

On the Sensitivity of Linear State-Space Systems

LOTHAR THIELE, MEMBER, IEEE

Abstract—This paper contains new measures to describe the transfer function sensitivity of state-space systems with respect to value and parameter perturbations. These measures are related to the newly defined generalized Gramian matrices. The value respecting the parameter variations contains the sensitivity at discrete frequency points, pole and zero sensitivities and the integral sensitivity as special cases. A general relation to the variance of the weighted output noise can be obtained in the case of a perturbed realization which is l_2 -scaled under a general non-white input process. The complete class of representations with minimum sensitivity and noise is given. The corresponding necessary and sufficient conditions lead to an analytic design of optimal state-space systems.

I. INTRODUCTION

ONE POSSIBILITY to convert linear systems or analytic functions into abstract mathematical processes consists in the state-space representations. They are characterized by the description of the outer behavior of physical systems with constant matrices and therefore, they establish a direct relation to the methods of linear algebra.

The state-space representation of a given linear system can be used in order to perform a pole-zero determination [1], a stability test by the use of Lyapunov functions, e.g., [2], a test concerning the passivity and losslessness [3] or a cascade factorization [4]. The solution of these problems using the methods of linear algebra is recommended from a numerical point of view as powerful software packages are available for the evaluation of standardized problems, e.g., [5]. It is well known [4], [6] that the numerical accuracy of the particular result depends crucially on the condition of the chosen state-space representation. Therefore, it is desirable to define appropriate measures which represent the sensitivity of the realized transfer function with respect to value perturbations (or internal noise sources) and to parameter variations. But up to now, results concerning the evaluation and analytic optimization of these properties are still sparse [6], [7].

Direct implementations of the signal flow graph which corresponds to the state equations received considerable attention in the past, e.g. [8], [9]. By the use of the equivalence transformation defined within this class of state-space representations it is possible to determine special realizations which are l_2 -scaled under a white noise input process. These realizations have the minimal norm of the output noise spectrum [8], [10], the minimum pole

sensitivity [9] or the minimum upper bound of an integral sensitivity measure [11], [12]. Moreover, normal realizations can be implemented to be stable under finite word-length effects [9]. As Kung has shown [13] that minimum noise realizations of arbitrary order satisfy $\|A\|_2 < 1$ where A denotes the system matrix they can be realized to be free of overflow and granularity limit cycles [14]. Efficient physical realizations of the state equations can be received by the use of vector arithmetic or a direct implementation of the continued matrix-vector product in form of a VLSI circuit. Modern concepts of VLSI design enable a direct implementation of algorithms of linear algebra and thus lead to a feasible realization of digital signal processors.

It can be stated that no results are available which generalize the above mentioned methods in the design of digital realizations to determine analytically state-space realizations which are optimal under the consideration of the actual surrounding of the system. In order to achieve a design of optimal state-space realizations according to practical aspects it is desirable to take into account an arbitrary input spectrum and the norm of the weighted output noise spectrum. In addition to the mentioned criteria of optimization it is useful to design realizations having minimum sensitivity of the transfer function with respect to parameter variations [15], [16]. In the case of a general VLSI circuit carrying out the continued matrix-vector product of the state equations it is also desirable to use a small coefficient word length. In [11], [12] only an upper bound of a special integral sensitivity measure could be minimized.

II. THE GENERALIZED GRAMIAN MATRICES

In this paper we are concerned with the state equations of a discrete time, single input, single output system of the form

$$\begin{bmatrix} x(k+1) \\ y(k) \end{bmatrix} = \begin{bmatrix} A & b \\ c & d \end{bmatrix} \begin{bmatrix} x(k) \\ u(k) \end{bmatrix} \quad (1)$$

with the input $u(k) \in \mathbb{C}$, the output $y(k) \in \mathbb{C}$, the state vector $x(k) \in \mathbb{C}^{n \times 1}$ and the time-invariant matrices $A \in \mathbb{C}^{n \times n}$, $b \in \mathbb{C}^{n \times 1}$, $c \in \mathbb{C}^{1 \times n}$, and $d \in \mathbb{C}$. In order to describe the global relation between the internal and external behavior of the state representation (1), the following z-transformed equations are given:

$$\begin{bmatrix} X(z) \\ Y(z) \end{bmatrix} = \begin{bmatrix} \Theta(z) & F(z) \\ G(z) & H(z) \end{bmatrix} \begin{bmatrix} zx(k=0) \\ U(z) \end{bmatrix} \quad (2)$$

Manuscript received May 22, 1985; revised November 18, 1985.

The author is with the Institute of Circuits and Systems, Department of Electrical Engineering, Technical University of Munich, D-8000 Munich, Germany.

IEEE Log Number 8607495.

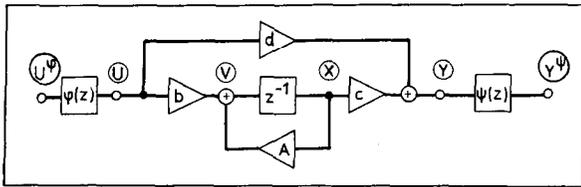


Fig. 2. Signal flowgraph interpretation of the generalized Gramian matrices.

procedure coincides in its interpretation with the method given in [8] for the synthesis of block optimal cascaded structures. The corresponding matrix equations can be solved using the iterative algorithm of Smith [18] or the efficient and accurate Bartels–Stewart algorithm [19].

The second procedure makes possible the determination of the matrices K^Φ and W^Ψ if the integrands in Definition 1 have to be taken into account only at l discrete frequency points z_k with the corresponding weights ϕ_k and ψ_k . In this case, the generalized Gramian matrices can be given to be

$$K^\Phi = \sum_{k=1}^l F(z_k) F(z_k)^h |\phi_k|^2,$$

$$W^\Psi = \sum_{k=1}^l G(z_k)^h G(z_k) |\psi_k|^2.$$

In the case of general weighting functions an algorithm can be developed which only requires the solution of $2n + 2$ integrals and that of two matrix equations [20].

III. ANALYSIS OF STATE-SPACE SYSTEMS

III.1. Variances of the Internal Signals

Now, we are looking for measures to represent the probability of overflow in the case of a statistical character of all signals. The subsequent analysis proceeds from a stationary, colored noise process with the noise power density spectrum $|\Phi(z)|^2$ which is an appropriate model of the actual environment of the system. By the use of Definition 1 it can easily be seen that the matrix K^Φ is the covariance matrix of the state vector $x(k)$ under a stationary input process with the spectral density $|\Phi(z)|^2$.

