

**Report No. UIUCDCS-R-2006-2678**

**UILU-ENG-2006-1705**

**Large-Scale Topology Optimization using Preconditioned Krylov  
Subspace Methods with Recycling**

by

**Shun Wang, Eric de Sturler, and Glaucio H. Paulino**

January 2006

# Large-scale topology optimization using preconditioned Krylov subspace methods with recycling

Shun Wang<sup>1</sup>, Eric de Sturler<sup>1</sup>, Glaucio H. Paulino<sup>2</sup>

<sup>1</sup>*Department of Computer Science*

<sup>2</sup>*Department of Civil and Environmental Engineering  
University of Illinois at Urbana-Champaign  
Urbana IL, 61801*

## SUMMARY

The computational bottleneck of topology optimization is the solution of a large number of linear systems arising in the finite element analysis. We propose fast iterative solvers for large three-dimensional topology optimization problems to address this problem. Since the linear systems in the sequence of optimization steps change slowly from one step to the next, we can significantly reduce the number of iterations and the runtime of the linear solver by recycling selected search spaces from previous linear systems. In addition, we introduce a MINRES (Minimum Residual method) version with recycling (and a short term recurrence) to make recycling more efficient for symmetric problems. Furthermore, we discuss preconditioning to ensure fast convergence. We show that a proper rescaling of the linear systems reduces the huge condition numbers that typically occur in topology optimization to roughly those arising for a problem with constant density. We demonstrate the effectiveness of our solvers by solving a topology optimization problem with more than a million unknowns on a fast PC.

KEY WORDS: topology optimization, three-dimensional analysis, iterative methods, Krylov methods, Krylov subspace recycling, preconditioning, large-scale computation

## 1. INTRODUCTION

The goal of topology optimization is to find a material distribution in terms of design variables such that a given objective function, e.g., compliance, is minimized subject to certain constraints. We give a brief introduction to topology optimization in the next section. To make topology optimization a truly effective tool in the design of large structures and complex materials, we must be able to use large three-dimensional models. Most work on topology optimization for continuum structures has emphasized developing new formulations and applications, designing suitable elements, studying existence and uniqueness issues, and solving (modest size) problems. However, the computational aspect of large-scale topology

---

Contract/grant sponsor: This work was supported, in part, by NSF grant DMR-0325939 through the Materials Computation Center at UIUC.

optimization, specifically the high cost of solving many large linear systems, has not received much attention. This is the focus of this paper.

The finite element analysis step in topology optimization requires the solution of a long sequence of linear systems of the type

$$\mathbf{K}(\boldsymbol{\rho}^{(i)})\mathbf{u}^{(i)} = \mathbf{f}, \quad (1)$$

where  $\mathbf{K}$  is the stiffness matrix as a function of the density distribution  $\boldsymbol{\rho}$  at the  $i$ th optimization step,  $\mathbf{f}$  is the load vector, and  $\mathbf{u}$  is the displacement vector. Currently, direct solvers are most commonly used, because of the very large condition numbers arising in topology optimization. Unfortunately, direct solvers cannot effectively handle large 3D problems, since their large storage and computational requirements make them prohibitively expensive. Iterative solvers have low storage requirements and the computational cost per iteration is small. Therefore, as long as convergence is reasonably fast, we can solve very large problems.

Iterative solvers offer a number of additional advantages compared with direct solvers. First, we do not need to solve very accurately in the early phase of the topology optimization process. Second, iterative solvers are easy to parallelize, which is important for very large problems. For instance, parallelization of topology optimization was studied in [1, 2]. Third, iterative solvers can use solutions from previous systems as starting guesses, which leads to smaller initial residuals. Last, for a sequence of linear systems that change slowly, we can reduce the total number of iterations by recycling subspaces of earlier search spaces [3, 4].

In topology optimization, the change in the design variables becomes small after the first few optimization steps. Therefore, the change in the system matrix  $\mathbf{K}(\boldsymbol{\rho})$  from one optimization step to the next is also small, and the Krylov subspace recycling methods introduced in [3] are likely to be effective. We give more background on recycling in Section 3.

In most structural problems, the matrices are symmetric but not necessarily positive definite. For example, in vibration problems symmetric indefinite matrices arise [5]. For such matrices, MINRES (Minimum Residual method) [6] is the method of choice. Therefore, we focus on MINRES in the present paper. We extend the idea of subspace recycling to MINRES and make it more efficient by exploiting symmetry and short term recurrences. We discuss the recycling MINRES in detail in Section 4.

To achieve fast convergence we do need to use preconditioning. As the material distribution in a structure is being optimized, some elements become nearly void (we set a small positive lower bound on the densities to avoid singularity). This makes the linear system very ill-conditioned. First, we show that the ill-conditioning is largely a problem of poor scaling. We reduce the condition number by 6 orders of magnitude by rescaling the system matrices with two diagonal matrices. Since diagonal scaling does not introduce numerical errors, this also mitigates the serious potential accuracy problems of ill-conditioning. Next, we combine the rescaling with other preconditioners to make the condition numbers of the linear systems even smaller. We discuss preconditioning in Section 5.

In Section 6, we give some implementation details of our methods. In Section 7, we analyze our methods for a model problem, and we present the numerical results to demonstrate the significant improvements our solvers achieve. In the last section, we provide the conclusions.

## 2. TOPOLOGY OPTIMIZATION

Topology optimization is a powerful structural optimization method that combines a numerical solution method, usually the finite element method (FEM), with an optimization algorithm to find the optimal material distribution inside a given domain [7, 8, 9, 10, 11, 12]. In designing the topology of a structure we determine which points of space should be material and which points should be void (i.e. no material). However, it is well known that an optimum result of topology optimization consists in a structure with intermediate (or composite) material. So, continuous values between 0 and 1 replace the discrete 0/1 numbers to represent the relative densities of the elements, while some form of penalization is used to steer the solution back to discrete 0/1 values [13]. The objective function is the compliance and there is a volume constraint. This is the basic setup of a topology optimization problem. We specify the problem mathematically as follows.

$$\begin{aligned}
 \min_{\boldsymbol{\rho}, \mathbf{u}} \quad & c(\boldsymbol{\rho}, \mathbf{u}) = \mathbf{u}^T \mathbf{K}(\boldsymbol{\rho}) \mathbf{u} & (2) \\
 \text{s.t.} \quad & \mathbf{K}(\boldsymbol{\rho}) \mathbf{u} = \mathbf{f} \\
 & 0 < \rho_0 \leq \rho_e \leq 1 \quad e = 1, 2, \dots, n_e \\
 & \int_{\Omega} \rho d\Omega \leq V,
 \end{aligned}$$

where  $c$  is the compliance,  $\mathbf{K}(\boldsymbol{\rho})$  is the stiffness matrix as a function of the density distribution  $\boldsymbol{\rho}$ ,  $\mathbf{u}$  and  $\mathbf{f}$  are the displacement vector and load vector,  $\rho_0$  is a chosen, small, positive lower bound for the density to avoid singularity of the stiffness matrix, and  $V$  is the total volume in use. The Solid Isotropic Material with Penalization (SIMP) method [13, 14] uses one design variable to represent the density in each element, while the recent method of Continuous Approximation of Material Distribution (CAMD) [15, 16] uses multiple variables per element. Since the focus of this paper is the linear solver in the finite element analysis (FEA), we use the SIMP method as a simple setup.

The basic scheme of topology optimization is described in Figure 1. First, we set up the geometry and the loading, and initialize the density distribution  $\boldsymbol{\rho}$ . Then, we start the optimization loop. We need a linear solver to solve the equilibrium equations  $\mathbf{K} \mathbf{u} = \mathbf{f}$  in the finite element analysis. In the sensitivity analysis, we compute the derivatives of the objective function  $\partial c / \partial \rho_e$ . After this, we can apply an optional low-pass filter to remedy the checkerboard problem [17]. The next step is the kernel of the optimization. There are various optimization algorithms that can be used for topology optimization. For instance, Optimality Criteria (OC) is a simple approach based on a set of intuitive criteria [18], while the Method of Moving Asymptotes (MMA) is a mathematical programming algorithm which is more robust and well established in theory [19]. Since this paper deals mainly with the FEA in topology optimization, the choice of the optimization method is less relevant for our discussion. Therefore, we choose the OC method for its simplicity. However, our Krylov subspace recycling method and preconditioning techniques are general and can be used with other optimization methods.

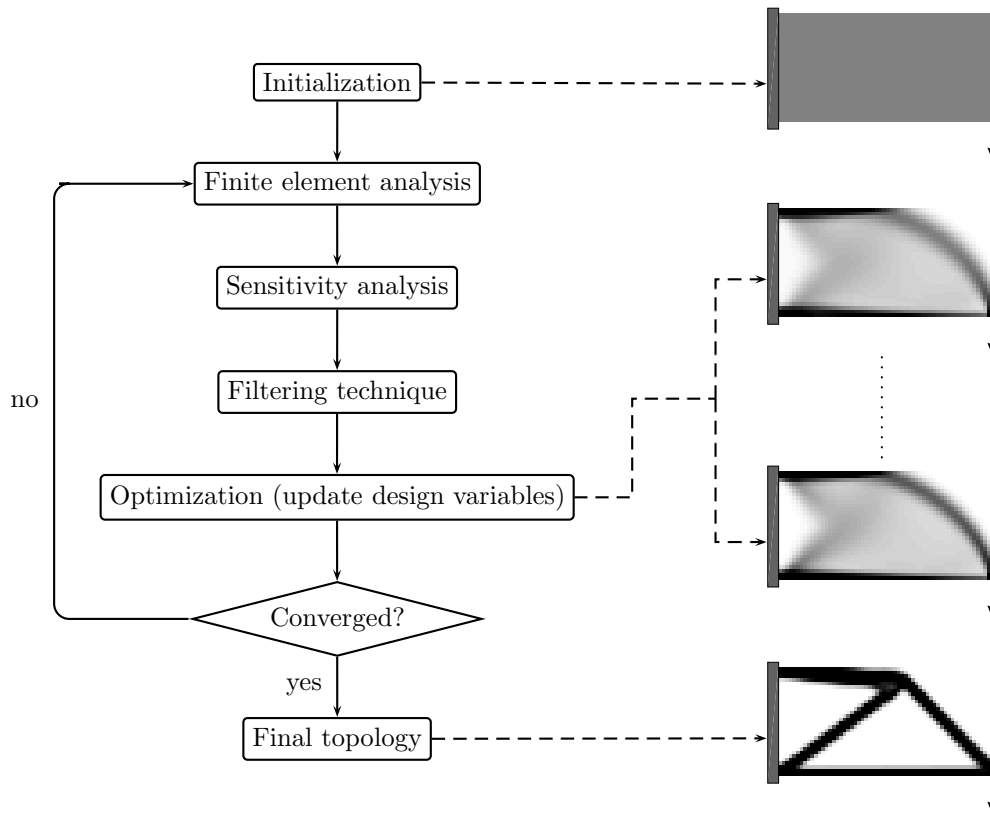


Figure 1. The general flow of computations for topology design.

