

# NEAR CAPACITY TOMLINSON-HARASHIMA PRECODING WITH UNIFORM STREAM PROPERTIES

Raphael Hunger, Maximilian Riemensberger, and Wolfgang Utschick

Associate Institute for Signal Processing  
Munich University of Technology, 80290 Munich, Germany  
hunger@tum.de

## ABSTRACT

We address the nonlinear transceiver design in a point-to-point MIMO system with *Tomlinson-Harashima precoding* (THP). By jointly optimizing both transmitter and receiver, capacity can be achieved up to the shaping loss and each stream can be decoded separately. In contrast to linear filtering, THP allows for uniform stream properties rendering bit-loading unnecessary and allows to span an arbitrarily high number of streams for the sake of a reduced cardinality modulation alphabet. Existing work studied either the decision feedback equalizer version with the nonlinearity located at the receiver or the perfect dirty paper precoding where the *geometric mean decomposition* (GMD) can be applied. We explicitly take the modulo operator into account leading to the fact that the *generalized triangular decomposition* has to be applied instead of the GMD.

## 1. INTRODUCTION

Recently, decision feedback aided nonlinear filtering was shown to achieve capacity in point-to-point systems with the nice characteristic of uniform stream properties and minimum sum mean square error [1, 2]. Thus, undesirable bit-loading and different coding on each stream, which is necessary for linear filtering to achieve capacity, become obsolete. The nonlinear structure can also be shifted to the transmitter: in [2], the authors also consider perfect *dirty paper* (DP) precoding for which stream-coding must be done jointly leading to a highly complex precoder. We deploy *Tomlinson Harashima Precoding* (THP) as a low-complexity practical implementation of DP. For the THP system, we show that capacity is achieved except for the shaping loss [3]. Moreover, a uniform stream property is achieved in the relevant SNR region and an arbitrarily high number of streams can be transmitted for the sake of a reduced cardinality modulation alphabet. However, the main difference to the work [1, 2] is that for THP, the so called *generalized triangular decomposition* (GTD) [4] has to be applied instead of the *geometric mean decomposition* (GMD) [5] which in contrast to the GMD not necessarily exists. The GTD must be used due to the fact that the modulo

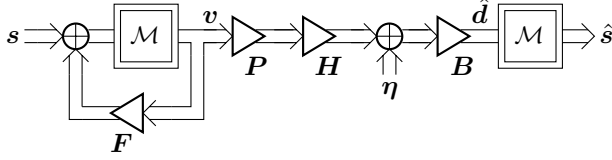
operator does not affect the last stream which therefore has a different variance. Both GTD and GMD are matrix factorizations where a matrix is decomposed into the product of a unitary matrix, an upper/lower triangular matrix with prescribed main diagonal, and another unitary matrix.

In Section 2, we discuss the system model and briefly review the *linear* capacity achieving filtering in Section 3. Afterwards, we show in Section 4 that the capacity achieving precoder in THP has enough degrees of freedom left to let all individual streams have uniform properties in the relevant SNR region. To this end, the GTD is utilized, the properties and construction of which are explained in Section 5. Finally, simulation results in Section 6 show that THP can clearly outperform the decision feedback aided system despite the power loss and the modulo loss since THP does not suffer from error propagation.

*Notation:* Matrices and vectors are denoted by upper and lower case bold italic letters, respectively. The operators  $\mathbb{E}[\cdot]$ ,  $(\cdot)^H$ ,  $(\cdot)^T$ ,  $\text{tr}(\cdot)$ ,  $\det(\cdot)$ , and  $[\cdot]_{i,j}$  denote expectation with respect to the noise  $\boldsymbol{\eta}$  and the signal  $\boldsymbol{v}$ , conjugate transposition, transposition, trace, determinant, and the matrix element in the  $i$ th row and  $j$ th column, respectively. The set  $\mathbb{T}^{K \times K} \subset \mathbb{C}^{K \times K}$  contains all  $K \times K$  strictly upper triangular matrices,  $\mathbb{A}$  represents the discrete or continuous modulation alphabet,  $\mathbb{S}_+^{K \times K} \subset \mathbb{C}^{K \times K}$  denotes the cone of the  $K \times K$  positive semidefinite matrices,  $\mathbb{D}^{K \times r} \subset \mathbb{R}_{+,0}^{K \times r}$  is the set of all (not necessarily square) diagonal matrices, and the set of all  $K \times K$  unitary matrices is denoted by  $\mathbb{U}^{K \times K} \subset \mathbb{C}^{K \times K}$ .  $\Re\{\cdot\}$ ,  $\Im\{\cdot\}$ , and  $j$  denote real and imaginary part of the complex argument, and the imaginary unit  $\sqrt{-1}$ , respectively.  $\succcurlyeq$  is a partial order on the proper cone of positive semidefinite matrices.

## 2. SYSTEM MODEL

We consider a point-to-point MIMO communication link as depicted in Figure 1 where the zero-mean data vector  $\boldsymbol{s} \in \mathbb{A}^K$  with identity covariance matrix contains the modulated symbols of the  $K$  streams. The feedback filter  $\boldsymbol{F} \in \mathbb{T}^{K \times K}$  vanishes in case of linear precoding whereas it is restricted to



