  $\hat{\sigma}_{ij}$  (47) calculated for the weighting function  $q(t) = 1(t-T)$ , it can be seen that only the first column has been altered. Obviously it holds that

$${}^r \hat{U}_k = {}^r U_k + (\hat{\sigma}_{k,r} - \sigma_{k,r} - \hat{\sigma}_{k,0} + \sigma_{k,0}) A_0 \quad (54)$$

When calculating the term in the parentheses, it is necessary to differentiate two cases:  $k > 0$  and  $k = 0$ . By substituting relations (46) and (47), we obtain in the first case the relation

$$k > 0, \hat{\sigma}_{k,r} - \sigma_{k,r} - \hat{\sigma}_{k,0} + \sigma_{k,0} = C_0 \sum_{v=1}^n \frac{C_v}{P_v} (1 - z_v) = x$$

which is independent of  $k$ . Similarly for  $k = 0$

$$\begin{aligned} \hat{\sigma}_{0,r} - \sigma_{0,r} - \hat{\sigma}_{0,0} + \sigma_{0,0} \\ = -C_0 \sum_{v=1}^n \frac{C_v}{P_v} (1 - z_v) + \sum_{v=1}^n C_v z_v \delta_v + C_0 \sum_{v=1}^n \frac{C_v}{P_v} - \theta(0) = \kappa_0 \end{aligned}$$

With this denotation the relation (54) may be rewritten in the form

$$\begin{aligned} k > 0, & \quad {}^r \hat{U}_k = {}^r U_k + \kappa A_0 \\ k = 0, & \quad {}^r \hat{U}_0 = {}^r U_0 + \kappa_0 A_0 \end{aligned} \quad (55)$$

For the verification of formulas (21) and (22) it will suffice to execute operations (51) and (49) with the relations (55), and to denote  $k_0 - k = \lambda$ .

The checking formula (20) can be verified by substituting relations (44) and (45) and by using the relation

$$\sigma_{i+r, j+r} - \sigma_{i, j} = rc$$

that follows from eqn (46).

In Case (b), with the transfer function  $S(p)$  having one zero pole, the continuously acting member of the system is astatic [ $s(\infty) = \infty$ ], and integrals (42) are not converging. It is possible to by-pass this difficulty, if the curve of the controlled variable is not represented as the superposition of unit-step responses, but as the superposition of responses to rectangular pulses. Otherwise the procedure of derivation is the same as in Case (a).

\* Eqn (52):

$$[\sigma] = \begin{bmatrix} \sigma_{0,r} & -\sigma_{0,0}; & \sigma_{0,1+r} & -\sigma_{0,1}; & \sigma_{0,2+r} & -\sigma_{0,2}; & \dots & \sigma_{0,n+r} & -\sigma_{0,n} \\ \sigma_{1,r} & -\sigma_{1,0}; & \sigma_{1,1+r} & -\sigma_{1,1}; & \sigma_{1,2+r} & -\sigma_{1,2}; & \dots & \sigma_{1,n+r} & -\sigma_{1,n} \\ \sigma_{2,r} & -\sigma_{2,0}; & \sigma_{2,1+r} & -\sigma_{2,1}; & \sigma_{2,2+r} & -\sigma_{2,2}; & \dots & \sigma_{2,n+r} & -\sigma_{2,n} \\ \vdots & & \vdots & & \vdots & & & \vdots & \\ \sigma_{n+L,r} & -\sigma_{n+L,0}; & \sigma_{n+L,1+r} & -\sigma_{n+L,1}; & \sigma_{n+L,2+r} & -\sigma_{n+L,2}; & \dots & \sigma_{n+L,n+r} & -\sigma_{n+L,n} \end{bmatrix} \quad (52)$$

References

- <sup>1</sup> RAGAZZINI, J. R. and FRANKLIN, G. F. *Sampled-data Control Systems*. 1958. New York; McGraw-Hill
- <sup>2</sup> JURY, E. I. *Sampled-data Control Systems*. 1958. New York; Wiley
- <sup>3</sup> TOU, J. T. *Digital and Sampled-data Control Systems*. 1959. New York; McGraw-Hill
- <sup>4</sup> STREJC, V. Ensuring reliability in complex automation by automatic digital computers. *Automatizace*. V (1962) 5

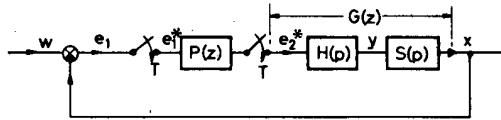


Figure 1

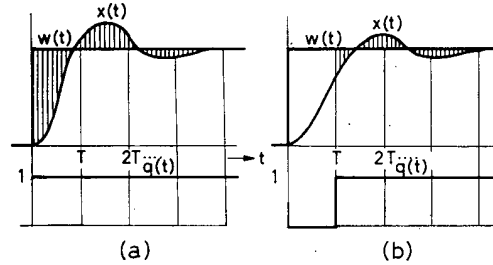


Figure 2

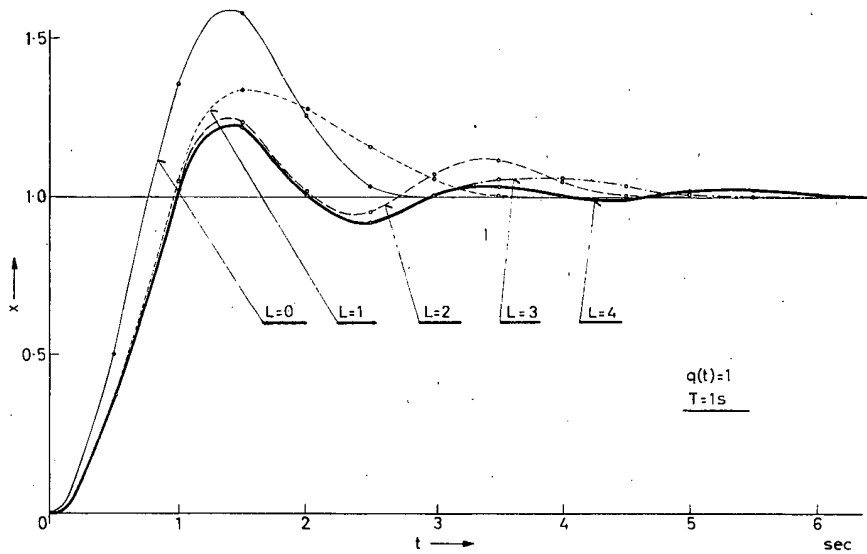


Figure 3

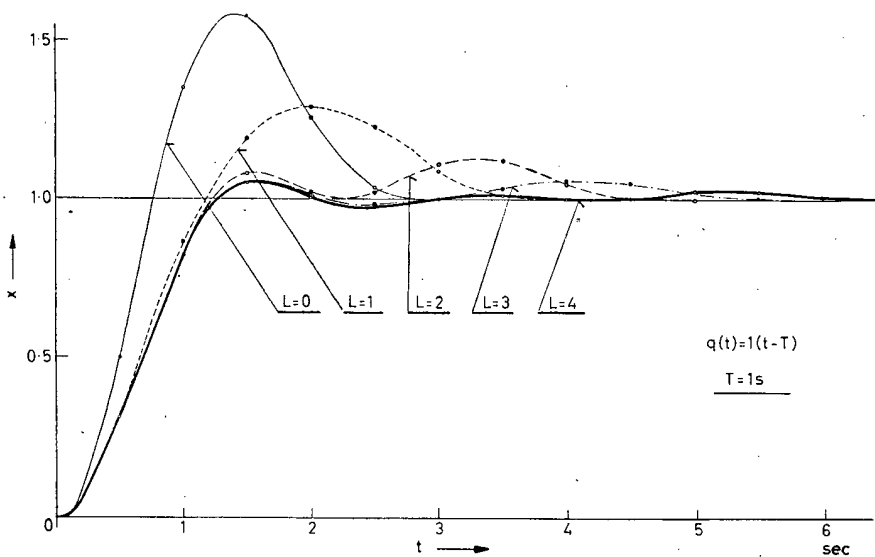


Figure 4

# The Dynamic Properties of Rectification Stations with Plate Columns

J. ZAVORKA

@zech-

The control of rectification stations, as carried out at the present time, is confined only to some control loops which are designed without any thorough theoretical consideration. As far as individual control diagrams are concerned, quite a number of them have been designed; for instance, see Anizinov<sup>1</sup>. The advantages and shortcomings of various connection schemes have been published by the respective authors, however, and the evaluation is mainly based on technical sense and experimental results. Information on the general operational analysis of rectification columns has been appearing only recently<sup>2, 3, 5, 7, 11</sup>.

In most of these papers the pressure and hold-up of the plate have been considered as constant quantities. Due to this, the validity of results is limited to cases with slow changes in the input quantities; for instance, changes in feed composition or changes occurring during the starting of the column. For rapidly changing input variables, for instance pressure, the results are erroneous. In view of these facts, an operational analysis was worked out by Voetter and Houtappel<sup>13</sup> where the pressure and hold-up of the plate were considered as variables. Starting from linearized equations the authors demonstrated that, nevertheless, the results hold for a rather wide range of input quantities. The same authors extended their study to ternary mixtures, and used digital computers for the calculation of dynamic properties. It has been found that the solutions of these problems are exceedingly time consuming with regard to the computer, and Rose and Williams<sup>5, 10, 11</sup>, attempted the modelling of the system on an analogue computer. However, these authors designed the model of the dynamics of the vapour phase as single-capacity members connected in series which does not correspond with reality. This deficiency has been eliminated by the work of Rijnsdorp and Maarleveld<sup>9</sup>, who succeeded in modelling a 32-plate column on an analogue computer built from passive elements especially for this purpose. The Bode frequency characteristics are the result of this work. As an example, one of these characteristics is shown in *Figure 1*. Obviously it cannot be evaluated, as the curve has no distinct straight sections to permit the determination of the respective intersects. Apart from this, it is not possible to agree with the assumption made by the authors in the equations describing the system, namely that the heat of evaporation is merely a function of pressure and independent of the composition of the mixture.

The aim of the present paper is to derive generally valid relationships for the computation of transfer functions for the individual input and output variables of the whole rectification station, to create in this way the possibility of comparing and assessing the advantages and shortcomings of various control diagrams, and to obtain the data necessary for the synthesis of

control loops and for the complex automation of rectification stations.

The task has been limited to rectification stations with plate columns for the separation of binary mixtures.

The purpose of the work is to determine the transfer functions of the system, which in turn determine the relationship between the input variables ( $N$ : the flow rate of the feed;  $X_N$ : the composition of the feed;  $P_k$ : the pressure in the condenser;  $G_1$ : the flow rate of the heating steam) and the output variables ( $A$ : the flow rate of the product;  $X_A$ : the composition of the product;  $B$ : the flow rate of the residue;  $X_B$ : the composition of the residue;  $P_0$ : pressure at the first plate of the column) and possibly between the concentration at some other plates.

The diagram of a rectification station with a plate column for the continuous separation of binary mixtures is shown in *Figure 2*.

For the investigation of dynamic properties let the rectifying station be divided into three sections shown by the dash line in the illustration. The first to be investigated is the independent rectifying column, the second section consists of the bottom of the column with the still, while the third section contains the top of the column, the condenser, the cooler and the condensate tank.

The rectifying column consists of plates that are to be considered as separate units with regard to function and construction. The diagram of a plate is shown in *Figure 3*. It can be seen that the plate may be acted upon by the following nine input variables:

$N$	The feed flow rate
$X_N$	The feed composition
$H_{l,N}$	The enthalpy of the feed
$V_{n-1}$	The flow rate of vapour from the plate below
$Y_{n-1}$	The concentration of this vapour
$H_{v,n-1}$	The enthalpy of this vapour
$L_{n+1}$	The reflux from the plate above
$X_{n+1}$	The composition of this reflux
$H_{l,n+1}$	The enthalpy of this reflux

By these variables changes are produced in nine output variables:

$M_{l,n}$	The liquid hold-up of the plate
$M_{v,n}$	The vapour hold-up of the plate
$P_n$	The pressure on the plate
$V_n$	The flow rate of vapour streaming from the plate
$L_n$	The reflux from the plate
$H_{v,n}$	The enthalpy of vapour streaming from the plate
$H_{l,n}$	The enthalpy of the reflux from the plate
$Y_n$	The composition of vapour
$X_n$	The composition of the reflux

127/2

The plate is described thus by a system of nine simultaneous equations which are now derived.

First, the material balance of the plate is set up.

$$\frac{dM_{l,n}}{d\tau} + \frac{dM_{v,n}}{d\tau} = L_{n+1} - L_n + V_{n-1} - V_n + N \quad (1)$$

By multiplying the individual terms by the corresponding concentrations the total material balance equation is transformed into the material balance of the more volatile component:

$$\frac{d(M_{l,n} \cdot X_n)}{d\tau} + \frac{d(M_{v,n} \cdot Y_n)}{d\tau} = L_{n+1} X_{n+1} - L_n \cdot X_n + V_{n-1} \cdot Y_{n-1} - V_n Y_n + N \cdot X_N \quad (2)$$

In accordance with the material balance equation it is possible to write the heat balance equation as follows:

$$\frac{d(M_{l,n} \cdot H_{l,n})}{d\tau} + \frac{d(M_{v,n} H_{v,n})}{d\tau} - V^* \cdot \frac{dP_n}{d\tau} = L_{n+1} H_{l,n+1} - L_n H_{l,n} + V_{n-1} \cdot H_{v,n-1} - V_n \cdot H_{v,n} + N \cdot H_N \quad (3)$$

The last term on the left-hand side of the equation (which represents the consideration given to the difference between the enthalpy of the vapour phase and its internal energy for which the equation holds) is neglected later with regard to the pressure changes being of the order of millimetres of water gauge.

The vapour flow rate depends on the square root of the pressure differential on two adjacent plates and on the density of the vapour. In view of the fact that the difference in pressure on two adjacent plates fluctuates within the range of 25–50 mm w.g., the influence of density may be neglected. The relationship between flow rate and pressure is then described by the equation

$$V_n^2 = k_0 \cdot (P_n - P_{n+1}) \quad (4)$$

The following relationship should be further investigated,

$$M_{l,n} = M_{l,n}(L_n)$$

By the application of relation

$$S_1 = 10 \cdot \left( \frac{10^{-5} \cdot \mu L}{1.773 \rho \cdot \gamma} \right)^{\frac{2}{3}}$$

one obtains

$$M_{l,n} = 10 \frac{S_1}{\xi} \left( \frac{10^{-5} \cdot \mu L}{1.773 \rho \cdot \gamma_l} \right)^{\frac{2}{3}} + K \quad (5)$$

Now consider the relationship between the concentration of the more volatile component in the vapours and the concentration of the more volatile component in the liquid during the state of equilibrium of both phases at the boiling point temperature of the binary mixture.

$$Y_n = Y_n(X_n) \quad (6)$$

The description of this relationship was attempted by a number of equations (Wohl, Scatchard-Hammer, Van Laar, Margules, symmetric<sup>4</sup>). However, they all contain constants that can be determined only experimentally. Due to this, and also due to their complexity, none of these equations has been accepted in practice. The effect of the composition of the liquid upon the composition of the vapours (established experimentally) is nor-

mally represented by the X–Y equilibrium diagram. This method of representation has been accepted for the following sections of this paper.

The remaining three equations are written in the form of general relations:

$$M_{v,n} = M_{v,n}(P_n) \quad (7)$$

$$H_{l,n} = H_{l,n}(P_n, X_n) \quad (8)$$

$$H_{v,n} = H_{v,n}(P_n, Y_n) \quad (9)$$

The system of the above-stated nine simultaneous equations describes one plate of the rectifying column. As interest here is only in the non-steady states of pressure, composition of the liquid phase and flow rate of the liquid phase, all other variables will be eliminated. The transfer functions of pressure, composition of the liquid phase and flow rate for one plate are obtained by the linearization of the equations or possibly by their transformation into differential equations, followed by the LW transformation and the arrangement of the equations. These transfer functions are used for drawing the partial block diagrams of one plate for the dynamic behaviour of the three variables. The block diagrams are shown in *Figure 4*. The overall block diagram of one plate is obtained by the interconnection of all three partial diagrams. The complete block diagram of the whole rectifying column is obtained by the interconnection of the block diagrams of the individual plates as shown in *Figure 5*. For the sake of clarity the multiplication constants are not shown in *Figure 5*. Now, it remains to conclude the block diagram of the column by the connections of the condenser and of the still.

The block diagram of the bottom section of the column (the first plate and still), and the block diagram of the top section of the column (the highest plate, condenser, cooler of the condensate, condensate tank and the piping) have been derived by a similar method as used for the derivation of the block diagram of the column proper. For the sake of brevity the respective procedures are omitted, and only their results are given in *Figures 6 and 7*.

The complete block diagrams of all sections of the rectifying station have been obtained so far. The description may serve as the source of some data for the modelling of the system. Owing to the high complexity of the diagram, a large number of integrating units will be required for the modelling and, therefore, it should be possible to model only the simplest stations with a small number of plates. For this reason the results of the preceding chapters have been subjected to a further theoretical analysis. The analysis follows the aim of simplifying the block diagram of the column proper so that it is suited for modelling, or so that it is possible to compute the transfer functions of the system. First of all it was necessary to determine the zones within which the values of individual design, physico-chemical and operational parameters can vary. Further the relations were to be stated that were required for the numerical solution of various terms occurring in the formulae for the time and multiplying constants. A quantitative analysis of the time and multiplying constants was made on the basis of these values and relations. The results obtained were used for certain simplifications of the formulae. Further, it appears that the dynamics of pressure and composition in the whole column are represented by block diagrams of the same structure (*Figure 8*). The diagram is formed by single-capacity members connected in series with feedbacks by-passing two members that follow behind. The output signals of this



chain are formed by the algebraic sum of the signals of three adjacent members and they link together the diagram of pressure and the diagram of composition.

The general analysis of this block diagram was made; a matrix calculation was used for deriving the matrices of the transfer functions of this block diagram as the functions of the number of the chain members (or of the member of the plates of the column). A further analysis was used for establishing the conditions at which the static value of the output signals of the above chain is equal to zero (the conditions are related to the magnitude of the multiplying constants), and the conditions at which it is possible also to neglect the dynamic value of the output signals (the conditions are related to the number of plates). It was proved by a further general procedure that the above-stated conditions are fulfilled by each column. Assume for an instant that, during the investigation of the dynamic properties of the distilling column, there is no interest in the non-steady states of pressure. Under this assumption, and owing to the former conclusions, it is possible to interrupt in the block diagram the connections of the pressure changes between the individual plates. This can be done because any disturbance entering any plate lying below or above the plate under investigation can influence neither the flow rate, nor the pressure, but only the pressure values at different points of the block diagram, or of the column, and these values are of no interest for the time being.

Now consider composition in the same way—supposing that one is not interested in the non-steady states of composition. Similarly, as in the case of pressures, the connections between individual plates may be interrupted. The block diagram is then transformed into the form shown in *Figure 9*. The values  $\varphi_p$  and  $\varphi_x$  are the sums of the input signals of the individual nodes of the block diagrams of the dynamics of pressure and composition respectively. Now the non-steady states of pressure and composition, that were formerly excluded from discussion, are considered. The partial block diagrams of pressure and composition respectively are easily attached to the diagram in *Figure 9* by introducing the signals  $\varphi_p$  and  $\varphi_x$  into the individual nodes of the block diagrams of pressure and composition respectively. The result is shown in *Figure 10*. The section of the block diagram bordered by the dot-and-dash lines corresponds with one plate of the rectifying column. By the solution of the system of equations written for all three nodes of the block diagram of one plate (naturally after the introduction of all multiplying constants) the transfer functions of all output variables of the plate are obtained. Finally, in the application of the transfer functions, it is possible to re-draw the block diagram shown in *Figure 6* into the final form according to *Figure 11*. This block diagram holds for a general column with any arbitrary parameters with regard to design, physico-chemical conditions and operation.

The block diagram shown in *Figure 11* together with the pertaining transfer functions and formulae for various constants and transfer functions, is the final product of the theoretical part of the work. These results make possible the computation of the transfer functions of a general rectifying station. During the solution of concrete problems a number of possible simplifications appeared that followed from the numerical evaluation of individual constants and plate transfer functions. It is not possible to prove the general validity of these simplifications. However, it may be assumed that they will be identical in most cases.

Further work<sup>16</sup> contains the practical computation of several transfer functions and step response curves of a concrete rectifying station on the basis of the results obtained from a general analysis. The necessary measurements were also made on this station in operation. After a comparison, the results of the computation were in very good agreement with the results of the measurements.

#### Nomenclature

$A$	Flow rate of the product (mol/sec)
$B$	Flow rate of the residue (mol/sec)
$C_n$	Multiplying constants
$c$	Specific heat of heating wall (kcal/kg °C)
$E$	Reflux ratio
$G$	Mass of the heating wall (kg)
$G_1$	Flow rate of the heating steam (kg/sec)
$H_l$	Enthalpy of the liquid (kcal/mol)
$H_N$	Enthalpy of the feed (kcal/mol)
$H_v$	Enthalpy of the vapour (kcal/mol)
$H_{01}$	Enthalpy of the heating steam (kcal/mol)
$H_{02}$	Enthalpy of the condensate from the still (kcal/mol)
$i$	Number of plates
$k$	Constants
$k$	Subscript of condenser
$L$	Reflux (mol/sec)
$L_{i+1}$	Reflux to the top (mol/sec)
$M_k$	Molar hold-up of the condenser (mol)
$M_l$	Liquid hold-up of the plate (mol)
$M_v$	Vapour hold-up of the plate (mol)
$N$	Feed flow rate (mol/sec)
$N$	Subscript of feed plate
$n$	Ordinal number of plate
$O/p/$	Transfer function of the still
$P$	Pressure (atm)
$P_0$	Pressure in the heating system of the still (atm)
$P_k$	Pressure in the condenser (atm)
$Q_1$	Heat flow to the heating wall (kcal/sec)
$Q_2$	Heat flow from wall to substance (kcal/sec)
$Q/b/$	Elementary transfer function of the still
$r$	Latent heat (kcal/mol)
$s_l$	Surface area of liquid hold-up (dm <sup>2</sup> )
$s_1$	Heating wall area on steam side (m <sup>2</sup> )
$s_2$	Heating wall area on liquid side (m <sup>2</sup> )
$s_1$	Height of liquid level on plate above the vapour nozzle of the bubble-cap (dm)
$T_{s\text{ str}}$	Mean temperature of heating wall (°C)
$T_{1\text{ str}}$	Mean temperature of heating wall on the steam side (°C)
$T_2$	Temperature of heating wall on the side of the heated substance (°C)
$U$	Free energy (kcal)
$U_{01}$	Free energy of the heating steam entering the still (kcal)
$U_{02}$	Free energy of condensate leaving the still (kcal)
$V$	Flow rate of vapour through column (mol/sec)
$V^*$	Volume (l)
$V^*_0$	Steam volume in the still heating system (l)
$X$	Concentration of the more volatile component in the liquid (mol %)
$X_A$	Concentration of the more volatile component in the product (mol %)
$X_{A0}$	Concentration of the more volatile condensate component after the condenser (mol %)
$X_{A1}$	Concentration of the more volatile product component in the cooler of condensate (mol %)
$X_{A2}$	Concentration of the more volatile component in the reflux (mol %)

127/4

$X_B$	Concentration of the more volatile component in the residue (mol %)
$X_N$	Concentration of the more volatile component in the liquid on the feed plate (mol %)
$Y$	Concentration of the more volatile component in the vapour (mol %)
$Y_N$	Concentration of the more volatile component in the vapour on the feed plate (mol %)
$\alpha_1$	Heat transfer coefficient steam-heating wall (kcal/m <sup>2</sup> h°C)
$\alpha_2$	Heat transfer coefficient heating wall-liquid (kcal/m <sup>2</sup> h°C)
$\gamma_l$	Specific gravity of liquid (kg/l)
$\gamma_v$	Specific gravity of vapour (kg/l)
$A(P)$	Elementary transfer function of the flow rate of the liquid phase.
$\mu$	Molecular weight
$\Xi(P)$	Elementary transfer function of concentration molar volume (dm <sup>3</sup> /mol)
$\Pi(P)$	Elementary transfer function of pressure
$\varrho$	Circumference of down-take pipe (dm)
$\tau$	Time (sec)
$\tau_d$	Transport lag (sec)
$\tau_l$	Time constant of the elementary transfer function of the flow rate of the liquid phase (sec)
$\tau_p$	Time constant of the elementary transfer function of pressure (sec)
$\tau_{px}$	Derivative time constant of the pressure-concentration link (sec)
$\tau_g$	Time constant of the elementary transfer function of the still (sec)
$\tau_k$	Time constant of the transfer function of the condensate (sec)
$\tau_{k3}$	Derivative time constant of the transfer function of the condenser (sec)
$\tau_x$	Time constant of the elementary transfer function of the concentration (sec)
$\tau_{xp}$	Derivative time constant of the concentration—pressure link (sec)
$\tau_z$	Time constant of the condensate tank (sec)

## References

- 1 ANIZINOV, J. V. *Automatičeskoje regulirovanie proces rektifikacii*. 1957. Moscow; Gostoptechizdat
- 2 ARMSTRONG, W. D., and WILKINSON, W. L. *Trans. Instn. chem. Engrs, Lond.* 35 (1957), 352
- 3 DAVIDSON, J. F. *Trans. Instn. chem. Engrs, Lond.* 34 (1956), 44
- 4 HÁLA, E., PICK, J., FRIED, V., and VILIM, O. *Rovnováha kapalina-pára*. 1955. Prague; NČSAV
- 5 HARNETT, R. T., ROSE, A., and WILLIAMS, T. J. *Industr. Engng. Chem.* 48 (1956), 1008
- 6 JACKSON, F. R., and PIGFORD, R. L. *Industr. Engng. Chem.* 48 (1956), 1020
- 7 KIRSCHBAUM, E. *Destilier- und Rektifiziertchnik*. 1950
- 8 MARSHALL, W. R., and PIGFORD, R. L. *The Applications of Differential Equations to Chemical Engineering*. 1947. University Delaware
- 9 RIJNSDORP, J. E., and MAARLEVELD, A. Use of electrical analogues in the study of the dynamic behaviour and control of distillation columns. *J. Symp. Instrument Comp. Develop. Plant Design*. London 11-13 (1959)
- 10 ROSE, A., JOHNSON, C. L., and WILLIAMS, T. J. *Industr. Engng. Chem.* 48 (1956), 1173
- 11 ROSE, A., and WILLIAMS, T. J. *Industr. Engng. Chem.* 47 (1955), 2284
- 12 ROSENBRICK, H. H. *Trans. Instn. chem. Engrs, Lond.* 35 (1957) 347
- 13 VOETTER, H. *Plant and Process Dynamic Characteristic*. 1957. London; Butterworths
- 14 YU-CHIN-CHU, BRENECKE, R. J., GETTY, R. J. and RAJINDRA, P. *Distillation Equilibrium Data*. 1950. New York; ■■■■
- 15 ZÁVORKA, J. Obecný analytický rozbor dynamických vlastností rektifikačních stanic s patrovými kolonami pro dělení binárních směsí. *ÚTIA ČSAV* 68 (September 1960)
- 16 ZÁVORKA, J. Výpoč et některých přenosu kolony 31 (provoz 03) ve Stalinových závodech. e srovnání s výsledky měření. *ÚTIA ČSAV*, 86 (September 1961)

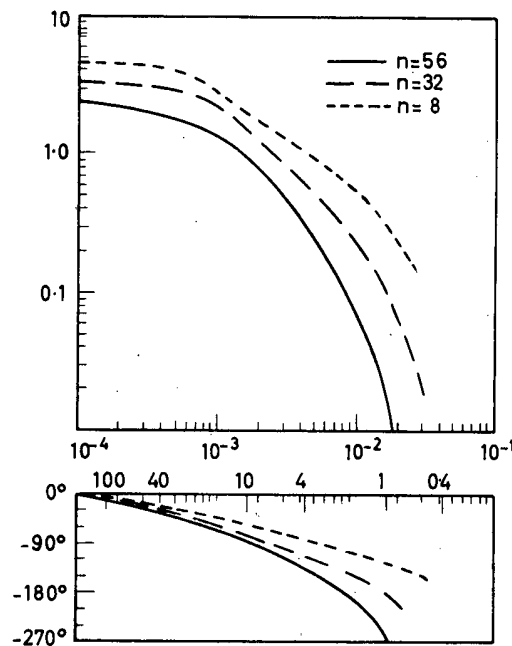


Figure 1

127/4

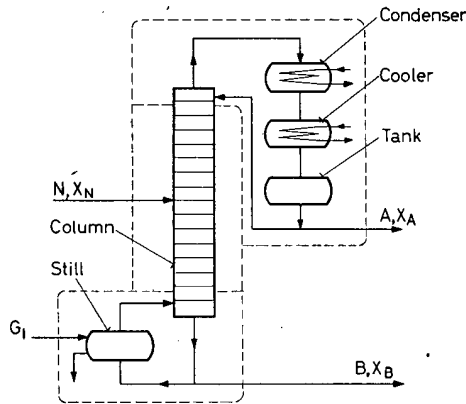


Figure 2

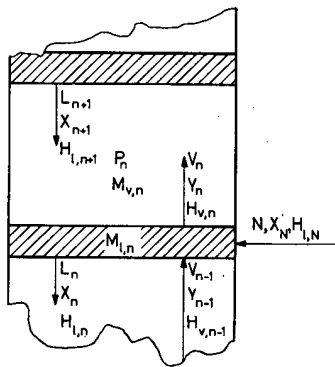


Figure 3

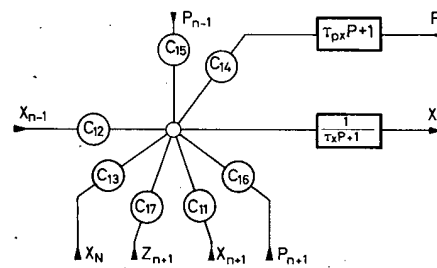
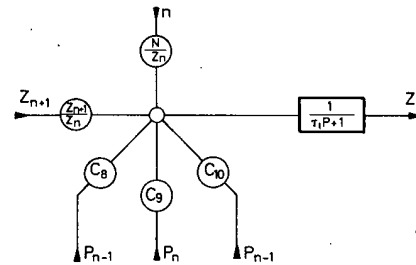
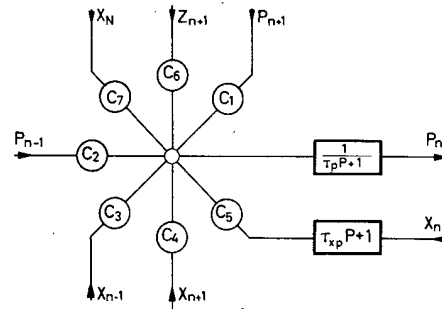


