

# PCCP

Accepted Manuscript



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

*Accepted Manuscripts* are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.

# The grid-based fast multipole method – a massively parallel numerical scheme for calculating two-electron interaction energies

Elias A. Toivanen<sup>a</sup>, Sergio A. Losilla<sup>b</sup>, Dage Sundholm<sup>\*c</sup>

Received Xth XXXXXXXXXXXX 20XX, Accepted Xth XXXXXXXXXXXX 20XX

First published on the web Xth XXXXXXXXXXXX 200X

DOI: 10.1039/b000000x

Algorithms and working expressions for a grid-based fast multipole method (GB-FMM) have been developed and implemented. The computational domain is divided into cubic subdomains, organized in a hierarchical tree. The contribution to the electrostatic interaction energies from pairs of neighboring subdomains are computed using numerical integration, whereas as the contributions from further apart subdomains are obtained using multipole expansions. The multipole moments of the subdomains are obtained by numerical integration. Linear scaling is achieved by translating and summing the multipoles according to the tree structure, such that each subdomain interacts with a number of subdomains that is almost independent of the size of the system. To compute electrostatic interaction energies of neighboring subdomains, we employ an algorithm which performs efficiently on general purpose graphics processing units (GPGPU). Calculations using one CPU for the FMM part and 20 GPGPUs consisting of tens of thousands of execution threads for the numerical integration algorithm show the scalability and parallel performance of the scheme. For calculations on systems consisting of Gaussian functions ( $\alpha = 6$ ) distributed as fullerenes from  $C_{20}$  to  $C_{720}$ , the total computation time and relative accuracy (ppb) are independent of system size.

## 1 Introduction

The computational speed has doubled every 18 month since the birth of the first computer. The exponential growth of the computational efficiency is called Moore's law.<sup>1</sup> However, the computational speed of the individual processors has already practically reached its maximum when using silicon-based technology, and the primary means to improve the performance is parallelization.

The parallelization of quantum chemistry codes is not a trivial task, especially when aiming at implementations that run efficiently on thousands of central processing units (CPU) or general purpose graphics processing units (GPGPU), although great advances have been made in the recent years.<sup>2–15</sup> For example, Yasuda and Maruoka<sup>16</sup> report a speed-up of four when using a GPGPU-accelerated version of their electron repulsion integral code, which is to be contrasted with the theoretical peak performance advertised by GPGPU vendors that suggest several orders of magnitude better results.

Setting technicalities such as memory access patterns and communication overheads aside, the software-hardware gap ultimately results from a mismatch between computer algorithms and the execution model of the computer hardware. In particular, maximal performance can never be attained if the focus is on parallelizing individual code segments interleaved by sequential parts, which is the core content of Amdahl's

<sup>a</sup> Department of Chemistry, A. I. Virtanens plats 1, P. O. Box 55, FI-00014 University of Helsinki, Finland. Email: [eatoivan@chem.helsinki.fi](mailto:eatoivan@chem.helsinki.fi)

<sup>b</sup> Department of Chemistry, A. I. Virtanens plats 1, P. O. Box 55, FI-00014 University of Helsinki, Finland. Present address: Department of Chemistry, University of Aarhus, DK-8000 Aarhus C, Denmark. Email: [sergio@chem.au.dk](mailto:sergio@chem.au.dk)

<sup>c</sup> Department of Chemistry, A. I. Virtanens plats 1, P. O. Box 55, FI-00014 University of Helsinki, Finland. E-mail: [Dage.Sundholm@helsinki.fi](mailto:Dage.Sundholm@helsinki.fi)

law.<sup>17</sup> It states that for massively parallel computers, the wall time it takes to perform a given calculations is determined by the computational time that is needed to perform the operations in the serial part of the code. The law seems to introduce severe limitations, because if 10% of the code is run serially the maximum speed-up is a factor of about 10. However, if the time needed for the serial part can be made more or less independent of the amount of parameters needed to describe the molecule, the total computational scaling becomes independent of the size of the molecule, as long as a sufficiently large computer is available. Alternatively, if the size of the system is kept the same, a higher accuracy is obtained at the same cost.

To be able to explore the computational capacity of future computers one has to design algorithms that run efficiently on massively parallel computers. Real-space grid-based methods, also known as numerical methods are well-suited for these architectures because the local support of the basis functions introduces new opportunities to divide the calculation into lots of independent parts that can be performed in parallel. Once numerical integration can be performed concurrently on sufficiently many cores, the performance is not going to be limited by the total amount of work but by the work per computational node. In addition, the flexibility of numerical approaches also renders accurate calculations with small basis-set truncation errors ( $\mu E_h$ ) feasible.<sup>18–24</sup> Thus, quantum chemistry calculations of the future employ local basis functions.

In this work, we discuss a fully numerical technique for computing electrostatic interaction energies between charge densities. In the context of *ab initio* electronic structure calculations, this is needed to compute large numbers of two-electron repulsion integrals. The evaluation of these integrals constitutes one of the computational bottlenecks of the self-consistent-field (SCF) calculations. Implementation of the FMM scheme in quantum chemistry codes that rely on atom-centered Gaussian-type basis sets significantly reduces the number of explicit integrals over the basis functions.<sup>25–29</sup> Instead of the  $\mathcal{O}(M^4)$  integrals in the naive case, the FMM-accelerated prescreening scheme computes all interactions with only  $\mathcal{O}(M)$  work, as demonstrated by e.g. Rudberg and Sałek.<sup>29</sup> We have adapted the FMM scheme originally conceived by Greengard and Rokhlin<sup>30</sup> and introduced into quantum chemistry codes by White and Head-Gordon *et al.*<sup>31,32</sup> to a fully numerical framework. In the grid-based fast multipole method (GB-FMM), the long-ranged and short-ranged contributions of the two-electron interactions are identified and treated differently. The long-ranged interactions are computed using an FMM scheme with numerically calculated multipole moments, whereas the short-ranged contributions are obtained by numerical integration as described in our previous work.<sup>15,33–36</sup> For a discussion of alternative parallelization techniques for Poisson solvers, we refer to the recent survey by García-Risueño *et al.*<sup>37</sup>

Even though the FMM scheme is surely a key component in the quest for linear scaling SCF calculations,<sup>38</sup> our interest in the FMM algorithm does not stem merely from asymptotic complexity arguments. Instead, we also tackle a somewhat different problem that is typical for grid-based methods, namely that the memory requirements for representing functions on a grid grow linearly with the volume of the system. The FMM scheme allows one to circumvent such limitations as it provides a natural framework for decomposing functions into local grids that can be distributed across several computational nodes opening the avenue for massive parallelization.

The article is outlined as follows. In Section 2, the GB-FMM scheme is thoroughly described. The bipolar series expansion of electrostatic interactions and the translation of multipole moments are discussed in Subsections 2.1 and 2.2. The partitioning of the computational domain and the grids are described in Subsections 2.3 and 2.4. The expressions for calculation the two-electron interaction energy using numerical integration combined with the GB-FMM scheme are derived in Subsections 2.5. The algorithm for calculating multipole moments is presented in Subsection 2.7, whereas the algorithm of the direct numerical integration of the near-field contributions to the two-electron interaction is described in Subsection 2.8. The

different components are put together in Subsection 2.9, where the final algorithm is given. The timings and accuracy of the calculations are presented in Section 3. The article ends with conclusions and a future outlook in Section 4.

## 2 The Grid-Based Fast Multipole Method

The ultimate goal is efficient calculations of the  $N_\rho(N_\rho + 1)/2$  distinct pair interactions resulting from  $N_\rho$  charge distributions, that is

$$U_{\mu\nu} = \int_{\mathbb{R}^3} \int_{\mathbb{R}^3} \rho_\mu(\mathbf{r}) \frac{1}{|\mathbf{r}' - \mathbf{r}|} \rho_\nu(\mathbf{r}') d^3r d^3r' \quad (1)$$

for  $1 \leq \mu \leq N_\rho$ ,  $1 \leq \nu \leq \mu$ .

In this Section, we present the different expressions and methods required to arrive to the final expression for the energy, Eq. (26). The algorithm to efficiently evaluate the interaction energy is then given in Subsection 2.9.

### 2.1 Bipolar expansion of electrostatic interactions

The mathematical basis of the fast multipole method is the bipolar series expansion<sup>39</sup> of the Coulomb potential:

$$\frac{1}{|\mathbf{r}' - \mathbf{r}|} = \sum_{l=0}^{\infty} \sum_{m=-l}^l \sum_{l'=0}^{\infty} \sum_{m'=-l'}^{l'} S_{lm}(\mathbf{r} - \mathbf{P}) T_{lm,l'm'}(\mathbf{Q} - \mathbf{P}) S_{l'm'}(\mathbf{r}' - \mathbf{Q}). \quad (2)$$

In Eq. (2),  $\mathbf{P}$  and  $\mathbf{Q}$  are two distinct ( $\mathbf{P} \neq \mathbf{Q}$ ) but arbitrary reference points and  $S_{lm}$  are real solid harmonics in Racah's normalization.<sup>40,41</sup>  $T_{lm,l'm'}(\mathbf{Q} - \mathbf{P})$  are the elements of the *interaction matrix*  $\mathbf{T}(\mathbf{Q} - \mathbf{P})$ . We provide explicit expressions for the elements in terms of real spherical harmonics in Appendix A.

