

First Principles Order N Calculations on Very Large Systems

M.J.Gillan, D.R.Bowler, C.M.Goringe,* and E.H.Hernández†
Physics Department, Keele University, Keele, Staffordshire, ST5 5BG, UK

We describe recent progress in developing practical first principles methods for which the computer effort is proportional to the number of atoms: linear scaling or $\mathcal{O}(N)$ methods. It is shown that the locality property of the density matrix gives a general framework for constructing such methods. We then outline some of the main technical problems which must be solved in order to develop a practical $\mathcal{O}(N)$ method based on density functional theory and the pseudopotential method. Recent progress in solving these problems is presented, and we show that the spatial cut-off distances needed to achieve good accuracy are small enough to make the calculations feasible. Parallel implementation of the $\mathcal{O}(N)$ methods in the CONQUEST code is outlined, and it is shown that the code exhibits excellent linear-scaling behaviour on test systems of several thousand atoms. It is pointed out that the most important remaining problem concerns the optimal strategy for seeking the ground state. It is argued that there are three different mechanisms of ill-conditioning which cause present search methods to be inefficient, and some partial solutions are suggested.

Published as “The Physics of Complex Liquids”, Proceedings of the International Symposium, 10-12 November 1997, Nagoya, Japan, ed. F.Yonezawa, K.Tsuji, K.Kaji, M.Doï and T.Fujiwara (World Scientific, 1998)

I. INTRODUCTION

First principles simulation based on density functional theory (DFT) and the pseudopotential method is now in widespread use in many areas of physics and chemistry. It is already playing an important role in the study of complex liquids, and there have been many first principles simulations of liquids including silicon¹, selenium² and alloys³. However, current calculations are limited to systems of no more than a few hundred atoms, because the computer time needed increases at least as the square of the number of atoms N . There has recently been an intensive effort⁴⁻¹⁴ to overcome this limitation by developing techniques for which the computer time is linearly proportional to N . We report here recent progress with these so-called $\mathcal{O}(N)$ or linear scaling techniques.

The reason for the N^2 dependence is easy to understand. The usual DFT techniques work with the eigenfunctions $\psi_i(\mathbf{r})$ of the Kohn-Sham Hamiltonian, which extend over the entire volume of the system. The number of ψ_i functions is proportional to N , and the amount of information in *each* ψ_i is proportional to the volume of the system (which is proportional to N), so that the total number of variables needed to describe the electrons increases as N^2 . For *very* large systems, the dependence becomes N^3 , because conventional methods require the calculation of the scalar products $\langle \psi_i | \psi_j \rangle$:

$$\langle \psi_i | \psi_j \rangle = \int d\mathbf{r} \psi_i(\mathbf{r}) \psi_j(\mathbf{r})^*, \quad (1)$$

and each of these N^2 quantities needs a computational time proportional to the volume. These bad dependencies on N arise because the ψ_i orbitals extend over the

whole system. But it has been recognised for a long time that the quantum state of the electrons in condensed matter can be described in terms of localised functions, and this insight provides the key to developing $\mathcal{O}(N)$ methods.

II. GENERAL FRAMEWORK FOR $\mathcal{O}(N)$

Several authors have stressed that the fundamental reason for the existence of $\mathcal{O}(N)$ techniques can be seen in the properties of the two particle density matrix¹³, $\rho(\mathbf{r}, \mathbf{r}')$. This can be defined in terms of the Kohn-Sham eigenfunctions as:

$$\rho(\mathbf{r}, \mathbf{r}') = \sum_i f_i \psi_i(\mathbf{r}) \psi_i(\mathbf{r}')^*, \quad (2)$$

where f_i is the occupation number of orbital i . But instead of considering the ψ_i as the primary quantities in terms of which $\rho(\mathbf{r}, \mathbf{r}')$ is defined, we regard $\rho(\mathbf{r}, \mathbf{r}')$ itself as the primary quantity²⁷. It is known that DFT can be formulated perfectly well in terms of $\rho(\mathbf{r}, \mathbf{r}')$ without any explicit mention of wavefunctions. The total energy E_{tot} can be expressed explicitly in terms of ρ , and the ground state is obtained by minimising the functional $E_{\text{tot}}[\rho(\mathbf{r}, \mathbf{r}')]$, subject to the conditions: (i) ρ is Hermitian; (ii) ρ is idempotent (i.e. its eigenvalues are all 0 or 1); (iii) ρ yields the correct number of electrons, N_{el} .