Figure 4

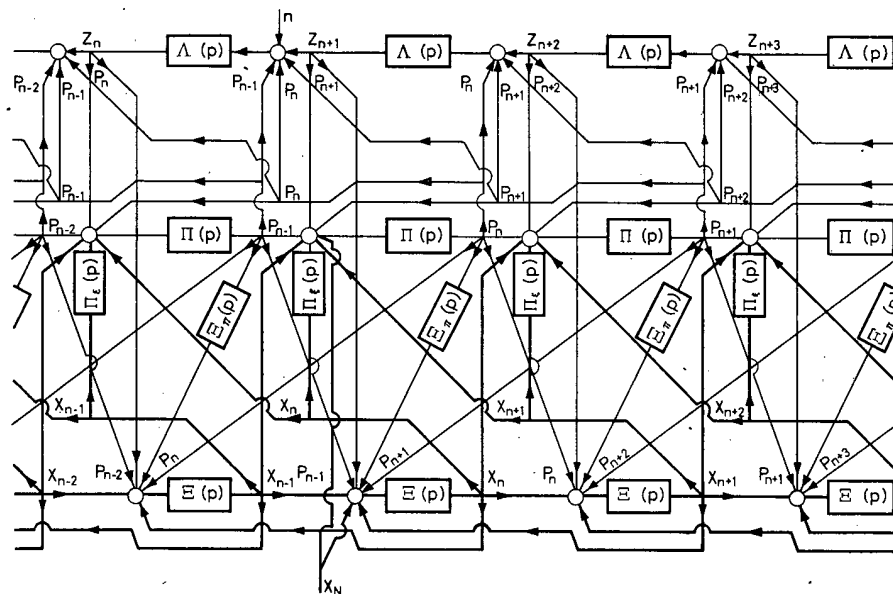


Figure 5

127/6

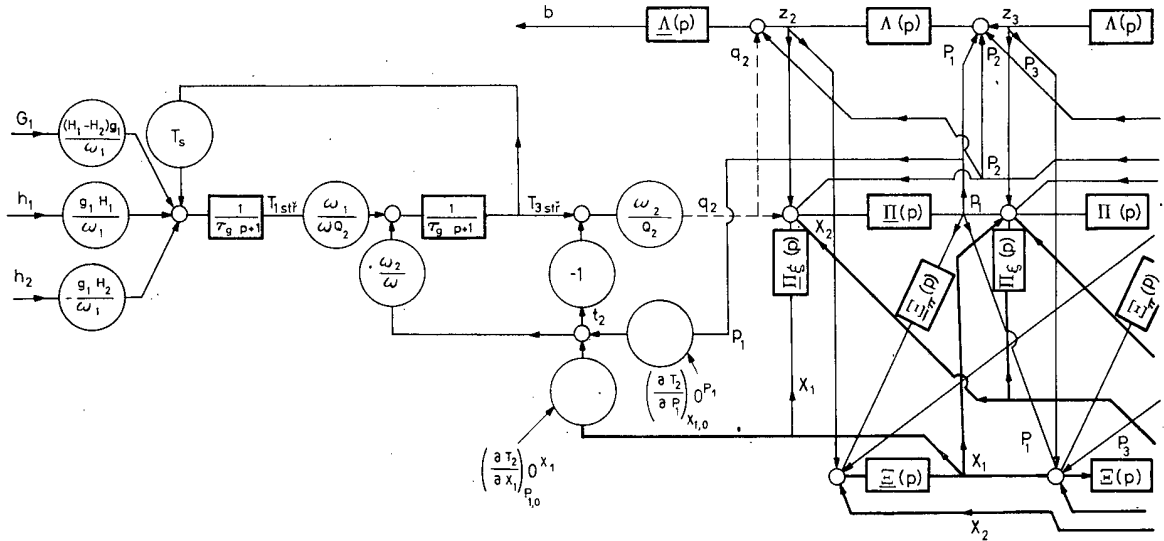


Figure 6

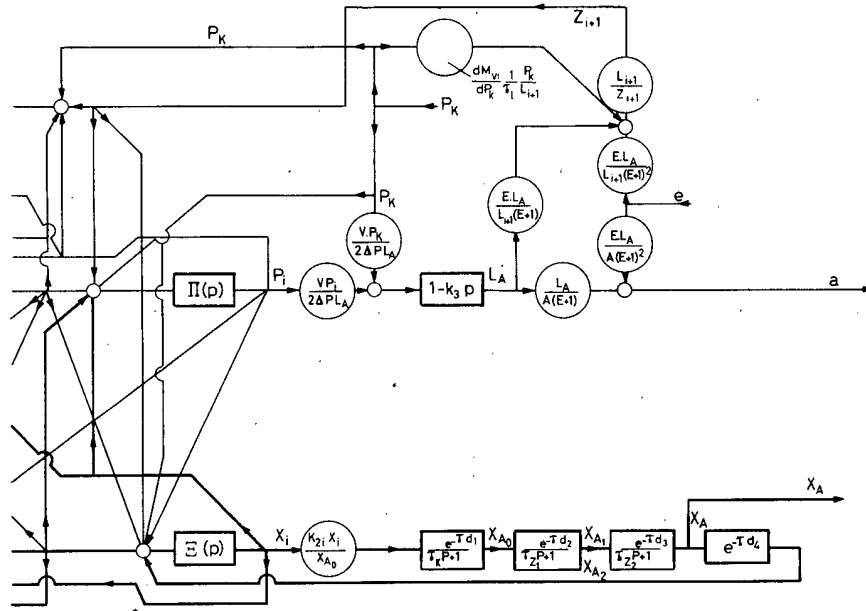


Figure 7

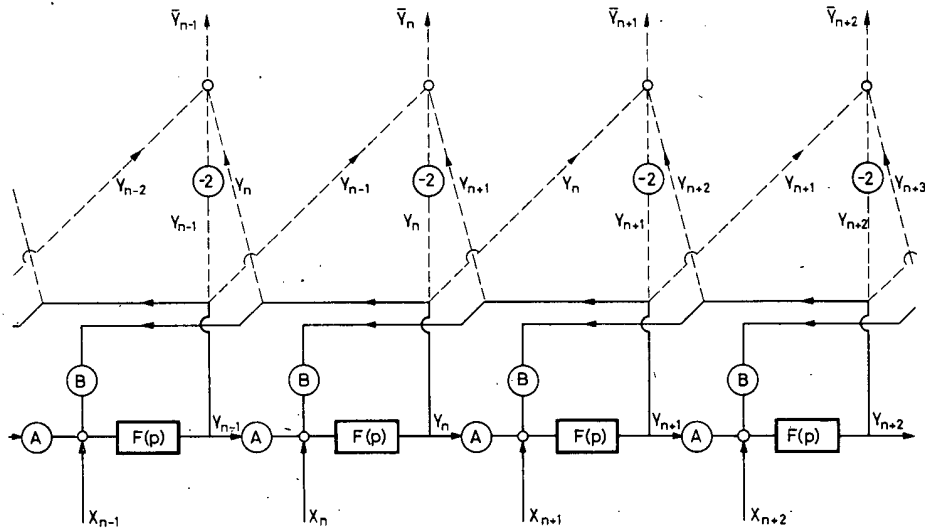


Figure 8

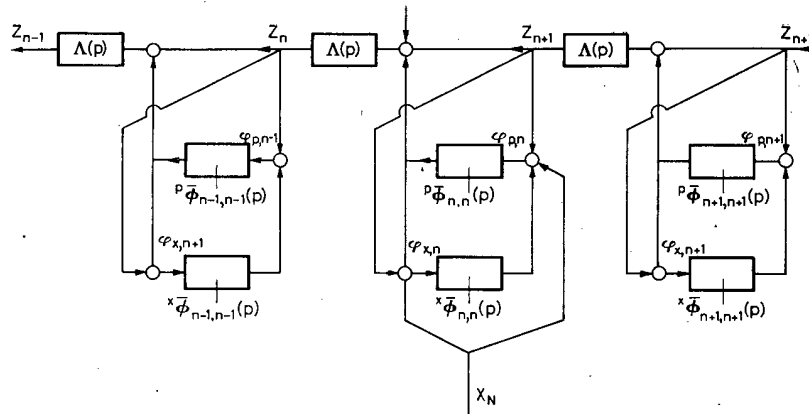


Figure 9

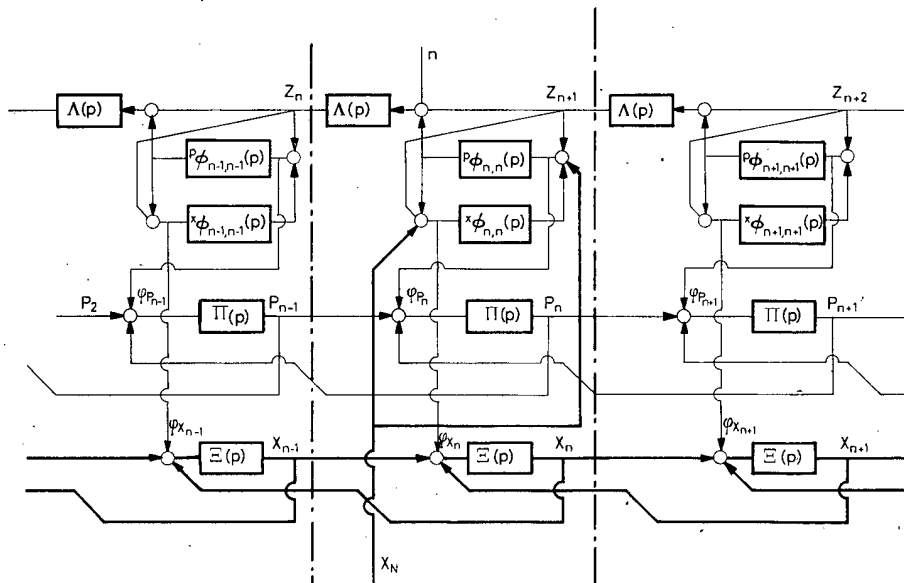


Figure 10

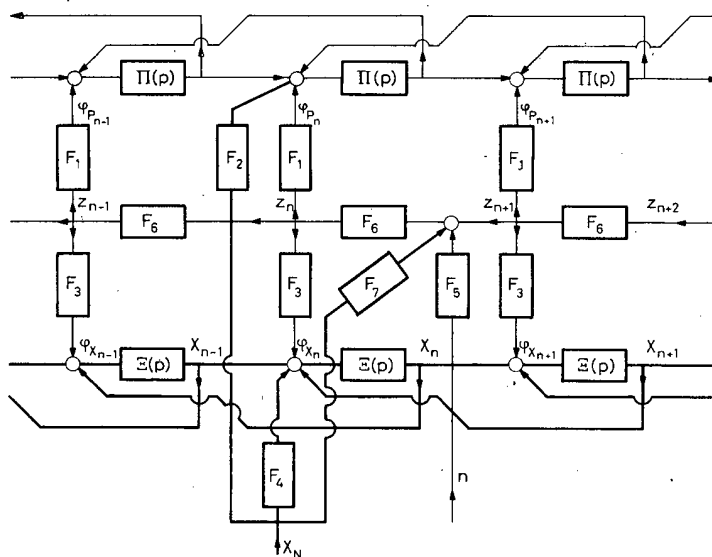


Figure 11

# Some Recent Results in the Computer Control of Energy Systems

T. VAMOS, S. BENEDIKT and M. UZSOKI — *Seen*

The results of automation in the Hungarian energy system were reported in earlier papers<sup>1, 2</sup>. The most important equipment realized is an automatic economic load dispatcher, based on new principles, calculating the effect of the network losses (differing from the usual solutions) from the actual network configuration. The special analogue computer of *Figures 1* and *2* calculates the matrix  $B$  using the well-known formula<sup>3, 4</sup>  $P_v = P B P$ , according to the scheme of *Figure 3*. The potentiometers  $k^g$  in section G of the figure have the values

$$k_i^g = \left| \frac{1 - j \frac{Q_i}{P_i}}{\bar{U}_i} \right|$$

characterizing the generators, where  $Q_i/P_i$  is the active and reactive power rate,  $U_i$  the generator terminal voltage, while

$$k_i^l = \left| \frac{I_i^f}{\sum_{r=1}^m I_i^g} \right|$$

in section L characterizes the proportions of the loads on the total loss, where  $I_{il}$  is the load current at the  $i$ th node;  $I_{ig}$  is the power station current at the  $i$ th node; and  $N$  is the direct current model of the actual network.

By a more detailed analysis it may be proved<sup>5</sup> that if the voltages corresponding to the generator powers  $P_i^g$  (without the network losses and formed similarly to the previous analogue computer solutions) are the driving voltages of the scheme's input, the voltages  $\Sigma B_{ik} P_k$ , being proportional to the network incremental losses, are obtained at the output, which provide with a suitable feedback the optimum load distribution considering the network losses. The papers quoted prove that accuracy of the method (considering the neglects) exceeds the practically reasonable limits.

The analogue single-purpose machine for such a solution, installed in the system control, takes into account, as compared with the former solutions, the active-reactive power proportions of the generator bus-bars, as well as variations in the load proportions.

With the part simulating the actual network (matrix  $N$ ), one obtains an adaptive system, resetting the economic load distribution by switching the elements of network  $H$  (that is, the actual lines, by hand, or automatically with remote control).

The experiences with the automatic load dispatcher have shown that the methods developed up to now for evaluating the economy (the incremental heat rate curves, plotted on statistical bases) are not satisfactory. In connection with this, the following problems have arisen:

(a) Continuous evaluation of the economy characteristics (efficiency, increment costs).

(b) Determination of the estimation periods for the data processing and economic load distribution, permitting filtering of the measurement uncertainties and other short cycle, transitory disturbances, but giving information about the effect of the system variations (e. g. fluctuations in connection with the frequency and power control).

(c) Measurement accuracy corresponding to better calculating and data processing possibilities and improvement of the sensing elements.

(d) Calculation of the transient phenomena effect (e. g. un-load, increase in load) in the automatic load distribution system.

(e) Problems of availability, probability of breakdown, objective judgement of the operation during partial disturbances, or unfavourable service conditions for an automatic dispatcher.

(f) The complex logical decision problems of the automatic energy system dispatcher control for searching the most favourable network connection manipulations.

Among the above problems (a) is generally solved, and a great number of power system data processors are operating. It is worth mentioning, that as regards development, these problems are not resolved. The endeavour for a practically perfect service safety, the complications of practice in connection with the electromechanical output equipment, the reasonable combination of the analogue and digital elements, and the development of a more reliable and cheaper annunciator system, rendering the whole apparatus less expensive, justify numerous new solutions.

Problem (b) must be regarded as the most open one. The digital instruments have generally a class accuracy of 0.1 per cent, while the digital computing technique is practically of absolute accuracy. At the same time the power system measuring and control instruments are of class 1–2 per cent, but in practice instruments and sensing devices for a higher accuracy can be reproduced, but these are accuracy limits under service conditions. Determination of the most important quantities, such as fluid and solid material flows, heat content, ash content, etc., leads to the greatest number of uncertainties, and here the measurement accuracy is 2–5 per cent. The error is increased with the data calculated from such uncertainly measured values, e. g. with the quotient formation necessary for the efficiency. This is the pivotal question and basic contradiction of the whole energetical optimization. We want to attain prospective efficiency improvements of 0.5–1 per cent, the sum of which may be in one country in an integrated average many millions, perhaps many tens of millions, of dollars per year, based on a measurement uncertainty of 1–5 per cent.

Up to now attention has been concentrated on the problem that has not yet been completely solved, i. e. continuous

and more accurate control of the coal heat content value, this being all the more justified in Hungary as the fuel quality varies considerably, and most of the power stations are thermal ones working with coal. Based on some former results<sup>6-8, 11</sup>, a definite improvement has been attained in the field of coal analysis by radioisotopes<sup>9, 10</sup>. We have succeeded in deriving a method permitting the continuous control of the heat content of the coal, at least within the accuracy of the laboratory calorimetric method, which is unacceptable from the statistical sampling point of view. The main experiences were as follows:

(a) In the case of thorough sample preparation a correlation of better than 0.9 may be reached with laboratory calorimetric control, which covers the uncertainty range of the laboratory method measurement accuracy.

(b) With mechanical sample preparation, the measurement accuracy is much influenced.

(c) In the case of considerably varying coal composition, a multi-ray method of different discrete energies can be advised, combined with a suitable, simple computer.

Generally it may be concluded that the next most important step in process automatization will not so much concern the automatic system itself, but rather the development of the quality analysing and quantity measuring devices and sensing elements.

In designing the process optimization a very important and insufficiently considered viewpoint is the determination of the estimation period of the characteristic to be optimized.

This view is especially clear when optimizing the efficiency, as efficiency samples taken for too short a time may lead, for example, to an efficiency value exceeding one after a former storage period, while sampling periods which are too long eliminate the possibilities of estimation of system variations that may be important for optimization. Determination of the ideal sampling period is complicated by the fact that in a boiler the transit and storage time constants of the single energy quantities are extremely different, changing even during the operation. The ideal accurate evaluation of the efficiency could be accepted only for a complete start-operation-stop period. The criterion of duration regarding the estimation period  $\tau$  is that the deviation between the efficiency calculated from the efficiency taken for the total operation time and from that taken for the partial times should not exceed the error  $\varepsilon$ , caused by the measurement inaccuracy, i.e.

$$\eta = \frac{\int_{t_b}^{t_e} x_{out} dt}{\int_{t_k}^{t_v} x_{in} dt} = \frac{1}{n} \sum_i^n \eta_i \pm \varepsilon = \frac{1}{n} \sum_{i=0}^n \frac{\int_{t_i}^{t_{i+1}} x_{out} dt}{\int_{t_i}^{t_{i+1}} x_{in} dt} \pm \varepsilon$$

where  $x_{in}$  = power quantity supplied during the measurement;  $x_{out}$  = power quantity taken out during the measurement;  $t_b$  = initial time of measurement;  $t_e$  = final time of measurement;  $\eta$  = process efficiency;  $\eta_i$  = efficiency of the  $i$ th partial time, and  $n$  = number of samples taken during the whole process.

The estimation period must be chosen as the shortest one meeting this accuracy criterion.

There may be several practical solutions among which the most simple is the working with a time  $\tau$  fixed by experience

on the basis of the above criterion. With the boilers used in Hungary there is an interval of 10–15 min, taking values into consideration only if deviation between the output and input energy levels is less than 3–5 per cent from the beginning to the end of the estimation. The greater variations are, in any case, to be processed separately. The other system adjusts adaptively the evaluation interval on the basis of the auto and cross correlations of the output and input energy characteristic. The numerical results show also, in an apparently entirely identical mode of operation and circumstances, efficiency changes of 2–6 per cent. This is partly due to the considerable variations in the fuel quality. A test made in Czechoslovakia<sup>12</sup> shows variations of  $\pm 10$  per cent for coal quality fluctuations within a very short time. The experiences in Hungary gave similar, or even worse, results, and coal quality fluctuates sometimes by minutes. The effect of the system power and frequency control on the change of efficiency is also most interesting, the load fluctuations having relatively rapid frequencies resulted in an efficiency deterioration of 2–3 per cent in some cases, against the same level steady state operation. The experience in Czechoslovakia justifies the introduction of a corrective control working on the basis of quick coal analysis, while that in Hungary demonstrates the necessity of sensing the effect of the relatively faster changes upon the efficiency.

From the foregoing it follows that the former view of the static load distribution is not satisfactory for calculating the economic load distribution, and the costs of the necessary alterations (heating, unload, switchover, etc.) must be considered.

The problem is clarified by the following example. A power station is operating with four identical boilers, each being loaded to 90 per cent. If the demand increases so that loading of the boilers is to be raised to 100 per cent, the alternative may be considered, i.e. starting a fifth boiler of similar capacity, as a consequence of which the single boilers may operate with 80 per cent load, generally the optimum efficiency level. In this case the expenses of the transients (start, possible later stop, loss of life due to manifold start and stop) must be compared with the savings of the more economical steady-state operation for the expected interval. These circumstances are taken into account already, though in a more simple way, in the present load distribution practice.

The former static load distribution methods are to be generalized to an optimum energetical programming, taking into consideration also the presumable changes. These methods start, as a rule, from Lagrange's method of constrained extrema and are calculated on the basis of the equal incremental costs. The generalized task is the typical case of the multi-step decision problem. On the basis of the power demand given, the system must be programmed in an optimal way, considering afterwards the transition to the power demand expected for subsequent periods and the optimal mode of operation on the new levels. Considering the calculating difficulties and practical demands, the programme was realized for two steps, consequently, besides the system performance level given, the search for the optimum is realized for the next two levels. Consideration of the second change provides information about the first alteration being justified. (In our example the heating up of the new boiler is made reasonable by the time elapsing until the next change and by the direction of the next variation.)

The optimum energetical programming must calculate the availability of the system and its units also, and therefore the probability factors must be considered not only when estimating the power demand, but also when calculating the available power system capacities and network interconnections. In the period to be planned, the system service conditions are characterized by the prospective capacity distributions of the individual units (power stations, machine units, etc.), that is, the probabilities of the available capacities as a function of time, and further, the probability cost values relative to these. These cost values are probability variables not only in the sense that they belong to probable power values, but they are also in themselves only probable values, e.g. the efficiency of the condensation machines is considerably dependent on the cooling water temperature, and consequently on the probable factors of the weather. The third important characteristic, from the optimization point of view, as mentioned earlier, is the excess cost of the transient states, this corresponding not to the expenses integral taken along the static time diagram of the given capacities, but generally exceeding it.

Consequently, for the predictive characterization of the availability, the following are needed. (1) Probability distribution of the capacities depending on the time and direction of transients, (2) the probability cost distributions belonging to these values, and (3) the time integrals of the expense distributions along time tables to be considered.

That is, the availability  $A$  is a set:

$$A = \{p_1 = f_1(P, t); p_2 = f_2(K, P, t); p_3 = f_3(P, K, \int K dt)\}$$

where  $p_1, p_2, p_3$  are the probabilities discussed above.

Accordingly, the availability  $A$  at instant  $t$  is the set of the possible power capacity values  $P$ , where to each value  $P$  belongs a value  $K$  (the costs of the service in steady state conditions) and to each curve  $P_i = P_i(t)$  belongs an integral cost curve  $\int K dt$ .

The task of optimization is as follows. The lines of constraint of the possible power capacities  $P_i = P_i(t)$  are given, that is, the boundary surfaces of a solution space of dimension  $n$ , the lower and upper power limits belonging to the individual units and changing in time. Given the probable system power  $\sum P_i(t) = P_x(t)$ . The trajectory of the vector  $P$  of  $n$  dimensions is to be determined (the vector characterizing the power output condition of the system units), the vector being, under the above conditions,

$$\min \sum_{i=1}^n \int_{t_0}^t K_i [P_{i, \text{opt}}(t)] dt$$

that is, providing the minimum of the cost integral taken along the trajectory. The line integral taken along the trajectory in the coordinate space  $P$  forms no conservative space, as the line integral is not independent of the path and the integrals taken along the closed curves (the cost of returning to the same power distribution) is not zero.

The probability influence of the availability and of the costs have been derived according to the following considerations:

(a) One determines for all equipments the operation time permitted on an experimental base, that may be considered—if there is no special fault indication—as a time of practically perfect safety. During this time the service costs of operating the apparatus in steady-state conditions correspond to the value calculated in general till now.

(b) At the beginning of operation (primary disorders) and over the service time permitted, the probability of outage is greater. Here a penalty tariff is stated, depending on the time and calculating from the former outfall statistics and from the probable economic consequences of the outfall.

(c) Similar penalty tariff is stipulated in case of some error signals (fault indications).

(d) For all important units the transient costs (the expenses of the transient conditions) obtained by experience or calculation are stored, adding to this in some cases the penalty tariff calculated from the disturbance danger relative to the transition.

The above data can be elaborated by the individual power station data processors with a relatively small storage and time requisition to data necessary for the load distribution. These are the curves corresponding to the classical increment cost curves, corrected by the penalty tariffs considering the availability, the possible time functions of the transients and the integral cost curves of the transient conditions. For power stations a relatively slow processing of about 10–20,000 data is needed and the communication of about 300–400 data with the central load dispatcher, as a result of the above calculation. The latter must be dispatched only in case of and to the degree of change. The knowledge of these 300–400 data per power station accomplishes the two-step optimizing programme mentioned earlier.

In this manner, with the aid of suitable power station data processors, by the otherwise available telemetering channels and by a central, medium size computer, energetical automatic optimization may be realized, which takes into account the economic consequences of the power system transient conditions and of its availability, and also the changes in production costs and efficiencies during operation.

The optimum system control referring to the whole power system does not make superfluous the optimization of the individual control circuits, which may be considered partly to be autonomous. Reference is made here, for example, to the control of coal pulverizers, which may be controlled directly by a continuous analyser of ash, assuring the given fuel quantity as a primary condition. As against the non-interacting control systems suggested recently by many authors, installing fixed matrix connections into the control circuits considered previously autonomous, we think to be rather practicable such semi-autonomous adaptive circuits, as the rigid functional connections give suitable results only under perfectly steady-state conditions (e.g. time constants), this condition being chiefly realized with boilers.

In the course of the dispatcher control automatization, the question arose to what extent the dispatcher work may be mechanized in addition to chart preparation and beyond the tasks of the continuous economic load distribution. This idea is supported by the fact that the switching, manipulating and failure suppression activity of the dispatcher control is motivated by subjective factors; extremely hazardous decisions must be made in a short time, and the presence of mind, momentary mood and luck of the dispatcher influence considerably his activity in this field. Mechanization of the task is complicated by the fact that the methods of judgement of the situations were partly subjective ones, based on the experiences and intuitive improvisation capabilities of the dispatcher, as there is no possibility for accurate analysis in the case of rapid decisions.



190/4

Accordingly, mechanization of the dispatcher control permits the application in practice of cybernetics in a narrow sense, and the adoption of recognition and heuristic search.

In the course of elaborating the problem the method of approximating the tasks step by step has been chosen, selecting a single logical task of the dispatcher control. It is seen, taking into account the present machine capacities, that the question arises as to how this task can be elaborated and, after solving this, the kind of further tasks that remain for the dispatcher to solve. By this one can remove from the total dispatcher's activity the parts having not been exactly formulated up to now and examine their weight in the total work and how to handle them. In any case, as more tasks are mechanized and separated from the dispatcher's control, the more time and possibility remain for accomplishing the part demanding the most complicated intellectual activity.

As a first task, estimation of the possible circuit diagram was examined from the overloading point of view. Similar calculations (load flow programmes in the network) have been made regularly for more than a decade on digital machines, but this was the first digital computer application in the power systems. At the same time the methods complying with computer requirements are not fully practicable for automatic control purposes, due to their other demands. Here the analogy of differences between the measuring instruments and the sensing elements of automatic control must be referred to. In sensing systems, however, detecting identical quantities using identical physical principles as measuring instruments the difference in their field of application, demands different approaches. For automatic control we have confined ourselves to a load flow computing method of an accuracy of 5-8 per cent, but being most rapid, providing the results for a 40 node network on a medium size machine in less than 1 sec. The storage capacity demanded is about 1,000 words over the programme. Otherwise the method was a generalization of the well-known method of current distribution factors and imaginary loads, reducing the evaluation of a network of  $n$  nodes to the solution of a complex, linear equation system of about  $i$  unknowns, if the calculated network differs from a basic configuration with  $i$  lines.

The next step was the determination of the optimum connection configuration of the connection manipulations (maintenance, disturbance) tested from a safety point of view. When elaborating the programme, the theory of games has been adopted, interpreting the dispatcher's work as a two-step game of two persons, a game against nature. The pure strategies of one of the players, i.e. the dispatcher, are the available connection manipulations, while those of the other player, that is, the nature, are the disturbances imaginable in the system. The game is two-stepped, first a favourable main network connection diagram is selected by the dispatcher, not knowing yet what kind of disturbances may arise during the validity of this circuit diagram. After this the nature 'moves', a possible disturbance ensues, and as a last step, a changeover strategy is chosen by the dispatcher which reduces the power limitation produced by the disturbance to the minimum. To adopt the minimax principle, a suitable pay-off criterion had to be found by which the elements of the game matrix may be filled and the optimal strategy may be evaluated. This criterion is established on the basis of the damage caused by the possible power outage and the weight functions formed by the outage probability. On the basis of

several considerations, the outage probability  $p$  is not directly applied for weighting, but this is done, however, with the relation

$$k = f(p_i) = \frac{1}{1 - \ln p_i}$$

so the criterion of the optimal game is:

$$\min \left( \max_j W_{ij} K_{ij} \frac{1}{1 - \ln p_j} \right)$$

where  $W_{ij}$  is the power outage caused by the  $i$ th dispatcher's strategy and the  $j$ th disturbance possibility (kWh),  $K_{ij}$  is the specific damage due to the above disturbance (S/kWh), and  $p_j$  is the probability of the  $j$ th disturbance.

The machine time for analysing a complete situation in the case of a medium size machine and of a starting position deviating not from the normal one but at most with the state of the four lines is about 4-5 min for a 40 node network, its storage demand being without programme about 800-1,000 words.

The availability of the network, and that of the power stations, may be considered along similar lines making use of the suggestions mentioned earlier, thus extending further the possibility of the objective evaluation of the network configuration. The programme evaluating the manipulations may include the data referring also to the stability. As examination of the stability conditions of a single situation demands considerable time even by a computer, the application of the pre-calculated, stored stability data, as well as the continuous processing of the data of the stability reserve indicators, are referred to here.

Control of the dispatcher by computers would not make superfluous the application of less complicated network automatics, such as protections, overswitch and backswitch automatics, etc.

It must be emphasized that in the field of the present summarizing report on the authors' developments and ideas, these are up to now mainly theoretical achievements calculated for a mathematical model, prepared for simulation on a digital computer. Their expediency and adaptability must be decided, however, by practice, for many technological and other realization difficulties must be overcome.

## References

- 1 UZSOKI, M. and VAMOS, T. Some questions regarding control of power systems. *Automatic and Remote Control*. 1961. London; Butterworths
- 2 VAMOS, T., UZSOKI, M. and BOROVSZKY, L. Novüj, nyepozredstvennyj, masinnüj szposzob ekonomicsnova raszpregyeljenyija nagruzki mezszdu elektrosztancijami i nyeszkojko voproszov szvjazannüch sz optimizaciej enyergoszisztyem. *Symposium of Automation of Large Energetical Units*. 1961; Prague.
- 3 KRON, G. Tensorial analysis of integrated transmission systems. The six basic reference frames 1. *Trans. Amer. Inst. elect. Engrs*, 70 Pt II. (1951) 1239
- 4 KIRCHMAYER, L. K. *Economic Operation of Power Systems*. 1958. New York; Wiley
- 5 UZSOKI, M. Uj, gépi módszer a gazdaságos teherelosztás számítására. *Colloquium of Automatic Control*. 1962. Budapest
- 6 NAUMOV, A. A. O primenyenii obratnovo rasszejannovo — izlucsenyija dja avtomaticheszkoje kontrolja szosztava szlozsnüch szred. *Avtomaticheszkoje upravlenyije*, pp. 152-159. 1959. Moscow; Izdatyelsztvo Acad. Nauk SSSR

- <sup>7</sup> PIVOVAROV, L. L. O primenyenii javleniya pogloscheniya  $\gamma$ -izlucseniya dlja avtomaticheskovo kontrolja szosztava mnogokomponentnuch szred.
- <sup>8</sup> DIJKSTRA, H. and SIESWERDA, B. S. Apparatus for continuous determination of the ash content of coal. *Int. Coal Prep. Conf.*, 1958; Liège
- <sup>9</sup> BISZTRAY-BALKU, S., Dr. LÉVAI, A., KAKAS, J., NAGY, M. and VARGA, K. Szenek fűtőértékének meghatározása radiológiai módszerrel. *Energia és Atomtechnika*, 6 (1960) 472

- <sup>10</sup> BISZTRAY-BALKU, S., KAKAS, J., NAGY, M., VARGA, K. and LÉVAI, A. Die Bestimmung des Heizwerts von Kohlen durch radioaktive Strahlung. *Isotopentechnik*, Nr. 5-6 (1960-61)
- <sup>11</sup> BELUGOU, P. and CONJÉANUD, P. The determination of the ash content of coals by means of x-rays. *1st Int. Coal Prep. Conf.*, 1950. Paris
- <sup>12</sup> IBLER, J. Trebovanija k regulirovaniju energeticszkich blokov sz tocski zrenija upravlenija energeticszkov szisztyémü. *Symposium of Automation of Large Energetical Units*. 1961. Prague

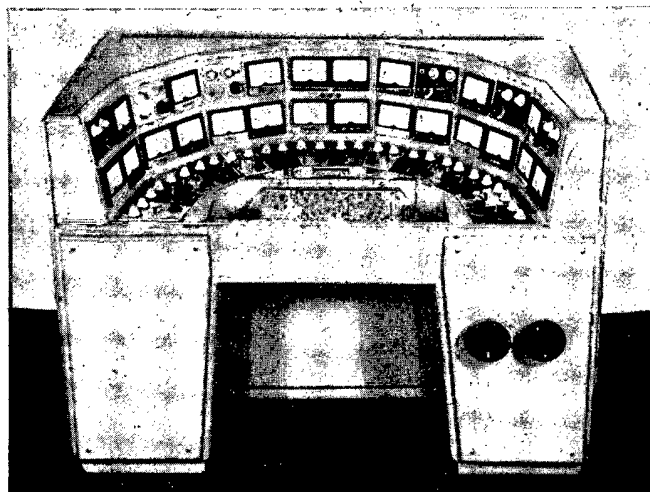


Figure 1

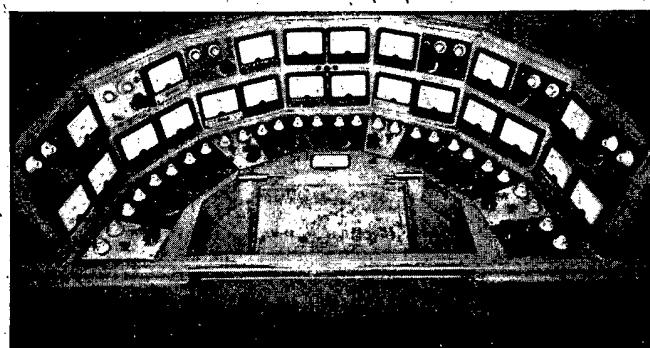


Figure 2

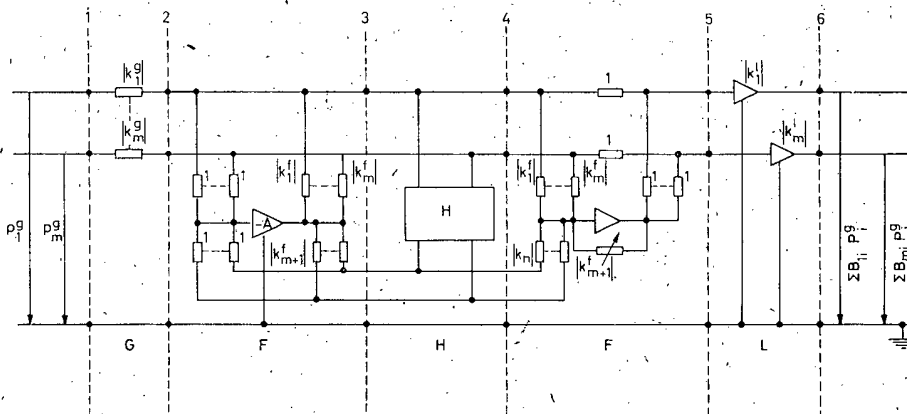


Figure 3

# The Problems, Operation and Calculation of a New Component to be Applied in Certain Control Circuits

O. BENEDIKT *— Hun*

## Introduction

The object of the paper is to describe the physical operation of a new component for the stabilization of oscillatory processes arising in certain control circuits, as well as to give account of a new practical method for calculating the parameters of this component. The component inspires a lively scientific interest, not merely because it can stabilize most effectively an otherwise entirely unstable control circuit in certain cases, but also from a theoretical respect, since there is no need to connect it to the external circuit of the machine. Moreover, without increasing the size of the machine to be stabilized, an effect is realized which up to now could be attained only by a relatively large set consisting of auxiliary devices, with an increase in machine size.

