

# Solution of indefinite linear systems using an LQ decomposition for the linear constraints

**Citation for published version (APA):**

Schilders, W. H. A. (2009). *Solution of indefinite linear systems using an LQ decomposition for the linear constraints*. (CASA-report; Vol. 0906). Technische Universiteit Eindhoven.

**Document status and date:**

Published: 01/01/2009

**Document Version:**

Publisher's PDF, also known as Version of Record (includes final page, issue and volume numbers)

**Please check the document version of this publication:**

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

**General rights**

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

[www.tue.nl/taverne](http://www.tue.nl/taverne)

**Take down policy**

If you believe that this document breaches copyright please contact us at:

[openaccess@tue.nl](mailto:openaccess@tue.nl)

providing details and we will investigate your claim.

# Solution of indefinite linear systems using an LQ decomposition for the linear constraints

Wil H.A. Schilders

---

## Abstract

In this paper, indefinite linear systems with linear constraints are considered. We present a special decomposition that makes use of the LQ decomposition, and retains the constraints in the factors. The resulting decomposition is of a structure similar to that obtained using the Bunch-Kaufman-Parlett algorithm. The decomposition can be used in a direct solution algorithm for indefinite systems, but it can also be used to construct effective preconditioners. Combinations of the latter with conjugate gradient type methods have been demonstrated to be very useful.

*Key words:* indefinite system, linear constraint, LQ decomposition, Bunch-Kaufman-Parlett, conjugate gradients, incomplete preconditioning

---

## 1. Introduction

In 1977, the seminal paper by Meijerink and Van der Vorst [20] on using incomplete factorisations to construct preconditioners drastically changed the view on the use of iterative solution methods for linear systems. Since then, many preconditioning techniques based upon this concept have been published, and shown to be extremely effective for solving challenging and large industrial problems.

In the original Meijerink-Van der Vorst paper, the preconditioner is based upon an incomplete Cholesky decomposition. In later publications, and for special situations, the use of an incomplete Crout decomposition was advocated, and in [13] it was shown that this can be used to obtain even more efficient methods.

For indefinite symmetric linear systems, the straightforward use of incomplete Cholesky or incomplete Crout decompositions may lead to problems

with zero pivots caused by the fact that eigenvalues are located on both ends of the real axis. However, if the indefinite systems are of a special form, a technique has been developed that overcomes this problem. This technique is now known as the Schilders factorization [1, 7, 8, 19], and it has been used extensively for constructing different families of preconditioners for constraint linear systems [8].

The method itself was already developed in 1999, but the ideas behind it have never been published. These ideas are based upon using explicitly the structure of the linear systems, in particular the fact that there are different types of unknowns. This turns out to be the basics of the method, and paves the way for the development of new classes of decomposition techniques. Interesting is the fact that the original idea stems from the use of these decompositions in the area of electronic circuit simulation. The ideas are not restricted to this class of problems, but much more widely applicable as will be shown in this paper.

In order to set the scene, we first give a brief overview of solution methods for indefinite systems in Section 2. Then, in Section 3 the main idea that has led to the Schilders factorisation is explained in detail. This is the most important section of the paper, and the basis for further development of methods. In Section 4 the idea is put into a more abstract mathematical context, so that it becomes apparent that  $LQ$  factorisations can be used to achieve the same results. Finally, Section 5 discusses the use of the decomposition for preconditioning purposes.

## 2. A brief account of solution methods for indefinite systems

Consider linear systems of the form

$$\begin{pmatrix} A & B \\ B^T & 0 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} b \\ c \end{pmatrix}, \quad (1)$$

where the  $n \times n$  matrix  $A$  is symmetric and positive definite, and the  $n \times m$  matrix  $B$  is of full rank. Throughout this paper, we shall assume that  $m \leq n$ . Note that, since  $B$  is of full rank, the coefficient matrix in (1), which we shall denote by  $\mathcal{A}$ , is a nonsingular matrix. It should be noted that in several papers the notation is somewhat different from ours, in the sense that the role of  $B$  and  $B^T$  is interchanged.

Systems of the form (1) occur frequently in applications, and also when using specific numerical methods. To show this, we first give a number of

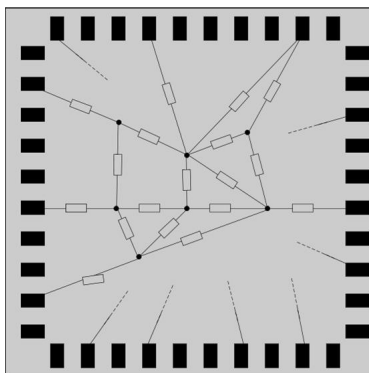


Figure 1: Resistor network

examples.

*Example 1.1*

Consider the use of the mixed finite element method for the discretisation of the problem

$$\nabla \cdot (a \nabla u) = f,$$

with suitable boundary conditions, and  $a = a(x, y) \geq \alpha > 0$ . The problem is reformulated as a system of first-order equations,

$$a^{-1} \sigma - \nabla u = 0,$$

$$-\nabla \cdot \sigma = -f.$$

Since the divergence and gradient operators are adjoints, the discretisation of this first-order system naturally leads to a system of the form (1). The resulting discrete problem is a "saddle point problem", and was analysed thoroughly in [2]. More information about mixed finite element methods, and the well known family of Raviart-Thomas mixed finite element spaces, can be found in [3, 22].

*Example 1.2*

Indefinite systems also occur quite naturally in the analysis of electronic circuits. Consider the network of resistors displayed in Figure 1. The voltage unknowns are associated with the nodes, whereas the currents are associated with the branches between nodes. The set of equations describing the

behaviour of this circuit is obtained by combining the so-called branch equations with the Kirchhoff laws for currents and voltages. Branch equations relate the voltage differences between nodes with the corresponding branch current. For example, a branch containing a resistor with value  $R$  will lead to a branch equation of the form

$$V_i - V_j - RI_{ij} = 0.$$

The set of all branch equations can, therefore, be written in the form

$$AI + BV = 0.$$

Kirchhoff's current law (KCL) states that, at each node in the network, the sum of all currents should be zero. Graph theoretical considerations lead to the conclusion that this can be formulated as

$$B^T I = 0,$$

thus demonstrating that the set of equations is of the form (1). This also holds for more general circuits, consisting of resistors, capacitors, inductors and nonlinear devices such as transistors and diodes [16].

Indefinite systems have attracted many researchers, and various approaches have been suggested to solve them. There are also some standard techniques. A straightforward method for solving the indefinite problem in (1) is direct elimination of the unknowns  $x$ :

$$x = A^{-1}b - A^{-1}By.$$

Substituting this in the second set of equations, leads to

$$B^T A^{-1}By = c - B^T A^{-1}b.$$

This approach is known as the range space method or the Schur complement method. At first glance it may look unattractive since, for sparse  $A$ , the matrix  $A^{-1}$  is full and hence the coefficient matrix  $B^T A^{-1}B$  is also a full matrix. However, in the special case of the Stokes problem,  $B$  and  $B^T$  are discrete versions of the gradient and divergence operator, whereas  $A$  is a discrete Laplace operator. Hence it is to be expected that  $A$ , in some sense, resembles the product of  $B$  and  $B^T$ , so that we may hope that the matrix  $B^T A^{-1}B$  is close to the identity, again, in some sense. This heuristic argument can be made more precise, and it can be shown that iterative methods

indeed perform well in this case. However, for more general problems the method often fails to provide a solution efficiently.

