

Parallelization of the variational fitting of the Coulomb potential

Patrizia Calaminici, Víctor D. Domínguez-Soria, Gerald Geudtner, Elizabeth Hernández-Marín and Andreas M. Köster

Departamento de Química, CINVESTAV

A.P. 14-740 México D.F. 07000, México

(Recibido: 4 de octubre de 2004; Aceptado: 15 de febrero de 2005)

The parallelization of the three-center electron repulsion integrals arising from the variational fitting of the Coulomb potential is presented. A scheme for dynamical load balancing of the corresponding loop structure is discussed. The implementation in the density functional theory program deMon using the message passing interface is described. The efficiency of the parallelization is analyzed by selected benchmark calculations.

Keywords: Density functional theory; DeMon program; Parallelization; Variational fitting; Coulomb potential

1. Introduction

Over the last decade the interest in density functional theory (DFT) methods has grown dramatically within computational chemistry. Numerous systematic studies have shown that DFT methods offer a promising alternative to the more traditional Hartree-Fock approach. Most interestingly, DFT methods include electron correlation at a Hartree like scaling. This is a fundamental difference to Hartree-Fock based methods, where the inclusion of correlation typically increases the computational effort considerably. The rate-limiting steps in today's DFT methods using localized atomic orbitals are the calculations of the two-electron repulsion integrals and the numerical integration of the exchange-correlation contribution.

The use of auxiliary functions for the calculation of these contributions has a long history in density functional theory (DFT) methods [1].

Many years ago Sambe and Felton [2] proposed the use of auxiliary functions in the framework of X_c implementations using the linear combination of Gaussian-type orbitals (LCGTO). Based on this early work Dunlap, Connolly and Sabin [3] introduced the variational fitting of the Coulomb potential in order to avoid the evaluation of four-center integrals in DFT methods. More recently, the so-called resolution of the identity has gained attention [4]. The working formulas are identical to the variational fitting of the Coulomb potential without constraints. The variational fitting of the Coulomb potential reduces the formal scaling to $N^2 \times M$, where N is the number of basis functions and M the number of auxiliary functions. Usually the number of auxiliary functions is three to five times the number of basis functions. It is obvious that this fitting can improve the efficiency of DFT methods considerably. Moreover, it has been shown in the past that the variational approximation of the Coulomb potential is robust and within the intrinsic accuracy of LCGTO-DFT methods [5]. This also includes the availability of accurate gradients and second order properties [6,7]. For the variational fitting of the Coulomb potential an approximated density is introduced. It is expanded in atom-centered auxiliary functions. The

expansion coefficients are obtained variationally by the minimization of the following self-interaction error [3]:

$$\varepsilon_2 = \iint \frac{|\rho(\mathbf{r}_1) - \tilde{\rho}(\mathbf{r}_1)| |\rho(\mathbf{r}_2) - \tilde{\rho}(\mathbf{r}_2)|}{|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2$$

Here, $\rho(\mathbf{r})$ and $\tilde{\rho}(\mathbf{r})$ are the orbital and approximated densities, respectively. In most implementations primitive or contracted Cartesian Gaussian functions are used for the expansion of the approximated density. Therefore, three-center electron repulsion integrals (ERIs) over Cartesian Gaussian functions have to be evaluated. In earlier implementations modified recurrence relations of Obara and Saika (OS) [8], which connect a given integral to others with lower angular momentum, have been used [9]. In a more recent implementation [10] a modified OS method [11] was combined with the method from McMurchie and Davidson (MD) [12] for the evaluation of the ERIs. In the MD method Hermite Gaussian functions [13] are used as intermediates to calculate the integrals over Cartesian Gaussian functions.

Very recently, the use of primitive Hermite Gaussian functions for the expansion of the auxiliary function density has been proposed [14]. It was also shown that this auxiliary function density can be used for the calculation of the exchange-correlation potential [15]. Because the approximated density is a linear combination of auxiliary functions the density calculation and, thus, the exchange-correlation potential calculation at each grid point becomes linear. Instead, the orbital density is calculated from basis function products as:

$$\delta(\vec{r}) = \sum_{\mu,\nu} P_{\mu\nu} \mu(\vec{r}) \nu(\vec{r}),$$

where $P_{\mu\nu}$ represents an element of the density matrix and $\mu(\vec{r})$ and $\nu(\vec{r})$ are contracted Gaussian functions. Therefore, the products of contracted Gaussian functions must be evaluated on the grid. Obviously, the use of the linear expanded approximated density represents a considerable simplification of the grid work.

