

Low-Complexity Near-Maximum-Likelihood Detection and Precoding for MIMO Systems using Lattice Reduction

Christoph Windpassinger Robert F. H. Fischer

Lehrstuhl für Informationsübertragung, Universität Erlangen–Nürnberg,

Cauerstraße 7/NT, 91058 Erlangen, Germany, <http://www.LNT.de>,

email: {windpass,fischer}@LNT.de

Abstract — We consider the lattice-reduction-aided detection scheme for 2×2 channels recently proposed by Yao and Wornell [11]. Using an equivalent real-valued substitute MIMO channel model their lattice reduction algorithm can be replaced by the well-known LLL algorithm, which enables the application to MIMO systems with arbitrary numbers of dimensions. We show how lattice reduction can also be favourably applied in systems that use precoding and give simulation results that underline the usefulness of this approach.

I. INTRODUCTION

In a recent publication by Yao and Wornell [11] a novel scheme for improved detection of signals transmitted over multiple-input/multiple-output (MIMO) systems was presented. The astonishing property of this scheme is that it results in error rate curves that parallel those for maximum-likelihood (ML) detection (with some penalty in power efficiency), at only a fraction of the complexity.

In the present work we show how their approach fits in the general (maximum-likelihood) lattice decoding framework of [1] and extend the work of [11], which presented an optimum algorithm for 2×2 complex MIMO systems based on Gaussian reduction [3], to higher-dimensional settings. The key point is the application of the (sub-optimum) basis reduction algorithm by A. K. Lenstra, H. W. Lenstra and L. Lovász (“LLL algorithm”, [9]). Note that this algorithm has also been used in connection with efficient near-ML decoding of differential space-time codes in [2]. Furthermore, we will show how this approach can be applied in precoding/preequalization schemes.

II. TRANSMISSION MODEL AND GENERAL FRAMEWORK

We consider the typical flat-fading MIMO transmission model

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{n}, \quad (1)$$

where usually the vectors \mathbf{x} , \mathbf{n} , \mathbf{y} and the matrix \mathbf{H} are given in the equivalent baseband, and hence are complex-valued. \mathbf{H} is the channel matrix, \mathbf{x} the vector of transmitted symbols, each chosen from some (finite) set \mathcal{A} and \mathbf{n} the additive white Gaussian noise vector. For simplicity we assume all vectors to be column vectors of dimension K , and \mathbf{H} a $K \times K$ matrix of complex transfer coefficients between transmit and receive antennas/subchannels.

We can equivalently write (1) as

$$\begin{bmatrix} \Re \mathbf{y} \\ \Im \mathbf{y} \end{bmatrix} = \begin{bmatrix} \Re \mathbf{H} & -\Im \mathbf{H} \\ \Im \mathbf{H} & \Re \mathbf{H} \end{bmatrix} \begin{bmatrix} \Re \mathbf{x} \\ \Im \mathbf{x} \end{bmatrix} + \begin{bmatrix} \Re \mathbf{n} \\ \Im \mathbf{n} \end{bmatrix}, \quad (2)$$

(where the \Re and \Im prefix denote the real and imaginary parts), which gives an equivalent $2K$ -dimensional real model of the form

$$\mathbf{y}_r = \mathbf{H}_r \mathbf{x}_r + \mathbf{n}_r, \quad (3)$$

with the obvious definitions of \mathbf{y}_r , etc. Note that \mathbf{H}^H translates to \mathbf{H}_r^T ($(\cdot)^H$ denotes the Hermitian, $(\cdot)^T$ the transpose of a matrix).

This step is essential since now the noiseless received signal $\mathbf{H}_r \mathbf{x}_r$ can be considered as a point of the lattice specified by the (generator) matrix \mathbf{H}_r , if \mathbf{x}_r is taken from the set of integers, i.e., $\mathbf{x}_r \in \mathcal{A}_r^{2K} \subset \mathbb{Z}^{2K}$. Only a subset of this lattice (the points within a boundary region around the origin) is actually used, as practical systems typically have only limited average and peak output power. The application to the commonly used QAM constellations, which are taken from translates of the integer lattice, is straightforward and will be discussed along the way.

