

A VLSI Optimised Parallel Tree Search for MIMO

Geoff Knagge¹, Graeme Woodward², Steven R. Weller¹, Brett Ninness¹

Abstract— Multiple Input-Multiple Output (MIMO) systems are of great interest due to their ability to significantly increase the capacity of wireless communications systems, but for these to be useful they must also be practical for implementation in very large scale integrated (VLSI) circuits. A particularly difficult part of these systems is the detector, where the maximum-likelihood (ML) solution cannot be directly implemented due to its exponential complexity.

Lattice decoders, such as the sphere search, exhibit near-ML performance with reduced complexity, but their application is still limited by computational requirements. Here, a number of optimisations are presented, designed to reduce the computational cost of the sphere search in the context of VLSI implementation for MIMO applications. We also propose parallel implementation strategies for such a detector, suitable for implementation in VLSI. This is then combined with a single-pass tree search approach and it is demonstrated that it can be designed so that the error-rate performance is not significantly impaired.

Index Terms— MIMO, Sphere Detector, Tree Search, VLSI, ASIC Implementation.

I. PROBLEM BACKGROUND

Multiple input multiple-output (MIMO) systems utilise spatial diversity between arrays of transmit and receive antennae to achieve high data rates [1]. Current research allows for data rates up to 28.8Mbps [2], however while further increases in data rate are theoretically possible, the practicality is limited. The very large scale integrated (VLSI) circuit in [2] is a 4x4 MIMO receiver using a QPSK constellation, but it is desirable to achieve a higher data rate.

There are two methods in which this can be done. While one such method is to increase the constellation size, [3] indicates that it is preferable to first increase the antennae dimensionality. However, this increases the size of the detection problem by an exponential order. This can be explained by examination of the following system model for a MIMO channel with n transmitters and m receivers:

$$\mathbf{y} = \mathbf{H}\mathbf{s} + \mathbf{n}. \quad (1)$$

Here, \mathbf{y} is an m -vector with each element representing the received despread sample from one antennae, \mathbf{s} is the n -vector of transmitted symbols, \mathbf{n} is a length m noise vector, and \mathbf{H} is an $m \times n$ matrix of channel coefficients between antennae.

The MIMO detection problem is to solve (1) for \mathbf{s} , with knowledge of the channel \mathbf{H} , the received symbol estimates \mathbf{y} , and that the elements of \mathbf{s} are known to be from a finite set

¹: With the School of Electrical Engineering and Computer Science, The University of Newcastle, Callaghan NSW 2308, Australia, e-mail: {gknagge, steve, brett}@ee.newcastle.edu.au, and partly supported by the Australian Research Council under Linkage Grant LP0211210. ²: With Agere Systems, Lvl 3, 11-17 Khartoum Rd, North Ryde, NSW 2113, Australia, e-mail: graemew@agere.com

of constellation points. The problem is in fact identical to the multiuser detection (MUD) problem, but with different meanings given to the variables [4]. The maximum-likelihood (ML) multiuser detector [5] involves finding

$$\tilde{\mathbf{s}} = \arg \min_{\mathbf{s} \in \Lambda} \|\mathbf{y} - \mathbf{H}\mathbf{s}\|^2, \quad (2)$$

where Λ is the set of possible decisions over all users.

This is a combinatorial optimisation problem that is NP-hard, and hence impractical for all but the smallest of problems. In particular, with n transmitters, each transmitting from a constellation of size 2^q , the complexity of the problem becomes $O(2^{qn})$. The MIMO receiver chip in [2] has a search space of $2^{2 \times 4} = 256$ and is able to perform a brute force search, but doubling the number of antennae to 8 increases the search space to 65536, which is beyond current feasible receiver designs.

Lattice decoding, especially the sphere search variant, is regarded as a very promising candidate for practical, high performance, near-optimal ML detection algorithms. It has recently received wide attention for its potential application to space-time decoding, and MUD for multi-carrier code-division multiple-access (MC-CDMA) [6] and direct sequence CDMA (DS-SS) [7] systems. An equivalent algorithm, the closest point search, is well described by [8], and the applications to MIMO channels have been studied in [9].

