

Dedicated to Professor Andy Majda on the occasion of his 60th birthday

## POLE-BASED APPROXIMATION OF THE FERMI-DIRAC FUNCTION

LIN LIN, JIANFENG LU, LEXING YING, AND WEINAN E

ABSTRACT. Two approaches for the efficient rational approximation of the Fermi-Dirac function are discussed: one uses the contour integral representation and conformal mapping and the other is based on a version of the multipole representation of the Fermi-Dirac function that uses only simple poles. Both representations have logarithmic computational complexity. They are of great interest for electronic structure calculations.

### 1. INTRODUCTION

Given an effective one-particle Hamiltonian  $\mathbf{H}$ , the inverse temperature  $\beta = 1/k_{\text{B}}T$  and the chemical potential  $\mu$ , the finite temperature single-particle density matrix of the system is given by the Fermi operator

$$(1) \quad \rho = 2(1 + \exp(\beta(\mathbf{H} - \mu)))^{-1} = 1 - \tanh\left(\frac{\beta}{2}(\mathbf{H} - \mu)\right),$$

where  $\tanh$  is the hyperbolic tangent function.

In the last decade or so, the development of accurate and numerically efficient representations of the Fermi operator has attracted a great deal of attention in the quest for linear scaling electronic structure methods based on effective one-electron Hamiltonians. These approaches have a numerical cost that scales linearly with  $N$ , the number of electrons, and thus hold the promise of making electronic structure analysis of large systems feasible. Achieving linear scaling for realistic systems is very challenging. Formulations based on the Fermi operator are appealing since this operator gives directly the single particle density matrix without the need for diagonalizing the Hamiltonian.

From a computational viewpoint, one main issue is that the right hand side of (1) is an operator-valued function. To evaluate this function, we have to replace or approximate it by something which can be computed directly without diagonalization. Obvious candidates are polynomial or rational approximations. Such an approach was first introduced by Baroni and Giannozzi [1] and Goedecker and co-workers [4, 5] (see also the review article [6]). Several improvements have been made since then, for example, in [2, 3, 9–11, 14]. These are put broadly under the umbrella of “Fermi operator expansion” (abbreviated as FOE).

From the viewpoint of efficiency, a major concern is the cost for representing the Fermi operator as a function of  $\beta\Delta E$  (for finite temperature) or  $\Delta E/E_g$  (for gapped systems) where  $\beta$  is the inverse temperature,  $\Delta E$  is the spectral width of the discretized Hamiltonian matrix and  $E_g$  is the spectrum gap of the Hamiltonian around the chemical potential. Consider a finite temperature gapless system for example, the cost of the FOE proposed by Goedecker *et al* scales as  $\beta\Delta E$ . The fast polynomial summation technique introduced by Head-Gordon *et al* [10, 11] reduces the cost to  $(\beta\Delta E)^{1/2}$ . The cost of the hybrid algorithm proposed by Parrinello *et*

$al$  in a recent preprint [2] scales as  $(\beta\Delta E)^{1/3}$ . The cost was brought down to logarithmic scaling  $\ln(\beta\Delta E)$  in [12] using a multipole representation of the Matsubara expansion of the Fermi-Dirac function.

The purpose of this article is to introduce two alternative rational expansions of the Fermi-Dirac function that use only simple poles and have computational cost that scales logarithmically. The first strategy is to use the contour integral and conformal mapping idea proposed recently in [8]. This will be presented in the next section. The other strategy is to borrow ideas from [15] and use a version of multipole expansion [12] that only involves simple poles. This will be discussed in Section 3. Numerical examples illustrating the efficiency and accuracy of the representations are discussed in Section 4.