### 3. KRYLOV SUBSPACE RECYCLING

Consider the linear system  $\mathbf{K}\mathbf{u} = \mathbf{f}$  and an initial guess  $\mathbf{u}_0$ . A Krylov subspace method, such as the Generalized Minimum Residual method (GMRES) [20], builds the Krylov subspace,  $\text{span}\{\mathbf{r}_0, \mathbf{K}\mathbf{r}_0, \mathbf{K}^2\mathbf{r}_0, \dots, \mathbf{K}^{m-1}\mathbf{r}_0\}$ , where  $\mathbf{r}_0 = \mathbf{f} - \mathbf{K}\mathbf{u}_0$ , and computes the optimal solution over that subspace. We use the Arnoldi recurrence [21] to obtain an orthonormal basis of the Krylov subspace:

$$\begin{aligned} \mathbf{v}_1 &= \mathbf{r}_0 / \|\mathbf{r}_0\|_2, \\ h_{i,i+1}\mathbf{v}_{i+1} &= \mathbf{K}\mathbf{v}_i - h_{i,i}\mathbf{v}_i - h_{i-1,i}\mathbf{v}_{i-1} - \dots - h_{1,i}\mathbf{v}_1, \end{aligned} \quad (3)$$

which in matrix form is written as

$$\mathbf{K}\mathbf{V}_m = \mathbf{V}_{m+1}\underline{\mathbf{H}}_m, \quad (4)$$

where the columns of  $\mathbf{V}_m$  are  $\mathbf{v}_1, \dots, \mathbf{v}_m$ ; the columns of  $\mathbf{V}_{m+1}$  are  $\mathbf{v}_1, \dots, \mathbf{v}_{m+1}$ ; and  $\underline{\mathbf{H}}_m$  is an  $(m+1) \times m$  upper Hessenberg matrix with coefficients  $\{h_{ij}\}$ .

For a symmetric matrix, we have  $\mathbf{K} = \mathbf{S}\mathbf{\Lambda}\mathbf{S}^T$ , where  $\mathbf{S}$  is an orthogonal matrix, and  $\mathbf{\Lambda}$  is a diagonal matrix whose coefficients are the eigenvalues of  $\mathbf{K}$ . The convergence rate of GMRES for a symmetric problem is bounded by [22, p. 206]

$$\frac{\|\mathbf{r}_m\|_2}{\|\mathbf{r}_0\|_2} \leq \min_{p_m \in \Pi_m^{(0)}} \max_{\lambda \in \lambda(\mathbf{K})} |p_m(\lambda)|, \quad (5)$$

where  $\Pi_m^{(0)}$  is the set of polynomials  $p_m$  of degree  $m$  such that  $p_m(0) = 1$ , and  $\lambda(\mathbf{K})$  is the set of eigenvalues of  $\mathbf{K}$ . So, the bound depends on the spectrum of the matrix. Therefore, if we remove an appropriate subset of the eigenvalues,  $M$ ,

$$\min_{p_m \in \Pi_m^{(0)}} \max_{\lambda \in \lambda(\mathbf{K})/M} |p_m(\lambda)|$$

can be significantly smaller than (5), and then the rate of convergence will be greatly improved. This is the motivation for recycling approximate invariant subspaces; other subspaces of the Krylov space may also be effective as a recycle space [4, 3]. Given the normalization condition,  $p_m(0) = 1$ , it is often effective to remove the eigenvalues close to the origin. This filtering of eigenvalues is achieved by including the corresponding (approximate) invariant subspace in the Krylov subspace over which we minimize. We typically recycle harmonic Ritz vectors with respect to the Krylov subspace to approximate an invariant subspace [3].

When solving the next linear system,  $\mathbf{K}(\boldsymbol{\rho}^{(i+1)})\mathbf{u}^{(i+1)} = \mathbf{f}$ , we include the recycle space as follows. We choose a basis for the recycle space to be the columns of a matrix  $\mathbf{U}$ , such that  $\mathbf{C} = \mathbf{K}\mathbf{U}$  and  $\mathbf{C}^T\mathbf{C} = \mathbf{I}$ . In addition, we adapt the Arnoldi process to make each new Krylov vector  $\mathbf{v}$  orthogonal to range( $\mathbf{C}$ ). This leads to the following recurrence:

$$\begin{aligned} (\mathbf{I} - \mathbf{C}\mathbf{C}^T)\mathbf{K}\mathbf{V}_m &= \mathbf{V}_{m+1}\underline{\mathbf{H}}_m \Leftrightarrow \\ \mathbf{K}\mathbf{V}_m &= \mathbf{C}\mathbf{C}^T\mathbf{K}\mathbf{V}_m + \mathbf{V}_{m+1}\underline{\mathbf{H}}_m, \end{aligned} \quad (6)$$

where  $\underline{\mathbf{H}}_m$  is still an  $(m+1) \times m$  upper Hessenberg matrix. Next, we compute the vector  $\boldsymbol{\varepsilon}_m = \mathbf{U}\mathbf{z}_m + \mathbf{V}_m\mathbf{y}_m$ , such that  $\mathbf{u}_m = \mathbf{u}_0 + \boldsymbol{\varepsilon}_m$  minimizes  $\|\mathbf{r}_m\|_2$ . This gives

$$\begin{aligned} \|\mathbf{r}_m\|_2 &= \left\| \mathbf{r}_0 - \mathbf{K}[\mathbf{U} \ \mathbf{V}_m] \begin{pmatrix} \mathbf{z}_m \\ \mathbf{y}_m \end{pmatrix} \right\|_2 \\ &= \left\| [\mathbf{C} \ \mathbf{V}_{m+1}] \left( \begin{pmatrix} \mathbf{C}^T\mathbf{r}_0 \\ \beta\mathbf{e}_1 \end{pmatrix} - \begin{bmatrix} \mathbf{I} & \mathbf{B}_m \\ \mathbf{0} & \underline{\mathbf{H}}_m \end{bmatrix} \begin{pmatrix} \mathbf{z}_m \\ \mathbf{y}_m \end{pmatrix} \right) \right\|_2 \\ &= \left\| \begin{pmatrix} \mathbf{C}^T\mathbf{r}_0 \\ \beta\mathbf{e}_1 \end{pmatrix} - \begin{bmatrix} \mathbf{I} & \mathbf{B}_m \\ 0 & \underline{\mathbf{H}}_m \end{bmatrix} \begin{pmatrix} \mathbf{z}_m \\ \mathbf{y}_m \end{pmatrix} \right\|_2, \end{aligned} \quad (7)$$

where  $\beta = \|(\mathbf{I} - \mathbf{C}\mathbf{C}^T)\mathbf{r}_0\|_2$  and  $\mathbf{B}_m = \mathbf{C}^T\mathbf{K}\mathbf{V}_m$ . This least square problem can be solved using the QR decomposition of  $\underline{\mathbf{H}}_m$ . This approach derives from the GCRO method [23] and is also used in the GCRODR method and GCROT with recycling [3, 4, 24].

An important issue for GMRES is that it relies (for general matrices) on a complete orthogonalization of the Krylov subspace. Therefore, as the Krylov subspace expands, the memory needed for the orthogonal basis vectors and the computational time for orthogonalization increase. As a result, normally restarting is required for GMRES, and we call the solution steps between two restarts a *cycle*. To mitigate the reduced convergence rate due to the loss of orthogonality caused by restarting, we use the recycle space immediately in the next cycle for the same system.

#### 4. RECYCLING MINRES

In topology optimization of structures, the system matrices are always symmetric. In most cases they are also positive definite. However, for some applications, e.g., topology design with dynamic vibrations, they can be indefinite [5]. So, in general, the MINRES method [6] is the most suitable iterative solver for topology optimization problems.

Both MINRES and GMRES minimize the two-norm of the residual over the Krylov subspace. The difference is that MINRES utilizes the symmetry of the matrix, and the resulting Lanczos three-term recurrence leads to significant reductions in memory requirements and computational cost.

We can use the matrices  $\mathbf{U}$  and  $\mathbf{C}$  defining the recycle space, obtained from solving previous linear systems, in the same way as in GCRODR. This leads to the same recurrence as in (6). However, the symmetry of  $\mathbf{K}$  implies the symmetry of  $H_m$ , the leading  $m \times m$  submatrix of  $\underline{\mathbf{H}}_m$ . Since  $\underline{\mathbf{H}}_m$  is also an upper Hessenberg matrix, this gives a tridiagonal  $\underline{\mathbf{H}}_m$ , which we will denote as  $\underline{\mathbf{T}}_m$  from now on. So, including the recycle space into the Krylov subspace does not affect the Lanczos recurrence of MINRES. Moreover, because of the three-term recurrence, we do not need to restart. So, in exact arithmetic, there is no need to use the recycle space generated during the solution of a linear system for solving that same system\*. As a consequence, we can derive a more efficient method for recycling for symmetric matrices. Although the Lanczos recurrence requires only the latest two basis vectors from the Krylov subspace (Lanczos vectors) for orthogonalization and no restarting is necessary, we need to retain the Lanczos vectors to be able to select a recycle space. To limit the memory requirements, we update the selected recycle space periodically. In this case, a *cycle* refers to the solution process between two updates of the recycle space.