Eq. (2) conveniently decouples coordinates  $\mathbf{r}$  and  $\mathbf{r}'$ , which significantly simplifies the evaluation of the interaction energy in Eq. (1). The pair interaction energy can now be written as:

$$U_{\mu\nu} = \sum_{l=0}^{\infty} \sum_{m=-l}^l \left( \int_{\mathbb{R}^3} S_{lm}(\mathbf{r} - \mathbf{P}) \rho_\mu(\mathbf{r}) d^3r \right) \sum_{l'=0}^{\infty} \sum_{m'=-l'}^{l'} T_{lm,l'm'}(\mathbf{Q} - \mathbf{P}) \left( \int_{\mathbb{R}^3} S_{l'm'}(\mathbf{r}' - \mathbf{Q}) \rho_\nu(\mathbf{r}') d^3r' \right), \quad (3)$$

where we identify the multipole moments (monopole, dipole, quadrupole, etc.) of the charge distributions

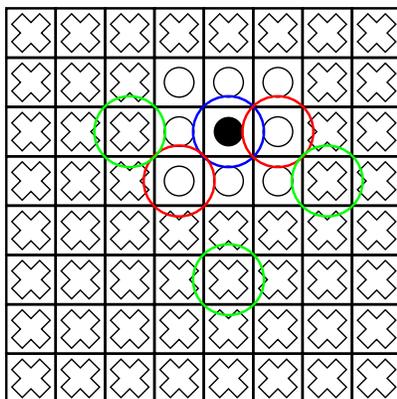
$$q_{\mu,lm}(\mathbf{P}) = \int_{\mathbb{R}^3} S_{lm}(\mathbf{r} - \mathbf{P}) \rho_\mu(\mathbf{r}) d^3r. \quad (4)$$

Organizing the multipole moments as a vector  $\mathbf{q}_\mu$ , Eq. (3) can be rewritten as

$$U_{\mu\nu} = \mathbf{q}_\mu^T(\mathbf{P}) \mathbf{v}_\nu(\mathbf{P}). \quad (5)$$

with the potential moments vector given by

$$\mathbf{v}_\nu(\mathbf{P}) = \mathbf{T}(\mathbf{Q} - \mathbf{P}) \mathbf{q}_\nu(\mathbf{Q}). \quad (6)$$



**Fig. 1** The blue sphere overlaps with the red spheres, but not with the green spheres. Hence, the interaction of the density enclosed in the box marked with a black circle with the boxes marked with white circles must be computed directly using Eq. (1). On the other hand, the interaction between the box with the black circle and the boxes with white crosses can be computed using the multipole expansion in Eq. (3).

The elements of the potential vectors  $\mathbf{v}_v(\mathbf{P})$  have a direct physical significance,<sup>31</sup> as they are the expansion coefficients in solid harmonics of the electric potential due to  $\rho_v(\mathbf{r})$  expanded around  $\mathbf{P}$ :

$$V_v(\mathbf{r}) = \sum_{l=0}^{\infty} \sum_{m=-l}^l v_{v,lm} S_{lm}(\mathbf{r} - \mathbf{P}). \quad (7)$$

The main limitation of the bipolar expansion in Eq. (2) is that it converges only if  $|\mathbf{r}' - \mathbf{Q}| + |\mathbf{r} - \mathbf{P}| < |\mathbf{Q} - \mathbf{P}|$ . In geometrical terms, this means that it must be possible to enclose the two charge densities in spheres which do not overlap, as shown in Fig. 1. Furthermore, in a practical implementation the expansion must be truncated at a certain order  $l_{\max}$ . However, the error diminishes systematically as  $l_{\max}$  is increased, implying that this is in general not a problem.

In the fast multipole method, the bipolar expansion is efficiently exploited by partitioning the charge densities into spatial subdomains, such that each of these subdomains overlaps only with a handful of the other subdomains, as described in the next Section. Grid-based methods have the advantage that non-overlapping domains are very easily identified due to the local support of the numerical basis functions.

## 2.2 Translation of multipole moments

An important property of multipole moments is that the expansion point can be shifted, without explicit recalculation of the multipole moments. The translation of the multipole moment from  $\mathbf{Q}$  to  $\mathbf{P}$  is given by

$$\mathbf{q}_\mu(\mathbf{P}) = \mathbf{W}(\mathbf{Q} - \mathbf{P})\mathbf{q}_\mu(\mathbf{Q}). \quad (8)$$

Expressions for the elements of the translation matrix  $\mathbf{W}$  can be found in Appendix A. The potential moments can be translated in a similar fashion using

$$\mathbf{v}_v(\mathbf{Q}) = \mathbf{W}^T(\mathbf{Q} - \mathbf{P})\mathbf{v}_v(\mathbf{P}). \quad (9)$$

## 2.3 Domain partitioning and box hierarchy

A *box* is here a three-dimensional Cartesian domain. The  $A$ -th box  $\Omega^{(A)}$  is defined as:

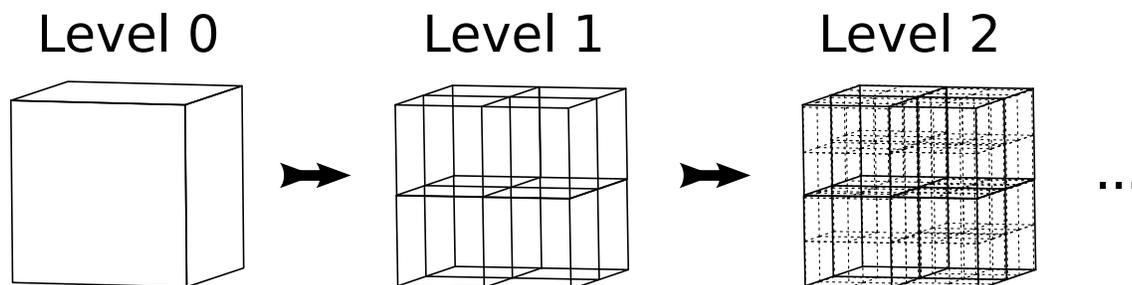


Fig. 2 An octree data structure is obtained by subdividing a three-dimensional domain recursively.

$$\Omega^{(A)} = [x_{\min}^{(A)}, x_{\max}^{(A)}] \times [y_{\min}^{(A)}, y_{\max}^{(A)}] \times [z_{\min}^{(A)}, z_{\max}^{(A)}], \quad (10)$$

The zeroth box,  $\Omega^{(0)}$  is the complete computational domain. The center of box  $A$  is the point  $\mathbf{C}^{(A)}$  with coordinates

$$\mathbf{C}^{(A)} = \left( \frac{x_{\min}^{(A)} + x_{\max}^{(A)}}{2}, \frac{y_{\min}^{(A)} + y_{\max}^{(A)}}{2}, \frac{z_{\min}^{(A)} + z_{\max}^{(A)}}{2} \right). \quad (11)$$

Each box is completely enclosed by a sphere centered at  $\mathbf{C}^{(A)}$  with radius

$$r^{(A)} = \frac{1}{2} \sqrt{(x_{\max}^{(A)} - x_{\min}^{(A)})^2 + (y_{\max}^{(A)} - y_{\min}^{(A)})^2 + (z_{\max}^{(A)} - z_{\min}^{(A)})^2}. \quad (12)$$

We will refer to  $r^{(A)}$  as the *extent* of box  $A$ . The enclosing spheres of boxes  $A$  and  $B$  overlap when  $|\mathbf{C}^{(A)} - \mathbf{C}^{(B)}| \leq r^{(A)} + r^{(B)}$ .

Boxes are recursively subdivided such that each box is the *parent* to eight *children* boxes, thus forming the octree data structure that is illustrated in Fig. 2. The indices of the children of box  $A$  are denoted by  $\text{Children}(A)$ . Likewise, the index of the parent of box  $A$  is  $\text{parent}(A)$ .