$$N_{el} = 2 \int d\mathbf{r} \rho(\mathbf{r}, \mathbf{r}). \quad (3)$$

The key property of $\rho(\mathbf{r}, \mathbf{r}')$ is that it decays to zero as $|\mathbf{r} - \mathbf{r}'| \rightarrow \infty$. The fundamental reason for this decay is the loss of quantum phase coherence between distant points. This means that the amount of information in $\rho(\mathbf{r}, \mathbf{r}')$ is proportional to N , and so much of the *apparent* information contained in a wavefunction description must be redundant.

Given these properties of $\rho(\mathbf{r}, \mathbf{r}')$, a general $\mathcal{O}(N)$ approach to the determination of the ground state is to

minimise E_{tot} with respect to ρ with the additional constraint:

$$\rho(\mathbf{r}, \mathbf{r}') = 0, \quad |\mathbf{r} - \mathbf{r}'| > R_c, \quad (4)$$

where R_c is some cut-off radius. This will lead to an upper bound to the true ground state energy, which will converge to the true value as R_c is increased. In developing practical methods, one cannot proceed exactly like this, as the ρ depends on two vector positions. Instead, we introduce the additional approximation that ρ is *separable*, i.e. that it can be expressed in the form:

$$\rho(\mathbf{r}, \mathbf{r}') = \sum_{i\alpha, j\beta} \phi_{i\alpha}(\mathbf{r}) K_{i\alpha, j\beta} \phi_{j\beta}(\mathbf{r}'); \quad (5)$$

this is equivalent to demanding that ρ have only a finite number of non-zero eigenvalues. The functions $\phi_{i\alpha}(\mathbf{r})$ are known as ‘support functions’, or sometimes as ‘localised orbitals’. The notation indicates that $\phi_{i\alpha}$ is the α^{th} support function on atom i . The matrix $K_{i\alpha, j\beta}$ is the density matrix in the representation of $\phi_{i\alpha}(\mathbf{r})$.

In the procedure we have proposed, the localisation of $\rho(\mathbf{r}, \mathbf{r}')$ expressed in equation 4 is replaced by the condition that the support functions be non-zero only within limited regions, which in practice are spheres of radius R_{reg} centred on the atoms. The $\phi_{i\alpha}$ can then be seen as closely related to the atomic-like orbitals that would be used in a tight binding description. The localisation of $\rho(\mathbf{r}, \mathbf{r}')$ also requires that the matrix $K_{i\alpha, j\beta}$ be subject to a spatial cutoff, so that $K_{i\alpha, j\beta} = 0$ when the separation of the atoms i and j exceeds some value.

The matrix $K_{i\alpha, j\beta}$ cannot be allowed to vary freely, since $\rho(\mathbf{r}, \mathbf{r}')$ is required to be idempotent, or at least ‘weakly idempotent’ (i.e. its eigenvalues must lie in the interval [0,1]). There are several ways of enforcing this condition, one of which is equivalent to the $\mathcal{O}(N)$ DFT scheme proposed by Mauri, Galli and Car⁴. The method we have proposed¹³ is based on the ‘purification’ technique of McWeeny¹⁵, recently used in tight binding calculations by Li, Nunes and Vanderbilt⁷. As explained elsewhere¹³, this requires the matrix K to be written as:

$$K = 3LSL - 2LSLSL \quad (6)$$

where $L_{i\alpha, j\beta}$ is an ‘auxiliary’ density matrix and $S_{i\alpha, j\beta}$ is the overlap matrix:

$$S_{i\alpha, j\beta} = \int d\mathbf{r} \phi_{i\alpha}(\mathbf{r}) \phi_{j\beta}(\mathbf{r}). \quad (7)$$

The localisation of $\rho(\mathbf{r}, \mathbf{r}')$ is then imposed as a spatial cutoff on $L_{i\alpha, j\beta}$:

$$L_{i\alpha, j\beta} = 0, \quad |\mathbf{R}_{i\alpha} - \mathbf{R}_{j\beta}| > R_L, \quad (8)$$

where R_i is the position of atom i and R_L is a cutoff radius.