Although the consideration of a general non-white input process is important from a practical point of view if, e.g., a digital realization is embedded in a larger system, often only rough bounds of the norm of the internal signals had been applied. Usually they are based on the Minkowski inequality

$$\sigma_{x_i}^2 = E\{|x_i(k)|^2\} \leq \max_{z \in \Gamma} \{|\Phi(z)|^2\} \|F_i(z)\|_2^2$$

where $E\{\cdot\}$ denotes the expectation operator. Now, the exact l_2 -norms can easily be computed using

$$\sigma_{x_i}^2 = k_{ii}^\Phi.$$

Thus in the case of a non-white input process an exact analysis and better design of state-space realizations can be

obtained concerning the sensitivity with respect to parameter and value variations [21].

III.2. Analysis of the Noise Behavior

In order to define a measure of the output error of a practical realization the following linear model of the internal deviations will be applied. A stationary, white vector process $\{v(k)\}$ with the corresponding covariance matrix Q is added to the state-space system according to Fig. 1 at the node 'V'. If not all parts of the noise spectrum at the output node 'Y' of the system are considered to be equally important or coloured noise sources with the noise power density spectrum $|\Psi(z)|^2$ are applied, a weighting function can be used which changes the output noise spectrum to

$$G(z)QG(z)^h |\Psi(z)|^2.$$

Therefore, the weighted output variance of the disturbed system can be evaluated to be

$$\sigma_y^2 = E\{|y^\Psi(k)|^2\} = \text{tr}(QW^\Psi).$$

In the case of uncorrelated noise sources at the node 'V' with $Q = I$, the weighted noise power gain g^Ψ can be given to be

$$g^\Psi = \text{tr}(W^\Psi). \quad (8)$$

III.3. Sensitivity of State-Space Representations

The sensitivity of the transfer function of a particular state-space model is an important criterion for the comparison of equivalent networks. In the case of a digital realization of (1) the quantization of the coefficients leads to a deviation of the nominal transfer characteristics as a finite number of different coefficient values causes a finite degree of accuracy of the transfer function. It is of particular importance to choose an original state-space model with low sensitivity if numerical algorithms for solving system theoretic problems are based upon this representation [4], [6], [7].

In the case of a state-space realization whose coefficients are exactly the elements of the state matrices A , b , c , and d the absolute sensitivities of the transfer function $H(z)$ with respect to variations of the coefficients can be given to be

$$S_{b_i}(z) = \frac{\partial H(z)}{\partial b_i} = c(zI - A)^{-1} e_i = G_i(z)$$

$$S_{c_i}(z) = \frac{\partial H(z)}{\partial c_i} = e_i(zI - A)^{-1} b = F_i(z)$$

$$S_{a_{ij}}(z) = \frac{\partial H(z)}{\partial a_{ij}} = c(zI - A)^{-1} e_i e_j^t (zI - A)^{-1} b$$

$$= G_i(z) F_j(z).$$

where e_i is the unit vector with i th element unity. The relation between sensitivities and partial transfer functions is well known and can be specialized on state-space realizations using a theorem given in [22].

In the following, only a statistical derivation of the defined sensitivity measure will be given, whereas a deterministic one could also be used. We assume that the element variations are small enough to utilize a linear approximation to the Taylor series expansion of $\Delta H(z)$:

$$\Delta H_A = \sum_{i=1}^n \sum_{j=1}^n S_{a_{ij}} \Delta a_{ij},$$

$$\Delta H_b = \sum_{i=1}^n S_{b_i} \Delta b_i,$$

$$\Delta H_c = \sum_{i=1}^n S_{c_i} \Delta c_i.$$

In the case of statistically independent variations of the coefficients with unit variances the frequency dependent variances of the transfer function can be obtained as

$$\Sigma_{\Delta h, b}^2(z) = E\{|\Delta H_b(z)|^2\} = \sum_{i=1}^n |S_{b_i}|^2 = G(z)G(z)^h$$

$$\Sigma_{\Delta h, c}^2(z) = E\{|\Delta H_c(z)|^2\} = \sum_{i=1}^n |S_{c_i}|^2 = F(z)^h F(z)$$

$$\begin{aligned} \Sigma_{\Delta h, A}^2(z) &= E\{|\Delta H_A(z)|^2\} = \sum_{i=1}^n \sum_{j=1}^n |S_{a_{ij}}|^2 \\ &= G(z)G(z)^h F(z)^h F(z). \end{aligned}$$

In order to give a measure which takes into account the weighted sensitivity behavior of the state-space representation in the whole frequency range the following definition is suggested:

Definition 2:

The following values are sensitivity measures of a state-space representation according to (1):

$$m_A^{\Phi\Psi} = \left\{ \frac{1}{2\pi j} \oint_{\Gamma} \Sigma_{\Delta h, A}^2(z) |\Phi(z)\Psi(z)| \frac{dz}{z} \right\}^2 = \|\Sigma_{\Delta h, A} \Phi\Psi\|_1^2$$

$$m_b^{\Psi} = \left\{ \frac{1}{2\pi j} \oint_{\Gamma} \Sigma_{\Delta h, b}^2(z) |\Psi(z)|^2 \frac{dz}{z} \right\} = \|\Sigma_{\Delta h, b} \Psi\|_2^2$$

$$m_c^{\Phi} = \left\{ \frac{1}{2\pi j} \oint_{\Gamma} \Sigma_{\Delta h, c}^2(z) |\Phi(z)|^2 \frac{dz}{z} \right\} = \|\Sigma_{\Delta h, c} \Phi\|_2^2$$

$$m^{\Phi\Psi} = m_A^{\Phi\Psi} + m_b^{\Psi} + m_c^{\Phi}.$$

The values of Definition 2 can be interpreted as statistical multiparameter sensitivity measures of the transfer function. As, e.g.,

$$m_b^{\Psi} = \left(\int_{-\pi}^{\pi} E\{|\Delta H_b(e^{jt})|^2\} |\Psi(e^{jt})|^2 dt \right)$$

and

$$m_A^{\Phi\Psi} = \left(\int_{-\pi}^{\pi} E^{1/2}\{|\Delta H_A(e^{jt})|^2\} |\Phi(e^{jt})\Psi(e^{jt})| dt \right)^2$$

the values of Definition 2 average the frequency dependent variances of the transfer function in the whole frequency range. The functions $\Phi(z)$ and $\Psi(z)$ enable the consideration of the variances in a designer-specified frequency band

or at some discrete frequency points. The differences in the definitions of $m_A^{\Phi\Psi}$ and, e.g., m_b^{Ψ} are caused by the analytic properties of $S_{a_{ij}}$ and S_{b_i} . The proposed measures can be optimized by an analytic procedure.