We use  $s$  to denote the maximum length of a cycle (and hence the maximum number of Lanczos vectors kept), and  $k$  to denote the number of linearly independent vectors selected for recycling. We use  $\text{RMINRES}(s, k)$  to indicate the recycling MINRES method with the parameters  $s$  and  $k$ . The matrix  $\mathbf{V}_j$  contains the Lanczos vectors generated in the  $j$ th cycle,  $\mathbf{V}_j = [\mathbf{v}_{(j-1)s+1}, \dots, \mathbf{v}_{js}]$ , and the matrix  $\overline{\mathbf{V}}_j = [\mathbf{v}_{(j-1)s}, \dots, \mathbf{v}_{js+1}]$  denotes  $\mathbf{V}_j$  extended with with one previous and one subsequence Lanczos vector. Then, for the  $j$ th cycle, the modified Lanczos process gives

$$(\mathbf{I} - \mathbf{C}\mathbf{C}^T)\mathbf{K}\mathbf{V}_j = \overline{\mathbf{V}}_j\overline{\mathbf{T}}_j,^\dagger \quad (8)$$

where  $\overline{\mathbf{T}}_j$  is the tridiagonal matrix  $\underline{\mathbf{T}}_j$  with an additional row corresponding to  $\mathbf{v}_{(j-1)s}$  at the top. The bottom row corresponds to  $\mathbf{v}_{js+1}$ . To be specific,  $\overline{\mathbf{T}}_j$  has the nonzero pattern shown in Figure 2.

Let  $\mathbf{U}_j$  give the basis of a subspace that is selected at the end of the  $j$ th cycle for the current linear system and used only in the solution of the next system. We compute this matrix from the whole search space available in memory at the end of the  $j$ th cycle. This includes the matrix  $\mathbf{U}$  that is recycled from the previous system and used in the current system, the matrix  $\mathbf{U}_{j-1}$

---

\*In floating point arithmetic, including the recycle space obtained from the current Krylov subspace may help remedy the loss of orthogonality that generally occurs due to rounding errors.

<sup>†</sup>Note that  $(\mathbf{I} - \mathbf{C}\mathbf{C}^T)\mathbf{K}\mathbf{x} = (\mathbf{I} - \mathbf{C}\mathbf{C}^T)\mathbf{K}(\mathbf{I} - \mathbf{C}\mathbf{C}^T)\mathbf{x}$  for  $\mathbf{x} \in \text{range}(\mathbf{C})^\perp$ , so that  $(\mathbf{I} - \mathbf{C}\mathbf{C}^T)\mathbf{K}$  is *symmetric* over  $\text{range}(\mathbf{C})^\perp$ .



After solving (13), we choose the  $k$  harmonic Ritz vectors with the (absolute) smallest harmonic Ritz values for recycling, and set  $\tilde{U}_j = \mathbf{W}_j \mathbf{P}_j$ , where the columns of  $\mathbf{P}_j$  are the chosen harmonic Ritz vectors. Now we have  $\tilde{C}_j = \mathbf{K} \tilde{U}_j = \tilde{\mathbf{W}}_j \tilde{\mathbf{H}}_j \mathbf{P}_j$ . To obtain  $\mathbf{C}_j$  with orthonormal columns, we compute the QR decomposition of  $\tilde{\mathbf{W}}_j$  (note that by construction almost all columns are already orthogonal),

$$\tilde{\mathbf{W}}_j = \hat{\mathbf{W}}_j \mathbf{G}_j, \quad (14)$$

and of  $\mathbf{G}_j \tilde{\mathbf{H}}_j \mathbf{P}_j$ ,

$$\mathbf{G}_j \tilde{\mathbf{H}}_j \mathbf{P}_j = \mathbf{Q}_j \mathbf{R}_j. \quad (15)$$

Next, we set

$$\mathbf{U}_j = \mathbf{W}_j \hat{\mathbf{P}}_j, \quad \mathbf{C}_j = \hat{\mathbf{W}}_j \mathbf{Q}_j = \tilde{\mathbf{W}}_j \hat{\mathbf{Q}}_j, \quad (16)$$

where  $\hat{\mathbf{P}}_j = \mathbf{P}_j \mathbf{R}_j^{-1}$ , and  $\hat{\mathbf{Q}}_j = \mathbf{G}_j^{-1} \mathbf{Q}_j$ . Then  $\mathbf{C}_j$  is orthogonal and  $\mathbf{K} \mathbf{U}_j = \mathbf{C}_j$ . The two QR decompositions (14–15) are cheap to compute, because  $\mathbf{G}_j$  has very few nonzeros whose positions are known in advance, and  $\mathbf{G}_j \tilde{\mathbf{H}}_j \mathbf{P}_j$  is a product of matrices of small dimensions.

Finally, to solve the generalized eigenvalue problem (13), we need the matrices  $\tilde{\mathbf{H}}_j^T \tilde{\mathbf{W}}_j^T \tilde{\mathbf{W}}_j \tilde{\mathbf{H}}_j$  and  $\tilde{\mathbf{H}}_j^T \tilde{\mathbf{W}}_j^T \mathbf{W}_j$ . We can simplify  $\tilde{\mathbf{W}}_j^T \tilde{\mathbf{W}}_j$  and  $\tilde{\mathbf{W}}_j^T \mathbf{W}_j$  as follows.

$$\tilde{\mathbf{W}}_j^T \tilde{\mathbf{W}}_j = \begin{bmatrix} \mathbf{I} & \mathbf{C}^T \mathbf{C}_{j-1} & \mathbf{0} \\ \mathbf{C}_{j-1}^T \mathbf{C} & \mathbf{I} & \mathbf{C}_{j-1}^T \bar{\mathbf{V}}_j \\ \mathbf{0} & \bar{\mathbf{V}}_j^T \mathbf{C}_{j-1} & \mathbf{I} \end{bmatrix}, \quad (17)$$

$$\tilde{\mathbf{W}}_j^T \mathbf{W}_j = \begin{bmatrix} \mathbf{C}^T \mathbf{U}_{j-1} & \mathbf{0} \\ \mathbf{C}_{j-1}^T \mathbf{U}_{j-1} & \mathbf{C}_{j-1}^T \mathbf{V}_j \\ \bar{\mathbf{V}}_j^T \mathbf{U}_{j-1} & \bar{\mathbf{I}} \end{bmatrix}, \quad (18)$$

where  $\bar{\mathbf{I}}$  is an extended identity matrix with an additional row of zeros at the top and at the bottom. We can simplify the computation of most blocks in these two matrices further.

$$\mathbf{C}^T \mathbf{C}_{j-1} = \mathbf{C}^T \tilde{\mathbf{W}}_{j-1} \hat{\mathbf{Q}}_{j-1} = [\mathbf{I} \ \mathbf{C}^T \mathbf{C}_{j-2} \ \mathbf{0}] \hat{\mathbf{Q}}_{j-1}, \quad (19)$$

$$\bar{\mathbf{V}}_j^T \mathbf{C}_{j-1} = \bar{\mathbf{V}}_j^T \tilde{\mathbf{W}}_{j-1} \hat{\mathbf{Q}}_{j-1} = \begin{bmatrix} 0 & \cdots & 0 & 1 & 0 \\ 0 & \cdots & 0 & 0 & 1 \\ \vdots & \ddots & \ddots & \vdots & \vdots \\ 0 & \cdots & \cdots & 0 & 0 \end{bmatrix} \hat{\mathbf{Q}}_{j-1}, \quad (20)$$

$$\mathbf{C}^T \mathbf{U}_{j-1} = \mathbf{C}^T \mathbf{W}_{j-1} \hat{\mathbf{P}}_{j-1} = [\mathbf{C}^T \mathbf{U}_{j-2} \ \mathbf{0}] \hat{\mathbf{P}}_{j-1}, \quad (21)$$

$$\mathbf{C}_{j-1}^T \mathbf{U}_{j-1} = \hat{\mathbf{Q}}_{j-1}^T (\tilde{\mathbf{W}}_{j-1}^T \mathbf{W}_{j-1}) \hat{\mathbf{P}}_{j-1}, \quad (22)$$

$$\mathbf{C}_{j-1}^T \mathbf{V}_j = \hat{\mathbf{Q}}_{j-1}^T \begin{bmatrix} \mathbf{C}^T \\ \mathbf{C}_{j-2}^T \\ \bar{\mathbf{V}}_{j-1} \end{bmatrix} \mathbf{V}_j = \hat{\mathbf{Q}}_{j-1}^T \begin{bmatrix} 0 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \cdots & 0 \end{bmatrix}. \quad (23)$$

From the above derivation, we can obtain  $\bar{\mathbf{V}}_j^T \mathbf{C}_{j-1}$  and  $\mathbf{C}_{j-1}^T \mathbf{V}_j$  simply from  $\hat{\mathbf{Q}}_{j-1}$  (known from the previous cycle), and we can compute  $\mathbf{C}^T \mathbf{C}_{j-1}$ ,  $\mathbf{C}^T \mathbf{U}_{j-1}$  and  $\mathbf{C}_{j-1}^T \mathbf{U}_{j-1}$  recursively

with  $2k^3$ ,  $2k^3$  and  $2k(k+s)(3k+s)$  flops, respectively. Therefore, the computation of these submatrices is very cheap. Only  $\overline{\mathbf{V}}_j^T \mathbf{U}_{j-1}$  must be computed explicitly by matrix-matrix product, which takes about  $2ksn$  flops. In summary, the cost of each update of the recycle space is about  $(12k^2 + 6ks + 6k + 4)n$  flops, ignoring the terms that do not have a factor  $n$ . Compared with the cost of MINRES, which is mainly determined by the matrix-vector product and the forward and backward solve for the preconditioner for each iteration, the overhead of the subspace selection is modest (see the timing results in Section 7).

