# Reduced forward operator for electromagnetic wave scattering problems

K.W.A. Van Dongen, C. Brennan and W.M.D. Wright

**Abstract:** The paper describes a reduced forward operator for solving electromagnetic scattering problems using a volume integral equation in conjunction with a conjugate gradient fast Fourier transform scheme. The reduction is obtained by decoupling of the interaction between the locations in the spatial computational domain at which there is non-zero contrast and those positions at which there is zero contrast. The decoupling is achieved by multiplication of the kernel by a diagonal matrix whose entries reflect the presence or absence of contrast at the associated point. Analysis supported by numerical experiments shows that the conjugate gradient algorithm applied to the reduced system converges more rapidly than when it is applied to the original system.

## 1 Introduction

The propagation and scattering of acoustic, elastodynamic or electromagnetic wave fields are described by integral equation formulations [1]. For the electromagnetic case, the equation formulates, in the temporal Laplace domain, the total electric wave field $\vec{E}^{tot}(\vec{r})$, at location $\vec{r}$ in the spatial computational domain $\mathbb{D}$, as the sum of an incident and a scattered electric wave field, where the scattered wave field is described by the convolution of a Green's tensor function $\underline{G}(\vec{r}, \vec{r}')$ with contrast sources. These contrast sources are the product of a contrast function $\chi(\vec{r}')$, defined by changes in the medium's electromagnetic parameters of permittivity and conductivity, and the local total electric wave field. The configuration shown in Fig. 1 shows that non-zero contrast sources reside in the region $\mathbb{D}' \subset \mathbb{D}$, whereas there is zero contrast in the region $\mathbb{D}'' = \mathbb{D} \backslash \mathbb{D}'$. For this configuration, the following integral equation holds [1]

$$\vec{E}^{\text{tot}}(\vec{r}) = \vec{E}^{\text{inc}}(\vec{r}) + \int_{\vec{r}' \in \mathbb{D}'} \underline{G}(\vec{r}, \vec{r}') \chi(\vec{r}') \vec{E}^{\text{tot}}(\vec{r}') \, dV(\vec{r}') \forall \vec{r} \in \mathbb{D} \tag{1}$$

For a numerical solution of the continuous integral equation (1) to be obtained, the unknown total wave field is written as a linear combination of $N$ basis functions. Application of the method of moments leads to the matrix equation

$$\boldsymbol{A}\boldsymbol{x} = \boldsymbol{b} \tag{2}$$
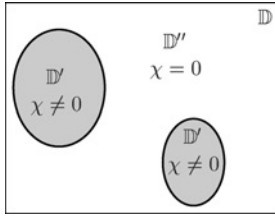
where the vector $\boldsymbol{b}$ contains information about the known incident wave field, and the vector $\boldsymbol{x}$ denotes the total wave field unknowns. The dense, complex-valued matrix $\boldsymbol{A}$ is of the form $\boldsymbol{I} - \boldsymbol{G}$, where $\boldsymbol{I}$ is an $N \times N$ identity matrix, and $\boldsymbol{G}$ contains coupling information between the basis functions used to represent the unknown total wave field $\boldsymbol{x}$. For practically sized problems, $\boldsymbol{A}$ cannot be readily inverted (or indeed explicitly stored), and we must turn to iterative methods to solve (2). Several techniques have been suggested in the literature for expediting the solution of these matrix equations. Examples include the adaptive integral equation method [2], precorrected FFT method [3], fast multipole method [4] and, the topic of this paper, the conjugate gradient fast Fourier transform method (CG-FFT) [5–9]. Specifically this paper analyses a procedure for improving the convergence of the CG-FFT and, in doing so, extends the results presented in [10] to the case of electromagnetic wave scattering.