The counterpart of the range space method described is the null space method. Here the variables  $y$  are eliminated from the system, and this is done as follows. Assume that a basis for the null space of  $B^T$  is formed by the columns of the matrix  $Z$ , so that  $B^T Z = 0$ . Then we can write

$$x = B\hat{y} + Zz,$$

where  $\hat{y}$  is a special solution satisfying  $B^T B\hat{y} = c$ , and  $z$  is as yet unknown. Substituting the expression for  $x$  in the first set of equations, we obtain

$$AZz + By = b - AB\hat{y}.$$

Multiplying this by  $Z^T$  and using the fact that  $Z^T B = 0$ , we find

$$Z^T AZz = Z^T b - Z^T AB\hat{y}.$$

The coefficient matrix looks much more attractive than the one obtained in the range space method, provided  $A$  is a sparse matrix. However, in order not to perturb the sparsity too much, one will have to take care that the matrix  $Z$  is also rather sparse. This means that a sparse basis for the null space has to be used. For certain problems, this is indeed possible. In electronic circuit simulation, and in electromagnetics, the elements of the null space have a physical meaning and are the closed (current) loops which can be found from the topology (of the network, or the mesh). The dependence on the topology means that the basis has to be constructed only once. In [29] this technique, which makes use of an old algorithm published by Alex Orden, is described in more detail.

In some cases, it is possible to avoid the indefiniteness of the system entirely, by modifying the numerical method. In [14] it was suggested to introduce Lagrange multipliers on the edges of elements, and to impose continuity via these new unknowns. This means that the space of basis functions for the fluxes is enlarged, allowing fluxes to be discontinuous in principle. The enlarged system of equations is now of the form

$$\begin{pmatrix} \hat{A} & \hat{B} & C \\ \hat{B}^T & 0 & 0 \\ C^T & 0 & 0 \end{pmatrix} \begin{pmatrix} \hat{x} \\ \hat{y} \\ \lambda \end{pmatrix} = \text{rhs},$$

where  $\widehat{A}$  and  $\widehat{B}$  are (local) block diagonal matrices (note this is again an indefinite system). The latter property implies that the unknowns  $\widehat{x}$  and  $\widehat{y}$  can locally be eliminated (in fact a rather simple application of the range space method) and expressed in terms of the Lagrange multipliers. Hence a system for  $\lambda$  is obtained. The resulting coefficient matrix is larger than the original matrix, but is usually rather sparse. The approach can be quite effective for practical problems. In [4, 21], the use of this method is demonstrated for semiconductor device simulation, and it is shown that the physical meaning of the Lagrange multipliers is similar to that of the unknowns  $x$ .

The foregoing discussion clearly demonstrates that there are various ways of solving indefinite systems, but it is also clear that the treatment is far from uniform. Of course, many attempts have been undertaken to present a more unified treatment. The paper by Rusten and Winther [23] is one of the first to present an in-depth analysis of saddle point problems. Since then, many research papers have appeared, and we refer the reader to the thorough review paper by Benzi, Golub and Liesen [1] to obtain an excellent overview of the developments.

An entirely new concept for solving indefinite systems was presented at the 1999 conference on preconditioning techniques in Minneapolis. Wathen [30] presented the idea to keep the constraints in the preconditioning matrix, whereas in [24] a similar result was obtained in an entirely different way. In a sense, the approach is comparable to the ideas underlying the modified ICCG method: retain properties of the original system in the preconditioning matrix. Although there is no rigorous mathematical proof, this general concept often proves itself to be very useful. It restrict solutions of the numerical problem to a subspace that already contains characteristics of the original problem. Especially in the case of saddle point problems originating from optimization, it is important to satisfy the constraints. Also in model order reduction, a relatively new field in numerical mathematics, the advantage of retaining structural properties is recognized, cf. the chapters by Freund, and by Bai et al. in [26].

The approach presented by Wathen was detailed further in [17]. In that paper the preconditioning matrices for the system (1) are of the form

$$\mathcal{G} = \begin{pmatrix} G & B \\ B^T & 0 \end{pmatrix}. \quad (2)$$

From the analysis in [17] it follows that it may be very beneficial to retain the constraints and to use these special preconditioning matrices: the eigenvalue

distribution is improved as far as their impact on the convergence of iterative solution techniques is concerned. In fact, the preconditioned system has at least  $2m$  eigenvalues equal to 1.

Similar results were obtained in [25], where an incomplete decomposition was used as the basis for a preconditioned iterative method. Here, it was also found that there are at least  $2m$  eigenvalues equal to 1. In addition, it was proved that the eigenvalues of the preconditioned system are all real and positive (this is also proved in [17], under the condition that  $Z^T AZ$  and  $Z^T GZ$  are positive definite). The preconditioning matrix is also of the form (2), but the advantage of this preconditioning technique is that a decomposition of the matrix  $\mathcal{G}$  is available. Clearly, this is important in view of the efficiency of the iterative method. In fact, it is possible to reformulate the method in such a way that a full decomposition of the matrix  $\mathcal{A}$  is obtained, which can then be used to directly solve the indefinite linear systems rather than iteratively. This is one of the main results of this paper, and will be discussed in Section 4. In order to better understand the reasons for this decomposition, we will summarize and motivate the incomplete decompositions of [25] in Section 3. In Section 5, we discuss the use of the incomplete decompositions as a basis for preconditioned iterative solution methods.

### 3. Incomplete preconditioning using $1 \times 1$ and $2 \times 2$ blocks

The idea for the decomposition technique originates from problems in the electronics industry. In the area of electronic circuit simulation, huge systems of equations must be solved. If resistors, capacitors and inductors are used, these systems are linear, but when diodes and transistors are part of the circuit, the systems become extremely nonlinear. Newton-type methods, often in combination with continuation methods, are used to solve the nonlinear problems, whence large linear systems are at the core of most circuit simulation software. A detailed discussion of electronic circuit simulation and mathematical techniques associated with it can be found in [16].

Important for the context of the present paper is that the systems involved are of the form (1). Virtually all known circuit simulation packages (both in-house codes like Pstar and Titan, and commercially available codes like Spectre and Spice) use direct solvers for such systems. The proprietary solver Pstar of NXP Semiconductors uses a hierarchical set-up and solution procedure, due to the natural hierarchy of electronic circuits that are often made up of standard building blocks.



We are interested in using iterative procedures for the solution of these linear systems originating from electronic circuit simulation. As these systems naturally contain two different types of unknowns, the idea came up to use both  $1 \times 1$  and  $2 \times 2$  pivots, and first use a special re-ordering scheme based upon these pivots before performing an incomplete decomposition. The idea turned out to be effective, and also generalizable to other systems containing different types of variables. Also, it turned out that the method can be cast into a much more general form, without having to explicitly mention the  $1 \times 1$  and  $2 \times 2$  pivots. However, before presenting this more general class of methods, we present in this section the original idea based upon a coupling of the current and voltage unknowns, as we feel that this may inspire similar ideas for other types of multi-variable problems. Furthermore, it reveals clearly why the approach is effective.