The following theorem relates the newly defined sensitivity measures to the generalized Gramian matrices of Definition 1.

Theorem 1:

The values of Definitions 1 and 2 are given. Then the relations

$$m_A^{\Phi\Psi} \leq \text{tr}(K^{\Phi}) \text{tr}(W^{\Psi}) \quad (9)$$

$$m_b^{\Psi} = \text{tr}(W^{\Psi}) \quad (10)$$

$$m_c^{\Phi} = \text{tr}(K^{\Phi}) \quad (11)$$

$$m^{\Phi\Psi} \leq \text{tr}(K^{\Phi}) \text{tr}(W^{\Psi}) + \text{tr}(K^{\Phi}) + \text{tr}(W^{\Psi}) \quad (12)$$

hold. Moreover, the equality signs in (9,12) are valid iff

$$\begin{aligned} \rho^2 |\Phi(z)|^2 F(z)^h F(z) &= |\Psi(z)|^2 G(z)G(z)^h, \\ z &\in \Gamma, \quad \rho \in \mathbb{R} \setminus \{0\}. \end{aligned}$$

Proof: At first, the relations (10) and (11) will be proved. By the use of m_b^{Ψ} of Definition 2 it can be concluded that

$$\begin{aligned} m_b^{\Psi} &= \frac{1}{2\pi j} \oint_{\Gamma} G(z)G(z)^h |\Psi(z)|^2 \frac{dz}{z} \\ &= \sum_{i=1}^n \frac{1}{2\pi j} \oint_{\Gamma} |G_i(z)\Psi(z)|^2 \frac{dz}{z} = \sum_{i=1}^n w_{ii}^{\Psi} = \text{tr}(W^{\Psi}). \end{aligned}$$

A similar proof can be given in order to obtain (11). The inequalities (9) and (12) are a consequence of the Cauchy-Schwartz inequality, as

$$\begin{aligned} m_A^{\Phi\Psi} &= \|\Sigma_{\Delta h, b} \Sigma_{\Delta h, c} \Phi\Psi\|_1^2 \\ &\leq \|\Sigma_{\Delta h, b} \Psi\|_2^2 \|\Sigma_{\Delta h, c} \Phi\|_2^2 \\ &= \text{tr}(W^{\Psi}) \text{tr}(K^{\Phi}). \end{aligned}$$

The equality sign is valid if and only if

$$\Sigma_{\Delta h, b}^2 |\Psi|^2 = \rho^2 \Sigma_{\Delta h, c}^2 |\Phi|^2, \quad z \in \Gamma, \quad \rho \in \mathbb{R} \setminus \{0\}. \quad \square$$

An appropriate choice of the weighting functions in Definition 2 leads to the specialization of the sensitivity measures on pole and zero sensitivities and on the sensitivity at discrete weighted frequency points. Therefore, common sensitivity measures can be made uniform and can be related to the Gramian matrices in the case of a state-space representation.

Let us suppose that the sensitivity of the considered state-space realization will be taken into account at the l discrete frequency points z_k with the corresponding weights w_k . The weighting functions are chosen so that the Gramian matrices can be given to be

$$K^{\Phi} = \sum_{k=1}^l F(z_k) F(z_k)^h |w_k|^2$$

$$W^{\Psi} = \sum_{k=1}^l G(z_k)^h G(z_k) |w_k|^2.$$

Then Definition 2 and Theorem 1 imply the following sensitivity measures:

$$\begin{aligned} m_A^{\Phi\Psi} &= \left\{ \sum_{k=1}^l \sum_{\Delta h, A} (z_k) |w_k|^2 \right\}^2 \\ &= \left\{ \sum_{k=1}^l \left(\sum_{i,j=1}^n \left| \frac{\partial H(z_k)}{\partial a_{ij}} \right|^2 \right)^{1/2} |w_k|^2 \right\}^2 \\ &\leq \text{tr}(K^\Phi) \text{tr}(W^\Psi) \end{aligned} \quad (13)$$

$$\begin{aligned} m_b^\Psi &= \sum_{k=1}^l \sum_{\Delta h, b}^2 (z_k) |w_k|^2 \\ &= \sum_{k=1}^l \sum_{i=1}^n \left| \frac{\partial H(z_k)}{\partial b_i} \right|^2 |w_k|^2 = \text{tr}(W^\Psi) \end{aligned} \quad (14)$$

$$\begin{aligned} m_c^\Phi &= \sum_{k=1}^l \sum_{\Delta h, c}^2 (z_k) |w_k|^2 \\ &= \sum_{k=1}^l \sum_{i=1}^n \left| \frac{\partial H(z_k)}{\partial c_i} \right|^2 |w_k|^2 = \text{tr}(K^\Phi). \end{aligned} \quad (15)$$

It will be shown now that the pole and zero sensitivities can be considered as a special case of (15). By the use of the pole-zero decomposition of a normalized transfer function according to

$$H(z) = \frac{\prod_{i=1}^p (z - \mu_i)}{\prod_{i=1}^q (z - \lambda_i)} = \frac{N(z)}{D(z)}$$

the following relation can be derived [23]:

$$\frac{\partial H(z)}{\partial e} = H(z) \left(\sum_{k=1}^q \frac{\partial \lambda_k / \partial e}{z - \lambda_k} - \sum_{k=1}^p \frac{\partial \mu_k / \partial e}{z - \mu_k} \right).$$