The boxes are grouped into *levels*, which we define recursively as:

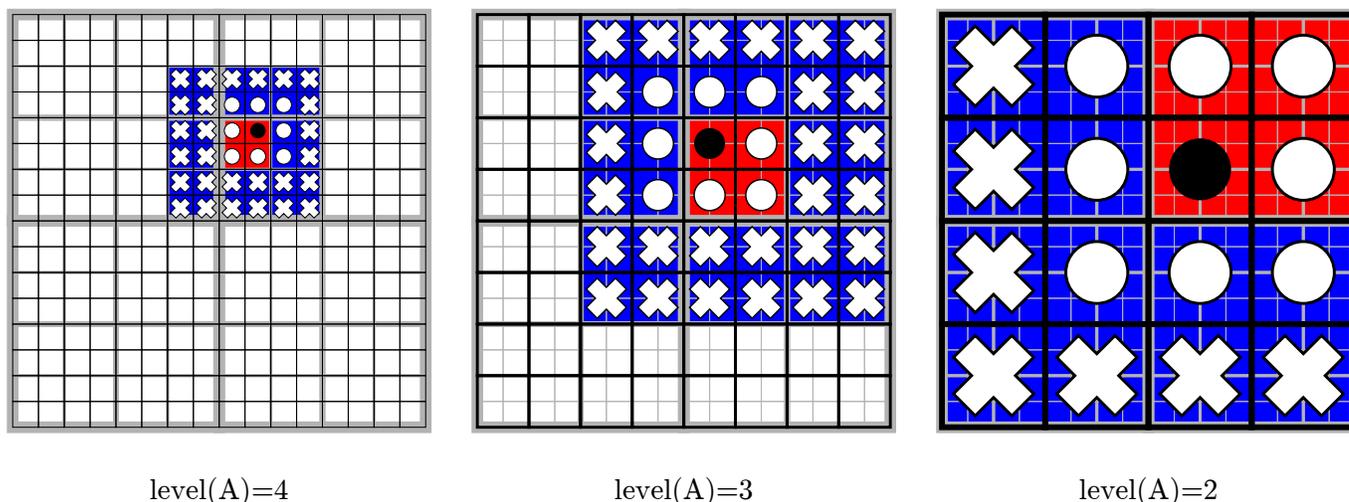
$$\text{level}(0) = 0, \quad (13)$$

$$\text{level}(A) = \text{level}(\text{parent}(A)) + 1. \quad (14)$$

The  $l$ -th level contains  $8^l$  boxes, and when the number of divisions is set to  $D_{\max}$ , the number of childless leaf nodes at the highest level is  $8^{D_{\max}}$ . In order to make an efficient use of the bipolar expansion, the vicinity of every box at all levels is divided into two groups: the *nearest neighbors* and the *local far field*. The nearest neighbors of box  $A$ ,  $\text{NN}(A)$ , are the boxes which belong to the same level as  $A$  and whose enclosing spheres overlap with that of  $A$ :

$$\text{NN}(A) = \left\{ B : \text{level}(A) = \text{level}(B), |\mathbf{C}^{(A)} - \mathbf{C}^{(B)}| \leq r^{(A)} + r^{(B)} \right\}. \quad (15)$$

Note that  $A$  is itself a member of  $\text{NN}(A)$ . The number of nearest neighbors ranges from 8 when  $A$  is in a corner to 27 when it is in the interior of the computational domain. The local far field of the  $A$ -th



**Fig. 3** Illustration of the domain partitioning for  $D_{\max} = 4$ . Each figure shows, for a selected box  $A$  (marked with a black circle):  $\text{NN}(A)$  (circles, both white and black),  $\text{parent}(A)$  (in red),  $\text{NN}(\text{parent}(A))$  (in blue) and  $\text{LFF}(A)$  (white crosses). The selected box (black circle) in each figure is the parent to the box marked with a black circle in the figure to its left.

box,  $\text{LFF}(A)$ , consists of all the children of the nearest neighbors of the parent of  $A$ , which are not nearest neighbors of  $A$  themselves:

$$\text{LFF}(A) = \left[ \bigcup_{B \in \text{NN}(\text{parent}(A))} \text{Children}(B) \right] \setminus \text{NN}(A). \quad (16)$$

There are  $4^3 - 3^3 = 37$  to  $6^3 - 3^3 = 189$  boxes in the local far field of any box. From this definition we note that, for boxes at levels 0 and 1, the local far field is empty, because box 0 has no parent and no nearest neighbors.

An important property of this separation is that  $\Omega^{(0)}$  is *exactly covered* by the boxes in  $\text{NN}(A) \cup \text{LFF}(A) \cup \text{LFF}(\text{parent}(A)) \cup \text{LFF}(\text{parent}(\text{parent}(A))) \cup \dots$ , for any box  $A$ . In this way, the complete domain can be partitioned into a hierarchy of non-overlapping domains,

$$\int_{\mathbb{R}^3} f(\mathbf{r}) d^3 r = \sum_{B \in \text{NN}(A)} \int_{\Omega(B)} f(\mathbf{r}) d^3 r + \sum_{B \in \text{LFF}(A)} \int_{\Omega(B)} f(\mathbf{r}) d^3 r + \sum_{B \in \text{LFF}(\text{parent}(A))} \int_{\Omega(B)} f(\mathbf{r}) d^3 r + \dots \quad (17)$$

The decomposition in Eq. (17) will always terminate with boxes at level two. The concepts defined in this section are illustrated in Fig. 3 in two dimensions, for  $D_{\max} = 4$ .

## 2.4 Grid representation

For each of the boxes at the highest level  $D_{\max}$ , an equidistant Cartesian grid is constructed as presented below. A more detailed discussion can be found elsewhere.<sup>42</sup>

The one-dimensional grids consist of  $N$  points separated by a distance  $h$ , the *grid step*, denoted by  $h$ . For example, for the  $x$  dimension, the points are  $x_1^{(A)} = x_{\min}^{(A)}, x_2^{(A)}, \dots, x_N^{(A)} = x_{\max}^{(A)}$ . A one-dimensional basis function is assigned to each point, *e.g.*  $\chi_j^{(A,y)}(y)$  is assigned to point  $y_j^{(A)}$ . The basis functions are polynomials of degree  $P$  with small support around their grid points. We employ sixth-degree polynomials ( $P = 6$ ) in the present work.

The three-dimensional basis set is an outer (tensorial) product of the one-dimensional local basis sets. The expansion coefficients are the values of the function in the corresponding grid points. Taking the density  $\rho_\mu(\mathbf{r})$  as an example, its expansion in the local basis functions within box  $A$  is given by

$$\rho_\mu(\mathbf{r}) \approx \sum_{ijk} (\rho_\mu)_{ijk}^{(A)} \chi_i^{(A,x)}(x) \chi_j^{(A,y)}(y) \chi_k^{(A,z)}(z). \quad (18)$$

with the expansion coefficients  $(\rho_\mu)_{ijk}^{(A)} = \rho_\mu(x_i^{(A)}, y_j^{(A)}, z_k^{(A)})$ .

As we have shown before,<sup>35</sup> this representation is adequate for smooth charge densities. For all-electron calculations, the presence of cusps at the nuclear positions requires a prohibitive number of points, which can be circumvented by explicitly representing those cusps, which is beyond the scope of the present work but described in detail in our previous work.<sup>42</sup>

The total number of grid points is given by  $(2^{D_{\max}} N)^3$ . Equivalently, this is  $(L/h)^3$ , where  $L$  is the total length of the computational domain  $\Omega^{(0)}$ . For a typical step of  $0.1 a_0$  and  $N = 200$ ,  $D_{\max} = 3$  allows treating systems with a volume of up to  $160 \times 160 \times 160 a_0^3$ .

## 2.5 The FMM energy expression

Using Eq. (17), the expression for the energy in Eq. (1) can be reorganized as

$$U_{\mu\nu} = \sum_{A:\text{level}(A)=D_{\max}} \int_{\Omega^{(A)}} \rho_\mu(\mathbf{r}) \left[ \sum_{B \in \text{NN}(A)} \int_{\Omega^{(B)}} \frac{1}{|\mathbf{r}' - \mathbf{r}|} \rho_\nu(\mathbf{r}') d^3 r' + \sum_{B \in \text{LFF}(A)} \int_{\Omega^{(B)}} \frac{1}{|\mathbf{r}' - \mathbf{r}|} \rho_\nu(\mathbf{r}') d^3 r' + \sum_{B \in \text{LFF}(\text{Parent}(A))} \int_{\Omega^{(B)}} \frac{1}{|\mathbf{r}' - \mathbf{r}|} \rho_\nu(\mathbf{r}') d^3 r' + \dots \right] d^3 r. \quad (19)$$

The contributions arising from the nearest neighbors are grouped together into the near-field interaction energy,  $U_{NF}^{(A)}$ . The remaining terms constitute the far-field interaction energy  $U_{FF}^{(A)}$ .

$$U_{\mu\nu} = \sum_{A:\text{level}(A)=D_{\max}} U_{NF}^{(A)} + U_{FF}^{(A)}. \quad (20)$$

$U_{NF}^{(A)}$  can be written more compactly as

$$U_{NF}^{(A)} = \int_{\Omega^{(A)}} \rho_{\mu}(\mathbf{r}) V_{\nu}^{(A)}(\mathbf{r}) d^3 r \quad (21)$$

where  $V_{\nu}^{(A)}(\mathbf{r})$ , which is the near-field potential at box  $A$ , is computed as a sum of contributions  $V_{\nu}^{(B,A)}(\mathbf{r})$ :

$$V_{\nu}^{(A)}(\mathbf{r}) = \sum_{B \in \text{NN}(A)} \int_{\Omega^{(B)}} \frac{1}{|\mathbf{r}' - \mathbf{r}|} \rho_{\nu}(\mathbf{r}') d^3 r' = \sum_{B \in \text{NN}(A)} V_{\nu}^{(B,A)}(\mathbf{r}). \quad (22)$$