We summarise the overall scheme: The ground state energy and density matrix of the system are determined

by minimising E_{tot} with respect to the $\phi_{i\alpha}(\mathbf{r})$ functions and the auxiliary matrix $L_{i\alpha, j\beta}$, subject to the spatial cutoffs R_{reg} and R_L . This gives an upper bound to the true E_{tot} , which is expected to go to the true value as R_{reg} and R_L are increased.

III. PRACTICAL QUESTIONS

The above general scheme can clearly be implemented in many ways. These are some of the practical questions that should be asked:

- What is the best way of representing the support functions $\phi_{i\alpha}(\mathbf{r})$, i.e. what basis functions should be used ?
- How should we search for the ground state, i.e. what strategy should be used to vary the $\phi_{i\alpha}(\mathbf{r})$ and the $L_{i\alpha, j\beta}$ in minimising E_{tot} ?
- How well does the method work in practice: what cutoffs R_{reg} and R_L are needed to achieve good convergence to the true ground state, and for what sizes of system does the $\mathcal{O}(N)$ scheme become more efficient than conventional methods ?
- How can the forces on the atoms be calculated ? This is a crucial question if the scheme is to be of any practical use !
- What is the right way to implement the scheme on parallel computers ? This is an important question, because $\mathcal{O}(N)$ schemes will show their true power for large systems, and this will usually require parallel machines.

Recently, answers to some of these questions have begun to emerge. In our own work, much has been learnt through the writing of the parallel $\mathcal{O}(N)$ code called CONQUEST (Concurrent Order N QUantum Electronic Simulation Technique)¹⁶. The next section summarises some of the findings.

IV. SOME PRACTICAL ANSWERS

A. Representation of the support functions

In considering this question, it is helpful to remember the lessons that have been learnt from the use of plane wave basis sets in conventional pseudopotential calculations. Two of the major advantages of plane waves are that: (i) Systematic convergence of the total energy with respect to basis set completeness is achieved by increasing a single parameter – the plane wave cutoff energy, E_{cut} ; (ii) they are free of bias, i.e. they are completely flexible, and no judgement has to be made about the kind of chemical bonding that will occur. If possible, the basis

set used in $\mathcal{O}(N)$ calculations should retain these advantages. In any case, we must certainly aim to achieve the accuracy normally expected of current plane wave calculations.

One possibility is to use the spherical analogue of plane waves, namely the product of spherical Bessel functions $j_l(r)$ and spherical harmonics $Y_l^m(\hat{\mathbf{r}})$, in each support region. The advantages of this representation have been discussed by Haynes and Payne¹⁷, but practical results have not yet been reported. An alternative is simply to represent the $\phi_{i\alpha}(\mathbf{r})$ by their values on a grid and to calculate matrix elements of the kinetic energy by finite differences. This method is well established to be practicable in non- $\mathcal{O}(N)$ calculations¹⁸ and has been studied in the $\mathcal{O}(N)$ context by us¹⁹ and very recently by Hoshi and Fujiwara²⁰. A third method is the B-spline basis set used in the CONQUEST code. In this scheme, the basis functions $\chi_s(\mathbf{r})$ in the expansion:

$$\phi_{i\alpha}(\mathbf{r}) = \sum_s b_{i\alpha s} \chi_s(\mathbf{r}) \quad (9)$$

are piecewise polynomial functions strictly localised on the points of a grid which is rigidly attached to each atom. Details of the B-spline scheme are reported elsewhere¹³.

B. Searching for the ground state

The only practical work on this question that we are aware of is our own work on the CONQUEST code. Basically, we use the standard conjugate gradients technique to minimise E_{tot} with respect to the basis coefficients $b_{i\alpha s}$ and the matrix elements $L_{i\alpha, j\beta}$. The gradients $\partial E_{\text{tot}}/\partial b_{i\alpha s}$ and $\partial E_{\text{tot}}/\partial L_{i\alpha, j\beta}$ needed to do this are readily calculated. Since two very different kinds of variables are involved, the search is organised as a double loop. In the inner loop, E_{tot} is minimised with respect to $L_{i\alpha, j\beta}$, with the support functions held fixed. In the outer loop, the $b_{i\alpha s}$ are varied.