Choosing the frequency points z_k according to $|z_k - \mu_k| = \delta$, ($k=1, \dots, p$) and $|z_k - \lambda_{k-p}| = \delta$, ($k=p+1, \dots, p+q$) with $\delta \rightarrow 0$ it can be concluded that

$$\left| \frac{\partial \mu_k}{\partial e} \right| = \left| \frac{z_k - \mu_k}{H(z_k)} \frac{\partial H(z_k)}{\partial e} \right| = \left| \frac{D(\mu_k)}{N'(\mu_k)} \frac{\partial H(z_k)}{\partial e} \right|, \quad k=1, \dots, p$$

$$\left| \frac{\partial \lambda_{k-p}}{\partial e} \right| = \left| \frac{z_k - \lambda_{k-p}}{H(z_k)} \frac{\partial H(z_k)}{\partial e} \right| = \left| \delta^2 \frac{D'(\lambda_{k-p})}{N(\lambda_{k-p})} \frac{\partial H(z_k)}{\partial e} \right|, \quad k=p+1, \dots, p+q$$

where $D'(z) = \partial D(z) / \partial z$ and $N'(z) = \partial N(z) / \partial z$. Therefore, using the weights

$$|w_k|^2 = |\tilde{w}_k|^2 \left| \frac{D(\mu_k)}{N'(\mu_k)} \right| \quad \text{for } 1 \leq k \leq p$$

and

$$|w_k|^2 = |\tilde{w}_k|^2 \left| \delta^2 \frac{D'(\lambda_{k-p})}{N(\lambda_{k-p})} \right|, \quad \text{for } p+1 \leq k \leq p+q$$

in (13) the following sensitivity measure is obtained which considers the weighted sensitivities of all poles and zeros with respect to all coefficients of the state matrix A :

$$\begin{aligned} m_A^{\Phi\Psi} &= \left\{ \sum_{k=1}^p |\tilde{w}_k|^2 \left(\sum_{i=1}^n \sum_{j=1}^n \left| \frac{\partial \mu_k}{\partial a_{ij}} \right|^2 \right)^{1/2} \right. \\ &\quad \left. + \sum_{k=p+1}^{p+q} |\tilde{w}_k|^2 \left(\sum_{i=1}^n \sum_{j=1}^n \left| \frac{\partial \lambda_{k-p}}{\partial a_{ij}} \right|^2 \right)^{1/2} \right\}^2 \\ &\leq \text{tr}(K^\Phi) \text{tr}(W^\Psi). \end{aligned}$$

By the use of the decompositions $F(z) = F_p(z)/D(z)$, $G(z) = G_p(z)/D(z)$ with the numerator polynomial vectors $F_p(z) \in \mathbb{C}^{n \times 1}$ and $G_p(z) \in \mathbb{C}^{1 \times n}$, the corresponding Gramian matrices can be evaluated according to

$$\begin{aligned} K^\Phi &= \sum_{k=1}^p F_p(\mu_k) F_p(\mu_k)^h \left| \frac{\tilde{w}_k^2}{N'(\mu_k) D(\mu_k)} \right| \\ &\quad + \sum_{k=p+1}^{p+q} F_p(\lambda_{k-p}) F_p(\lambda_{k-p})^h \left| \frac{\tilde{w}_k^2}{N(\lambda_{k-p}) D'(\lambda_{k-p})} \right| \\ W^\Psi &= \sum_{k=1}^p G_p(\mu_k)^h G_p(\mu_k) \left| \frac{\tilde{w}_k^2}{N'(\mu_k) D(\mu_k)} \right| \\ &\quad + \sum_{k=p+1}^{p+q} G_p(\lambda_{k-p})^h G_p(\lambda_{k-p}) \left| \frac{\tilde{w}_k^2}{N(\lambda_{k-p}) D'(\lambda_{k-p})} \right|. \end{aligned}$$

IV. RELATION BETWEEN SENSITIVITY AND NOISE

It is known for a long time that relations exist between sensitivity and noise measures, e.g., [24], [25]. A generalization is possible by the use of the previously given results. In the case of state-space representations with internal signals which are scaled under a coloured input process, there is a linear dependence between the weighted sensitivity measures of Definition 2 and the weighted noise measure of (8).

Theorem 2:

A scaled state-space representation according to (1) satisfies $\sigma_{x_i}^2 = 1$, ($i=1, \dots, n$) under a stationary input process with the spectral density $|\Phi(z)|^2$. Then, the relations

$$m_A^{\Phi\Psi} \leq n g^\Psi, \quad (16)$$

$$m_b^\Psi = g^\Psi \quad (17)$$

$$m_c^\Phi = n, \quad (18)$$

$$m^{\Phi\Psi} \leq n + (n+1) g^\Psi \quad (19)$$

hold where the equality signs in (16) and (19) are valid if and only if

$$\rho^2 |\Phi(z)|^2 F(z)^h F(z) = |\Psi(z)|^2 G(z) G(z)^h,$$

$$z \in \Gamma, \quad \rho \in \mathbb{R} \setminus \{0\}.$$

Proof: The proof of this theorem is a simple consequence of Theorem 1 if the relations $\sigma_{x_i}^2 = k_{ii}^\Phi = 1 \rightarrow \text{tr}(K^\Phi) = n$ and $g^\Psi = \text{tr}(W^\Psi)$ are taken into account. \square

This new relation will lead to the determination of the whole class of optimal state-space representations which are distinguished by simultaneous minimal weighted sensitivity and noise.

V. THE CLASS OF OPTIMAL STATE-SPACE REPRESENTATIONS

In order to derive the necessary and sufficient conditions for the optimality of a state-space representation the following two inequalities are useful. Although their proofs have already been given in [11], [26] these relations will be presented again because of their importance in this connection.

Theorem 3:

$K^\Phi \in \mathbb{C}^{n \times n}$, $W^\Psi \in \mathbb{C}^{n \times n}$ are two hermitian, positive semi-definite matrices. Then the relations

$$\text{tr}(K^\Phi) \text{tr}(W^\Psi) \geq \left(\sum_{i=1}^n \mu_i \right)^2 \tag{20}$$

$$\text{tr}(K^\Phi) + \text{tr}(W^\Psi) \geq 2 \sum_{i=1}^n \mu_i, \quad \mu_i \geq 0 \tag{21}$$

hold where $\{\mu_i^2\}$ are the eigenvalues of the matrix $K^\Phi W^\Psi$. The equality sign in (20) is valid iff $W^\Psi = \rho^2 K^\Phi$, $\rho \in \mathbb{R} \setminus \{0\}$ and in (21) iff $W^\Psi = K^\Phi$.

By the use of (20) it is as a generalization of the results given in [11], [12] now possible to define the class of state-space realizations which are scaled under an arbitrary stationary input process and which are characterized by the minimum possible weighted noise power gain g^Ψ .