The circuits to be stabilized by the component in question are control circuits, in which the newly developed electrical amplifier 'autodyne' is applied to maintain the load current at a constant value (e.g. for the automatic charging of accumulator batteries, for automatic welding, for supplying motors in series, etc.).

In his paper 'The New Electrical Amplifier', presented at the 1st IFAC Congress the author gave a general report on the theoretical bases, the main application field and control circuit connections of the autodyne, mentioning the autodyne for the above purpose only in short. Csáki, Fekete and Borka, however, referred to the experimental test of another kind of autodyne, namely an autodyne maintaining the output voltage constant.

The following shows in detail the characteristics of the transient phenomena in the autodyne maintaining the load current, as the task and problems of the new stabilizing component of the control circuit of this machine may be understood only in this relation.

## Comparison of the Stability Criteria of the Autodynes Controlling Voltage and Current

The operation of all autodynes working as amplifiers is based, independently of their concrete connection, upon the physical phenomenon that the spatial fundamental harmonic  $\phi_{1\text{res}}$  of the main flux of a converter (Figure 1) may theoretically take up any spatial position (in a different state of equilibrium) with a suitable arrangement of the split poles and a synchronous speed  $n_0$  of the rotor. At the same time, this flux is produced by a magnetizing excitation, the direction of which is set automatically to the flux direction. (Regarding the problems dealt with below, this magnetizing excitation is of no practical importance and therefore is not shown in the figures.) In consequence, with the appearance of a small positive, or negative excitation of  $\pm \Delta AW'$  in the control winding  $W_y$ , the control

torque  $\mp \Delta M$  produced by this excitation and the main flux, and also the small rotor lag or lead caused by the main flux, change considerably the spatial position of the flux  $\phi_{1\text{res}}$ . At the same time, the internal phase voltage  $\bar{U}_{1\text{res}}$ , being in equilibrium with the terminal voltage vector  $\bar{E}_{1\text{res}}$ , can be displaced between the limits of  $\beta = 0$  and  $\beta = 180^\circ$ , while the output voltage  $U$  is varying continuously between the limits  $\pm U_{\text{max}}$ . If the output voltage  $U$  of the amplifier realized, or another control circuit parameter depending on  $U$ , is fed back negatively, then the parameter may be maintained automatically at a constant value. For example in Figure 1 an autodyne is shown stabilizing the output voltage  $U$  to the value of the control voltage  $U_y$ .

In the publications of the USSR Academy of Sciences Technical Section, Energetics and Automatics, No. 2., 1962, the author examined the transient phenomenon taking place in the autodyne controlling voltage (Figure 1), using a different simplifying supposition and neglecting the relatively small rotor resistances.

The characteristic equation of the control circuit, using the operator calculus, yields

$$A' + pB' + p^2C' + p^3D' = 0 \quad (1)$$

As a stability criterion, the following relation is obtained:

$$1 \geq \frac{W_y C_4 (C_1 C_2 n_0 + x_\phi)}{a \cdot r_y \cdot C_1 \cdot C_2} \lambda_5 + \frac{W_y \cdot C_4 \cdot C_6}{a \cdot r_y \cdot C_1 \cdot C_2} x_\phi \quad (2)$$

The quantities  $A'$ ,  $B'$ ,  $C'$ ,  $D'$ ,  $C_1$ ,  $C_2$ ,  $C_4$  and  $C_6$  are constants depending on the machine dimensions.

The physical meaning of these two formulas may be illustrated briefly as follows. Suppose the synchronous speed  $n_0$  of the rotor is decreased to a value  $n$  hardly deviating from  $n_0$  (Figure 2), as a consequence of which the vectors  $\bar{\phi}_{1\text{res}}$  and  $\bar{E}_{1\text{res}}$  rotate by a small angle  $\Delta\beta$  anticlockwise. Meanwhile  $U$  is increased by  $\Delta U$ , and a control current  $\Delta I_y$  arises, producing an excitation  $\Delta AW'$  downwards. The resulting accelerating torque  $\Delta M$  is greatest when the vectors  $\bar{\phi}_{1\text{res}}$  and  $\bar{E}_{1\text{res}}$  reach their dotted upper limit position. At the central position of the two vectors shown by the full line  $\Delta\beta = 0$ . Evidently,  $\Delta AW'$  causes the vectors to oscillate freely around their central position, and such oscillations would appear if the values  $B'$  and  $D'$  in eqn (1) were equal to zero.

Nevertheless, in addition to the voltage  $\Delta U$ , the control winding is effected also by the voltage induced by the increment of the direct axis component  $\Delta\phi_1'$  of the flux  $\bar{\phi}_{1\text{res}}$  (resulting from the rotation of  $\bar{\phi}_{1\text{res}}$ ), which lags the increment  $\Delta\phi_1'$  by  $90^\circ$ . The additional control current being formed evidently establishes

195/2

an excitation upwards when  $\bar{\phi}_{1\text{res}}$  has returned to its central position, thus the excitation has a damping effect. This effect is represented in eqn (1) by the damping term  $p B'$ . The unit of the stabilizing flux  $\Delta \phi_1'$  corresponds to the left-hand side of eqn (2) (the right-hand side is, for the time being, zero).

The voltages induced in the control winding by the fluxes proportional to the current  $\Delta I_y$  must now be considered.

The excitation  $\Delta A W'$  produced by  $\Delta I_y$ , being proportional to it and arising in the rotor winding, is short-circuited by the a.c. network, since according to the converter theory, all excitation arising in addition to the magnetizing excitation is cancelled out almost completely as a consequence of the compensating current  $\Delta I_1'$ . Nevertheless, because of the leakage reactance  $x_\phi$  of the phase winding, the excitation of the current  $\Delta I_1'$  is smaller by few per cent than  $\Delta A W'$ , and consequently a small flux difference  $\Delta \phi_x$  appears. As this, compared with the flux  $\Delta \phi_1'$ , is of opposite direction, it induces a current in the control winding, which reduces the damping effect of the current induced by the flux  $\Delta \phi_1'$ . At the right-hand side of eqn (2) the second term, proportional to  $x_\phi$ , corresponds to the flux  $\Delta \phi_x$ .

In addition to this, the leakage flux  $\Delta \phi_\lambda$  must also be considered. This is caused by the current  $\Delta I_y$ , which passes through the control winding. To this corresponds the first term, proportional to the leakage factor  $\lambda_5$ , at the right-hand side of eqn (2), which proves that the stability degree is now even lower. The effect of fluxes  $\Delta \phi_x$  and  $\Delta \phi_\lambda$  is represented in eqn (1) by the term  $D' p^3$ .

The right-hand side of eqn (2) is proportional to the expression  $W_y/a \cdot r_y$ , where  $W_y$  is the number of turns of the control winding,  $a$  is the number of parallel branches of the control winding, and  $r_y$  is the resistance of this circuit. This expression is obviously proportional to the cross section of one turn of the control winding. The greater this is the greater is the increment of the current  $\Delta I_y$  corresponding to the angle  $\Delta \beta$ , as well as the value of  $\Delta A W'$ , and, evidently, the steady-state control accuracy, also. On the other hand, the value  $\Delta \phi_x + \Delta \phi_\lambda$  is increasing together with  $\Delta I_y$ . However, owing to the fact that at a given value of  $\Delta \beta$ ,  $\Delta \phi_1'$  remains constant, it may be concluded physically—as shown mathematically by eqn (2)—that as the cross section increases, the stability is reduced. If the sum  $\Delta \phi_x + \Delta \phi_\lambda$  were equal to  $\Delta \phi_1'$ , then obviously no voltage would be induced in the control winding and free oscillations would again arise, while the two sides of eqn (2) would be equal.

In practice this never occurs, as a satisfactory control accuracy may be realized by small values of  $W_y/a \cdot r_y$ , at which the stability limit is very great.

The case is quite different with an autodyne used for controlling the load current to a constant value, e.g. in spite of the variation in the internal voltage  $E_A$  of an accumulator (Figure 3). To demonstrate this question theoretically in a more simple way, compare Figure 3 with Figure 1.

At first sight the difference is great. Actually, to the winding  $W'$  (working in this case instead of  $W_y$ ) the loading current  $I$  is fed back, not the voltage  $U$ . Further, in this winding, not two voltages ( $U$  and  $U_y$ ), but two excitations are compared, that is, the excitation  $I W'$  with the excitation  $i_p W_p$  of the continuously controllable regulating current  $i_p$ . Consequently, instead of the law  $U = U_y$  the law  $I = i_p W_p/W'$  is valid here.

However, examining the problem of the transient phenomena, an important analogy of principle may be observed between the

two connections at once, as in this question the magnitude of the current  $i_p$  is practically of no importance and it may be made equal to zero. In this case, however, the connection of Figure 3 does not differ in any respect from that of Figure 1, as  $E_A$  may be regarded as the given control voltage, while the control winding is connected to  $E_A$  and to the voltage  $U$ . In consequence the factors illustrated in Figure 2, affecting the stability, may be distinguished also in the autodyne shown in Figure 3 and if  $\Delta \phi_1' = \Delta \phi_x + \Delta \phi_\lambda$ , free oscillations also arise here. Moreover, it may be seen that in this case the presence of winding  $W_p$  may not practically cause any deviation either, as the fluxes mentioned pass also through this winding and so to the latter, and if the fluxes balance each other mutually, no voltage is induced. Accordingly under the same conditions as have produced eqn (2), a stability criterion corresponding theoretically to eqn (2) must also be obtained. From this, however, follows the interesting fact mentioned below.

While, in the case of Figure 2, the control winding forms a shunt winding, and consequently the cross section of its turns is very small; with an autodyne maintaining its load current at a constant value, the cross section is very large, because the  $W'$  is series connected. This means, however, that the right-hand side of eqn (2) is, in this case, incomparably greater, i.e., there is an actual danger of oscillations arising. This has in fact occurred in practice at an early stage in the development of the autodynes.

It is to be considered that (compared with Figure 2) in the case of Figure 3 the resistance of the rotor may not be neglected with respect to the actually small resistance of winding  $W'$ . As the current  $\Delta I$  must now overcome the resistance of winding  $W'$ , in addition to series-connected resistance  $\Sigma R$ , the effect of  $W_y/a \cdot r_y$  will be somewhat smaller. It is clear, however, that if this term is replaced by  $W'/\Sigma R$ , being physically analogous, the latter will still be incomparably greater, than  $W_y/a \cdot r_y$  in the case of the autodyne controlling its output voltage to a constant value. So it is proved that the autodyne shown in Figure 3 can perform its task only if provision is made for its stability by some supplementary means.

#### Problems Concerning the Development of a Suitable Stabilizing Device and the Way Leading to the Solution

The auxiliary devices for stabilizing circuits, in which the loading current is to be maintained at a constant value, are theoretically known. This is obtained as follows (Figure 4).

Assume the autodyne operates just at the limit of lability, as a consequence of which sinusoidal currents  $\Delta I$  are superposed on the current  $I$ . These would induce sinusoidal voltages in the transformer  $T$ , the primary coil of which is series connected with the load circuit. If the capacity of this voltage is increased by the amplifier  $A$  and the stabilizing winding  $W_c$  is joined to windings  $W_p$  and  $W'$  of Figure 3, with a suitable connection there arises in  $W'$  an excitation leading in time with regard to the excitation  $\Delta I W'$  and proportional to it. In this way effect of fluxes  $\Delta \phi_x + \Delta \phi_\lambda$  could be theoretically reduced by well-known means.

Nevertheless, this arrangement has several great disadvantages. The additional winding increases the machine dimensions. Further, through the application of auxiliary devices, the service safety is reduced. It must also be taken into account that the dimensions of the transformer  $T$  are considerably increased, because its primary coil must be dimensioned for the total

load current  $I$  and saturation of the iron core of the transformer by the load current  $I$  must be avoided. The other stabilizing devices of the classical control technique to be adopted here have similar disadvantages.

Accordingly, it has become essential to seek a novel device for additional stabilization, permitting elimination of the disadvantages mentioned above.

Actually, it has been proved that the physical processes corresponding to *Figure 4* may be realized without the application of a transformer or amplifier, while the required additional winding may be placed in the machine in a way which does not reduce the useful winding area.

This problem is to be solved step by step as follows. (1) In order to spare the primary coil and flux of the transformer  $T$ , instead of this flux another existing flux, already in the autodyne and being proportional to the current  $\Delta I$ , is applied to produce a current in the winding to be placed in the autodyne, playing the role of the secondary coil. This current must lag behind  $\Delta I$ . (2) To amplify the effect of this current, a generated voltage proportional to it is established in the autodyne. (3) To eliminate also winding  $W_s$ , this generated voltage is established in the winding  $W'$  itself, that is, between the main brushes.

Meanwhile, two difficult problems arise. First, (3) obviously necessitates that the winding sought should operate in the direct axis of the autodyne. But then it is inductively coupled with the winding  $W'$ , which eliminates the effect wanted, i.e. only a single suitable additional generated voltage should affect this winding. On the other hand, the following problem arises. If the new winding is placed in the direct axis, then the fluxes  $\Delta \phi_x, \Delta \phi_y$ , being proportional to current  $\Delta I_y$ , will pass through it, and also the flux  $\Delta \phi_1'$ .

It is already known, that free oscillations arise, when the sum of these fluxes is zero, in which case no current is induced in the winding, and therefore the desired effect does not arise at the occurrence of the free oscillations. From this it follows that the tested winding should fulfil the following, apparently contradictory conditions: on the one hand, the magnetic effect of the current arising in it should fall into the direct axis of the machine, but, on the other hand, the direct axis flux  $\Delta \phi_1'$  should not be enclosed by the winding, consequently, the flux enclosed by it and being proportional to the current  $\Delta I$  must not exercise any effect in the direction of the direct axis.

This problem may be solved by a special shape of the tested winding and also by the winding having a particular physical function.

### Physical Operation and Method of Calculation of the Stabilizing Winding

In view of the fact that the autodynes as *Figures 1* and *3* have a practically analogous behaviour regarding the transient phenomena, the following consideration should be valid also in the case shown in *Figure 1*. Therefore, instead of  $\Delta I$  and  $W'$ , the physically similar symbols  $\Delta I_y$  and  $W_y$  shall be applied.

The stabilizing winding, as shown in *Figure 5*, has the shape of a figure eight and is placed, according to *Figure 6*, to the pole shoes of the half-pole I and II belonging to the pole pitch. Thus the condition that they should not be inductively coupled with the winding  $W_y$ , is fulfilled. The condition, that the flux, proportional to the current  $\Delta I_y$  and enclosed by the winding,

does not exert any effect in the direct axis of the machine, may be fulfilled on the basis of the following consideration.

As is known, the compensation current  $\Delta I_1'$ , corresponding to the excitation  $\Delta I_y W_y$ , is proportional to the current  $\Delta I_y$ . In the airgaps below the half-poles the induction of the flux produced by  $\Delta I_1'$  corresponds evidently, within a pole pitch  $\tau_p$ , to the ordinates of curve 1-2-3-4-5-6-7-8-9-10 in *Figure 7*.

If, everywhere, constant inductions of the flux of the same magnitude are represented with the aid of line 1-11-12-13-5-6-14-15-16-10, it becomes apparent that the area 12-3-4-13-12 is equal to the difference of areas 2-11-17-2 and 17-3-12-17. As a result of this, the part of the area 3-4-13-12-3 of the flux proportional to  $\Delta I_y$ , as shown in *Figure 5*, enters the half-pole through one half of the stabilizing winding and leaves on the side of its other half, that is, it is twofold inductively coupled with this coil. Obviously, the situation is just the same in the other half-poles. If the total flux being established is denoted by  $\Delta \phi_8$  and the ordinates of curves 2-4, 7-9 by  $\Delta B(x)$ , then, adopting the above symbols

$$\Delta \phi_8 = K \left[ \begin{array}{l} x = \frac{\tau_p}{4}(1 + \alpha) \\ 2l \int \Delta B(x) dx - \frac{C_5 \cdot \Delta I_1'}{2} \\ x = \frac{\tau_p}{4} \end{array} \right] \quad (3)$$

$$\begin{array}{l} x = \frac{\tau_p}{4}(1 + \alpha) \\ \text{and } l \int \Delta B(x) dx = \frac{C_5 \cdot \Delta I_1'}{2} \\ x = \frac{\tau_p}{4}(1 - \alpha) \end{array} \quad (4)$$

where  $l$  is the active length,  $C_5 \Delta I_1'$  is the flux produced by  $\Delta I_1'$ , and  $K$  is the factor considering the saturation. (The cause of  $C_5$  being constant in spite of the saturation is explained in the paper mentioned previously.)

It follows from eqn (3) and (4) that

$$\Delta \phi_8 = \frac{1 - \cos \frac{\pi \cdot \alpha}{4}}{\sin \frac{\pi \cdot \alpha}{4}} \cdot \frac{C_5 \Delta I_1'}{2} \quad (5)$$

On the other hand

$$\Delta I_1' = K_1 \cdot \Delta I_y W_y \quad (6)$$

where, as is known,  $K_1$  is a constant depending on  $x_p$ . The flux  $\Delta \phi_8$  induces a current of

$$\Delta I_8 = \frac{1}{r_8} \frac{d\Delta \phi_8}{dt} \quad (7)$$

in the stabilizing winding, where  $r_8$  is the resistance of the winding. For the sake of simplicity, the inducing effect of the stray magnetic field of the winding is neglected here. As shown by theory and practice, this is permissible, because the frequency of the free oscillations is insignificant.

Thus, up to now the secondary coil of the transformer  $T$  has been replaced by a winding corresponding to *Figure 5*, while the transformer primary coil and its iron core become

195/4

superfluous. The production of the generated voltage mentioned in (3) between the main brushes of the rotor, will be attempted with the aid of current  $\Delta I_g$ .

At first sight this seems to be impossible. Namely, the excitation  $\Delta \theta_g$  produced by current  $\Delta I_g$  is obviously

$$\Delta \theta_g = \Delta I_g \quad (8)$$

that is, of the same magnitude but of opposite direction in the two half-parts of the figure-of-eight winding. Thus, the excitation is divided on one pole pitch according to the line 15-1-2-3-4-5-6-7-8-9-10-11-12-13-14-17 of *Figure 8*. As the induction  $\Delta B_{g1}$  established by  $\Delta \theta_g$  is distributed practically in a similar way, and as a result of this the total flux arising on one pole pitch yields

$$\Delta \phi_g = \int_{x=0}^{x=\tau_p} \Delta \theta_g dx = 0 \quad (9)$$

it may be concluded that the flux produced by the current  $\Delta I_g$  of the winding cannot produce the generated voltage required in the d. c. winding of the rotor.

The problem can be solved if it is considered that the excitation  $\Delta \theta_g$  must have a positive fundamental harmonic in the case of the distribution in *Figure 8*, because the positive areas 4-5-6-7 and 8-9-10-11 are closer to the central line than the negative areas 1-2-3-4 and 11-12-13-14. The harmonic analysis proves that the amplitude value of the fundamental harmonic  $\Delta \theta_{g1}$  is

$$\Delta \theta_{g1} = \Delta \theta_g \frac{8}{\pi \sqrt{2}} \left( 1 - \cos \frac{\pi \cdot \alpha}{2} \right) \quad (10)$$

Accordingly, this excitation has an inducing effect on the phase winding of the rotor and consequently is practically eliminated by a compensating current  $\Delta I_{g1}$ , because  $\Delta I_{g1}$  is proportional to the excitation,  $\Delta \theta_{g1}$  having produced it, i. e.

$$\Delta I_{g1} = K_2 \Delta \theta_{g1} \quad (11)$$

where  $K_2$  is another constant depending on the value of  $x_\phi$ . The induction produced by the excitation of sinusoidal distribution of this current  $\Delta I_{g1}$  is evidently distributed in the same way as the induction produced by the excitation of the current  $\Delta I_g$ , but in the opposite direction. Consequently, the integral of this induction taken within the section  $\tau_p$  establishes the flux formed by the current  $\Delta I_{g1}$ , which produces the generated voltage

wanted between the main brushes, the latter operating opposite to the voltage induced by  $\Delta \phi_x$  and  $\Delta \phi_\lambda$  in  $W_y$ . Considering the fact that this voltage is proportional to  $\Delta I_{g1}$ , as well as eqn (6), (5), (7), (8), (10) and (11), the generated voltage may be made equal to

$$p \frac{K_8 \cdot \Delta I_y(p)}{r_8}$$

where  $K_8$  is constant. On the other hand, the voltage produced by the fluxes  $\Delta \phi_x$  and  $\Delta \phi_\lambda$  is obtained as  $p D \cdot \Delta I_y(p)$ , where  $D$  is constant. Performing the substitution

$$K'_8 = K_8 (C_1 \cdot n_0 \cdot C_2 + x_\phi) \frac{\theta \cdot \omega}{p_\pi} \quad (12)$$

where  $\omega$  is the angular frequency of the rotor,  $\theta$  is the moment of inertia of the rotor, and  $p_\pi$  is the number of pairs of poles, the characteristic equation is, after all

$$A' + pB' + p^2 \cdot C' + p^3 \cdot \left( D' - \frac{K'_8}{r_8} \right) = 0 \quad (13)$$

The stability criterion:

$$1 \geq \frac{W_y}{a \cdot r_y} \left[ \lambda_5 \frac{C_4 \cdot (C_1 C_2 n_0 + x_\phi)}{C_1 C_2} + x_\phi \cdot \frac{C_4 C_6}{C_1 C_2} \right] - \frac{K''_8}{r_8} \quad (14)$$

where  $K''_8$  is a constant, in which  $W_y$  does not figure. It is recognized that the stabilizing winding is actually in possession of the effect demanded, as for instance the term comprising  $p^3$ , reducing the stability and, in an analogous way, the right-hand side of eqn (14) may be decreased most effectively, if the cross section of the winding, i. e.  $1/r_8$ , is suitably increased. With extremely high values of  $W_y$  the first term of the right-hand side is increased and there is no place in the machine for giving a cross section so large to the stabilizing winding that would suffice for a sensible decrease of the first term. Therefore, in the cases illustrated in *Figure 1*, that is, in the autodyne controlling the voltage to a constant value, this winding has not been applied. Nevertheless, in the cases of *Figure 3*, where the value  $W'$  taking the place of  $W_y$ , is small, calculation shows that the right-hand side of eqn (14) will be zero with such small cross sections, which (with the stabilizing winding set on the pole shoes) has practically no effect upon the machine dimensions.

The autodyne of serial production, provided with such a winding and maintaining the load current at a constant value, would operate without the above-mentioned winding far within the unstable range and prove itself entirely stable in practice.

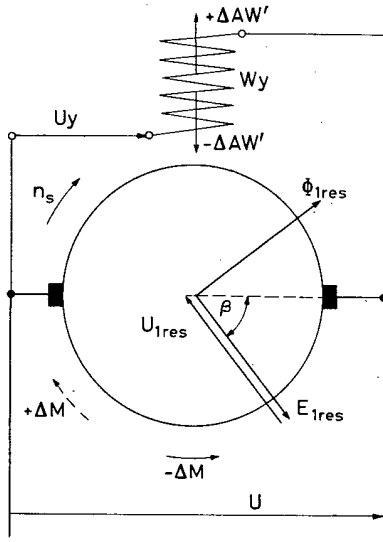


Figure 1.

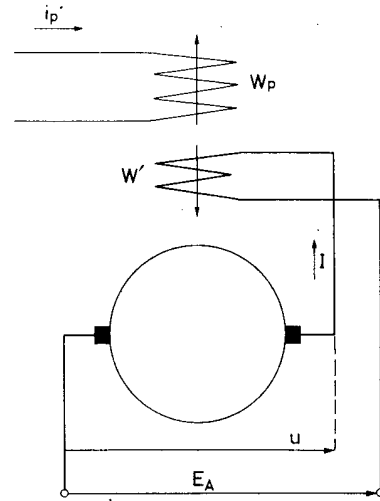


Figure 3.

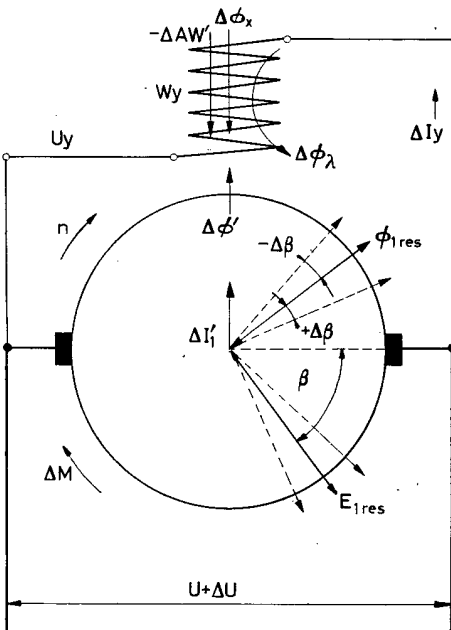


Figure 2.

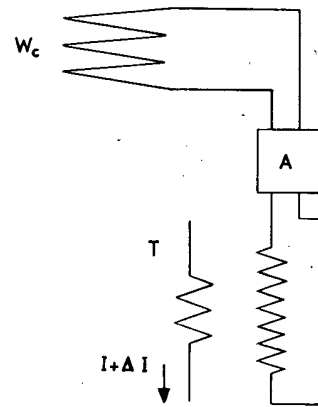


Figure 4.

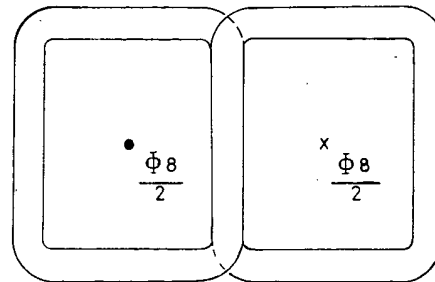


Figure 5.

195/6

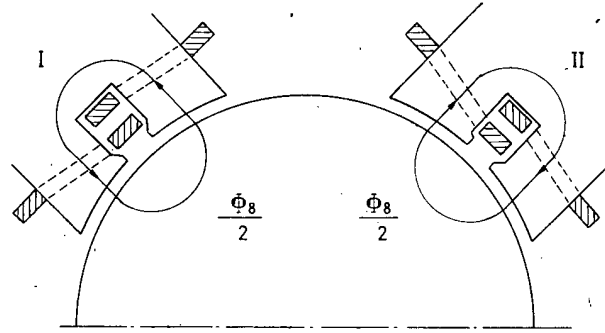


Figure 6.

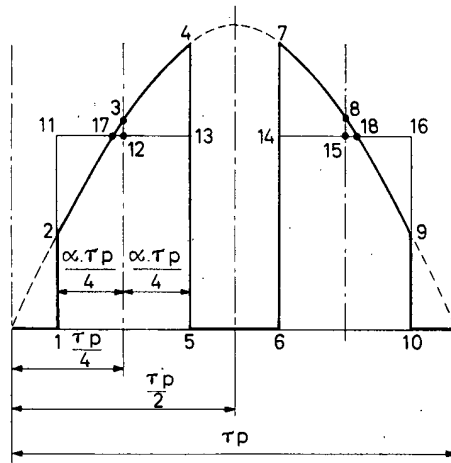


Figure 7.

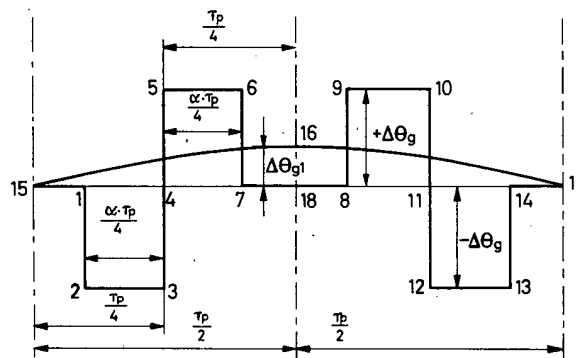


Figure 8.

195/6



# Optimization of Non-linear Random Control Processes

R. KULIKOWSKI\*

- Pal

## Introduction

In the theory of optimum control systems it is usually assumed that the plant differential or operator equations are completely known to the controller. In such cases, through the application of known optimization techniques, the optimum control signal can be derived by an analogue or digital computer and applied to the plant during any time interval. However, there are known systems such as chemical plants and aircraft whose differential equations are not known completely to the controller because of environmental changes, ageing, etc. In many systems of this kind the best that can be accomplished is to construct a multistage optimizing process which converges to the optimum control. All the necessary information for the construction of such a process can be obtained by observing outputs of the plant at every stage for known inputs. Applying this approach to non-linear, zero-memory plants, the gradient of the performance measure can be determined and known iteration methods, based on the gradient concept (such as steepest descent, non-linear programming, contracting iterations, Newton method, etc.) can be applied<sup>1</sup>.

If it is desired to extend these methods for the case of non-linear and random plants having memory (i.e., possessing inertial elements), with the object of obtaining a stable optimizing process, one should first define and determine experimentally the generalized gradient of the performance measure and then construct the convergent iteration process. It will be shown that these problems can be solved successfully, at least in the case of certain classes of non-linear, inertial, random plants, by using some concepts of non-linear and probabilistic functional analysis. However, since the writer realizes that one of the main purposes of a short technical paper is to present the arguments and results in a form which is understandable for the majority of engineers, an attempt has been made to avoid abstract formulations. The more delicate, formal questions are therefore explained in *Remarks I, II, III* which can be omitted during the first reading.

## Assumptions

(1) It will be assumed that the controller generates signals  $x(t)$  which may be subject to certain constraints, such as volume or energy constraints, i. e.

$$\int_0^T |x(t)|^p dt \leq L = \text{const. where } p=1,2 \quad (1)$$

or amplitude constraints, i. e.

$$\max |x(t)| \leq M = \text{const. etc.} \quad (2)$$

\* This research was partially supported by the National Science Foundation under Grant NSFG-14514.

The controller can observe the output  $y(t)$  of the plant for every  $x(t)$  applied to the input by the feedback loop (see *Figure 1*).

(2) We shall also assume that the form of the output-input relation of the controlled system can be described with sufficient accuracy by a non-linear, twice differentiable, integral operator. This operator for example, may be of the polynomial type:

$$y = A(x) = A_0(t) + \sum_{i=1}^{\infty} \int_0^T \dots \int_0^T k_i(t; \tau_1 \dots \tau_i) x(\tau_1) \dots x(\tau_i) d\tau_1 \dots d\tau_i \quad (3)$$

where the kernels  $k_i$  and the function  $A_0(t)$  are generally unknown to the controller.

The differential  $dA(x, h)$  of the operator  $A(x)$ , which is an extension of the usual concept of the differential of a function, can be defined as

$$dA(x, h) = \lim_{\gamma \rightarrow 0} \frac{1}{\gamma} \{A[x(t) + \gamma h(t)] - A[x(t)]\} = \frac{d}{d\gamma} A[x(t) + \gamma h(t)]_{\gamma=0} \quad (4)$$

where  $h(t)$  is an arbitrary function subject to the same constraints as  $x(t)$ . We assume also that it is possible to determine the approximate value of (4) experimentally by observing the outputs of the plant for  $x(t)$  and  $x(t) + \gamma h(t)$  and computing<sup>3</sup>

$$\frac{1}{\gamma} \{A[x(t) + \gamma h(t)] - A[x(t)]\} \approx dA(x, h) \quad (5)$$

where  $\gamma$  is a sufficiently small number.

(3) It is assumed that a performance measure  $F(x)$  is given

$$F(x) = \int_0^T G[x, y, y_d] dt \quad (6)$$

where  $G[x, y, y_d]$  is a known, twice differentiable function of the arguments  $x, y$ .