**Fig. 1.** Nonlinear point-to-point system model.

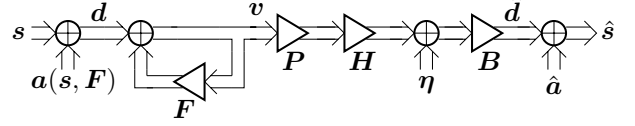
be *strictly* upper<sup>1</sup> triangular (with zeros on its main diagonal) when THP is applied. The modulo operator  $\mathcal{M}(\cdot) : \mathbb{C}^K \rightarrow \mathbb{V}^K$  with  $\mathbb{V} = \{z \in \mathbb{C} \mid -\tau \leq 2\Re\{z\} < \tau, -\tau \leq 2\Im\{z\} < \tau\}$  is an element-wise many-to-one mapping from the entire  $K$ -dimensional complex hyper-plane into the half-open complex hyper-cube by an element-wise addition of integer multiples of its argument such that its image lies in  $\mathbb{V}^K$ . We make the common statistical assumptions on the output  $\mathbf{v} \in \mathbb{V}^K$  of the modulo operator [6], i.e., we assume that the covariance matrix  $\mathbf{C}_v = \mathbb{E}[\mathbf{v}\mathbf{v}^H] \in \mathbb{S}_+^{K \times K}$  is diagonal which means that the individual entries are uncorrelated. Furthermore, the last stream is not affected by the modulo operator when the modulo constant  $\tau$  is chosen sufficiently large, hence its variance remains one. In case of an  $M$ -ary QAM modulation with  $M = 4^n, n \in \mathbb{N}$ , the first  $K - 1$  diagonal entries of  $\mathbf{C}_v$  are assumed to have variance  $\sigma_v^2 = \frac{M}{M-1}$ , see [6], if the respective feedback is *active*. This follows from the assumption that  $[v]_{1,1}, \dots, [v]_{K-1,1}$  are uniformly distributed over  $\mathbb{V}$ . The signal  $\mathbf{v}$  is linearly precoded by  $\mathbf{P} \in \mathbb{C}^{N_{\text{Tx}} \times K}$  and propagates over the frequency flat channel  $\mathbf{H}^{N_{\text{Rx}} \times N_{\text{Tx}}}$ . At the receiver side, zero-mean additive white Gaussian noise  $\boldsymbol{\eta} \in \mathbb{C}^{N_{\text{Rx}}}$  is superimposed and the receive filter  $\mathbf{B}^{K \times N_{\text{Rx}}}$  generates  $\hat{\mathbf{d}} \in \mathbb{C}^K$  representing an estimate for the virtual signal  $\mathbf{d} \in \mathbb{A}^K + \tau\mathbb{Z}^K + j\tau\mathbb{Z}^K \subset \mathbb{C}^K$ . In the end, the modulo operator generates the estimated symbol vector  $\hat{\mathbf{s}}$  by performing the remapping onto  $\mathbb{V}^K$ .

Fig. 2 shows a modulo operator free representation of the transmitter with the additive signal  $\mathbf{a}(\mathbf{s}, \mathbf{F})$  from a  $K$ -dimensional infinite lattice  $\tau\mathbb{Z}^K + j\tau\mathbb{Z}^K$  such that  $\mathbf{v} \in \mathbb{V}^K$  remains the same. As the modulo operators avoid simple relations between  $\mathbf{s}$  and  $\hat{\mathbf{s}}$ , all our optimizations are based on the virtual signals  $\mathbf{d}$  and  $\hat{\mathbf{d}}$ .

### 3. CAPACITY ACHIEVING LINEAR FILTERING REVISITED

If the precoder is restricted to act linearly, the feedback filter  $\mathbf{F} = \mathbf{0}$  is inactive and all modulo operators are removed. In the following, we shortly review the capacity achieving linear filtering from an information theoretic point of view and from

<sup>1</sup>In general, any symmetrically permuted strictly upper triangular matrix may be chosen since causality and realizability are ensured by means of the resulting matrix structure. We choose the strictly upper triangularity due to the fact that the general triangular decomposition, which we will make use of later, was introduced for upper triangular matrices first.



**Fig. 2.** Modulo-operator free precoder and receiver representation.

a signal processing point of view and highlight key properties of the resulting transmission chain.

#### 3.1. Information Theoretic Point of View

Capacity is defined as the maximum mutual information between the transmit signal  $\mathbf{x} = \mathbf{P}\mathbf{s} \in \mathbb{C}^{N_{\text{Tx}}}$  and the receive signal  $\mathbf{y} = \mathbf{H}\mathbf{P}\mathbf{s} + \boldsymbol{\eta} \in \mathbb{C}^{N_{\text{Rx}}}$ , and maximization is done with respect to the input distribution of  $\mathbf{x}$  and its covariance  $\mathbf{C}_x = \mathbb{E}[\mathbf{x}\mathbf{x}^H] = \mathbf{P}\mathbf{P}^H \in \mathbb{S}_+^{N_{\text{Tx}} \times N_{\text{Tx}}}$  subject to an average power constraint  $\text{tr}(\mathbf{C}_x) \leq P_{\text{Tx}}$ :

$$\begin{aligned} \max_{\mathbf{C}_x} \log_2 \det \left( \mathbf{I}_{N_{\text{Tx}}} + \mathbf{C}_x \mathbf{H}^H \mathbf{C}_\eta^{-1} \mathbf{H} \right) \\ \text{s.t.: } \text{tr}(\mathbf{C}_x) \leq P_{\text{Tx}} \text{ and } \mathbf{C}_x \succcurlyeq \mathbf{0} \end{aligned} \quad (1)$$

First of all, Gaussian signaling has to be chosen such that the utility in (1) represents not only an upper bound, but also the truly achievable sum-rate. Telatar [7] came up with the well known result that the input covariance  $\mathbf{C}_x$  needs to have the same eigenspace as  $\mathbf{H}^H \mathbf{C}_\eta^{-1} \mathbf{H}$ . Moreover, the water-filling policy guarantees optimum power allocation. The optimum covariance  $\check{\mathbf{C}}_x$  directly follows from the KKT conditions assigned to (1) and can compactly be expressed as

$$\check{\mathbf{C}}_x = [\mu^{-1} \mathbf{I}_{N_{\text{Tx}}} - (\mathbf{H}^H \mathbf{C}_\eta^{-1} \mathbf{H})^{-1}]_{\perp}, \quad (2)$$

where  $\mu \ln 2$  is the Lagrangian factor which is chosen such that the constraint  $\text{tr}(\check{\mathbf{C}}_x) \leq P_{\text{Tx}}$  is fulfilled with equality. The operator  $[\cdot]_{\perp}$  performs the orthogonal projection of its Hermitian argument onto the proper cone of positive semidefinite matrices by setting all negative eigenvalues to zero [8]. It can be thought of as the multi-dimensional extension of the  $\max(0, \cdot)$  operator known from scalar water-pouring. However, this compact notation already tells us which eigenspace needs to be chosen and the optimum power allocation follows as well, cf. [9].