## 5. PRECONDITIONING FOR TOPOLOGY OPTIMIZATION

The convergence rate of Krylov methods for a symmetric matrix depends only on the spectrum of the matrix. In fact, the ratio between the absolute largest and smallest eigenvalue governs a worst-case upper bound on the convergence rate. In large-scale finite element simulations in physics and engineering, the linear systems tend to be ill-conditioned. In topology optimization, this problem is exacerbated by the wide range of magnitudes of the element densities.

Ill-conditioning creates two problems for numerical simulation. First, ill-conditioning may seriously affect the accuracy of the computed solution. Second, the convergence of iterative methods is poor for ill-conditioned problems. The second problem is generally addressed by proper preconditioning. In principle, preconditioning does not alleviate the potential accuracy problem, because a preconditioner that is effective for an ill-conditioned matrix has to be fairly ill-conditioned itself. This leads to two multiplications by ill-conditioned matrices in each iteration (or three for two-sided preconditioning), which may lead in turn to serious accumulation of numerical errors. However, in certain cases the accuracy problem can be relieved by properly scaling the problem. We show that this is the case for topology optimization. This leads to a preprocessing step and a preconditioning step (or two preconditioners depending on one's view).

In the next section, we discuss the preprocessing and preconditioner that we used for our numerical experiments. We illustrate the idea of rescaling from a mechanical point of view for a 1D problem in Section 5.2. Borrvall and Petersson [1] suggested, without further discussion, that the condition number of the matrix can be as large as the ratio of maximum to minimum density. We show that this ratio provides only a lower bound on the condition number and that the actual condition number typically is even larger. The actual conditioning is a combination of this ratio and the conditioning of a corresponding problem with constant density.

### 5.1. Preconditioning

The following analysis shows how ill-conditioned the stiffness matrices can be.

The two-norm condition number of a matrix  $\mathbf{K}$  can be defined as

$$\kappa(\mathbf{K}) = \frac{\max_{\|\mathbf{u}\|=1} \|\mathbf{K}\mathbf{u}\|}{\min_{\|\mathbf{u}\|=1} \|\mathbf{K}\mathbf{u}\|}. \quad (24)$$

Since

$$\min_{\|\mathbf{u}\|=1} \|\mathbf{K}\mathbf{u}\| \leq \|\mathbf{K}\mathbf{e}_l\| = \|\mathbf{k}_l\| \leq \max_{\|\mathbf{u}\|=1} \|\mathbf{K}\mathbf{u}\|, \quad \text{for any } l = 1, \dots, n, \quad (25)$$

where  $\mathbf{k}_l$  is the  $l$ th column of  $\mathbf{K}$  and  $\mathbf{e}_l$  is the Cartesian basis vector with the  $l$ th coefficient

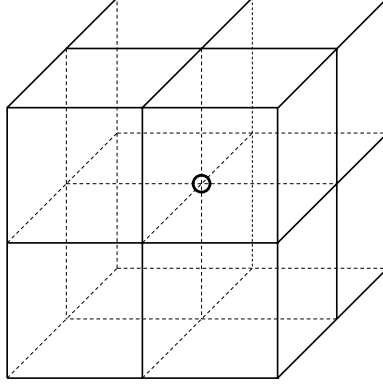


Figure 3.  $N_l$ : the set of elements associated with the  $l$ th d.o.f. indicated by the circle in the middle.

equal to 1, we have

$$\kappa(\mathbf{K}) \geq \frac{\|\mathbf{k}_{l_1}\|}{\|\mathbf{k}_{l_2}\|}, \quad \text{for any } l_1, l_2 = 1, \dots, n. \quad (26)$$

In topology optimization, a column of the stiffness matrix is given by

$$\mathbf{k}_l = \sum_{e \in N_l} \rho_e^p \mathbf{L}_e^T \mathbf{K}_0 \mathbf{L}_e \mathbf{e}_l, \quad (27)$$

where  $\mathbf{K}_0$  is the unit element stiffness matrix,  $\mathbf{L}_e$  is the local-to-global transformation matrix, and  $N_l$  is the set of elements that are associated with the  $l$ th d.o.f.. These usually form a  $2 \times 2 \times 2$  block in the 3D mesh (see Figure 3). If the blocks associated with d.o.f.  $l_1$  and  $l_2$  are solid and void respectively, namely  $\rho_e = 1$  for  $e \in N_{l_1}$  and  $\rho_e = \rho_0$  for  $e \in N_{l_2}$ , we have

$$\mathbf{k}_{l_1} = \sum_{e \in N_{l_1}} \mathbf{L}_e^T \mathbf{K}_0 \mathbf{L}_e \mathbf{e}_{l_1}, \quad (28)$$

$$\mathbf{k}_{l_2} = \rho_0^p \sum_{e \in N_{l_2}} \mathbf{L}_e^T \mathbf{K}_0 \mathbf{L}_e \mathbf{e}_{l_2}. \quad (29)$$

Then, assuming that the elements are uniform and isotropic, we have

$$\kappa(\mathbf{K}) \geq \frac{\|\mathbf{k}_{l_1}\|}{\|\mathbf{k}_{l_2}\|} = \frac{1}{\rho_0^p}. \quad (30)$$

For  $\rho_0 = 10^{-3}$  and  $p = 3$ , which are commonly used in topology optimization, the condition number of the stiffness matrix will be greater than  $10^9$  when solid and void areas begin to appear in the design domain.

Note that this analysis provides only a lower bound on the condition number, and that structures from homogeneous material can also have large condition numbers. However, the analysis suggests that, to a significant degree, the ill-conditioning comes from the poor scaling of the material densities over the design domain. We can understand this intuitively as follows.

A change in an algebraic degree of freedom, say the Cartesian basis vector  $\mathbf{e}_j$ , associated with a nodal basis function in a region with very small density corresponds to a displacement that requires a very small amount of energy ( $\mathbf{e}_j^T \mathbf{K}(\boldsymbol{\rho}) \mathbf{e}_j$  small). However, that same change in an algebraic degree of freedom,  $\mathbf{e}_i$ , associated with a nodal basis function in a region with large density corresponds to a displacement of the same magnitude that requires a large amount of energy ( $\mathbf{e}_i^T \mathbf{K}(\boldsymbol{\rho}) \mathbf{e}_i$  large). Since for symmetric  $\mathbf{K}$

$$\kappa(\mathbf{K}(\boldsymbol{\rho})) \geq \frac{\mathbf{e}_i^T \mathbf{K}(\boldsymbol{\rho}) \mathbf{e}_i}{\mathbf{e}_j^T \mathbf{K}(\boldsymbol{\rho}) \mathbf{e}_j},$$

this shows that the system is inherently ill-conditioned. Therefore, we expect that we can reduce the ill-conditioning due to the large variation in density by scaling the linear system such that changes of equal magnitude in algebraic degrees of freedom yield equal changes in energy ( $\mathbf{e}_i^T \mathbf{K}(\boldsymbol{\rho}) \mathbf{e}_i = \mathbf{e}_j^T \mathbf{K}(\boldsymbol{\rho}) \mathbf{e}_j$  for all  $i$  and  $j$ ). Since this is the case for a problem with homogeneous density, we expect that this scaling reduces the condition number of the stiffness matrix to roughly that for a similar problem with homogeneous density. Indeed, in general we obtain a condition number that is slightly better than that for a problem with constant density. Alternatively, in light of (30) we may want to scale the linear system such that all columns have the same norm. In the next section, we discuss the effects of rescaling for a simple 1D problem with heterogeneous density.

We propose to rescale the stiffness matrices  $\mathbf{K}$  by multiplying with a diagonal matrix on both sides (for symmetry),

$$\tilde{\mathbf{K}} = \mathbf{D}^{-1/2} \mathbf{K} \mathbf{D}^{-1/2}, \quad (31)$$

where the entries of the diagonal matrix  $\mathbf{D}$  are either the diagonal coefficients of  $\mathbf{K}$  or the absolute column sums of  $\mathbf{K}$ , i.e.  $d_i = \|\mathbf{k}_i\|_1$ . In Figure 4 we compare the condition numbers of stiffness matrices that arise in topology optimization for a model problem on a  $18 \times 6 \times 3$  mesh with the condition numbers of the rescaled stiffness matrices. The model problem is the same as that used in the numerical results section. The condition numbers of the stiffness matrices quickly rise to about  $10^{11}$  after only a few optimization steps. However, the condition numbers of the rescaled matrices remain at about the same level as those at the beginning (at about  $10^5$ ).

To obtain rapid convergence for iterative methods, it is important to further reduce the condition number after rescaling by more general preconditioning techniques. In our numerical experiments, we use an incomplete Cholesky decomposition with zero fill-in of the rescaled stiffness matrix as a preconditioner [25],

$$\tilde{\mathbf{K}} = \mathbf{D}^{-1/2} \mathbf{K} \mathbf{D}^{-1/2} \approx \mathbf{L} \mathbf{L}^T. \quad (32)$$

Finally, we note that diagonal scaling does not decrease the relative accuracy of the matrix coefficients, and hence such scaling leads to a real improvement in the worst case numerical error in the computed solution. The second type of preconditioning, using the incomplete Cholesky decomposition, improves the rate of convergence, but will typically not affect the accuracy of the computed solutions. Since the Cholesky decomposition may fail for a very ill-conditioned matrix, we explicitly rescale the stiffness matrix before we compute the incomplete Cholesky decomposition.