An optimum detector performs maximum-likelihood detection, i.e., calculates

$$\hat{\mathbf{x}}_r = \underset{\mathbf{x}_r \in \mathcal{A}_r^{2K}}{\operatorname{argmin}} \|\mathbf{y}_r - \mathbf{H}_r \mathbf{x}_r\|^2, \quad (4)$$

where we assume (as always) that the detector has perfect knowledge of the channel state, \mathbf{H}_r . The size of the search space, $|\mathcal{A}_r^{2K}|$, prohibits the use of maximum-likelihood detection in practical systems (especially for large constellations and numbers of dimensions).

In [1] a general procedure for efficient lattice decoding, which performs the minimization in (4) with respect to \mathbb{Z}^{2K} , was presented. Conceptually it consists of the following steps (refer to [1] for details):

1. reduce lattice basis
2. perform closest-point search in reduced lattice
3. transform result to original lattice

Using this procedure we can perform (near) maximum-likelihood detection in MIMO systems. However, the closest point search, Step 2, (see also, e.g., [10]) is still extremely time-consuming for large constellations and dimensions K . It is therefore interesting to consider sub-optimum but less complex variants of this procedure.

It turns out that the first step of reducing the lattice basis has an essential impact on the number of operations required in the closest-point search [1]. For “small” dimensions (< 15) the

speedup due to the application of the LLL algorithm for Step 1 is up to an order of magnitude, and similar to that possible with the more complex Korkine-Zolotareff (KZ) reduction, as simulations shown in [1] confirm. For larger-dimensional lattices KZ reduction is superior and provides a speedup of up to two orders of magnitude.

The new idea in [11] is to combine the lattice reduction preprocessing step with “traditional” low-complexity detectors with the hope that the detection performance of these is improved as well. In the light of this framework that approach seems natural. In [11] an optimum reduction algorithm for a 2×2 system is given. For higher dimensions we expect the complexity of the optimum reduction to be prohibitive, and as our interest mainly lies in medium-size MIMO systems (e.g., 4×4 or 8×8), we will concentrate on the application of LLL reduction in the sequel.

III. THE LLL REDUCTION OF A LATTICE

As shown above, the noiseless received points in the communication scenario correspond to points of the lattice $\mathbf{H}_r \mathbb{Z}^{2K}$.

Lattice (basis) reduction [9, 7, 2] optimizes the generating matrix of the lattice to obtain a “nicer” description of the lattice. It obtains

$$\mathbf{H}_{\text{red}} = \mathbf{H}_r \mathbf{P}, \quad (5)$$

where \mathbf{P} is a matrix with integer entries that has determinant 1, i.e., \mathbf{P}^{-1} also contains only integer entries (“unimodular matrix”). The matrix \mathbf{H}_{red} gives a “reduced” basis for the same lattice, i.e., $\mathbf{H}_r \mathbb{Z}^{2K} \equiv \mathbf{H}_{\text{red}} \mathbb{Z}^{2K}$. Using the LLL algorithm for basis reduction, the following properties of the columns of $\mathbf{H}_{\text{red}} = [\mathbf{h}_1, \dots, \mathbf{h}_K]$ can be achieved [2]:

$$\|\mathbf{h}_i\|^2 \leq 2\|\mathbf{h}_{i+1}\|^2 \quad (6)$$

$$\frac{|\mathbf{h}_i^\top \mathbf{h}_j|}{\|\mathbf{h}_i\|^2} \leq \frac{1}{2}, \quad j > i, \quad (7)$$