We note that the inner loop is identical to the ground state search in a self-consistent tight binding calculation. It has been stressed recently²¹ that in this context invariance with respect to linear transformations of the $\phi_{i\alpha}(\mathbf{r})$ imposes certain natural constraints, which should be respected. The implication of this is that in a non-orthogonal basis set, the correct metric must be applied to ensure correct gradients (i.e. a fully contravariant gradient $\partial E_{\text{tot}}/\partial L_{i\alpha, j\beta}$ must be added to the contravariant density matrix; this involves applying the metric S^{-1} to the gradient calculated by Nunes and Vanderbilt²²). Our practical experience with the CONQUEST code is that this procedure is markedly more efficient than a naive search along the gradients of $\partial E_{\text{tot}}/\partial L_{i\alpha, j\beta}$.

Although we find that these search methods generally work, we are not satisfied with their efficiency, and we return to this question in Section V.

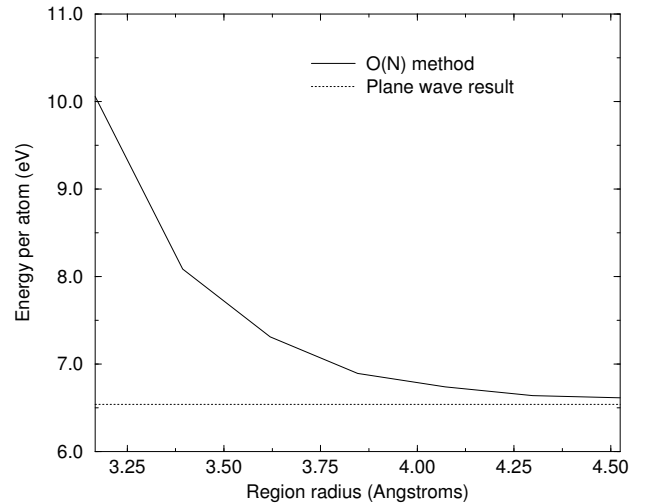


FIG. 1: The total energy per atom versus region radius (R_{reg}) for silicon as calculated by CONQUEST. The value from the plane wave code CASTEP is shown for comparison.

C. Dependence on the spatial cutoffs

We have already reported preliminary results on the dependence of the calculated ground state energy on the cutoffs R_{reg} and R_L , which suggested that accurate results are obtained with quite modest cutoffs. But these were mainly based on a model local pseudopotential and so were not fully realistic. We report here new tests using standard non-local pseudopotentials^{23,24}. The tests are done on perfect crystals of silicon.

In testing the dependence of E_{tot} on R_{reg} , we set R_L equal to infinity, which is equivalent to exact diagonalisation. For comparison, we have also done calculations with the standard plane wave code CASTEP²⁵ using precisely the same pseudopotential and other parameters. Figure 1 shows the calculated total energy as a function of R_{reg} for Si and Ge. The results show that E_{tot} converges to the correct value extremely rapidly once R_{reg} is greater than 4 Å. For this radius, each support region contains 17 neighbouring atoms, and the calculations are perfectly manageable.

Our tests on R_L were done with $R_{\text{reg}} = 2.715$ Å, and the results for Si are shown in Figure 2. Rather accurate convergence to the $R_L = \infty$ value is obtained for $R_L \geq 8$ Å, which again is acceptable. No value is shown for exact diagonalisation because of technical difficulties in performing comparisons.

The conclusion from these tests is that the practical values of the spatial cutoffs needed to achieve good accuracy are encouragingly small. In Section III, we asked for what system size $\mathcal{O}(N)$ becomes more efficient than conventional methods. We believe that it is too soon to answer this question. The main reason is that the efficiency of the ground state search technique in $\mathcal{O}(N)$ is

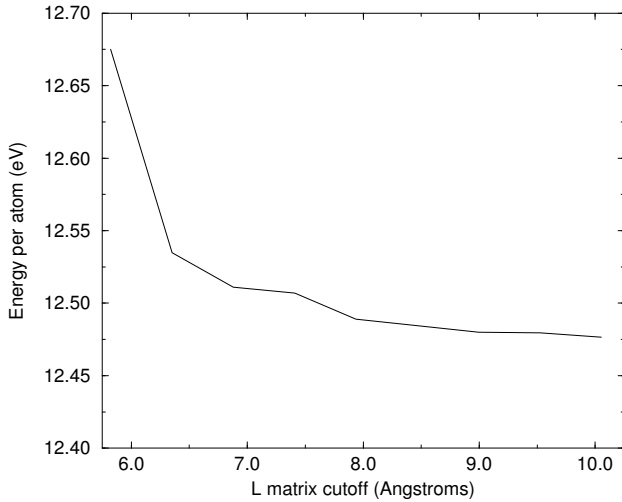


FIG. 2: The total energy per atom versus density matrix cutoff (R_L) as calculated by CONQUEST for silicon. The region radius was 2.715 Å.

still rather poor, as discussed in Section V.