Because all the density pairs contributing to the far-field energy are non-overlapping,  $U_{FF}^{(A)}$  can be calculated by exploiting the bipolar expansion in Eq. (3) as

$$U_{FF}^{(A)} = \left( \mathbf{q}_{\mu}^{(A)} \right)^T \mathbf{v}_{\nu}^{(A)} \quad (23)$$

where  $\mathbf{q}_{\mu}^{(A)}$  are the multipole moments of the part of  $\rho_{\mu}(\mathbf{r})$  contained within  $\Omega^{(A)}$ ,

$$q_{\mu,lm}^{(A)} = \int_{\Omega^{(A)}} S_{lm}(\mathbf{r} - \mathbf{C}^{(A)}) \rho_{\mu}(\mathbf{r}) d^3 r, \quad (24)$$

and the far-field potential moment vector  $\mathbf{v}_{\mu}^{(A)}$  is defined recursively as

$$\mathbf{v}_{\mu}^{(A)} = \sum_{B \in \text{LFF}(A)} \mathbf{T}(\mathbf{C}^{(B)} - \mathbf{C}^{(A)}) \mathbf{q}_{\mu}^{(B)} + \mathbf{W}^T (\mathbf{C}^{(A)} - \mathbf{C}^{(\text{parent}(A))}) \mathbf{v}_{\mu}^{(\text{parent}(A))}. \quad (25)$$

The final expression for the energy is then given by

$$U_{\mu\nu} = \sum_{A: \text{level}(A)=D_{\max}} \int_{\Omega^{(A)}} \rho_{\mu}(\mathbf{r}) V_{\nu}^{(A)}(\mathbf{r}) d^3 r + \left( \mathbf{q}_{\mu}^{(A)} \right)^T \mathbf{v}_{\nu}^{(A)}. \quad (26)$$

## 2.6 Parallelization

The linear scaling nature of the fast multipole method is apparent from Eq. (26): for a given set of far field potential vectors and near field potentials, the energy is evaluated in a single loop over the boxes. For an analysis establishing linear scaling complexity also in the potential construction step, we refer to the work by White and Head-Gordon.<sup>31</sup> Here, we focus on presenting the traditional FMM algorithm in a maximally parallel fashion.

Each box  $A$  is assigned to a node, denoted with  $A$ . If the GB-FMM scheme is parallelized to the fullest, the number of computational nodes will be  $8^2 + 8^3 + \dots + 8^{D_{\max}}$ . We will refer to the  $8^{D_{\max}}$  nodes at the highest level as *near-field nodes*. It is not required that all nodes are physically different. For example, in the benchmarks shown in Section 3, all non near-field nodes are in fact the same.

The algorithms required for the calculation of pair energies are represented in Fig. 4 to Fig. 8. The operations, which are run concurrently by all nodes, are represented within rectangles. The near-field nodes must often carry out additional steps, which are indicated separately. Grey rectangles indicate steps that require some type of inter-node communication.

In order to carry out the necessary operations, each node  $A$  stores the center of its own box,  $\mathbf{C}^{(A)}$ , and the centers of all its children, its parent and all boxes in  $\text{LFF}(A)$ . In addition, each near-field node must also

All near-field nodes  $A$  concurrently do

$$\begin{aligned}
 & a_{ijk} = (\rho_\mu)_{ijk}^{(A)} \forall i, j, k \\
 & \text{for all } 0 \leq u \leq l_{\max} \text{ do} \\
 & \quad b_{jk} = \sum_i I_i^{(A,x)} a_{ijk} \forall j, k \\
 & \quad \text{for all } 0 \leq v \leq l_{\max} - u \text{ do} \\
 & \quad \quad c_k = \sum_j I_j^{(A,y)} b_{jk} \forall k \\
 & \quad \quad \text{for all } 0 \leq w \leq l_{\max} - u - v \text{ do} \\
 & \quad \quad \quad d = \sum_k I_k^{(A,z)} c_k \\
 & \quad \quad \quad \text{for all } (l, m) \text{ such that } C_{uvw}^{lm} \neq 0 \text{ do} \\
 & \quad \quad \quad \quad q_{\mu,lm}^{(A)} \leftarrow q_{\mu,lm}^{(A)} + C_{uvw}^{lm} d \\
 & \quad \quad \quad \quad \text{(Note: } d = \int_{\Omega(A)} x^u y^v z^w \rho_\mu(\mathbf{r}) dx dy dz) \\
 & \quad \quad \quad c_k \leftarrow \tilde{z}_k c_k \forall k \\
 & \quad \quad b_{jk} \leftarrow \tilde{y}_j b_{jk} \forall j, k \\
 & \quad a_{ijk} \leftarrow \tilde{x}_i a_{ijk} \forall i, j, k
 \end{aligned}$$

**Fig. 4** Algorithm for the numerical integration of the multipole moments of the boxes with  $\text{level}(A) = D_{\max}$ .

store the grids of its nearest neighbors (including its own), and be able to calculate or read the expansion coefficients  $(\rho_\mu^{(A)})_{ijk}$ , which are the expansion coefficients of all the input densities in box  $A$ .

## 2.7 Evaluation of multipole moments

At the highest division level, the multipoles are integrated using the expression in Eq. (24). The fact that the multipole moments have to be computed by numerical integration may seem daunting, especially when the multipole moments of Gaussian primitives can be evaluated more or less trivially.<sup>32</sup> However, it turns out that the tensorial nature of the local basis functions leads to a very efficient algorithm, as the multipole moment integrals always separate to products of one-dimensional integrals, many of which can be reused.

The real solid harmonics can be expressed using Cartesian coordinates as

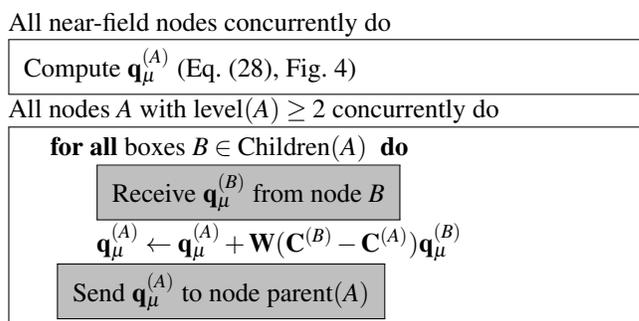
$$S_{lm}(r, \theta, \varphi) = \sum_{uvw} C_{uvw}^{lm} x^u y^v z^w. \quad (27)$$

where  $C_{uvw}^{lm}$  are the transformation coefficients.<sup>41</sup> The multipole moments in the Cartesian representation are easily computed with the following expression:

$$q_{\mu,lm}^{(A)} = \sum_{uvw} C_{uvw}^{lm} \sum_k I_k^{(A,z)} \tilde{z}_k^w \sum_j I_j^{(A,y)} \tilde{y}_j^v \sum_i I_i^{(A,x)} (\rho_\nu)_{ijk}^{(A)} \tilde{x}_i^u, \quad (28)$$

where the shifted coordinates are  $\tilde{x}_i = x_i^{(A)} - C_x^{(A)}$ ,  $\tilde{y}_j = y_j^{(A)} - C_y^{(A)}$  and  $\tilde{z}_k = z_k^{(A)} - C_z^{(A)}$ . The quantities  $I_i^{(A,x)}$ ,  $I_j^{(A,y)}$  and  $I_k^{(A,z)}$  are one-dimensional integrals over the basis functions, e.g.

$$I_i^{(A,x)} = \int_{x_{\min}^{(A)}}^{x_{\max}^{(A)}} \chi_i^{(A,x)}(x) dx. \quad (29)$$



**Fig. 5** Algorithm for calculating  $\mathbf{q}_\mu^{(A)}$ , the multipole moments for all boxes at all levels greater or equal to two for a charge density  $\rho_\mu$ .

The most expensive operations appear in the outermost loops, with a computational cost that is asymptotically proportional to  $\mathcal{O}((l_{\max} + 1)N^3)$ , which is a substantial improvement as compared to the cost of  $\mathcal{O}((l_{\max} + 1)^2N^3)$  for the naive implementation.

For the boxes at lower levels, the multipole moments can be computed recursively from the multipoles of their children using Eq. (8):

$$\mathbf{q}_\mu^{(A)}(\mathbf{C}^{(A)}) = \sum_{B \in \text{Children}(A)} \mathbf{W}(\mathbf{C}^{(B)} - \mathbf{C}^{(A)})\mathbf{q}_\mu^{(B)}(\mathbf{C}^{(B)}). \quad (30)$$

The algorithm to compute the multipole moments at all boxes is outlined in Fig. 5. Initially, all near-field nodes must integrate the multipole moments in their respective boxes as described above. Then, the boxes at lower levels can recursively reconstruct their multipole moments. For each of its 8 children, each node receives the  $(l_{\max} + 1)^2$  entries of the multipole moment vector, which is then multiplied by the  $(l_{\max} + 1)^2 \times (l_{\max} + 1)^2$  translation matrix, and the result is accumulated into  $\mathbf{q}_\mu^{(A)}$ . In the FMM literature, this step is often referred to as the *upward pass*.