We review the convergence behaviour of the CG algorithm and identifies the CG-FFT as a practical technique to reduce the computational requirement of each iteration. The FFT is used to compute efficiently the convolution of the source terms with the Green's function. A computational drawback of the CG-FFT method is the need to introduce 'dummy' unknowns in region $\mathbb{D}''$ to pad out the regular lattice structure. These dummy unknowns do not contribute to the scattered field but do adversely affect the convergence of the CG algorithm. We outline a simple measure by which the influence of the dummy unknowns can be hidden from the CG minimisation procedure without compromising the ability to perform rapid matrix vector multiplications with the FFT. This involves pre-multiplying both sides of (2) by a diagonal matrix whose values reflect the presence or absence of contrast in the associated basis function domain. We refer to the pre-multiplied $\boldsymbol{A}$ matrix as the reduced operator. Although this technique has been suggested before in the literature (see, for example, [7]), a rigorous analysis of its convergence has not to date been presented. We introduce such an analysis and show why the CG algorithm applied to the reduced operator must lead to a more rapid convergence than when it is applied to the full operator. Numerical results support this analysis.

## 2 Convergence of CG-FFT scheme

The CG inversion scheme [11] is a method for iteratively finding a solution for the linear system of size $N \times N$ presented in (2). The CG method cannot be applied to the

*IET Sci. Meas. Technol.*, 2007, **1**, (1), pp. 57–62

57

**Fig. 1** *Spatial volume $\mathbb{D}$ contains spatial region $\mathbb{D}'$, in which contrast function $\chi \neq 0$, and spatial region $\mathbb{D}''$, in which contrast function $\chi = 0$*

matrix equation (2) directly, as the matrix $A$ is not Hermitian positive-definite. Instead, the CG algorithm is applied to the equation

$$A^{\dagger}Ax = A^{\dagger}b \qquad (3)$$

where $A^{\dagger}$ is the conjugate transpose of $A$, and, hence, $A^{\dagger}A$ is a Hermitian positive-definite matrix. This technique (application of CG to $A^{\dagger}A$) is also known as 'CG applied to normal equations' or CG-NE, for which the solution scheme is shown in Table 1. The CG-NE algorithm starts with an initial guess $x_0$ and proceeds by constructing optimum correction vectors along a series of orthogonal search directions. The action of the algorithm at the $n$th step can be equivalently thought of as minimisation of the following error functional [12]

$$E_n = \sum_{i=1}^{N} |\langle u_i, r_0 \rangle|^2 [R_n(\lambda_i)]^2 \qquad (4)$$

where $u_i$ are the orthonormal eigenvectors of the matrix $A^{\dagger}A$, and $\lambda_i$ are the corresponding (real-valued) eigenvalues. $R_n$ is an $n$th-order polynomial whose coefficients are determined by the correction vectors; $r_0$ is the initial residual

$$r_0 = Ax_0 - b \qquad (5)$$

Thus the $n$th iteration implicitly involves the choice of an appropriate $n$th-order polynomial $R_n$ such that the error functional (4) is minimised. Once the iteration number $n$ reaches the number of independent eigenvalues of $A^{\dagger}A$, it is possible to choose a polynomial that is zero at each $\lambda_i$, and the error is reduced to zero. (In practice, accumulation of floating point roundoff errors due to finite arithmetic

**Table 1: Conjugate gradient method applied to normal equations using Polak–Ribière update directions**

initial step:
$\quad d_0 = 0$
$\quad x_0 = 0$
$\quad r_0 = Ax_0 - b$
for $j = 1, 2, \ldots$
$\quad \beta_j = \langle A^{\dagger}r_{j-1}, A^{\dagger}r_{j-1} - A^{\dagger}r_{j-2} \rangle / \|A^{\dagger}r_{j-2}\|^2$
$\quad d_j = -A^{\dagger}r_{j-1} + \beta_j d_{j-1}$
$\quad \alpha_j = -\langle Ad_j, r_{j-1} \rangle / \|Ad_j\|^2$
$\quad x_j = x_{j-1} + \alpha_j d_j$
$\quad r_j = Ax_j - b$

computer precision causes the residual to lose its accuracy gradually, and cancellation error causes the vectors to lose their A-orthogonality [13]. As a consequence, an error of exactly zero is never actually attained.) In practice, the algorithm is terminated after far fewer steps, once an acceptable threshold level of accuracy is achieved. How rapidly this threshold is reached depends on how spread out the eigenvalues of $A^{\dagger}A$ are. The CG will converge more rapidly when applied to linear systems with clustered eigenvalues as, in that case, it is possible to find a polynomial of lower order that possesses zeros near all of the independent eigenvalues, thus yielding a low value of $E_n$. Conversely, it is more difficult to force a low-order polynomial through a widely separated set of eigenvalues such that the polynomial is almost zero at all of them.