## 2. RATIONAL EXPANSIONS BASED ON CONTOUR INTEGRAL

Our first approach is an adaptation of the ideas proposed recently in [8] based on contour integral representation and conformal mapping. Let us first briefly recall the main idea of [8]. Consider a function  $f$  that is analytic in  $\mathbb{C}\setminus(-\infty, 0]$  and an operator  $\mathbf{A}$  with spectrum in  $[m, M] \subset \mathbb{R}^+$ , one wants to evaluate  $f(\mathbf{A})$  using a rational expansion of  $f$  by discretizing the contour integral

$$(2) \quad f(\mathbf{A}) = \frac{1}{2\pi i} \int_{\Gamma} f(z)(z - \mathbf{A})^{-1} dz.$$

The innovative technique in [8] was to construct a conformal map that maps the stripe  $S = [-K, K] \times [0, K']$  to the upper half (denoted as  $\Omega^+$ ) of the domain  $\Omega = \mathbb{C}\setminus((-\infty, 0] \cup [m, M])$ . This special map from  $t \in S$  to  $z \in \Omega^+$  is given by

$$(3) \quad z = \sqrt{mM} \left( \frac{k^{-1} + u}{k^{-1} - u} \right), \quad u = \operatorname{sn}(t) = \operatorname{sn}(t|k), \quad k = \frac{\sqrt{M/m} - 1}{\sqrt{M/m} + 1}.$$

Here  $\operatorname{sn}(t)$  is one of the Jacobi elliptic functions and the numbers  $K$  and  $K'$  are complete elliptic integrals whose values are given by the condition that the map is from  $S$  to  $\Omega^+$ .

Applying the trapezoidal rule with  $Q$  equally spaced points in  $(-K + iK'/2, K + iK'/2)$ ,

$$(4) \quad t_j = -K + \frac{iK'}{2} + 2\frac{(j - \frac{1}{2})K}{Q}, \quad 1 \leq j \leq Q,$$

we get the quadrature rule (denote  $z_j = z(t_j)$ )

$$(5) \quad f_Q(\mathbf{A}) = \frac{-4K\sqrt{mM}}{\pi Qk} \operatorname{Im} \sum_{j=1}^Q \frac{f(z_j)(z_j - \mathbf{A})^{-1} \operatorname{cn}(t_j) \operatorname{dn}(t_j)}{(k^{-1} - \operatorname{sn}(t_j))^2}.$$

Here  $\operatorname{cn}$  and  $\operatorname{dn}$  are the other two Jacobi elliptic functions in standard notation and the factor  $\operatorname{cn}(t_j) \operatorname{dn}(t_j)(k^{-1} - \operatorname{sn}(t_j))^{-2} \sqrt{mM}/k$  comes from the Jacobian of the function  $z(t)$ .

It is proved in [8] that the convergence is exponential in the number of quadrature points  $Q$  and the exponent deteriorates only logarithmically as  $M/m \rightarrow \infty$ :

$$(6) \quad \|f(\mathbf{A}) - f_Q(\mathbf{A})\| = \mathcal{O}(e^{-\pi^2 Q / (\log(M/m) + 3)}).$$

To adapt the idea to our setting with the Fermi-Dirac function or the hyperbolic tangent function, we face with two differences: First, the  $\tanh$  function has singularities on the imaginary axis. Second, the operator we are considering,  $\beta(\mathbf{H} - \mu)$ , has spectrum on both the negative and positive axis.

**2.1. Gapped case.** We first consider the case when the Hamiltonian  $\mathbf{H}$  has a gap in its spectrum around the chemical potential  $\mu$ , such that  $\text{dist}(\mu, \sigma(\mathbf{H})) = E_g > 0$ . Physically, this will be the case when the system is an insulator.

Let us consider  $f(z) = \tanh(\frac{\beta}{2}z^{1/2})$  acting on the operator  $\mathbf{A} = (\mathbf{H} - \mu)^2$ . Now,  $f(z)$  has singularities only on  $(-\infty, 0]$  and the spectrum of  $\mathbf{A}$  is contained in  $[E_g^2, E_M^2]$ , where

$$E_M = \max_{E \in \sigma(\mathbf{H})} |E - \mu|.$$

We note that obviously  $E_M \leq \Delta E$ . Hence we are back in the same scenario as considered in [8] except that we need to take care of different branches of the square root function when we apply the quadrature rule.