#### 3.2. Signal Processing Point of View

Solution (2) is important from an information theoretic point of view. However, from a signal processing perspective, it is more relevant to know how to realize the precoder  $\mathbf{P}$  instead of its covariance  $\mathbf{C}_x = \mathbf{P}\mathbf{P}^H$  and how to design the receiver. In conjunction with an MMSE-type *receiver*  $\mathbf{B}$  recovering  $\mathbf{s}$ , the utility in (1) can also be expressed as  $-\log_2 \det \mathbf{C}_e$  (see [10]) with

$$\mathbf{C}_e = (\mathbf{I}_K + \mathbf{P}^H \mathbf{H}^H \mathbf{C}_\eta^{-1} \mathbf{H} \mathbf{P})^{-1}$$

denoting the covariance matrix of the error vector  $e = s - \hat{s}$ .

In order to find the optimum transmit filter, we introduce the *sorted* reduced eigenvalue decomposition  $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^H$  of  $\mathbf{H}^H\mathbf{C}_\eta^{-1}\mathbf{H}$  with  $\mathbf{U} \in \mathbb{C}^{N_{\text{Tx}} \times r}$  containing the first  $r$  unit-norm eigenvectors and  $\mathbf{\Lambda} \in \mathbb{D}^{r \times r}$  containing the  $r$  non-zero eigenvalues  $\lambda_1, \dots, \lambda_r$  in non-increasing order. Setting  $\check{\mathbf{P}} = \mathbf{U}\check{\Phi}\mathbf{S}$ , where  $\check{\Phi} \in \mathbb{D}^{r \times K}$  is diagonal but not necessarily square with entries  $\check{\Phi}_1, \dots, \check{\Phi}_{\min(r,K)} \in \mathbb{R}$ , the water-filling policy leads to a power allocation according to  $\check{\Phi}_k^2 = \max(0, \frac{1}{\mu} - \frac{1}{\lambda_k})$ ,  $k \in \{1, \dots, \min(K, r)\}$ . In order to let the covariance matrix  $\mathbf{C}_x = \mathbf{P}\mathbf{P}^H$  achieve the rank of the optimum covariance matrix  $\check{\mathbf{C}}_x$  from (2), the number of streams  $K$  has to be chosen at least as large as the rank of the optimum covariance matrix  $d = \text{rank}(\check{\mathbf{C}}_x) \leq N_{\text{Tx}}$ . Otherwise, capacity cannot be achieved. Finally, a unitary matrix  $\mathbf{S} \in \mathbb{U}^{K \times K}$  can be chosen as a degree of freedom as it does not change the obtained sum-rate.

### 3.3. Key Properties of the Transmission Chain

Choosing  $\mathbf{U}$  and  $\check{\Phi}$  as in the previous section, capacity is achieved. The resulting impacts on the transmission chain are now summarized:

1.) For  $\mathbf{S} = \mathbf{I}_K$ , the number of streams  $K$  has to be chosen equal to the dimension  $d$  of the subspace where transmission takes place. If  $K < d$ , capacity cannot be achieved, if  $K > d$ ,  $K - d$  streams get zero power allocated. So for  $\mathbf{S} = \mathbf{I}_K$  and  $K = d$ , the error covariance matrix  $\mathbf{C}_e = (\mathbf{I}_K + \check{\Phi}^T \mathbf{\Lambda} \check{\Phi})^{-1}$  gets diagonalized and all  $K$  streams can be decoded *separately*, which has an enormous practical relevancy and drastically reduces complexity. However, the  $K$  streams have different individual MSEs and hence different signal-to-noise ratios in general. The MSE of stream  $k$  reads as  $\varepsilon_k = \frac{\mu}{\lambda_k}$  and obviously depends on the associated eigenvalue  $\lambda_k$ . As a consequence, both transmitter and receiver must be capable of handling different coding schemes and different coding rates in order to achieve capacity. Furthermore, different modulation schemes must be available if the theoretical limit which is based on Gaussian signaling shall be approached by a practical QAM modulation scheme. Whereas one stream may require QPSK, another stream might demand for 64QAM in order to come close to the Gaussian limit.

2.) Identical stream MSEs and SINRs, i.e., equal diagonal elements of the error covariance matrix  $\mathbf{C}_e$ , can be achieved by choosing  $\mathbf{S}$  as a DFT or Hadamard matrix [11]. For this case,  $K > d$  would also make sense, since power for the additional  $K - d$  streams is allocated by the unitary matrix  $\mathbf{S}$ . Despite the fact that capacity is still obtained, the price one has to pay is that the error covariance matrix is no longer diagonal. Since all streams are coupled then, they have to be decoded jointly, leading to an intractable complexity.

3.) The number of transmitted streams  $K$  is therefore limited by the dimension  $d = \text{rank}(\check{\mathbf{C}}_x)$  and reaches the rank of the channel  $r = \text{rank}(\mathbf{H}) \leq \min(N_{\text{Tx}}, N_{\text{Rx}})$  in the high

SNR region. No more streams can be spanned for the sake of a reduced cardinality of the modulation alphabet if separate stream decoding is preferred.

## 4. OPTIMUM THP TRANSCEIVER DESIGN

### 4.1. Precoder Structure for Near-Capacity Transmission

For THP, the modulo operator limits the real and imaginary part of each entry in  $\mathbf{v} \in \mathbb{V}^K$  and hence prevents Gaussian signaling. The induced *shaping loss* [3] leads to the fact that the channel capacity can only be obtained up to this 1.53 dB shaping loss [3]. Nonetheless, we focus on the maximization of the upper bound. Inserting the MMSE receiver, depending on  $\mathbf{F}$  and  $\mathbf{P}$ , into the covariance matrix  $\mathbf{C}_e$  of the error signal  $e = \mathbf{d} - \hat{\mathbf{d}}$  yields

$$\mathbf{C}_e = (\mathbf{I} - \mathbf{F})(\mathbf{P}^H \mathbf{H}^H \mathbf{C}_\eta^{-1} \mathbf{H} \mathbf{P} + \mathbf{C}_v^{-1})^{-1} (\mathbf{I} - \mathbf{F})^H. \quad (3)$$