The sphere search can be described as a variation of a tree search, making use of simplifications to greatly reduce the search space to only a few percent of the points considered in the full ML problem. However, the algorithm and the associated preprocessing is still computationally intensive, and optimisations need to be found to further reduce its complexity.

To address this problem, this paper proposes a variation in the sphere search technique, and investigates strategies that may be used to assist in the VLSI implementation of such a detector. Section II introduces the sphere search algorithm and its computational requirements. In Section III, our proposed parallel search strategy and single-pass simplification is described, and Section IV describes how this allows for a simplified architecture. Section V presents the relative complexity comparisons of this technique against standard algorithms. Finally, the paper is concluded with a summary and evaluation of these proposals.

II. STANDARD SPHERE SEARCH ALGORITHM

The optimal solution, $\tilde{\mathbf{s}}$, in (2) may be expressed as [9]:

$$\tilde{\mathbf{s}} = \arg \min_{\mathbf{s} \in \Lambda} (\mathbf{s} - \hat{\mathbf{s}})^H \mathbf{H}^H \mathbf{H} (\mathbf{s} - \hat{\mathbf{s}}). \quad (3)$$

Here, Λ is the lattice of possible decisions over all transmitters, \mathbf{H}^H denotes the conjugate transpose of \mathbf{H} , and $\hat{\mathbf{s}}$ is the unconstrained ML estimate of \mathbf{s} , given by

$$\hat{\mathbf{s}} = (\mathbf{H}^H \mathbf{H})^{-1} \mathbf{H}^H \mathbf{y}. \quad (4)$$

The sphere search has the additional concept of a radius, r , which is a threshold that defines the maximum distance around the search centre, $\hat{\mathbf{s}}$, that will be searched. The problem becomes one of solving (3) for cases $\mathbf{s} \in \Lambda$ that satisfy

$$(\mathbf{s} - \hat{\mathbf{s}})^H \mathbf{H}^H \mathbf{H} (\mathbf{s} - \hat{\mathbf{s}}) \leq r^2. \quad (5)$$

By utilising either a Cholesky or QR decomposition, an upper triangular \mathbf{U} can be obtained such that $\mathbf{U}^H \mathbf{U} = \mathbf{H}^H \mathbf{H}$, with the added constraint that the diagonals of \mathbf{U} are non-negative, so that (5) may be rewritten as

$$\sum_{i=1}^K \left| u_{ii} (s_i - \hat{s}_i) + \sum_{j=i+1}^K u_{ij} (s_j - \hat{s}_j) \right|^2 \leq r^2, \quad (6)$$

where u_{ij} is the (i, j) -th element of \mathbf{U} .

The upper triangular nature of \mathbf{U} allows the optimisation problem to be structured as a tree search, with each transmitter representing one level of the tree, and the branches representing a choice of one of the constellation points available for each transmitter. Associated with each branch is a cost contribution, represented by one term of the outer summation in (6). Each leaf then represents the entire collection of decisions, and has a cost that is the sum of the cost contributions associated with each of the branches taken to reach that leaf from the root of the tree.

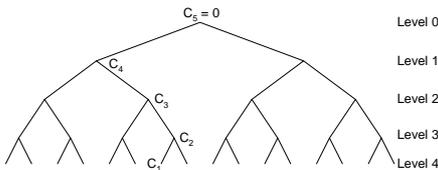


Fig. 1. Example of cost allocations for tree search decisions. Note that “Level 0” does not actually exist in the search, since the first decision is represented by “Level 1”.

Fig. 1 presents an example for a $K = 4$ antennae problem, with binary decisions for each transmitter. The first decision is represented by level 1, which corresponds to the last row of \mathbf{U} . Defining cost C_{K+1} as 0, the cost for the node at level $K+1-n$ of the tree is

$$C_n = C_{n+1} + \left| u_{nn} (s_n - \hat{s}_n) + \sum_{j=n+1}^K u_{nj} (s_j - \hat{s}_j) \right|^2. \quad (7)$$

It is apparent from (7) that all descendants of a given node will have a cost that is not smaller than the cost of that parent node. Therefore, if a given node has a cost greater than the current radius, all of its descendants will also have a greater cost, and so the tree may be pruned at that point. It is via this pruning that the sphere search significantly reduces the search space and therefore the complexity of the detection problem.