Thus, in this section, we restrict ourselves to a special class of matrices  $B$ , namely those having the following properties:

$$B_{i,j} \in \{-1, 0, 1\} \quad \forall 1 \leq i \leq n, 1 \leq j \leq m.$$

We also assume that each row of  $B$  contains at most two non zero elements, which are of opposite sign:

$$\sum_{i=1}^m |B_{i,j}| \leq 2,$$

$$-1 \leq \sum_{i=1}^m B_{i,j} \leq 1.$$

As before, we assume that  $\text{rank}(B) = m$ . Matrices of this type are related to the so-called incidence matrices whose entries are 0 or 1. In fact, the matrices we are considering are differences of two incidence matrices.

Now let  $P : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$  be a permutation with the property that

$$B_{P(i),i} \neq 0,$$

and

$$B_{P(i),j} = 0 \quad \text{for } j > i. \quad (3)$$

In fact,  $B$  is permuted to lower trapezoidal form, meaning that the top  $m \times m$  part is lower triangular. Such a permutation  $P$  does not necessarily exist for all matrices considered in this paper. However, it is easy to show that for

matrices  $B$  of the above form there exist a row permutation  $P$  and a column permutation  $S$  such the permuted  $B$  is lower trapezoidal. Here we will assume that  $S(i) = i$ , but the generalization to  $S(i) \neq i$  is straightforward.

Next we define the permutation matrix  $Q$  by

$$Q = (\mathbf{e}_{P(1)}, \mathbf{e}_{n+1}, \dots, \mathbf{e}_{P(m)}, \mathbf{e}_{n+m}, \mathbf{e}_{P(m+1)}, \dots, \mathbf{e}_{P(n)}),$$

where  $\mathbf{e}_i \in R^{n+m}$  is the  $i$ -th unit vector. After permutation of rows and columns, we obtain the matrix

$$\tilde{A} = Q^T A Q,$$

Note that the vector of unknowns changes from  $(x_1, \dots, x_n, y_1, \dots, y_n)^T$  to  $(x_{P(1)}, y_1, \dots, x_{P(m)}, y_m, x_{P(m+1)}, \dots, x_{P(n)})^T$ .

In order to find a suitable preconditioning technique for the original indefinite system, we first transform it and propose an incomplete decomposition for the system with coefficient matrix  $\tilde{A}$ . After having found this decomposition, the preconditioning matrix is transformed back. The preconditioning matrix  $\tilde{M}$  for the transformed system is cast into the form

$$\tilde{M} \equiv (\tilde{L} + \tilde{D})\tilde{D}^{-1}(\tilde{L} + \tilde{D})^T,$$

where

$$\tilde{L} = \left( \begin{array}{ccc|ccc} 0 & \cdots & 0 & 0 & \cdots & 0 \\ \tilde{L}_{2,1} & \ddots & \vdots & \vdots & & \vdots \\ \vdots & \ddots & 0 & 0 & \cdots & 0 \\ \hline \tilde{L}_{m+1,1} & \cdots & \tilde{L}_{m+1,m} & 0 & \cdots & 0 \\ \vdots & & \vdots & \vdots & \ddots & \vdots \\ \tilde{L}_{n,1} & \cdots & \tilde{L}_{n,m} & \cdots & \tilde{L}_{n,n-1} & 0 \end{array} \right),$$

where  $\tilde{L}_{i,j}$  is a  $2 \times 2$  block for  $1 \leq j < i \leq m$ , a  $1 \times 1$  block whenever  $m \leq j < i \leq n$ , and a  $1 \times 2$  block in all other cases. We shall use the notation

$$\tilde{L} = \text{"lower"}(\tilde{A}).$$

Also,

$$\tilde{D} = \left( \begin{array}{ccc|ccc} \tilde{D}_1 & & & & & \\ & \ddots & & & & \\ & & \tilde{D}_m & & & \\ \hline & & & \tilde{d}_{m+1} & & \\ & & & & \ddots & \\ & & & & & \tilde{d}_n \end{array} \right).$$

When  $1 \leq j \leq m$ , we find that

$$\tilde{L}_{i,j} = \begin{pmatrix} A_{P(i),P(j)} & B_{P(i),j} \\ B_{i,P(j)}^T & 0 \end{pmatrix}.$$

The matrices  $\tilde{D}_1, \dots, \tilde{D}_m$  and the scalars  $\tilde{d}_{m+1}, \dots, \tilde{d}_n$  are required to be such that

$$\text{"diag"} \left( (\tilde{L} + \tilde{D})\tilde{D}^{-1}(\tilde{L} + \tilde{D})^T \right) = \text{"diag"}(\tilde{A}), \quad (4)$$

where the operation "*diag*" is defined as follows:

$$\text{"diag"}(\tilde{A}) \equiv \left( \begin{array}{ccc|ccc} \tilde{A}_{1,1} & \tilde{A}_{1,2} & & & & \\ \tilde{A}_{2,1} & \tilde{A}_{2,2} & & & & \\ & & \ddots & & & \\ & & & \tilde{A}_{2m-1,2m-1} & \tilde{A}_{2m-1,2m} & \\ & & & \tilde{A}_{2m,2m-1} & \tilde{A}_{2m,2m} & \\ \hline & & & & & \tilde{A}_{2m+1,2m+1} \\ & & & & & \ddots \\ & & & & & \tilde{A}_{n,n} \end{array} \right).$$

The scalars  $\tilde{d}_{m+1}, \dots, \tilde{d}_n$  do not necessarily exist for all symmetric positive definite (spd)  $A$  and general  $B$ , because the recurrence may break down at a zero pivot:

$$\begin{aligned} \tilde{d}_{m+1} &= A_{P(m+1),P(m+1)}, \\ \tilde{d}_i &= A_{P(i),P(i)} - \sum_{j=1}^{i-1} \frac{(A_{P(j),P(j)})^2}{\tilde{d}_j}, \quad m+2 \leq i \leq n. \end{aligned}$$

This is similar to the standard  $ILLU(0)$  preconditioner that is guaranteed to exist for  $M$ -matrices, but not for general spd matrices.

The diagonal  $2 \times 2$  blocks  $\tilde{D}_i$  for  $1 \leq i \leq m$  turn out not to be singular, and can even be proved to have a very special structure, as is shown in the following lemma.

**Lemma 3.1.** *There exist  $\tilde{d}_1, \dots, \tilde{d}_m$  such that, for  $1 \leq i \leq m$ ,*

$$\tilde{D}_i = \begin{pmatrix} \tilde{d}_i & B_{P(i),i} \\ B_{i,P(i)}^T & 0 \end{pmatrix}.$$

Proof:

The proof proceeds by induction. It is easily verified that

$$\tilde{D}_1 = \begin{pmatrix} A_{P(1),P(1)} & B_{P(1),1} \\ B_{1,P(1)}^T & 0 \end{pmatrix},$$

so that  $\tilde{d}_1 = A_{P(1),P(1)}$ . Now assume that  $\tilde{D}_1, \dots, \tilde{D}_{i-1}$  are of the desired form (where  $2 \leq i \leq m$ ). Then  $\tilde{D}_i$  is determined by the equation

$$\begin{pmatrix} A_{P(i),P(i)} & B_{P(i),i} \\ B_{i,P(i)}^T & 0 \end{pmatrix} = \tilde{D}_i + \sum_{j=1}^{i-1} \tilde{L}_{i,j} \tilde{D}_j^{-1} \tilde{L}_{i,j}^T.$$