Theorem 4:

A scaled state-space description according to (1) satisfies

$$\sum_{i=1}^n \sigma_{x_i}^2 = n \tag{22}$$

under a stationary input process with the spectral density $|\Phi(z)|^2$. Then the measure g^Ψ takes its absolute minimum value

$$g^\Psi = \frac{1}{n} \left(\sum_{i=1}^n \mu_i \right)^2$$

iff

$$W^\Psi = \left(\frac{1}{n} \sum_{i=1}^n \mu_i \right)^2 K^\Phi.$$

Proof:

The assumption of the above theorem can be formulated as follows:

$$\sum_{i=1}^n \sigma_{x_i}^2 = \sum_{i=1}^n k_{ii}^\Phi = \text{tr}(K^\Phi) = n.$$

Using this constraint in (20), it can be concluded that

$$g^\Psi = \text{tr}(W^\Psi) \geq \frac{1}{n} \left(\sum_{i=1}^n \mu_i \right)^2.$$

As according to Theorem 3 the equality sign is valid iff $W^\Psi = \rho^2 K^\Phi$ and, therefore, $g^\Psi = \text{tr}(W^\Psi) = n\rho^2$, one ob-

tains the condition

$$\rho^2 = 1/n^2 \left(\sum_{i=1}^n \mu_i \right)^2. \quad \square$$

The l_2 -dynamic range constraint of (22) can be satisfied, e.g., by the n equations $\sigma_{x_i}^2 = 1$, ($i = 1, \dots, n$). Theorem 4 is more general than the corresponding statement of Mullis and Roberts [8], as the input signal is not restricted to have a white noise characteristic and the noise measure contains an arbitrary weighting function.

Now the class of sensitivity optimal state-space realizations will be defined. It will be proved that the general measures of Definition 2 take their absolute minimum values if the considered representation is in this class. Contrary to former results, e.g., [11], [12], [25], not only an upper bound of sensitivities but their exact values are concerned in the following.

Theorem 5:

A real state-space representation according to (1) and a weighting function $|\Phi(z)|^2$ are given such that $r(K^\Phi) = r(W^\Phi) = n$. Then

$$m_A^\Phi \geq \left(\sum_{i=1}^n \mu_i \right)^2 \tag{23}$$

$$m_b^\Phi + m_c^\Phi \geq 2 \sum_{i=1}^n \mu_i \tag{24}$$

$$m^\Phi \geq \left(\sum_{i=1}^n \mu_i \right)^2 + 2 \sum_{i=1}^n \mu_i \tag{25}$$

where the equality sign in (23) is valid iff $W^\Phi = \rho^2 K^\Phi$, $\rho \in \mathbb{R} \setminus \{0\}$ and those in (24,25) are valid iff $W^\Phi = K^\Phi$.

The proof of this theorem is given in the Appendix.

By the use of Theorem 5 it is possible to find a direct relation (which therefore dispenses with bounds) between round-off noise and sensitivity in the case of an optimal realization with proportional generalized Gramian matrices W^Φ and K^Φ . Moreover, with $|\Phi(z)|^2 = |\Psi(z)|^2 = 1$ it is shown that the roundoff noise optimal realizations according to [8], [10] which satisfy $W = \rho^2 K$ possess the absolute minimum value of the integral sensitivity measures of Definition 2. It is proved that the balanced realization with $K = W = \text{diag}(\mu_i)$ has the minimum sensitivity and thus is recommended as a starting realization if numerical algorithms are applied to rational functions.

VI. TRANSFORMATION TO OPTIMAL REPRESENTATIONS

As it has been shown in Theorems 4 and 5 that realizations with proportional Gramian matrices satisfy some optimality conditions it is desirable to define the class by the use of the continuous equivalence transformation given in (3). The following definition turns out to be the central point of the optimization procedure:

Definition 3:

A given state-space representation is Φ, Ψ -balanced iff the corresponding Gramian matrices according to Defini-

tion 1 satisfy

$$K^\Phi = W^\Psi = \text{diag}(\mu_i), \quad \mu_i > 0, \quad i=1, \dots, n.$$

This definition generalizes the well-known balanced property which plays an essential role in the results given by Moore [17] and which has been used implicitly by Mullis and Roberts [8]. Definition 1 leads to the conclusion that the partial transfer functions $F_i(z)$ (and $G_i(z)$) of a Φ, Ψ -balanced realization according to Definition 3 are orthogonal with respect to the weighting functions $\Phi(z)$ (and $\Psi(z)$). Therefore, the transfer function $H(z)$ is decomposed into linear combinations of these orthogonal functions in the case of a Φ, Ψ -balanced state-space description.

Now, the whole class of optimal state-space representations in terms of a set of transformation matrices will be determined. A numerically well conditioned algorithm to determine a balanced realization has been given by Laub [27]. The matrix $T_D = LUM^{-1/2}$ transforms according to (3) so that $K_D^\Phi = W_D^\Psi = M$ with the Cholesky factorization $K^\Phi = LL^h$ and the symmetric eigenvalue problem $U^h(L^hW^\Psi L)U = M^2$, $M = \text{diag}(\mu_i)$. Now, the matrix

$$T_0 = |\rho|^{1/2} T_D Q \in \mathbb{C}^{n \times n}, \quad QQ^h = I, \quad \rho \in \mathbb{R} \setminus \{0\}$$

transforms to all possible realizations with $W_0^\Psi = \rho^2 K_0^\Phi$. This property can be proved as follows:

$$\begin{aligned} W_0^\Psi = \rho^2 K_0^\Phi &\leftrightarrow |\rho| Q^h W_D^\Psi Q = \rho^2 |\rho|^{-1} Q^{-1} K_D^\Phi (Q^{-1})^h \\ &\leftrightarrow (QQ^h) M (QQ^h) = M \leftrightarrow QQ^h = (QQ^h)^{-1} \\ &\leftrightarrow QQ^h = I. \end{aligned}$$

The third equivalence can be shown by the use of Lemma 1 which is given in the Appendix.

It has been proved that all elements of the complete class of optimal representations with $K^\Phi = \rho^2 W^\Psi$ are connected by the orthogonal group of transformation matrices with $T=Q$. This degree of freedom can be used in order to optimize additional properties of a state-space representation. In the case of Gaussian input process with the spectral density function $|\Phi(z)|^2$, the scaled realization has equal state variances and thus equal probabilities of the occurrence of overflow at any state. In this case, the relations $\sigma_{x_i}^2 = k_{ii}^\Phi = \kappa^2$, ($i=1, \dots, n$) hold, where the real number κ controls the probability of overflow of the internal signals. The noniterative algorithm of Hwang [10] can be applied in order to construct the desired orthogonal matrix Q which consists of $(n-1)$ -coordinate rotations and which distributes the weighted state vector norm equally among the states.