Table 1. Number of effective shells N_{shell}^{eff} and their percent distribution $\omega_1, \omega_2, \omega_3$ and ω_4 according to the dynamical ERI load balancing for the first 6 SCF cycles of $C_{72}H_{146}$. The corresponding time steps [sec] and normalized time steps are also given.

	SCF cycle					
	1	2	3	4	5	6
N_{shell}^{eff}	28502	47086	45482	44144	43222	42930
ω_1	0.23244	0.24284	0.23554	0.23743	0.23641	0.23576
ω_2	0.23707	0.23970	0.23627	0.23666	0.23507	0.23626
ω_3	0.23700	0.24305	0.23529	0.23605	0.23533	0.23661
ω_4	0.29349	0.27441	0.29290	0.28985	0.29318	0.29137
t_1	33.65	60.80	56.95	56.41	55.21	54.42
t_2	34.77	59.83	57.31	56.55	54.78	54.57
t_3	34.28	60.92	57.22	56.34	54.76	54.75
t_4	37.60	55.25	58.01	55.53	55.40	54.48
$t_{1,norm}$	0.23984	0.25676	0.24816	0.25090	0.25078	0.24938
$t_{2,norm}$	0.24783	0.25266	0.24973	0.25152	0.24883	0.25007
$t_{3,norm}$	0.24433	0.25726	0.24934	0.25059	0.24874	0.25089
$t_{4,norm}$	0.26800	0.23332	0.25278	0.24699	0.25165	0.24966

However, the number of auxiliary functions is usually three to five times the number of basis functions. This could deteriorate the performance of the approximated density calculation on the grid because the calculation of the exponential function is computationally not negligible. Therefore, we are using in our implementation primitive Hermite Gaussian functions which are joined together in sets, sharing the same exponents. As an example, a d auxiliary function set contains ten primitive Hermite Gaussians, one s , three p and six d functions, all with the same exponent. Thus, the number of exponents that have to be evaluated at each grid point is for an auxiliary function density of this structure much smaller than for the corresponding orbital density. Moreover, by using Hermite Gaussian auxiliary functions the Hermite polynomial recurrence relations [13] can be used for the function calculation on the grid. With this approach the computational demand for the numerical integration of the exchange correlation energy and potential becomes secondary with respect to the calculation of the ERIs. Therefore, we decided to parallelize the calculation of these integrals.

In this paper we describe the basic strategy for the parallelization of the three-center electron repulsion integral calculation using FORTRAN and the Message Passing Interface (MPI). In the following section the strategy for the parallelization of the DFT code deMon [16] using MPI is discussed in general. In Sec. 3 the parallelization of the three-center electron repulsion is presented. Benchmark results are discussed in Sec. 4 and concluding remarks are drawn in the last section.

2. Parallelization strategy

The parallelized routines of a program usually have to exchange data. To handle this communication the widely used Message Passing Interface (MPI) is employed in

deMon. Under MPI each parallel execution is called a task. It should be noted that such a task is not necessarily connected to an own node or processor. In principle several tasks may run on the same processor using MPI. However, for our discussion here we can assume that each task is connected to its own CPU.

The transfer of data between tasks is initiated by calls to MPI routines. The calls to these system routines would lead to an error in the link step during the generation of the serial deMon executable on a single processor computer because of the absence of the necessary libraries. Therefore, the CALL statements for the MPI routines are replaced by CALL statements of interface routines. Two sets of these interface routines exist within the deMon source code, one for the serial version and another for the parallel version. For the serial version these routines are dummy routines without any code to be executed. For the parallel version they include the calls to the MPI routines. This strategy has the advantage that only few routines have to be exchanged in order to compile a parallel or serial version of deMon. Obviously, this simplifies code maintenance and ensures the integrity of both, parallel and serial, versions. It also allows an easy change of the messaging system, e.g. from MPI to shmem or Parallel Virtual Machine (PVM).