By the induction hypothesis and the fact that  $B_{P(j),j}^2 = 1$  for all  $1 \leq j \leq m$ , we find that

$$\tilde{D}_j^{-1} = \begin{pmatrix} 0 & B_{P(j),j} \\ B_{j,P(j)}^T & -\tilde{d}_j \end{pmatrix}.$$

Hence,

$$\tilde{L}_{i,j} \tilde{D}_j^{-1} \tilde{L}_{i,j}^T = \begin{pmatrix} 2A_{P(i),P(j)} B_{P(j),j} B_{P(i),j} - \tilde{d}_j B_{P(i),j}^2 & B_{P(i),j} B_{P(j),j} B_{P(j),i} \\ B_{P(i),j} B_{P(j),j} B_{P(j),i} & 0 \end{pmatrix}.$$

Due to (3) we have that  $B_{P(j),i} = 0$ , and we conclude that

$$\tilde{L}_{i,j} \tilde{D}_j^{-1} \tilde{L}_{i,j}^T = \begin{pmatrix} 2A_{P(i),P(j)} B_{P(j),j} B_{P(i),j} - \tilde{d}_j B_{P(i),j}^2 & 0 \\ 0 & 0 \end{pmatrix}.$$

So,

$$\tilde{D}_i = \begin{pmatrix} \tilde{d}_i & B_{P(i),i} \\ B_{i,P(i)}^T & 0 \end{pmatrix},$$

with

$$\tilde{d}_i = A_{P(i),P(i)} + \sum_{j=1}^{i-1} B_{P(i),j}^2 \tilde{d}_j - 2A_{P(i),P(j)} B_{P(j),j} B_{P(i),j}.$$

Hence, the lemma is proved.

Note that there is at most one  $j \in \{1, \dots, n\}$ ,  $j \neq i$ , such that  $B_{P(i),j} \neq 0$ . Denote this number by  $j(i)$ . Then we have  $j(i) \leq i - 1$  and

$$\tilde{d}_i = A_{P(i),P(i)} + \tilde{d}_{j(i)} - 2A_{P(i),P(j(i))} B_{P(j(i)),j(i)} B_{P(i),j(i)}.$$

Lemma 1 tells us that the blocks in  $\tilde{D}$  are of the same structure as the  $2 \times 2$  blocks in the upper left part of  $\tilde{M}$ . Hence, the following corollary is not surprising.

**Corollary 3.2.** *Let  $\tilde{L}$  and  $\tilde{D}$  be determined as described in the foregoing and suppose that the scalars  $\tilde{d}_{m+1}, \dots, \tilde{d}_n$  defined by (4) exist. Then*

$$Q(\tilde{L} + \tilde{D})\tilde{D}^{-1}(\tilde{L} + \tilde{D})^T Q^T = \begin{pmatrix} G & B \\ B^T & 0 \end{pmatrix}.$$

for some matrix  $G$ .

Proof:

Define  $l(C)$  as the strictly lower triangular part of a matrix  $C$ , and  $P(i) = i$  for all  $i$ . Then we obtain

$$L = Q\tilde{L}Q^T = \begin{pmatrix} l(A_{22}) & 0 & l(B_1) \\ A_{21} & l(A_{22}) & B_2 \\ 0 & 0 & 0 \end{pmatrix},$$

$$D = Q\tilde{D}Q^T = \begin{pmatrix} D_1 & 0 & \text{diag}(B_1) \\ 0 & D_2 & 0 \\ \text{diag}(B_1) & 0 & 0 \end{pmatrix},$$

where

$$D_1 = \text{diag}(\tilde{d}_1, \dots, \tilde{d}_m),$$

$$D_2 = \text{diag}(\tilde{d}_{m+1}, \dots, \tilde{d}_n).$$

Multiplying out  $(L + D)D^{-1}(L + D)^T$  then gives the result.

The corollary demonstrates that the preconditioning matrix is in exactly the same form as suggested by Keller et al [17], i.e. it is a so-called constrained preconditioner. Even more importantly, the corollary shows that this preconditioner is obtained in factorized form. Thus, we have found a way to construct constraint preconditioners that are easily inverted. This observation has sparked much research into constraint preconditioners for saddle point problems.

It should also be noted that in [18] experiments with similar use of  $1 \times 1$  and  $2 \times 2$  pivots have been carried out for indefinite systems. In the aforementioned paper, an ILU decomposition for indefinite systems, based on a Crout decomposition, is employed. The paper contains many interesting numerical results.

#### 4. A general decomposition for indefinite matrices

The technique described in the previous section is based on properties of the matrix  $B$ . In fact, it was assumed that  $B$  is an incidence matrix with only few non-zero entries. Such matrices can be put into lower trapezoidal form, meaning that the top  $m \times m$  part is lower triangular. For more general  $B$ , a similar treatment is possible by making use of LQ decompositions. To this end, we write

$$\Pi B = \widehat{B}Q,$$

where  $\Pi$  is an  $n \times n$  permutation matrix,  $Q$  is an  $m \times m$  orthogonal matrix, and  $\widehat{B}$  is of lower trapezoidal form. Furthermore we require that the top  $m \times m$  part of  $\widehat{B}$  is nonsingular. Such decompositions are always possible, and many software routines are available. Actually, the matrix  $Q$  can be obtained as the product of a permutation matrix and a number of matrices describing Givens rotations.

Now define

$$Q = \begin{pmatrix} \Pi & 0 \\ 0 & Q \end{pmatrix},$$

and let

$$\widehat{A} = \Pi A \Pi^T.$$

Then

$$Q A Q^T = \begin{pmatrix} \widehat{A} & \widehat{B} \\ \widehat{B}^T & 0 \end{pmatrix}.$$

The matrix  $\widehat{B}$  is now of a form similar to that in Section 3, and the following holds (see also Theorem 4.2 in [11]):

**Lemma 4.1.** *Let  $\widehat{A}$  and  $\widehat{B}$  be as in the foregoing, and write  $\widehat{B}^T = (\widehat{B}_1, \widehat{B}_2)^T$  where  $\widehat{B}_1$  is the  $m \times m$  top part of  $\widehat{B}$ . Then there exist an  $m \times m$  diagonal matrix  $D_1$ , an  $(n - m) \times (n - m)$  diagonal matrix  $D_2$ , an  $m \times m$  strictly lower triangular matrix  $L_1$ , an  $(n - m) \times (n - m)$  strictly lower triangular matrix  $L_2$ , and an  $(n - m) \times m$  matrix  $M$ , such that*

$$\begin{pmatrix} \widehat{A} & \widehat{B} \\ \widehat{B}^T & 0 \end{pmatrix} = \begin{pmatrix} \widehat{B}_1 & 0 & L_1 \\ \widehat{B}_2 & I_{n-m} + L_2 & M \\ 0 & 0 & I_m \end{pmatrix} \begin{pmatrix} D_1 & 0 & I_m \\ 0 & D_2 & 0 \\ I_m & 0 & 0 \end{pmatrix} \begin{pmatrix} \widehat{B}_1^T & \widehat{B}_2^T & 0 \\ 0 & I_{n-m} + L_2^T & 0 \\ L_1^T & M^T & I_m \end{pmatrix} \quad (5)$$

Proof:

Working out the expression in the right hand side, and writing

$$\widehat{A} = \begin{pmatrix} \widehat{A}_{11} & \widehat{A}_{12} \\ \widehat{A}_{21} & \widehat{A}_{22} \end{pmatrix},$$

we find that the following relations must be satisfied:

$$\widehat{B}_1 D_1 \widehat{B}_1^T + \widehat{B}_1 L_1^T + L_1 \widehat{B}_1^T = \widehat{A}_{11}, \quad (6)$$

$$\widehat{B}_1 D_1 \widehat{B}_2^T + \widehat{B}_1 M^T + L_1 \widehat{B}_2^T = \widehat{A}_{12}, \quad (7)$$

$$\widehat{B}_2 D_1 \widehat{B}_1^T + \widehat{B}_2 L_1^T + M \widehat{B}_1^T = \widehat{A}_{21}, \quad (8)$$

$$(I_{n-m} + L_2) D_2 (I_{n-m} + L_2)^T + \widehat{B}_2 D_1 \widehat{B}_2^T + \widehat{B}_2 M^T + M \widehat{B}_2^T = \widehat{A}_{22}. \quad (9)$$

Multiplying equation (6) from the left by  $\widehat{B}_1^{-1}$  and from the right by  $\widehat{B}_1^{-T}$  yields

$$D_1 + L_1^T \widehat{B}_1^{-T} + \widehat{B}_1^{-1} L_1 = \widehat{B}_1^{-1} \widehat{A}_{11} \widehat{B}_1^{-T}.$$

Thus, the matrices  $D_1$ ,  $L_1$  can be found from the expressions:

$$D_1 = \text{diag}(\widehat{B}_1^{-1} \widehat{A}_{11} \widehat{B}_1^{-T}),$$

$$L_1 = \widehat{B}_1 \text{lower}(\widehat{B}_1^{-1} \widehat{A}_{11} \widehat{B}_1^{-T}).$$

Note that we have explicitly used the fact that  $\widehat{B}_1$  is lower triangular here!

Having found  $D_1$  and  $L_1$ , the matrix  $M$  is simply obtained from either (7) or (8), to give

$$M = (\widehat{A}_{21} - \widehat{B}_2 L_1^T) \widehat{B}_1^{-T} - \widehat{B}_2 D_1.$$

It remains to show that matrices  $L_2$  and  $D_2$  exist such that (9) is satisfied. To this end, we first observe that

$$\begin{aligned} \widehat{B}_2 M^T &= \widehat{B}_2 \widehat{B}_1^{-1} (\widehat{A}_{12} - L_1 \widehat{B}_2^T) - \widehat{B}_2 D_1 \widehat{B}_2^T, \\ M \widehat{B}_2^T &= (\widehat{A}_{21} - \widehat{B}_2 L_1^T) \widehat{B}_1^{-1} \widehat{B}_2^T - \widehat{B}_2 D_1 \widehat{B}_2^T, \end{aligned}$$

by virtue of (7) and (8). Substituting this in (9), and making use of the expressions for  $D_1$  and  $L_1$ , we find that the following must hold:

$$(I_{n-m} + L_2) D_2 (I_{n-m} + L_2)^T = \widehat{A}_{22} + \widehat{B}_2 \widehat{B}_1^{-1} \widehat{A}_{11} \widehat{B}_1^{-T} \widehat{B}_2^T - \widehat{B}_2 \widehat{B}_1^{-1} \widehat{A}_{12} - \widehat{A}_{21} \widehat{B}_1^{-T} \widehat{B}_2^T.$$

In other words.  $D_2$  and  $L_2$  are to be found from the expression

$$(I_{n-m} + L_2) D_2 (I_{n-m} + L_2)^T = \begin{pmatrix} -\widehat{B}_2 \widehat{B}_1^{-1} & I_{n-m} \end{pmatrix} \begin{pmatrix} \widehat{A}_{11} & \widehat{A}_{12} \\ \widehat{A}_{21} & \widehat{A}_{22} \end{pmatrix} \begin{pmatrix} -\widehat{B}_1^{-T} \widehat{B}_2^T \\ I_{n-m} \end{pmatrix}, \quad (10)$$

which is possible because  $\widehat{A}$  is a positive definite, symmetric matrix. This completes the proof.

A straightforward consequence of this lemma is the following decomposition theorem for indefinite matrices:

**Theorem 4.2.** *Let  $A$  be an  $n \times n$  symmetric, positive definite matrix,  $B$  an  $n \times m$  matrix of full rank,  $m \leq n$ , and set*

$$\mathcal{A} = \begin{pmatrix} A & B \\ B^T & 0 \end{pmatrix}.$$

*Then there exist an  $n \times n$  permutation matrix  $\Pi$ , an  $m \times m$  orthogonal matrix  $Q$ , an  $m \times m$  diagonal matrix  $D_1$ , an  $(n-m) \times (n-m)$  diagonal matrix  $D_2$ , an  $m \times m$  strictly lower triangular matrix  $L_1$ , an  $(n-m) \times (n-m)$  strictly lower triangular matrix  $L_2$ , and an  $(n-m) \times m$  matrix  $M$ , such that  $\Pi B Q^T$  is lower trapezoidal and*

$$\mathcal{A} = Q \mathcal{L} \mathcal{D} \mathcal{L}^T Q^T, \quad (11)$$



where

$$\mathcal{Q} = \begin{pmatrix} 0 & \Pi^T \\ Q^T & 0 \end{pmatrix},$$

$$\mathcal{L} = \begin{pmatrix} I_m & 0 & 0 \\ L_1 & \widehat{B}_1 & 0 \\ M & \widehat{B}_2 & I_{n-m} + L_2 \end{pmatrix},$$

$$\mathcal{D} = \begin{pmatrix} 0 & I_m & 0 \\ I_m & D_1 & 0 \\ 0 & 0 & D_2 \end{pmatrix},$$

where  $\widehat{B}_1$  is the top  $m \times m$  part of  $\Pi B Q^T$ , and  $\widehat{B}_2$  is the lower  $(n-m) \times m$  part of the same matrix.

Proof:

Using Lemma 1, a decomposition of the form (5) is found. With a simple permutation of rows and columns, we find

$$\begin{pmatrix} 0 & 0 & I_m \\ I_m & 0 & 0 \\ 0 & I_{n-m} & 0 \end{pmatrix} \begin{pmatrix} \widehat{B}_1 & 0 & L_1 \\ \widehat{B}_2 & I_{n-m} + L_2 & M \\ 0 & 0 & I_m \end{pmatrix} \begin{pmatrix} 0 & I_m & 0 \\ 0 & 0 & I_{n-m} \\ I_m & 0 & 0 \end{pmatrix} = \begin{pmatrix} \widehat{I}_m & 0 & 0 \\ L_1 & \widehat{B}_1 & 0 \\ M & \widehat{B}_2 & I_{n-m} + L_2 \end{pmatrix}.$$

The proof now follows from the observation that

$$\begin{pmatrix} 0 & 0 & I_m \\ I_m & 0 & 0 \\ 0 & I_{n-m} & 0 \end{pmatrix} \begin{pmatrix} D_1 & 0 & I_m \\ 0 & D_2 & 0 \\ I_m & 0 & 0 \end{pmatrix} \begin{pmatrix} 0 & I_m & 0 \\ 0 & 0 & I_{n-m} \\ I_m & 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & I_m & 0 \\ I_m & D_1 & 0 \\ 0 & 0 & D_2 \end{pmatrix},$$

and (take care of the dimensions of  $Q$ )

$$\begin{pmatrix} \Pi^T & 0 \\ 0 & Q^T \end{pmatrix} \begin{pmatrix} 0 & I_m & 0 \\ 0 & 0 & I_{n-m} \\ I_m & 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & \Pi^T \\ Q^T & 0 \end{pmatrix}.$$

*Remark 2.1* Note the resemblance of the decomposition in (11) with the Bunch-Kaufman-Parlett decomposition described in [5, 6, 15]<sup>1</sup>. The structure of the decomposition is similar, the difference being that the permutation

---

<sup>1</sup>Historical note: by pure coincidence, the famous paper by Bunch and Kaufman [6] follows immediately, in the same volume of *Mathematics of Computation*, the equally famous paper by Meijerink and Van der Vorst [20].

matrix in the BKP method is now a more general orthogonal matrix. Note also that the matrix  $\mathcal{L}$  is a lower triangular matrix.