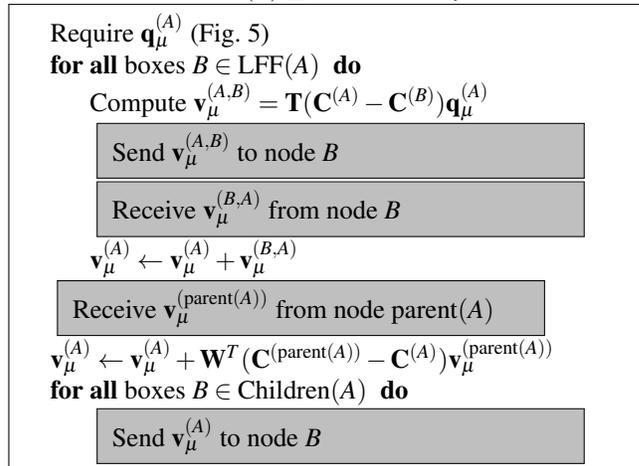
Once the multipole moments have been obtained, the corresponding potential vectors  $\mathbf{v}_\mu^{(A)}$  are computed using the algorithm described in Fig. 6. Each node computes its contribution to the potential in each box  $B$  of  $\text{LFF}(A)$ . This computational step consists of a series of constructions of interaction matrices and the corresponding matrix-vector multiplications with a cost of  $\mathcal{O}(l_{\max}^4)$ . The resulting potential vectors, containing  $(l_{\max} + 1)^2$  entries each, are then sent to the appropriate nodes  $B$ . In this way, each node receives all the necessary contributions which are added up to form the local far-field contribution to  $\mathbf{v}_\mu^{(A)}$ . Then, starting from the lowest level, which is level 2 in practice as the contributions from levels 0 and 1 vanish, because they have no local far field neighbors. Each node sends  $\mathbf{v}_\mu^{(A)}$  to its children, which translate the contribution by means of Eq. (9) and accumulate it into their own  $\mathbf{v}_\mu^{(A)}$ , which is commonly known as the *downward pass*.

## 2.8 Calculation of the direct interactions

For the direct interactions, we use a low-rank separated representation of the Coulomb potential based on the well-known integral expression by Boys and Singer<sup>43,44</sup>:

$$\frac{1}{r_{12}} = \frac{2}{\sqrt{\pi}} \int_0^\infty e^{-t^2 r_{12}^2} dt \approx \sum_{p=1}^R \omega_p e^{-t_p^2 r_{12}^2} + \frac{\pi}{t_f^2} \delta(\mathbf{r}_1 - \mathbf{r}_2). \quad (31)$$

All nodes  $A$  with  $\text{level}(A) \geq 2$  concurrently do



**Fig. 6** Algorithm for calculating  $\mathbf{v}_\mu^{(A)}$ , the potential multipole moments at every box.

We refer to our previous work for further details on how to obtain the quadrature weights  $\omega_p$  and points  $t_p$ .<sup>35</sup> In this work, the operator rank is  $R = 50$  and the chosen quadrature parameters are  $t_l = 2$ ,  $t_f = 50$ ,  $N_{\text{lin}} = 25$ , and  $N_{\text{log}} = 25$ . The parameters are chosen such that the accuracy of the near-field interactions is not a bottleneck. Thus, in practical applications the total number of quadrature points can be significantly reduced.

Calculating the contributions to the potential of the nearest neighbors in Eq. (22) is done as a series of tensor contractions:

$$\left(V_\mu^{(B,A)}\right)_{ijk} \approx \sum_{p=1}^R \omega_p \sum_{k'=1}^{N_z} O_{kk'}^{z,p} \sum_{j'=1}^{N_y} O_{jj'}^{y,p} \sum_{i'=1}^{N_x} O_{ii'}^{x,p} \left(\rho_\mu^{(B)}\right)_{i'j'k'} + \frac{\pi}{t_f^2} \left(\rho_\mu^{(B)}\right)_{ijk} \quad (32)$$

The linear transformation is executed very efficiently using GPGPUs, as we recently have shown.<sup>15</sup> The performance obtained on a single Nvidia K40 card is between 0.5 and 1 TFLOPs depending on the size of the matrices  $\mathbf{O}^{\xi,p}$ .

The matrix elements of the  $\mathbf{O}^{\xi,p}$  matrices are given by

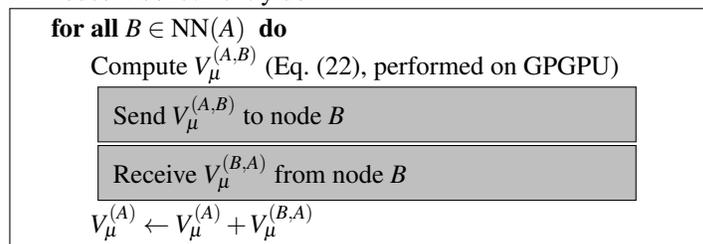
$$O_{ii'}^{\xi,p} = \int_{\xi_{\min}^{(B)}}^{\xi_{\max}^{(B)}} e^{-t_p^2(\xi - \xi_i^{(A)})^2} \chi_{i'}^{(B,\xi)}(\xi) d\xi \quad (33)$$

The algorithm to compute the near-field potentials is summarized in Fig. 7. Each near-field node computes all contributions in their neighbors boxes and sends the data to the neighbor nodes. The computation itself is carried out very efficiently on the GPGPU. The inter-node communication can be expensive: the majority of the nodes have to send and receive  $27N^3$  coefficients, which for typical values of  $N = 100$ – $200$  amounts from hundreds of MB to over 1 GB for 64-bit floating point numbers.

The near-field energy contribution to the energy in Eq. (26) is then computed as

$$\int_{\Omega^{(A)}} \rho_\mu(\mathbf{r}) V_v^{(A)}(\mathbf{r}) d^3r = \sum_k I_k^{(A,z)} \sum_j I_j^{(A,y)} \sum_i I_i^{(A,x)} \left(\rho_\mu^{(A)}\right)_{ijk} \left(V_v^{(A)}\right)_{ijk}, \quad (34)$$

All nodes  $A$  concurrently do



**Fig. 7** Algorithm for calculating  $V_{\mu}^{(A)}(\mathbf{r})$ , the nearest-neighbor contribution to the potentials in all boxes at the highest level.

which can be obtained efficiently using the algorithm for integrating the multipole moments with  $l_{\max} = 0$  shown in Fig. 4.

## 2.9 Calculation of pair interactions

We are now ready to put together all the pieces introduced in this Section into an overall algorithm that computes all pair energies. The final algorithm is presented in Fig. 8. In the initialization step, all nodes compute the translation and interaction matrices that they need. Near-field nodes also compute the necessary Coulomb matrices and the integral vectors. Then, for each input density, each node computes the necessary potential vectors ( $\mathbf{v}_{\mu}^{(A)}$ ) and the potential contributions ( $\mathbf{V}_{\mu}^{(A)}(\mathbf{r})$ ) of the near-field nodes. For each density  $\rho_{\nu}(\mathbf{r})$  with  $1 \leq \nu \leq \mu$ , the node computes its contribution to the interaction energy  $U_{\mu\nu}$ , either as the dot product  $(\mathbf{q}_{\mu}^{(A)})^T \mathbf{v}_{\nu}^{(A)}$  or using the tensor contraction in Eq. (34) for the near-field nodes. All contributions are added up into a master node.

## 3 Results

### 3.1 Benchmark calculations on fullerene systems

The accuracy and parallel performance of the GB-FMM scheme was analyzed in a series of benchmark calculations on fullerene structures that were obtained from Mitsuho Yoshida's fullerene library database.<sup>45</sup> We constructed model densities from linear combinations of Gaussian functions centered on the nuclear coordinates of the fullerene structures:

$$\rho(\mathbf{r}) = \left(\frac{1}{\pi}\right)^{3/2} \sum_K q_K e^{-\alpha_K(\mathbf{r}-\mathbf{R}_K)^2}. \quad (35)$$

The densities were partitioned to the boxes at the highest level of division. For example, when  $D_{\max} = 3$ , the density was represented as an array of size  $8^3 = 512$  containing function objects. These objects contained in turn the information about the near field grids and the numerical meshes. The use of densities of the form in Eq. (35) permitted us to determine the accuracy of the GB-FMM scheme. The so-called self-interaction energy of the charge density (Eq. (1) in the case  $\mu \equiv \nu$ ) was computed and compared with its analytical value. The analytical value was obtained as

All nodes  $A$  concurrently do

All nodes:	Near-field nodes:
<b>for all <math>B \in \text{LFF}(A)</math> do</b> Compute $\mathbf{T}(\mathbf{C}^{(A)} - \mathbf{C}^{(B)})$ <b>for all <math>B \in \text{Children}(A)</math> do</b> Compute $\mathbf{W}(\mathbf{C}^{(B)} - \mathbf{C}^{(A)})$ Compute $\mathbf{W}^T(\mathbf{C}^{(\text{parent}(A))} - \mathbf{C}^{(A)})$	<b>for all <math>B \in \text{NN}(A)</math> do</b> Compute $\mathbf{O}^{\xi,p}$ (Eq. (33)) Compute $\mathbf{I}^{\xi}$ (Eq. (29))

**for all  $\mu : 1 \leq \mu \leq N_p$  do**

All nodes  $A$  concurrently do

All nodes:	Near-field nodes:
Compute $\mathbf{v}_{\mu}^{(A)}$ (Fig. 6)	Compute $V_{\mu}^{(A)}$ (Fig. 7)

**for all  $\nu : 1 \leq \nu \leq \mu$  do**

All nodes  $A$  concurrently do

Compute  $\mathbf{q}_{\nu}^{(A)}$  (Fig. 5)

All near-field nodes  $A$  concurrently do

$$U_{\mu\nu}^{(A)} = \mathbf{v}_{\mu}^{(A)} \cdot \mathbf{q}_{\nu}^{(A)} + \int_{\Omega^{(A)}} V_{\mu}^{(A)} \rho_{\nu}^{(A)} d^3r$$

All near-field nodes reduce:

$$U_{\mu\nu} \leftarrow U_{\mu\nu} + U_{\mu\nu}^{(A)}$$

**Fig. 8** The grid-based FMM algorithm for the calculation of all pair interaction energies from  $N_p$  charge densities.

$$U = \sum_K q_K^2 \sqrt{\frac{2\alpha_K}{\pi}} + 2 \sum_{K>J} \frac{q_K q_J \operatorname{erf}\left(R_{JK} \frac{\alpha_J \alpha_K}{\alpha_J + \alpha_K}\right)}{R_{JK}}, \quad (36)$$

where  $R_{JK} = |\mathbf{R}_J - \mathbf{R}_K|$ .