3. Parallelization of three-center electron repulsion integrals

As already described in the introduction, the calculation of the three-center ERIs becomes the most computational demanding step if the exchange-correlation potential is calculated with the auxiliary function density. In the direct self-consistent field (SCF) approach these integrals have to be calculated twice in each SCF step. One time for the construction of the Kohn-Sham matrix elements,

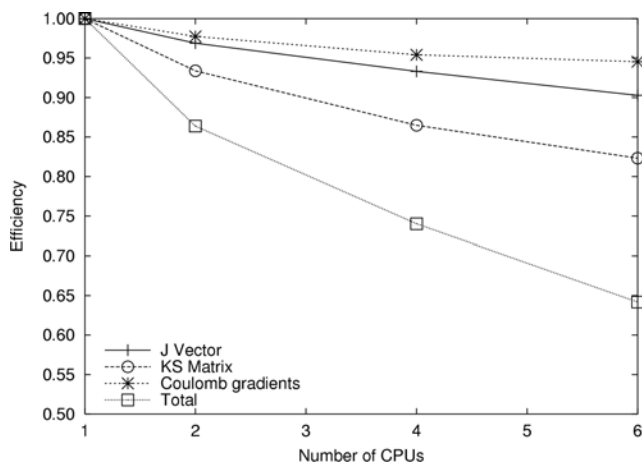


Figure 1. Efficiency, as described in the text, of the parallelized Coulomb vector (J Vector), the Kohn-Sham (KS) matrix, the Coulomb gradients and the complete calculation depending of number of CPUs. As benchmark system the geometry optimization of the alkene molecule $C_{72}H_{146}$ was used.

$$K_{\mu\nu} = H_{\mu\nu} + \sum_k \langle \mu\nu || k \rangle (x_k + z_k),$$

and another time for the calculation of the fitting coefficients,

$$x_l = \sum_k G_{lk} J_k, \quad J_k = \sum_{\mu,\nu} P_{\mu\nu} \langle \mu\nu || k \rangle.$$

Here $H_{\mu\nu}$ and $P_{\mu\nu}$ denote elements of the core Hamiltonian and the density matrix. The fitting coefficients of the approximated density and the exchange-correlation potential are represented by x_k and z_k , respectively. To avoid unnecessary complication in the presentation we have dropped terms arising from the normalization of the approximated density (see [15] for more details). Here we will focus on the parallel computation of the ERIs $\langle \mu\nu || k \rangle$. In this notation, μ and ν refer to contracted Cartesian Gaussian functions which represent the basis set of the calculation. The primitive Hermite Gaussian functions used for the auxiliary function expansion are denoted by k . The symbol $||$ stands for the Coulomb operator $1/|r_1 - r_2|$.

In order to achieve a good efficiency for the parallelization it is necessary to distribute the work, i.e. the calculation of the ERIs, homogeneously over the tasks. For this reason the parallelization of the innermost loop over the auxiliary functions would be too coarse to achieve a good load balancing of the tasks. It is, therefore, more efficient to parallelize over the shells or shell combinations. Here a shell defines a set of contracted basis functions that share the same exponents and contraction pattern, e.g. s , p or d orbitals. Each shell combination possesses a set of orbital products $\mu\nu$. The number of these products depends on the quantum number of the orbitals in the two

shells. The different number of orbital products in the various shell combinations as well as the screening of ERIs based on the maximum overlap integrals of the shell combinations makes the workload balancing a formidable task. Moreover, the dynamical screening of shell combinations based on density matrix differences in the direct SCF procedure forbids any static workload balancing. This situation calls for a dynamical load balancing as described now.