As an example, consider a chemical plant (for instance, a reactor, distillation column, etc.) described by the positive operator  $A(x)$  [which is non-negative for any  $x(t)$ ]. The amount of steam, fuel or electrical energy delivered to the plant within the time interval  $[0, T]$  will be equal to  $\int_0^T |x(t)|^p dt$ , where  $p = 1$  or  $2$ . The output product obtained in time  $T$  from the plant will be  $\int_0^T A(x) dt$ . Then as the cost of running the plant in the time  $T$ , we can take the performance expression

$$F(x) = \lambda_1 \int_0^T |x(t)|^p dt - \lambda_2 \int_0^T A(x) dt \quad (7)$$

283/2

where  $\lambda_1, \lambda_2 =$  positive coefficients which express the cost in the accepted currency.

As the next example consideration can be given to an autopilot-controller, which minimizes the integral error between the desired  $y_d(t)$  and the actual  $y = A(x)$  path angle of an aircraft:

$$F(x) = \lambda \int_0^{+T} w(t) |y_d(t) - A(x)|^p dt \quad (8)$$

where  $w(t) =$  the given weighting function and  $x(t)$  is subject to the amplitude or integral constraints (2) or (1).

In the case where we want to minimize the final deflection  $\varepsilon(T) = y_d(T) - y(T)$  and its derivatives  $\varepsilon^{(i)}(t)|_{t=T}$ , i.e., when

$$F(x) = \lambda \sum_{i=0}^n \lambda_i \left| \frac{d^i}{dt^i} \varepsilon(t) \right|_{t=T} \quad (9)$$

the weighting function

$$w(t) = \sum_{i=0}^n \lambda_i (-1)^i \delta^{(i)}(T-t) \quad \text{and} \quad p=1$$

should be substituted into (8). The problem becomes more complicated when one wants to minimize the time  $T$  subject to the constraints  $\varepsilon^{(i)}(T) = 0$ , and (1) or (2).

(4) In the general case  $A(x)$  may be a random operator, i.e., for the same  $x(t)$  it may be the case that  $y(t) = A(x)$  is a random function. Therefore, in the performance measures (7), (8), (9) the expected values will be assumed, i.e.  $E\{A(x)\}$  instead of  $A(x)$ .

In many cases  $y_d(t)$  is not known *a priori* and the transient term  $A_0(t)$  caused by the non-zero initial conditions of  $A(x)$  is not known as well. Therefore in the case of (8), (9) it will be assumed that the function  $y_d(t) - A_0(t)$  can be predicted so that it will be known, at least approximately, in the interval  $[0, T]$  and  $A(x)$  will not depend on the initial conditions. In the case of (7) the output  $y(t)$  is the sum of the processes due to  $x(t)$  acting within  $[0, T]$  and  $x(t)$  acting in the past, and this last term contributes to the output.

Now the goal can be formulated, and it is necessary to find a control signal  $x(t)$  which will minimize the performance measure  $F(x)$ . In order to solve this problem one has to determine the conditions of optimality and construct the optimizing process which will converge to the best  $x(t)$ .

*Remark 1.* Speaking more precisely, one wants to minimize the twice-weakly differentiable functional  $F(x)$ , determined on the open or closed sphere of  $L^p[0, T]$  space

$$\|x\| = \left\{ \int_0^T |x(t)|^p dt \right\}^{1/p} \leq R, \quad 1 \leq p$$

The norm  $\|x\|$  in the space  $L^\infty$  should be defined as the so-called 'essential maximum' or

$$\|x\| = \inf_E \left\{ \sup_{t \in [0, T] - E} |x(t)| \right\}, \quad \text{mes } E = 0$$

which is, roughly, equivalent to (2). The functional (9) should be regarded as the so-called Schwartz distribution or generalized

$$dA(x, h) = n \int_0^T k_1(t-\tau) \left[ \int_0^T k_2(\tau-\tau_1) x(\tau_1) d\tau_1 \right]^{n-1} \int_0^T k_2(\tau-\tau_1) h(\tau_1) d\tau_1 d\tau \quad (12)$$

function. The  $\delta^{(i)}(t)$  functions can then be defined as the limits of weakly converging linear functionals.

The concept of a random operator is based on the notion of the so-called generalized random variable<sup>2</sup>. Usually, in control theory, random phenomena are described by random numbers or stochastic processes, which are, roughly, random numbers for any fixed time moments. It is known that the random numbers can be defined in the axiomatic way as the mapping of the space of events into the space of real numbers. It is possible to extend the notion of random numbers to the generalized random variable, which is a Borel measurable mapping of the space of events into some topological or metric space (in our case only, a sphere of  $L^p[0, T]$  space). More precisely<sup>2</sup>, let  $(\Omega, \mathcal{S})$  be a measurable space and  $X$  a non-empty metric space with the  $\sigma$ -algebra  $\mathcal{Z}$  of all Borel subsets of the space  $X$ . Then the mapping  $V$  of the space  $\Omega$  into  $X$  is called a generalized random variable if the inverse image under the mapping  $V$  of each Borel set  $B$  belongs to the  $\sigma$ -algebra  $\mathcal{S}$ , or in symbols, if  $\{\omega : V(\omega) \in B\} : B \in \mathcal{Z} \subset \mathcal{S}$ .

The random operator, which can be denoted by  $A(\omega, x)$ ,  $\omega \in \Omega$ , can be defined as the operator which for every fixed  $x$  is a generalized random variable.

The expected value of  $A(\omega, x)$  can be defined as the Bochner integral over the space  $\Omega$ :

$$E\{A(x)\} \Delta = \int_{\Omega} A(\omega, x) d\mu(\omega)$$

where  $\mu$  is the probability measure, i.e. a non-negative, countable, additive, real set function with the property  $\mu(\Omega) = 1$ . It is assumed that  $E\{A\}$  exists and the expectation sign will be treated as a linear operator acting from the random variable space into the output signal space  $Y$ .

#### Conditions of Optimality

When  $x(t)$  is optimum, any variation  $\gamma h(t)$  of  $x(t)$  should not decrease  $F(x)$ . For example, taking  $G[x, y, y_d] = \lambda x^2(t) - g[y, y_d]$  one can express this condition in the form:

$$\begin{aligned} dF(x, h) &= \frac{d}{dy} F[x + \gamma h] \Big|_{\gamma=0} \\ &= 2\lambda \int_0^T x(t) h(t) dt - \int_0^T \frac{dg}{dy}[y, y_d] \frac{d}{dy} A[x + \gamma h] \Big|_{\gamma=0} dt \\ &= 2\lambda \int_0^T x(t) h(t) dt - \int_0^T g'[y, y_d] dA(x, h) dt = 0 \end{aligned} \quad (10)$$

assuming that the second differential  $d^2/d\gamma^2 F[x + \gamma h]_{\gamma=0}$  is positive for all  $h(t)$ .

It is more convenient to formulate this condition in a form which does not depend upon the arbitrary function  $h(t)$ . If an example is taken of the operator

$$A(x) = \int_0^T k_1(t-\tau) \left[ \int_0^T k_2(\tau-\tau_1) x(\tau_1) d\tau_1 \right]^n d\tau \quad (11)$$

which has the following differential

one can substitute (12) into (10) and by interchanging the integration order there is obtained

$$dF(x, h) = \int_0^T h(t) dt \{2\lambda x(t) - dA^*(x, g')\} = 0$$

where

$$\boxed{\text{Eqn (13)}}^*$$

Then it can be observed that  $dF(x, h) = 0$  for every  $h(t)$  if

$$f(x) = 2\lambda x(t) - dA^*(x, g') = 0 \quad (14)$$

The operator  $f(x)$  will be called the gradient of  $F(x)$ , [ $f(x) = \text{grad } F(x)$ ] because it can be regarded as a generalization of the notion of gradient as commonly thought of in analytic geometry.

When the gradient  $f(x)$  of  $F(x)$  in the neighbourhood of a certain  $x(t) = x_0(t)$  is known, it is possible to express the decrease of  $F(x)$  along the trajectory  $x_0(t) + \alpha [x(t) - x_0(t)]$  (where  $\alpha$  is changing from  $\alpha = 0$  up to  $\alpha = 1$ ) by the mean value

$$\bar{f}(x_0, x) = \int_0^1 d\alpha f\{x_0(t) + \alpha [x(t) - x_0(t)]\}$$

of the gradient along this trajectory. Indeed, one obtains (see Remark II) the following inequality:

$$|F(x_0) - F(x)| \leq \left\{ \int_0^T |\bar{f}(x_0, x)|^q dt \right\}^{1/q} \left\{ \int_0^T |x_0(t) - x(t)|^p dt \right\}^{1/p} \quad (15)$$

which becomes an equality when the two arguments  $\bar{f}(x_0, x)$  and  $x_0(t) - x(t)$  are adjuncts, i.e.

$$|x_0(t) - x(t)|^p = \text{const.} \cdot |\bar{f}(x_0, x)|^q, \quad p^{-1} + q^{-1} = 1$$

Then from (15) it follows that the best change or variation of the control signal should be subordinate to the mean gradient of performance measure.

*Remark II.* It is assumed that the weak differential  $dF(x, h)$  of  $F(x)$  is a linear functional with respect to  $h$  and therefore it can be written in the scalar product form  $dF(x, h) = [f(x), h]$ , where  $h, x \in L^p[0, T]$ ,  $f(x) \in L^q[0, T]$ .

To prove inequality (15) let us observe that for every number  $\alpha \in [0, 1]$  we have

$$\begin{aligned} \frac{d}{d\alpha} F[x_0 + \alpha(x - x_0)] &= dF[x_0 + \alpha(x - x_0), x - x_0] \\ &= \{f[x_0 + \alpha(x - x_0)], x - x_0\} \end{aligned}$$

Integrating this relation we obtain:

$$F(x) - F(x_0) = \int_0^1 \{d\alpha f[x_0 + \alpha(x - x_0)], (x_0 - x)\}$$

Applying the Hölder inequality, the 'maximum principle' is expressed by formula (15).

The necessary condition for a minimum of  $F(x)$  can be

written in the form  $\text{grad } F(x) = \theta$ , where  $\|\theta\| = 0$ ,<sup>4</sup> and for the sufficient condition the following formula is obtained.

$$d^2F(x, h, h) \geq \gamma (\|h\|) \|h\|$$

where  $\gamma(z)$  is a non-negative function having the property  $\lim_{z \rightarrow \infty} \gamma(z) = \infty$ .

In the case of conditional minimums it must be assumed that the functionals are strongly differentiable or, what is equivalent, that the weak differentials are continuous with respect to  $x^4$ . When, for instance, it is required to minimize a certain  $F_1(x)$  subject to the condition  $F_2(x) = c = \text{const.}$ , then, for the necessary condition, the following equation is obtained<sup>4</sup>.

$$\text{grad } F_1(x) = \lambda \text{grad } F_2(x)$$

where  $\lambda$  is a number and, in addition, at point  $x$  one has  $\|\text{grad } F_2(x)\| > 0$ .

In the case when the time  $T$  should be minimized subject to the constraints  $\varepsilon^{(i)}(T) = 0, i = 0, 1, \dots, n$ , in a closed sphere of  $L^p[0, T]$  space  $\|x\| \leq R$ , the problem can be solved in two independent steps:

(1) Fix  $T$  and solve the conditional optimization problem in an open sphere of  $L^p[0, T]$ , by minimizing the functional  $F(x) = \|x\| + \sum_{i=0}^n \lambda_i \varepsilon^{(i)}(T)$ , where  $\lambda_i = \text{constant multipliers}$  determined by the constraints:  $\varepsilon^{(i)}(T) = 0$ .

(2) Assuming that the norm of the solution of (1) depends monotonically on  $T$ , the minimum  $T$  which satisfies the condition  $\|x\| \leq R$  is found.

### Optimizing Processes

When  $A(x)$  is unknown one cannot solve equation (14) and find the best  $x(t)$  in the first interval  $[0, T]$ . But it is sometimes possible to construct an optimizing process  $x_n(t), n = 0, 1, 2, \dots$  in the consecutive intervals  $[nT, (n+1)T]$ , which converges to the best control signal. Consider, for example, the problem of minimizing (8) which is equivalent to the solution of the equation  $y_d(t) - A(x) = 0$ , or the equivalent equation

$$x = x + \kappa [y_d(t) - A(x)] = T(x) \quad (16)$$

where  $\kappa$  is a number. This equation can be solved by the iteration

$$x_{n+1}(t) = T[x_n], \quad n = 0, 1, \dots \quad (17)$$

where  $x_0(t)$  is an arbitrary function, provided the process converges, i.e. the integral distance between  $x_{n+1}$  and  $x_n$  is smaller than the distance between  $x_n$  and  $x_{n-1}$

$$\begin{aligned} & \left\{ \int_0^T |x_{n+1}(t) - x_n(t)|^p dt \right\}^{1/p} \\ &= \left\{ \int_0^T |T(x_n) - T(x_{n-1})|^p dt \right\}^{1/p} \\ &\leq \beta \left\{ \int_0^T |x_n(t) - x_{n-1}(t)|^p dt \right\}^{1/p} \quad (18) \end{aligned}$$

where  $\beta < 1$ .

\* Eqn (13):

$$dA^*(x, g') = n \int_0^T k_2(\tau - t) \left[ \int_0^T k_2(\tau - \tau_1) x(\tau_1) d\tau_1 \right]^{n-1} \int_0^T k_1(\tau_1 - \tau) g'(\tau_1) d\tau_1 d\tau \quad (13)$$

283/4

Assuming that condition (18) is satisfied for every  $x$  (which sometimes can be accomplished by choosing the proper value of  $x$ ) one can construct the sequence of functions  $x_n(t)$  by applying  $x_0(t)$  to the plant, observing  $A(x_0)$ , and computing  $x_1(t) = x_0(t) + x[y_d - A(x_0)]$ , etc. The smaller  $\beta$  is, the faster this process converges to the best  $x(t)$ . Of course, a faster converging optimizing process can be constructed if one has more information about the plant. When the plant changes slowly in time this information can be collected by observing the outputs  $y_i = A(x_i)$  for known inputs  $x_i$  and interpolating the plant operator by the polynomial operator (3), or equalizing the differentials of (3) to the differentials of the plant, determined experimentally. However, one cannot apply this approach in the case when the plant characteristics vary fast in time, because all the information collected in the past becomes obsolete in the future. Therefore it is usually better to use such iteration processes which require a minimum amount of information at every stage of optimization.

In the case when  $A(x)$  is a random operator and it is necessary to find the best  $x(t)$  with respect to the expected value, i.e., if one wants to solve the equation  $x = E\{T(x)\} = S(x)$ , use can be made of an iteration scheme similar to (17) provided  $T(x)$  satisfies certain additional conditions.

*Remark III.* Namely<sup>2</sup>, let  $\chi$  be a separable Banach space and  $T_1, T_2, \dots$  a sequence of weakly independent, weakly equally distributed continuous random operators mapping the Cartesian product  $\Omega \times \chi$  into the space  $\chi$ , i.e.,

$$\mu \left\{ \bigcap_{i=1}^n [\omega: T_i(\omega, x) \in B_i] \right\} = \prod_{i=1}^n \mu[\omega: T_i(\omega, x) \in B_i], \quad B_i \in \mathcal{B},$$

$$\mu[\omega: T_i(\omega, x) \in B] = \mu[\omega: T_k(\omega, x) \in B]$$

and satisfying the conditions:

(a) for every  $x \in \chi$  there exists the Bochner integral:

$$S(x) = \int_{\Omega} T_1(\omega, x) d\mu(\omega)$$

(b) there exists a number  $\beta < 1$  such that for every  $x_1, x_2 \in \chi$  and  $n = 1, 2, \dots$

$$\mu[\omega: \|T_n(\omega, x_1) - T_n(\omega, x_2)\| \leq \beta \|x_1 - x_2\|] = 1 \quad (19)$$

Choosing the generalized random variable  $V_1$  arbitrarily and defining the mapping  $V_{n+1}$  ( $n = 1, 2, \dots$ ) of the space  $\Omega$  into  $\chi$  for every  $\omega \in \Omega$  by the formula  $V_{n+1}(\omega) = S_n[\omega, V_n(\omega)]$ , where the mapping  $S_n$  of the Cartesian product  $\Omega \times \chi$  into  $\chi$  is defined by the formula

$$S_n(\omega, x) = \frac{1}{n} \sum_{i=1}^n T_i(\omega, x) \quad (20)$$

Then there exists a unique point  $\bar{x} \in \chi$ , such that  $S(\bar{x}) = \bar{x}$  and the sequence of generalized random variables converges strongly, almost surely (with probability one), to the fixed point  $\bar{x}$ .

The assumption (19) of this theorem can be relaxed, as was shown by Hans<sup>2</sup>, by assuming

$$\mu[\omega: \|T_n(\omega, x) - T_n(\omega, \bar{x})\| \leq \beta \|x - \bar{x}\|] = 1$$

or

$$\|S(x) - \bar{x}\| \leq \beta \|x - \bar{x}\|, \quad \beta < 1$$

The equation (20) can also be substituted by:

$$\bar{S}_n(\omega, x) = \frac{1}{k_n} \sum_{i=1}^{k_n} T_{j_n+i}(\omega, x)$$

where  $k_1, k_2, \dots, j_1, j_2, \dots$  are two sequences of positive integers:

$$j_1 = 1, j_{n+1} = \sum_{i=1}^n k_i + 1, n = 1, 2, \dots, \sum_{i=1}^{\infty} [1/k_n] < \infty$$

Then each realization of the process can be used only once and the control  $x$  is changed less frequently the further one proceeds.

Now check whether a similar iterative approach can be applied in the more general case of the solution of (14). Assuming that  $dA^*(x, g')$  satisfies condition (18), one can observe that this can be accomplished if we can determine experimentally the functions  $dA^*[x_n, g'(x_n)]$ , for every function  $x_n(t)$  and  $g'(x_n)$ . As an example consider the plant described by (11) for  $n = 1$  and  $k_1(t) = k_2(t) = 0$  for  $t < 0$ . The differentials (12) and (13) become linear subordinate operators, i.e.,

$$dA(x, g') = \int_0^t k(t-\tau) g'(\tau) d\tau,$$

$$dA^*(x, g') = \int_t^T k(\tau-t) g'(\tau) d\tau$$

where

$$k(t) = \int_0^t k_1(t-\tau) k_2(\tau) d\tau, \quad k(t) = 0 \text{ for } t < 0$$

Then it is easy to determine the function  $f^*(t) = dA^*[x, g'(\tau)]$  from  $f(t) = dA[x, g'(T-\tau)]$ , which can be determined experimentally by reversing in time the input  $g'(T-\tau)$  and output  $f(T-t)$ . Indeed,

$$\begin{aligned} f^*(t) &= \int_t^T k(\tau-t) g'(\tau) d\tau \\ &= \int_0^{T-t} k(T-t-\tau) g'(T-t-\tau) d\tau = f(T-t) \end{aligned} \quad (21)$$

In the case of non-linear operators, e.g. (11) for  $n > 1$ , a similar relation holds only for certain types of non-linear operators. Assuming, for example,  $k_1(t) = k_2(t) = k(t) = k(-t)$  and  $f^*(t) = dA^*[x(\tau), g'(\tau)]$ ,  $f(t) = dA[x(T-\tau), g'(T-\tau)]$  it can be proved that the gradient  $f^*(t)$  can be obtained from the differential  $dA[x, g']$  by reversing in time the inputs and outputs; i.e.  $f^*(t) = f(T-t)$ . In the case when  $k(t)$  is symmetrical rather with respect to a certain time instant  $t = T_0$ , than  $t = 0$ , which can be regarded as the delay (or the slope of phase characteristics of the linear parts of the non-linear operator) one can find  $f^*(t)$  in an analogous way from the relation  $f^*(t) = f(T + T_0 - t)$ . The same approach can also be applied to plants described by the sum of:

(1) Linear, delayed by  $T_0$ , operator:

$$\int_0^{t-T_0} k(t-T_0-\tau) x(\tau) d\tau, \quad k(t) = 0, t \leq 0$$

(2) Non-linear, delayed by  $T_0$ , operator:  $\phi[x(-T_0 + t)]$ , where  $\phi(x)$  is a non-linear function.

(3) Non-linear operator of the general type (3) which does

not change when substituting  $t = T + T_0 - t$ ,  $\tau_i = T - \tau_i$ , and interchanging the integration order.

When not sure whether a particular plant belongs to the class for which the gradient can be determined by reversing inputs and outputs, one can test the required property experimentally for every input  $x(t)$ , using the following criterion:

$$\int_0^T h_1(t) dA[x(\tau), h_2(\tau)] dt \\ = \int_0^T h_2(t) d\bar{A}[x(T-\tau), h_1(T-\tau)] dt \quad (22)$$

where  $d\bar{A}[x, h]$  denotes the reversed in time  $dA[x, h]$  operator and  $h_1(t)$ ,  $h_2(t)$  are arbitrary functions.

It can be proved that for plants which satisfy equation (22) and for which the operator  $dA[x, g']$  or  $1/\gamma \{A[x + \gamma g'] - A(x)\}$  satisfies (18) (which can be tested experimentally) the operator  $d^*A[x, g']$  (which is equal to the input-output reversed in time  $dA[x, g']$ ) will also satisfy (18), thus assuring that iteration processes of the type (17) will converge to the best  $x(t)$  when  $\gamma \rightarrow 0$  for  $n \rightarrow \infty$ .

A more general method of identification of  $dA^*[x, g']$  can be constructed using the relation

$$\int_0^T g'(t) dA[x, h] dt = \int_0^T h(t) dA^*[x, g'] dt$$

which connects  $dA[x, h]$  and  $dA^*[x, g']$ .

Indeed, the numbers

$$a_i = \frac{1}{T} \int_0^T h_i(t) dA^*[x_k, g'] dt, \quad k, i = 1, 2, \dots$$

where  $h_i(t)$  are orthogonal, i.e.

$$\frac{1}{T} \int_0^T h_i(t) h_j(t) dt = 0, \quad i \neq j \\ = 1, \quad i = j$$

can be regarded as coefficients of the expansion of the function  $dA^*[x_k, g']$  into the series

$$dA^*[x_k, g'] = \sum_{i=1}^{\infty} a_i h_i(t)$$

Every coefficient  $a_i$  can be written as

$$a_i = \frac{1}{T} \int_0^T g'(t) dA[x_k, h_i] dt \\ = \frac{1}{T} \int_0^T g'(t) \lim_{\gamma \rightarrow 0} \left\{ \frac{A[x_k + \gamma h_i] - A[x_k]}{\gamma} \right\} dt$$

where  $dA[x_k, h_i]$  can be determined experimentally.

By assuming  $h_i t = T\delta(t - t_i)$  it is possible to identify  $dA^*[x_k, g']$  at any desired time moment  $t_i$ .

This method, generally speaking, requires infinite time for complete determination of  $dA^*[x_k, g']$ ,  $k = 1, 2, \dots$ . However, when the orthogonal functions  $h_i(t)$  are properly chosen a few terms of  $a_i h_i(t)$  can provide a good approximation to  $dA^*[x_k, g']$ .

When the output noise is present  $a_i$  are random variables. However it is possible to minimize the corresponding, R.M.S. error or the so-called average risk using Bayes estimates of  $a_i$ .

In this case it is also possible to improve the performance of the controller by collecting and utilizing all the past information about the plant characteristics:  $dA^*[x_k, g']$ .

It should be noted that when the gradient of the performance measure is determined experimentally many other methods, such as steepest descent or Newton generalized process, can also be constructed and applied for the plant optimization.

#### Example

For the sake of simplicity consider the non-random plant described by the operator

$$A(x) = \int_0^t k(t-\tau) x(\tau) d\tau - \varepsilon [x(t)]^2 \quad (23)$$

and the controller which minimizes the cost (7) for  $p = 2$ . The optimal iteration process corresponding to (14) is

$$x_{n+1}(t) = \frac{1}{2\lambda} dA^*[x_n(t), g'(x_n)] = \frac{1}{2\lambda} dA^*[x_n(t), 1(t)] \\ = \frac{1}{2\lambda} \int_t^T k(\tau-t) 1(\tau) d\tau - \frac{\varepsilon}{\lambda} x_n(t) \lambda = \lambda_1/\lambda_2 \quad (24)$$

It can be shown that the plant satisfies (22) and that the process converges if  $\beta = |\varepsilon/\lambda| < 1$ .

$$\left( d^2 F(x, h, h) = 2(\lambda_1 + \varepsilon\lambda_2) \int_0^T h^2(t) dt > 0 \text{ if } \lambda_1 + \varepsilon\lambda_2 > 0 \right)$$

Substituting  $x_0(t) = 0$  into (24) one gets  $x_1(t) = 1/2\lambda dA^*[0, 1(t)]$ . This function can be determined by applying the step function  $\gamma 1(t)$  to the plant and reversing in time the response of the plant, which is multiplied by  $(2\lambda\gamma)^{-1}$ .

For the succeeding iterations we get

$$x_n(t) = x_1(t) \left[ 1 - \frac{\varepsilon}{\lambda} + \left( \frac{\varepsilon}{\lambda} \right)^2 - \dots - \left( \frac{\varepsilon}{\lambda} \right)^{n-1} \right]$$

and

$$x(t) = \lim_{n \rightarrow \infty} x_n(t) = \frac{x_1(t)}{1 + \frac{\varepsilon}{\lambda}} \quad (25)$$

If  $\varepsilon$  were known the best  $x(t)$  could be found for the first interval by (25). When it is not known, or is changing, one can still observe the jumps  $\varepsilon/\lambda \cdot x_{n-1}(0)$ , at the beginning of every interval, and determine the value of  $\varepsilon$ . This optimizing process is shown in Figure 2 for the case  $k(t) = \alpha e^{-\alpha t}$ ,  $\alpha = 3/T$ ,  $\lambda = 1/2$ ,  $\varepsilon = 1/4$ . It is interesting to observe that the optimizing process assumes the scanning form similar to the scanning in the so-called extremum controllers. In the general case one cannot use (25) and in order to find  $x_2(t)$  should reverse in time  $1/2\lambda dA[x_1(T-\varepsilon), 1(\tau)]$ ; a procedure which is shown in Figures 3(a) and (b). One applies  $x_1(T-t)$  to the plant in the first interval and  $x_1(T-t) + \gamma 1(t)$  in the third interval; then finds [see Figure 3(b)]  $1/2\lambda\gamma \{A_I[x_1(T-\tau) + \gamma 1(t)] - A_{III}[x_1(T-\tau)]\}$  and reverses it in time to obtain  $x_2(t)$  etc. In order to utilize all the intervals we can also apply the same  $x_1(T-\tau)$  in the even intervals as shown by dotted line in Figure 3(a). When it is observed that the initial conditions of  $A(x)$  at the beginning of the adjacent intervals are the same (i.e. when the steady-

state process is obtained) [see Figure 3(c)] the step signal  $\gamma 1(t)$  can be applied and the differential  $dA[x_1(T-\tau), 1(\tau)]$  can be determined.

It should be noticed that the number of intervals which are necessary for the determination of the differentials can be reduced in the case when one has two identical plants or a model of the plant, e.g. two chemical reactors or two or more cylinders of the same combustion engine which are subject to the same physical conditions. In this case the differential  $dA[x_n, 1]$  for every  $x_n$  can be determined without the equalization of initial conditions. The scanning period  $T$  should be as short as possible, but not shorter than the settling time of the plant because it would not be possible to determine all the information contained in the transient process. It should be observed that in the general case of performance measure (6)

the determination of the gradient is more complicated because one must reverse in time also  $g'(x_n)$  and in many cases should predict the desired state  $y_d$ .

#### References

- <sup>1</sup> FELDBAUM, A. A. *Computers in Automatic Systems* (in Russian). 1959. Moscow; G.I.F.M.L.
- <sup>2</sup> HANŠ, O. Random fixed points theorem. *Trans. Inst. Prague Conf. on Information Theory, Statistical Decision Functions, Random Processes*. Prague, 1957
- <sup>3</sup> KULIKOWSKI, R. On optimization of time-varying, inertial and nonlinear control systems. *Bull. Acad. Pol. Sci. Ser. Tech. No. 8*, (in Russian). 1961
- <sup>4</sup> WEINBERG, M. M. *Variational Methods in Nonlinear Operators* (in Russian). 1956. Moscow; G.I.T.T.L.

#### Summary

283

There are known controlled systems such as chemical plants or aeroplanes whose differential equations are not known completely to the controller because of environmental changes, ageing, etc. In many systems of this kind the best which can be accomplished is to construct a multistage optimizing process which converges to the optimum control. In the case of zero-memory, non-linear plants processes of this kind can be constructed if the gradient of the performance measure is known or can be determined experimentally. In this paper an extension of this method is considered for the case of non-linear plants having memory and changing randomly in time. In the two introductory

parts of the paper the assumptions and the necessary and sufficient conditions of optimality are formulated. It is shown that at any stage of optimization the best change of the input signal should be adjoint to the mean value of the gradient of the performance measure. Then an optimizing process, based on the so-called fixed point theorem, is constructed. It is shown that for certain classes of non-linear, random plants all the necessary information about the generalized gradient can be obtained experimentally. As an example an optimizing process which minimizes the cost of input energy and maximizes the output gain of a non-linear plant has been constructed and discussed.

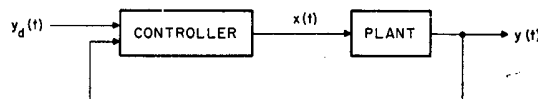


Figure 1. The optimizing control system

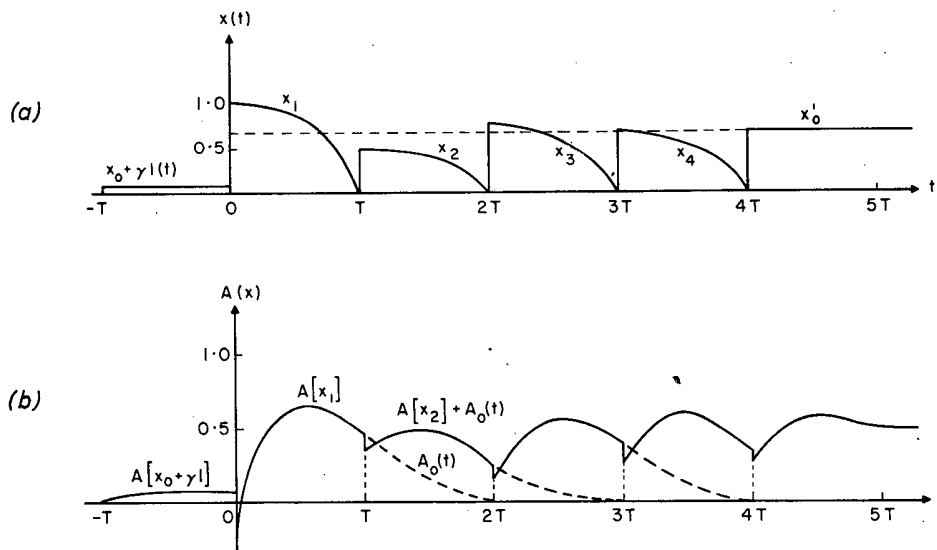


Figure 2. Construction of optimizing processes

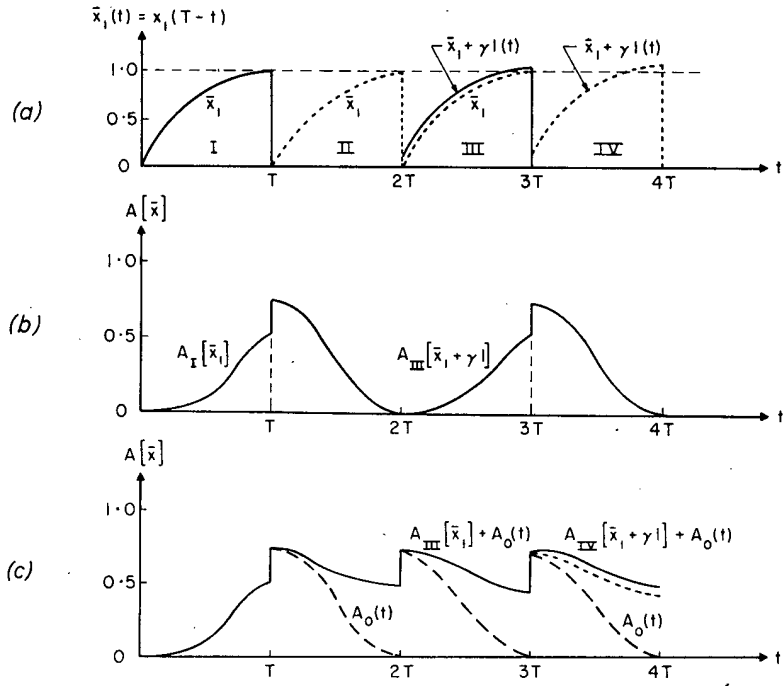


Figure 3. Construction of optimizing processes

# Adaptive Control for a System with a Finite number of States

S. PASHKOVSKII - Pol

## Introduction

In this article a system with incomplete information concerning the medium is considered. Problems of this kind are encountered in engineering, economics, and in systems of mass maintenance. In the systems of control with incomplete information regarding the behaviour of the medium, irregular and inaccurate controlling solutions may be adopted, which results in great losses. In connection with this the development of such an algorithm for controlling solutions, which would rapidly reduce the number of inaccurate and costly solutions, represents an important problem. The system, which will realize this algorithm, may be called the automatic system of control.