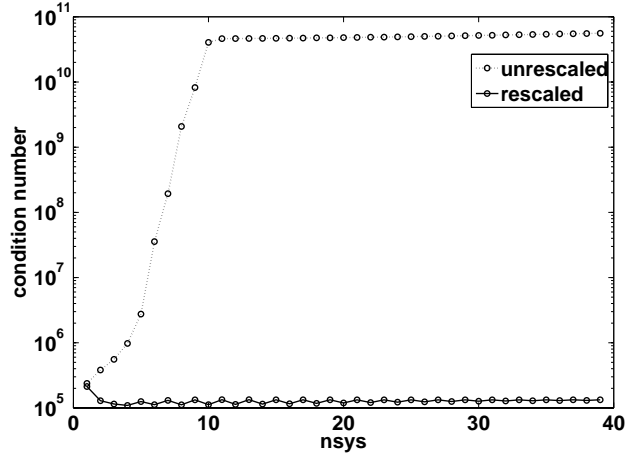


Figure 4. Condition numbers of stiffness matrices and rescaled stiffness matrices for the model problem in Figure 7 on a  $18 \times 6 \times 3$  mesh.

### 5.2. Rescaling for a 1D elasticity problem

We use an idealized 1D elasticity problem with piece-wise constant modulus of elasticity to explain the idea of rescaling. Consider the following problem.

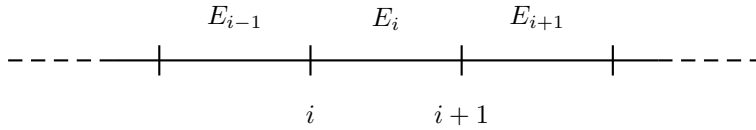


Figure 5. Piecewise constant modulus of elasticity  $E_i$ .

Find  $u(x)$  with boundary conditions  $u(0) = 0$  and  $u(1) = 1$ , such that

$$\mathbf{a}(u, v) \equiv \int_0^1 E(x) u_x v_x dx = 0, \quad \text{with } E(x) \geq E_0 > 0, \quad (33)$$

for all  $v$  with  $v(0) = v(1) = 0$ . Furthermore, following the typical case of topology optimization, we assume that  $E$  is piecewise constant and varies over a large range of values.

For simplicity, we discretize the problem using piecewise linear nodal basis functions and a mesh with equal length elements. This yields the following linear system.

$$\begin{bmatrix} E_1 + E_2 & -E_2 & & & \\ -E_2 & E_2 + E_3 & -E_3 & & \\ & \ddots & \ddots & \ddots & \\ & & -E_{n-2} & E_{n-2} + E_{n-1} & \end{bmatrix} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_{n-1} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}. \quad (34)$$

We can write this system of equations as follows (note  $(Eu_x)_x = 0 \Leftrightarrow Eu_x = \text{constant}$ ).

$$E_i(u_i - u_{i-1}) - E_{i+1}(u_{i+1} - u_i) = 0, \quad \text{for } i = 1, \dots, n-1,$$

where we have used  $u_0 = 0$  and  $u_n = 1$ . Introducing the difference matrix

$$\mathbf{D}_1 = \begin{bmatrix} 1 & & & & \\ -1 & 1 & & & \\ & & \ddots & \ddots & \\ & & & -1 & 1 \end{bmatrix}$$

and the diagonal matrix  $\mathbf{\Omega} = \text{diag}(E_1, E_2, \dots, E_{n-1})$ , we can write (34) as

$$(\mathbf{D}_1^T \mathbf{\Omega} \mathbf{D}_1 + E_n \mathbf{e}_{n-1} \mathbf{e}_{n-1}^T) \mathbf{u} = E_n \mathbf{e}_{n-1}. \quad (35)$$

For a problem with constant modulus of elasticity,  $E$ , this equation gives

$$E (\mathbf{D}_1^T \mathbf{D}_1 + \mathbf{e}_{n-1} \mathbf{e}_{n-1}^T) \mathbf{u} = E \mathbf{e}_{n-1}, \quad (36)$$

where  $\mathbf{D}_1^T \mathbf{D}_1 + \mathbf{e}_{n-1} \mathbf{e}_{n-1}^T$  is the well-known tridiagonal matrix with coefficients  $[-1 \ 2 \ -1]$ .

Next, we want to demonstrate two issues. Comparing (35) with (36), it is clear that extreme ill-conditioning in (35) must arise from the scaling introduced by  $\mathbf{\Omega}$ . First, we demonstrate that this leads to a condition number (bound) that is roughly the product of the condition number of the constant elasticity problem and the condition number of  $\mathbf{\Omega}$ . Second, we show how a proper (re)scaling brings the condition number down to that for the constant elasticity case, if the solution is properly defined. We note that for general choices of  $\mathbf{\Omega}$  there may be no diagonal scaling that reduces the condition number. For example, in 1D if we have two non-adjacent ‘holes’ the displacement for material in between the holes is not properly defined (as the modulus of elasticity goes to zero), since there is no connection to any point with a fixed displacement. In higher dimensions this is rarely a problem, as the topology optimization algorithm leads to energetically favorable solutions that do not have such anomalies.

Below, we need the following well-known result for symmetric positive definite matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$  and  $\alpha, \beta \in \mathbb{R}^+$  [26, pp. 338-9]. Let  $\mathbf{A}$  and  $\mathbf{B}$  be such that for all  $\mathbf{u} \neq 0$

$$\alpha \leq \frac{\mathbf{u}^T \mathbf{A} \mathbf{u}}{\mathbf{u}^T \mathbf{B} \mathbf{u}} \leq \beta. \quad (37)$$

Then

$$\kappa(\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2}) \leq \frac{\beta}{\alpha}, \quad (38)$$

where  $\kappa$  denotes the condition number.

Using (37–38), we can bound the condition number of the matrix in (35) as follows.

$$\kappa(\mathbf{D}_1^T \mathbf{\Omega} \mathbf{D}_1 + E_n \mathbf{e}_{n-1} \mathbf{e}_{n-1}^T) = \frac{\max_{\|\mathbf{u}\|=1} \mathbf{u}^T (\mathbf{D}_1^T \mathbf{\Omega} \mathbf{D}_1 + E_n \mathbf{e}_{n-1} \mathbf{e}_{n-1}^T) \mathbf{u}}{\min_{\|\mathbf{u}\|=1} \mathbf{u}^T (\mathbf{D}_1^T \mathbf{\Omega} \mathbf{D}_1 + E_n \mathbf{e}_{n-1} \mathbf{e}_{n-1}^T) \mathbf{u}}. \quad (39)$$

Let  $\mathbf{y} = \mathbf{D}_1 \mathbf{u}$  and hence  $\mathbf{u} = \mathbf{D}_1^{-1} \mathbf{y}$ . Then

$$\begin{aligned} \mathbf{u}^T \mathbf{D}_1^T \mathbf{\Omega} \mathbf{D}_1 \mathbf{u} + E_n \mathbf{u}^T \mathbf{e}_{n-1} \mathbf{e}_{n-1}^T \mathbf{u} &= (\mathbf{D}_1 \mathbf{u})^T \mathbf{\Omega} (\mathbf{D}_1 \mathbf{u}) + E_n u_{n-1}^2 \\ &= \frac{\mathbf{y}^T \mathbf{\Omega} \mathbf{y} + E_n (y_1 + \dots + y_{n-1})^2}{(\mathbf{D}_1^{-1} \mathbf{y})^T (\mathbf{D}_1^{-1} \mathbf{y})}. \end{aligned} \quad (40)$$

Let  $E_{\min} = \min_i(E_i)$ ,  $E_{\max} = \max_i(E_i)$ , and let  $\lambda_{\min}$  be the smallest eigenvalue of the matrix  $\mathbf{D}_1\mathbf{D}_1^T$  and  $\lambda_{\max}$  its largest eigenvalue. Furthermore, note that  $\mathbf{D}_1\mathbf{D}_1^T$  and  $\mathbf{D}_1^T\mathbf{D}_1$  have the same eigenvalues. Since  $\mathbf{y}$  appears quadratically in both the numerator and the denominator, we can assume  $\mathbf{y}$  to be normalized. Then, (40) gives

$$E_{\min}\lambda_{\min} \leq \frac{\mathbf{y}^T\boldsymbol{\Omega}\mathbf{y} + E_n(y_1 + \dots + y_{n-1})^2}{(\mathbf{D}_1^{-1}\mathbf{y})^T(\mathbf{D}_1^{-1}\mathbf{y})} \leq nE_{\max}\lambda_{\max}, \quad (41)$$

which finally gives

$$\kappa(\mathbf{D}_1^T\boldsymbol{\Omega}\mathbf{D}_1 + E_n\mathbf{e}_{n-1}\mathbf{e}_{n-1}^T) \leq n\frac{E_{\max}}{E_{\min}}\frac{\lambda_{\max}}{\lambda_{\min}} = n\kappa(\boldsymbol{\Omega})\kappa(\mathbf{D}_1^T\mathbf{D}_1). \quad (42)$$

Next, we show that scaling a problem (without non-adjacent ‘holes’) reduces the condition number of the linear system to roughly that of a problem with constant elasticity. Let

$$\mathbf{S} = \text{diag}(E_1 + E_2, E_2 + E_3, \dots, E_{n-1} + E_n). \quad (43)$$

We have

$$\begin{aligned} & \frac{\mathbf{u}^T(\mathbf{D}_1^T\boldsymbol{\Omega}\mathbf{D}_1 + E_n\mathbf{e}_{n-1}\mathbf{e}_{n-1}^T)\mathbf{u}}{\mathbf{u}^T\mathbf{S}\mathbf{u}} = \\ & \frac{E_1u_1^2 + E_2(u_1 - u_2)^2 + \dots + E_{n-1}(u_{n-2} - u_{n-1})^2 + E_nu_{n-1}^2}{E_1u_1^2 + E_2(u_1^2 + u_2^2) + \dots + E_{n-1}(u_{n-2}^2 + u_{n-1}^2) + E_nu_{n-1}^2} = \\ & \frac{(E_1 + E_2)u_1^2 + \dots + (E_{n-1} + E_n)u_{n-1}^2 - 2(E_2u_1u_2 + \dots + E_{n-1}u_{n-2}u_{n-1})}{(E_1 + E_2)u_1^2 + \dots + (E_{n-1} + E_n)u_{n-1}^2} = \\ & 1 - \frac{2(E_2u_1u_2 + \dots + E_{n-1}u_{n-2}u_{n-1})}{(E_1 + E_2)u_1^2 + \dots + (E_{n-1} + E_n)u_{n-1}^2} = \\ & 1 - \frac{2(E_2u_1u_2 + \dots + E_{n-1}u_{n-2}u_{n-1})}{E_1u_1^2 + E_2(u_1^2 + u_2^2) + \dots + E_{n-1}(u_{n-2}^2 + u_{n-1}^2) + E_nu_{n-1}^2}. \end{aligned} \quad (44)$$

It is easy to see that the maximum of (44) is bounded by 2. The condition number therefore depends mainly on the minimum of (44). We consider three examples.