The decomposition presented in Theorem 1 can be used for the direct solution of indefinite systems of the form (1). Roughly speaking, the algorithm entails the following steps:

1. determine orthogonal matrices  $\Pi$ ,  $Q$  which transform  $B$  into the lower trapezoidal matrix  $\widehat{B}$
2. transform the matrix  $A$  by forming  $\Pi A \Pi^T$
3. determine the matrices  $D_1$ ,  $L_1$ , and  $M$
4. perform a Cholesky decomposition of the matrix

$$\begin{pmatrix} -\widehat{B}_2 \widehat{B}_1^{-1} & I_{n-m} \end{pmatrix} \Pi A \Pi^T \begin{pmatrix} -\widehat{B}_1^{-T} \widehat{B}_2^T \\ I_{n-m} \end{pmatrix},$$

leading to the matrices  $D_2$  and  $L_2$ .

Fortunately, the transformation of  $A$  performed in step 2. involves a permutation matrix, so that the sparsity of  $\Pi A \Pi^T$  is the same as for  $A$ . If  $B$  is an incidence matrix, then we know that  $Q$  is a simple permutation too. Depending on the specific type of problem, it may be possible to construct the matrices  $\Pi$  and  $Q$  using only topological information, just as in the case discussed in Section 3. For indefinite systems obtained after discretisation of partial differential equations, the sparsity of the matrix  $B$  depends on the type of (finite) elements used. If higher order elements are used, there will be more non-zero elements in  $B$ . If the indefinite system describes an optimisation problem, the matrix  $B$  will be sparse since constraints usually couple only a few variables; problems in which all constraints contain all variables are not to be expected.

It is interesting to have a closer look at the block diagonal matrix  $\mathcal{D}$ , since this matrix contains essential information about the eigenvalues.

- the matrix  $D_2$  has  $n - m$  positive (real) eigenvalues
- the matrix

$$\begin{pmatrix} 0 & I_m \\ I_m & D_1 \end{pmatrix}$$

has  $m$  positive and  $m$  negative (real) eigenvalues

Hence, the indefiniteness of the original matrix  $\mathcal{A}$  is fully reflected in the matrix  $\mathcal{D}$ . The lower and upper triangular factors have unit eigenvalues, as is to be expected.

## 5. Preconditioning and iterative solution techniques

Although originally we set out to construct incomplete preconditioners for the indefinite systems occurring in electronic circuit simulation, the foregoing sections clearly show that, in fact, we have obtained a very general way of constructing exact decompositions of saddle point matrices. Hence, the decomposition in Theorem 4.2 can be used for a direct solution of the indefinite system (1).

However, the discussion in Sections 3 and 4 also leads to another, extremely interesting and valuable observation. This essential observation was made originally by Dollar, Gould and Wathen, and further elaborated in [8, 9, 10, 11, 12]. They noted that the factorization in Theorem 1 leads to a constraint preconditioner for all choices of the matrices  $D_1$ ,  $D_2$ ,  $L_1$ ,  $L_2$ , and  $M$ ! In other words, no matter what these matrices are, the resulting product  $QLDL^TQ^T$  will always be of the form (2).

Using the foregoing observation, it is rather simple to generate a wealth of constraint preconditioners, and the thesis [8] contains many families of these so-called implicit preconditioners. This terminology reflects the fact that, implicitly, always a constrained preconditioner is found, without having to explicitly calculate the matrices  $D_1$ ,  $D_2$ ,  $L_1$ ,  $L_2$ , and  $M$ . One could make choices for a number of these matrices, and calculate the remaining matrices explicitly. Or, alternatively, make specific choices for all of these matrices. The main question is, of course, how such preconditioners will perform in practice. Once again, this is nicely summarized in the aforementioned papers.

Hence, it is clear that the decomposition technique discussed in the previous sections can also be used as the basis for preconditioned iterative solution methods. Both in Section 4 and in [17] it has been demonstrated that it is a good idea to use preconditioners which retain the constraints, whence we restrict ourselves to preconditioning matrices of the form

$$\mathcal{G} = \begin{pmatrix} G & B \\ B^T & 0 \end{pmatrix}.$$

There are several criteria for preconditioners, an important one being that the matrix used for preconditioning is easily inverted. By virtue of Theorem 1, this is the case for  $\mathcal{G}$ , since we can write

$$\mathcal{G} = QL_GD_GL_G^TQ^T,$$

with

$$\mathcal{L}_G = \begin{pmatrix} \widehat{I}_m & 0 & 0 \\ L_{G,1} & \widehat{B}_1 & 0 \\ M_G & \widehat{B}_2 & I_{n-m} + L_{G,2} \end{pmatrix},$$

$$\mathcal{D} = \begin{pmatrix} 0 & I_m & 0 \\ I_m & D_{G,1} & 0 \\ 0 & 0 & D_{G,2} \end{pmatrix}.$$

Clearly, the matrix  $\mathcal{Q}$  is the same as in the previous section, since it only depends on the matrix  $B$ .

Motivated by the results in Section 4, the use of an incomplete factorisation is most appealing. This means that the matrices  $D_{G,1}$ ,  $L_{G,1}$ ,  $M_G$ ,  $L_{G,2}$ , and  $D_{G,2}$  should be approximations to the corresponding elements of the decomposition of  $A$ , which we shall denote by  $D_{A,1}$ ,  $L_{A,1}$ ,  $M_A$ ,  $L_{A,2}$ , and  $D_{A,2}$ , respectively. We observe that the calculation of the first three of these matrices is rather straightforward. Furthermore, working out the product  $(\mathcal{Q}\mathcal{L}_G\mathcal{D}_G\mathcal{L}_G^T\mathcal{Q}^T)^{-1}\mathcal{Q}\mathcal{L}_A\mathcal{D}_A\mathcal{L}_A^T\mathcal{Q}^T$ , we find that the product is a full matrix for which we can not easily find the eigenvalues. For these reasons, we shall assume the following:

$$\begin{aligned} D_{G,1} &= D_{A,1}, \\ L_{G,1} &= L_{A,1}, \\ M_G &= M_A. \end{aligned}$$

A straightforward calculation then shows that

$$\mathcal{Q}^T (\mathcal{Q}\mathcal{L}_G\mathcal{D}_G\mathcal{L}_G^T\mathcal{Q}^T)^{-1} \mathcal{Q}\mathcal{L}_A\mathcal{D}_A\mathcal{L}_A^T\mathcal{Q}^T \mathcal{Q} = \begin{pmatrix} I_m & 0 & X \\ 0 & I_m & Y \\ 0 & 0 & Z \end{pmatrix},$$

where  $X$ ,  $Y$  are not further specified, and

$$Z = (I_{n-m} + L_{G,2})^{-T} D_{G,2}^{-1} (I_{n-m} + L_{G,2})^{-1} (I_{n-m} + L_{A,2}) D_{A,2} (I_{n-m} + L_{A,2})^T.$$

This proves the following lemma.