More specifically, we construct the contour and quadrature points  $z_j$  in the  $z$ -plane using parameters  $m = E_g^2$  and  $M = E_M^2$ . Denote  $g(\xi) = \tanh(\beta\xi/2)$ ,  $\xi_j^\pm = \pm z_j^{1/2}$ , and  $\mathbf{B} = \mathbf{H} - \mu$ . The quadrature rule is then given by

$$(7) \quad g_Q(\mathbf{B}) = \frac{-2K\sqrt{mM}}{\pi Qk} \text{Im} \left( \sum_{j=1}^Q \frac{g(\xi_j^+) (\xi_j^+ - \mathbf{B})^{-1} \text{cn}(t_j) \text{dn}(t_j)}{\xi_j^+ (k^{-1} - \text{sn}(t_j))^2} + \sum_{j=1}^Q \frac{g(\xi_j^-) (\xi_j^- - \mathbf{B})^{-1} \text{cn}(t_j) \text{dn}(t_j)}{\xi_j^- (k^{-1} - \text{sn}(t_j))^2} \right),$$

where the factors  $\xi_j^\pm$  in the denominator come from the Jacobian of the map from  $z$  to  $\xi$ . The number of poles to be inverted is  $N_{\text{pole}} = 2Q$ . After applying (6), we have a similar error estimate for  $g(\mathbf{B})$

$$(8) \quad \|g(\mathbf{B}) - g_Q(\mathbf{B})\| = \mathcal{O}(e^{-\pi^2 Q / (2 \log(E_M/E_g) + 3)}).$$

In Fig. 1, a typical configuration of the quadrature points is shown. The x-axis is taken to be  $E - \mu$ . We see that in this case the contour consists of two loops, one around the spectrum below the chemical potential and the other around the spectrum above the chemical potential.

Note further that as the temperature goes to zero, the Fermi-Dirac function converges to the step function:

$$(9) \quad \eta(\xi) = \begin{cases} 2, & \xi \leq 0, \\ 0, & \xi > 0. \end{cases}$$

Therefore, the contribution of the quadrature points  $\xi_j^+$  on the right half plane ( $\text{Re } \xi_j^+ > 0$ ) is negligible when  $\beta$  is large. In particular, for the case of zero temperature, one may choose only the quadrature points on the left half plane. The quadrature formula we obtain then becomes

$$(10) \quad \eta_Q(\mathbf{B}) = \frac{-4K\sqrt{mM}}{\pi Qk} \text{Im} \left( \sum_{j=1}^Q \frac{(\xi_j^- - \mathbf{B})^{-1} \text{cn}(t_j) \text{dn}(t_j)}{\xi_j^- (k^{-1} - \text{sn}(t_j))^2} \right).$$

The number of poles to be inverted is then  $N_{\text{pole}} = Q$ .