## Formulation of Problem

Let  $U = (u_1, u_2, \dots, u_z)$  be the set of actions at one's disposal. These actions can occur only at a certain interval of time. On each step only one action,  $u(n) \in U$ , can occur.

Let  $X^* = (x_1^*, x_2^*, \dots, x_z^*)$  be the set of events, which can be received by receivers  $R$ . The event occurring on the  $n$ th step will be denoted by  $x^*(n)$ .

Let  $X = (x_1, x_2, \dots, x_z)$  be the set of events, which can be received by receivers  $L$ . The event occurring on the  $n$ th step will be denoted by  $x(n)$ .

In addition there is the criterion function  $S(x^*(n), x(n+1))$  which determines for each step the occurring events. This function is represented in the form of *Table 1*.

From *Table 1* it follows that events  $x^*(n) = x(n+1)$  are the events desired. For any other event the 'penalty' represented by number  $r$  must be paid.

Event  $x^*$  may be regarded as the request received on the given step, and event  $x$ , as the realization of that request. When

the realization is identical with the request there are no losses. Otherwise losses  $r$  occur.

Receivers  $R$  receive events  $x^*$  from medium  $A$ . The processes in medium  $A$  which have an effect on received events  $x^*$  may be described only in the form of a probability. In this particular problem it was assumed that  $P(x^* = x_i^*) = p^* = 1/z$  when  $i = 1, 2, 3, \dots, z$ . This means that at each stage any event  $x^* \in X^*$  may occur with an equal probability.

Receivers  $L$  receive events  $x$  from medium  $B$ . The processes of medium  $B$  may be influenced by permissible actions  $u \in U$ . Nothing is known about the mechanism of the effect of processes of medium  $B$  and of adopted action on events  $x$ , except that such an effect does exist. The structure of the mathematical model which will be used for the finding of a connection between the adopted action and the received event  $x(n)$  is represented by the matrix for the probabilities of transition

$$(p_{ij}^k) \quad i, j, k = 1, 2, 3, \dots, z \quad (1)$$

where  $p_{ij}^k = P(x_i(n) \rightarrow x_j(n+1))$  is the probability of occurrence of the event  $x(n+1) = x_j$  when  $x(n) = x_i$ ;  $u(n) = u_k$ . In the problem considered  $p_{ij}^k$  are slowly changing unknown numbers.

The aim is that, under the above-defined conditions, the losses obtained should be at a minimum. This is the general aim of action of an organized system. In this system a stochastic process takes place and, therefore, the mentioned aim should be regarded as the realization of a minimum of mean expected losses. In connection with this the problem of automatic control is to produce on each step such actions for which the mean expected losses will be at a minimum.

The block diagram of an organized system is shown in *Figure 1*. It may be assumed that in the given system there exists a series of actions, which solves the basic problem. This series of actions, for which a minimum of mean expected losses is obtained, will be called the 'decisive (determinative) strategy'. In the given system there may be several decisive strategies. With the change in the operating conditions of the system there is a change in the decisive strategy.

The information regarding the behaviour of the medium is incomplete in the system. Therefore it is impossible to determine directly the decisive strategy. In connection with this a new, additional problem for the automatic control is created. It is, thus, necessary to control the actions in such a way that the decisive strategy is obtained with the minimum of additional losses. This additional problem is reduced to the finding of the decisive strategy.

*Table 1*

$x^*(n)$	$x(n+1)$			
	$x_1$	$x_2$	...	$x_z$
$x_1^*$	0	$r$		$r$
$x_2^*$	$r$	0		$r$
.	.			.
.	.			.
.	.			.
$x_z^*$	$r$	$r$	...	0



284/2

**The Solution of the Problem for Known Probabilities  $p_{ij}^k$** 

First of all, automatic control for the known probabilities  $p_{ij}^k$  will be considered. For this, the algorithm for the working out of the decisive strategy will be determined.

The mean expected losses will be determined as the losses on the  $N$ th step of the path, where  $N$  can be as large a number as desired.

For the determination of the corresponding algorithm the method of dynamic programming will be used. First, the losses over one step will be calculated for the following conditions.

(a) The initial condition for the approaching step is known.

$$\begin{aligned} x^*(n) &= x_i^* \\ x(n) &= x_i \end{aligned}$$

Under these conditions the mathematical expectation for losses over a single step is determined by the formula:

$$v(x_i, x_i^*, u_k) = \sum_{j=1}^z p_{ij}^k S(x_i^*; x_j(n+1)) = (1 - p_{ii}^k) \cdot r \quad (2)$$

(b) The initial condition is given in the form:

$$x(n) = x_i$$

$p_i^* = \frac{1}{z}$  is the probability of occurrence of any  $x_i^* \in X^*$ . Under these conditions the mathematical expectation for the losses over a given step is determined by the formula:

$$v^*(x_i, u_k) = \sum_{i=1}^z p_i^* (1 - p_{ii}^k) \cdot r = \frac{z-1}{z} \cdot r \quad (3)$$

From (3) it is seen that for the unknown initial condition, the mathematical expectation for losses over a step does not depend on the probabilities of transition  $p_{ij}^k$ .

On the basis of the method of dynamic programming it is possible to write down the following equation:

$$\begin{aligned} V_{n+N}(x_i(n), x_i^*(n)) &= \min_{u_k \in U} \left\{ v(x_i, x_i^*, u_k) \right. \\ &\left. + \sum_{h=1}^z p_h^* \sum_{j=1}^z p_{ij}^k V_{n+N-1}[x_j(n+1), x_h^*(n+1)] \right\} \quad (4) \end{aligned}$$

In this equation the events  $x^*(n+1)$ ,  $x^*(n+2)$ , ...,  $x^*(n+N)$  are given only in the form of probability  $p^* = \frac{1}{z}$ . Therefore, the second portion of the right-hand side of eqn (4) may be represented in the form:

$$\sum_{h=1}^z p_h^* \sum_{j=1}^z V_{n+N-1}(x_j(n+1), x_h^*(n+1)) = \frac{z-1}{z} (N-1) r \quad (5)$$

Under these conditions

$$V_{n+N}(x_i(n), x_i^*(n)) = \min_{u_k \in U} \left\{ (1 - p_{ii}^k) \cdot r + \frac{z-1}{z} (N-1) r \right\} \quad (6)$$

From this it is evident that the optimum solution is that solution which gives the minimum value for the mathematical expectation of losses over a step and this is clear instinctively. For the same probability of occurrence of event  $x_i^* \in U$  satisfied over each step, there is no point in planning the actions over  $N$  steps.

On the basis of the given reasoning, the following algorithm for the operation of the decisive strategy is adopted. For each step such actions are adopted, for which the mathematical expectation for step losses is at a minimum.

**Solution of Problem for Unknown Values of  $p_{ij}^k$** 

For unknown probabilities  $p_{ij}^k$  it is impossible to apply the algorithm determined above. It is necessary to develop a new algorithm in order to find a decisive strategy. One of the methods for finding this is to use the information regarding the medium obtained during the time of operation of the system, and the gradual approach to the unknown strategy.

The results, obtained during the time of operation of automatic control, will be represented in the form:

$$\left( \frac{v_{ij}^k}{m_j^k} \right) \quad i, j, k = 1, 2, 3, \dots, z \quad (7)$$

where  $m_i^k$  is the number of adopted actions  $u_k$  for the initial condition  $x_i$ , and  $v_{ij}^k$  is the number of obtained transitions from  $x_i(n)$  to  $x_j(n+1)$ , for  $m_i^k$  experiments.

These results will be used for the determination of unknown probabilities  $p_{ij}^k$ . The determination of the unknown values of  $p_{ij}^k$  will be made by means of reliable intervals. For each value of  $v_{ij}^k/m_i^k$  it is possible to calculate the reliable interval  $(P_{ijH}^k; P_{ijB}^k)$ , where  $P_{ijH}^k$  is the lower limit, and  $P_{ijB}^k$  the upper limit, of the interval.

The limits of intervals may be calculated from known expressions or they can be obtained from tables<sup>2</sup>.

The reliable interval determines the set of the hypothetically possible actual values of  $p_{ij}^k$ . With a high degree of reliability it can be assumed that the actual value of probability  $p_{ij}^k$  will be found in the above-defined interval.

Let the initial condition for the approaching step be:

$$\left. \begin{aligned} x^*(n) &= x_i^* \\ x(n) &= x_i \\ \frac{v_{ii}^k}{m_i^k} &\rightarrow (P_{iiH}^k, P_{iiB}^k) \end{aligned} \right\} \quad (8)$$

The working out of the decisive strategy represents the general problem of the system which consists in the control of actions. From this it follows that for a given initial condition it is necessary to choose action  $u_k \in U$  for which  $p_{ii}^k$  is at a maximum. However, since one knows only the reliable intervals, it is not possible to make a direct choice. In connection with this the following algorithm for the choice of action  $u_k \in U$ , is adopted: to choose such  $u_k \in U$  actions for which, for a given initial condition, there is a hope that probability  $p_{ii}^k$  has the maximum value. This is identical with the method based on the choice of an action, for which there is a hope that the expected losses over a single step will be at a minimum.

It should be pointed out that the upper limit of the reliable interval  $P_{iiB}^k$  when  $k = 1, 2, 3, \dots, z$  represents the basis for the choice of the action. From this it follows that it is necessary to choose such values of  $u_k$ , for which the upper limit of the reliable interval has the maximum value. The result of the action will either confirm the correctness of choice or, in the case of a negative result, decrease the upper limit of the interval, which in

the following intervals gives the possibility for the choice of another value for  $u_k$ . This method guarantees a sufficiently quick convergence of the actions being chosen towards the decisive strategy.

*Example 1*—In the given example the set of events  $X^*$  consists of three events ( $x_1^*$ ,  $x_2^*$ ,  $x_3^*$ ). Medium B is described by means of graphs, shown in *Figure 2*. The results of actions of the system are shown on the graph (*Figure 3*). On this graph the deviations of the actual actions from the decisive strategy are seen. This system was investigated.

The results obtained indicate a rapid convergence of the actions towards the decisive strategy.

*Example 2*—In this example the behaviour of medium B has changed. Medium B is represented by the graph of *Figure 4*. The results obtained are shown in *Figure 5*. In this case also, a rapid convergence towards the sought strategy is obtained,

**References**

- 1 FELDBAUM, A. A. Information storage in closed systems of automatic control. *Izv. Akad. Nauk SSSR, OTN, Energ. i Avt.* 4 (1961)
- 2 YANKO, YA. *Mathematical-statistical Tables*, 1961 ■■■■
- 3 BELLMANN, R. *Dynamic Programming*. 1957. New York; ■■■■
- 4 BUSK, R. and MOSTELLER, F. *Stochastic Models for Learning*. 1955. ■■■■

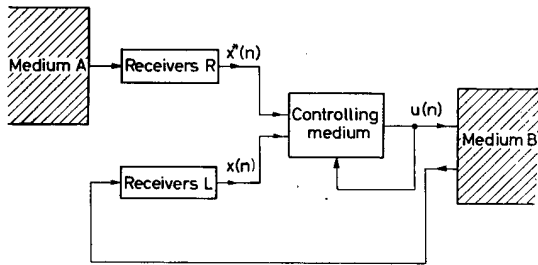


Figure 1

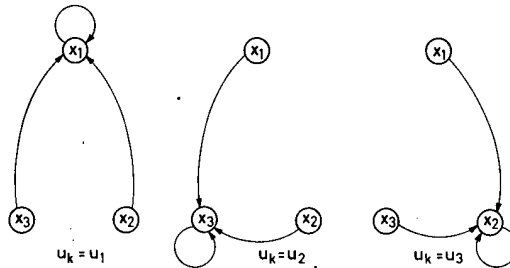


Figure 2

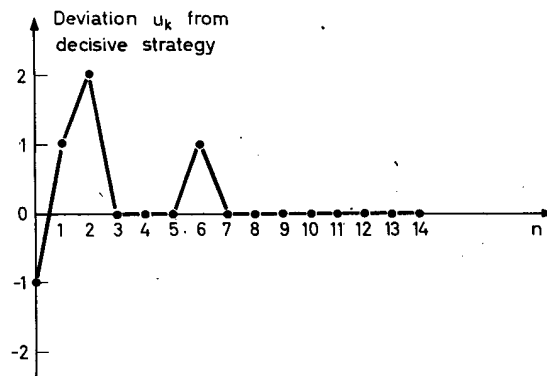


Figure 3

284/4

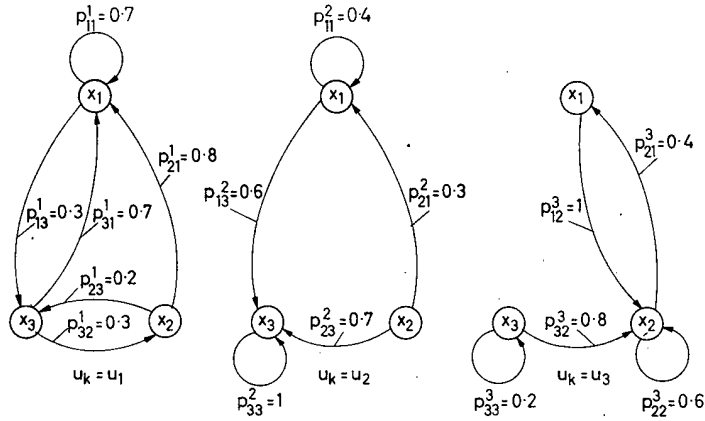


Figure 4

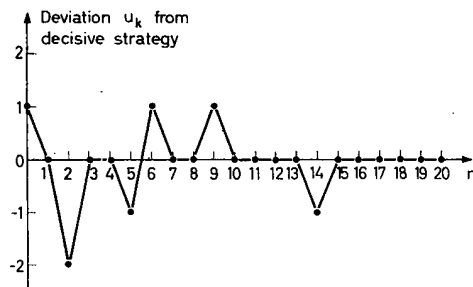


Figure 5

# The Inverse Problem of Integral Square Estimation of Transient Responses

W. JAROMINEK *Pol*

## Introduction

This paper is a continuation of the author's work devoted to the investigation of automatic control systems by means of determinant indices of stability margin<sup>1</sup>. One of the principal problems considered in that reference is the inverse stability problem of linear systems; that is, the problem of obtaining expressions enabling the characteristic equation for prescribed values of the indices of stability margin to be established.

The paper is devoted to the inverse of integral square estimation of transient responses. The inverse integral square estimation may be considered to constitute the inverse problem of quality of a transient process in a linear system. The solution of this problem is obtained by introducing the notion of spectrum of the integral square estimation. The expressions obtained enables the determination, in a unique form, of the transfer functions corresponding to a prescribed qualitative evaluation according to the integral square estimation of transient responses, thus being particularly useful for the synthesis of one- and multi-loop systems.

This work was done under the direction of the Academician B. N. Petrov to whom the author wishes to express his gratitude for many valuable remarks and suggestions.

## 1. The Inverse Stability Problem of Linear Systems

The starting point of the present paper is the inverse stability problem of linear systems. The integral square estimation expressed in the form proposed by A. Krasovskii and called, in what follows, the Krasovskii integral criterion or the Krasovskii evaluation<sup>2, 3</sup>, is the assumed estimation of quality. For the sake of comparability of results the normalized Krasovskii evaluations  $\bar{J}_n^{(m)}$  are considered. One has

$$\bar{J}_n^{(m)} = \bar{J}_n^{(m)} [\bar{F}(p)] \quad (1)$$

$$\bar{F}(p) = \frac{\bar{B}(p)}{\bar{A}(p)} = \frac{b_0 p^m + b_1 p^{m-1} + \dots + b_{m-1} p + 1}{p^n + a_1 p^{n-1} + \dots + a_{n-1} p + 1} \quad (2)$$

where  $\bar{F}(p)$  is the normalized transfer function and  $n > m \geq 0$ .

The consideration of normalized Krasovskii evaluation and normalized transfer functions  $\bar{F}(p)$  does not affect the generality of the assumptions.

It has been shown<sup>1, 4, 5</sup> that the Markov stability criterion enables a solution of the inverse stability problem of linear systems to be obtained. The generalized notion of determinant indices of stability margin has also been introduced<sup>1, 4, 5</sup>; the indices will be denoted by *SMI* (Stability Margin Indices).

The determination of the values of the coefficients of the characteristic equation corresponding to arbitrary values of the *SMI* is obtained according to the developed method<sup>1</sup> by intermediate determination of the Markov parameters. To omit the intermediate stage (the determination and calculation of Markov parameters), which is specially convenient in the case of synthesis of linear systems based on the qualitative Krasovskii's integral criterion, a new method has been developed for establishing characteristic equations, corresponding to any prescribed conditions concerning the *SMI*<sup>6, 7</sup>. It presents a new and independent solution of the inverse stability problem.

## 2. Expansion of the Coefficients of the Characteristic Equation in Terms of *SMI*

The new solution of the inverse stability problem consists in expansion of the coefficients of the characteristic equation in terms of determinant indices of stability margin and, in particular cases, in terms of the Hurwitz or Markov determinants or Routh parameters. As an example of the expansion of the coefficients of the normalized characteristic equation in terms of Hurwitz determinants, mention should be made of *Table 2*, Reference 1. To generalize the results obtained there to the case of any degree 'n' write the characteristic equation  $A_n(p) = 0$  in the following form<sup>7</sup>:

$$A_n(p) = p^n + a_{1,n} p^{n-1} + a_{2,n} p^{n-2} + \dots + a_{k,n} p^{n-k} + \dots + a_{n,n} \quad (3)$$

By considering the sequence of Routh's matrices corresponding to successive values of the degree  $n$  and the equivalent sequence of Hurwitz matrices, it can be shown that the coefficients  $a_{k,n}$  ( $k = 1, 2, \dots, n$ ) of the characteristic polynomial (3) can be expressed in a unique form in terms of Hurwitz determinants<sup>7</sup>. In particular, the following expansions of the coefficients  $a_{k,n}$  are obtained in terms of Hurwitz determinants  $\Delta_{k,n} \equiv \Delta_k$ :

$$\left. \begin{aligned} a_{1,n} &= \frac{\Delta_1}{\Delta_0} \\ a_{2,n} &= \frac{\Delta_2}{\Delta_1} + \sum_{i=1}^{n-2} \frac{\Delta_{i-1}}{\Delta_i} \cdot \frac{\Delta_{i+2}}{\Delta_{i+1}} \\ a_{3,n} &= \frac{\Delta_3}{\Delta_1} + \frac{\Delta_1}{\Delta_0} \sum_{i=2}^{n-2} \frac{\Delta_{i-1}}{\Delta_i} \cdot \frac{\Delta_{i+2}}{\Delta_{i+1}} \end{aligned} \right\} \quad (4)$$

$$a_{2k-1, 2k} = \frac{\Delta_{2k-1}}{\Delta_{2k-2}} + \frac{\Delta_{2k}}{\Delta_{2k-1}} \cdot \sum_{i=1}^{k-1} \frac{\Delta_{2i-1}}{\Delta_{2i-2}} \cdot \frac{\Delta_{2i-1}}{\Delta_{2i}} \quad (5)$$

287/2

$$a_{2k, 2k+1} = \frac{\Delta_{2k}}{\Delta_{2k-1}} + \frac{\Delta_{2k+1}}{\Delta_{2k}} \left( \frac{\Delta_0}{\Delta_1} + \sum_{i=1}^{k-1} \frac{\Delta_{2i}}{\Delta_{2i-1}} \cdot \frac{\Delta_{2i}}{\Delta_{2i+1}} \right) \quad (6)$$

$$a_{n,n} = \frac{\Delta_n}{\Delta_{n-1}} \quad (7)$$

Table 1 has been prepared on the basis of eqns (4)-(7).

In the general case the expansion of the coefficients in terms of Hurwitz determinants is expressed by the following algorithm<sup>7</sup>:

$$a_{k,n} = a_{k,n-1} + \frac{\Delta_{n-3}\Delta_n}{\Delta_{n-2}\Delta_{n-1}} a_{k-2,n-2} \quad (8)$$

where  $a_{i,s} \equiv 1$  in the case of  $i = 0$  and  $0 \leq s \leq n$   
 $a_{\alpha,\alpha} \equiv 1$  in the case of  $\alpha = 0, -1, -2, \dots$   
 $a_{i,s} \equiv 0$  in the case of  $i < 0$  or  $i > s$

The expressions for the expansion of the coefficients  $a_{k,n}$  can be most easily obtained by means of the recurrence equations

$$A_k(p) = pA_{k-1}(p) + \frac{\Delta_{k-3}\Delta_k}{\Delta_{k-2}\Delta_{k-1}} A_{k-2}(p) \quad (9)$$

where  $A_k(p)$  is a polynomial of degree  $k$  and  $\Delta_i \equiv 1$  for  $i = 0, -1, -2, \dots$ . The recurrence equation (9) holds for  $1 \leq k \leq n$  (Reference 7). In particular, in the case of  $k = 1, 2$  one obtains  $A_0(p) = A_{-1}(p) \equiv 1$ .

Analogous expressions may be derived for the remaining forms of the *SMI*<sup>7</sup>.

Equations (4)-(9) enable the inverse stability problem of linear systems to be easily solved. The selection of appropriate values of the *SMI* should be done on the basis of a suitable qualitative criterion of transient responses.

### 3. The Transformed Krasovskii Integral Criterion

As an estimation of quality of transient responses assume the Krasovskii integral criterion  $\bar{J}_n^{(m)}$ . It has been shown<sup>1, 5, 8</sup> that as a result of a suitable transformation the integral square estimation  $\bar{J}_n^{(0)}$  can be expressed in a simple manner in terms of the indices of stability margin. The transformed evaluation  $\bar{J}_n^{(0)}$  takes, when Hurwitz determinants are used, the form

$$\bar{J}_{2k}^{(0)} = \frac{1}{2} \left( \frac{a_{2k-1}}{a_{2k}} + \frac{a_0}{a_1} + \frac{\Delta_2\Delta_2}{\Delta_1\Delta_3} + \frac{\Delta_4\Delta_4}{\Delta_3\Delta_5} + \dots + \frac{\Delta_{2(k-1)}\Delta_{2(k-1)}}{\Delta_{2k-3}\Delta_{2k-1}} \right) \quad (10)$$

$$= \frac{1}{2} \left( \frac{a_{2k-1}}{a_{2k}} + \frac{\Delta_0}{\Delta_1} + \sum_{i=2}^k \frac{\Delta_{2(i-1)}}{\Delta_{2i-3}} \cdot \frac{\Delta_{2(i-1)}}{\Delta_{2i-1}} \right)$$

$$\bar{J}_{2k+1}^{(0)} = \frac{1}{2} \left( \frac{a_{2k}}{a_{2k+1}} + \frac{\Delta_1\Delta_1}{\Delta_0\Delta_2} + \frac{\Delta_3\Delta_3}{\Delta_2\Delta_4} + \frac{\Delta_5\Delta_5}{\Delta_4\Delta_6} + \dots + \frac{\Delta_{2k-1}\Delta_{2k-1}}{\Delta_{2k-2}\Delta_{2k}} \right)$$

$$= \frac{1}{2} \left( \frac{a_{2k}}{a_{2k+1}} + \sum_{i=1}^k \frac{\Delta_{2i-1}}{\Delta_{2i-2}} \cdot \frac{\Delta_{2i-1}}{\Delta_{2i}} \right) \quad (11)$$

In order to obtain a complete transformation of the evaluation  $\bar{J}_n^{(0)}$  make use of the relations (5)-(7) and express the ratios  $a_n - 1/a_n$  to the coefficients of the characteristic equation  $A_n(p)$  in terms of Hurwitz determinants. It is found that

$$\frac{a_{2k-1}}{a_{2k}} = \sum_{i=1}^k \frac{\Delta_{2i-1}}{\Delta_{2i-2}} \cdot \frac{\Delta_{2i-1}}{\Delta_{2i}} \quad (12)$$

$$\frac{a_{2k}}{a_{2k+1}} = \frac{\Delta_0}{\Delta_1} + \sum_{i=1}^k \frac{\Delta_{2i}}{\Delta_{2i-1}} \cdot \frac{\Delta_{2i}}{\Delta_{2i+1}} \quad (13)$$

Observe that the evaluations (10) and (11) have different forms in case of even ( $n = 2k$ ) and odd ( $n = 2k + 1$ ) degrees  $n$  of the characteristic equation. Substituting in (10) the expression (12) and in (11) the expression (13) and performing an appropriate change of summation indices one obtains a single general expression

$$\bar{J}_n^{(0)} = \frac{1}{2} \sum_{i=1}^n \frac{\Delta_{i-1}}{\Delta_{i-2}} \cdot \frac{\Delta_{i-1}}{\Delta_i} \quad (14)$$

where  $\Delta_0 \equiv \Delta_{-1} \equiv 1$  is arbitrarily assumed.

The expression (14) is the transformed Krasovskii evaluation expressed exclusively in terms of indices of stability margin\*. It holds for both even and odd degrees; that is, for any degree  $n$  of the characteristic equation.

The above transformation of the Krasovskii evaluation may be considered as a transition from one set of independent variables to another. The independent variables of the first set are the coefficients of transfer function; those of the other, the *SMI*. Further investigations show that this transformation is of essential importance chiefly because the *SMI* furnish much more necessary information on the control system than the transfer function coefficients. It is also of importance that the new expressions of the integral square estimation take a much simpler analytical form, which is essential for the synthesis of control systems.

To generalize the results obtained to systems of the non-zero class ( $m \neq 0$ ;  $n > m > 0$ ) consider some of the relations between the Krasovskii determinants and the *SMI*.

### 4. Expansion of the Krasovskii Determinants in Terms of *SMI*

In the general case the normalized<sup>†</sup> integral square estimation takes the form

$$\bar{J}_n^{(m)} = \frac{1}{2} \sum_{\alpha=0}^m B_{m-\alpha}^{(m)} \cdot \frac{\Delta_{m-\alpha}^{(n)}}{\Delta_n} - b_{m-1} \quad (15)$$

where

$$B_{m-\alpha}^{(m)} = b_{m-\alpha}^2 - 2b_{m-\alpha+1}b_{m-\alpha-1} + 2b_{m-\alpha+2}b_{m-\alpha-2} + \dots + 2(-1)^{m-\alpha}b_m b_{m-2\alpha} \quad (16)$$

for  $\alpha = 0, 1, 2, \dots, m$  and  $b_m \equiv 1$ ;  $b_k \equiv 0$  ( $k < 0$ ;  $k > m$ ).

The expression of the normalized evaluation (15) in terms of *SMI* requires, above all, the expansion of the Krasovskii determinants  $\Delta_{m-\alpha}^{(n)}/\Delta_n$  in terms of *SMI*<sup>7, 8</sup>. The elements of these determinants are exclusively the coefficients of the characteristic equation, therefore the unique expansion  $\Delta_{m-\alpha}^{(n)}/\Delta_n$  in terms of *SMI* may be done on the basis of eqns (4)-(9). As an example a few of the relations obtained are quoted:

\* Other forms of the transformed Krasovskii evaluation may be found in Reference 7.

† The normalized evaluation  $\bar{J}_n^{(m)}$  corresponds to the normalized transfer function  $(\bar{F})_p$ , for which one has  $a_0 = a_n = b_m \equiv 1$ .

Table 1. Expansion of the coefficients of the characteristic equation in terms of indices of stability margin of the group H (the general case)

Coefficients of the characteristic equation							
$n$	$a_0$	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$
1	1	$a_1 = \frac{\Delta_1}{\Delta_0}$					
2	1	$a_1 = \frac{\Delta_1}{\Delta_0}$	$a_2 = \frac{\Delta_2}{\Delta_1}$				
3	1	$a_1 = \frac{\Delta_1}{\Delta_0}$	$a_2 = \frac{\Delta_2}{\Delta_1} + \frac{\Delta_0 \Delta_3}{\Delta_1 \Delta_2}$	$a_3 = \frac{\Delta_3}{\Delta_2}$			
4	1	$a_1 = \frac{\Delta_1}{\Delta_0}$	$a_2 = \frac{\Delta_2}{\Delta_1} + \frac{\Delta_0 \Delta_3}{\Delta_1 \Delta_2} + \frac{\Delta_1 \Delta_4}{\Delta_2 \Delta_3}$	$a_3 = \frac{\Delta_3}{\Delta_2} + \frac{\Delta_1}{\Delta_0} \cdot \frac{\Delta_1 \Delta_4}{\Delta_2 \Delta_3}$	$a_4 = \frac{\Delta_4}{\Delta_3}$		
5	1	$a_1 = \frac{\Delta_1}{\Delta_0}$	$a_2 = \frac{\Delta_2}{\Delta_1} + \frac{\Delta_0 \Delta_3}{\Delta_1 \Delta_2} + \frac{\Delta_1 \Delta_4}{\Delta_2 \Delta_3} + \frac{\Delta_2 \Delta_5}{\Delta_3 \Delta_4}$	$a_3 = \frac{\Delta_3}{\Delta_2} + \frac{\Delta_1}{\Delta_0} \left( \frac{\Delta_1 \Delta_4}{\Delta_2 \Delta_3} + \frac{\Delta_2 \Delta_5}{\Delta_3 \Delta_4} \right)$	$a_4 = \frac{\Delta_4}{\Delta_3} + \frac{\Delta_5}{\Delta_4} \left( \frac{\Delta_0}{\Delta_1} + \frac{\Delta_2 \Delta_2}{\Delta_1 \Delta_3} \right)$	$a_5 = \frac{\Delta_5}{\Delta_4}$	
6	1	$a_1 = \frac{\Delta_1}{\Delta_0}$	$a_2 = \frac{\Delta_2}{\Delta_1} + \frac{\Delta_0 \Delta_3}{\Delta_1 \Delta_2} + \frac{\Delta_1 \Delta_4}{\Delta_2 \Delta_3} + \frac{\Delta_2 \Delta_5}{\Delta_3 \Delta_4} + \frac{\Delta_3 \Delta_6}{\Delta_4 \Delta_5}$	$a_3 = \frac{\Delta_3}{\Delta_2} + \frac{\Delta_1}{\Delta_0} \left( \frac{\Delta_1 \Delta_4}{\Delta_2 \Delta_3} + \frac{\Delta_2 \Delta_5}{\Delta_3 \Delta_4} + \frac{\Delta_3 \Delta_6}{\Delta_4 \Delta_5} \right)$	$a_4 = \frac{\Delta_4}{\Delta_3} + \frac{\Delta_5}{\Delta_4} \left( \frac{\Delta_0}{\Delta_1} + \frac{\Delta_2 \Delta_2}{\Delta_1 \Delta_3} + \frac{\Delta_3 \Delta_2}{\Delta_1 \Delta_4} + \frac{\Delta_0 \Delta_3 \Delta_3}{\Delta_1 \Delta_2 \Delta_4} \right)$	$a_5 = \frac{\Delta_5}{\Delta_4} + \frac{\Delta_6}{\Delta_5} \left( \frac{\Delta_1 \Delta_1}{\Delta_0 \Delta_2} + \frac{\Delta_3 \Delta_3}{\Delta_2 \Delta_4} \right)$	$a_6 = \frac{\Delta_6}{\Delta_5}$

Notes: (1) In the case of normalized characteristic equation one should substitute  $a_n = \frac{\Delta_n}{\Delta_{n-1}} \equiv 1$ .

(2)  $n$  is the degree of the characteristic equation.

(3)  $\Delta_0 \equiv 1$  (assumed arbitrarily).

$$\frac{\Delta_m^{(n)}}{\Delta_n} = \sum_{i=1}^n \frac{\Delta_{i-1}}{\Delta_{i-2}} \cdot \frac{\Delta_{i-1}}{\Delta_i} \quad (17)$$

$$\frac{\Delta_{m-1}^{(n)}}{\Delta_n} = \frac{\Delta_{n-2}}{\Delta_{n-1}} \quad (18)$$

$$\frac{\Delta_{m-2}^{(n)}}{\Delta_n} = \frac{\Delta_{n-3}}{\Delta_{n-1}} \quad (19)$$

$$\frac{\Delta_{m-3}^{(n)}}{\Delta_n} = \frac{\Delta_{n-4}}{\Delta_{n-2}} + \frac{\Delta_{n-3}}{\Delta_{n-2}} \frac{\Delta_{n-3}}{\Delta_{n-1}} \quad (20)$$

.....