When a leaf node is evaluated, it may be added to a “leading candidates list” of an arbitrary fixed length n , which contains the best n leaves found to date. This list is used, after the search

is complete, to generate soft information about the decision for each transmitter, as described in Section IV-C. Once the leading candidates list is filled, the radius becomes equal to the cost of the highest cost leaf in that list. Further additions to the list will result in that highest cost leaf being discarded, and the radius being adjusted to match the new highest cost leaf within the list.

III. PROPOSED PARALLEL SINGLE-PASS SEARCH

To make the processing of larger problems more practical, a distributed processing approach may be used to consider multiple paths of the tree at one time in an attempt to generate the leading candidates list in a shorter time. A key requirement of this approach is that all parallel searchers act on the same level of the tree at any one time, providing a simple solution to memory contention issues. This means that all searchers will require the same row of the decomposed matrix, and will also be utilising the same cells of that matrix at any given time.

In the parallel scheme considered in Fig. 2, each of the two parallel searchers evaluates the cost of its two children ($[A, B]$ and $[C, D]$ respectively) on each step. The two most promising of these children are then chosen and assigned to the distributed entities on the next step of the search, while the others are stored for later consideration.

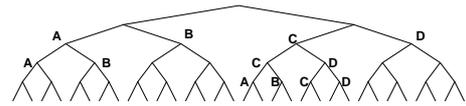


Fig. 2. On each level, the children of the current nodes (labeled A-D) are evaluated by the two parallel searchers. The best two are then selected as the current nodes for the next iteration.

Once the leaves of the tree are reached, a traditional sphere search back-tracks up the tree to evaluate the validity of the other unexplored paths. These may in turn find additional leaf nodes, some of which may be more promising than the current ones, and the process then continues until all remaining branches are either explored or discarded as non-promising.

To aid in evaluating the computation penalties or gains from increasing the amount of parallelism, simulations were conducted on a multiuser detection problem to count the number of leaves added to the promising candidates list as the parallelism increased. The key result, shown in Fig. 3, is that as the number of parallel searchers increased, the number of candidates found converged to the target list size. That is, the first set of candidates encountered were kept as the final result. Backtracking yielded few, if any, better candidates.

Hence, the proposed “single-pass” approach simply instantiates an appropriate amount of parallelism, and terminates the search as soon as the first set of leaf nodes is reached, thus entirely eliminating the need for back-tracking. While [10] alludes to a similar underlying idea, the approach and implementation strategies presented here allow for a more modular and flexible design.

Using the parallel search strategy described, the first set of candidates produced by the searchers are very likely to be “good”, even though some better candidates may have been

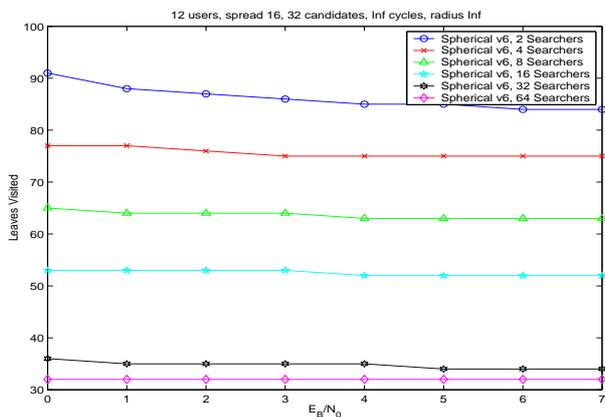


Fig. 3. Comparison of the number of candidates found as the number of parallel searchers increases.

missed. Since the only purpose of these candidates is to produce soft information, as described in section IV-C, it is sufficient to do this as long as there are enough candidates generated. A larger number of parallel searchers will provide better soft information, but will incur an increase in complexity. However, as described in section IV, an increase in the number of searchers does not necessarily involve significant duplication of hardware.

