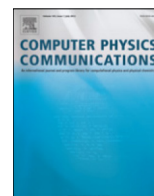




ELSEVIER

Contents lists available at ScienceDirect

Computer Physics Communications

journal homepage: www.elsevier.com/locate/cpc

Parallel 3-dim fast Fourier transforms with load balancing of the plane waves

Xingyu Gao^{a,b,c}, Zeyao Mo^{a,b,c}, Jun Fang^{b,c}, Haifeng Song^{a,b,c}, Han Wang^{b,c,*}^a Laboratory of Computational Physics, Huayuan Road 6, Beijing 100088, PR China^b Institute of Applied Physics and Computational Mathematics, Fenghao East Road 2, Beijing 100094, PR China^c CAEP Software Center for High Performance Numerical Simulation, Huayuan Road 6, Beijing 100088, PR China

ARTICLE INFO

Article history:

Received 29 January 2016

Received in revised form

23 May 2016

Accepted 2 July 2016

Available online xxxx

Keywords:

First-principles calculation

Kohn–Sham equation

Plane wave

FFT

Load balancing

ABSTRACT

The plane wave method is most widely used for solving the Kohn–Sham equations in first-principles materials science computations. In this procedure, the three-dimensional (3-dim) trial wave functions' fast Fourier transform (FFT) is a regular operation and one of the most demanding algorithms in terms of the scalability on a parallel machine. We propose a new partitioning algorithm for the 3-dim FFT grid to accomplish the trade-off between the communication overhead and load balancing of the plane waves. It is shown by qualitative analysis and numerical results that our approach could scale the plane wave first-principles calculations up to more nodes.

© 2016 Elsevier B.V. All rights reserved.

1. Introduction

In the context of Density Functional Theory (DFT), solving the Kohn–Sham equation is the most time-consuming part of the first-principles materials science computations [1–3]. The plane wave method, which is a widely used numerical approach [4], could lead to a large-scale dense algebraic eigenvalue problem. This problem is usually solved by iterative diagonalization methods such as Davidson's [5], RMM-DIIS [3], LOBPCG [6], Chebyshev polynomial filtering subspace iteration [7], etc. The elementary operation of the iteration methods is the matrix–vector multiplication. Since the large-scale dense matrix is not suitable for explicit assembly, we realize the matrix–vector multiplication by applying the Hamiltonian operator on trial wave functions. The local term of the effective potential is one part of the Hamiltonian operator. In order to compute its action in a lower time complexity, we perform 3-dim FFT twice on one trial wave function in each matrix–vector multiplication.

There are three features to make the trial wave function's FFT one of the most demanding algorithms to scale on a parallel machine. The first is the moderate sized FFT grid rather than a large

one. The ratio of computation to communication of the parallel 3-dim FFT is of order $\log N$ where N , the single dimension of the FFT grid, is usually $\mathcal{O}(10^2)$ in most first-principles calculations of bulk materials. The second is the accumulated communication overhead led by many execution times corresponding to a number of wave functions. Thousands of FFTs may be executed at each step of iterative diagonalization. The third is the all-to-all communication required by the data transposes. This could limit the parallel scaling due to the large number of small messages in the network resulting in competition as well as latency issues.

It has already been recognized that making fewer and larger messages can speed up parallel trial wave functions' FFTs. The hybrid OpenMP/MPI implementation [8,9] can lead to fewer and larger messages compared to a pure MPI version. And a blocked version [9] performs a number of trial wave functions' FFTs at the same time to aggregate the message sizes and reduce the latency problem.

In first-principles calculations, we should consider not only the parallel scaling of trial wave functions' FFTs, but also the load balancing of intensive computations on the plane waves that expand the wave functions. The workload of these computations are inhomogeneously distributed on a standard 3-dim FFT grid. Thus a greedy algorithm is usually used to optimize the load balancing. However, this algorithm results in global all-to-all communications across all the processors, thus the latency overhead would grow in proportion to the number of processors

* Corresponding author at: Institute of Applied Physics and Computational Mathematics, Fenghao East Road 2, Beijing 100094, PR China.

E-mail address: wang_han@iapcm.ac.cn (H. Wang).

<http://dx.doi.org/10.1016/j.cpc.2016.07.001>

0010-4655/© 2016 Elsevier B.V. All rights reserved.

and might contribute substantially to the total simulation time. Haynes et al. [10] present a partitioning approach for the 3-dim FFT grid that minimizes the latency cost. Their method depends critically on the Danielson–Lanczos Lemma [11] and requires a particular data distribution, which limits the possibilities to improve the load balancing of the plane waves.

In this paper, we propose a new partitioning method for the 3-dim FFT grid, with which we need independent local all-to-all communications for each data transpose rather than one global all-to-all communication. With this communication pattern preserved, we develop the method to improve the load balancing by adjusting the data distribution among working processors. By numerical examples, we show that although its load balancing is not as perfect as that of the greedy algorithm, the new approach can be more favorable for parallel scaling by making the fewer and larger messages. Hence, we are allowed to accomplish the trade-off between the load balancing of the plane waves and communication overhead in the trial wave functions' FFTs. And such a trade-off could scale the plane wave first-principles calculations up to more nodes. With the proposed partitioning method, we design a compact parallel 3-dim FFT to reduce the amount of calculations and passing messages without loss of accuracy.

The rest of this paper is organized as follows. In Section 2, we explain the elemental role of trial wave functions' FFTs in the plane wave method. In Section 3, we introduce the greedy algorithm for load balancing of the plane waves and analyze the resulting communication cost. In Section 4, we describe the new partitioning algorithms and implementations. In Section 5, we show the numerical results. The last section gives concluding remarks.