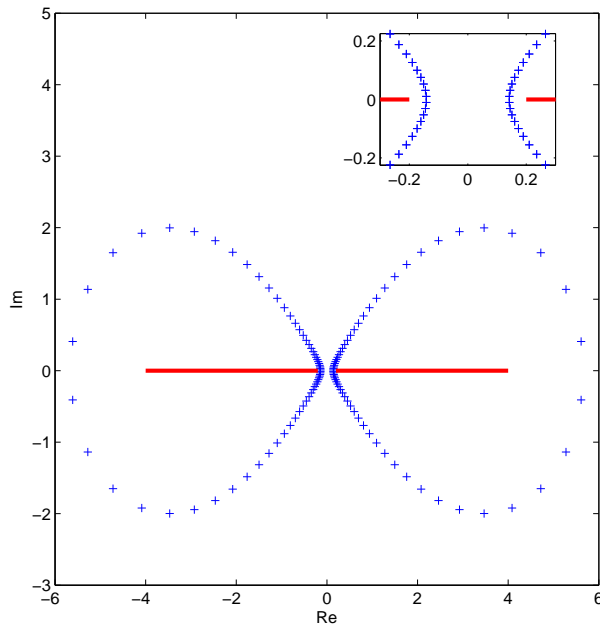


FIGURE 1. A typical configuration of the poles on a two-loop contour.  $Q = 30$ ,  $E_g = 0.2$ ,  $E_M = 4$  and  $\beta = 1000$ . The red line indicates the spectrum. The inset shows the poles close to the origin. The x-axis is  $E - \mu$  with  $E$  the eigenvalue of  $\mathbf{H}$ . The poles with negative imaginary parts are not explicitly calculated.

We show in Fig. 2 a typical configuration of the set of quadrature points. Only one loop is required compared with Fig. 1.

**2.2. Gapless case.** The more challenging case is when the spectrum of  $\mathbf{H}$  does not have a gap, *i.e.*,  $E_g = 0$ . Physically, this corresponds to the case of metallic systems. In this case, the construction discussed in the last subsection does not work.

To overcome this problem, we note that the hyperbolic tangent function  $\tanh(\frac{\beta}{2}z)$  is analytic except at poles  $(2l-1)\pi/\beta i$ ,  $l \in \mathbb{Z}$  on the imaginary axis. Therefore, we could construct a contour around the whole spectrum of  $\mathbf{H}$  which passes through the imaginary axis on the upper half plane between the origin and  $\pi/\beta i$  and also on the lower half plane between the origin and  $-\pi/\beta i$ . Thus, we will have a dumbbell shaped contour as shown in Fig. 3.

To be more specific, let us first construct the contour and quadrature points  $z_j$  in the  $z$ -plane as in the last subsection using parameters  $m = \pi^2/\beta^2$  and  $M = E_M^2 + \pi^2/\beta^2$ . Denote  $\xi_j^\pm = \pm(z_j - \pi^2/\beta^2)^{1/2}$ ,  $g = \tanh(\beta\xi/2)$  and  $\mathbf{B} = \mathbf{H} - \mu$ .

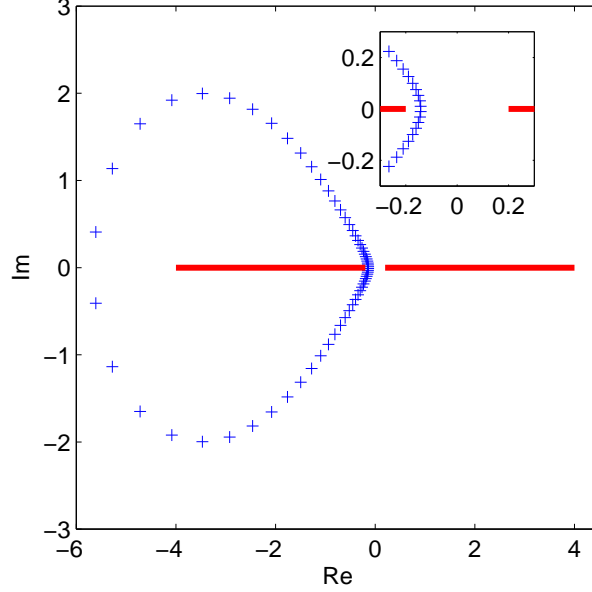


FIGURE 2. A typical configuration of the poles for zero temperature ( $\beta = \infty$ ).  $Q = 30$ ,  $E_g = 0.2$  and  $E_M = 4$ . The red line indicates the spectrum. The inset zooms into the poles that is close to the origin. The x-axis is  $E - \mu$  with  $E$  the eigenvalue of  $\mathbf{H}$ . The poles with negative imaginary parts are not explicitly calculated.