Let the eigenvalues of $A^{\dagger}A$ be ordered such that

$$\lambda_1(A^{\dagger}A) \geq \lambda_2(A^{\dagger}A) \geq \cdots \geq \lambda_N(A^{\dagger}A) \qquad (6)$$

The rate at which the CG converges thus depends on the ratio $R$, which measures the eigenvalue spread of $A^{\dagger}A$, that is $R$ is defined by

$$R = \frac{\lambda_1(A^{\dagger}A)}{\lambda_N(A^{\dagger}A)} \qquad (7)$$

Note that the eigenvalues of $A^{\dagger}A$ are related to the singular values of $A$ by

$$\sigma_i(A) = [\lambda_i(A^{\dagger}A)]^{1/2} \quad \text{for } i = 1, \ldots, N \qquad (8)$$

where $\sigma_i(A)$ are the singular values of $A$ and are ordered as follows

$$\sigma_1(A) \geq \sigma_2(A) \geq \cdots \geq \sigma_N(A) \qquad (9)$$

Consequently, it is trivial to express $R$ as a function of $\sigma_i(A)$

$$R = \frac{\sigma_1^2(A)}{\sigma_N^2(A)} \qquad (10)$$

$R$ is thus the square of the condition number of $A$.

The condition number of $A$ thus determines how many iterative steps are necessary before the CG algorithm will converge. Each iterative step involves a matrix-vector multiplication and thus involves $\mathcal{O}(N^2)$ computations. The volume integral equation lends itself to solution with a CG-FFT method where each iterative step involves $\mathcal{O}(N \log_2 N)$ computations. Application of the CG-FFT method necessitates the uniform discretisation of a rectilinear volume enclosing the scatterers. Under these conditions, the matrix $G$ can be decomposed into a diagonal matrix containing contrast information and a matrix with a cyclical structure containing Green's function information. This cyclical structure facilitates its efficient premultiplication of an arbitrary vector using the FFT. The use of the FFT to efficiently compute the convolutions inherent in the integral equation formulation reduces the computation time but necessitates the introduction of dummy unknowns at locations where there is no contrast. These dummy unknowns produce no scattered field and have no influence on the value of the other unknowns. They are introduced merely to support a uniform lattice structure that allows the use of the FFT. The CG-FFT makes no differentiation between these dummy unknowns and the unknowns lying within the scatterers themselves. This has an adverse effect on the convergence of the CG-FFT, as the CG minimisation procedure must solve for both unknowns alike. The following sections show

how the dummy unknowns can be effectively hidden from the CG minimisation procedure, so that faster convergence can be guaranteed without compromising the ability for the FFT to be used to compute the convolutions.

## 3 Reduced forward operator

Let the set $\mathbb{N}'$ represent the indices of the $N'$ basis functions whose domains possess non-zero contrast. The complementary set $\mathbb{N}''$ represents the index of the $N'' = N - N'$ basis functions whose domains possess no contrast. The presence or absence of contrast in a subdomain is manifested in the structure of $A$. Specifically, the columns $A_{mn''}$ for each $n'' \in \mathbb{N}''$ have all zero entries except for the diagonal entry, which is equal to 1. Equation (2) can thus be written as

$$
\begin{pmatrix}
A_{11} & \cdots & & \cdots & 0 & \cdots & & \cdots & A_{1N} \\
\vdots & \ddots & & & \vdots & & & & \vdots \\
\vdots & & & & 0 & & & & \vdots \\
A_{n''1} & \cdots & A_{n''(n''-1)} & 1 & A_{n''(n''+1)} & \cdots & & & A_{n''N} \\
\vdots & & & & 0 & & & & \vdots \\
\vdots & & & & \vdots & & \ddots & & \vdots \\
A_{N1} & \cdots & & \cdots & 0 & \cdots & & \cdots & A_{NN}
\end{pmatrix}
\begin{pmatrix}
x_1 \\ \vdots \\ \vdots \\ x_{n''} \\ \vdots \\ \vdots \\ x_N
\end{pmatrix}
= 
\begin{pmatrix}
b_1 \\ \vdots \\ \vdots \\ b_{n''} \\ \vdots \\ \vdots \\ b_N
\end{pmatrix}
\tag{11}
$$