Now, the methods presented will be compared with previous results on state-space realizations [8]–[12] in order to show the application of the proposed theory for filter design. The matrix K^Φ can be used to scale any starting realization with respect to an arbitrary input spectrum. To this end, a diagonal and orthogonal transformation matrix has proven to be useful as both procedures do not disturb the stability under nonlinear operations. Therefore, especially minimum noise realizations [8], [10], [11] or normal

realizations [9] can be optimally l_2 -scaled such that considerably better roundoff noise properties are obtained in comparison to the usual scaling procedure [21]. The state-space realization according to Theorem 4 has the minimum noise power gain whereby the difference to [8], [10] depends on its transfer function and on the input spectrum which has been considered here [21]. The proposed sensitivity measures of Definition 2 take into account all coefficients of the state-space representation. As an extension of [11], [12], [25] arbitrary weighting functions make possible a matching of the expected deviations of the transfer function to the given specifications. Contrary to previously known results [11], [12], [21], [25] the sensitivity itself can be minimized instead of its upper bound.

VII. CONCLUDING REMARKS

It has been shown in the two Theorems 4 and 5 that all optimal realizations satisfy $K^\Phi \propto W^\Psi$. These representations are connected by an orthogonal transformation and a transformation with a scalar. One special element of this class of realizations is the Φ, Ψ -balanced representation of Definition 3. It is pointed out in this section that numerous important digital filter realizations are elements of this class and therefore share the condition $K^\Phi \propto W^\Psi$.

1) A transfer function $H(z)$ with $\|H(z)\|_\infty < 1$ can be embedded in a paraunitary transfer function matrix. It has been shown by several authors that a corresponding state-space realization can be chosen to be orthogonal, which yields $K=W=I$. Therefore, this special class of state-space wave digital filters [28] or state-space orthogonal filters belongs to the class of optimal representations if it is considered as a class of multivariable systems. Moreover, if these systems are regarded as scalar ones, neglecting one input/output pair, it can be shown that they are optimal according to the measures of (13)–(15). In this case, the frequency points have to be located at the zeroes of the attenuation with $|H(z_k)|=1$, which yields $K^\Phi \propto W^\Phi$.

2) It is known that the scaled, round-off noise optimal realizations given by Hwang [10] and Mullis and Roberts [8] are characterized by the condition $K = \rho^2 W$. Therefore, these realizations are optimal in the sense of the Theorems 4 and 5 with $|\Phi(z)|^2 = |\Psi(z)|^2 = 1$.

3) It can be shown that a parallel realization of normal second and first order blocks is identical to a pole balanced representation which satisfies $K^\Phi = W^\Phi = \text{diag}(\mu_i)$ with the frequency points $z_k = \lambda_k$.

4) The most general class of optimal state-space representations uses general weighting functions $|\Phi(z)|^2$ and $|\Psi(z)|^2$. It has been shown by the use of examples [21], [26] that the corresponding state-space realizations have better numerical properties than those designed by other methods, if the designed system is considered to be embedded in an actual environment.

APPENDIX

Now, the proof of Theorem 5 will be given. Let us first proof (24). As according to Theorem 1 $m_b^\Phi = \text{tr}(W^\Phi)$ and $m_c^\Phi = \text{tr}(K^\Phi)$ one obtains $m_b^\Phi + m_c^\Phi = \text{tr}(K^\Phi) + \text{tr}(W^\Phi)$

which yields under consideration of Theorem 3 the desired result.

Lemma 1:

A nonsingular Hermitian matrix $P \in \mathbb{C}^{n \times n}$ and a positive definite Hermitian matrix $K \in \mathbb{C}^{n \times n}$ are given. Then

$$K = PKP \leftrightarrow PK = KP, \quad P = P^{-1}.$$

Proof:

$$PK = KP, \quad P = P^{-1} \rightarrow K = PKP$$

$$K = PKP \rightarrow K = U\Lambda_p U^h K U\Lambda_p U^h$$

$$\rightarrow (U^h K U) = \Lambda_p (U^h K U) \Lambda_p$$

$$\rightarrow \Lambda_p = \text{diag}(\pm 1) \quad \text{as } U^h K U > 0$$

$$\text{causes } (U^h K U)_{ii} > 0$$

$$\rightarrow P = P^{-1}.$$

There we use the eigenvalue decomposition $P = U\Lambda_p U^h$ with $UU^h = I$ and $\Lambda_p = \text{diag}(\lambda_{p_i})$. \square

Lemma 2:

A real state-space representation and a weighting function $|\Phi(z)|^2$ are given so that $r(K^\Phi) = r(W^\Phi) = n$. Then

$$W^\Phi = \rho^2 K^\Phi \rightarrow G(z)G(z)^h = \rho^2 F(z)^h F(z), \\ z \in \Gamma, \quad \rho \in \mathbb{R} \setminus \{0\}.$$

Proof:

The proof uses a statement that has been given in [3]. In the case of a real state-space representation there exists a unique symmetric nonsingular matrix $P \in \mathbb{R}^{n \times n}$ with

$$\begin{bmatrix} P & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} A & b \\ c & d \end{bmatrix} = \begin{bmatrix} A & b \\ c & d \end{bmatrix}' \begin{bmatrix} P & 0 \\ 0 & 1 \end{bmatrix}$$

and

$$\begin{bmatrix} P & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \Theta(z) & F(z) \\ G(z) & H(z) \end{bmatrix} = \begin{bmatrix} \Theta(z) & F(z) \\ G(z) & H(z) \end{bmatrix}' \begin{bmatrix} P & 0 \\ 0 & 1 \end{bmatrix}.$$