The quadrature rule takes the following form

$$(11) \quad g_Q(\mathbf{B}) = \frac{-2K\sqrt{mM}}{\pi Qk} \operatorname{Im} \left( \sum_{j=1}^Q \frac{g(\xi_j^+) (\xi_j^+ - \mathbf{B})^{-1} \operatorname{cn}(t_j) \operatorname{dn}(t_j)}{\xi_j^+ (k^{-1} - \operatorname{sn}(t_j))^2} + \sum_{j=1}^Q \frac{g(\xi_j^-) (\xi_j^- - \mathbf{B})^{-1} \operatorname{cn}(t_j) \operatorname{dn}(t_j)}{\xi_j^- (k^{-1} - \operatorname{sn}(t_j))^2} \right).$$

When apply the quadrature formula, the number of poles to be inverted is  $N_{\text{pole}} = 2Q$ . Fig. 3 shows a typical configuration of quadrature points for  $Q = 30$ . The map  $\xi(z) = (z - \pi^2/\beta^2)^{1/2}$  maps the circle in the  $z$ -plane to a dumbbell-shaped contour (put two branches together).

Actually, what is done could be understood as follows. Similar to [8], we have constructed a map from the rectangular domain  $[-3K, K] \times [0, K']$  to the upper half of the domain

$$U = \{z \mid \operatorname{Im} z \geq 0\} \setminus ([-E_M, E_M] \cup i[\pi/\beta, \infty)).$$

The map is carried out in three steps, shown in Fig. 4. The first two steps use the original map constructed in [8], however with extended domain of definition. First,

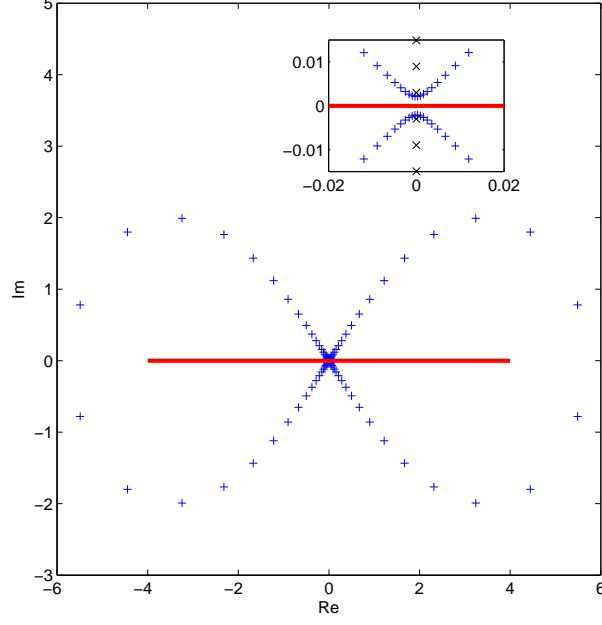


FIGURE 3. A typical configuration of the poles on a dumbbell-shaped contour.  $Q = 30$ ,  $E_g = 0$ ,  $E_M = 4$  and  $\beta = 1000$ . The inset zooms into the part close to the origin. The red line indicates the spectrum. The black crosses indicate the positions of the poles of  $\tanh$  function on the imaginary axis. The poles with negative imaginary parts are not explicitly calculated.

the Jacobi elliptic function

$$(12) \quad u = \operatorname{sn}(t) = \operatorname{sn}(t|k), \quad k = \frac{\sqrt{M/m} - 1}{\sqrt{M/m} + 1}$$

maps the rectangular domain to the complex plane, with the ends mapping to  $[1, k^{-1}]$  and the middle vertical line  $-K + i[0, K']$  to  $[-k^{-1}, -1]$ . Then, the Möbius transformation

$$(13) \quad z = \sqrt{mM} \left( \frac{k^{-1} + u}{k^{-1} - u} \right)$$

maps the complex plane to itself in such a way that  $[-k^{-1}, -1]$  and  $[1, k^{-1}]$  are mapped to  $[0, m]$  and  $[M, \infty]$ , respectively. Finally, the shifted square root function