where we have made explicit the structure of the $n''$th column and row for some $n'' \in \mathbb{N}''$. Similarly structured rows and columns will be present in the matrix for all other $n \in \mathbb{N}''$. Hence, the unknowns $x_{n''}$ for $n'' \in \mathbb{N}''$ do not affect the values of the unknowns at any other location in space. It is therefore desirable to remove their influence from the CG minimisation procedure. The CG process should strive to compute the unknowns only at locations where contrast is present. These values can subsequently be used to compute the fields in $\mathbb{D}''$, as desired. We propose removing the influence of the zero-contrast unknowns by pre-multiplying both sides of (2) with a diagonal matrix $\tilde{I}$ whose diagonal elements reflect the presence or absence of contrast, that is

$$
\tilde{I}_{mm} = \begin{cases} 1 & \forall \{m \in \mathbb{N}'\} \\ 0 & \forall \{m \in \mathbb{N}''\} \end{cases}
\tag{12}
$$

Note that this pre-multiplication does not compromise the ability for the FFT to be used for rapid matrix-vector multiplication as it merely introduces a trivial extra multiplication by a diagonal matrix at each iteration.

Premultiplication by $\tilde{I}$ results in a matrix equation with increased sparsity

$$
\begin{pmatrix}
A_{11} & \cdots & \cdots & 0 & \cdots & \cdots & A_{1N} \\
\vdots & \ddots & & \vdots & & & \vdots \\
\vdots & & & 0 & & & \vdots \\
0 & \cdots & 0 & 0 & 0 & \cdots & 0 \\
\vdots & & & 0 & & & \vdots \\
\vdots & & & \vdots & & \ddots & \vdots \\
A_{N1} & \cdots & \cdots & 0 & \cdots & \cdots & A_{NN}
\end{pmatrix}
\begin{pmatrix}
x_1 \\ \vdots \\ \vdots \\ x_{n''} \\ \vdots \\ \vdots \\ x_N
\end{pmatrix}
$$
$$
= 
\begin{pmatrix}
b_1 \\ \vdots \\ \vdots \\ 0 \\ \vdots \\ \vdots \\ b_N
\end{pmatrix}
\tag{13}
$$

where we have again made explicit the structure of the $n''$th column and row for some $n'' \in \mathbb{N}''$. Again, we note that similarly structured rows and columns will be present in the matrix for all $n \in \mathbb{N}''$. We can write this reduced matrix equation as

$$
\tilde{A}\tilde{x} = \tilde{b}
\tag{14}
$$

with $\tilde{A} = \tilde{I}A$ and $\tilde{b} = \tilde{I}b$. Although the condition number of $\tilde{A}$ is technically infinite, we note that $\tilde{x}_n = x_n$, $\forall n \in \mathbb{N}'$; that is, pre-multiplication by the $\tilde{I}$ matrix will not affect the solution at the locations where contrast exists. The reason for this is that the equations involving unknowns at locations where no contrast is present are effectively ignored, as the CG algorithm will not generate search directions in the null space that the pre-multiplication by $\tilde{I}$ introduces into $\tilde{A}$. The infinite condition number merely reflects the fact that the CG scheme will be unable to update the initial values chosen for $x_i$, $i \in \mathbb{N}''$. We stress that this does not present a problem, as any starting value for components $x_i$, $i \in \mathbb{N}''$ trivially satisfy the relevant row equation in (13) and as such will not contribute to the error.

Although it seems reasonable to remove the influence of these unknowns from the CG minimisation procedure, it is not fully clear what effect such a removal will have on its convergence. The CG procedure applied to (13) should converge exactly in at most $N'$ steps, where $N'$ is the number of non-zero eigenvalues of $\tilde{A}$ and is less than the at most $N$ steps required for an exact CG solution of (11). However, in practice, we do not seek an exact solution, preferring to terminate the procedure once an acceptable threshold error has been reached. It is not obvious that the convergence of the CG applied to (13) should be guaranteed to be better at each step of the iterative process than that of the CG applied to (11) and thus that an acceptable threshold will be reached more quickly by use of the reduced operator. Numerical results presented in [10] suggest that this indeed is the case, but no analytic proof was offered. The following section presents such an analysis of the convergence properties of the reduced operator and shows that the CG will indeed converge faster at each step.