i.e., the vectors are sorted in length and are roughly orthogonal. (The length constraint (6) depends on a parameter used in the LLL algorithm; by varying it the factor 2 can be improved to $2/\sqrt{3}$ at the expense of somewhat higher complexity). An indication of the orthogonality of the reduced matrix is its condition number (defined as $\text{cond}(\mathbf{H}) \triangleq \|\mathbf{H}_r\|_2 \|\mathbf{H}_r^{-1}\|_2$, using the spectral matrix norm $\|\cdot\|_2$ [8]), which describes its behaviour with respect to inversion. Orthogonal matrices are “perfectly conditioned” with condition number 1, while matrices which are nearly singular have large condition numbers. As an illustration we have plotted the pdfs of the condition numbers of random 8×8 real-valued matrices and their reduced counterparts in Fig. 1 (entries as in (2), corresponding to complex Gaussian i.i.d. random variables of variance 1). Not only is the spread in the condition numbers of the LLL-reduced matrices much smaller, but their average value is considerably smaller as well.

IV. LATTICE-REDUCTION-AIDED DETECTION

As the basis change does not change the lattice, we can now interpret the noiseless received signal points as points in the lattice described by \mathbf{H}_{red} . Since the matrix \mathbf{H}_{red} has much

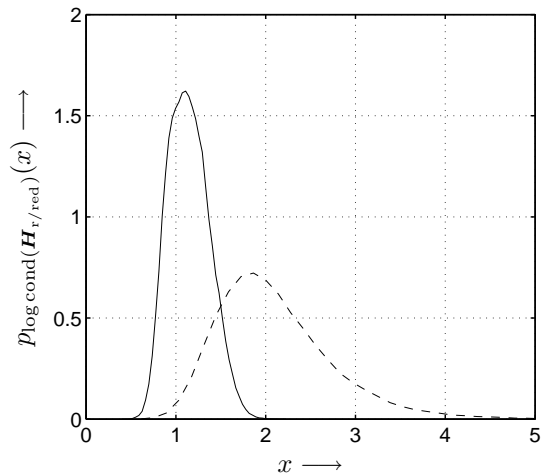


Figure 1: Probability density functions of the natural logarithm of $\text{cond}(\mathbf{H}_r)$ (dashed) and $\text{cond}(\mathbf{H}_{\text{red}})$ (solid) for 8×8 matrices.

“nicer” properties than \mathbf{H}_r , for instance with respect to inversion as shown by the results in Fig. 1, it is “easier” to detect the transmitted symbol when noise is present using a simple low-complexity detector. Having found these estimates, we can change the lattice basis to obtain $\hat{\mathbf{x}}_r$.

Using the LLL reduction of \mathbf{H}_r given in (5), we can write $\mathbf{H}_{\text{red}}^{-1} = \mathbf{P}^{-1} \mathbf{H}_r^{-1}$, and if we apply this matrix to the received vector \mathbf{y}_r , i.e., perform linear equalization of \mathbf{H}_{red} , we obtain

$$\mathbf{y}'_r = \mathbf{H}_{\text{red}}^{-1} \mathbf{y}_r \quad (8)$$

$$= \mathbf{H}_{\text{red}}^{-1} \mathbf{H}_r \mathbf{x}_r + \mathbf{H}_{\text{red}}^{-1} \mathbf{n}_r \quad (9)$$

$$= \mathbf{P}^{-1} \mathbf{x}_r + \mathbf{H}_{\text{red}}^{-1} \mathbf{n}_r. \quad (10)$$

We see that the signal \mathbf{y}'_r contains the desired signal \mathbf{x}_r plus noise $\mathbf{H}_{\text{red}}^{-1} \mathbf{n}_r$. Since the columns of \mathbf{H}_{red} are “rather orthogonal”, only relatively small noise enhancement and coloring is present.