The first example considers the case of constant modulus of elasticity. The second example demonstrates that the case of a bar with variable modulus (solid bar with a ‘hole’) leads to about the same condition number after scaling as the case of a bar with constant modulus (homogeneous). The third example shows that for the hypothetical case of a 1D bar with two non-adjacent ‘holes’ scaling cannot remove the actual singularity.

#### *Example 1: Constant Modulus*

For a constant modulus of elasticity, (44) leads to

$$\begin{aligned} & \frac{\mathbf{u}^T(E\mathbf{D}_1^T\mathbf{D}_1 + E\mathbf{e}_{n-1}\mathbf{e}_{n-1}^T)\mathbf{u}}{\mathbf{u}^T\mathbf{S}\mathbf{u}} \\ & = 1 - \frac{2u_1u_2 + \dots + 2u_{n-2}u_{n-1}}{u_1^2 + (u_1^2 + u_2^2) + \dots + (u_{n-2}^2 + u_{n-1}^2) + u_{n-1}^2}. \end{aligned} \quad (45)$$

The minimum for (45) is obtained for  $u_i = \sin(\pi ih)$  which gives  $u_{i-1}u_i \approx u_i^2$  and minimizes the influence of the terms  $Eu_1^2$  and  $Eu_{n-1}^2$ . This leads to a condition number for the preconditioned system of  $O(h^{-2})$ .

*Example 2: Variable Modulus*

Now consider a problem with a ‘hole’ at the end of the bar;  $E_i = 1$  for  $i = 1, \dots, n - 5$ , where  $n \gg 5$ , and  $E_i = \varepsilon$  (small) for the remaining elements. The minimum for (44) is obtained for a vector  $u$  such that  $|u_i| = O(1)$  for  $i = 1, \dots, n - 5$  and  $|u_i| = O(\varepsilon)$  for the remaining elements. After substituting for the  $E_i$  in (44) and dropping the  $u_i$  terms that are  $O(\varepsilon)$ , we need to minimize the following expression

$$1 - \frac{2u_1u_2 + \dots + 2u_{n-6}u_{n-5}}{u_1^2 + (u_1^2 + u_2^2) + \dots + (u_{n-6}^2 + u_{n-5}^2) + \varepsilon u_{n-5}^2}. \quad (46)$$

Comparing (46) to (45), we see that this minimization problem is essentially the same as the one for the constant modulus (with a few terms of small magnitude dropped). Therefore, the resulting condition number will be about the same.

*Example 3: Hypothetical Case*

Finally, consider a hypothetical problem of a 1D bar with two non-adjacent ‘holes’. Let  $n = 5$ , and let  $E_1 = E_3 = E_5 = 1$  and  $E_2 = E_4 = \varepsilon$ . Now taking  $u_1 = u_4 = 0$  and  $u_2 = u_3 = 1$  in (44) gives

$$\begin{aligned} \min_{\mathbf{u} \neq \mathbf{0}} \frac{\mathbf{u}^T (\mathbf{D}_1^T \boldsymbol{\Omega} \mathbf{D}_1 + E_n e_{n-1} e_{n-1}^T) \mathbf{u}}{\mathbf{u}^T \mathbf{S} \mathbf{u}} &\leq \\ 1 - \frac{2E_3 u_2 u_3}{E_2 u_2^2 + E_3 (u_2^2 + u_3^2) + E_4 u_3^2} &= \\ 1 - \frac{2}{\varepsilon + 2 + \varepsilon} &= \frac{\varepsilon}{1 + \varepsilon}. \end{aligned} \quad (47)$$

So, (47) can be made arbitrarily small, and therefore the condition number  $\kappa(\mathbf{D}_1^T \boldsymbol{\Omega} \mathbf{D}_1 + E_n e_{n-1} e_{n-1}^T)$  can be made arbitrarily large.

## 6. SOME IMPLEMENTATION ISSUES

We have developed our C/C++ code in an object-oriented fashion, since we intend to make the RMINRES solver available as public domain software. This makes it easy to integrate the software in other packages, and also facilitates maintenance and extension, while retaining high efficiency. We store sparse matrices in compressed sparse row format (CSR). The (column) vectors in  $\mathbf{U}$ ,  $\mathbf{C}$ ,  $\mathbf{V}_j$ , and so on, are stored as one-dimensional arrays linked by a one-dimensional array of pointers. The memory required by the system matrix and the incomplete Cholesky factor is linear in the number of unknowns,  $n$ , since the number of nonzero coefficients per row is never greater than 81. The RMINRES method requires only matrix-vector multiplications, dot products, vector updates, forward and backward solves with the incomplete Cholesky factors, and the incomplete Cholesky decomposition itself. All of these operations have linear computational cost.

The small matrices, e.g., the matrices in (15), are all stored as dense matrices in column-wise ordering (F77 format), so that dense matrix routines from LAPACK and BLAS can be used. For the generalized eigenvalue problem (13) we use the LAPACK routine DSYGV. For convenience we use the CLAPACK library, which provides an interface for C programs

Table I. Three discretizations used for the example in Figure 6.

problem	mesh size	#unknowns (in simulation)	solution time	optimization steps	iterative solver
small	$36 \times 12 \times 6$	9,360	0.1h	142	RMINRES(100,10)
medium	$84 \times 28 \times 14$	107,184	2.4h	139	RMINRES(100,10)
large	$180 \times 60 \times 30$	1,010,160	45.7h	130	RMINRES(200,10)

to LAPACK. However, since CLAPACK routines call the corresponding LAPACK routines, we still need to adhere to F77 storage formats. The computational cost of the work with these small matrices is negligible. Finally, we note that the computational cost is significantly reduced by taking the simplifications in (17)–(23) into account.

## 7. NUMERICAL RESULTS

We demonstrate the performance of the MINRES iterative solver with the recycling and preconditioning techniques discussed in Sections 4 and 5 on a large 3D design problem. We also provide some analysis for the selection of the linear solver parameters.

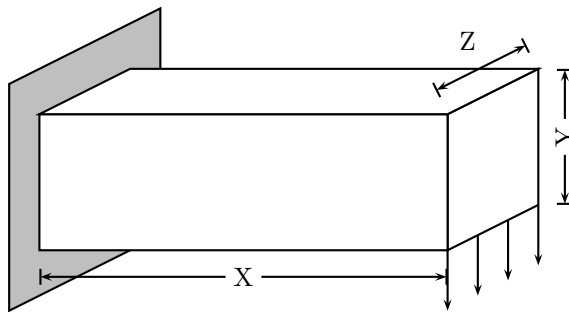


Figure 6. Design problem: finding optimal material distribution in a hexahedron with the left end fixed and a distributed load applied on the right bottom edge. ( $X : Y : Z = 3 : 1 : 1$ ).

Figure 6 shows our model problem. The volume fraction is 50%, and the radius of the filter is 10% of  $Y$ . We use continuation on the density penalization, ranging from 1 to 3 with increments of 0.5. As stated before, we use the OC method as the optimization algorithm. The convergence criterion is that the maximum change in the design variables is less than 0.01. We test three discretizations of increasing mesh resolution. Exploiting the symmetry of the problem, we model and simulate only half of the domain. For each test case, Table I lists the mesh size (for half of the domain), the number of unknowns, the overall solution time, the number of optimization steps, and the parameters used for the iterative solver. Figure 7 shows the final topologies.

First, we analyze the convergence properties of RMINRES for several parameter choices on