The expansions of the remaining expressions of the form

$$\frac{\Delta_{m-\alpha}^{(n)}}{\Delta_n}$$

can be represented in a similar manner. Consider the sequence

$$\left\{ \frac{\Delta_{m-\alpha}^{(n)}}{\Delta_n} \right\}$$

of these expressions:

$$\frac{\Delta_m^{(n)}}{\Delta_n}, \frac{\Delta_{m-1}^{(n)}}{\Delta_n}, \frac{\Delta_{m-2}^{(n)}}{\Delta_n}, \dots, \frac{\Delta_{m-\alpha}^{(n)}}{\Delta_n}, \dots, \frac{\Delta_{m-m}^{(n)}}{\Delta_n} \quad (21)$$

It can be shown that the structure of the equation obtained for the expansion of each particular expression

$$\frac{\Delta_{m-\alpha}^{(n)}}{\Delta_n},$$

in terms of the *SMI* is independent of the degrees  $n$  and  $m$  of the transfer function polynomials and depends only the ordinal number  $\alpha$  in the sequence

$$\left\{ \frac{\Delta_{m-\alpha}^{(n)}}{\Delta_n} \right\}$$

This property is very useful for the generalization of the considered problem for the case of  $m > 0$ .

**5. The Optimum Integral Estimation  $\bar{J}_n^{(m)}$  in the Sense of *SMI***

The value of the evaluation (15) depends on the distribution of poles and zeros of the transfer function (2). Assume that in the general case the distribution of the zeros is independent of that of the poles. Then, the coefficients  $B^{(m)}_{m-\alpha}$  are also independent of the coefficients of the characteristic equation and cannot, in general, be expressed in terms of *SMI*. For any assigned distribution of transfer function poles there exists only one distribution of zeros of the polynomial  $\bar{B}(p)$  in the numerator of the transfer function, which, for the given assumptions, corresponds to the minimum value of the evaluation  $\bar{J}_n^{(m)}$ . Such a distribution of zeros will be called, in what follows, optimum in relation to the *SMI*. The determination of the corresponding optimum polynomial  $\bar{B}(p) = \bar{B}(p)_{opt}$  will be called the optimization of the integral square estimation  $\bar{J}_n^{(m)}$  in the sense of *SMI*<sup>7</sup>.

The determination of the values of the coefficients of the optimum polynomial  $\bar{B}(p)_{opt}$  reduces to that of the extremum (minimum) value of a function of many independent variables. To do this one must equal to zero the partial derivatives of the evaluation  $\bar{J}_n^{(m)}$  with respect to the coefficients of the polynomial  $\bar{B}(p)$ . One has

$$\frac{\partial \bar{J}_n^{(m)}}{\partial b_0} = 0; \frac{\partial \bar{J}_n^{(m)}}{\partial b_1} = 0; \dots; \frac{\partial \bar{J}_n^{(m)}}{\partial b_k} = 0; \dots; \frac{\partial \bar{J}_n^{(m)}}{\partial b_{m-1}} = 0 \quad (22)$$

Solving successively for each  $m$  the system of  $m$  equations  $\partial \bar{J}_n^{(m)} / \partial b_k \equiv 0$  ( $k = 0, 1, 2, \dots, m - 1$ ) the values of the coefficients of the polynomial  $\bar{B}(p)_{opt}$  will be obtained, i.e., the optimum in the sense of *SMI*. Thus, for instance, in the case of  $m = 3$  one has:

$$b_0 = \frac{\Delta_{n-3}}{\Delta_{n-4}}, b_1 = \frac{\Delta_{n-2}}{\Delta_{n-3}}, b_2 = \frac{\Delta_{n-1}}{\Delta_{n-2}} + \frac{\Delta_{n-3}}{\Delta_{n-4}} \frac{\Delta_{n-3}}{\Delta_{n-2}}, b_3 = 1 \quad (23)$$

Table 2 contains expressions for the coefficients of optimum polynomials  $\bar{B}(p)_{opt}$ , obtained as a result of solution of the set of eqns (22) for a few successive values of the degree  $m$  of the polynomial  $\bar{B}(p)$ .

The integral evaluation  $\bar{J}_n^{(m)}$ , that satisfies the set of conditions (22), will be called optimum in the sense of stability margin and denoted by  $\bar{J}_n^{(m)}_{opt}$ . Optimum evaluations in the sense of *SMI*, have a number of valuable properties. Some of them will be considered below. Of particular importance is the fact that for full analytical description of the evaluation  $\bar{J}_n^{(m)}_{opt}$  the *SMI* are required only.

**6. The Two Equivalent Forms of the Integral Evaluation  $\bar{J}_n^{(m)}$**

In the general case the integral evaluation  $\bar{J}_n^{(m)}$  does not satisfy the optimum conditions (22), and therefore it cannot be expressed in terms of the *SMI* only. This follows directly from the assumption, that the coefficients of the polynomial  $\bar{B}(p)$  are independent of the coefficients of the characteristic equation  $\bar{A}(p)$ . In this connection try to separate in the integral evaluation  $\bar{J}_n^{(m)}$  a component depending exclusively on the *SMI* from another component in which the influence of the polynomial  $\bar{B}(p)$  is taken into account. The introduction of the *SMI* and the notion of optimum conditions in the sense of *SMI* enables two new equivalent forms of the integral evaluation  $\bar{J}_n^{(m)}$  to be established, that is:

$$\bar{J}_n^{(m)} = \bar{J}_n^{(0)} + M_n^{(m)} \quad (24)$$

and

$$\bar{J}_n^{(m)} = \bar{J}_n^{(m)}_{opt} + \Delta M_n^{(m)} \quad (25)$$

A detailed analysis of expressions (24) and (25) will be shown later. Now one is satisfied with the statement that for the determination of the first components, that is  $\bar{J}_n^{(0)}$  and  $\bar{J}_n^{(m)}_{opt}$ , only *SMI* are needed. To find the remaining components, that is  $M_n^{(m)}$  and  $\Delta M_n^{(m)}$ , the knowledge of the polynomial  $\bar{B}(p)$  is also needed. In particular, the component  $M_n^{(m)}$  expresses the increase of the evaluation  $\bar{J}_n^{(m)}$  due to the fact that the polynomial increase  $B(p) = \bar{B}(p) - b_m$  has been taken into consideration, and the component  $\Delta M_n^{(m)}$  is the increase due to the introduction of the polynomial  $\Delta B(p) = \bar{B}(p) - \bar{B}(p)_{opt}$  in the numerator of the transfer function  $\bar{F}(p)$ . For further investigation form (25) will be of particular use.

Table 2. Coefficients  $b_i$  ( $i = 0, 1, 2, \dots, m$ ) of the polynomials  $B(p) = B(p)_{\text{opt}}$ , satisfying the optimum conditions in the sense of the indices of stability margin  $\left(q_i = \frac{\Delta_{i+1}}{\Delta_i}\right)$

m	Coefficients of the polynomial $B(p) = b_0 p^m + b_1 p^{m-1} + \dots + b_{m-1} p + b_m = B(p)_{\text{opt}}$					
	$b_0$	$b_1$	$b_2$	$b_3$	$b_4$	$b_5$
1	$q_{n-2}$	1				
2	$q_{n-3}$	$q_{n-2}$	1			
3	$q_{n-4}$	$q_{n-3}$	$q_{n-2} + \frac{q_{n-4}}{q_{n-3}}$	1		
4	$q_{n-5}$	$q_{n-4}$	$q_{n-3} + \frac{q_{n-5}}{q_{n-3}} + \frac{q_{n-5}}{q_{n-4}} \cdot q_{n-2}$	$q_{n-2} + \frac{q_{n-4}}{q_{n-3}}$	1	
5	$q_{n-5}$	$q_{n-5}$	$q_{n-4} + \frac{q_{n-6}}{q_{n-3}} + \frac{q_{n-6}}{q_{n-4}} \cdot q_{n-2} + \frac{q_{n-6}}{q_{n-5}} \cdot q_{n-3}$	$q_{n-3} + \frac{q_{n-5}}{q_{n-3}} + \frac{q_{n-5}}{q_{n-4}} \cdot q_{n-2}$	$q_{n-2} + \frac{q_{n-4}}{q_{n-3}} + \frac{q_{n-6}}{q_{n-5}}$	1

Notes: (1)  $n$  is the degree of the characteristic equation  $A(p) = 0$   
 (2)  $m$  is the degree of the polynomial  $B(p)$ ;  $0 \leq m < n$   
 (3)  $B(p)$  is the polynomial in the numerator of the normalized  
 $t$  transfer function  $\bar{F}(p) = \frac{A(p)}{B(p)}$

7. The Primary Spectrum of the Integral Evaluation  $\bar{J}_n^{(m)}$

Under the term of primary spectrum of the integral evaluation  $\bar{J}_n^{(m)}$  one will understand the expression

$$R_n = R_n(r_1, r_2, r_3, \dots, r_n) \quad (26)$$

The elements of the spectrum  $R_n$  are  $r_1, r_2, \dots, r_n$ . They are related to the SMI by the formulae

$$r_i = \frac{q_{i-2}}{q_{i-1}} = \frac{\Delta_{i-1} \cdot \Delta_{i-1}}{\Delta_{i-2} \cdot \Delta_i} = \frac{S_{i-1}^* S_{i-1}^*}{S_{i-2}^* S_i^*}; \quad (i=1, 2, \dots, n) \quad (27)$$

where  $q_i$  are the Routh parameters,  $\Delta_i$  the Hurwitz determinants, and  $S_i^*$  the Markov determinants ( $\Delta_i = a_1^i \cdot S_i^*$ ).

Example:

$$r_1 = \frac{\Delta_0}{\Delta_1}; r_2 = \frac{\Delta_1 \Delta_1}{\Delta_0 \Delta_2}; r_3 = \frac{\Delta_2 \Delta_2}{\Delta_1 \Delta_3}; \dots, r_n = \frac{\Delta_{n-1} \cdot \Delta_{n-1}}{\Delta_{n-2} \cdot \Delta_n}$$

Knowing the values of the elements  $r_1, r_2, \dots, r_n$  the values of the corresponding indices of stability margin can easily be determined:

Routh parameters  $q_s$

$$\frac{1}{q_{k-1}} = r_1 r_2 \dots r_k = \prod_{i=1}^k r_i \quad (k=1, 2, \dots, n) \quad (28)$$

Hurwitz determinants  $\Delta_k$

$$\frac{1}{\Delta_k} = r_1^k r_2^{k-1} r_3^{k-2}, \dots, r_k = \prod_{\alpha=1}^k r_\alpha^{k+1-\alpha} \text{ for } k=1, 2, \dots, n \quad (29)$$

Markov determinants  $S_k^*$ ; ( $S_1^* \equiv 1$ )

$$\frac{1}{S_k^*} = r_2^{k-1} r_3^{k-2}, \dots, r_k = \prod_{\alpha=2}^k r_\alpha^{k+1-\alpha} \text{ for } k=2, 3, \dots, n \quad (30)$$

The primary spectrum  $R_n$  determines uniquely the first components of the forms (24) and (25) of the integral evaluation  $\bar{J}_n^{(m)}$ . In particular, by virtue of eqns (14) and (27), one can write at once

$$\bar{J}_n^{(0)} = \frac{1}{2} (r_1 + r_2 + \dots + r_n) = \frac{1}{2} \sum_{i=1}^n r_i \quad (31)$$

It can also be shown<sup>7</sup> that when the optimum conditions (22) are satisfied, the expression of the evaluation  $\bar{J}_n^{(m)}_{\text{opt}}$  takes the following exceptionally simple form

$$\bar{J}_n^{(m)}_{\text{opt}} = \frac{1}{2} (r_1 + r_2 + \dots + r_{n-m}) = \frac{1}{2} \sum_{i=1}^{n-m} r_i \quad (32)$$

where  $0 \leq m < n$ .

From (32) it follows that evaluation  $\bar{J}_n^{(m)}_{\text{opt}}$  depends only on the first  $n - m$  elements of the spectrum  $R_n$  and is invariant in relation to the remaining ones. Thus, the elements  $r_1, r_2, \dots, r_{n-m}$  will be called weight (influence) elements and the remaining ones, that is  $r_{n-(m-1)}, r_{n-(m-2)}, \dots, r_n$ , independent or free ones\*. Observe that although the independent elements of the spectrum  $R_n$  show no influence on the value of the evaluation  $\bar{J}_n^{(m)}_{\text{opt}}$ ,

\* In the case of normalized transfer function the condition  $\prod_{i=1}^n r_i \equiv 1$  should be satisfied.



they influence the character of the transient response. This is a separate problem and is not dealt with in this paper.

On the basis of eqns (27)–(30) the stability of a control system can easily be analysed. From this analysis it follows that if a control system is stable, all the elements of the spectrum  $R_n$  are positive. If, in addition, the system is physically real, these elements are bounded. The spectrum  $R_n$ , of which all the elements are different from zero and positive will be called 'essentially positive'.

Another interesting property of the spectrum  $R_n$  is now shown. It is known that the stability margin of the system is greater for greater values of SMI<sup>3,4</sup>, Hurwitz determinants, for instance. This means that the stability margin is greater for smaller values of elements of the spectrum  $R_n$ .

On the basis of the above results and considerations the following cardinal properties of the spectrum  $R_n$  can be formulated:

*Property I:* In order that a linear control system with the characteristic equation  $A(p) = p^n + a_1 p^{n-1} + \dots + a_n = 0$  is stable and physically real it is necessary and sufficient that the primary spectrum  $R_n = R_n(r_1, r_2, \dots, r_n)$  of the evaluation  $\bar{J}_n^{(m)}$  of this system is essentially positive and bounded, that is  $0 < r_i < \infty$  for  $i = 1, 2, \dots, n$ .

*Property II:* The primary spectrum  $R_n$  characterizes the transient performance in a linear control system of the order  $n$  and class  $m$ , because the sum of its weight elements  $r_i$  ( $i = 1, 2, \dots, n - m$ ) determines the value of the evaluation  $\bar{J}_n^{(m)}$  satisfying the optimum conditions in the sense of stability margin (SMI), that is

$$\bar{J}_{n \text{ opt}}^{(m)} = \frac{1}{2} \sum_{i=1}^{n-m} r_i$$

for  $m = 0, 1, 2, \dots, n - 1$ .

**8. The Secondary Spectrum of the Integral Evaluation  $\bar{J}_n^{(m)}$**

The components  $M_n^{(m)}$  and  $\Delta M_n^{(m)}$  in expressions (24) and (25) for the evaluations  $\bar{J}_n^{(m)}$  depend in the general case on the spectrum  $R_n$  and the polynomial  $\bar{B}(p) = b_0 p^m + b_1 p^{m-1} + \dots + b_{m-1} p + b_m$  or the equivalent polynomial  $C(p) \equiv \bar{B}(p)$ , where

$$C(p) = c_m p^m + c_{m-1} p^{m-1} + \dots + c_1 p + c_0 \equiv \bar{B}(p); (c_0 = b_m = 1) \tag{33}$$

The task now is to find a set of  $m$  parameters such that their structure contains as much information as possible on the transient performance in a control system and would enable the determination, in a unique form, of the values of the coefficients  $c_\alpha = b_{m-\alpha}$  and the easy computation of the component  $M_n^{(m)}$  or  $\Delta M_n^{(m)}$  of  $\bar{J}_n^{(m)}$ . To this aim consider the partial derivatives

$$w_i = \frac{\partial M_n^{(m)}}{\partial c_i} = \sum_{\alpha=i}^m \frac{\partial M_n^{(m)}}{\partial c_\alpha} \text{ for } i = 1, 2, \dots, m \tag{34}$$

One has, in the case of odd  $i$ :

$$w_i = -1 + \sum_{\alpha=i}^{E[\frac{m+1}{2}]} (-1)^{i+\alpha} \cdot c_{2\alpha-1} \cdot \frac{\Delta_{m-i}^{(n)}}{\Delta_n}$$

and

$$w_{2i+1} = \sum_{\alpha=i}^{E[\frac{m+1}{2}]} (-1)^{(i+\alpha)+1} \cdot c_{2\alpha-1} \cdot \frac{\Delta_{m-(i+\alpha)}}{\Delta_n} \tag{35}$$

where  $1 \leq l \leq E[m - 1/2]$ . One finds, for even  $i$ :

$$w_{2l} = (-1)^l \cdot \frac{\Delta_{m-l}^{(n)}}{\Delta_n} + \sum_{i=1}^{E[\frac{m}{2}]} (-1)^{(i+l)} \cdot c_{2i} \cdot \frac{\Delta_{m-(i+l)}^{(n)}}{\Delta_n} \tag{36}$$

for  $1 \leq l \leq k = E[m/2]$ .

The expressions (35) constitute a set of equations for the odd coefficients  $c_\alpha$  of the polynomial (33) and the expressions (36) are a set of equations for even coefficients  $c_\alpha$  of that polynomial. The principal determinants of these systems will be denoted by  $W_m^{(1)}$  and  $W_m^{(2)}$  respectively, where

$$W_m^{(1)} = \begin{vmatrix} \frac{\partial w_1}{\partial c_1} & \frac{\partial w_1}{\partial c_3} & \dots \\ \frac{\partial w_3}{\partial c_1} & \frac{\partial w_3}{\partial c_3} & \dots \\ \dots & \dots & \dots \end{vmatrix} = \begin{vmatrix} w_{11} & w_{13} & w_{15} & \dots \\ w_{31} & w_{33} & w_{35} & \dots \\ w_{51} & w_{53} & w_{55} & \dots \\ \dots & \dots & \dots & \dots \end{vmatrix} \tag{37}$$

$$W_m^{(2)} = \begin{vmatrix} \frac{\partial w_2}{\partial c_2} & \frac{\partial w_2}{\partial c_4} & \dots \\ \frac{\partial w_4}{\partial c_2} & \frac{\partial w_4}{\partial c_4} & \dots \\ \dots & \dots & \dots \end{vmatrix} = \begin{vmatrix} w_{22} & w_{24} & w_{26} & \dots \\ w_{42} & w_{44} & w_{46} & \dots \\ w_{62} & w_{64} & w_{66} & \dots \\ \dots & \dots & \dots & \dots \end{vmatrix} \tag{38}$$

or

$$\left. \begin{aligned} W_m^{(1)} &= \frac{\partial(w_1, w_3, \dots, w_{2k-1}, \dots)}{\partial(c_1, c_3, \dots, c_{2k-1}, \dots)} = \frac{1}{q_{n-2} q_{n-3}, \dots, q_{n-k}}; \\ & \quad k = 2 E \left[ \frac{1+m}{2} \right] \\ W_m^{(2)} &= \frac{\partial(w_2, w_4, \dots, w_{2k}, \dots)}{\partial(c_2, c_4, \dots, c_{2k}, \dots)} = \frac{1}{q_{n-2} q_{n-3}, \dots, q_{n-l}}; \\ & \quad l = 2 E \left[ \frac{2+m}{2} \right] \end{aligned} \right\} \tag{39}$$

From eqns (37)–(39) it follows that the determinants  $W_m^{(1)}$  and  $W_m^{(2)}$  are the Jacobians of the transformation. The elements  $W_{ij}$  of these Jacobians are Krasovskii determinants

$$\frac{\Delta_{m-\alpha}^{(n)}}{\Delta_n}$$

with appropriate signs and

$$w_{ij} = \frac{\partial w_i}{\partial c_j} = \frac{\partial^2 M_n^{(m)}}{\partial c_i \partial c_j} \text{ (} i, j = 1, 2, 3, \dots, m \text{)} \tag{40}$$

Assume that the system is stable and its spectrum  $R_n$  is invariable (constant); then, assume also that the Jacobians  $W_m^{(1)}$  and  $W_m^{(2)}$  have, in agreement with (39), constant values different from zero and positive. From the analysis it follows<sup>11</sup> that in this case all the necessary and sufficient conditions are satisfied for the transformation considered to be homeomorphic. It follows that the representation of the set of parameters  $w_i$  ( $i = 1, 2, \dots, m$ ) in an  $m$ -dimensional space  $L^{(m)}$  on the set of parameters  $c_i$  ( $i = 1, 2, \dots, m$ ) in an  $m$ -dimensional space  $D^{(m)}$  is one-to-one, and that the homeomorphic representation of a space region is a space region and the representation of an arc

is an arc. The set of the parameters  $w_i$  ( $i = 1, 2, \dots, m$ ) will be called the secondary spectrum of the integral evaluation  $\bar{J}_n^{(m)}$  and denoted by

$$w_m = w_m(w_1, w_2, \dots, w_m) \quad (41)$$

If the values of the elements  $w_i$  of the spectrum  $w_m$  are known, it is easy to calculate all the coefficients  $c_\alpha$  of the polynomial (33). To do this it suffices to solve in relation to  $c_\alpha$  the matrix equations\*:

$$\|W_m^{(1)}\| \cdot \|C_m^{(1)}\| = \|V_m^{(1)}\| \text{ and } \|W_m^{(2)}\| \cdot \|C_m^{(2)}\| = \|V_m^{(2)}\| \quad (42)$$

The spectrum  $w_m$  is called positive, zero or negative if all its elements  $w_i$  ( $i = 1, 2, \dots, m$ ) are, respectively, positive or zero or negative. A spectrum  $w_m$  may also be of a mixed type. In particular, from the solution of eqns (42) it follows that if the spectrum  $w_m$  is zero, the set of eqns (22) is satisfied. This important feature of the spectrum  $w_m$  concerning the optimum evaluation  $\bar{J}_n^{(m)}$  in the sense of SMI can be expressed in the form of the following.

*Property of the spectrum  $w_m$ :* In order that the Krasovskii integral evaluation  $\bar{J}_n^{(m)}$  should satisfy the optimum conditions in the sense of stability margin (SMI) it is necessary and sufficient, that its secondary spectrum  $w_m$  is zero; that is,  $w_i \equiv 0$  ( $i = 1, 2, \dots, m$ ).

Now pass to another form of the spectrum  $w_m$  connected with the increment  $\Delta M_n^{(m)}$  of the evaluation  $\bar{J}_n^{(m)}$ . For this purpose the coefficients of the polynomial  $C(p)$  should first be represented in the form

$$c_i = c_{i \text{ opt}} + h_i \quad (i = 1, 2, \dots, m) \quad (43)$$

where  $c_{i \text{ opt}}$  satisfy the optimum conditions in the sense of SMI and expand  $M_n^{(m)}$  in Taylor's series for functions of more than one independent variable

$$\begin{aligned} M_n^{(m)}(c_{1 \text{ opt}} + h_1, c_{2 \text{ opt}} + h_2, \dots, c_{m \text{ opt}} + h_m) \\ = M_{n \text{ opt}}^{(m)} + \frac{dM_n^{(m)}}{1!} + \frac{d^2M_n^{(m)}}{2!} + \dots + \frac{d^{k-1}M_n^{(m)}}{(k-1)!} + R_k \end{aligned} \quad (44)$$

In the general case the derivatives  $d^v M_n^{(m)}$  and the rest  $R_k$  of the expansion (44) are

$$d^v M_n^{(m)} = \left( \frac{\partial M_n^{(m)}}{\partial c_1} h_1 + \frac{\partial M_n^{(m)}}{\partial c_2} h_2 + \dots + \frac{\partial M_n^{(m)}}{\partial c_m} h_m \right)^v \quad (45)$$

and

$$R_k = \frac{d^k M_n^{(m)}}{k!} \quad (46)$$

where the derivatives  $d^v M_n^{(m)}$  for  $v < k$  should be determined at the point  $Q_{\text{opt}} = Q_{\text{opt}}(c_{1 \text{ opt}}, c_{2 \text{ opt}}, \dots, c_{m \text{ opt}})$  and the rest  $R_k$  at an intermediate point  $(C_{1 \text{ opt}} + \theta h_1, c_{2 \text{ opt}} + \theta h_2, \dots, c_{m \text{ opt}} + \theta h_m)$ , where  $0 < \theta < 1$ . One obtains

$$dM_n^{(m)} = \sum_{i=1}^m \frac{\partial M_n^{(m)}}{\partial c_i} h_i = \sum_{i=1}^m w_i h_i \quad (47)$$

$$R_k = R_2 = \frac{1}{2} \left( \sum_{i=1}^m \frac{\partial M_n^{(m)}}{\partial c_i} h_i \right)^2 = \frac{1}{2} \sum_{i=1}^m w_{ii} h_i^2 + \sum_{\substack{i=1, j=2 \\ (i < j)}}^m w_{ij} h_i h_j \quad (48)$$

\* The solution of eqn (42) are given in Table.

where

$$w_{ij} = \frac{\partial^2 M_n^{(m)}}{\partial c_i \partial c_j}$$

The expansion of the function  $M_n^{(m)}$  in Taylor's series, taking into account eqns (47) and (48), will now be written

$$M_n^{(m)} = M_{n \text{ opt}}^{(m)} + dM_n^{(m)} + R_2 \quad (49)$$

or

$$\Delta M_n^{(m)} = M_n^{(m)} - M_{n \text{ opt}}^{(m)} = dM_n^{(m)} + R_2$$

In agreement with the optimum theorem of the evaluation  $\bar{J}_n^{(m)}$  the point  $Q_{\text{opt}}(c_{1 \text{ opt}}, c_{2 \text{ opt}}, \dots, c_{m \text{ opt}})$  corresponds to a zero spectrum  $w_m$ . In other words the derivative (47) is at this point equal to zero, i.e.

$$dM_n^{(m)}(Q_{\text{opt}}) = 0 \quad (50)$$

It can easily be shown that the second partial derivatives  $w_{ij}$  do not depend on the choice of the intermediate point. Therefore the expression of the component  $\Delta M_n^{(m)}$  of the evaluation  $\bar{J}_n^{(m)}$  takes the following very simple form

$$\Delta M_n^{(m)} = R_2 = \frac{1}{2} \sum_{i=1}^m w_{ii} h_i^2 + \sum_{\substack{i=1, j=2 \\ (i < j)}}^m w_{ij} h_i h_j \quad (w_{i, i+1} \equiv 0) \quad (51)$$

The partial derivatives  $w_{ii}$  and  $w_{ij}$  are Krasovskii determinants taken with an appropriate sign, therefore they depend on the spectrum  $R_n$  only. Analogous considerations show that the representation of the set of parameters  $w_i$  into the set of parameters  $h_i$  ( $i = 1, 2, \dots, m$ ) is also homeomorphic; that is, one-to-one. In this connection the parameters  $h_i$  will be taken as elements of the second, equivalent form of the secondary spectrum  $w_m$ ; that is,

$$w_m(h_1, h_2, \dots, h_m) = w_m(w_1, w_2, \dots, w_m) \quad (52)$$

If the spectrum  $w_m$  is zero; that is,  $h_i \equiv 0$  ( $i = 1, 2, \dots, m$ ); then  $\Delta M_n^{(m)} = 0$  and  $\bar{J}_n^{(m)} = \bar{J}_{n \text{ opt}}^{(m)}$ .

The secondary spectrum  $w_m$  has a number of properties facilitating the qualitative analysis of the influence of distribution of zeros of the transfer function on the transient performance<sup>7, 12</sup>. Thus, for instance, a positive or negative spectrum  $w_m$  shows that the corresponding fluctuations of transient responses are greater or less than the same fluctuations for the evaluation  $\bar{J}_n^{(m)}$  opt.

## 9. The Inversion of the Integral Evaluation $\bar{J}_n^{(m)}$ by means of the Spectra $R_n$ and $w_m$

To estimate the transient performance in a control system various integral criteria have found broad application. This is done most often by a comparative method. The less is the value of the integral evaluation chosen, the higher is the quality of the transient response. In this connection various methods have been developed for investigation of the relation between a change of values of selected transfer function parameters and the corresponding change of the value of the integral evaluation. Of the best known and most widely used, mention should be made of methods of minimizing the integral evaluation in relation to one or a few parameters: graphoanalytic methods of determining the minimum evaluation and the method of successive trials and approximations.

The aim is to obtain analytically a new solution of this problem using the integral square criterion of transient performance which was called the inverse problem of the integral square estimation  $\bar{J}_n^{(m)}$  and which could also be called the inverse problem of transient performance. This is a problem encountered particularly in the synthesis of linear system.

Under the name of inversion of the integral square estimation  $\bar{J}_n^{(m)}$  one will understand the determination of the normalized transfer function for a prescribed value of the normalized Krasovskii integral evaluation  $\bar{J}_n^{(m)}$ . The solution of the inversion problem of the integral criterion  $\bar{J}_n^{(m)}$  has been obtained by introducing the notions of the spectra  $R_n$  and  $w_m$  defined above. It should be explained that in the general case the evaluation  $\bar{J}_n^{(m)}$  is a multivalued function. For any assigned value of the evaluation  $\bar{J}_n^{(m)}$  an infinite number of various linear transfer functions can be made to correspond. A different case is that where the evaluation  $\bar{J}_n^{(m)}$  is expressed in terms of the spectra  $R_n$  and  $w_m$ , the correspondence between a transfer function and these spectra being now one-to-one. In this connection, if the inverse evaluation  $\bar{J}_n^{(m)}$  is spoken of one always means the inverse evaluation  $\bar{J}_n^{(m)}$  expressed in terms of definite spectra  $R_n$  and  $w_m$ .

Assume that the spectra of the Krasovskii evaluation  $\bar{J}_n^{(m)}$ , primary  $R_n = R_n(r_1, r_2, \dots, r_n)$  and secondary  $w_m = w_m(h_1, h_2, \dots, h_m)$ , are known. They contain full information on the transfer function  $F(p) = C(p)/A(p)$  of the control channel under consideration and much information on the transient performance in this control channel.

The transfer function will be determined by inversion of the evaluation  $\bar{J}_n^{(m)}$  expressed in terms of the spectra  $R_n$  and  $w_m$ ; that is, by inverse transformation of the spectra  $R_n$  and  $w_m$ . For this purpose determine first the characteristic equation corresponding to the spectrum  $R_n$ .

The method of determining the characteristic equation (polynomial)  $A(p) = p^n + a_1 p^{n-1} + \dots + a_{n-1} p + a_n$  is an elementary one. The values of all the elements  $r_1, r_2, \dots, r_n$  of the spectrum  $R_n$  being known, it is easy, on the basis of (29), to determine, for instance, the values of all the Hurwitz determinants and then to make use of Table 1 which enables the values of the coefficients of the characteristic polynomial  $A(p)$  to be found directly.

In order to avoid the intermediate stage of computing the Hurwitz determinants, Table 3 has been prepared, containing expansions of the coefficients of the characteristic equation directly in terms of the elements  $r_i$  ( $i = 1, 2, \dots, n$ ) of the primary spectrum  $R_n$ . In this case the coefficients  $a_{k,n}$  may be determined by means of the algorithm

$$a_{k,n} = a_{k,n-1} + \frac{a_{k-2,n-2}}{r_{n-1} \cdot r_n} \quad (53)$$

or by means of the recurrence equations

$$A_k(p) = p A_{k-1}(p) + \frac{A_{k-2}(p)}{r_{k-1} \cdot r_k} \quad (k = 1, 2, \dots, n) \quad (54)$$

The expressions (53) and (54) are equivalent to eqns (8) and (9)<sup>7</sup>. In the case of low degrees  $n$  on the characteristic polynomial it is more convenient to use Table 3. In the case of high degrees  $n$  eqns (53) or (54) are better suited for computation.

The method for finding the polynomial  $c(p)$  in the numerator

of the transfer function  $F(p) = C(p)/A(p)$  is also elementary; it consists in representing the polynomial  $C(p)$  as a sum of two polynomials

$$C(p) = c_m p^m + c_{m-1} p^{m-1} + \dots + c_1 p + 1 = C(p)_{\text{opt}} + \Delta C(p) \quad (55)$$

where

$$\left. \begin{aligned} C(p)_{\text{opt}} &= c_{m \text{ opt}} p^m + c_{m-1 \text{ opt}} p^{m-1} + \dots + c_{1 \text{ opt}} p + 1 \\ \Delta C(p) &= h_m p^m + h_{m-1} p^{m-1} + \dots + h_1 p \end{aligned} \right\} \quad (56)$$

The coefficients  $c_{\alpha \text{ opt}}$  ( $\alpha = 1, 2, \dots, m$ ) of the polynomial  $C(p)_{\text{opt}}$  can easily be obtained from Table 2 and the spectrum  $R_n$ . One has only to take into consideration the relation  $c_{\alpha \text{ opt}} = b_{m-\alpha \text{ opt}}$  between the coefficients of the equivalent polynomial  $C(p)_{\text{opt}}$  and  $B(p)_{\text{opt}}$ . The coefficients  $h_i$  ( $i = 1, 2, \dots, m$ ) are known if the spectrum  $w_m$  is known. As a consequence the transfer function  $F(p) = C(p)/A(p)$  is uniquely determined. The spectra  $R_n$  and  $w_m$  enable one to estimate simultaneously the stability of a system and the transient performance in agreement with the principles studied in Sections 6-8.

The method based on the inverse integral criterion  $\bar{J}_n^{(m)}$  and the spectra of the evaluation  $\bar{J}_n^{(m)}$  may be used successfully for analysis and synthesis of linear systems of automatic control. In the first case the transfer function is known, therefore the corresponding spectra of the evaluation  $\bar{J}_n^{(m)}$  can be found easily; in particular, the primary spectrum  $R_n$  can rapidly be determined by using, for instance, the Markov criterion or finding, by successive elimination, the elements  $r_i$  ( $i = 1, 2, \dots, n$ ) directly from the expansion of the characteristic equation in terms of these elements. The next stage is that of correction of the spectra obtained.