Since the matrix \mathbf{P}^{-1} contains only integer entries, we have $\mathbf{P}^{-1} \mathbb{Z}^{2K} \equiv \mathbb{Z}^{2K}$, and consequently we can quantize \mathbf{y}'_r to \mathbb{Z}^{2K} . The estimates corresponding to the original signal points can now be obtained by

$$\hat{\mathbf{x}}_r = \mathbf{P} Q_{\mathbb{Z}^{2K}} \{\mathbf{y}'_r\}, \quad (11)$$

where $Q_{\mathbb{Z}^{2K}} \{\cdot\}$ denotes the quantization operation to the $2K$ -dimensional integer lattice. For the integer lattice \mathbb{Z}^{2K} , lattice quantization is identical to simple per-component quantization, i.e., $Q_{\mathbb{Z}^{2K}} \{\mathbf{y}'_r\} = [Q_{\mathbb{Z}} \{y'_{r,1}\}, \dots, Q_{\mathbb{Z}} \{y'_{r,2K}\}]^\top$.

As this quantization does not regard the boundary region of the constellation used for \mathbf{x}_r , the points obtained in $\hat{\mathbf{x}}_r$ stem from an extended version of the original constellation, and, in a final step, points that happen to lie outside the boundary region of the original constellation have to be assigned to the nearest point within the boundary region.

The commonly used QAM constellations are centered on the origin of the complex plane, and consequently square constellations that have sizes corresponding to powers of 2 do not include the origin, e.g., for 4-QAM: $\mathcal{A} = \{\pm \frac{1}{2} \pm j \frac{1}{2}\}$. That is,

$\mathbf{H}_r \mathbf{x}_r$ is a translate of a lattice, specifically,

$$\mathbf{x}_r = \tilde{\mathbf{x}}_r + \left[\frac{1}{2}, \dots, \frac{1}{2}\right]^T, \quad \tilde{\mathbf{x}}_r \in \mathbb{Z}^{2K}. \quad (12)$$

Then

$$\mathbf{y}'_r = \mathbf{P}^{-1} \tilde{\mathbf{x}}_r + \mathbf{P}^{-1} \left[\frac{1}{2}, \dots, \frac{1}{2}\right]^T + \mathbf{H}_{\text{red}}^{-1} \mathbf{n}_r, \quad (13)$$

and a comparison with (10), (11) shows that by

$$\hat{\mathbf{x}}_r = \mathbf{P} \mathcal{Q}_{\mathbb{Z}^{2K}} \left\{ \mathbf{y}'_r - \mathbf{P}^{-1} \left[\frac{1}{2}, \dots, \frac{1}{2}\right]^T \right\} + \left[\frac{1}{2}, \dots, \frac{1}{2}\right]^T \quad (14)$$

estimates for these shifted lattice constellations are obtained.

As the inverse of the channel matrix can be expressed as $\mathbf{H}_r^{-1} = \mathbf{P} \mathbf{H}_{\text{red}}^{-1}$, the procedure outlined above says that the receiver (1) compensates for the channel corresponding to \mathbf{H}_{red} by means of linear equalization, (2) quantizes the output to the constellation lattice, and (3) applies \mathbf{P} to obtain the corresponding point in the original constellation lattice basis. Hence we can also use a step-by-step decision feedback approach to compensate for $\mathbf{H}_{\text{red}}^{-1}$ [11].

The V-BLAST decision feedback detection algorithm [6] can be formulated starting from the decomposition $\mathbf{H}_{\text{red}} = \mathbf{W}^{-1} \mathbf{B} \mathbf{P}_{\text{vb}}$, where \mathbf{P}_{vb} is the permutation matrix corresponding to the optimum sorting order of the subchannels, and \mathbf{W} is the interference suppression filter chosen such that the matrix \mathbf{B} is lower-triangular with unit diagonal (e.g., obtained from the QR-decomposition). Forming $\mathbf{y}'_{\text{vb}} = \mathbf{W} \mathbf{y} = \mathbf{B} \mathbf{P}_{\text{vb}} \mathbf{P}^{-1} \mathbf{x}$, due to the structure of \mathbf{B} we can quantize the first component and subtract its interference into the other dimensions using the entries in the first column of \mathbf{B} . The same procedure is applied for the second, third, etc., components. The estimate for the original \mathbf{x} is obtained by reordering the components using \mathbf{P}_{vb}^T and changing the lattice basis using \mathbf{P} .