$$(14) \quad \xi = (z - m)^{1/2}$$

maps the complex plane to the upper-half plane (we choose the branch of the square root such that the lower-half plane is mapped to the second quadrant and the upper-half plane is mapped to the first quadrant), in such a way that  $[0, m]$  is sent to  $i[0, \sqrt{m}]$  and  $[M, \infty)$  is sent to  $(-\infty, -\sqrt{M - m}] \cup [\sqrt{M - m}, \infty)$ . The map can be extended to a map from  $[-7K, K] \times [0, K']$  to the whole  $U$ , in this case, the  $z$ -plane becomes a double-covered Riemann surface with branch point at  $m$ .

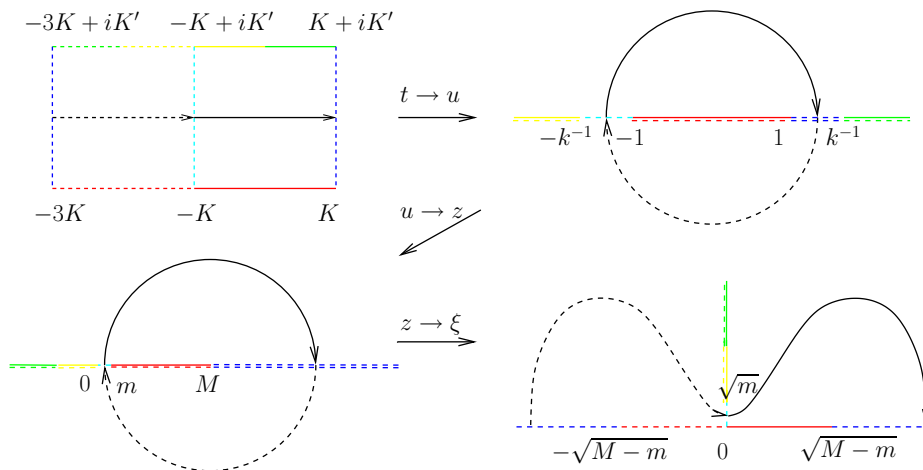


FIGURE 4. The map from the rectangular domain  $[-3K, K] \times [0, K']$  to the upper-half of the domain  $U$ . The map is constructed in three steps:  $t \rightarrow u \rightarrow z \rightarrow \xi$ . The boundaries are shown in various colors and line styles.

Since the function  $g$  is analytic in the domain  $U$ , the composite function  $g(t) = g(\xi(z(u(t))))$  is analytic in the stripe in the  $t$ -plane, and therefore, the trapezoidal rule converges exponentially fast. Using a similar analysis that leads to (6), it can be shown that

$$(15) \quad \|g(\mathbf{B}) - g_Q(\mathbf{B})\| = \mathcal{O}(e^{-CQ/\log(\beta E_M)}),$$

where  $C$  is a constant.

We remark that the construction proposed in this subsection also applies to the gapped case. In practice, if the temperature is high (so that  $\beta$  is small) or the gap around the chemical potential is small (in particular, for gapless system), the contour passing through the imaginary axis will be favorable; otherwise, the construction in the last subsection will be more efficient.

### 3. RATIONAL APPROXIMATIONS BASED ON THE MULTIPOLE EXPANSION

Another strategy for obtaining an efficient rational approximation for the Fermi-Dirac function for finite temperature is based on the multipole expansion, proposed recently in [12]. Let us first recall the construction of the multipole representation.