At the beginning of the SCF an effective number of shell combinations N_{shell}^{eff} are calculated. In this way the initial screening is taken into account. The ERIs are homogeneously distributed according to N_{shell}^{eff} and the number of orbital products in each shell. For each shell combination the inner auxiliary function loop is processed within the assigned task. The computational time for the ERI calculation in the construction of the Kohn-Sham matrix and the Coulomb vector J for each task is then measured. Based on these time stamps a redistribution of the effective number of shell combinations for each task is performed according to the following formula:

$$\omega_I^{(n+1)} = \frac{\omega_I^{(n)}}{t_I^{(n)} \sum_J \frac{\omega_J^{(n)}}{t_J^{(n)}}}$$

Here $\omega_I^{(n)}$ represents the percentage of calculated N_{shell}^{eff} by the task I in the n^{th} SCF cycle divided by 100, $t_I^{(n)}$ denotes the corresponding time stamp for task I in the n^{th} SCF cycle and N is the total number of available tasks in the calculation. After about 3 SCF cycles a stable and very efficient load balancing is achieved. We refer to this approach as dynamical ERI load balancing.

In order to be more explicit we now present an example for the dynamical load balancing from a calculation on four CPUs. In the first step N_{shell}^{eff} was 28502. The first percentage values were $\omega_1^{(1)}=0.232$, $\omega_2^{(1)}=0.237$, $\omega_3^{(1)}=0.237$ and $\omega_4^{(1)}=0.293$. This results in the local N_{shell}^{eff} ranges on the 1st CPU from 1 to 6625, on the 2nd CPU from 6626 to 13382, on the 3rd CPU from 13383 to 20137 and on the 4th CPU from 20138 to 28502. The corresponding time stamps for the J vector calculation were $t_1=33.65$ sec, $t_2=34.77$ sec, $t_3=34.28$ sec and $t_4=37.60$ sec. From these timing the normalized time stamps

$$t_{I,norm}^{(n)} = \frac{t_I^{(n)}}{\sum_J t_J^{(n)}}$$

are obtained as $t_{1,norm}=0.240$, $t_{2,norm}=0.248$, $t_{3,norm}=0.244$ and $t_{4,norm}=0.268$. With these time stamps and the above formula for $\omega_I^{(n+1)}$ new percentages for the distribution of

N_{shell}^{eff} in the 2nd SCF cycle are calculated. The results are $\omega_1^{(2)}=0.243$, $\omega_2^{(2)}=0.240$, $\omega_3^{(2)}=0.243$ and $\omega_4^{(2)}=0.274$.

In the 2nd SCF cycle N_{shell}^{eff} is 47086. The number is much larger than in 1st SCF cycle because the used tight-binding starting density matrix is much sparser as the *ab-initio* density matrix. The distribution of this number of effective shell combinations according to the $\omega_l^{(2)}$ values leads to the following $t_{l,norm}$ values: $t_{1,norm}=0.257$, $t_{2,norm}=0.253$, $t_{3,norm}=0.257$ and $t_{4,norm}=0.233$. The corresponding time stamps lead to new $\omega_l^{(3)}$ values for the 3rd SCF cycle. In Table 1 the complete set of the different values for the first 6 SCF cycles are shown. One can see that in the 3rd SCF cycle the $t_{l,norm}$ values are about 0.25 which means that each of the 4 CPUs is doing 1/4 of the work for the calculation of the ERIs even if the N_{shell}^{eff} value changes in every SCF cycle. Because of the similar loop structure for the calculation of the ERIs for the Kohn-Sham (KS) matrix construction a very similar scheme for parallelization of this calculation is used. The same scheme is also applied for the parallelization of the Coulomb integral gradients.

4. Benchmark results

As benchmark system we used the geometry optimization of the alkene molecule C₇₂H₁₄₆. The DZVP basis sets [17] and the local VWN [18] functional were used. In this system the time for the calculation of the parts related to the three center integrals, e.g. the construction of the KS matrix, the Coulomb vector J and the Coulomb gradients, represent about 71% of the total CPU time on one processor. The parts related to the calculation of the exchange correlation energy are only about 10% of the total CPU time. The system was calculated on a cluster with Intel® Xeon™ 2.4 GHz CPUs on 1, 2, 4 and 6 CPUs. The results of the efficiency of the parallelized version for the different number of CPUs are shown in Fig. 1. The efficiency E_N is defined as the quotient of calculation time on one CPU, t_1 , and the product of the calculation time on N CPUs, t_N , and the number of used CPUs, N

$$E_N = \frac{t_1}{Nt_N} .$$

Therefore, the efficiency E_N multiplied by the number of CPUs gives the speedup.