It can easily be derived from (6), (7) that the Gramian matrices satisfy the matrix equations

$$AK^\Phi A' + Q_K = K^\Phi, \quad A'W^\Phi A + Q_W = W^\Phi$$

where $PQ_K P = Q_W$. If general nonrational weighting functions are applied the same relations hold [20]. Now, the following proof can be given:

$$AK^\Phi A' + Q_K = K^\Phi, \quad W^\Phi = A'W^\Phi A + Q_W$$

$$\rightarrow PAK^\Phi A'P + PQ_K P = PK^\Phi P, \quad W^\Phi = A'W^\Phi A + Q_W$$

$$\rightarrow A'(PK^\Phi P)A + Q_W = PK^\Phi P, \quad W^\Phi = A'W^\Phi A + Q_W$$

$$\rightarrow PK^\Phi P = \rho^2 K^\Phi \quad \text{because of the uniqueness}$$

$$\text{of (6), (7) and } W^\Phi = \rho^2 K^\Phi$$

$$\rightarrow P = \rho^2 P^{-1} \quad \text{because of Lemma 1}$$

$$\rightarrow G(z)G(z)^h = F(z)' P P F(z)^* \quad \text{as } G(z) = F(z)' P$$

$$\rightarrow G(z)G(z)^h = \rho^2 F(z)^h F(z) \quad \text{as } G(z)G(z)^h \in \mathbb{R}.$$

Now, by the use of these lemmas, the proof of Theorem 5 can be formulated. With $\tilde{W}^\Phi = \tilde{T}' W^\Phi \tilde{T}$, $\tilde{K}^\Phi = \tilde{T}^{-1} K^\Phi (\tilde{T}^{-1})'$ and $\tilde{R} = \tilde{T} \tilde{T}'$ it will be shown, that a realization with $\tilde{W}^\Phi = \rho^2 \tilde{K}^\Phi$ satisfies

$$m_A^{\Phi\Phi}(\tilde{R}) = \left\{ \frac{1}{2\pi j} \oint_{\Gamma} (F(z)^h \tilde{R}^{-1} F(z) G(z) \tilde{R} G(z)^h)^{1/2} \cdot |\Phi(z)|^2 \frac{dz}{z} \right\}^2 \\ = \left(\sum_{i=1}^n \mu_i \right)^2$$

and that a realization whose functional $(m_A^{\Phi\Phi}(\tilde{R}))^{1/2}$ is minimal satisfies $\tilde{W}^\Phi = \rho^2 \tilde{K}^\Phi$, $\rho \in \mathbb{R} \setminus \{0\}$.

The first part can be proved as follows:

$$\tilde{W}^\Phi = \rho^2 \tilde{K}^\Phi \rightarrow m_A^{\Phi\Phi}(\tilde{R}) = \text{tr}(\tilde{K}^\Phi) \text{tr}(\tilde{W}^\Phi)$$

because of (9) and Lemma 2

$$\rightarrow m_A^{\Phi\Phi}(\tilde{R}) = \left(\sum_{i=1}^n \mu_i \right)^2 \quad \text{because of Theorem 3.}$$

In order to show the second part, the Jacobian matrix of the functional $(m_A^{\Phi\Phi}(\tilde{R}))^{1/2}$ is given:

$$D(\tilde{R}) = \frac{1}{4\pi j} \oint_{\Gamma} \left\{ k(\tilde{R}) \tilde{R}^{-1} F(z) F(z)^h \tilde{R}^{-1} - k(\tilde{R})^{-1} G(z)^h G(z) \right\} |\Phi(z)|^2 \frac{dz}{z}$$

with

$$k(\tilde{R}) = \left(\frac{G(z) \tilde{R} G(z)^h}{F(z)^h \tilde{R}^{-1} F(z)} \right)^{1/2}.$$

Now it remains to show that $D(\tilde{R}) = 0 \rightarrow \tilde{W}^\Phi = \rho^2 \tilde{K}^\Phi$, $\rho \in \mathbb{R} \setminus \{0\}$ because $D(\tilde{R}) = 0$ leads to the global minimum of the functional (as

$$m_A^{\Phi\Phi}(\tilde{R}) = \left(\sum_{i=1}^n \mu_i \right)^2$$

and as the functional satisfies $m_A^{\Phi\Phi}(R) > 0$, $R > 0$).

Let us choose $W^\Phi = \rho^2 K^\Phi$ and, therefore, $D(I) = 0$ as

$$D(I) = \frac{1}{4\pi j} \oint_{\Gamma} \left\{ k(I) F(z) F(z)^h - k(I)^{-1} G(z)^h G(z) \right\} \cdot |\Phi(z)|^2 \frac{dz}{z} \\ = \frac{1}{4\pi j} \left[\oint_{\Gamma} \rho F(z) F(z)^h |\Phi(z)|^2 \frac{dz}{z} - \oint_{\Gamma} \frac{1}{\rho} G(z)^h G(z) |\Phi(z)|^2 \frac{dz}{z} \right] = 0.$$

\square There we use the fact that $k(I) = \rho$ because of Lemma 2.

Now, the proof can be formulated as follows:

$$\begin{aligned}
 D(\tilde{R}) = 0 &\rightarrow \tilde{R}D(\tilde{R}) - D(I) = 0 \quad \text{as } D(I) = 0 \\
 &\rightarrow \frac{1}{4\pi j} \oint_{\Gamma} \left\{ F(z)F(z)^h (k(\tilde{R})\tilde{R}^{-1} - \rho I) \right. \\
 &\quad \left. + \left(\frac{1}{\rho} I - k(\tilde{R})^{-1}\tilde{R} \right) G(z)^h G(z) \right\} |\Phi(z)|^2 \frac{dz}{z} = 0 \\
 &\rightarrow (k(\tilde{R})\tilde{R}^{-1} - \rho I) \left(\frac{1}{\rho} I - k(\tilde{R})^{-1}\tilde{R} \right) \leq 0 \\
 &\quad \text{as } r(K^\Phi) = r(W^\Phi) = n \\
 &\rightarrow (k(\tilde{R})\tilde{R}^{-1} - \rho I) \left(\frac{1}{\rho} I - k(\tilde{R})^{-1}\tilde{R} \right) = 0 \\
 &\quad \text{as } (k(\tilde{R})\tilde{R}^{-1} - \rho I) \left(\frac{1}{\rho} I - k(\tilde{R})^{-1}\tilde{R} \right) \geq 0 \\
 &\rightarrow \frac{1}{\rho} k(\tilde{R})\tilde{R}^{-1} = I \\
 &\rightarrow \tilde{T}' = \frac{1}{\rho} k(\tilde{R})\tilde{T}^{-1} \\
 &\rightarrow \tilde{W}^\Phi = \tilde{T}'W^\Phi\tilde{T} = \frac{1}{\rho} k(\tilde{R})\tilde{T}^{-1}\rho^2 K^\Phi (T^{-1})' \frac{1}{\rho} k(\tilde{R}) \\
 &\rightarrow \tilde{W}^\Phi = k(\tilde{R})^2 \tilde{K}^\Phi. \quad \square
 \end{aligned}$$

ACKNOWLEDGMENT

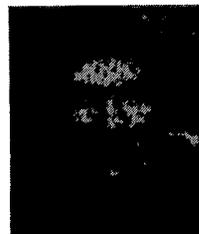
The author is grateful to Professor Saal for his continued encouragement and support on the subject of this paper. He also thanks the reviewers for many valuable comments and suggestions.