In the case of the QAM constellations, we can take an approach similar to (14), i.e., form

$$\tilde{\mathbf{y}}'_{\text{vb}} = \mathbf{y}'_{\text{vb}} - \mathbf{B} \mathbf{P}_{\text{vb}} \mathbf{P}^{-1} \left[\frac{1}{2}, \dots, \frac{1}{2}\right]^T \quad (15)$$

before starting the V-BLAST algorithm with quantization to \mathbb{Z}^{2K} , and add $\left[\frac{1}{2}, \dots, \frac{1}{2}\right]^T$ to the result after reordering and basis change.

An even simpler approach than the inversion of \mathbf{H}_{red} would be to use the matched filter receiver on \mathbf{y}_r , i.e., $\mathbf{D} \mathbf{H}_{\text{red}}^T$, where \mathbf{D} is a diagonal scaling matrix that ensures a transfer function of 1 per component. This is the approach used in [2], however in the present setting its performance is not satisfactory.

Simulation results for a 4×4 scenario using 4-QAM in all components of \mathbf{x} are shown in Figure 2. We have plotted the average symbol error rate $\overline{\text{SER}}$ over the ratio of average received energy per bit \bar{E}_b to one-sided power spectral density N_0 (cf., e.g., [5]). For comparison curves for maximum-likelihood detection (MLD) as well as linear equalization (LE) and the V-BLAST detection scheme are shown. Note that the diversity order achieved by lattice-reduction-aided linear equalization (LR-LE) is the same as that of MLD, namely 4, while that of V-BLAST is the same as that of LE, which is 1 since the number of receiving antennas equals the number of transmit antennas. Lattice reduction combined with the V-BLAST algorithm (LR-VB) gives some additional gain over

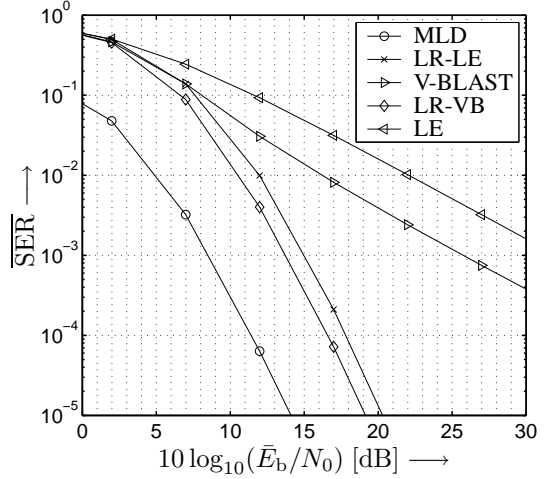


Figure 2: Simulation results for lattice-reduction-aided detection in a 4×4 system using 4-QAM compared to other strategies.

LR-LE, but the significant part of the gain comes from the increased diversity order by using the lattice reduction step.

V. LATTICE-REDUCTION-AIDED PRECODING

Instead of performing the LLL reduction on \mathbf{H}_r we can also perform LLL reduction on \mathbf{H}_r^T to obtain

$$\mathbf{H}_{\text{red}} = \mathbf{P} \mathbf{H}_r. \quad (16)$$

This relationship can be used as a starting point for a lattice-reduction-aided linear preequalization scheme, assuming perfect channel state information at the transmitter (for a review of other popular precoding schemes the reader is referred to [4]): by preprocessing the transmit signal with $\mathbf{H}_{\text{red}}^{-1}$, we obtain

$$\mathbf{y}_r = \mathbf{H}_r \mathbf{H}_{\text{red}}^{-1} \mathbf{x}_r + \mathbf{n}_r, \quad (17)$$

and the estimate for \mathbf{x}_r is given by

$$\hat{\mathbf{x}}_r = \mathbf{P} \mathcal{Q}_{\mathbb{Z}^{2K}} \{ \mathbf{y}_r \}. \quad (18)$$

Here the property of \mathbf{H}_{red} of consisting of close-to-orthogonal rows leads to a significant reduction in required transmit power for the linearly preequalized transmit signal $\mathbf{H}_{\text{red}}^{-1} \mathbf{x}_r$ compared to that of conventional linear preequalization via \mathbf{H}_r^{-1} .