A. Performance

To demonstrate that the single-pass approach does not significantly impact error-rate performance, a software model of a flat fading MIMO system has been used. Independent $m \times n$ Rayleigh flat fading channel coefficients were generated, with a Doppler frequency corresponding to mobility of 200km/h assuming a 2GHz carrier frequency. A DS-CDMA system with a chip rate of 3.84Mcps was assumed, to roughly approximate a 3G cellular link.

Fig. 4 shows that in a 4×4 MIMO system, as few as 4 searchers is sufficient to generate a near optimal BER curve, however more searchers would be required to generate useful soft information.

An 8×8 system is shown in Fig. 5, a 12×12 system is shown in Fig. 6, and a 16×16 system is shown in Fig. 7. These figures show that an increasing number of searchers are required as the size of the system increases. However, as is discussed in the next section, increasing the number of searchers by this order does not dramatically increase the complexity when calculation sharing strategies are used.

The algorithm is generally suitable to any system where the sphere tree search would achieve optimal results. A near optimal result can always be obtained, however it involves a trade-off against how much parallelism is feasible for a specific application. This will depend on the acceptable BER performance, the processing time allowed, and the limits on the complexity of the VLSI circuit.

IV. ARCHITECTURAL STRATEGIES

The implementation benefits of the single-pass approach are very significant in comparison to other sphere search architectures, such as the one described in [11]. Since the search is

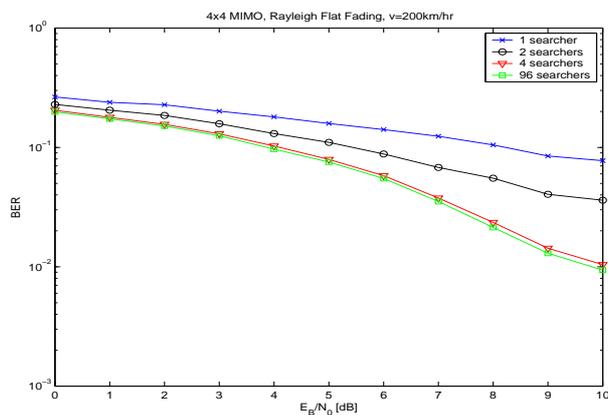


Fig. 4. Performance of single-pass approach on a 4×4 MIMO system.

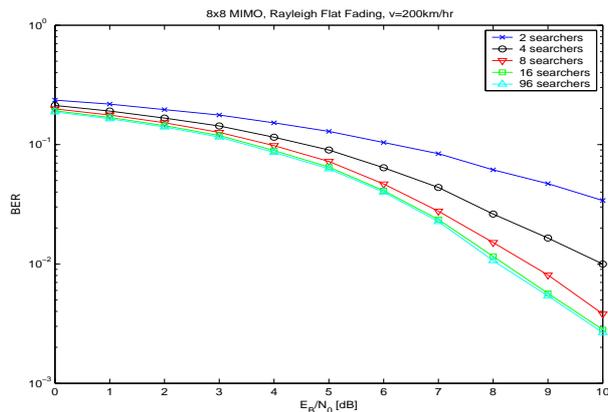


Fig. 5. Performance of single-pass approach on an 8×8 MIMO system.

terminated as soon as the first leaf nodes are reached, no back-tracking is required and so no stack structure is necessary to remember branches that were not followed. In addition, there no longer needs to be any concept of a radius, tree pruning strategy, nor any complete sorting of the leading candidates list. Thus, by removing the need for back-tracking, the associated overhead costs that previously limited the feasibility of implementation have been eliminated. In addition, our search strategy allows other optimisations to achieve a feasible MIMO detector in systems with large amounts of antennae.

While a detailed architecture is beyond the scope of this paper, Fig. 8, shows that implementation simply consists of a number of parallel search engines. The content of these engines is illustrated in Fig. 9. Each engine implements a series of simple node cost calculations, which are optimised as follows.