D. Forces on atoms

In the conventional plane wave technique, the force on each atom is simply the Hellmann-Feynman force, i.e. the force exerted by the electrons in the ground state associated with the current ionic positions (plus, of course, the Coulombic interaction between ionic cores). This relies on the fact that the basis set does not depend on the ionic positions. In the $\mathcal{O}(N)$ technique used in CONQUEST, the B-spline basis functions $\chi_s(\mathbf{r})$ move with the ions, and this gives rise to an additional contribution to the force, known as the Pulay contribution²⁶. If the calculation is well converged with respect to basis set completeness, then the Pulay correction is small, but it is nonetheless essential to include it, in order to ensure exact consistency between the total energy and the forces. As will be described in more detail elsewhere, the Pulay contribution is straightforward to calculate. This means that the relaxation of the system to mechanical equilibrium and the generation of time-dependent ionic trajectories will be feasible in $\mathcal{O}(N)$ DFT calculations.

E. Parallel implementation

The essence of $\mathcal{O}(N)$ is that the system can be separated into independent spatial regions. This means that $\mathcal{O}(N)$ is ideally suited to parallel implementation, with each processor being responsible for a set of atoms and/or spatial regions. The way this is done in the parallel CONQUEST code is described in detail elsewhere¹⁶, so we give here only a brief summary.

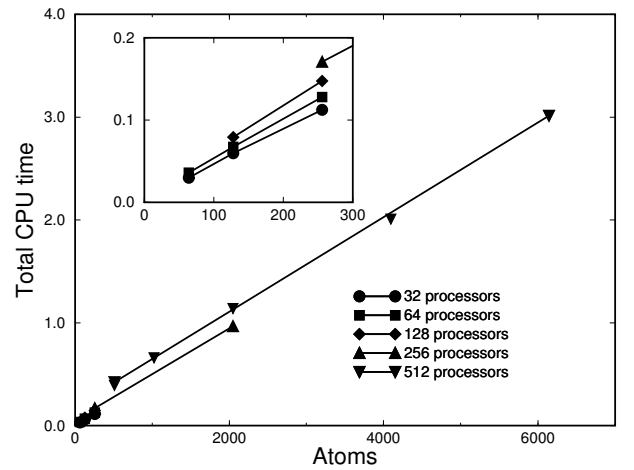


FIG. 3: The total CPU time taken for silicon systems of differing size, using the CONQUEST code on the Cray T3D.

Each processor is given a three-fold responsibility. First, it is in charge of a certain group of atoms. This means that it holds the basis coefficients $b_{i\alpha s}$, the derivatives $\partial E_{\text{tot}}/\partial b_{i\alpha s}$, and also the rows of all matrix elements such as $S_{i\alpha, j\beta}$ corresponding to these atoms. It is responsible for performing the transforms from basis coefficients to grid values $\phi_{i\alpha}(\mathbf{r}_l)$ for these atoms, and also for doing matrix multiplications needed to produce matrix rows associated with its atoms. Second, each processor is in charge of a domain of integration grid points \mathbf{r}_l , and has the job of calculating contributions to matrix elements coming from sums over this domain. It also has responsibility for the electron density and Kohn-Sham potential on its domain of points. Third, the processor is responsible for doing part of the spatial Fourier transforms needed in calculating the Hartree potential. In practice, this means that it deals with a set of columns of grid points in the x, y or z directions. The processors switch between their responsibilities in a synchronised manner, and communication of data between them is needed when this happens.

We have made extensive tests of the scaling properties of the CONQUEST code. It is important to stress that there are two completely different types of scaling. The first concerns the way the CPU time increases as the size of the simulated system increases for a *fixed* number of processors; we call this ‘intrinsic scaling’. The second concerns the way the CPU time changes when a *given* simulated system is treated on varying numbers of processors.