The Gaussian amplitudes were set to the nuclear charge of carbon,  $q_K \equiv 6$ , and the exponents were set to unity,  $\alpha_K \equiv 1$ . These parameters ensured that the charge was reproduced and, more importantly, that all model charge densities were globally smooth.

The GB-FMM scheme was implemented in Fortran and parallelized with MPI and cuBLAS, the CUDA-implementation of the BLAS library. One single-core physical node was responsible for computing the far-field potentials with the FMM, while the near field potential calculations were divided among hybrid nodes, each equipped with a GPGPU card (NVIDIA Tesla 40K cards that are controlled by Intel Xeon E5-2620-v2 CPUs). All parallel runs were performed on the Taito GPGPU cluster of the CSC computing center.

### 3.2 Accuracy

The fullerene  $C_{60}$  was picked as the test structure for studying how the accuracy of the GB-FMM scheme depends on the grid step  $h$ , the truncation of the multipole expansion  $l_{\max}$ , and the depth of the octree  $D_{\max}$ . The domain size was adjusted for ensuring that the density was negligible at the domain boundaries. In the case of  $C_{60}$ , the cubic domain had a side length of  $24a_0$ . The values  $D_{\max} = 2, 3, 4$  and  $l_{\max} = 5, 7, 10, 12, 15$  were considered. The grid step was set to  $h = a_0/2^i$  with  $i = 1, 2, 3, 4$ .

In Fig. 9, the absolute error in the self-interaction energy is plotted against the grid refinement parameter  $i$ . The error is seen to be saturated below one millihartree in the case of  $D_{\max} = 2$ , while the  $10^{-5}E_h$  level is reached with  $D_{\max} = 3$  and  $D_{\max} = 4$ . These absolute energies translate to relative accuracies in the range of 1 to 10 ppb, which agree well with our previous results.<sup>34,35</sup>

The error in the energy saturates with respect to both the truncation of the bipolar multipole expansion  $l_{\max}$  and the grid step  $h$ , which implies that the remaining error results from the quadrature parameters of the direct integration, *e.g.*, the discrete representation of the Coulomb operator.

### 3.3 Performance

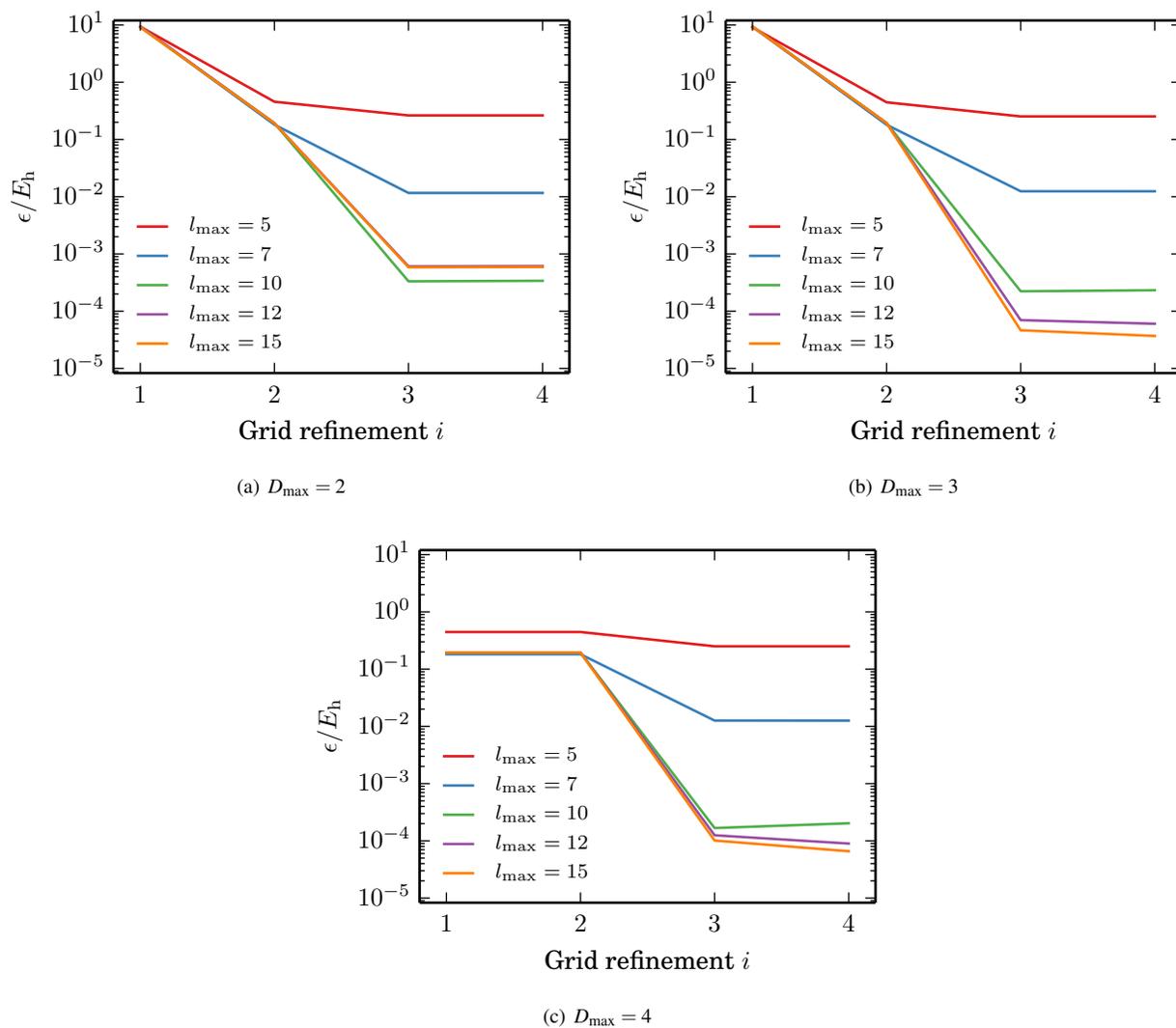
The evaluation of the electrostatic potential is by far the most time consuming step in the GB-FMM scheme. As the near-field and far-field potentials are computed concurrently by construction, the timing is

$$t_{\text{pot}} = \max(t_{\text{NF}}, t_{\text{FF}}). \quad (37)$$

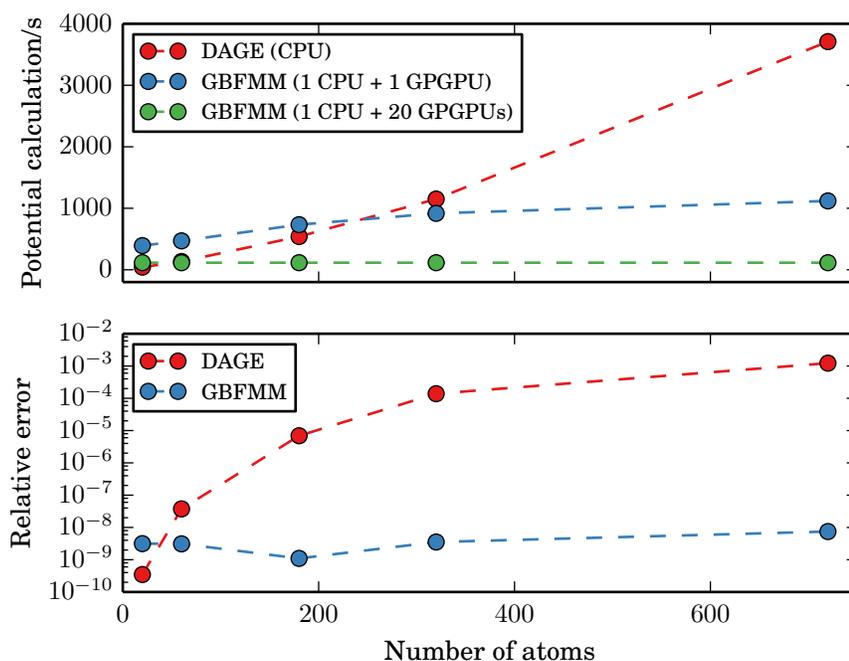
The near-field contribution is  $8^{D_{\max}} t_{\text{DAGE}}$ , where  $t_{\text{DAGE}}$  is the time for constructing a single near field potential with the direct integration. The far-field contribution consists of evaluating the  $8^{D_{\max}}$  potential vectors, *i.e.*, the expansion coefficients of the far-field potentials.