**Lemma 5.1.** *Assume that  $D_{G,1} = D_{A,1}$ ,  $L_{G,1} = L_{A,1}$ , and  $M_G = M_A$ . Then the matrix  $\mathcal{G}^{-1}\mathcal{A}$  has  $2m$  eigenvalues 1, and the remaining  $n - m$  eigenvalues are equal to the eigenvalues of the generalized eigenvalue problem*

$$(I_{n-m} + L_{A,2}) D_{A,2} (I_{n-m} + L_{A,2})^T x = \lambda (I_{n-m} + L_{G,2}) D_{G,2} (I_{n-m} + L_{G,2})^T \quad (12)$$

We conclude that the problem of finding a suitable preconditioner for the indefinite problem is equivalent to finding a suitable preconditioner for linear systems involving the positive definite coefficient matrix

$$\begin{pmatrix} -\widehat{B}_2\widehat{B}_1^{-1} & I_{n-m} \end{pmatrix} \begin{pmatrix} \widehat{A}_{11} & \widehat{A}_{12} \\ \widehat{A}_{21} & \widehat{A}_{22} \end{pmatrix} \begin{pmatrix} -\widehat{B}_1^{-T}\widehat{B}_2^T \\ I_{n-m} \end{pmatrix}. \quad (13)$$

This is not surprising, as we can see from the following reasoning. Making use of the orthogonal transformation matrix  $Q$  in the  $LQ$  decomposition of  $B$ , we can write the system (1) as

$$\begin{pmatrix} \widehat{A} & \widehat{B} \\ \widehat{B}^T & 0 \end{pmatrix} \begin{pmatrix} \widehat{x} \\ \widehat{y} \end{pmatrix} = \begin{pmatrix} \widehat{b} \\ \widehat{c} \end{pmatrix},$$

with  $\widehat{x} = \Pi x$ ,  $\widehat{b} = \Pi b$ ,  $\widehat{y} = Qy$ , and  $\widehat{c} = Qc$ . Following the notation of Section 3, this is equivalent to the system

$$\begin{pmatrix} \widehat{A}_{11} & \widehat{A}_{12} & \widehat{B}_1 \\ \widehat{A}_{21} & \widehat{A}_{22} & \widehat{B}_2 \\ \widehat{B}_1^T & \widehat{B}_2^T & 0 \end{pmatrix} \begin{pmatrix} \widehat{x}_1 \\ \widehat{x}_2 \\ \widehat{y} \end{pmatrix} = \begin{pmatrix} \widehat{b}_1 \\ \widehat{b}_2 \\ \widehat{c} \end{pmatrix}.$$

Since  $B_1$  is nonsingular, both  $\widehat{x}_1$  and  $\widehat{y}$  can be eliminated, so that a reduced system is obtained in terms of the unknown  $\widehat{x}_2$ :

$$(\widehat{A}_{22} - \widehat{A}_{21}\widehat{B}_1^{-T}\widehat{B}_2^T - \widehat{B}_2\widehat{B}_1^{-1}\widehat{A}_{12} + \widehat{B}_2\widehat{B}_1^{-1}\widehat{A}_{11}\widehat{B}_1^{-T}\widehat{B}_2^T)\widehat{x}_2 = \widehat{b}_2 + (\widehat{B}_2\widehat{B}_1^{-1}\widehat{A}_{11} - \widehat{A}_{21})\widehat{B}_1^{-T}\widehat{c} - \widehat{B}_2\widehat{B}_1^{-1}\widehat{b}_1, \quad (14)$$

where the coefficient matrix is the same as that in (13). This completes the argument.

Because of the foregoing observation, the iterative solution of systems of the form (1) may be performed in the following form:

1. determine a permutation matrix  $\Pi$  and an orthogonal matrix  $Q$  which transform  $B$  into the lower trapezoidal matrix  $\widehat{B}$
2. transform the matrix  $A$  by forming  $\Pi A \Pi^T$
3. determine the matrices  $D_1$ ,  $L_1$ , and  $M$
4. perform an incomplete decomposition of the matrix

$$\begin{pmatrix} -\widehat{B}_2\widehat{B}_1^{-1} & I_{n-m} \end{pmatrix} \Pi A \Pi^T \begin{pmatrix} -\widehat{B}_1^{-T}\widehat{B}_2^T \\ I_{n-m} \end{pmatrix},$$

leading to the matrices  $D_{G,2}$  and  $L_{G,2}$

5. iteratively solve the system (14) using the incomplete preconditioning matrix obtained in the previous step
6. calculate the remaining components  $\hat{x}_1$  and  $\hat{y}$  of the solution vector
7. transform the solution vector back to the original variables using the orthogonal matrix  $Q$  and the permutation matrix  $\Pi$

Clearly, the simplest possible preconditioning is obtained when assuming that  $L_{G,2} \equiv 0$ . In that case, we require

$$D_{G,2} = \text{diag} \left( -\hat{B}_2 \hat{B}_1^{-1} \quad I_{n-m} \right) \begin{pmatrix} \hat{A}_{11} & \hat{A}_{12} \\ \hat{A}_{21} & \hat{A}_{22} \end{pmatrix} \begin{pmatrix} -\hat{B}_1^{-T} \hat{B}_2^T \\ I_{n-m} \end{pmatrix}.$$

Dollar [8] has performed extensive research on suitable choices for these implicit factorization preconditioners, and a wealth of numerical results is available, also in [9, 10, 11, 12]. The results clearly demonstrate the potential of constrained preconditioning.

It should be noted that, despite the fact that the preconditioned system is non-symmetric in general, it is possible to use the conjugate gradient method for their solution. This is possible if we assume that an 'inner product'

$$[x, y] \equiv x^T \mathcal{G} y$$

is used that is based upon the preconditioning matrix  $\mathcal{G}$ . Such point of view for preconditioned CG is clearly explained in [28]. Of course, if we choose the wrong starting vector for the CG process, we may immediately end up with a failing CG process. However, in practical cases, it has turned out to be a very useful and certainly feasible method. This is mainly due to the fact that the preconditioned system has eigenvalues that are all located in the right half plane. This is not surprising if we look at the form of the preconditioner, which is very similar to that of the original matrix. In fact, the preconditioner has been constructed in such a way that negative eigenvalues of the original matrix are 'compensated' by negative eigenvalues of the preconditioning matrix, in such a way that the product matrix has eigenvalues with positive real parts. This observation clearly demonstrates that structure preservation is essential.