In an implementation of CONQUEST on the Cray T3D, both types of scaling turn out to be excellent. As an illustration of the intrinsic scaling, we show in Figure 3 the total CPU time required per iteration for silicon crystals containing from 64 to 6144 atoms. (Here, total CPU time means the number of processors multiplied by the CPU time per processor.) The results show that for

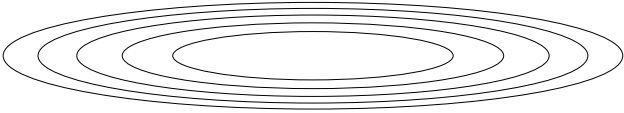


FIG. 4: A function with elongated surfaces of constant f

a fixed number of processors, the total CPU time is almost exactly proportional to the number of atoms, and this confirms the correctness of the underlying theory. We also find that for a given simulated system, the total CPU time increases only very weakly with the number of processors, which means that the fraction of time taken in communications is very small.

V. ILL-CONDITIONING PROBLEMS

In spite of the very encouraging findings summarised above, there remains a question that is still not completely solved: the strategy for finding the ground state. The search procedure outlined in Section IV B is generally successful, but it sometimes requires many iterations, particularly when the region radius R_{reg} is large. This means that practical calculations are rather inefficient. The problem concerns variations of the $\phi_{i\alpha}$, since convergence with respect to the $L_{i\alpha,j\beta}$ matrix in the inner loop is generally rapid. The large number of iterations is needed in the outer loop where the $\phi_{i\alpha}$ are varied.

The cause of the problem is ill-conditioning. This is an extremely generally phenomenon which afflicts minimisation problems in many areas of science, and occurs when the function being minimised has a wide range of curvatures. Suppose we need to locate the minimum of some general function $f(x_1, x_2, \dots, x_N)$ which depends on the set of variables $\{x_i\}$. Let $C_{ij} \equiv \partial^2 f / \partial x_i \partial x_j$ be the curvature matrix (sometimes called the Hessian) evaluated at the minimum. If the eigenvalues λ_n of C_{ij} span a wide range – the ratio between the largest and smallest eigenvalues $\lambda_{\text{max}}/\lambda_{\text{min}}$ is large – then the surfaces of constant f are very elongated (see Figure 4), and conventional techniques such as conjugate gradients become very inefficient. In fact it is known that the number of iterations needed by conjugate gradients is proportional to $(\lambda_{\text{max}}/\lambda_{\text{min}})^{1/2}$.

Although ill-conditioning is a very widespread phenomenon, its causes are specific to the problem at hand. In order to overcome the problem, it is essential to understand these causes. In fact, it has been recognised for many years that conventional first principles calculations can suffer from ill-conditioning, and its causes are already well understood. It turns out that the ill-conditioning we encounter in $\mathcal{O}(N)$ calculations is closely related to that found in conventional calculations, so it will be useful to spend a few moments recalling some well known facts.

In the usual plane wave techniques, the total energy E_{tot} has to be minimised with respect to the Kohn-Sham

orbitals ψ_i , which are represented in a plane wave expansion:

$$\psi_i = \sum_{\mathbf{G}} c_{i\mathbf{G}} \exp(i\mathbf{G} \cdot \mathbf{r}). \quad (10)$$

For many systems, particularly metals, it is common to allow partial occupation of the orbitals, so that orbital ψ_i has occupation number f_i . We then have $E_{\text{tot}} = E_{\text{tot}}(\{c_{i\mathbf{G}}\}, \{f_i\})$. The function E_{tot} has high curvatures associated with variation of $c_{i\mathbf{G}}$ at high wavevector \mathbf{G} . The reason for this is simply that the kinetic energy E_{kin} of the electrons is

$$E_{\text{kin}} = 2 \sum_i f_i \sum_{\mathbf{G}} \frac{\hbar^2 G^2}{2m} |c_{i\mathbf{G}}|^2, \quad (11)$$

so that the curvature is proportional to G^2 . Since this kind of ill-conditioning comes from the variation of curvature with length scale, we refer to this as ‘length scale ill-conditioning’. The problem can be cured by conventional preconditioning methods.