The curves in Fig. 1 show for the Coulomb vector J and for the Coulomb gradients a very good efficiency. With 6 CPUs the efficiency for these parts of the code is over 90% which represents a speedup of about 5.5. The efficiency for the calculation of the ERIs for construction of the KS matrix is less good. Nevertheless, a speedup with 6 CPUs of about 5 is reached. Because of different screening techniques used for the ERIs in the construction of the KS matrix it is more difficult for the dynamic load balancing mechanism to achieve an even workload here. As Fig. 1 shows, the efficiency of the total CPU time is less good

then for the single parts discussed. This is due to other parts in the program that are less efficiently parallelized. Still an overall speedup of 4 for 6 CPUs is achieved. The efficiency values for larger number of CPUs will decrease because the parts each CPU has to calculate will become smaller and smaller. As a result, the dynamic load balancing becomes less efficient. As the system size increases an opposite effect is observed. This means that the efficiency of the parallelization increases with system size for a fixed number of CPUs.

5. Conclusion

An efficient algorithm for the parallelization of the calculation of three-center electron repulsion integrals in the framework of the DFT code deMon is presented. This algorithm has the characteristic to adapt to changing conditions during the SCF cycles. This is achieved by a dynamical load balancing. Because similar situations also occur in the calculation of the exchange correlation potential this algorithm can be used here, too. Due to the implementation only a few routines have to be changed in order to compile a parallel or a serial version of the used code. The efficiency of the parts of the program in which this algorithm is used is with 6 CPUs about 85% or higher.

Acknowledgments

Financial support from the CONACYT projects 36037-E and 40379-F is gratefully acknowledged.

References

- [1] E.J. Baerends, D.E. Ellis, P. Ros, Chem. Phys. **2**, 41 (1973)
- [2] H. Sambe, R.H. Felton, J. Chem. Phys. **62**, 1122 (1975).
- [3] B.I. Dunlap, J.W.D. Connolly, J.R. Sabin, J. Chem. Phys. **71**, 4993 (1979); J.W. Mintmire, B.I. Dunlap, Phys. Rev. A **25**, 88 (1982).
- [4] O. Vahtras, J. Almlöf, M.W. Feyereisen, Chem. Phys. Lett. **213**, 514 (1993).
- [5] B.I. Dunlap, J. Mol. Struct. (THEOCHEM) **529**, 37 (2000).
- [6] B.I. Dunlap, J. Andzelm, Phys. Rev. A **45**, 81 (1992)
- [7] A. Komornicki, G. Fitzgerald, J. Chem. Phys. **98**, 1398 (1993).
- [8] S. Obara, A. Saika, J. Chem. Phys. **84**, 3963 (1986).
- [9] J. Andzelm, E. Wimmer, J. Chem. Phys. **96**, 1280 (1992).
- [10] A.M. Köster, J. Chem. Phys. **104**, 4114 (1996).
- [11] M. Head-Gordon, J.A. Pople, J. Chem. Phys. **89**, 5777 (1988).
- [12] L.E. McMurchie, E.R. Davidson, J. Comput. Phys. **26**, 218 (1978).
- [13] V.R. Saunders, in *Methods in Computational Physics*, edited by G.H.F. Diercksen and S. Wilson (Reidel, Dordrecht, 1983).
- [14] A.M. Köster, J. Chem. Phys. **118**, 9943 (2003)
- [15] A.M. Köster, J.U. Reveles, J.M. del Campo, J. Chem. Phys. **121**, 3417 (2004).
- [16] A.M. Köster, P. Calaminici, R. Flores-Moreno, G. Geudtner, A. N. Goursot, T. Heine, F. Janetzko, S. Patchkovskii, J.U. Reveles, A. Vela and D.R. Salahub, deMon, NRC, Canada (2004).

[17] N. Godbout, D.R. Salahub, J. Andzelm, E. Wimmer, *Can. J. Phys.* **70**, 560 (1992).

[18] S.H. Vosko, L. Wilk, M. Nusair, *Can. J. Phys.* **58**, 1200 (1980).