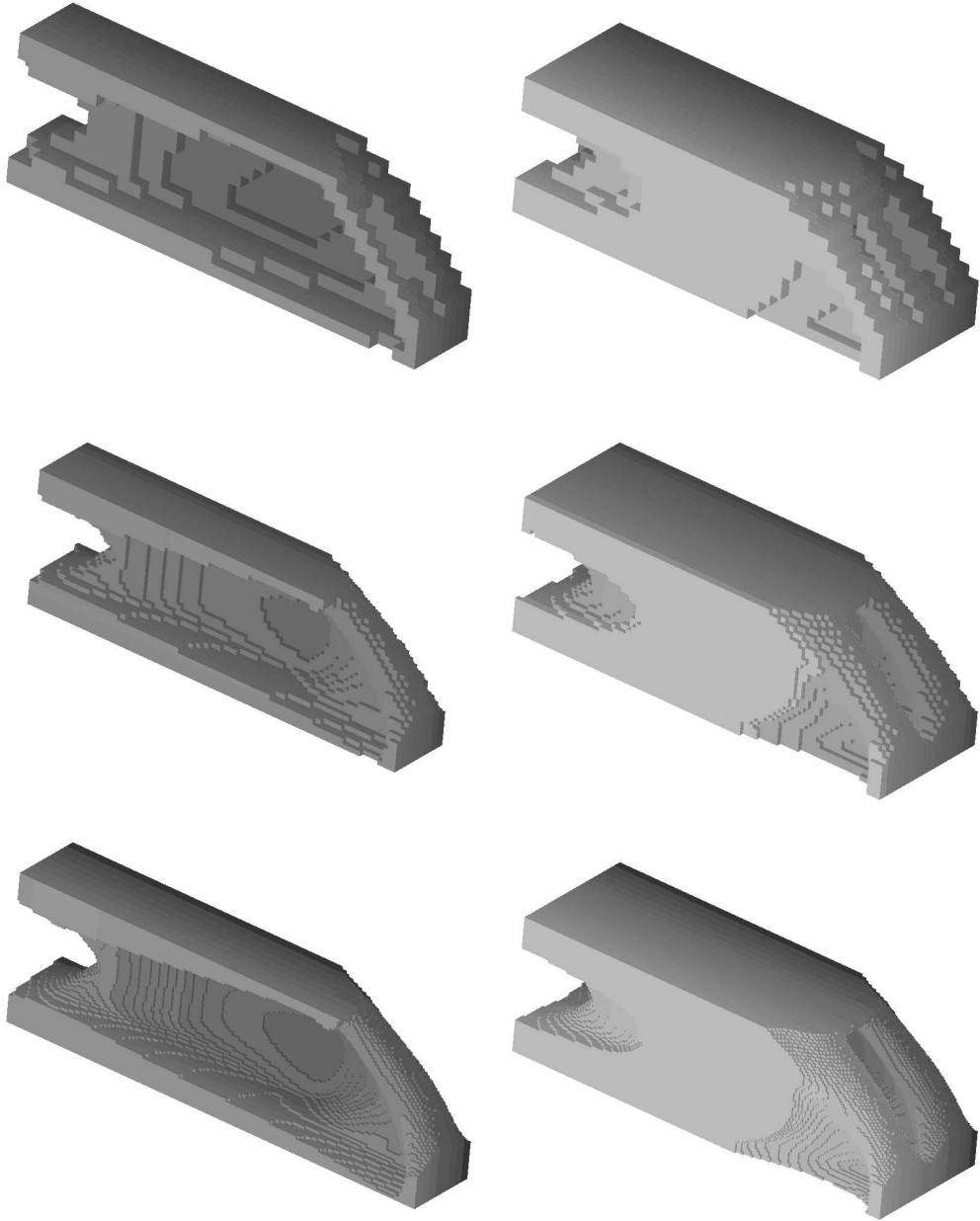


Figure 7. Final topologies of the model problem. Left: half domain; right: full domain. Top row: small mesh; middle row: medium size mesh; bottom row: large mesh.

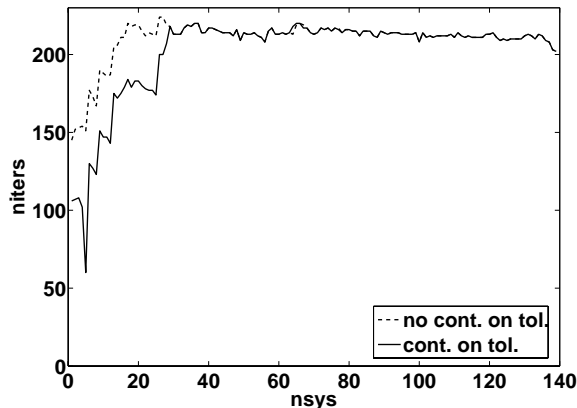


Figure 8. Reduction in the number of iterations using a relaxed tolerance for the linear solver (MINRES without recycling).

the medium size ( $84 \times 28 \times 14$ ) mesh. The number of optimization steps to compute the optimal design is 139, requiring the solution of 139 linear systems.

As mentioned in the first section, we can vary the tolerance for the iterative solver, since less accurate finite element solutions are sufficient at the beginning of the topology optimization process. So, we can also apply a continuation approach to the tolerance of the linear solver, which reduces the number of iterations in the early phase of the optimization process, as shown in Figure 8. The jumps in the iteration counts correspond to the steps where the tolerance of the linear solver is decreased or the penalization parameter  $p$  is increased. We update the solver tolerance and the penalization parameter every time the maximum change of the design variables drops below 0.1. Finally, we note that allowing a higher tolerance for the linear solver in the beginning of the optimization process did not affect the number of optimization steps required.

Next, we consider the parameters that govern the recycling for the MINRES solver, namely  $k$ , the dimension of the subspace that is recycled from one linear system to the next, and  $s$ , the maximum dimension of the Krylov subspace kept to periodically update the approximate invariant subspace that will be recycled. We carry out two sets of experiments to analyze the effects of varying these two parameters.

In the first set of experiments, we fix  $k = 10$  and vary  $s$ . Figure 9 compares the number of iterations and computation time for each linear system, for several choices of  $s$ . In the first few optimization steps, the topology changes significantly, and the effect of recycling is modest. After this, the recycling approach greatly reduces the number of iterations to solve each linear system. We see that if we keep a larger Krylov subspace to update the approximate invariant subspace, the recycling becomes more effective in reducing iteration counts. Since the dimension of the recycle space itself does not change, this suggests that we obtain a more accurate approximation to the invariant subspace this way. This reduction in iterations significantly reduces the computation time for RMINRES, in spite of the computational

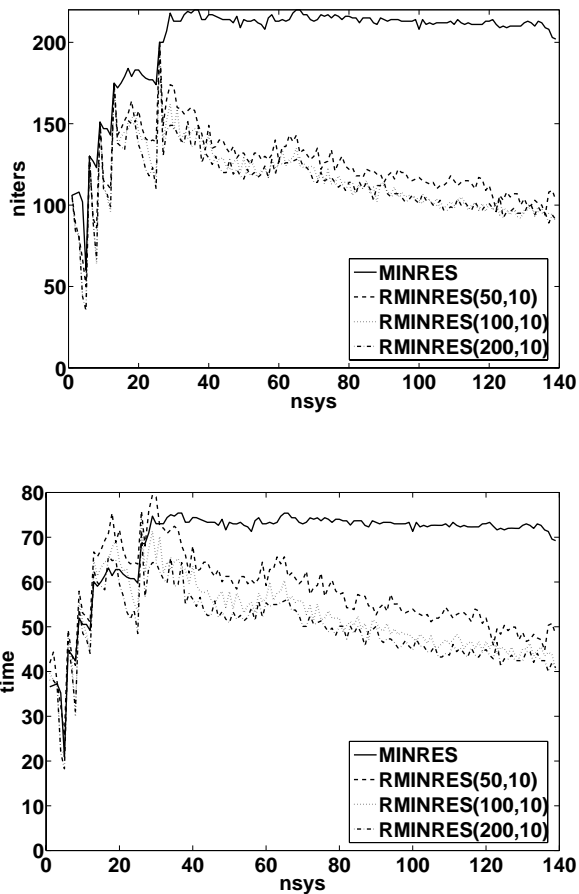


Figure 9. Number of iterations (niters) and time (seconds) of  $\text{RMINRES}(s, k)$  with fixed  $k = 10$  and varying  $s$

overhead from the orthogonalizations against the recycle space and from the updates of the recycle space. Towards the end of the optimization process, recycling leads to a 40% reduction in computation time and a 50% reduction in iterations. Notice that increasing  $s$  beyond 100 has limited effect since  $\text{RMINRES}$  rarely takes more than 100 iterations for this problem. However, for harder problems, e.g., for finer meshes, the solver may not converge so fast. In that case, larger values for  $s$  can be helpful. Note that increasing  $s$  does not increase the computational cost of  $\text{RMINRES}$ . The only limit on  $s$  is the memory size.

In the second set of experiments, we fix  $s = 100$  and vary  $k$ . The parameter  $k$  affects both the computational cost per iteration, specifically the number of orthogonalizations and the cost of subspace selection, and the total number of iterations for the solver. There is a tradeoff between these two factors, and in Figure 10 we compare the number of iterations and computational

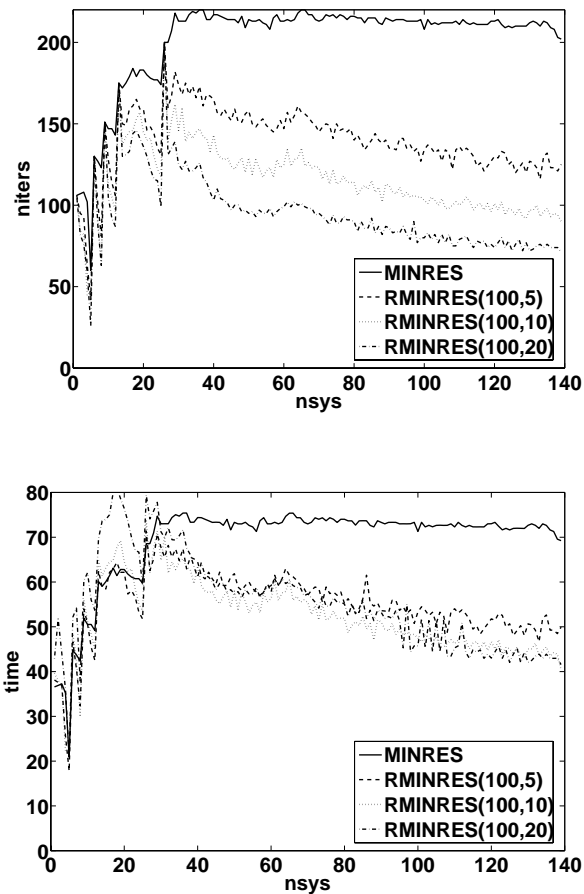


Figure 10. Number of iterations (niters) and time (seconds) of  $\text{RMINRES}(s, k)$  with varying  $k$  and fixed  $s = 100$

time for several values of  $k$ . Increasing  $k$  leads to a significant improvement in the convergence rate; towards the end we obtain a factor 3 reduction in the number of iterations. Time-wise, we obtain a 40% improvement. We also see that the computation time is not overly sensitive to the choice of  $k$ .