Conventional techniques can also suffer from a second type of ill-conditioning, associated with the invariance of E_{tot} under unitary transformations of the orbitals. If all f_i are zero or unity, the E_{tot} is exactly invariant under transformations:

$$\psi_i \rightarrow \psi'_i = \sum_j U_{ij} \psi_j \quad (12)$$

where U_{ij} is unitary. This invariance implies that some of the eigenvalues of the Hessian vanish. But if the f_i deviate from zero or unity, the exact invariance is broken, and the vanishing eigenvalues of the Hessian acquire small positive values. It is their smallness that causes the ill-conditioning. We refer to this mechanism as ‘superposition ill-conditioning’. In conventional techniques, this is usually cured by a method known as sub-space rotation.

There is yet another cause of ill-conditioning. When variable occupation numbers are employed, orbitals whose energies are well above the Fermi energy will have very small values of f_i . Variations of the corresponding ψ_i will therefore have little effect on E_{tot} , so that the curvatures will again be small. Since orbitals having small f_i are almost redundant, we call this mechanism ‘redundancy ill-conditioning’.

All these three types of ill-conditioning can also cause trouble in $\mathcal{O}(N)$ techniques. It is clear, for example, that the support functions $\phi_{i\alpha}$ can vary on different length scales, so that length scale ill-conditioning is inevitable. This will not cause serious problems, and will be overcome by conventional preconditioning techniques.

Superposition ill-conditioning, associated with linear mixing of the support functions, is more interesting. Two kinds should be distinguished. The first consists of mixing of different $\phi_{i\alpha}$ on the same atom. It is readily shown that this leaves E_{tot} exactly invariant, and cannot cause

trouble. The second consists of mixing of $\phi_{i\alpha}$ on different atoms. Since the $\phi_{i\alpha}$ are constrained to be zero outside their regions, this second kind is, strictly, impossible. However, for large region radii there are variations which respect this constraint, while consisting almost exactly of linear superpositions of $\phi_{i\alpha}$ on different atoms. The small curvatures of E_{tot} arising from these variations give rise to superposition ill-conditioning. Our present belief is that this problem will not be difficult to overcome. The reason is that the characteristic variations responsible for the ill-conditioning can be calculated, and this will make it possible to precondition them.

Finally, we comment on redundancy ill-conditioning. We have noted that in conventional techniques this occurs when the number of orbitals exceeds the sum of the occupation numbers (i.e. half the electron number in spin-paired calculations). An analogous problem will afflict $\mathcal{O}(N)$ when the number of $\phi_{i\alpha}$ is greater than half the electron number. This will not always happen, because for many systems the number of $\phi_{i\alpha}$ can be taken equal to half the electron number. But for other systems it will be essential, or at least desirable, to work with a larger number of $\phi_{i\alpha}$. Systems consisting of group IV elements are a case in point, because it will generally be natural to take four $\phi_{i\alpha}$ on each atom, one corresponding to the valence s -orbital and the other three to p -orbitals. Once again, we believe that preconditioning will allow us to overcome this problem, but detailed techniques have yet to be formulated.

VI. PROSPECTS

The developments presented here give reason for great optimism about the future potential of $\mathcal{O}(N)$ DFT tech-

niques. We have shown how the properties of the density matrix allow one to give a very general framework for constructing such techniques. The detailed methods we have implemented in the CONQUEST code represent only one possible way of doing this, and other ways will need to be investigated. We have pointed to a number of technical problems that must be overcome in constructing practical $\mathcal{O}(N)$ techniques, and we have shown that solutions to most of these problems are now available. However, some of these may only be interim solutions. We believe, for example, that the question of how best to represent the support functions will need considerable further investigation before any consensus will be reached. The same goes for parallel implementation. We have outlined one way of doing this, and have shown that this works well for systems of several thousand atoms. However, we have done this on only one kind of machine (the Cray T3D), and it may well be that implementation on other machines (e.g. vector parallel machines) will raise new questions. Finally, we have pointed out that there are unsolved questions about the right way to search for the ground state. The ill-conditioning problems that we have described will need deeper study.

Perhaps the most important conclusion is that $\mathcal{O}(N)$ DFT calculations definitely *work*. The spatial cutoff distances required are small enough to make the calculations perfectly feasible. Moreover, the $\mathcal{O}(N)$ behaviour of the calculations is actually demonstrated in practice. Encouraged by these results, our group is now working towards the application of these techniques to complex large-scale problems, including nanostructures on semiconductor surfaces.