*IET Sci. Meas. Technol., Vol. 1, No. 1, January 2007*

59

## 4 Convergence analysis for reduced forward operator

For the purposes of analysing the convergence properties of the reduced operator, we remove the zero rows and zero columns of the matrix equation (13). The equation can be effectively reduced to one involving a matrix $\hat{A}$ of size $N' \times N'$, given by

$$\hat{A} = \begin{pmatrix} \hat{A}_{11} & \cdots & \hat{A}_{N'1} \\ \vdots & \ddots & \vdots \\ \hat{A}_{1N'} & \cdots & \hat{A}_{N'N'} \end{pmatrix} \quad (15)$$

Application of the same reductions to $\tilde{x}$ and $\tilde{b}$ results in the vectors $\hat{x}$ and $\hat{b}$ with the corresponding matrix equation

$$\hat{A}\hat{x} = \hat{b} \quad (16)$$

The analysis of the behaviour of the CG algorithm as applied to (13) is identical to its analysis when applied to (16), and we will use the reduced operator of (16) for our convergence analysis. We stress that (13) is preferred in practical implementation as it implicitly retains the uniform lattice structure that permits the retention of the FFT for matrix-vector multiplications. Solving (16) exactly with a CG will only require at most $N'$ iterations, with $N' < N$. What is not clear is how quickly the algorithm will converge at each step. If we apply the ideas of Section 2, we know that this depends on the ratio $\hat{R}$ given by

$$\hat{R} = \frac{\sigma_1^2(\hat{A})}{\sigma_{N'}^2(\hat{A})} \quad (17)$$

Note that the denominator equals $\sigma_{N'}^2(\hat{A})$, as the matrix of interest is now of size $N' \times N'$, with $N' < N$. Let $R$ be the square of the condition number of the original system, as given by (10). Next, we will prove that $R \geq \hat{R}$ and thus that the reduced matrix has a more clustered set of eigenvalues and that its use leads to a more rapid convergence of the CG algorithm.

It is clear that matrix $\hat{A}$ is a subset of the $N' \times N$-dimensional matrix $B$, where $B$ is

$$B = \begin{pmatrix} \hat{A}_{11} & \cdots & \hat{A}_{1N'} & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & & \vdots \\ \hat{A}_{N'1} & \cdots & \hat{A}_{N'N'} & 0 & \cdots & 0 \end{pmatrix} \quad (18)$$

Therefore the singular values of $\hat{A}$ are the same as the singular values of $B$ (Lemma 4 in [14]), and, hence

$$\sigma_1(B) = \sigma_1(\hat{A}) \quad (19)$$

$$\sigma_{N'}(B) = \sigma_{N'}(\hat{A}) \quad (20)$$

As $B$ is a subset of $A$, obtained by the removal of $(N - N')$ rows and shifting of the columns with only zero entries towards the end of the matrix (this shifting of the columns has no effect on the singular values of $A$), the following is true (Lemma 3 in [14]):

$$\sigma_1(A) \geq \sigma_1(B) = \sigma_1(\hat{A}) \quad (21)$$

$$\sigma_N(A) \leq \sigma_{N'}(B) = \sigma_{N'}(\hat{A}) \quad (22)$$

Therefore the following holds for the ratios $R$ and $\hat{R}$:

$$R = \frac{\sigma_1^2(A)}{\sigma_N^2(A)} \geq \frac{\sigma_1^2(\hat{A})}{\sigma_{N'}^2(\hat{A})} = \hat{R} \quad (23)$$

which shows that the reduced matrix has a more compact eigenvalue spread and thus will lead to a more rapid convergence when used in a CG scheme.