A. Calculation of Node Costs

From (7), the cost of any node consists of the cost of its parent, plus the result of a multiplication between one row of the decomposed matrix, U , and the vector of the difference between the search centre and the current estimate. The number of multiplications required may seem large, but can be greatly reduced by careful examination of the algorithm.

The key features that allow a simplified implementation strategy are the choice of either a binary or quaternary tree

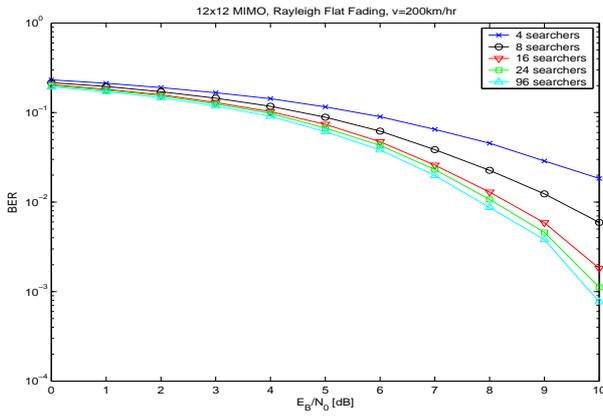


Fig. 6. Performance of single-pass approach on a 12x12 MIMO system.

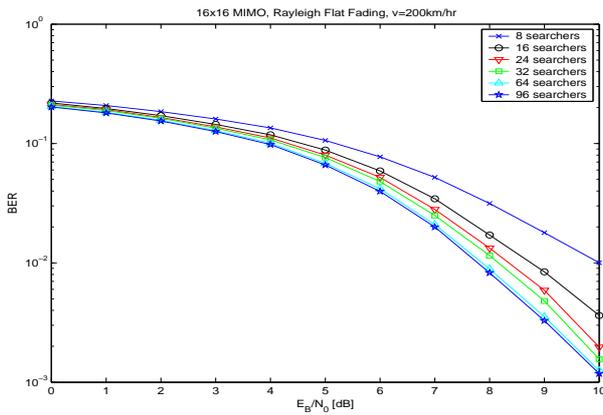


Fig. 7. Performance of single-pass approach on a 16x16 MIMO system.

(i.e. BPSK or QPSK constellations), and by requiring that all searchers process the same level of the tree at any given time.

Since all searchers are processing nodes on the same level of the tree, then from the expansion

$$u_{nj}(s_n - \hat{s}_n) = u_{nj}s_n - u_{nj}\hat{s}_n$$

it can be observed that only the first term is unique to each searcher. Furthermore, since the search tree is binary, the multiplication becomes a trivial matter of calculating $\pm u_{nj}$. The latter term is common to all searchers, and so only needs to be calculated once.

By examination of (7), it can be seen that the entire inner summation is common to the children of a given node. The computational cost then consists of

- A group of multiplications and additions common to all searchers;

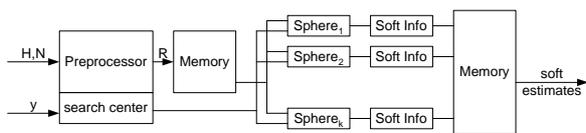


Fig. 8. A possible one-pass architecture, containing a number of parallel search engines to obtain required throughput.

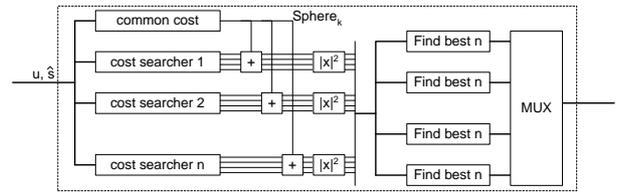


Fig. 9. Architecture of a single search engine.

- A group of additions for the first child of each searcher;
- A single addition for subsequent children of each searcher;
- Two squaring operations to find $|x|^2$ for each child.

Similar optimisations are also possible for a tree with quaternary decisions, where the decisions are of the form $\pm 1 \pm j$ with a precomputed scaling of $\sqrt{2}$.