-
- * Present address:Key Center for Microscopy and Microanalysis, Madsen Building (F09), University of Sydney, NSW 2006.
- † Present address:Fisica Teorica, Facultad de Ciencias, Universidad de Valladolid, Real de Burgos S/N, Valladolid 47011, Spain.
- ¹ I. Stich, R. Car and M. Parrinello, Phys. Rev. Lett. **63**, 2240 (1989).
- ² F. Kirchhoff, M. J. Gillan, J. M. Holender, G. Kresse and J. Hafner, J. Phys.: Condens. Matter **8**, 9353 (1996).
- ³ F. Kichhoff, J. M. Holender and M. J. Gillan, Phys. Rev. B **54**, 190 (1996).
- ⁴ F.Mauri, G.Galli and R.Car, Phys. Rev. B **47**, 9973 (1993).
- ⁵ F.Mauri and G.Galli, Phys. Rev. B **50**, 4316 (1994).
- ⁶ J.Kim, F.Mauri and G.Galli, Phys. Rev. B **52**, 1640 (1995).
- ⁷ X.-P.Li, R.W.Nunes and D.Vanderbilt, Phys. Rev. B **47**, 10891 (1994).
- ⁸ W.Kohn, Int. J. Quant. Chem. **56**, 229 (1995).
- ⁹ W.Kohn, Phys. Rev. Lett. **76**, 3168 (1996).
- ¹⁰ P.Ordejón, D.A.Drabold, R.M.Martin and M.P.Grumbach, Phys. Rev. B **51**, 1456 (1995).
- ¹¹ P.Ordejón, E.Artacho and J.M.Soler, Phys. Rev. B **53**, 10441 (1996).
- ¹² S.Goedecker and L.Colombo, Phys. Rev. Lett. **73**, 122 (1994).
- ¹³ E.Hernández, M.J.Gillan and C.M.Goringe, Phys. Rev. B **53**, 7147 (1996).
- ¹⁴ D.R.Bowler, M.Aoki, C.M.Goringe, A.P.Horsfield and D.G.Pettifor, Modell. Simul. in Mater. Sci. Eng. **5**, 199 (1997).
- ¹⁵ R.McWeeny, Rev. Mod. Phys. **32**, 335 (1960).
- ¹⁶ C.M.Goringe, E.Hernández, M.J.Gillan and I.J.Bush, Comp. Phys. Commun. **102**, 1 (1997).
- ¹⁷ P.D.Haynes and M.C.Payne, Comp. Phys. Commun. **102**, 17 (1997).
- ¹⁸ J.R.Chelikowsky, N.Troullier and Y.Saad, Phys. Rev. Lett. **72**, 1240 (1994); E.L.Briggs, D.J.Sullivan and J.Bernholc, Phys. Rev. B **54**, 14326 (1996).
- ¹⁹ E.Hernández and M.J.Gillan, Phys. Rev. B **51**, 10157 (1995).
- ²⁰ T.Hoshi, M.Arai and T.Fujiwara, Phys. Rev. B **52**, R5459 (1995); T.Hoshi and T.Fujiwara, unpublished (1997).

- ²¹ C.A.White, P.Maslen, M.S.Lee and M.Head-Gordon, Chem. Phys. Lett. **276**, 133 (1997).
- ²² R.W.Nunes and D.Vanderbilt, Phys. Rev. B **50**, 17611 (1994).
- ²³ G.Kerker, J. Phys. C **13**, 189 (1980).
- ²⁴ L.Kleinman and D.M.Bylander, Phys. Rev. Lett. **4**, 1425 (1982).
- ²⁵ M.C.Payne, M.P.Teter, D.C.Allan, T.A.Arias and J.D.Joannopoulos, Rev. Mod. Phys. **64**, 1045 (1992).
- ²⁶ P.Pulay, Mol. Phys. **17**, 197 (1969); M.Scheffler, J.P.Vigneron and G.B.Bachelet, Phys. Rev. B **31**, 6541 (1985).
- ²⁷ We are interested here in the pseudopotential approach to DFT, so that the ψ_i are actually the valence pseudo-wavefunctions