REFERENCES

- [1] P. M. v. Dooren, "The generalized eigenstructure problem in linear system theory," *IEEE Automat. Contr.* vol. AC-26, pp. 111-129, Aug. 1981.
- [2] S. Barnett and C. Storey, *Matrix Methods in Stability Theory*. London, UK: Thomas Nelson, 1970.
- [3] B. D. O. Anderson and S. Vongpanitlerd, *Network Analysis and Synthesis*. Englewood Cliffs, NJ: Prentice-Hall, 1973.
- [4] P. M. v. Dooren and P. Dewilde, "Minimal cascade factorization of real and complex rational transfer matrices," *IEEE Trans. Circuits Syst.*, vol. CAS-28, pp. 390-400, May 1981.
- [5] B. Garbow, *Matrix Eigensystem Routines—Eispack Guide Extension*. New York: Springer Verlag, 1977.
- [6] P. M. v. Dooren, "Numerical linear algebra: An increasing interest in linear system theory," in *Proc. European Conf. Circuit Theory Design*, The Hague, Netherlands, pp. 243-251, 1981.
- [7] A. S. Willsky, "Relationships between digital signal processing and control and estimation theory," *Proc. IEEE*, vol. 66, pp. 996-1017, Sept. 1978.
- [8] C. T. Mullis and R. A. Roberts, "Synthesis of minimum round-off noise fixed point digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-23, pp. 551-561, Sept. 1976.
- [9] C. W. Barnes, "Round-off noise and overflow in normal digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-26, pp. 154-159, Mar. 1979.
- [10] S. Y. Hwang, "Minimum uncorrelated unit noise in state-space digital filtering," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-25, pp. 273-281, Aug. 1977.

- [11] L. Thiele, "Design of sensitivity and round-off noise optimal state-space discrete systems," *Int. J. Circuit Theory Appl.* vol. 12, pp. 39-46, Jan. 1984.
- [12] V. Tavsanoğlu and L. Thiele, "Optimal design of state-space digital filters by simultaneous minimization of sensitivity and round-off noise," *IEEE Trans. Circuits Syst.*, vol. CAS-31, pp. 884-888, Oct. 1984.
- [13] S. Kung, "A new identification and model reduction algorithm via singular value decompositions," *Proc. 12th Annual Asilomar Conf. on Circuits Systems*, pp. 705-714, 1978.
- [14] W. Mills and C. T. Mullis, "Digital filter realizations without overflow oscillations," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-26, pp. 334-338, Aug. 1978.
- [15] J. B. Knowles and E. M. Olcayto, "Coefficient accuracy and digital filter response," *IRE Trans. Circuit Theory*, vol. CT-15, pp. 31-41, Mar. 1968.
- [16] E. Avenhaus, "On the design of digital filters with coefficients of limited word length," *IEEE Trans. Audio Electroacoust.*, vol. AU-20, pp. 206-212, Aug. 1972.
- [17] B. C. Moore, "Principal component analysis in linear systems: Controllability, observability and model reduction," *IEEE Trans. Automat. Contr.*, vol. AC-26, pp. 17-32, Feb. 1981.
- [18] R. A. Smith, "Matrix equation $AX + XB = C$," *SIAM J. Appl. Math.*, vol. 16, pp. 198-201, 1968.
- [19] R. H. Bartels and G. W. Stewart, "Solution to the matrix equation $AX + XB = C$," *Comm. ACM*, vol. 15, pp. 820-826, Sept. 1972.
- [20] L. Thiele, "Balanced model reduction in time and frequency domain," in *Proc. Int. Symp. Circuits and Systems 1985*, Kyoto, Japan, pp. 345-348.
- [21] L. Thiele, "Generalized Gramian matrices and their applications in digital filters," in *Digital Signal Processing-84*, V. Cappellini, A. G. Constantinides (Eds.), North-Holland: Elsevier, pp. 13-17, 1984.
- [22] A. Fettweis, "A general theorem for signal-flow networks with applications," *Arch. Elektron. Uebertragung*, vol. 25, pp. 557-561, Dec. 1971.
- [23] H. Ur, "Root locus properties and sensitivity relations in control systems," *IEEE Trans. Automat. Contr.*, vol. AC-5, pp. 57-65, Jan. 1960.
- [24] A. Fettweis, "On the connection between multiplier word length limitations and roundoff noise for digital filters," *IEEE Trans. Circuit Theory*, vol. CT-19, pp. 486-491, Sept. 1972.
- [25] L. B. Jackson, "Round-off noise bounds derived from coefficient sensitivities for digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-28, pp. 481-485, Aug. 1976.
- [26] L. Thiele, "On the approximation and realization of linear systems by the use of methods of linear algebra," (in German) Ph.D. dissertation, Tech. Univ. Munich, Munich, Germany, 1985.
- [27] A. J. Laub, "Computation of 'balancing' transformations," in *Proc. 1980 JACC*, San Francisco, CA, p. FA8-E, 1980.
- [28] A. Fettweis, "Digital circuits and systems," *IEEE Trans. Circuits Syst.*, vol. CAS-31, pp. 31-48, Jan. 1984.

✱



Lothar Thiele (S'83-M'86) was born in Aachen, Germany on April 7, 1957. He received the Diplom-Ingenieur and Dr.-Ing. degrees in electrical engineering from Technical University of Munich, Munich, Germany, in 1981 and 1985, respectively.

Since 1981, he has been with the Institute on Circuits and Systems of the Technical University of Munich. His research and teaching interests include methods of linear algebra in linear system theory and theoretical aspects of VLSI.

Dr. Thiele is a member of NTG (Nachrichtentechnische Gesellschaft, Germany).