For the usual QAM constellations the same shifting operation as in the lattice-reduction-aided detection scheme needs to be applied.

Note that just as decision-feedback equalization in the guise of the V-BLAST algorithm is possible in the lattice-reduction-aided detection scenario, Tomlinson-Harashima precoding for MIMO channels [5] can be combined with the lattice-reduction approach to obtain a non-linear precoding scheme.

Simulation results for a 4×4 system using 4-QAM are shown in Fig. 3. We compare linear preequalization (LPE) based on \mathbf{H}_r with lattice-reduction-aided linear preequalization (LR-LPE). In this simulation transmit power was allowed to vary according to the channel realization, i.e., it is high if

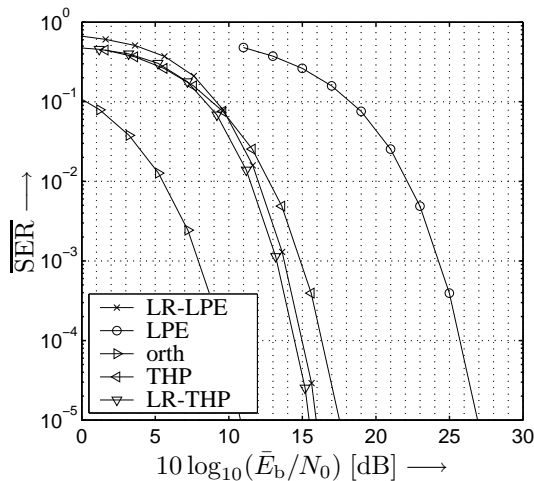


Figure 3: Simulation results for lattice-reduction-aided pre-equalization without a fixed power constraint in a 4×4 system using 4-QAM compared to standard linear pre-equalization.

some subchannels are “bad” and lower if the subchannels are “good”. A significant power efficiency gain of approximately 10 dB can be observed. Since we only look at average transmit power, we have in fact “infinite diversity” and the resulting SER curves are just the curve for the AWGN channel shifted by the average increase of transmit power. For reference we have given the curve that is achieved by a system that has 4 independent, i.e., orthogonal, fading channels with diversity order 4 each and linear pre-equalization (note that it shows only slightly worse performance than transmission over the AWGN channel, where a SER of 10^{-5} is achieved at roughly 10 dB E_b/N_0). Tomlinson-Harashima precoding (labeled THP in the figure) also shows a significant reduction in required transmit power, and combined with lattice-reduction-aided precoding (LR-THP) a small gain with respect to lattice-reduction-aided precoding is achieved.

If we impose a fixed transmit power constraint, i.e., only vary the distribution of the transmit power to the dimensions in order to equalize their error rates, but use the same average transmit power for every realization of the MIMO channel matrix \mathbf{H} , the gain becomes even more dramatic. The corresponding simulation results are shown in Fig. 4. Now the curve for LPE degrades to diversity order 1, as does linear equalization at the receiver side, while LR-LPE still approaches the full diversity order of 4. For comparison the curve achieved by a system with 4 independent fading channels with diversity order 4 each and linear pre-equalization is shown as well. The Tomlinson-Harashima precoding (THP) curves also show the expected behaviour: plain THP shows a significantly increased power efficiency compared to linear precoding, but diversity order 1, while combined with lattice-reduction-aided precoding the full diversity order is visible (LR-THP).