## 5 Results

Testing of the presented reduced forward operator was carried out for a spatial domain of size $3\frac{1}{2}\lambda \times 3\frac{1}{2}\lambda \times (1/2)\lambda$, with cubic spatial elements of volume $(\lambda/4)^3$, with $\lambda$ the wavelength of the 1.0 GHz monochromatic wave field. The spatial domain with medium parameters similar to those of breast tissue encloses one or three tumours with dimensions $\lambda \times \lambda \times (1/2)\lambda$; see Fig. 2. The relative permittivity $\epsilon_r$ of the background medium equals $\epsilon_r^{bg} = 18.23$ and that of the object $\epsilon_r^{obj} = 89$, whereas the conductivity equals $\sigma^{bg} = 0.068$ S m$^{-1}$ and $\sigma^{obj} = 1.3$ S m$^{-1}$, respectively.

For both configurations, the matrices $A$ and $\hat{A}$ have been computed. In Figs. 3a and c, the presence of non-zero entries in $A$ is indicated by black points in the image. It clearly shows the sparsity $S$ of the $A$ matrix, where $S$ is defined as

$$S = 100 \times \left(1 - \frac{N_{\text{non-zero}}}{N^2}\right) \quad (24)$$

Hence, in the presence of one object, $S = 96\%$, and, for three objects, $S = 90\%$. Figs. 3b and d show $|\hat{A}|$ for the same systems, obtained by removal of the rows/columns associated with the dummy unknowns. The condition number of all four matrices shown in Figs. 3a−d equals $1.39 \times 10^8$, 31.07, $4.4 \times 10^8$ and 36.97, respectively. We note the significant reduction in condition number of the reduced system $\hat{A}$ when compared with the unreduced system $A$. This is reflected in the greatly improved convergence rates observed and described later in this section. For this configuration, both the full and the reduced forward problems were solved with the CG scheme shown in Table 1, with Polak−Ribière update directions. Note that the gradient divergence present in the Greens operator is computed numerically by the mid-point rule. Consequently, the values of the vector potential, that
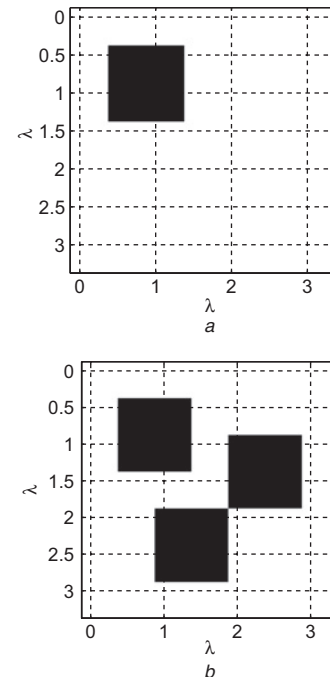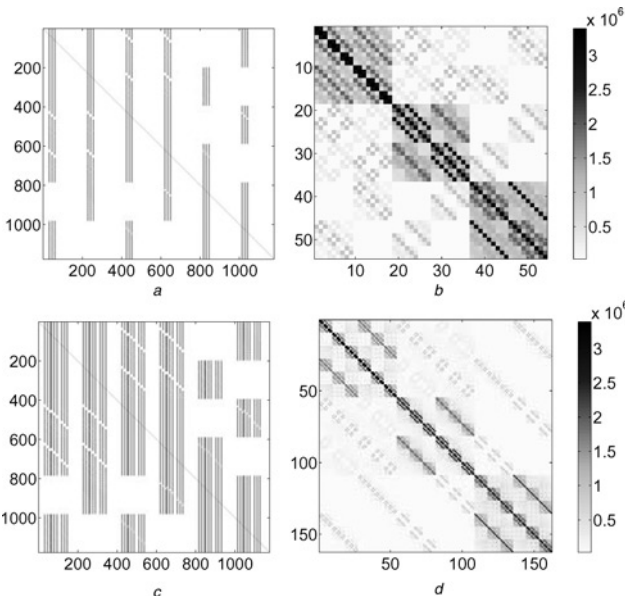


**Fig. 2** *Cross-sections of spatial domain $\mathbb{D}$ enclosing tumour(s)*

60

*IET Sci. Meas. Technol., Vol. 1, No. 1, January 2007*

**Fig. 3** *Presence of non-zero elements in matrix $A$ and absolute values of matrix $\hat{A}$, for one and three tumours*