B. Searcher Implementation

The single-pass approach allows for a very simple architecture of searcher hardware, allowing a large number of searchers to be easily instantiated. The additions required for node cost calculations can be performed in parallel as previously described, and it is not necessary to maintain a stack of unexplored options. The multiplications required to find the final cost of each node may either be performed in parallel, or time multiplexed, depending on the number of multipliers that may be feasibly implemented.

The only point where the number of searchers introduces possible complexity concerns is in the sorting of evaluated children on each level of the tree. With m parallel searchers, only the m best children are required. While standard full-sorting methods such as the bubble sort may be used, it is also acceptable to find the m best children in any order.

C. Generation of Soft Information

The required output of the sphere search is a soft decision for each transmitter's symbol, with the sign of the value representing the decision and the magnitude representing the reliability. Generally, a likelihood ratio of probabilities is required:

$$\text{LR}(\mathbf{y}) = \frac{P(s_k = -1|\mathbf{y})}{P(s_k = +1|\mathbf{y})}. \quad (8)$$

In a sphere list search, these probabilities can be determined directly from the leading candidates list. We note that the difference between the squared Euclidean distance in (2), and the leaf cost minimised in (3) is a constant, Δ . By defining the cost of a given leaf as d_s and applying Bayes' rule,

$$p(\mathbf{y}|\mathbf{s}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{\Delta}{2\sigma^2}} e^{-\frac{d_s^2}{2\sigma^2}}. \quad (9)$$

The probability of a "1" being transmitted by a particular transmitter is equal to the sum of the probabilities of all of the combinations containing a "1" for that given transmitter, and similarly for a "-1". It is not necessary to calculate the constant term $\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{\Delta}{2\sigma^2}}$, since it cancels when the likelihood ratio is computed.

If the costs of the best n solutions are known, then the others may be estimated by observing that their cost is at least as high as that of the worst known point. This value can then be substituted in place of the unknown costs. Alternatively, these unknown results may be ignored completely, since their contribution is likely to be relatively small. The variance σ^2 may be also approximated to a convenient constant without significantly affecting the performance.

The soft output log-likelihood ratio associated with the k th transmitter can then be determined by

$$\text{LLR}_k = \ln \frac{P(s_k = 1|y)}{P(s_k = -1|y)} \quad (10)$$

$$= \ln P(s_k = 1|y) - \ln P(s_k = -1|y) \quad (11)$$

A hard decision can be determined from the soft outputs by simply recording the sign of the output, with the magnitude representing the relative confidence of the decision.

V. COMPLEXITY

When considered for VLSI implementation, detection algorithms are generally compared in terms of the number of multiplications required per detection operation. In the case of the sphere search, this involves the calculation of (7) as the search progresses down the tree.

The sum of $u_{ij}\hat{s}_j$ is constant regardless of which branches are taken, and the calculation of $u_{ij}s_j$ is trivial in the case where s_j is taken from a binary or quaternary constellation. Therefore, the number of multiplications in the searching operation consists entirely of:

- p calculations of $|x|^2$ on each level, where p is the number of parallel evaluations at each level of the tree.
- $\frac{K(K+1)}{2}$ multiplications to evaluate $u_{ij}\hat{s}_j$. Since most of these are full complex multiplications, they are considered to be twice as complex as the $|x|^2$ calculations.

If the single pass approach is used, then the total amount of multiplications is at most $Kz + K(K+1)$, where z is the product of the size of the constellation and the number of parallel searchers. This analysis ignores the preprocessing requirements, such as performing the Cholesky decomposition, which are not changed by our proposed method and are approximately $O(K^3)$ in complexity. We demonstrate the feasibility of the preprocessing in [12]

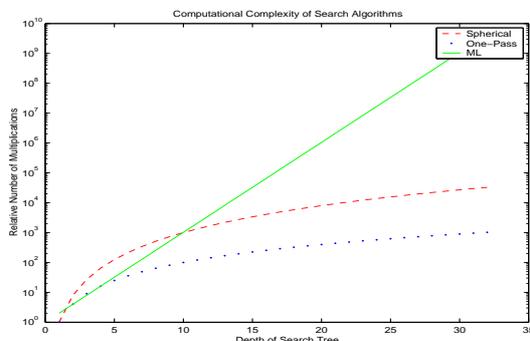


Fig. 10. Comparison of relative complexity of the various alternative detectors.