Finally, we compare a state-of-the-art sparse direct solver with our iterative method, including recycling and the continuation on the solver tolerance, the initial scaling, and the preconditioner. We use the multifrontal, supernodal Cholesky factorization from the TAUCS package [27]. We ran the comparison on a PC with an AMD Opteron™252 2.6GHz 64-bit processor, 8GB RAM of memory, and the SuSE Linux system. Table II lists the computation times of both the direct solver and the iterative solver for a single linear system, for several problem sizes. For the direct solver, the data include the time for the Cholesky decomposition,

Table II. Run time comparison (seconds) of direct and iterative solvers. The direct solver is the multifrontal, supernodal Cholesky factorization in TAUCS; the iterative solver is RMINRES with rescaling and incomplete Cholesky preconditioner and the continuation on the solver tolerance. These timings were obtained on a PC with an AMD Opteron <sup>TM</sup>252 2.6GHz 64-bit processor, 8GB RAM of memory, and the SuSE Linux system.

problem size	#unknowns in simulation	direct solver time (s)			iterative solver time (s)		
		decompose	solve	<b>total</b>	min	max	<b>average</b>
small	9,360	0.95	0.01	0.96	0.94	2.25	1.66
medium	107,184	179.0	0.3	179.3	21.2	71.9	50.5
large	1,010,160	21241	4904	26154	254	1546	1170

and for the forward and backward substitution to solve one system. For the iterative solver, to make a fair comparison the data include the maximum, minimum and average run time taken over all linear systems. We note that the run time of the direct solver seems to scale worse than quadratic, whereas the iterative solver scales slightly worse than linear. Since the matrix changes at every optimization step, we cannot reuse the Cholesky factors of the direct solver and a new factorization must be computed every optimization step.

## 8. CONCLUSIONS

In this paper, we investigate iterative solvers for the equilibrium equations in topology optimization. We address the main problem, the extreme ill-conditioning, by two approaches. First, we rescale the system to reduce the ill-conditioning due to the variation in the density. Second, we use an incomplete Cholesky preconditioner for the resulting linear system. In general, the rescaling improves the accuracy of the solution as well as the convergence rate. The Cholesky preconditioner improves convergence rate, but in general has no effect on the accuracy. Exploiting the slowly changing nature of the linear systems arising in topology optimization, the RMINRES method, which recycles selected subspaces, leads to a further significant reduction in iterations and run time.

We benchmark our methods for a design problem on a sequence of increasingly finer meshes. The largest problem has more than a million unknowns, resulting in a smooth solution. With the proposed methods, we can solve these problems on a single PC in a reasonably short time.

## APPENDIX

### LIST OF SYMBOLS

$n$	the size of the linear system
$m$	the current number of iterations of the iterative method
$k$	the dimension of the subspace to be recycled
$s$	the maximum number of Lanczos vectors kept for updating the recycle space

$j$	the index of the update cycle
$a(u, v)$	bilinear form of the 1D idealized elasticity problem
$c$	compliance
$E, E_i$	modulus of elasticity
$\mathbf{f}$	load vector
$\mathbf{K}, \mathbf{K}^{(i)}$	global stiffness matrix
$\mathbf{K}_0$	element stiffness matrix
$K_{ij}$	$(i, j)$ coefficient of $\mathbf{K}$
$\mathbf{k}_j$	the $j$ th column of $\mathbf{K}$
$p$	penalization parameter
$\mathbf{u}, \mathbf{u}^{(i)}$	displacement vector
$V$	total volume
$\rho, \rho^{(i)}$	density distribution (the design variables of topology optimization problems)
$\rho_0$	positive lower bound on the density to avoid singularity
$\Omega$	diagonal matrix of material constants
$\mathbf{C}, \mathbf{C}_j$	orthogonal basis of $\mathbf{K}\mathbf{U}$
$\underline{\mathbf{H}}_m, \mathbf{H}_m$	$(m + 1) \times m$ Hessenberg matrix generated in the Arnoldi iteration, and the first $m$ rows of $\underline{\mathbf{H}}_m$
$\mathbf{r}_0, \mathbf{r}_m$	residual
$\underline{\mathbf{T}}_j, \overline{\mathbf{T}}_j$	tridiagonal matrix generated by Lanczos iteration and its extended version with an additional top row
$\mathbf{U}, \mathbf{U}_j$	basis for the recycle space
$\mathbf{V}_m$	first $m$ vectors of the Krylov subspace, $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m$
$\mathbf{V}_j, \underline{\mathbf{V}}_j$	the Krylov subspace generated in $j$ th cycle, $\mathbf{v}_{(j-1)s+1}, \dots, \mathbf{v}_{js}$ , and its extended version, $\mathbf{v}_{(j-1)s}, \dots, \mathbf{v}_{js+1}$
$\kappa(\cdot)$	condition number of a matrix
$(\theta, \mathbf{w})$	an harmonic Ritz pair of $\mathbf{K}$ with Ritz value $\theta$ and Ritz vector $\mathbf{w}$

## ABBREVIATIONS

GMRES	Generalized Minimum Residual method [20]
MINRES	Minimum Residual method [6]
GCRODR	Generalized Conjugate Residual method with inner Orthogonalization and Deflated Restarting [3]
GCRO	Generalized Conjugate Residual method with inner Orthogonalization [23]
GCROT	Generalized Conjugate Residual method with inner Orthogonalization and optimal Truncation [24]

## REFERENCES

1. Borrvall T, Petersson J. Large-scale topology optimization in 3D using parallel computing. *Computer Methods in Applied Mechanics and Engineering* 2001; **190**:6201–6229.
2. Vemaganti K, Lawrence WE. Parallel methods for optimality criteria-based topology optimization. *Computer Methods in Applied Mechanics and Engineering* 2005; **194**:3637–3667.
3. Parks ML, de Sturler E, Mackey G, Johnson DD, Maiti S. Recycling Krylov subspaces for sequences of linear systems. Tech. rep., UIUC DCS-R-2004-2421, available from <http://www-faculty.cs.uiuc.edu/~sturler>, 2004.
4. Kilmer ME, de Sturler E. Recycling subspace information for diffuse optical tomography. *SIAM Journal on Scientific Computing* 2005; **to appear**, available from <http://www-faculty.cs.uiuc.edu/~sturler>.
5. Sigmund O, Jensen J. Systematic design of phononic band gap materials and structures. *Philosophical Transactions of the Royal Society London, Series A (Mathematical, Physical and Engineering Sciences)* 2003; **361**:1001–1019.
6. Paige CC, Saunders MA. Solution of sparse indefinite systems of linear equations. *SIAM Journal on Numerical Analysis* 1975; **12**:617–629.
7. Bendsoe MP, Sigmund O. *Topology Optimization: Theory, Methods and Applications*. Springer-Verlag: Berlin, 2003.
8. Sigmund O. Topology optimization: a tool for the tailoring of structures and materials. *A special issue of the Philosophical Transactions of the Royal Society: Science into the next Millennium (Issue III, Mathematics, Physics and Engineering)* 2000; **358**(1765):211–228.
9. Rozvany GIN. Aims, scope, methods, history and unified terminology of computer-aided topology optimization in structural mechanics. *Structural and Multidisciplinary Optimization* 2001; **21**:90–108.
10. Mackerle J. Topology and shape optimization of structures using FEM and BEM: A bibliography (1999–2001). *Finite Elements in Analysis and Design* 2003; **39**:243–253.
11. Paulino GH, Page III RC, Silva ECN. A java-based topology optimization program with web access: nodal design variable approach. In *Proceedings of the 6th World Congress on Structural and Multidisciplinary Optimization (WCSMO6), Rio de Janeiro, Brazil, May 30 – June 3, 2005*. International Society for Structural and Multidisciplinary Optimization, 2005; .
12. Stump FV, Paulino GH, Silva ECN. Material distribution design of functionally graded rotating disks with stress constraint. In *Proceedings of the 6th World Congress on Structural and Multidisciplinary Optimization (WCSMO6), Rio de Janeiro, Brazil, May 30 – June 3, 2005*. International Society for Structural and Multidisciplinary Optimization, 2005; .
13. Bendsoe MP. Optimal shape design as a material distribution problem. *Structural Optimization* 1989; **1**:193–202.
14. Bendsoe MP, Sigmund O. Material interpolation schemes in topology optimization. *Archives of Applied Mechanics* 1999; **69**(9-10):635–654.
15. Matsui K, Terada K. Continuous approximation of material distribution for topology optimization. *International Journal for Numerical Methods in Engineering* 2004; **59**:1925–1944.
16. Paulino GH, Silva ECN. Design of functionally graded structures using topology optimization. *Materials Science Forum* 2005; **492** – **493**:435 – 440.
17. Sigmund O. On the design of compliant mechanisms using topology optimization. *Mechanics of Structures and Machines* 1997; **25**:495–526.
18. Rozvany GIN. *Structural Design Via Optimality Criteria*. Kluwer Academic Publishers Group: Dordrecht, 1989.
19. Svanberg K. The method of moving asymptotes: a new method for structural optimization. *International Journal for Numerical Methods in Engineering* 1987; **24**(2).
20. Saad Y, Schultz MH. GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems. *SIAM Journal on Scientific and Statistical Computing* 1986; **7**(3):856–869.
21. Arnoldi WE. The principle of minimized iterations in the solution of the matrix eigenvalue problem. *Quarterly of Applied Mathematics* 1951; **9**:17–29.
22. Saad Y. *Iterative Methods for Sparse Linear Systems, 2nd edition*. SIAM: Philadelphia, 2003.
23. de Sturler E. Nested Krylov methods based on GCR. *Journal of Computational and Applied Mathematics* 1996; **67**:15–41.
24. de Sturler E. Truncation strategies for optimal Krylov subspace methods. *SIAM Journal on Numerical Analysis* 1999; **36**:864–889. Electronically available from <http://epubs.siam.org>.
25. Meijerink J, Van der Vorst H. An iterative solution method for linear equations systems of which the coefficient matrix is a symmetric  $M$ -matrix 1977; .
26. Axelsson O, Barker VA. *Finite Element Solution of Boundary Value Problems*, vol. 35 of *Classics in Applied Mathematics*. SIAM: 3600 University City Science Center, Philadelphia, PA 19104-2688, 2001.

27. Toledo S. *Taucs: A Library of Sparse Linear Solvers*. School of Computer Science, Tel-Aviv University, September 2003.