*a* Matrix $A$ for one tumour
*b* Matrix $\hat{A}$ for one tumour
*c* Matrix $A$ for three tumours
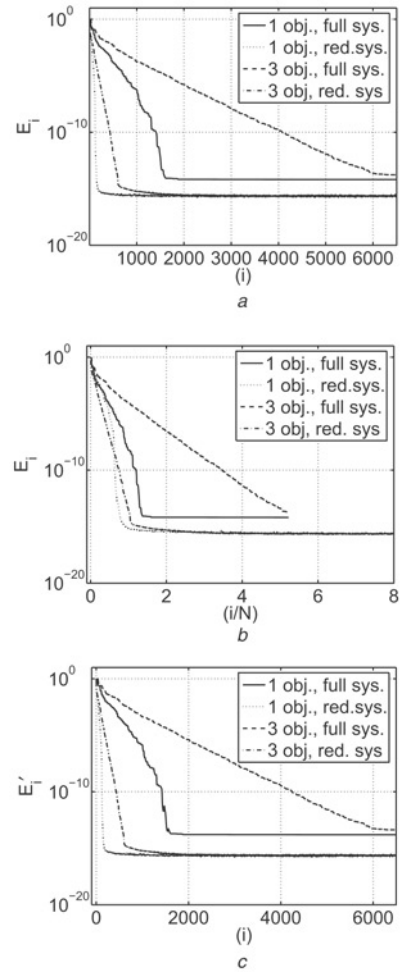*d* Matrix $\hat{A}$ for three tumours

is the convolution of the Greens scalar function with the contrast sources, at the edge of the computational domain (as opposed to the edge of the scatterer) are only used to compute the gradient divergence and, hence, the electric wave field at points just inside the boundary. Consequently, to compute the electric field in $14 \times 14 \times 2$ points, we computed the vector potential in $16 \times 16 \times 4$ points.

To illustrate the improvement in the convergence rate, we computed the normalised error $E_i$ after each iteration,

$$E_i = \frac{\|Ax_i - b\|_2}{\|b\|_2} \qquad (25)$$

where $x_i$ is the approximate solution at the $i$th iteration step, and $\|\cdot\|_2$ denotes the $L_2$-norm. In Fig. 4*a*, this error functional is plotted as a function of the number of iterations, where the solid line refers to a functional based on $Ax = b$ with only one object present, and the broken line is used when three objects are present, whereas the dotted line refers to the functional based on $\hat{A}\hat{x} = \hat{b}$, with only one object present, and the dashed-dotted line is the case where three objects are present. These results show that a clear reduction in the number of iterations needed is obtained by employment of the knowledge about the locations of zero contrast to form a reduced operator. In the presence of three objects, only a tenth of the number of iterations is required to obtain the same accuracy. Note also the effect of the finite numerical precision of the computer. This is indicated by the fact that the error functionals flatten out after a reduction of 16 times the order of magnitude and do not become exactly zero.

The same limitation in numerical precision is observed in Fig. 4*b*. Here, the same error functionals are shown, but now as a function of the number of iterations divided by the number of unknowns, that is $i/N$ and $i/N'$, respectively. Based on the knowledge that the number of iterations required to solve the problem exactly equals the number of independent eigenvalues, we would expect that all lines would cross the line $E_i \simeq 10^{-16}$ at $i/N \simeq 1$. This is clearly not the case for the situation where the full system with

**Fig. 4** *Error functionals $E_i$ and $E_i'$*

—— Functionals based on $Ax = b$ with one object present
– – – Functionals based on $Ax = b$ with three objects present
······ Functionals based on $\hat{A}\hat{x} = \hat{b}$ with one object present
-·-·- Functionals based on $\hat{A}\hat{x} = \hat{b}$ with three objects present

three objects is solved, where the poor condition number has exacerbated the problem of finite precision and led to extremely slow convergence. The reduced system for the three-object problem shows no such problem, however.