2. The role of trial wave functions' FFT

In this section, we explain the elemental role of trial wave functions' FFTs in solving the Kohn–Sham equation using a plane wave basis set.

In the pseudopotential (norm-conserving [12] or ultrasoft [13] pseudopotential) setting or the projector augmented wave (PAW) [14,15] approach, the pseudo wave function $\tilde{\psi}_i$ satisfies the Kohn–Sham equation which looks like

$$\left(-\frac{1}{2}\Delta + V_{loc} + V_{nl}\right)\tilde{\psi}_i = \epsilon_i S \tilde{\psi}_i, \quad (1)$$

where $-\frac{1}{2}\Delta$ is the kinetic energy operator, V_{loc} the local potential, V_{nl} the nonlocal term, and S the overlapping operator. In the case of the norm-conserving pseudopotential, S could simply be interpreted as the identity operator. In this manuscript, we refer to the pseudo wave function simply as the wave function.

We use always the periodic boundary condition and expand the wave function in plane waves:

$$\tilde{\psi}_{nk}(\mathbf{r}) = \sum_{\mathbf{G}} \tilde{\psi}_{nk}(\mathbf{G}) e^{-i(\mathbf{k}+\mathbf{G})\cdot\mathbf{r}}, \quad (2)$$

where the \mathbf{k} 's are vectors sampling the first Brillouin zone, n is an index of the energy level with given \mathbf{k} , and \mathbf{G} 's are the reciprocal lattice vectors. The expansion (2) only includes the plane waves satisfying

$$|\mathbf{k} + \mathbf{G}| < \sqrt{2E_{cut}} \equiv G_{cut}. \quad (3)$$

In the plane wave method, the Hamiltonian matrix is not assembled explicitly. Instead, iterative diagonalization techniques are employed together with the implicit matrix–vector multiplication that is realized as the action of the Hamiltonian operator on the trial wave functions. It is noticed that the local potential is diagonal in the real space. In order to obtain efficiently the action of the

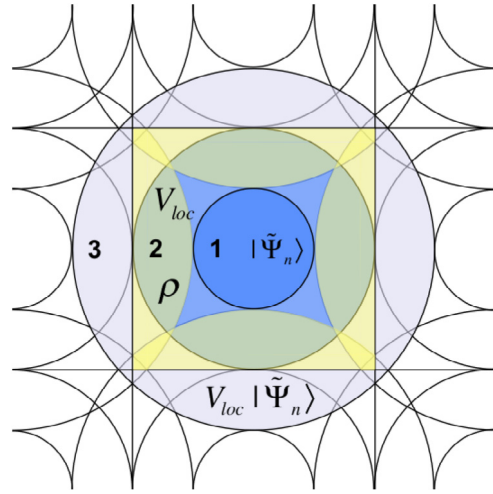


Fig. 1. A two dimensional sketch of the wrap-around errors in the reciprocal space. The wave function $|\Psi\rangle$ is sampled within a sphere with radius G_{cut} (the innermost circle 1). The charge density ρ and the local potential V_{loc} are defined inside a sphere with radius $2G_{cut}$ (circle 2). We would require a sphere with radius $3G_{cut}$ to accurately estimate the operation of the local potential on the trial wave function. If we apply a smaller FFT grid with only $2G_{cut}$, the artificial wrap-around error between $2G_{cut}$ and $3G_{cut}$ would occur and be folded back into circle 2 due to the periodicity. Hence, it is sufficient to approximate the wave functions correctly in circle 1.

local potential on the trial wave function, we should first transform $\tilde{\psi}_{nk}(\mathbf{G})$ to the real space representation $\tilde{\psi}_{nk}(\mathbf{r})$ by one FFT, multiply it with the local potential term, and then transform the product back to the reciprocal space. Consequently, two 3-dim FFTs are required by each action on a trial wave function.

3. The load balancing issue and the greedy algorithm

3.1. The load balancing issue

As mentioned in the previous section, the plane waves are truncated at cut-off radius G_{cut} . Since charge density ρ is the sum of squares of the wave functions, the corresponding cut-off radius for the charge density is $2G_{cut}$. The cut-off radius of the local potential V_{loc} can be regarded the same as that of ρ , because V_{loc} is a functional of ρ . Thus, the cut-off radius of $V_{loc}\tilde{\psi}_{nk}$ is $3G_{cut}$. As illustrated in Fig. 1, it is sufficient to take the FFT grid with only $2G_{cut}$ in order to prevent the wave functions from the wrap-around error.

On one hand, the operation of the local potential on trial wave functions are computed with 3-dim FFTs on the standard grid determined by $2G_{cut}$. On the other hand, we carry out some intensive computations, the time complexities of which are in proportion to the number of plane waves within the cut-off radius G_{cut} , including the assembly of the matrix on the subspace, the orthogonalization of wave functions, and the actions of other parts of the Hamiltonian operator. And the workload of these calculations are not homogeneously distributed on the FFT grid. Therefore, one should consider not only the parallel scaling of FFTs, but also the load balancing issue of the intensive plane wave computations.

3.2. The greedy algorithm

One 3-dim FFT consists of three successive sets of 1-dim FFTs along the x , y and z directions. For each set of 1-dim FFTs, the data layout should guarantee that each processor holds complete columns of data along the FFT direction. Therefore, there are three data layouts of the 1-dim FFTs along the x , y and z directions. We

Download English Version:

<https://daneshyari.com/en/article/4964558>

Download Persian Version:

<https://daneshyari.com/article/4964558>

[Daneshyari.com](https://daneshyari.com)