Again, the obtained rate can be expressed as a function of the determinant of  $\mathbf{C}_e$ , see [10]. However, as mentioned in Section 2, the covariance matrix  $\mathbf{C}_v$  of  $\mathbf{v}$  is no longer an identity matrix and the obtained rate, neglecting the shaping loss, now reads as  $\log_2 \det(\mathbf{C}_v) - \log_2 \det(\mathbf{C}_e)$ . Due to the fact that  $\mathbf{F}$  is strictly upper triangular, the determinant of  $\mathbf{C}_e$  does not depend on  $\mathbf{F}$ , i.e.,  $\det(\mathbf{I}_K - \mathbf{F}) = 1$ . This means that, no matter how  $\mathbf{F}$  is chosen, capacity up to the shaping loss can always be reached by properly choosing  $\mathbf{P}$ . The optimum precoder  $\check{\mathbf{P}}$  follows from a determinant maximization problem [12] with covariance matrix  $\mathbf{C}_v$  and reads as

$$\check{\mathbf{P}} = \mathbf{U}\check{\Phi}\mathbf{S}\mathbf{C}_v^{-\frac{1}{2}}, \quad (4)$$

where  $\mathbf{U}$  and  $\check{\Phi}$  are defined as in Section 3.2 and  $\mathbf{S} \in \mathbb{U}^{K \times K}$  is again an arbitrary unitary matrix. Summing up, capacity-achieving precoding, except from the non-Gaussianity, does not uniquely determine the feedforward filter  $\mathbf{P}$ , a unitary matrix  $\mathbf{S}$  remains as degree of freedom. Moreover, the feedback filter  $\mathbf{F}$  does not have any influence on the throughput.

### 4.2. MSE Minimization and Uniform Stream Properties

Defining the *arithmetic* MSE  $\varepsilon_A$  as the sum MSE via

$$\varepsilon_A = \mathbb{E}[\|\mathbf{d} - \hat{\mathbf{d}}\|_2^2] = \text{tr}(\mathbb{E}[\mathbf{e}\mathbf{e}^H]) = \text{tr}(\mathbf{C}_e),$$

it becomes evident that  $\varepsilon_A$  corresponds to the trace of the error covariance matrix  $\mathbf{C}_e$ . A lower bound on the arithmetic MSE can be derived from the trace-determinant inequality

$$\varepsilon_A = \text{tr}(\mathbf{C}_e) \geq K \sqrt[K]{\det(\mathbf{C}_e)} = K \sqrt[K]{\det(\varepsilon_G)} \quad (5)$$

relating the *arithmetic mean* to the *geometric mean*. Here,  $\varepsilon_G = \det(\mathbf{C}_e)$  denotes the *geometric* MSE. The lower bound is obtained if, and only if,  $\mathbf{C}_e$  is a scaled identity matrix. The

authors in [1,2] were the first to introduce this framework. Instead of directly minimizing the sum-MSE in a decision feedback system at the receiver, they minimized the lower bound, i.e., the geometric MSE, and showed that the minimum lower bound can also be achieved by means of the *geometric mean decomposition* (GMD) [5]. The GMD is a matrix factorization where an arbitrary matrix is decomposed into the product of a unitary matrix, an upper triangular matrix with *identical* main diagonal entries, and a second unitary matrix. But since we employ THP at the transmitter and not decision feedback at the receiver, the signal covariance matrix  $C_v$  is not a scaled identity matrix. This follows from the fact that the power loss does not affect that last stream which is not influenced by the modulo operator. As we will see later, we have to decompose a matrix in a similar way except that the main diagonal entries of the upper triangular matrix have (different) prescribed values now. As a consequence, the GMD cannot be applied and the more powerful *generalized triangular decomposition* (GTD) [4] must be utilized. Its existence and properties as well as a very fast implementation are discussed in Section 5.

As  $\hat{P}$  from (4) achieves capacity up to the shaping loss, it maximizes  $-\log_2 \det C_e$  and thus minimizes the geometric mean  $\varepsilon_G = \det(C_e)$ . In order to let the arithmetic mean  $\varepsilon_A = \text{tr}(C_e)$  merge with this minimum lower bound, the error covariance matrix  $C_e$  has to boil down to a scaled identity matrix. Plugging the optimum precoder  $\hat{P}$  into the covariance expression (3), we get

$$C_e = (\mathbf{I}_K - \mathbf{F})C_v^{\frac{1}{2}}\mathbf{S}^H\mathbf{D}^{-2}\mathbf{S}C_v^{\frac{1}{2}}(\mathbf{I}_K - \mathbf{F})^H, \quad (6)$$

where  $\mathbf{D} = (\mathbf{I}_K + \Phi^T \mathbf{A} \Phi)^{\frac{1}{2}}$  is diagonal and positive definite. If a unitary  $\mathbf{S}$  and a strictly upper triangular  $\mathbf{F} \in \mathbb{T}^{K \times K}$  exist such that  $C_e = \sigma_\varepsilon^2 \mathbf{I}_K = \sigma_\varepsilon^2 \mathbf{Q}^H \mathbf{Q}$  with unitary  $\mathbf{Q} \in \mathbb{U}^{K \times K}$ , i.e.,  $C_e$  becomes a scaled identity, the arithmetic mean achieves its global minimum. To this end, a unitary decomposition<sup>2</sup>

$$\mathbf{D} = \mathbf{Q}(\mathbf{I}_K - \mathbf{F})\sigma_\varepsilon^{-1}C_v^{\frac{1}{2}}\mathbf{S}^H = \mathbf{Q}\mathbf{R}\mathbf{S}^H \quad (7)$$

should exist with unitary  $\mathbf{S}$  and  $\mathbf{Q}$  and an upper triangular matrix

$$\mathbf{R} = \sigma_\varepsilon^{-1}(\mathbf{I}_K - \mathbf{F})C_v^{\frac{1}{2}}, \quad (8)$$

the diagonal of which corresponds to the diagonal of  $C_v^{\frac{1}{2}}$  divided by  $\sigma_\varepsilon$  and therefore has different entries on its main diagonal. That's why the GMD cannot be used and the GTD has to be employed. In Section 5, exact conditions for its existence are presented. From simulations, we find that these requirements are only violated for very low SNR values below -10 dB. More precisely, the common MSE level  $\sigma_\varepsilon^2$  cannot grow above one, see Section 5. As  $\mathbf{R}$  follows from two unitary transformations applied to  $\mathbf{D}$  (cf. Eq. 7), and since both  $\mathbf{R}$  and  $\mathbf{D}$  have positive real-valued main diagonal entries,  $\mathbf{Q}$