Finally, the error functionals shown in Fig. 4*c* are defined as

$$E_i' = \frac{\|Ax_i - b\|_{2,\chi}}{\|b\|_{2,\chi}} \qquad (26)$$

with $\|x\|_{2,\chi}$ being the $L_2$-norm of a vector $x$, where locations with zero contrast are excluded from the spatial domain and, hence, the norm. Comparison of these error functionals with the error functionals presented in Fig. 4*a*, shows that both sets of error functionals are almost identical. Hence, the largest contribution to the error functional originates from the fields at the locations of non-zero contrast.

## 6 Conclusions

A reduced operator has been presented that significantly reduces the number of steps needed for acceptable convergence of the CG-FFT when applied to electromagnetic scattering problems. The reduced operator hides the influence of the unknowns located at points of zero contrast from the CG minimisation procedure, thereby lowering the condition number of the system and greatly improving the convergence rate. Numerical simulations support the analytical findings.

# 7 References

1 de Hoop, A.T.: 'Handbook of radiation and scattering of waves: acoustic waves in fluids, elastic waves in solids, electromagnetic waves' (Academic Press, London, 1995)

2 Zhang, Z.Q., and Liu, Q.H.: 'A volume adaptive integral equation method (VAIM) for 3D inhomogeneous objects', *IEEE Antennas Wirel. Propag. Lett.*, 2002, **1**, (1), pp. 102–105

3 Nie, X.C., Yuan, N., Li, L.W., Gan, Y.B., and Yeo, T.S.: 'A fast combined field volume integral equation solution to EM scattering by 3D dielectric objects of arbitrary permittivity and permeability', *IEEE Trans. Antennas Propag.*, 2006, **54**, (3), pp. 961–969

4 Lu, C.C.: 'A fast algorithm based on volume integral equation for analysis of arbitrarily shaped dielectric radomes', *IEEE Trans. Antennas Propag.*, 2003, **51**, (3), pp. 606–612

5 Zhang, Z.Q., Liu, Q.H., Xiao, C., Ward, E., Ybarra, G., and Joines, W.T.: 'Microwave breast imaging: 3D forward scattering simulation', *IEEE Trans. Biomed. Eng.*, 2003, **50**, (10), pp. 1180–1189

6 Zhang, Z.Q., Liu, Q.H., and Xu, X.M.: 'RCS computation of large inhomogeneous objects using a fast integral equation solver', *IEEE Trans. Antennas Propag.*, 2003, **51**, (3), pp. 613–618

7 Peterson, A.F., Ray, S.L., Chan, C.H., and Mittra, R.: 'Numerical implementations of the conjugate gradient method and the CG-FFT for electromagnetic scattering' in Sarkar, T.K. (Ed.): ''PIER 5, Application of conjugate gradient method to electromagnetics and signal analysis', Elsevier, 1991, pp. 241–300

8 Sarkar, T.K., Arvas, E., and Rao, S.: 'Application of FFT and the conjugate gradient method for the solution of electromagnetic radiation from electrically large and small conducting bodies', *IEEE Trans. Antennas Propag.*, 1986, **34**, (5), pp. 635–640

9 Zwamborn, A.P.M., and van den Berg, P.M.: 'A weak form of the conjugate gradient FFT method for plate problems', *IEEE Trans. Antennas Propag.*, 1991, **39**, (2), pp. 224–228

10 van Dongen, K.W.A., Brennan, C., and Wright, W.M.D.: 'A reduced forward operator for acoustic scattering problems'. Proc. IEE Irish Signals and Systems Conf., Dublin, Ireland, 1–2 September 2005, pp. 294–299

11 Kleinman, R.E., and van den Berg, P.M.: 'Iterative methods for solving integral equations' in Sarkar, T.K. (Ed.): 'PIER 5, Application of conjugate gradient method to electromagnetics and signal analysis', Elsevier, 1991, pp. 67–102

12 Peterson, A.F., Ray, S.L., and Mittra, R.: 'Computational methods for electromagnetics' (Wiley-IEEE Press, 1997)

13 Golub, G.H., and van Loan, C.F.: 'Matrix computations' (John Hopkins University Press, 3rd edn.)

14 Mathias, R.: 'Two theorems on singular values and eigenvalues', *Am. Math. Monthly*, 1990, **97**, (1), pp. 47–50

62

*IET Sci. Meas. Technol., Vol. 1, No. 1, January 2007*