In the case of synthesis it is necessary to know the general structure of the transfer function (that is  $n$  and  $m$  must be known) and the requirements concerning the evaluation  $\bar{J}_n^{(m)} = \bar{J}_n^{(m)}_{\text{opt}} + \Delta M_n^{(m)}$ . Correct choice of the weight elements  $r_i$  ( $i = 1, 2, \dots, n - m$ ) and the elements  $h_i$  of the spectrum  $w_m$  is of particular importance. It follows that the problem of correct choice of the spectra  $R_n$  and  $w_m$  is essential for the synthesis of linear systems and the necessary correction of such system.

## 10. Minimization of the Integral Evaluation $\bar{J}_n^{(m)}$

Some additional data on the primary spectrum  $R_n$  may be obtained by minimizing the integral evaluation  $\bar{J}_n^{(m)}$ . The minimization of the evaluation  $\bar{J}_n^{(m)} = \bar{J}_n^{(m)}_{\text{opt}} + \Delta M_n^{(m)}$ , expressed in terms of the *SMI* may be divided into two stages. The first consists in obtaining the optimum integral estimation  $\bar{J}_n^{(m)}$  in the sense of *SMI*, discussed in Section 5. As result the component  $\Delta M_n^{(m)}$  vanishes.

The next stage consists in minimizing the component  $\bar{J}_n^{(m)}_{\text{opt}}$  in relation to a selected group of *SMI*, for instance

$$\begin{aligned} H &= H(\Delta_1, \Delta_2, \dots, \Delta_n) \\ R &= R(q_0, q_1, \dots, q_{n-1}) \\ M &= M(\bar{q}_0, \bar{q}_1, \dots, \bar{q}_{n-1}) \end{aligned} \quad (57)$$

where  $H, R, M$  are groups of determinant indices of stability margin in the sense of the criterion of Hurwitz, Routh and

Table 3. Expansion of the coefficients of the characteristic equation in terms of the elements  $r_i$  ( $i = 1, 2, \dots, n$ ) of the primary spectrum  $R_n$  of the evaluation  $J_n^{(n)}$

$n$	Coefficients of the characteristic equation $A(p) = a_0 p^n + a_1 p^{n-1} + \dots + a_n$						
	$a_0$	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$
1	1	$\frac{1}{r_1}$					
2	1	$\frac{1}{r_1}$	$a_2 = \frac{1}{r_1 r_2}$				
3	1	$\frac{1}{r_1}$	$\frac{1}{r_1} + \frac{1}{r_2 r_3}$	$a_3 = \frac{1}{r_1^2} \left( \frac{1}{r_2 r_3} \right)$			
4	1	$\frac{1}{r_1}$	$\frac{1}{r_1 r_2} + \frac{1}{r_2 r_3} + \frac{1}{r_3 r_4}$	$\frac{1}{r_1} \left( \frac{1}{r_2 r_3} + \frac{1}{r_3 r_4} \right)$	$a_4 = \frac{1}{r_3 r_4} \left( \frac{1}{r_1 r_2} \right)$		
5	1	$\frac{1}{r_1}$	$\frac{1}{r_1 r_2} + \frac{1}{r_2 r_3} + \frac{1}{r_3 r_4} + \frac{1}{r_4 r_5}$	$\frac{1}{r_1} \left( \frac{1}{r_2 r_3} + \frac{1}{r_3 r_4} + \frac{1}{r_4 r_5} \right)$	$\frac{1}{r_3 r_4} \left( \frac{1}{r_1 r_2} \right) + \frac{1}{r_4 r_5} \left( \frac{1}{r_1 r_2} + \frac{1}{r_2 r_3} \right)$	$a_5 = \frac{1}{r_1 r_4 r_5} \left( \frac{1}{r_2 r_3} \right)$	
6	1	$\frac{1}{r_1}$	$\frac{1}{r_1 r_2} + \frac{1}{r_2 r_3} + \frac{1}{r_3 r_4} + \frac{1}{r_4 r_5} + \frac{1}{r_5 r_6}$	$\frac{1}{r_1} \left( \frac{1}{r_2 r_3} + \frac{1}{r_3 r_4} + \frac{1}{r_4 r_5} + \frac{1}{r_5 r_6} \right)$	$\frac{1}{r_3 r_4} \left( \frac{1}{r_1 r_2} \right) + \frac{1}{r_4 r_5} \left( \frac{1}{r_1 r_2} + \frac{1}{r_2 r_3} \right) + \frac{1}{r_5 r_6} \left( \frac{1}{r_1 r_2} + \frac{1}{r_2 r_3} + \frac{1}{r_3 r_4} \right)$	$a_5 = \frac{1}{r_1 r_4 r_5} \left( \frac{1}{r_2 r_3} \right) + \frac{1}{r_1 r_5 r_6} \left( \frac{1}{r_2 r_3} + \frac{1}{r_3 r_4} \right)$	$a_6 = \frac{1}{\prod_{i=1}^6 r_i}$

Notes: (1)  $J_n^{(n)}$  is the Krasovskii integral evaluation.  
 (2) In the case of normalized characteristic equation the condition  $\prod_{i=1}^n r_i \equiv 1$  should be satisfied.

## 287 / 10

Markov, respectively. Of course, the result of the minimization process depends on the choice of the group of independent variables  $H$ ,  $R$  or  $M$ .

As a result of investigations a new property of the spectrum has been found, characteristic of the conditional minimum of the evaluation  $\bar{J}_n^{(m)}$ . This is as follows:

*Property III:* An integral evaluation  $\bar{J}_n^{(m)}$  minimized with respect to an appropriate group of  $SMI$  reaches a conditional minimum if its secondary spectrum  $w_m$  is zero and the weight elements of the primary spectrum  $R_n$  constitute an arithmetic progression of which the difference is  $\Delta r_i$  and  $\Delta r_i = \kappa/(n-m)\Delta_1$ ;  $r_{i+1} = r_i + \Delta r_i$  ( $i = 1, 2, \dots, n-m$ ;  $m > 0$ ).

The conditional minimum of the evaluation  $\bar{J}_n^{(m)}$  ( $m > 0$ ) is expressed by the equation (Jarominek<sup>7</sup>),

$$\bar{J}_n^{(m)}(\kappa)_{\min} = \frac{2(n-m) + \kappa(n-m-1)}{4\Delta_1} > 0 \quad (58)$$

In the particular case of

$$\kappa = -1 \text{ one has } \bar{J}_n^{(m)}(H)_{\min} = \bar{J}_n^{(m)}(\kappa = -1)_{\min} = \frac{n-m+1}{4\Delta_1} \quad (59)$$

and

$$\kappa = 0 \text{ one has } \bar{J}_n^{(m)}(R)_{\min} = \bar{J}_n^{(m)}(\kappa = 0)_{\min} = \frac{n-m}{2\Delta_1} \quad (60)$$

### 11. The Problem of Correct Distribution of Zeros of the Transfer Function

The secondary spectrum  $w_m$  characterizes indirectly the distribution of zeros of the transfer function, which has a strong influence on the value of the evaluation  $\bar{J}_n^{(m)}$ . As a criterion of correct distribution of the zeros of the transfer function in relation to the distribution of the poles, the value of the component  $\Delta M_n^{(m)} \geq 0$  of the evaluation  $\bar{J}_n^{(m)} = \bar{J}_n^{(m)}_{\text{opt}} + \Delta M_n^{(m)}$  is taken. The component  $\Delta M_n^{(m)}$  depends on both spectra  $R_n$  and  $w_m$ . The spectrum  $w_m$ , which is non-zero in general, will be correctly chosen to fit the spectrum  $R_n$  if the component  $\Delta M_n^{(m)}$  does not exceed a certain fixed value.

The upper bound of the value of the component  $\Delta M_n^{(m)}$  can be determined assuming as a principle the maximum utilization of the polynomial  $C(p)$  of degree  $m$ . For this will be used the first equivalent form of the integral criterion  $\bar{J}_n^{(m)}$ , that in

$$\bar{J}_n^{(m)} = \bar{J}_n^{(0)} + M_n^{(m)} \text{ and } \bar{J}_n^{(m)} = \bar{J}_n^{(0)} + M_n^{(m)} \quad (61)$$

On the basis of appropriate considerations the following can be written:

$$\lim \inf. M_n^{(m)} = M_n^{(m)}_{\text{opt}} \leq M_n^{(m)} \leq \lim \sup. M_n^{(m)} = M_n^{(m-1)}_{\text{opt}} \quad (62)$$

It follows that

$$0 \leq M_n^{(m)} - M_n^{(m)}_{\text{opt}} = \Delta M_n^{(m)} \leq M_n^{(m-1)} - M_n^{(m)}_{\text{opt}} \quad (63)$$

On the other hand it is easy to show that

$$M_n^{(m-1)} - M_n^{(m)}_{\text{opt}} = \frac{1}{2} r_{n-m+1} \quad (64)$$

It is inferred that  $\Delta M_n^{(m)}$  should be contained between the limits

$$0 \leq \Delta M_n^{(m)} \leq \frac{1}{2} r_{n-m+1} \quad (65)$$

The secondary spectrum  $w_m$  is therefore correctly chosen if condition (65) is satisfied.

The above method, based on the solution of the inverse stability problem and the inverse transformation of the integral square estimation  $\bar{J}_n^{(m)}$  may also be successfully used for the analysis and synthesis of multi-loop linear systems of automatic control. In the case of synthesis it enables the parameters following the requirements concerning the static and dynamic characteristics of the system to be chosen<sup>7</sup>. It may also be helpful for the investigation of non-linear systems in the applicability limits of the Liapunov theorems<sup>10</sup>, and for the investigation of some adaptive systems.

### References

- JAROMINEK, W. Investigating linear systems of automatic control by means of determinant stability margin indices. *Automatic and Remote Control*. 1960. London; Butterworths
- KRASOVSKII, A. A. *Integral Evaluations and the Selection of Parameters of Automatic Control Systems*. 1954. Moscow; Mashgiz
- POPOV, E. P. *Dinamika Sistem Avtomaticheskovo Regulirovania*. 1954. Gostekhizdat
- JAROMINEK, W. Ob ekvivalentnosti kriteria ustoichivosti po Routh'u, Hurwitz'u i po Markov'u. *DAN U.S.S.R.*, T. 130, No. 6 (1960)
- JAROMINEK, W. Primenenie pokazatelei zapasa ustoichivosti po opredelitelam dla issledovaniia lineinykh sistem avtomaticheskogo regulirovania (1959)
- JAROMINEK, W. Razlozhenie koeficientov kharakteristicheskogo uravenia po determinantnym pokazateliam zapasa ustoichivosti (in press)
- JAROMINEK, W. Issledovanie lineinykh stationarnykh sistem metodom spektrov integralnykh ocenok (1961)
- JAROMINEK, W. Ob odnom sluchae nakhozhenia minimuma kvadraticnoi integralnoi ocenki. *Izv. Vyzshikh Uchebnykh Zavedenii, razdel Electromech.* No. 13 (1959)
- JAROMINEK, W. O zapase ustoichivosti lineinykh sistem avtomaticheskogo regulirovania. *Izv. Vyzshikh Uchebnykh Zavedenii, razdel Electromech.*, No. 8 (1959)
- LETOV, A. M. Sostoianie problemy ustoichivosti v teorii avtomaticheskovo regulirovania, T. I. *Izd. AN U.S.S.R.* (1955)
- LEJA, F. Rachunek rozniczkowy, calkowy, *PWN*, Warsaw (1954)
- PETROV, B. N. *Sviaz Mezhdru Kachestvom Perekhodnogo Processa i Pazpredeleniem nulei i Poliusov Peredachnoi Funkcii v sb; Teoria Avtomaticheskogo Regulirovania*. 1954. Moscow; Mashgiz

# Axiomatization of the Theory of Simplification of Combinational Automata

GR. C. MOISIL *Rum*

It is intended to derive a calculus which the axioms have to satisfy in order that the simplifying method, given by Quine for the II-dipoles with contacts and relays, should be valid. Quine's method has been thoroughly investigated by many researchers whose contributions are important, especially J. McCluskey and J. Paul Roth. The axioms given here show that this method may be used in many other instances, such as that of circuits with triodes, transistors, cryotrons, or with three positional relays (as in the case of polarized relays or of real operating of ordinary relays) and with multipositional contacts (as in the case of selectors and codified relays).

I. The researches of Quine and of his successors are related to the two expansions in Boole series

$$f(x_1, \dots, x_n) = \bigcup_{\alpha} f(\alpha_1, \dots, \alpha_n) L_{\alpha_1}(x_1), \dots, L_{\alpha_n}(x_n) \quad (1)$$

$$f(x_1, \dots, x_n) = \prod_{\alpha} [f(\bar{\alpha}_1, \dots, \bar{\alpha}_n) \cup L_{\alpha_1}(x_1) \cup \dots \cup L_{\alpha_n}(x_n)] \quad (2)$$

where

$$L_0(z) = \bar{z}, \quad L_1(z) = z \quad (3)$$

The formulae are of the following type:

$$f(x_1, \dots, x_n) = \Omega [c_{\alpha_1, \dots, \alpha_n} \theta L_{\alpha_1}(x_1) \theta \dots \theta L_{\alpha_n}(x_n)] \quad (4)$$

where  $\theta$  and  $\omega$  represent two operations with any number whatever of variables and where expressions such as

$$\begin{aligned} \Omega z_i &= z_1 \omega \dots \omega z_r \\ \Theta z_i &= z_1 \theta \dots \theta z_r \end{aligned} \quad (5)$$

with  $r > 1$  have a meaning, namely eqns (1) and (2) are eqn (4) if

	$\omega$	$\theta$	$c_{\alpha_1, \dots, \alpha_n}$	
I	$\cup$	$\cdot$	$f(\alpha_1, \dots, \alpha_n)$	(6)
II	$\cdot$	$\cup$	$f(\bar{\alpha}_1, \dots, \bar{\alpha}_n)$	

It has already been shown that six other expansion formulae of the eqn (4) type are valid:

	$\omega$	$\theta$	$c_{\alpha_1, \dots, \alpha_n}$	
III	$\cup$	$\top$	$f(\bar{\alpha}_1, \dots, \bar{\alpha}_n)$	(7)
IV	$\Pi$	$\perp$	$f(\alpha_1, \dots, \alpha_n)$	
V	$\perp$	$\cup$	$f(\bar{\alpha}_1, \dots, \bar{\alpha}_n)$	
VI	$\top$	$\Pi$	$f(\alpha_1, \dots, \alpha_n)$	
VII	$\perp$	$\perp$	$f(\alpha_1, \dots, \alpha_n)$	
VIII	$\top$	$\top$	$f(\bar{\alpha}_1, \dots, \bar{\alpha}_n)$	

The functions  $\top$  and  $\perp$  are Sheffer's functions of several variables

$$\begin{aligned} z_1 \top \dots \top z_r &= \bar{z}_1 \dots \bar{z}_r \\ z_1 \perp \dots \perp z_r &= \bar{z}_1 \cup \dots \cup \bar{z}_r \end{aligned} \quad (8)$$

The interpolation formula of Lagrange in GF(2) is of the same type; as a matter of fact there are here two Lagrange interpolation formulae for GF(2) (IX and its dual X):

	$\omega$	$\theta$	$c_{\alpha_1, \dots, \alpha_n}$	
IX	$+$	$\cdot$	$f(\alpha_1, \dots, \alpha_n)$	(9)
X	$\nabla$	$\cup$	$f(\bar{\alpha}_1, \dots, \bar{\alpha}_n)$	

The functions  $+$  and  $\nabla$  of several variables are defined recurrently

$$\begin{aligned} z_1 + \dots + z_n &= (z_1 + \dots + z_{n-1}) + z_n \\ z_1 \nabla \dots \nabla z_n &= (z_1 \nabla \dots \nabla z_{n-1}) \nabla z_n \end{aligned} \quad (10)$$

while

$$\begin{aligned} \alpha + \beta &= \alpha \bar{\beta} \cup \bar{\alpha} \beta \\ \alpha \nabla \beta &= \alpha \bar{\beta} \cup \alpha \beta \end{aligned} \quad (11)$$

In all these cases,  $f(\alpha_1, \dots, \alpha_n)$  is 0 or 1 and eqn (4) reduces itself to the function generated by an expression

$$f(x_1, \dots, x_n) = f_{\mathcal{E}}(x_1, \dots, x_n) \quad (12)$$

where  $\mathcal{E}$  is an expression, yielding the definition:

I. The expressions are sequences of letters of the following form

$$\Omega (z_{h1} \theta \dots \theta z_{lm_n}) \quad (13)$$

If  $r > 1$ ,  $\Omega$  is defined by eqn (4). If  $r = 1$ ,  $\Omega z_h$  is  $t$  if  $\omega$  is  $\cup, \cdot, +, \nabla$  and is  $\bar{t}$  if  $\omega$  is  $\top$  or  $\perp$ ; in (13),  $z_{ij}$  will be substituted by  $L_{\alpha}(x_{\beta})$  and therefore by  $x_{\beta}$  or  $\bar{x}_{\beta}$ .

II. Between the expressions, a relation of equivalence  $=$  may take place, satisfying the following conditions if

$$\mathcal{A}_i = \mathcal{D}_i$$

then

$$\Omega \mathcal{A}_i = \Omega \mathcal{D}_i$$

$$\Theta \mathcal{A}_i = \Theta \mathcal{D}_i$$

Evidently,  $\cup, \dots, \nabla$  satisfy condition II.

304/2

II. The application of Quine's method is based on the formulae

$$\begin{aligned}x \cup \bar{x} &= 1 \\ 1x &= x\end{aligned}\quad (14)$$

Therefore

$$\begin{aligned}z_0 z_1 \dots z_r \cup \bar{z}_0 z_1 \dots z_r \cup t_1 \cup \dots \cup t_s \\ = (z_0 \cup \bar{z}_0) z_1 \dots z_r \cup t_1 \cup \dots \cup t_s \\ = 1 z_1 \dots z_r \cup t_1 \cup \dots \cup t_s \\ = z_1 \dots z_r \cup t_s \cup \dots \cup t_s\end{aligned}\quad (15)$$

In order to apply this method, the terms  $z_0 z_1 \dots z_r, \bar{z}_0 z_1 \dots z_r$  must be brought to be neighbours and the variables must be arranged in a definite order. Therefore, it will be assumed that

III.  $\Omega$  and  $\Theta$  are commutative; that is to say if  $\pi$  is a permutation of the indexes 1, ..., r, then

$$\begin{aligned}z_{\pi(1)} \theta \dots \theta z_{\pi(r)} &= z_1 \theta \dots \theta z_r \\ z_{\pi(1)} \omega \dots \omega z_{\pi(r)} &= z_1 \omega \dots \omega z_r\end{aligned}\quad (16)$$

This property allows the expression to take the form

$$(z_0 \theta z_1 \theta \dots \theta z_r) \omega (\bar{z}_0 \theta z_1 \theta \dots \theta z_r) \omega t_1 \omega \dots \omega t_s$$

The commutativity is valid for all the operations given as examples:  $\cup, \cdot, \top, \perp, +, \nabla$ . Yet, in order to make the simplification, it is not necessary that all the steps in eqn (15) could be made. It is sufficient that

IV. The following equality be true

$$\begin{aligned}(z_0 \theta z_1 \theta \dots \theta z_r) \omega (\bar{z}_0 \theta z_1 \theta \dots \theta z_r) \omega t_1 \omega \dots \omega t_s \\ = (z_1 \theta \dots \theta z_r) \omega t_1 \omega \dots \omega t_s\end{aligned}\quad (17)$$

It is important to emphasize that for all the pairs of operations  $\omega, \theta$ , from eqns (6), (7) and (9), eqn (14) remains true. That is so much more remarkable, since the various steps made in eqn (15), such as the associativity of  $\cup$ , the distributivity of  $\cdot$ , with regard to  $\cup$ , etc. are not valid for some of these pairs (I-X) of  $\omega, \theta$  operations.

III. This first stage of simplification is valid for:

(a) the dipoles  $\Pi$  with contacts, as well in the normally disjunctive form ( $\omega = \cup, \theta = \cdot$ ) as in the normally conjunctive form ( $\omega = \cdot, \theta = \cup$ );

(b) the diode circuits, in the same cases;

(c) the triode circuits, of the two following forms

$$\omega = \cup, \quad \theta = \top \quad (\text{form III})$$

$$\omega = \top, \quad \theta = \top \quad (\text{form VIII})$$

(d) the transistor circuits of the eight forms I-VIII;

(e) the transistor circuits of the form IX, X;

(f) the cryotron circuits of the following forms

$$\omega = \perp, \quad \theta = \perp \quad (\text{form VII})$$

$$\omega = \top, \quad \theta = \top \quad (\text{form VIII})$$

IV. In the classical case, the following simplification is made

$$\begin{aligned}xyz \cup \bar{x}yz \cup x\bar{y}z \cup xy\bar{z} \\ = xyz \cup \bar{x}yz \cup xyz \cup x\bar{y}z \cup xyz \cup xy\bar{z} \\ = (x \cup \bar{x})yz \cup (y \cup \bar{y})xz \cup (z \cup \bar{z})xy \\ = yz \cup xz \cup xy\end{aligned}\quad (18)$$

by virtue of the idempotence law

$$z \cup z = z \quad (19)$$

To indulge in this type of computation, it is necessary to assume that

V. The following equality is true

$$\begin{aligned}\mathcal{A}_0 \omega \mathcal{A}_0 \omega \mathcal{A}_1 \omega \dots \omega \mathcal{A}_r \\ = \mathcal{A}_0 \omega \mathcal{A}_1 \omega \dots \omega \mathcal{A}_r\end{aligned}$$

This property is valid for the operations  $\cup, \cdot, \top, \perp$ , but it is not valid for  $+$  and  $\nabla$ .

VI. It is known that in the classical case, there can be the following type of simplification

$$\begin{aligned}xy \cup \bar{y}z \cup xz = xy \cup \bar{y}z \cup xyz \cup x\bar{y}z \\ = (xy \cup xyz) \cup (\bar{y}z \cup x\bar{y}z) \\ = xy \cup \bar{y}z\end{aligned}\quad (21)$$

The problems arising from this type of simplification constitute the originality of Quine's method.

A start is made with an expression such as eqn (13) where the  $z_{i_h}$  have been replaced by  $L_{\alpha}(x_{\beta})$  as in the expressions provided by eqn (4).

An expression of the form

$$\mathcal{A} = L_{\alpha_1}(x_{\alpha_1}) \theta \dots \theta L_{\alpha_r}(x_{\alpha_r}) \quad (22)$$

is called a simple expression.

If  $\mathcal{A}, \mathcal{D}$  are simple expressions

$$\mathcal{A} \propto \mathcal{D} \quad (23)$$

provided that

$$\{L_{\alpha_1}(x_{\alpha_1}), \dots, L_{\alpha_r}(x_{\alpha_r})\} \supset \{L_{\beta_1}(x_{\beta_1}), \dots, L_{\beta_s}(x_{\beta_s})\} \quad (24)$$

where the inclusion is considered in the sense of the set theory.

It is obvious that, on the basis of principles I-V, it can be deduced from eqn (23), that\*

$$\begin{aligned}L_{\alpha_1}(x_{\alpha_1}) \theta \dots \theta L_{\alpha_r}(x_{\alpha_r}) \theta L_{\beta_1}(x_{\beta_1}) \theta \dots \theta L_{\beta_s}(x_{\beta_s}) \\ = L_{\alpha_1}(x_{\alpha_1}) \theta \dots \theta L_{\alpha_r}(x_{\alpha_r})\end{aligned}\quad (25)$$

Since eqn (23) is equivalent to eqn (24), it is easy to deduce that the relation  $\propto$  between the simple expressions is a relation of partial order, i.e.

\* It can be seen that this equality cannot be written as

$$\mathcal{A} \theta \mathcal{D} = \mathcal{A}$$

since  $\theta$  is not associative (in particular  $\top$  and  $\perp$  are not associative).

$$\mathcal{A} \propto \mathcal{A}$$

if  $\mathcal{A} \propto \mathcal{D}$  and  $\mathcal{D} \propto \mathcal{L}$ , then  $\mathcal{A} \propto \mathcal{L}$

if  $\mathcal{A} \propto \mathcal{D}$  and  $\mathcal{D} \propto \mathcal{A}$ , then  $\mathcal{A} = \mathcal{D}$

The relation of contiguity will be defined between two simple expressions of the form

$$\mathcal{A} = L_{\alpha_1}(x_1)\theta \dots \theta L_{\alpha_{s_1-1}}(x_{s_1-1})\theta L_{\alpha_{s_1+1}}(x_{s_1+1})\theta \dots \theta L_{\alpha_{t-1}}(x_{s_t-1})\theta L_{\alpha_{t+1}}(x_{s_t+1})\theta \dots \theta L_{\alpha_n}(x_n) \quad (26)$$

$$\mathcal{D} = L_{\beta_1}(x_1)\theta \dots \theta L_{\beta_{t-1}}(x_{s_t-1})\theta L_{\beta_{t+1}}(x_{s_t+1})\theta \dots \theta L_{\beta_n}(x_n)$$

(the missing letters  $x_{s_1}, \dots, x_{s_t}$  are the same in  $\mathcal{A}$  and  $\mathcal{D}$ ) by

$$|\alpha_1 - \beta_1| + \dots + |\alpha_n - \beta_n| = 1 \quad (27)$$

VI. It is proposed to simplify methodically the expression

$$\mathcal{E} = \mathcal{M}_1 \omega \dots \omega \mathcal{M}_t \quad (28)$$

where  $\mathcal{M}_i$  are simple expressions with  $n$  letters  $x$ .

If the simple expressions  $\mathcal{M}_i$  and  $\mathcal{M}_j$  are compared in case they are contiguous, if the letters  $x$  (third principle) are re-ordered, and if use is made of eqn (17) (fourth principle); a simple expression  $\mathcal{M}_{(ij)}$  is formed with  $n - 1$  letters. One thus obtains several simple expressions  $\mathcal{M}_{(ij)_1}, \dots, \mathcal{M}_{(ij)_{t_1}}$  with  $n - 1$  letters. By virtue of the fifth principle

$$\mathcal{E} = \mathcal{M}_1 \omega \dots \omega \mathcal{M}_t \omega \mathcal{M}_{(ij)_1} \omega \dots \omega \mathcal{M}_{(ij)_{t_1}}$$

Compare the simple expressions  $\mathcal{M}_{(ij)_p}$  and  $\mathcal{M}_{(ij)_q}$  with  $n$  letters; if two of them are contiguous, simplify a letter according to the fourth principle and obtain expressions  $\mathcal{M}_{(ijk)_1}, \dots, \mathcal{M}_{(ijk)_{t_2}}$  with  $n - 2$  letters. The operation is reiterated as many times as possible and it will yield  $\mathcal{E}$  in the following form:

$$\mathcal{E} = \mathcal{M}_1 \omega \dots \omega \mathcal{M}_t \omega \mathcal{M}_{(ij)_1} \omega \dots \omega \mathcal{M}_{(ij)_{t_1}} \omega \mathcal{M}_{(ijk)_1} \omega \dots \quad (29)$$

Among all these  $\mathcal{M}$ , the name of prime implicant is given to  $\mathcal{M}_\alpha$  if, from

$$\mathcal{M}_\alpha \propto \mathcal{M}_\beta$$

can be deduced  $\mathcal{M}_\alpha = \mathcal{M}_\beta$ . Let

$$B_1, \dots, B_v \quad (30)$$

be the totality of the prime implicants.

It is obvious that if, in the process of going from eqn (28) to (29), the fifth principle is not used, except when strictly necessary

$$\mathcal{E} = B_1 \omega \dots \omega B_v \quad (31)$$

but the example of eqn (29) shows that even this expression can be simplified to

$$\mathcal{E} = B_{\mu_1} \omega \dots \omega B_{\mu_e} \quad (32)$$

which contains only some of the prime implicants.

For this end perform again the simplification process described at the beginning of this paragraph. Starting from eqn (28) and considering each  $\mathcal{M}_i$ ; let  $B_{j_1}, \dots, B_{j_t}$  be those prime implicants for which

$$\mathcal{M}_i \propto B_j$$

There is at least one among the  $B_j$  which satisfies this condition. Consider the propositions  $p_j =$  the prime implicant  $B_j$  appears in eqn (32). Since the simple expression  $\mathcal{M}_i$  appears in eqn (28), the proposition

$$q_i = p_{j_1} v \dots v p_{j_t} \quad (33)$$

must be true.

Forming the proposition

$$P = q_1 \& \dots \& q_t \quad (34)$$

By expanding eqn (34)

$$P = R_1 v \dots v R_w \quad (35)$$

where  $R_i$  are propositions of the form

$$p_{\sigma_1} \& \dots \& p_{\sigma_t} \quad (36)$$

which means that

$$\mathcal{E} = B_{\sigma_1} \omega \dots \omega B_{\sigma_t} \quad (37)$$

However,  $B_{\sigma_1}, \dots, B_{\sigma_t}$  may be non-different, and therefore the application of the fifth principle will lead to a simplified form of  $\mathcal{E}$ .

Such a form corresponds at each  $R_i$ .

VII. The simplification problem may be stated not only for the canonical forms

$$[L_{\alpha_1}(x_1)\theta \dots \theta L_{\alpha_n}(x_n)] \omega \dots \omega [L_{\alpha_{r_1}}(x_1)\theta \dots \theta L_{\alpha_{r_n}}(x_n)]$$

but also for other normal forms, such as

$$[L_{\alpha_1}(x_{\beta_1})\theta \dots \theta L_{\alpha_{r_1}}(x_{\beta_{r_1}})] \omega \dots \omega [L_{\alpha_{r_1}}(x_{\beta_{r_1}})\theta \dots \theta L_{\alpha_{r_p}}(x_{\beta_{r_p}})]$$

since, by applying eqn (4), the missing variables can be introduced into every simple expression.

A parallel is drawn between the classical example exposed at the beginning of section V and another similar one.

*Classic example*

$$\begin{aligned} \mathcal{E} &= xy \cup \bar{y}z \cup xz = xyz \cup xy\bar{z} \cup \bar{y}z \cup xz \\ &= xyz \cup xy\bar{z} \cup x\bar{y}z \cup \bar{x}\bar{y}z \cup xz \\ &= xyz \cup xy\bar{z} \cup x\bar{y}z \cup \bar{x}\bar{y}z \cup xyz \cup x\bar{y}z \\ &= xyz \cup xy\bar{z} \cup x\bar{y}z \cup \bar{x}\bar{y}z \end{aligned}$$

*New example*

$$\begin{aligned} \mathcal{E}^* &= (x \perp y) \perp (\bar{y} \perp z) \perp (x \perp z) = (x \perp y \perp z) \\ &\quad \perp (x \perp y \perp \bar{z}) \perp (\bar{y} \perp z) \perp (x \perp z) \\ &= (x \perp y \perp z) \perp (x \perp y \perp \bar{z}) \perp (x \perp \bar{y} \perp z) \\ &\quad \perp (\bar{x} \perp \bar{y} \perp z) \perp (x \perp z) \\ &= (x \perp y \perp z) \perp (x \perp y \perp \bar{z}) \perp (x \perp \bar{y} \perp z) \perp (\bar{x} \perp \bar{y} \perp z) \\ &\quad \perp (x \perp y \perp z) \perp (x \perp \bar{y} \perp z) \\ &= (x \perp y \perp z) \perp (x \perp y \perp \bar{z}) \perp (x \perp \bar{y} \perp z) \perp (\bar{x} \perp \bar{y} \perp z) \end{aligned}$$



304/4

The expression considered is of the form

$$\mathcal{E} = M_1 \cup M_2 \cup M_3 \cup M_4$$

$$\mathcal{E}^* = M_1^* \cup M_2^* \cup M_3^* \cup M_4^*$$

with

$$M_1 = xyz$$

$$M_1^* = x \perp y \perp z$$

$$M_2 = xy\bar{z}$$

$$M_2^* = x \perp y \perp \bar{z}$$

$$M_3 = x\bar{y}z$$

$$M_3^* = x \perp \bar{y} \perp z$$

$$M_4 = \bar{x}\bar{y}z$$

$$M_4^* = \bar{x} \perp \bar{y} \perp z$$

According to the fifth principle  $(M_1, M_2), (M_2, M_3), (M_3, M_4)$  are introduced

$$M_{(12)} = xy$$

$$M_{(12)}^* = x \perp y$$

$$M_{(13)} = xz$$

$$M_{(13)}^* = x \perp z$$

$$M_{(34)} = \bar{y}z$$

$$M_{(34)}^* = \bar{y} \perp z$$

The prime implicants are

$$B_1 = M_{(12)}$$

$$B_1^* = M_{(12)}^*$$

$$B_2 = M_{(13)}$$

$$B_2^* = M_{(13)}^*$$

$$B_3 = M_{(34)}$$

$$B_3^* = M_{(34)}^*$$

$$\mathcal{E} = B_1 \cup B_2 \cup B_3 \text{ is eqn (29)}$$

$$M_1 \subset B_1$$

$$M_1^* \subset B_1^*$$

$$M_1 \subset B_2$$

$$M_1^* \subset B_2^*$$

$$M_2 \subset B_1$$

$$M_2^* \subset B_1^*$$

$$M_3 \subset B_2$$

$$M_3^* \subset B_2^*$$

$$M_3 \subset B_3$$

$$M_3^* \subset B_3^*$$

$$M_4 \subset B_3$$

$$M_4^* \subset B_3^*$$

$(p_1vp_2) \& p_1 \& (p_2vp_3) \& p_3 = p_1 \& p_3$  is true

$$\mathcal{E} = B_1 \cup B_3$$

$$\mathcal{E}^* = B_1^* \cup B_3^*$$

$$= xy \cup \bar{y}z$$

$$= (x \perp y) \perp (\bar{y} \perp z)$$

To sum up, in order to apply Quine's method, it is necessary that the first to the fifth principles be valid.