In Fig. 10, we report timings for calculation of the the potentials for the  $C_{20}$ ,  $C_{60}$ ,  $C_{180}$ ,  $C_{320}$  and  $C_{720}$  systems together with the absolute relative error. The domain sizes varied between  $18a_0 \times 19a_0 \times 18a_0$  and  $59a_0 \times 56a_0 \times 56a_0$ . The GB-FMM specific parameters were  $D_{\max} = 3$  and  $l_{\max} = 15$ . The grid step and the quadrature parameters were set to  $0.1a_0$  and  $N_{\text{lin}} = N_{\text{log}} = 25$ , respectively.

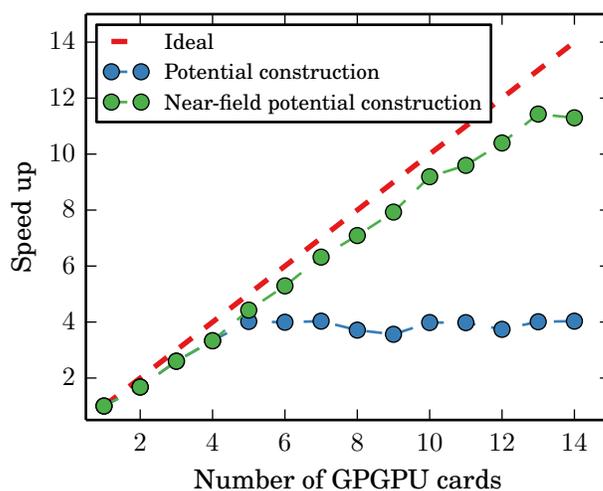
While the cost of the direct integration of the potential on a single-core CPU is competitive with the GB-FMM part for small systems, the GB-FMM scheme scales significantly better with the size of the



**Fig. 9** Accuracy of the GB-FMM with respect to grid quality, truncation of the multipole expansion and the depth of the octree. The grid refinement parameter  $i$  gives the grid step from  $h = a_0/2^i$ .



**Fig. 10** Time of the potential calculation as a function of system size. The data points correspond to the fullerene structures  $C_{20}$ ,  $C_{60}$ ,  $C_{180}$ ,  $C_{320}$  and  $C_{720}$ .



**Fig. 11** Speed up of potential calculation in  $C_{60}$  as a function of GPGPU cards. While the near field potential calculation parallelizes well by construction, the Amdahl's law limit is reached already with 4 to 5 GPGPU cards as the serial far-field potential calculation forms a bottleneck.

system, both in terms of accuracy and performance. In the GB-FMM calculation, a constant, ppb-level relative accuracy is obtained for all system sizes, whereas the accuracy of the direct numerical integration deteriorates with increasing size of the domain.

The results obtained for the GB-FMM scheme can be expected as the accuracy of the scheme should not depend on the size of the computational domain, because the entire domain is never simultaneously considered. In addition, as all the direct numerical integration parts of the calculations within the GB-FMM scheme are performed on domain sizes of a couple of atomic units per side, the GB-FMM should retain the accuracy characteristic (0.1–1 ppb) of the direct numerical integration for that size of the grid.

When the near field potential calculation is distributed among 20 GPGPU nodes, the GB-FMM scheme is clearly superior to the fully numerical integration method and exhibits effectively constant scaling of the computational time with respect to system size, which is a consequence of reaching the limit of Amdahl's law. When all GPGPU nodes performing the numerical integration complete their task faster than the CPU performing GB-FMM calculation, the running time of the potential calculation is dominated by the single-core GB-FMM part of the calculation. The calculation can be speeded up by also distributing the GB-FMM part on several CPUs or GPGPUs.

Fig. 11 shows in detail how the limit of Amdahl's law is reached for the  $C_{60}$  system. The employed parameters were the same as used in the calculations reported in Fig. 10. With  $l_{\max} = 15$  and  $D_{\max} = 3$ , the GB-FMM part becomes the bottleneck when the near field potential calculations are divided among four to five nodes. The calculations of the near-field potentials scale nearly ideally as expected.

#### 4 Conclusions and outlook

A grid-based fast multipole method (GB-FMM) scheme for calculating two-electron interaction energies has been developed and implemented. The computational domain is divided into subdomains that can be assigned to nearest neighbors and more distant subdomains. The far-field contributions to the two-electron interaction energies are obtained by employing the bipolar series expansion of the Coulomb operator, whereas the near-field contributions are obtained by fully numerical integration of the corresponding Coulomb energy integral. In the GB-FMM scheme, the multipole moments are calculated for the individual subdomains and combined to multipole moments for groups of subdomains. The series expansions of the charge density and the electrostatic potential are combined to yield the electrostatic interaction energy of separated non-overlapping charge densities. The numerical integration algorithm has been adapted for general purpose graphics processing (GPGPU) units. The performance of the computational scheme has been demonstrated by calculations using one CPU for the serial part including the GB-FMM and up to 20 GPGPUs for the numerical integration of the Coulomb energy. For the largest GPGPU cluster, the parallel part of the code is faster than the serial one and since the computing time of the serial part of the code is more or less independent of the system size, the computational time is found to be independent of the size of the system. The present calculations show that numerical calculations using local basis functions can be made to run very efficiently on massively parallel computers, because for local basis functions the computational domain can be easily subdivided into smaller units whose interactions can be computed independently. The algorithm yields a computational method that formally scales linearly with the system size, whereas the massively parallel computer architecture renders calculations whose wall time is independent of the system size feasible. The scalability, performance and accuracy of the present numerical calculations on the GPGPU cluster suggest that quantum chemistry calculations of the future will most likely be made this way using local basis functions.

---

## Acknowledgment

This work has been supported by the Academy of Finland through projects (137460 and 275845) and its Computational Science Research Programme (LASTU/258258) and by the Magnus Ehrnrooth Foundation. CSC – the Finnish IT Center for Science – is acknowledged for computer time. We thank Jonas Jusélius for support with the development and maintenance of the numerical code.

## Appendix

### A The interaction and translation matrices

Let us introduce the auxiliary functions

$$A_l^m = \begin{cases} l! & m = 0 \\ (-1)^m \sqrt{\frac{(l-m)!(l+m)!}{2}} & m > 0 \\ (-1)^m A_l^{-m} & m < 0 \end{cases} \quad (38)$$

$$B_l^m = \begin{cases} 0 & m = 0 \\ (-1)^m \sqrt{\frac{(l-m)!(l+m)!}{2}} & m > 0 \\ (-1)^{m+1} B_l^{-m} & m < 0 \end{cases} \quad (39)$$

$$C_l^m = \frac{(-1)^m}{\sqrt{(1 + \delta_{m0})(l+m)!(l-m)!}}, \quad (40)$$

$$T_{lm,jk} = \begin{cases} (-1)^j C_l^m C_j^k \left[ A_{l+j}^{m+k} I_{l+j}^{|m+k|} + (-1)^k A_{l+j}^{m-k} I_{l+j}^{|m-k|} \right] & m \geq 0, k \geq 0 \\ (-1)^j C_l^m C_j^k \left[ B_{l+j}^{m-k} I_{l+j}^{-|m-k|} - (-1)^k B_{l+j}^{m+k} I_{l+j}^{-|m+k|} \right] & m \geq 0, k < 0 \\ (-1)^j C_l^m C_j^k \left[ (-1)^{m+1} B_{l+j}^{m+k} I_{l+j}^{-|m+k|} - (-1)^{k+m} B_{l+j}^{m-k} I_{l+j}^{-|m-k|} \right] & m < 0, k \geq 0 \\ (-1)^j C_l^m C_j^k \left[ (-1)^m A_{l+j}^{m-k} I_{l+j}^{|m-k|} - (-1)^{m+k} A_{l+j}^{m+k} I_{l+j}^{|m+k|} \right] & m < 0, k < 0. \end{cases} \quad (41)$$