<sup>2</sup>Remember that  $\mathbf{D} = (\mathbf{I}_K + \Phi^T \mathbf{A} \Phi)^{\frac{1}{2}}$  is diagonal.

and  $\mathbf{S}$  can be chosen real-valued as well. Additionally,  $\mathbf{R}$  and  $\mathbf{D}$  have the same determinant from which we can compute the common MSE level

$$\sigma_\varepsilon^2 = \sqrt[2K]{\prod_{k=1}^K [C_v]_{k,k} \prod_{k=1}^K [D]_{k,k}^{-2}} = \sqrt[2K]{\sigma_v^{2(K-1)} \prod_{k=1}^K [D]_{k,k}^{-2}}, \quad (9)$$

where we exploited the fact that the first  $K - 1$  diagonal entries of  $C_v$  are  $\sigma_v^2$  and the last entry equals one. Given the unfavorable conditions at very low SNR that the GTD does not exist, the minimum MSE cannot be achieved in combination with a balancing of all individual MSEs. As soon as  $\sigma_\varepsilon^2$  is one, i.e., the sum-MSE reaches  $K$ , a further reduction of the transmit power prevents the GTD to exist, and consequently, all streams cannot be balanced any more. Interestingly, it is possible to switch off the last stream such that its MSE is one, and balance the remaining streams  $1, \dots, K - 1$  at an MSE level below one. Doing so, the  $K - 1$ st stream is precoded linearly, and therefore, the  $K - 1$ st diagonal entry in  $C_v$  reduces from  $\sigma_v^2$  to one, as the feedback for the linearly precoded stream  $K - 1$  is disabled. This change in the variance leads to the fact that the MSE level of the remaining streams changes, cf. (9). If this level has dropped below  $[D]_{K,K}^{-2}$ , the remaining  $K - 1$  streams can be balanced at this level, otherwise, the  $K - 1$ st stream will also be switched off, changing the variance of the  $K - 2$ nd stream to one, and so on. If only a single stream is active in the end, the covariance matrix  $C_v$  has been transformed to the identity matrix, and the modulo operators become obsolete. From this procedure it becomes obvious that the sum-MSE can never grow above the number of streams  $K$  no matter how small the transmit power  $P_{Tx}$  gets.

### 4.3. Key Properties of the Nonlinear Transmission

The following statements hold under the assumption that the GTD exists.

1.) Except for the shaping loss due to the non-Gaussianity of  $v$ , the precoder achieves capacity. At high SNRs, the modulo operator at the receiver side does not significantly change the Gaussianity of the receive noise. These propositions remain valid even if the GTD does not exist.

2.) The error covariance matrix  $C_e$  reduces to a scaled identity matrix, meaning that all streams have the same properties, especially the same MSEs and the same rates. Hence, no bit-loading is necessary which is the major advantage compared to capacity achieving linear filtering. Moreover, each stream can be decoded separately which has a drastically lower complexity than a joint decoding of all streams simultaneously.

3.) The number  $K$  of active streams must be larger than or equal to the dimension  $d = \text{rank}(\check{C}_x)$  of the optimum covariance matrix from (2). Capacity is achieved for an arbitrary high number of streams. Larger stream numbers facili-

tate the use of small cardinality modulation alphabets making this approach very important with respect to the practical implementation. Clearly, the MSE per stream is increased when increasing the number of streams.

4) In contrast to decision-feedback equalization (DFE) based nonlinear filtering [1,2], THP does not suffer from error propagation. Thus, considerable gains can be achieved in the mid-SNR region despite the modulo loss implicating the generation of new nearest neighbors and despite the power loss ( $\mathbf{C}_v \neq \mathbf{I}_K$ ) since the crucial assumption of perfect decision in DFE is not fulfilled in this SNR region, see the simulation results in Section 6.

## 5. THE GENERALIZED TRIANGULAR DECOMPOSITION

### 5.1. Existence and Uniqueness

Taking a closer look at (7), we find from  $\mathbf{R} = \mathbf{Q}^H \mathbf{D} \mathbf{S}$  that the diagonal matrix  $\mathbf{D}$  comprises the singular values of the matrix  $\mathbf{R}$ , which itself has the eigenvalues on its diagonal, since it is upper triangular. Weyl [13] proved that the singular values multiplicatively majorize<sup>3</sup> the eigenvalues of a matrix, and Horn [14] stated that an upper triangular matrix  $\mathbf{R}$  with prescribed diagonal and a specific singular value set exists, if the diagonal is majorized by the singular values, see [5]. In our context, this means that

$$\begin{aligned} \prod_{k=1}^{\ell} [\mathbf{R}]_{k,k}^2 &\leq \prod_{k=1}^{\ell} [\mathbf{D}]_{k,k}^2, \forall \ell : 1 \leq \ell < K \\ \prod_{k=1}^K [\mathbf{R}]_{k,k}^2 &= \prod_{k=1}^K [\mathbf{D}]_{k,k}^2, \end{aligned} \quad (10)$$

because the diagonal entries of both  $\mathbf{R}$  and  $\mathbf{D}$  are already sorted in an non-increasing fashion. For  $\mathbf{D}$ , this follows from the water-pouring policy and the fact that we made use of the *sorted* eigenvalue decomposition. The second row in (10) is always fulfilled due to the choice of  $\sigma_\varepsilon^2$  in (9). Since the entries in  $\mathbf{D}$  are non-increasing and the first  $K-1$  entries in  $\mathbf{C}_v$  are identical, we may replace (10) in conjunction with (8) by the necessary and sufficient condition

$$\frac{1}{\sigma_\varepsilon^2} \geq [\mathbf{D}]_{K,K}^2 \quad (11)$$

for the existence of the GTD. Note that  $[\mathbf{D}]_{K,K}^{-2} \leq 1$  and hence  $\sigma_\varepsilon^2 \leq 1$  must hold. In the following, we show that the transmit power  $P_{\text{Tx}}$  can always be chosen so small that the GTD does not exist. To this end, we assume that the transmit power is already so small that  $d = 1$  holds, i.e.,  $\Phi_1^2 = P_{\text{Tx}}$  and  $\Phi_k = 0 \forall k \neq 1$ . Moreover,  $K = 2$  streams shall be transmitted. Then  $[\mathbf{D}]_{1,1}^2 = 1 + \lambda_1 P_{\text{Tx}}$ ,  $[\mathbf{D}]_{2,2}^2 = 1$ , and  $\sigma_\varepsilon^2 = \sqrt{\frac{\sigma_v^2}{1 + \lambda_1 P_{\text{Tx}}}}$  hold. Eq. 10 now requires that  $\frac{\sigma_v^2}{\sigma_\varepsilon^2} \leq [\mathbf{D}]_{1,1}^2$ , which is violated for  $P_{\text{Tx}} < \lambda_1^{-1}(\sigma_v^2 - 1)$ .