VIII. On other occasions, the another has drawn attention to the fact that the multiplicative elements occur in circuits with contacts and relays. A few examples follow.

(a) In real operating conditions, the armature with contacts does not change suddenly from the attracted or the repulsed position. There exists also an intermediate position, in which the normally open contacts as well as the normally closed ones are open (the 'break before make' relays) (Figure 1) or else the normally open contacts as well as the normally closed ones are closed (the 'make before break' relay) (Figure 2).

(b) The polarized relays whose armatures possess three possible positions, namely, resting, attraction and repulsion positions.

(c) The codified relays: examples (Figures 3, 4, 5) of codified four-positional relays, taken from the book of Keister-Ritchie-Washburn, and the example given by Ivanin (Figure 6).

(d) The 'step by step' searcher or selector.

To each element a number of contacts can be associated, namely:

(i) In real operation of the relays  $X$  of the 'break before make' type, there exist normally open contacts  $\varphi_1^0(X)$  and normally closed contacts  $\varphi_0^0(X)$ ; in real operation of the 'make before break' relays  $X$  there are also normally open contacts  $\varphi_1^s(X)$  and normally closed contacts  $\varphi_0^s(X)$ .

(ii) The polarized relays  $X$  have contacts for the attraction position  $\varphi_1(X)$  and contacts for the repulsion position  $\varphi_2(X)$ ; some of the polarized relays also have, contacts for the resting position  $\varphi_0(X)$ , but some others, with an unstable neuter, lack such contacts.

(iii) The codified relays  $X$  possess several types of contacts.

(iv) The selector  $S$  with  $\nu$  steps has the brush contacts  $\varphi_0(S), \dots, \varphi_{\nu-1}(S)$ .

IX. To each  $n$ -positional element, two sets of  $n$  elements are associated:

(a) The ring of residue classes modulo  $n$

$$\mathcal{I}|(n) = (0, 1, \dots, n-1) \quad (38)$$

with two operations: the addition and multiplication modulo  $n$  denoted by  $+$  and  $\cdot$ .

(b) The  $n$ -valent Lukasiewicz algebra

$$L_n = \left( 0, \frac{1}{n-1}, \dots, \frac{n-2}{n-1}, 1 \right) \quad (39)$$

with the natural order relation and with the operations

$$\begin{aligned} a \cup b &= \max(a, b) \\ a \cap b &= \min(a, b) \end{aligned} \quad (40)$$

The Lagrange functions are denoted by  $L_\alpha(x)$

$$\begin{aligned} L_\alpha(\alpha) &= 1 \\ L_\alpha(\beta) &= 0, \alpha \neq \beta \end{aligned} \quad (41)$$

with  $\alpha \in \mathcal{T}/(n)$  respectively  $\alpha \in L_n$ . The dual functions  $\bar{L}_\alpha(x)$  are introduced with

$$\begin{aligned} \bar{L}_\alpha(\alpha) &= 0 \\ \bar{L}_\alpha(\beta) &= 1, \alpha \neq \beta \end{aligned} \quad (42)$$

There is a Lagrange interpolation formula in  $\mathcal{T}/(n)$

$$f(x_1, \dots, x_n) = \sum f(\alpha_1, \dots, \alpha_n) L_{\alpha_1}(x_1) \dots L_{\alpha_n}(x_n) \quad (43)$$

and two interpolation formulae in  $L_n$

$$f(x_1, \dots, x_n) = \bigcup_{\alpha} [f(\alpha_1, \dots, \alpha_n) \cap L_{\alpha_1}(x_1) \cap \dots \cap L_{\alpha_n}(x_n)] \quad (44)$$

$$f(x_1, \dots, x_n) = \bigcap_{\alpha} [f(\alpha_1, \dots, \alpha_n) \cup \bar{L}_{\alpha_1}(x_1) \cup \dots \cup \bar{L}_{\alpha_n}(x_n)] \quad (45)$$

Giving

$$[L_\alpha(x)]^2 = L_\alpha(x) \quad (46)$$

$$L_\alpha(x) L_\beta(x) = 0, \alpha \neq \beta \quad (47)$$

$$L_\alpha(x) \cap L_\beta(x) = 0, \alpha \neq \beta \quad (48)$$

$$L_1(x) + \dots + L_n(x) = 1 \quad (49)$$

$$L_1(x) \cup \dots \cup L_n(x) = 1 \quad (50)$$

To  $\omega$  and  $\theta$  in eqn (13), the following values can be given

	$\omega$	$\theta$	
XI	+	·	
XII	∪	∩	
XIII	∩	∪	(51)

In  $\mathcal{T}/(n)$  respectively in  $L_n$ , principles I, II, III, and V are satisfied for the substitutions XI, XII, XIII of  $\omega, \theta$ .

To the fourth principle, must be substituted

IV\*. The following equality is true

$$\begin{aligned} & [L_0(z) \theta y_1 \theta \dots \theta y_r] \omega [L_1(z) \theta y_1 \theta \dots \theta y_r] \omega \\ & \dots \omega [L_{n-1}(z) \theta y_1 \theta \dots \theta y_r] \omega t_1 \omega \dots \omega t_s \\ & = (y_1 \theta \dots \theta y_r) \omega t_1 \omega \dots \omega t_s \end{aligned}$$

Principles I, II, III, IV\*, V allow application of Quine's method to multipositional elements.

In order to have a better understanding of the method one should introduce also the notions of formula of structure, function of work and functional equivalence in these general cases.

### Bibliography

The author has published the following volumes

- <sup>1</sup> *Teoria algebrica a mecanismelor automate* (Algebraic theory of switching circuits). 1959; Bucuresti; Editura tehnică
- <sup>2</sup> *Scheme cu comanda directă cu contacte și relee* (Combinational switching circuits with contacts and relays). 1959. București; Ed. Acad. Rep. populare Romine.
- <sup>3</sup> *Funcționarea în mai mulți timpi a schemelor cu relee ideale* (Sequential ideal operation of relays switching circuits). București, Ed. Acad. Rep. populare Romine.
- <sup>4</sup> *Circuite cu transistori* (Transistor switching circuits). 2 Vols. 1962. Ed. Acad. Rep. populare Romine.  
The interpolation formulae have been taken from
- <sup>5</sup> *Simplificarea circuitelor cu tuburi electronice, cu transistori și criotroni* (Simplifying the circuits containing electronic tubes, transistors and cryotrons). *Revue Math. pures et appl.* (1959) 497  
For the electronic tubes circuits see 5 and the authors work
- <sup>6</sup> ———— theory of ———— with electronic tubes, physic-mathematical copy. *German Acad. Science* 4 (37) (1961) p. 7  
For circuits with transistors see 4, 5, 6 and the author's work
- <sup>7</sup> Sur la théorie algébrique des circuits logiques à transistors. *Automatisme*, 7 (1962) 136  
For the cryotron circuits, see 5  
For the real operating of circuits with ordinary switching relays, see 1 and Gh. Ivanin's work
- <sup>8</sup> A supra teoriei algebrice a contactelor multipoziționale și aplicațiile ei la studiul contactelor reale (On the algebraic theory of multi-positional contacts and its application to the study of real contacts) *Bul. sti. Acad. Repub. rom. sec. sti. mat. fiz.* 7 (1955) 231.  
See also the author's work with Ivanin
- <sup>9</sup> Asupra funcționării schemelor cu butoni reali (On the operation of circuits with real buttons). *Bul. sti. Acad. Repub. rom. sec. sti. mat. fiz.* 7 (1955) 33, see also 1 and
- <sup>10</sup> Sur l'application des logiques à trois valeurs à l'étude des schémas à contacts et relais, *Act. Congr. int. Automat., Paris*, 14-24 June 1956, p. 48
- <sup>11</sup> Aplicațiile logicei trivalente în studiul funcționării reale a schemelor cu contacte și relee (The applications of trivalent logic to the study of real operating of circuits with contacts and relays). *Bul. matem. Soc. sti. mat. fiz. Repub. rom.* 1 (49) (1957) 197; *Bul. math. Soc. sci. math. phys. Roum.* 1 (49) (1957) 147
- <sup>12</sup> Sinteza schemelor cu contacte și relee în funcționare reală (Synthesis of circuits with contacts and relays, under real operating conditions). *Bull. matem. Soc. sci. math. phys. Roum.* 3 (51) (1959) 65  
A volume on the real operating conditions is in the press  
For circuits with polarized relays see 1 and the author's works
- <sup>13</sup> Sur la synthèse des schémas à relais polarisés. *German Academy of Science* 2 (1957) 121
- <sup>14</sup> Sur la théorie algébrique des mécanisme automatiques. Synthèse des schémas à relais polarisés. *Ber. int. matem. Koll., Dresden*, 22-27 November 1955, Aktuelle Probleme der Rechentechnik. Deutscher Verlag der Wissenschaft. 1957. Berlin  
For circuits with codified relays see 1 and the author's works

<sup>15</sup> Logica matematică și tehnica modernă. Logicele cu mai multe valori și circuitele cu contacte și relee (Mathematical logic and modern technique. The logics with several values and the circuits with contacts and relays). *Probleme filosofice ale științelor naturii*. 1960. ISRS Acad. Rep. populare Romine  
For circuits with selectors, see 1 and Gh. Ivanin's works

<sup>16</sup> Sinteza schemelor în care intră selectori (Synthesis of circuits with selectors). *Bul. sti. Acad. Repub. rom., ser. mat. fiz.* 8 (1956) 489; *Automat. Telemekh., Moscow* 19 (1958) 855  
<sup>17</sup> Sur un type de problème concernant les schémas à sélecteurs. *Acta Logica, Bucharest* (1958) 187

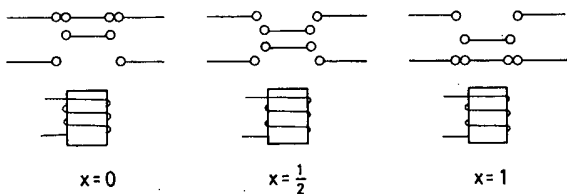


Figure 1

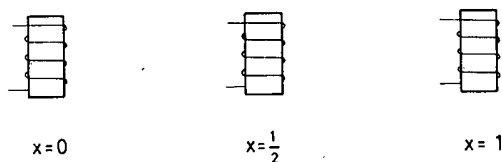
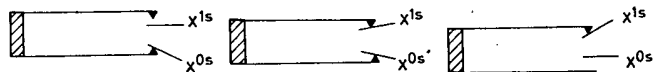


Figure 2

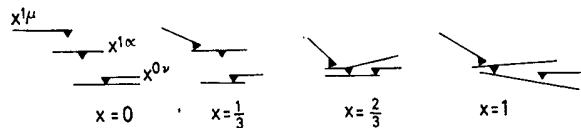


Figure 3

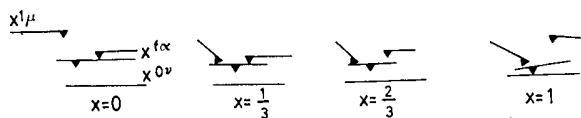


Figure 4

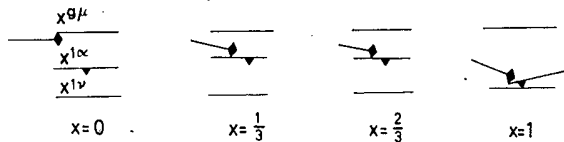


Figure 5

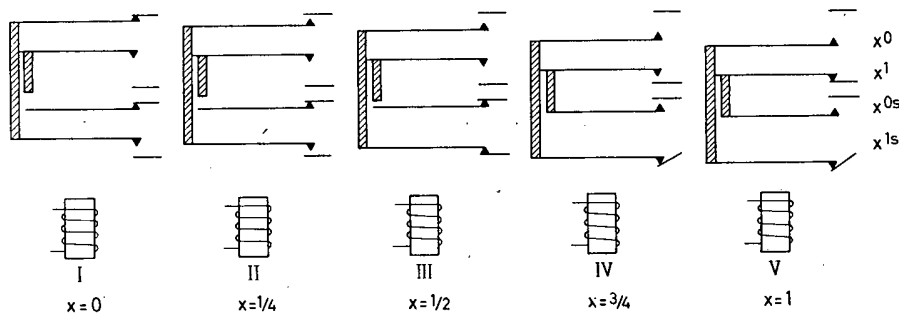


Figure 6

# Outline of a Control Theory of Prosthetics

R. TOMOVIC *- yug -*

## Introduction

Prosthetics has been treated, until recently, mainly as a branch of medicine; its relations with engineering were limited to mechanical and some electrical studies. The progress of automatic control has brought up the question of its application to prosthetics<sup>1-3</sup>. First achievements in this direction are very promising yet they were limited in scope. The intention here is to examine, on a much broader basis, the role in prosthetics to be played by automatic control theory. Instead of studying a definite prosthetic device, attention is directed at general scientific principles involved in the design of prosthetics. It is shown that automatic control can contribute in an important way to laying down a sound basis for improved design of prosthetic devices; prosthetic devices here being defined as applied to human extremities.

There is no doubt that many aspects of hand and foot prosthetics are different; but, looking from a deeper point of view, there are also important aspects in common. These common problems become specially evident when using control theory approach. In addition, the common ties relating the prosthetics of human extremities and remotely controlled manipulators or vehicles come clearly to the foreground in this way. A deeper insight into the common problems of all these fields has hardly existed. However, the task of handling materials in hostile environment is currently getting more and more important so that the solution of this problem, in itself, is of considerable interest. Thus, in many instances where a prosthetic device is mentioned in this paper it should be extended to remotely controlled manipulators as well.

## The Control Problem

The problem of controlling a prosthetic device can be treated in a general way. However, to make the understanding easier, consider the hand and arm control. Being even more specific, the process of lifting an object of arbitrary shape will be analysed. A closer examination shows that in the above action three different levels of control can be found. The first loop involves visual feedback and takes care of hand positioning with relation to the object. Since the arm consists of several mechanically independent units (upper arm, elbow, lower arm, wrist joint) which have their own degrees of freedom, the need arises to coordinate the movements of individual arm parts to perform the positioning action as a whole. Finally, when the hand has been brought into touch with the object, the grasping action can start.

From the control point of view, the grasping action can be divided into the following phases: (a) hand adjustment to the arbitrary shape of the object, (b) locking of the hand in the hold position, and (c) adjustment of the pressure to keep the object in the hand.

Again, it should be remembered that this is not the explanation of the biological control system, which must be treated by other means. At best, only intuitive ideas of how to look for a better explanation of biological phenomena may be gathered in this way.

Consider first problem (a), i.e., automatic hand adjustment to the shape of an object. It is well known that this problem has hardly been solved in the existing prostheses or remotely controlled manipulators. Some thought has been given to the use of electronic computers for this purpose. In certain automatic production lines the objects of strictly limited shapes must be handled. In such a case a stored programme computer directing the manipulator may represent the solution, but it is clear that this solution becomes easily obsolete if shapes are varied to a larger extent or if they are not known in advance. Using this example it will be shown that by treating the problem from a completely different point of view the general case of objects of arbitrary shapes can be solved in a simple way. The basis for this solution will be found in communication and control theory.

When studying a control system one usually begins with the block diagram explaining its structural set-up. Thus a servo-mechanism may be considered as having the following elements: input, amplifier, stage, actuator. This is represented in *Figure 1*. Assuming that the artificial hand with its fingers represents a positioning servo-mechanism<sup>1</sup> one can ask what is the basic difference compared with the diagram of *Figure 1*. In *Figure 2* the diagram of prosthesis control is shown. In contrast to *Figure 1* it is seen that now the control signal source is linked via the communication channel to the actuator. In addition to remotely sent control signals, there is also a local feedback loop with input signals produced on the spot. With the existing engineering knowledge it is not difficult to reproduce mechanically the form and movements of the hand, but the real problem is how to supply adequate control signals. As is known, the supply of control signals by muscular movements and by bio-electrical means was not very successful. With all the improvements in prosthesis design only a few elementary hand movements could therefore be reproduced. The adequate supply of control signals for prostheses and manipulators is still a very important problem to be solved.

A closer study of the role of skin sensitivity to pressure, temperature and other stimuli may give a hint for the solution. Returning to *Figure 2* one can easily understand that it is highly desirable to obtain maximum hand flexibility with minimum signals supplied by the remote source. In the case of long communication channels, this will mean a reduced channel capacity if remote handling of materials is in question, or a reduced burden on the part of the amputee if prosthesis control is concerned. In order to discuss the question in a more precise manner, consider the set of signals  $S_1$  which must be

390/2

provided by the control source in order to position the fingers. Taking each finger as a separate automatic positioning system, the source must provide five continuously varying signals

$$e_s = f_s(x), \quad s = 1, 2, \dots, 5$$

where  $x$  is the parameter defining finger position. For simplicity reasons the finger is considered as a dynamic system with one degree of freedom, i.e., as a rigid body with no lateral movements and phalanges. Even in this case the set of control signals

$$S_1 = \{e_1, e_2, e_3, e_4, e_5\} \quad (1)$$

is quite complex. Since the word complexity of the set  $S_1$  is of intuitive nature it needs additional explanation. Remembering that the control of the prosthesis is also a communication problem, the complexity of  $S_1$  will be measured by the information content of signals  $e_s$ . Designate the information content of  $e_s$  by  $i_s$ , so that the information content  $I_1$  corresponding to  $S_1$  is

$$I_1 = \sum_{i=1}^5 i_s$$

An explicit value of  $I_1$  is not needed here. Remember only that human hand control consists of 24 different muscle groups.

It is clear that the prosthesis control problem cannot be solved in a satisfactory manner by conscious control signals. Such control is in evident contrast with the basic design condition for prostheses, and manipulators to keep  $I_1$  as low as possible. The solution of keeping  $I_1$  low by reducing hand flexibility is naturally not acceptable since it badly limits the performances of the prosthesis. In the absence of a better solution this has been done in the existing models. Thus a new approach is needed. The first results, taking into account the requirement that  $I_1 = I_{\min}$ , can be found in previous papers<sup>2</sup>. The fundamental idea is to keep  $I_1$  low without affecting hand performances. The problem has been solved by dividing the control signals into two sets:  $S_1$  and  $S_2$ . The signals  $S_1$  are centrally or remotely produced and transmitted *via* the communication channel, while signals  $S_2$  are locally generated, i.e., at the receiver end. The information content of  $S_2$  is  $I_2$  and the total information available

$$I = I_1 + I_2 \quad (2)$$

The simple fact of dividing control signals in  $S_1$  and  $S_2$  allows for a great reduction of channel capacity, and consequently keeps the 'burden' on the central control source low without affecting hand performance.

Eqn (2) needs explanation, namely, the control signals  $S_2$  should be generated in such a way that the required adaptation to the shape of objects is obtained. In order to understand how this can be done two new concepts must be introduced. In the first place a topologically equivalent mechanical system of hand and finger movements is needed. The aim is to obtain a simple and symmetric mechanical structure equivalent to the human hand with regard to its capacity to handle objects of various shapes. *Figure 3* shows such an equivalent and symmetric model, consisting of five elastic segments which can be rotated around the central ring. An elastic segment is required for holding objects against a rigid segment, with two or more sections rotating around individual joints. Each segment is provided with a fixed cable along which a central force  $P$  can be applied. Actually all five cables may represent five branches

of a central cable so that the model is activated by the application of just one force  $P$ :

$$P_0 = 5P$$

The force  $P_0$  is directed perpendicularly to the plane of the drawing. It is further supposed for simplicity reasons that the joints of all segments lay on the perimeter of the circle with radius  $R$ .

The elastic mechanical model of *Figure 3* is, for study purposes, equivalent to the human hand, and it is easy to see that objects of arbitrary shape can be grasped by this system if one first assumes that a ball, of radius  $r$ , is placed in the centre of the mechanical model. The only condition which should be satisfied in order that the object remains in the 'hand' is

$$r \leq R \quad (3)$$

If the friction between surfaces is assumed, the condition (3) becomes less strict.

The problem of holding the object of arbitrary shape with the model of *Figure 3* can always be reduced topologically to condition (3). In the case of irregular shape the radius  $r$  in (3) means the radius of the smallest sphere described around the object. It should be remarked that the uneven disposition of segments along the perimeter of the central ring allows for holding of objects of elongated shapes like pencils, for instance.

Another new concept, which in grasping actions helps to reduce  $I_1$  in eqn (2), is the sensitivity of the actuator to external stimuli. For instance, the instant of touching an object with hand prosthesis must be recorded not only by visual signals but also by pressure sensitive elements. The application of pressure-sensitive elements to prostheses is quite simple<sup>3</sup>. However, this new concept facilitates greatly the control problem by reducing the information content of the central control unit. The application of sensory elements to prostheses and remotely controlled manipulators adds actually a new local information source which can be used for object identification or local motor control. The information content of signal source  $I_2$  is thus increased while  $I_1$  is kept low. This redistribution of the information content of control signals  $I_1$  and  $I_2$  is not affecting hand performance but saves channel capacity and reduces the need for frequent intervention of the central control unit. The hand being demonstrated at this conference handles, therefore, objects of completely arbitrary shapes requiring, however, only one bit of information being produced by the amputee.

### Object Identification

In the previous paragraph it has been explained how the special mechanical structure of the actuator of the positioning servo-mechanism simplifies the remote control of the prosthetic device or manipulator. The coverage of the control part of the servo-mechanism by sensory elements served the same purpose. The considerations here will be limited to pressure sensitive elements, although the basic conclusions apply to temperature, radioactive or other type of sensory transducers.

The difference between a sensory element and an ordinary transducer should be clearly defined. The basic characteristic of a transducer is to establish a one-to-one correspondence between two different physical quantities. In most cases, and this will be understood here, the output of the transducer is

electrical quantity, voltage or current. The equation of the transducer may be written in the following form

$$e = f(p) \quad (4)$$

where  $e$  is the voltage, and  $p$  the pressure in our case. A sensory element or surface as understood here, differs in some important aspects from the above definition of the transducer. The symbolic representation of the pressure sensitive surface is seen in *Figure 4*. The surface represented in *Figure 4* should be understood as a piece of the 'skin' with pressure sensitive cells. Each cell, upon touch, gives a voltage output proportional to  $e_{rs}/j = f(p)$ . It is not important that all functions  $f(p)$  be strictly identical. A practical version of such a sensory surface can be realized in different ways.

The difference between the conventional transducer and sensory surface can best be grasped by establishing the equation of the sensory surface

$$e = f(p, x, y) \quad (5)$$

When compared with eqn (4), one sees immediately that the sensory elements provide, in addition to intensity of the stimulus, information about the spot of its application. Thus, from the point of view of information content, new dimensions are added when a set of transducers is geometrically ordered in space.

Eqn (5) needs a refinement which is quite important. Namely, the set of transducers is discrete so that the equation of the sensory surface corresponds exactly to the following form:

$$e = f(p, r \Delta x, s \Delta y) \quad \begin{matrix} 0 \leq r \leq n \\ 0 \leq s \leq m \end{matrix} \quad (6)$$

An important conclusion obtained from eqn (6) is that the resolution rate in  $x$  and  $y$  is finite. This fact corresponds with the actual situation in biological systems where resolution rates of sensory elements are always finite. How this fact allows extraction of important informations about the object held in the hand, is now explained; only informations flowing directly through the communication channel of *Figure 2* are in mind, and not those which can be obtained, for instance, by direct or remote visual examination of the object.

The first kind of object identification made possible by sensory elements eqn (6) regards the shape. When the artificial hand, covered with the pressure sensitive surface, is closed around an object, a one-to-one correspondence between the electrical waveform and the shape of the object is established. This fact can best be understood by taking two characteristic geometric forms. It is assumed that objects to be identified by the artificial hand have circular and rectangular cross-sections as represented in *Figure 5(a)* and *(b)*. Associated waveforms for the two types of shapes are seen in *Figure 6*. It has also been assumed that  $r$  is variable but  $s = s_0$ , is fixed. The restriction is not important. If different  $y$  sections of the hand are taken then the waveforms of *Figure 6* become functions of  $s$  as well. They may or may not be identical, depending on the fact if the object keeps cross section unchanged along  $y$  axis. The correspondence of electrical waveforms in *Figure 6* with object shapes in *Figure 5* is evident from eqn (6). Namely, in the case of the circular cross section more or less the whole surface is equally exposed to pressure. Thus,  $e = \text{const.}$  for all  $r$ . The constant voltage output for circular cross section requires an even hand surface, but a slightly uneven hand surface will not

affect the object identification. In this case the parts of the hand which are not in contact with the object will provide for a gap in *Figure 6(a)*. Since there is a large amount of redundancy in this identification process, the general information obtained will be adequate even if an ideally flat sensory surface is not assumed. In the case of cross section with the edges, however, there will be, ideally, two subsets of  $r$

$$\begin{aligned} R_1 &= r_\alpha \quad s = s_0 \\ R_2 &= r_\beta \quad s = s_0 \\ R_1 + R_2 &= R = \{r_\alpha, r_\beta\} \end{aligned} \quad (7)$$

The location of  $r_\alpha$  corresponds geometrically to the spots of hand contacts with the edges of the object, and  $r_\beta$  is the complementary set of  $R$  with respect to  $R_1$ . Now, the equation of the characteristic waveform reads

$$\begin{aligned} \text{if } r \in R_1 \quad e &= 1 \\ \text{if } r \in R_2 \quad e &= 0 \end{aligned}$$

Examination shows how this information regarding object shape can be sent back to the central control place by the communication channel of *Figure 2*. The problem is technically trivial since one needs a two-dimensional scanning system. To make it clear, refer to the sensitive hand surface in *Figure 4*. Since for the purpose of pure shape identification only the distinction between activated and non-activated spots is important, the output of the artificial hand surface is a binary matrix  $e_{rs} = 1$  or  $0$ , according to eqn (7). The transmission to the remote control place is simply solved, for instance, by a magnetic core selection matrix.

A further interesting tactile information may be obtained from the pressure sensitive surface. That is, if the number of object edges is increased the distinction between circular and polygonal shapes will be lost. If the resolution rate is correspondingly increased, i.e.,  $\Delta x \rightarrow 0$ , one will be able to map into the electrical form the roughness of the surface with which the hand is in contact. Thus, the waveform of *Figure 6(b)* contains both the information about shape or roughness of the object depending on the resolution of the sensory surface, i.e., the magnitude of  $\Delta x$ . Although these two tactile effects are distinct from the sensory point of view, mathematically they are equivalent; the only difference being the order of magnitude of  $\Delta x$ . Actually both effects are the consequence of the discrete structure of the sensory surface. An important condition for practical realization of such a discrete pressure sensitive surface is a high resolution rate of individual transducer elements. This implies the mechanical isolation of the adjacent elements so that they can react in a distinct way, although being geometrically close. One is led therefore to the design of very thin elastic surfaces.

At the beginning of the paper it was outlined that the principles exposed here have general significance. Besides their theoretical value of giving mathematical insight into the problem of remote object identification without visual feedback, there are other fields of application. Namely, in the existing foot and leg prostheses the role of the shape identification of the ground for control purposes has been completely neglected. However, the application of the pressure sensitive discrete surface of *Figure 4* allows easily the coordination of different phases of human gait according to which part (front or back) of the foot is in touch with the ground; further, hitting of

obstacles can easily be detected in the electrical form and used for control purposes as well. The idea of object identification by sensory elements exposed here can therefore be exploited for variety of control purposes.

In the design of hand prostheses evidently there is no need for object identification by tactile feedback since it is more simple to use visual information for this purpose. However, in the remote handling problems, due to the fact that communication capacities may become critical, the relative importance of tactile and visual feedback may change. It is hard to give a precise evaluation since all the existing designs have relied exclusively on visual feedback (television). A general selection criterion is not possible since the application conditions must be taken into account. However, for simple remote identification problems (size, shape, weight) sensitive surfaces may serve the purpose. It should be remembered that according to Figure 6 the tactile feedback needs just a few  $y$  lines to be sent over communication channel. Thus a great reduction of channel capacity is possible in certain instances. In other instances it may occur that a combined identification system represents the best solution. As has been written, it is not the intention to discuss the absolute merits of visual or tactile information feedback, but to stress the fact that more general identification methods when designing remote handling control systems should be used.

**Conclusions**

In this paper several questions have been raised. First of all,

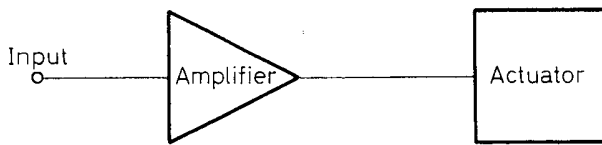


Figure 1

the importance of improving the actuators used in servo-mechanisms has been shown. It is proposed to solve this problem in an unconventional way by using special mechanical shapes of the actuator covered with sensitive surface. Such an approach has the merit of showing the transition phases of a conventional positioning servo-mechanism to an artificial hand.

In addition, the identification problem of automatic control theory is presented in a new way. The notion of the transfer function for linear systems, or other methods of identification for non-linear systems are in current use. However, in the future development of automatic control systems many situations may arise where the identification problem cannot be solved satisfactorily by the existing methods. Object identification by shape, surface characteristics and other ways such as those occurring in biological systems will also be needed in engineering systems.

Looking at the identification problem in engineering from a broader point of view allows the synthesis of new cybernetic control systems which can duplicate functions of biological structures in a very efficient way.

**References**

- <sup>1</sup> TOMOVIC, R. Human hand as a feedback system, *Automatic and Remote Control*. 1960, 624-628. London; Butterworths
- <sup>2</sup> KOBRINSKI *et al.* Problems of bioelectric control, *Automatic and Remote Control*. 1960. London; Butterworths
- <sup>3</sup> TOMOVIC, R., and BONI, G. An adaptive artificial hand, *Trans. IRE*, vol. PGAC (1962)

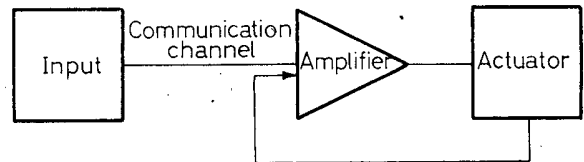


Figure 2

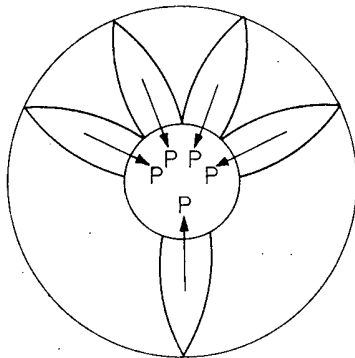


Figure 3

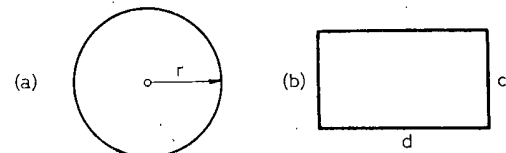


Figure 5

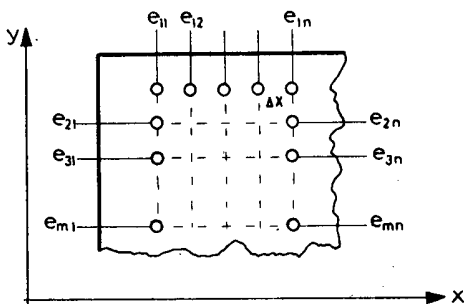


Figure 4

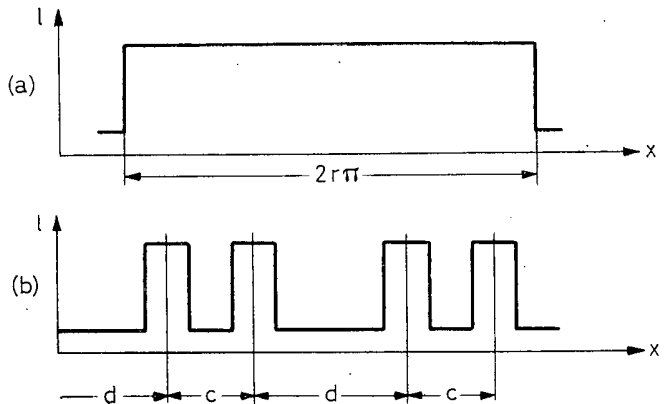


Figure 6

STAT

**Page Denied**