## 6. Conclusion

In this paper, we have elaborated the ideas underlying the Schilders' factorization. It has been demonstrated that a Bunch-Kaufman-Parlett like

strategy can be employed with an a priori known structure of the pivots. The concept has been generalized, and has led to a decomposition method for symmetric indefinite matrices of a special form. The method can readily be extended to the non-symmetric case (this has been done in [7, 19]), and also for non-zero lower right hand blocks the ideas can be used to obtain factorizations. In addition to the exact decompositions, the method has also been used to generate implicit factorization preconditioners, and these have been shown to be very effective. For numerical results, the reader is referred to [8, 9, 10, 11, 12].

## 7. Acknowledgement

The author is grateful to Henk van der Vorst, who has been an advisor to our company for 20 years. Within that period, we have greatly benefitted from the extensive knowledge that Henk has, both via his in-depth knowledge of linear algebra, and his extensive network of colleagues. Henk, thank you very much for this extremely exciting, interesting, rewarding and fruitful period!

## References

- [1] M. Benzi, G.H. Golub and J. Liesen, Numerical solution of saddle point problems, *Acta Numerica*, vol. 14, pp. 1-137 (2005)
- [2] F. Brezzi, *On the existence, uniqueness and approximation of saddle-point problems arising from Lagrangian multipliers*, RAIRO Anal. Numer., vol. 8, pp. 129-151 (1974)
- [3] F. Brezzi and M. Fortin, *Mixed and hybrid finite element methods*, Springer Series In Computational Mathematics, vol. 15 (1991)
- [4] F. Brezzi, L.D. Marini and P. Pietra, *Two-dimensional exponential fitting and applications to semiconductor device equations*, Publ. no. 597, Consiglio Nazionale Delle Ricerche, Pavia, Italy (1987)
- [5] J.R. Bunch and B.N. Parlett, Direct methods for solving symmetric indefinite systems of linear equations, *SIAM J. Matrix Anal. Appl.*, vol. 8, pp. 639-655 (1971)

- [6] J.R. Bunch and L. Kaufman, Some stable methods for calculating inertia and solving symmetric linear systems, *Math. Comp.*, vol. 31, pp. 163-179 (1977)
- [7] Z.-H. Cao, A class of constraint preconditioners for nonsymmetric saddle point problems, *Numer. Math.*, vol. 103, pp. 47-61 (2006)
- [8] H.S. Dollar, *Iterative linear algebra for constrained optimization*, DPhil (PhD) thesis, Oxford University Computing Laboratory (2005)
- [9] H.S. Dollar, N.I.M. Gould and A.J. Wathen, On implicit-factorization constraint preconditioners, in: *Large scale nonlinear optimization*, G. di Pillo and M. Roma (eds.), Springer Verlag, Heidelberg, Berlin, New York, pp. 61-82 (2006)
- [10] H.S. Dollar, N.I.M. Gould, W.H.A. Schilders and A.J. Wathen, Implicit-factorization preconditioning and iterative solvers for regularized saddle-point systems, *SIAM J. Matrix Anal. Appl.*, vol. 28, pp. 170-189 (2006)
- [11] H.S. Dollar and A.J. Wathen, Approximate factorization constraint preconditioners for saddle-point matrices, *SIAM J. Sci. Comput.*, vol. 27, pp. 1555-1572 (2006)
- [12] H.S. Dollar, N.I.M. Gould, W.H.A. Schilders and A.J. Wathen, Using constraint preconditioners with regularized saddle-point problems, *Comput. Optim. Appl.*, vol. 36 (2-3), pp. 249-270 (2008)
- [13] S.C. Eisenstat, Efficient implementation of a class of preconditioned conjugate gradient methods, *SIAM J. Sci. Stat. Comp.*, vol. 2, pp. 1-4 (1981)
- [14] B.X. Fraeijs de Veubeke, Displacement and equilibrium models in the finite element method, in: *Stress Analysis*, O.C. Zienkiewicz and G. Hollister (eds.), John Wiley, New York (1965)
- [15] G.H. Golub and C.F. Van Loan, *Matrix computations*, Johns Hopkins Studies in the Mathematical Sciences, Johns Hopkins University Press, Baltimore, MD, third ed. (1996)
- [16] M. Günther, J. ter Maten and U. Feldmann, *Modeling and discretization of circuit problems*, in: Numerical Methods in Electromagnetics,



- W.H.A. Schilders and E.J.W. ter Maten (eds.), vol. XIII of Handbook of Numerical Analysis, Elsevier (2005)
- [17] C. Keller, N.I.M. Gould and A.J. Wathen, Constraint preconditioning for indefinite linear systems, *SIAM J. Matrix Anal. Appl.*, vol. 21, pp. 1300-1317 (2000)
- [18] N. Li and Y. Saad, Crout versions of ILU factorization with pivoting for sparse symmetric matrices, *Electr. Trans. Num. Anal.*, vol. 20, pp. 75-85 (2005)
- [19] Y. Lin and Y. Wei, A note on constraint preconditioners for nonsymmetric saddle point problems, *Numer. Lin. Alg. Appl.*, vol. 14, pp. 659-664 (2007)
- [20] J.A. Meijerink and H.A. van der Vorst, An iterative solution method for linear systems of which the coefficient matrix is a symmetric M-matrix, *Math. Comp.*, vol. 31, pp. 148-162 (1977)
- [21] S.J. Polak, W.H.A. Schilders and H.D. Couperus, A finite element method with current conservation, in: *Proceedings SISPAD 1988 Conference*, G. Baccarani and M. Rudan (eds.), pp. 453-462 (1988)
- [22] R.A. Raviart and J.M. Thomas, *A mixed finite element method for 2nd order elliptic problems*, in: *Mathematical Aspects of the Finite Element Method*, Lecture Notes in Mathematics, vol. 606, Springer-Verlag, New York, pp. 292-315 (1977)
- [23] T. Rusten and R. Winther, A preconditioned iterative method for saddlepoint problems, *SIAM J. Matrix Anal. Appl.*, vol. 13, pp. 887-904 (1992)
- [24] W.H.A. Schilders and H.A. van der Vorst, Preconditioning techniques for indefinite linear systems with applications to circuit simulation, in: *Proc. Int. Conf. on Preconditioning Techniques for Large Sparse Matrix Problems in Industrial Applications*, June 10-12, 1999, Minneapolis (1999)
- [25] W.H.A. Schilders, *A preconditioning technique for indefinite linear systems*, RANA report 18, Eindhoven University of Technology (2000)

- [26] W.H.A. Schilders, H.A. van der Vorst and J. Rommes, *Model order reduction: theory, research aspects and applications*, Mathematics in Industry series, vol. 13, Springer-Verlag, Berlin (2008)
- [27] H.A. van der Vorst, *Preconditioning by incomplete decompositions*, ACCU Report no. 32, Utrecht (1982)
- [28] H.A. van der Vorst, *Iterative Krylov methods for large linear systems*, Cambridge University Press, Cambridge, UK (2003)
- [29] A.J.H. Wachtters and W.H.A. Schilders, Simulation of EMC behaviour, in: *Numerical Methods in Electromagnetics*, W.H.A. Schilders and E.J.W. ter Maten (eds.), vol. XIII of Handbook of Numerical Analysis, Elsevier (2005)
- [30] A. Wathen, Preconditioning constrained systems, in: *Proc. Int. Conf. on Preconditioning Techniques for Large Sparse Matrix Problems in Industrial Applications*, June 10-12, 1999, Minneapolis (1999)