<sup>3</sup>In this context, majorization is related to an non-increasing ordering.

Similar to the GMD [5], the GTD is not unique if it exists. In particular, infinitely many decompositions exist but the Frobenius norm of all those matrices  $\mathbf{R}$  is identical to  $\|\mathbf{D}\|_F$ . One type of invariance are unitary diagonal matrices: Let  $\mathbf{D} = \mathbf{Q} \mathbf{R} \mathbf{S}^H$  be a GTD decomposition of the diagonal matrix  $\mathbf{D}$ , and let  $\mathbf{V}$  and  $\mathbf{W}$  be unitary diagonal matrices. With  $\tilde{\mathbf{R}} = \mathbf{V} \mathbf{R} \mathbf{V}^H$ ,  $\tilde{\mathbf{Q}} = \mathbf{W} \mathbf{Q} \mathbf{V}^H$ , and  $\tilde{\mathbf{S}} = \mathbf{W} \mathbf{S} \mathbf{V}^H$ , we find that  $\mathbf{D} = \mathbf{W} \mathbf{D} \mathbf{W}^H = \tilde{\mathbf{W}} \mathbf{Q} \tilde{\mathbf{R}} \tilde{\mathbf{S}}^H \mathbf{W}^H = \tilde{\mathbf{Q}} \tilde{\mathbf{R}} \tilde{\mathbf{S}}^H$  is also a valid GTD decomposition of  $\mathbf{D}$ . Furthermore, a completely different off-diagonal structure in  $\mathbf{R}$  can be obtained for  $K > 2$  by changing the order selection of the Givens rotations, see the following section.

### 5.2. Efficient and Stable Computation of the GTD

Independently from but similar to [2,4], we found a low complexity and extremely stable way to compute the GTD when trying to recude the complexity of the GMD from quartic order in [15] to quadratic order. This goal was achieved by the use of Givens rotations as in [2]. We interpret the unitary matrices  $\mathbf{S}$  and  $\mathbf{Q}$  as the product of  $K-1$  real-valued Givens rotations  $\mathbf{Q}_i$  and  $\mathbf{S}_i$ , respectively:

$$\mathbf{S} = \prod_{i=1}^{K-1} \mathbf{S}_i, \quad \text{and} \quad \mathbf{Q} = \prod_{i=1}^{K-1} \mathbf{Q}_i.$$

A (real-valued) Givens rotation  $\mathbf{G}_{j,k}$  is a rank-two perturbation of the  $K$ -dimensional identity matrix with  $[\mathbf{G}_{j,k}]_{j,j} = [\mathbf{G}_{j,k}]_{k,k} = a$ ,  $[\mathbf{G}_{j,k}]_{j,k} = b$ , and  $[\mathbf{G}_{j,k}]_{k,j} = -b$  holds. Imposing the constraints  $a, b \in \mathbb{R}$ ,  $-1 \leq a \leq 1$ ,  $-1 \leq b \leq 1$ , and  $a^2 + b^2 = 1$  leads to the fact that  $\mathbf{G}_{j,k}$  is orthogonal. We assume that  $j < k$ . Note that  $\mathbf{G}_{j,k}$  performs a counterclockwise rotation in the  $j, k$  plane with  $a$  and  $b$  corresponding to the cosine and the sine of the rotation angle [16].

Multiplying  $\mathbf{G}_{j,k}$  from the right-hand side (RHS) onto a matrix has influence only on the columns  $j$  and  $k$  whereas a left-hand side (LHS) multiplication has impact on rows  $j$  and  $k$  only. Hence, we can focus on rows  $j$  and  $k$  and columns  $j$  and  $k$  separately. Extracting the four entries of the LHS and RHS multiplication of two different Givens rotations applied to the matrix  $\mathbf{D}$  yields

$$\begin{aligned} \begin{bmatrix} a & b \\ -b & a \end{bmatrix} \begin{bmatrix} z_j & 0 \\ 0 & z_k \end{bmatrix} \begin{bmatrix} c & d \\ -d & c \end{bmatrix} &= \\ \begin{bmatrix} acz_j - bdz_k & adz_j + bcz_k \\ -(bcz_j + adz_k) & -bdz_j + acz_k \end{bmatrix}. \end{aligned} \quad (12)$$

$z_j$  and  $z_k$  are the  $j$ th and the  $k$ th diagonal entries of the matrix  $\mathbf{Z}$  which is initialized with the matrix  $\mathbf{D}$  and to which the Givens rotations are applied such that after  $N-1$  rotations,  $\mathbf{Z} = \mathbf{R}$  holds. Since the upper triangular structure of  $\mathbf{Z}$  must be preserved, the lower left entry in (12) has to be zero. W.l.o.g., we set the upper left<sup>4</sup> element to the desired

<sup>4</sup>It is also possible to set the lower right element to the desired value. However, the order selection changes then.

value  $\beta_j = \sigma_\varepsilon^{-1} [C_\sigma^{\frac{1}{2}}]_{j,j}$ , see (8). The set of equations

$$\begin{aligned} acz_j - bdz_k &= \beta_j, & bcz_j + adz_k &= 0 \\ a^2 + b^2 &= 1, & c^2 + d^2 &= 1 \end{aligned}$$

leads to the solutions

$$c = \sqrt{\frac{\beta_j^2 - z_k^2}{z_j^2 - z_k^2}}, \quad d = -\sqrt{\frac{z_j^2 - \beta_j^2}{z_j^2 - z_k^2}}, \quad a = \frac{z_j c}{\beta_j}, \quad b = -\frac{z_k d}{\beta_j}.$$