Conversion to a real-valued formulation is once again effected. We introduce a second set of auxiliary functions

$$\alpha_l^m = \begin{cases} \frac{1}{l!} & m = 0 \\ (-1)^m \frac{1}{\sqrt{2(l+m)!(l-m)!}} & m > 0 \\ (-1)^m \alpha_l^{-m} & m < 0 \end{cases} \quad (42)$$

$$\beta_l^m = \begin{cases} 0 & m = 0 \\ (-1)^m \frac{1}{\sqrt{2(l+m)!(l-m)!}} & m > 0 \\ (-1)^{(m+1)} \beta_l^{-m} & m < 0 \end{cases} \quad (43)$$

$$\gamma_l^m = (-1)^m \sqrt{(2 - \delta_{m0})(l-m)!(l+m)!}, \quad (44)$$

whereby the translation matrix is given by

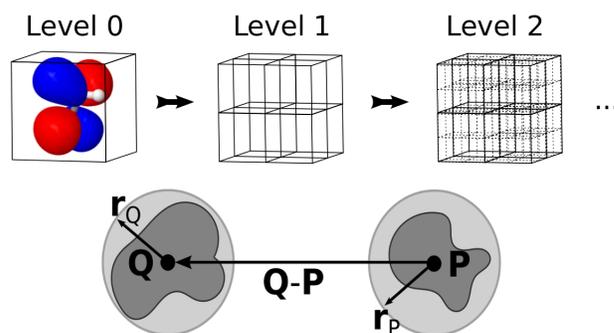
$$W_{lm,jk} = \begin{cases} \left(\frac{1}{2}\right) \delta_{k0} \left(\frac{\gamma_l^m}{\gamma_j^k}\right) \left[\alpha_{l-j}^{m-k} S_{l-j}^{|m-k|} + (-1)^k \alpha_{l-j}^{m+k} S_{l-j}^{|m+k|}\right] & m \geq 0, k \geq 0 \\ \left(\frac{\gamma_l^m}{\gamma_j^k}\right) \left[(-1)^k \beta_{l-j}^{m-k} S_{l-j}^{-|m-k|} - \beta_{l-j}^{m+k} S_{l-j}^{-|m+k|}\right] & m \geq 0, k < 0 \\ \left(\frac{1}{2}\right) \delta_{k0} \left(\frac{\gamma_l^m}{\gamma_j^k}\right) \left[(-1)^{m+1} \beta_{l-j}^{m-k} S_{l-j}^{-|m-k|} - (-1)^{m+k} \beta_{l-j}^{m+k} S_{l-j}^{-|m+k|}\right] & m < 0, k \geq 0 \\ \left(\frac{\gamma_l^m}{\gamma_j^k}\right) \left[(-1)^{m-k} \alpha_{l-j}^{m-k} S_{l-j}^{|m-k|} - (-1)^m \alpha_{l-j}^{m+k} S_{l-j}^{|m+k|}\right] & m < 0, k < 0. \end{cases} \quad (45)$$

## References

- 1 G. E. Moore, *Electron. Mag.*, 1965, 38.
- 2 A. V. Titov, I. S. Ufimtsev, N. Luehr and T. J. Martinez, *J. Chem. Theory Comput.*, 2013, **9**, 213–221.
- 3 K. Bhaskaran-Nair, W. Ma, S. Krishnamoorthy, O. Villa, H. J. J. van Dam, E. Aprà and K. Kowalski, *J. Chem. Theory Comput.*, 2013, **9**, 1949–1957.
- 4 J. E. Stone, J. C. Phillips, P. L. Freddolino, D. J. Hardy, L. G. Trabuco and K. Schulten, *J. Comp. Chem.*, 2007, **28**, 2618–2640.
- 5 K. Yasuda, *J. Comp. Chem.*, 2008, **29**, 334–342.
- 6 J. A. Anderson, C. D. Lorenz and A. Travesset, *J. Comp. Phys.*, 2008, **227**, 5342–5359.
- 7 L. Vogt, R. Olivares-Amaya, S. Kermes, Y. Shao, C. Amador-Bedolla and A. Aspuru-Guzik, *J. Phys. Chem. A*, 2008, **112**, 2049–2057.
- 8 M. S. Friedrichs, P. Eastman, V. Vaidyanathan, M. Houston, S. Legrand, A. L. Beberg, D. L. Ensign, C. M. Bruns and V. S. Pande, *J. Comp. Chem.*, 2009, **30**, 864–872.
- 9 J. van Meel, A. Arnold, D. Frenkel, S. Portegies Zwart, and R. Belleman, *Mol. Simul.*, 2008, **34**, 259–266.
- 10 A. Asadchev, V. Allada, J. Felder, B. M. Bode, M. S. Gordon and T. L. Windus, *J. Chem. Theory Comput.*, 2010, **6**, 696–704.
- 11 R. Olivares-Amaya, M. A. Watson, R. G. Edgar, L. Vogt, Y. Shao and A. Aspuru-Guzik, *J. Chem. Theory Comput.*, 2010, **6**, 135–144.
- 12 B. G. Levine, D. N. LeBard, R. DeVane, W. Shinoda, A. Kohlmeyer and M. L. Klein, *J. Chem. Theory Comput.*, 2011, **7**, 4135–4145.
- 13 X. Wu, A. Koslowski and W. Thiel, *J. Chem. Theory Comput.*, 2012, **8**, 2272–2281.
- 14 S. A. Maurer, J. Kussmann and C. Ochsenfeld, *J. Chem. Phys.*, 2014, **141**, 051106.
- 15 S. Losilla, M. A. Watson, A. Aspuru-Guzik and D. Sundholm, *J. Chem. Theory Comput.*, 2015.
- 16 K. Yasuda and H. Maruoka, *Intern. J. Quant. Chem.*, 2014, **114**, 543–552.
- 17 G. M. Amdahl, Proc. of the AFIPS '67 Spring Joint Computer Conference, 1967, pp. 483–485.
- 18 C. F. Fischer, *Hartree-Fock Method for Atoms. A Numerical Approach*, Wiley, New York, USA, 1977.
- 19 T. Koga, S. Watanabe, K. Kanayama, R. Yasuda and A. J. Thakkar, *J. Chem. Phys.*, 1995, **103**, 3000–3005.
- 20 L. Laaksonen, P. Pyykkö and D. Sundholm, *Comput. Phys. Rep.*, 1986, **4**, 315–344.
- 21 J. Kobus, L. Laaksonen and D. Sundholm, *Comput. Phys. Commun.*, 1996, **98**, 346–358.

- 
- 22 F. Jensen, *Theoret. Chim. Acta*, 2005, **113**, 187–190.
- 23 F. A. Bischoff, R. J. Harrison and E. F. Valeev, *J. Chem. Phys.*, 2012, **137**, 104103.
- 24 F. A. Bischoff and E. F. Valeev, *J. Chem. Phys.*, 2013, **139**, 114106.
- 25 J. M. Pérez-Jordá and W. Yang, *J. Chem. Phys.*, 1997, **107**, 1218–1226.
- 26 C. H. Choi, K. Ruedenberg and M. S. Gordon, *J. Comp. Chem.*, 2001, **22**, 1484–1501.
- 27 M. A. Watson, P. Sałek, P. Macak and T. Helgaker, *J. Chem. Phys.*, 2004, **121**, 2915–2931.
- 28 H. Dachsels, *J. Chem. Phys.*, 2010, **132**, 119901.
- 29 E. Rudberg and P. Sałek, *J. Chem. Phys.*, 2006, **125**, 84106.
- 30 L. Greengard and V. Rokhlin, *J. Comp. Phys.*, 1987, **73**, 325–348.
- 31 C. A. White and M. Head-Gordon, *J. Chem. Phys.*, 1994, **101**, 6593–6605.
- 32 C. A. White, B. G. Johnson, P. M. Gill and M. Head-Gordon, *Chem. Phys. Letters*, 1994, **230**, 8–16.
- 33 D. Sundholm, *J. Chem. Phys.*, 2005, **122**, 194107.
- 34 J. Jusélius and D. Sundholm, *J. Chem. Phys.*, 2007, **126**, 094101.
- 35 S. A. Losilla, D. Sundholm and J. Jusélius, *J. Chem. Phys.*, 2010, **132**, 024102.
- 36 S. A. Losilla and D. Sundholm, *J. Chem. Phys.*, 2012, **136**, 214104.
- 37 P. García-Risueño, J. Alberdi-Rodríguez, M. J. Oliveira, X. Andrade, M. Pippig, J. Muguerza, A. Arruabarrena and A. Rubio, *J. Comp. Chem.*, 2014, **35**, 427–444.
- 38 J. Kussmann, M. Beer and C. Ochsenfeld, *WIREs Comput. Mol. Sci.*, 2013, **3**, 614–636.
- 39 B. C. Carlson and G. S. Rushbrooke, *Math. Proc. Camb. Phil. Soc.*, 1950, **46**, 626–633.
- 40 E. Steinborn and K. Ruedenberg, *Adv. Quantum Chem*, 1973, **7**, 1–81.
- 41 T. Helgaker, P. Jørgensen and J. Olsen, *Molecular Electronic-structure Theory*, Wiley, New York, USA, 2000.
- 42 S. Losilla, M. Mehine and D. Sundholm, *Mol. Phys.*, 2012, **110**, 2569–2578.
- 43 K. Singer, *Proc. R. Soc. A*, 1960, **258**, 412–420.
- 44 S. F. Boys, *Proc. R. Soc. A*, 1960, **258**, 402–411.
- 45 S. Weber, <http://www.jcrystal.com/steffenweber/gallery/Fullerenes/Fullerenes.html>, (Mitsuho Yoshida's Fullerene Library).

## Graphical abstract



A grid-based fast multipole method have been developed for calculating two-electron interaction energies for non-overlapping charge densities.