Fig. 10 shows the general trend of the computational complexity of each type of tree search technique. Due to its exhaustive search, the standard ML detector can only be used on very small problems, with K transmitters requiring 2^K cost calculations. The sphere search algorithm can be configured to have an approximate complexity of $O(K^3)$ [9]. Hence, it can handle significantly more transmitters, particularly if parallelism is exploited in a hardware implementation.

The single-pass approach offers not only computational savings, but the complexity of the hardware required is also significantly reduced. In particular, by removing the need for back-tracking, there is no need for a stack to store unexplored options, and associated control mechanisms.

VI. CONCLUSION

This paper has described a low complexity tree search algorithm that is feasible for VLSI application to MIMO detection. The results have demonstrated that the number of searchers needed to obtain a good result is not excessive, and is dependent on the number of antennae in the system. The calculation sharing strategies proposed in this paper will significantly reduce the complexity of the implementation of the parallel searchers, allowing more to be implemented. Furthermore, the single-pass approach also removes the need to generate and store information that would otherwise be required for back-tracking in the search tree.

REFERENCES

- [1] G. J. Foschini and M. J. Gans, "On limits of wireless communications in a fading environment when using multiple antennas," *Wireless Personal Communications*, vol. 6, pp. 311–335, March 1998.
- [2] D. Garrett, L. Davis, S. ten Brink, B. Hochwald, and G. Knagge, "A 28.8 mb/s 4x4 mimo 3g high-speed downlink packet access receiver with normalized least mean square equalization," in *IEEE International Solid-State Circuits Conference (ISSCC)*, pp. 420–421, Feb 2004.
- [3] L. M. Davis, D. C. Garrett, G. K. Woodward, M. A. Bickerstaff, and M. F. J., "System architecture and ASICs for a MIMO 3GPP-HSDPA receiver," in *57th IEEE Semiannual Vehicular Technology Conference, 2003. VTC 2003-Spring*, vol. 2, pp. 818–822, Apr 2003.
- [4] G. Knagge, G. Woodward, B. Ninness, and S. Weller, "An optimised parallel tree search for multiuser detection with VLSI implementation strategy," in *IEEE GLOBECOM*, Dec 2004.
- [5] S. Verdú, "Minimum probability of error for asynchronous Gaussian multiple-access channels," vol. 32, pp. 85–96, Jan. 1986.
- [6] L. Brunel, "Multiuser detection techniques using maximum likelihood sphere decoding in multicarrier CDMA systems," *IEEE Trans. Wireless Commun.*, vol. 3, pp. 949–957, May 2004.
- [7] L. Brunel and J. J. Boutros, "Lattice decoding for joint detection in direct-sequence CDMA systems," *IEEE Trans. Inform. Theory*, vol. 49, pp. 1030–1037, September 2003.
- [8] E. Agrell, T. Eriksson, A. Vardy, and K. Zager, "Closest point search in lattices," *IEEE Transactions on Information Theory*, vol. 48, pp. 2201–2213, Aug 2002.
- [9] B. Hochwald and S. ten Brink, "Achieving near-capacity on a multiple-antenna channel," *IEEE Trans. Information Theory*, vol. 51, pp. 389–399, Mar 2003.
- [10] K. Wong, C. Tsui, R. Cheng, and W. Mow, "A VLSI architecture of a k -best lattice decoding algorithm for MIMO channels," in *IEEE International Symposium on Circuits and Systems*, pp. III–273 – III–276, May 2002.
- [11] B. Widdup, G. Woodward, and G. Knagge, "A highly-parallel VLSI architecture for a list sphere detector," in *International Conference on Communications (ICC 2004)*, vol. 5, pp. 2720–2725, Jun 2004.
- [12] G. Knagge, L. Davis, G. Woodward, and S. R. Weller, "VLSI preprocessing techniques for MUD and MIMO sphere detection," in *Australian Communications Theory Workshop (AusCTW) 2005*, to appear, Feb 2005.