VI. CONCLUSIONS AND FURTHER WORK

We have shown that the performance of lattice-reduction-aided detection and precoding is superior to that of “conventional” detectors or precoding schemes. In the detection

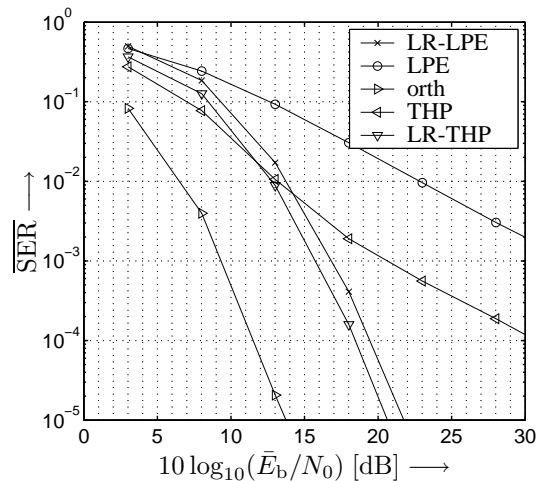


Figure 4: Simulation results for lattice-reduction-aided linear pre-equalization for fixed transmit power in a 4×4 system using 4-QAM compared to linear pre-equalization.

setting lattice-reduction-aided schemes approach the performance of maximum-likelihood detection (MLD); but while MLD complexity is exponential in the number of dimensions and constellation size, the complexity of the LLL algorithm is polynomial in the number of dimensions, and it is only required once for each transmission burst that experiences the same channel matrix \mathbf{H} .

In further work, the LLL algorithm, which was applied here in a “black box” fashion, should be studied in more detail. In particular lattice reduction according to other criteria, as mentioned in [1], as well as efficient implementation should be looked into.

REFERENCES

- [1] E. Agrell, T. Eriksson, A. Vardy, and K. Zeger. Closest point search in lattices, *IEEE Transactions on Information Theory*, vol. 48, no. 8, pp. 2201–2214, Aug. 2002.
- [2] K. L. Clarkson, W. Sweldens, and A. Zheng. Fast Multiple Antenna Differential Decoding, *IEEE Transactions on Communications*, Vol. 29, No. 2, pp. 253–261, 2001.
- [3] H. Cohen. *A Course in Computational Algebraic Number Theory*. Springer Verlag, Berlin, Germany, 1993.
- [4] R. F. H. Fischer. *Precoding and Signal Shaping for Digital Transmission*, John Wiley & Sons, New York, 2002.
- [5] R. F. H. Fischer, C. Windpassinger, A. Lampe, and J. B. Huber. Space-Time Transmission using Tomlinson-Harashima Precoding, *4th Intern. ITG Conf. on Source and Channel Coding*, Berlin, January 2002.
- [6] G. J. Foschini, G. D. Golden, R. A. Valenzuela, and P. W. Wolniansky. Simplified Processing for High Spectral Efficiency Wireless Communication Employing Multi-Element Arrays, *IEEE Journal on Selected Areas in Communications*, JSAC–17, pp. 1841–1852, November 1999.
- [7] J. von zur Gathen and J. Gerhard. *Modern Computer Algebra*, Cambridge University Press, Cambridge, UK, 2nd edition, 2002.
- [8] R. A. Horn and C. R. Johnson. *Matrix analysis*, Cambridge University Press, Cambridge, UK, 1985.
- [9] A. K. Lenstra, H. W. Lenstra, L. Lovász. Factoring polynomials with rational coefficients, *Math. Ann.*, 261:515–534, 1982.
- [10] H. Vikalo and B. Hassibi. Modified Fincke-Pohst Algorithm for Low-Complexity Iterative Decoding over Multiple Antenna Channels, in *Proceedings of IEEE ISIT*, Lausanne, Switzerland, July 2002.
- [11] H. Yao and G. W. Wornell. Lattice-Reduction-Aided Detectors for MIMO Communication Systems, in *Proceedings of IEEE Globecom 2002*, Taipei, Taiwan, November 2002.