The verification of the constraints  $0 \leq 1 - b^2 \leq 1$  and  $0 \leq d^2 \leq 1$  tells us how to choose the indices  $j$  and  $k$ :

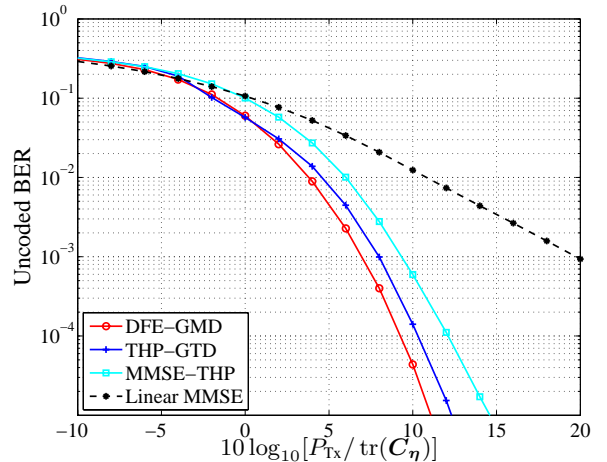
$$z_j \leq \beta_j \leq z_k \quad \text{or} \quad z_j \geq \beta_j \geq z_k. \quad (13)$$

Apparently, we need to find  $k$  such that the desired value  $\beta_j$  is enclosed by the two diagonal elements  $z_k$  and  $z_j$ . Starting with  $j = 1$ , we search for an index  $k > j$  for which (13) is fulfilled. If the majorization criterion (10) is fulfilled, such a  $k$  always exists. After the application of the Givens rotations to  $\mathbf{Z}$ ,  $j$  is increased by one and a  $k > j$  satisfying (13) has to be found again. This procedure terminates after  $K - 1$  steps and during the last step,  $j = K - 1$  and  $k = K$  hold since the determinants of  $\mathbf{D}$  and  $\mathbf{R}$  are the same. The overall complexity to compute both  $\mathbf{Q}$ ,  $\mathbf{S}$ , and  $\mathbf{R}$  is only  $\mathcal{O}(K^2)$ .

## 6. SIMULATION RESULTS

In our simulation results, we compare the THP GTD-based joint transmitter and receiver design (cross marker) presented in this paper, with the joint transmitter and receiver GMD-based decision feedback equalization (circle marker) from [1, 2]. Both versions guarantee uniform stream properties (the GTD requires the common MSE level to be below one for stream balancing). For purposes of comparison, we also simulated the nonlinear MMSE THP (square marker) from [17] and the linear MMSE filter (dashed curve) from [18] where the receive filters are forced to be scaled identity matrices. Clearly, these filters only exist for  $K = N_{\text{Rx}}$  but they do not balance the individual streams' MSEs. Nonetheless, we plot stream-averaged uncoded bit-error ratios (BERs).

Fig. 3 shows a QPSK transmission of  $K = 4$  streams in a system with  $N_{\text{Tx}} = 4$  transmit and  $N_{\text{Rx}} = 4$  receive antennas averaged over uncorrelated channels. Linear MMSE precoding (dashed curve) [18] with a non-cooperative receiver is clearly inferior to nonlinear THP-based precoding or nonlinear DFE-based equalization. Both the GMD-based and the GTD-based schemes clearly outperform the MMSE THP filter [17] with a non-cooperative receiver being a scaled identity matrix. There are two reasons why the DFE-based system performs better than the GTD-based one. First, the generation of next neighbors due to the modulo operator is quite crucial for QPSK, since every symbol now has four nearest neighbors instead of two. Second, the variance  $\sigma_v^2$  of the first three

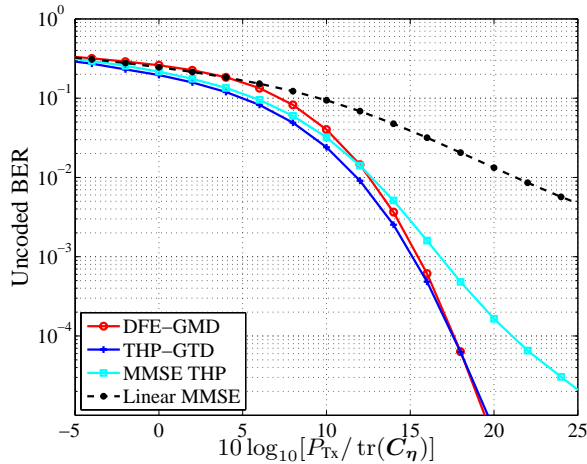


**Fig. 3.** QPSK transmission of  $K = 4$  streams over a system with  $N_{\text{Tx}} = 4$  transmit and  $N_{\text{Rx}} = 4$  receive antennas (uncorrelated channel).

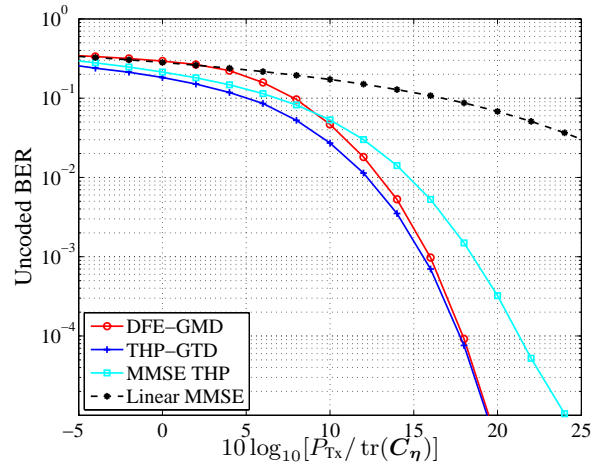
streams raises from one to  $\sigma_v^2 = \frac{M}{M-1} = \frac{4}{3}$  for QPSK bringing an increased MSE level  $\sigma_\varepsilon^2$  with it compared to the case when  $\sigma_v^2 = 1$  would hold, see the definition of  $\sigma_\varepsilon^2$  in Section 4.2.

This changes when we switch the modulation alphabet to 16-QAM, see Fig. 4. All other parameters are left unchanged. First, we notice that all filters perform slightly worse. In the previous figure, the transmit signal-to-noise ratio  $P_{\text{Tx}} / \text{tr}(C_\eta)$  was too small to see that the MMSE THP with non-cooperative receiver (square marker) from [17] flattens due to the fact that it does not achieve the full diversity order. In Fig. 4, this can already be anticipated, and in fact, the square marker curve achieves the diversity order  $1 + \min(N_{\text{Tx}}, N_{\text{Rx}}) - K = 1$ , since one stream is precoded linearly. By contrast, both the GMD (circle marker) and therefore also the GTD based (cross marker) version have the full diversity order, see [2]. Another interesting observation is the intersection between the DFE-GMD curve (circle marker) and the MMSE THP curve (square marker). It results from the error propagation of the DFE in the small and moderate SNR range. The assumption of perfect interference subtraction is clearly violated in this region. Furthermore, the generation of next neighbors due to the modulo operator is not so severe for the THP-GTD (cross marker) in case of 16-QAM, since only the outer symbol get new nearest neighbors. In addition, the power loss induced by  $\sigma_v^2 > 1$  is very small for 16-QAM, since  $\sigma_v^2 = \frac{M}{M-1} = \frac{16}{15}$ . Both effects are not as detrimental as the error propagation in case of decision feedback equalization. For an uncoded BER of  $10^{-1}$ , a gain of 2 dB can be achieved if THP is applied instead of the DFE.

Even larger gains can be obtained when the channel becomes correlated as in Fig. 5. Here,  $K = 8$  streams are transmitted from a  $N_{\text{Tx}} = 8$  antenna sender to a  $N_{\text{Rx}} = 8$  antenna receiver. Due to the correlated scenario, error propa-



**Fig. 4.** 16-QAM transmission of  $K = 4$  streams over a system with  $N_{Tx} = 4$  transmit and  $N_{Rx} = 4$  receive antennas (uncorrelated channel).



**Fig. 5.** 16-QAM transmission of  $K = 8$  streams over a system with  $N_{Tx} = 8$  transmit and  $N_{Rx} = 8$  receive antennas (correlated channel).

gation turns out to be even more severe. For an uncoded BER of  $10^{-1}$ , 3 dB can be gained by the THP-GTD compared to the DFE-GMD, and if  $10^{-2}$  is the desired operation point, we can still save 1 dB.

## 7. REFERENCES

- [1] F. Xu, T. N. Davidson, J.-K. Zhang, and K. M. Wong, "Design of Block Transceivers with Decision Feedback Detection," *IEEE Transactions on Signal Processing*, vol. 54, no. 3, pp. 965–978, March 2006.
- [2] Y. Jiang, J. Li, and W. W. Hager, "Uniform Channel Decomposition for MIMO Communications," *IEEE Transactions on Signal Processing*, vol. 53, no. 11, November 2005.
- [3] Wei Yu, David P. Varodayan, and John M. Cioffi, "Trellis and convolutional precoding for transmitter-based interference pre-subtraction," *IEEE Transactions on Communications*, vol. 53, no. 7, pp. 1220–1230, 2005.
- [4] Y. Jiang, W. W. Hager, and J. Li, "The Generalized Triangular Decomposition," Accepted for publication in *Mathematics of Computation*, 2005.
- [5] Y. Jiang, W. W. Hager, and J. Li, "The Geometric Mean Decomposition," *Linear Algebra and its Application*, vol. 396, pp. 373–384, 2005.
- [6] Robert F. Fischer, *Precoding and Signal Shaping for Digital Transmission*, John Wiley & Sons, Inc., New York, NY, USA, 2002.
- [7] Emre Telatar, "Capacity of multi-antenna gaussian channels," *European Transactions on Telecommunications*, vol. 10, no. 6, pp. 585–596, November/December 1999.
- [8] N. J. Higham, "Matrix Nearness Problems and Applications," in *Applications of Matrix Theory*, M. J. C. Gover and S. Barnett, Eds. 1989, pp. 1–27, Oxford University Press.
- [9] R. Hunger, D. A. Schmidt, and W. Utschick, "Sum Capacity and MMSE for the MIMO Broadcast Channel without Eigenvalue Decompositions," *Submitted to ISIT 2007*.
- [10] P. Stoica, Y. Jiang, and J. Li, "On MIMO Channel Capacity: An Intuitive Discussion," *IEEE Signal Processing Mag.*, pp. 83–84, May 2005.
- [11] D. Palomar, J. Cioffi, and M. Lagunas, "Joint Tx-Rx beamforming Design for Multicarrier MIMO Channels: A Unified Framework for Convex Optimization," *IEEE Transactions on Signal Processing*, vol. 51, no. 9, pp. 2381–2401, 2003.
- [12] H. S. Witsenhausen, "A Determinant Maximization Problem Occurring in The Theory of Data Communication," *Siam Journal on Applied Mathematics*, vol. 29, no. 3, pp. 515–522, November 1975.
- [13] H. Weyl, "Inequalities between Two Kinds of Eigenvalues of a Linear Transformation," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 35, pp. 408–411, 1949.
- [14] A. Horn, "On the Eigenvalues of a Matrix with Prescribed Singular Values," in *Proceedings of the American Mathematical Society*, February 1954, vol. 5, pp. 4–7.
- [15] Jian Kang Zhang, Aleksandar Kavčić, and Kon Max Wong, "Equal-Diagonal QR Decomposition and its Application to Precoder Design for Successive-Cancellation Detection," January 2005.
- [16] Gene H. Golub and Charles F. Van Loan, *Matrix Computations*, The Johns Hopkins University Press, 1991.
- [17] M. Joham, J. Brehmer, and W. Utschick, "MMSE approaches to multiuser spatio-temporal Tomlinson-Harashima precoding," in *Proc. ITG SCC'04*, January 2004, pp. 387–394.
- [18] M. Joham, W. Utschick, and J. A. Nossek, "Linear Transmit Processing in MIMO Communications Systems," *IEEE Transactions on Signal Processing*, vol. 53, no. 8, pp. 2700–2712, August 2005.