Using the Matsubara representation (pole expansion) of the Fermi-Dirac function, the density matrix is given by

$$(16) \quad \rho = 1 - 4\text{Re} \sum_{l=1}^{\infty} \frac{1}{\beta(\mathbf{H} - \mu) - (2l-1)\pi i}.$$

The summation in (16) can be seen as a summation of residues contributed from the poles  $\{(2l-1)\pi i\}$ , with  $l$  a positive integer, on the imaginary axis. This suggests looking for a multipole expansion of the contributions from the poles, as was done in the fast multipole method (FMM) [7]. To do so, we use a dyadic grouping of

the poles, in which the  $n$ -th group contains terms from  $l = 2^{n-1}$  to  $l = 2^n - 1$ , for a total of  $2^{n-1}$  terms. We decompose the summation in (16) accordingly. Let  $x = \beta(\mathbf{H} - \mu)$ . Then

$$(17) \quad \sum_{l=1}^{\infty} \frac{1}{x - (2l-1)\pi i} = \sum_{n=1}^{\infty} \sum_{l=2^{n-1}}^{2^n-1} \frac{1}{x - (2l-1)\pi i} = \sum_{n=1}^{\infty} S_n.$$

The basic idea is to combine the simple poles into a set of multipoles at  $l = l_n$ , where  $l_n$  is taken as the midpoint of the interval  $[2^{n-1}, 2^n - 1]$

$$(18) \quad l_n = \frac{3 \cdot 2^{n-1} - 1}{2}.$$

Then the  $S_n$  term in the above equation can be written as

$$(19) \quad \begin{aligned} S_n &= \sum_{l=2^{n-1}}^{2^n-1} \frac{1}{x - (2l-1)\pi i - 2(l-l_n)\pi i} \\ &= \sum_{l=2^{n-1}}^{2^n-1} \frac{1}{x - (2l_n-1)\pi i} \sum_{\nu=0}^{\infty} \left( \frac{2(l-l_n)\pi i}{x - (2l_n-1)\pi i} \right)^{\nu} \\ &= \sum_{l=2^{n-1}}^{2^n-1} \frac{1}{x - (2l_n-1)\pi i} \sum_{\nu=0}^{P-1} \left( \frac{2(l-l_n)\pi i}{x - (2l_n-1)\pi i} \right)^{\nu} \\ &\quad + \sum_{l=2^{n-1}}^{2^n-1} \frac{1}{x - (2l-1)\pi i} \left( \frac{2(l-l_n)\pi i}{x - (2l_n-1)\pi i} \right)^P. \end{aligned}$$

Using the fact that  $x$  is real, the second term in (19) can be bounded by

$$\sum_{l=2^{n-1}}^{2^n-1} \left| \frac{1}{x - (2l-1)\pi i} \right| \left| \frac{2(l-l_n)\pi i}{x - (2l_n-1)\pi i} \right|^P \leq \sum_{l=2^{n-1}}^{2^n-1} \frac{1}{|(2l-1)\pi|} \left| \frac{2(l-l_n)}{2l_n-1} \right|^P \leq \frac{1}{2\pi} \frac{1}{3^P}.$$

Therefore, if we approximate the sum  $S_n$  by the first  $P$  terms, the error decays exponentially fast with  $P$ :

$$(20) \quad \left| S_n(x) - \sum_{l=2^{n-1}}^{2^n-1} \frac{1}{x - (2l_n-1)\pi i} \sum_{\nu=0}^{P-1} \left( \frac{2(l-l_n)\pi i}{x - (2l_n-1)\pi i} \right)^{\nu} \right| \leq \frac{1}{2\pi} \frac{1}{3^P}.$$

The above analysis is of course standard from the view point of the fast multipole method [7]. The overall philosophy is also similar: given a preset error tolerance, one selects the value of  $P$ , the number of terms to retain in  $S_n$ , according to (20).

Moreover, the remainder of the sum in (16) from  $l = M_{\text{pole}} + 1$  to  $\infty$  has an explicit expression

$$(21) \quad \operatorname{Re} \sum_{l=M_{\text{pole}}+1}^{\infty} \frac{1}{2x - (2l-1)i\pi} = \frac{1}{2\pi} \operatorname{Im} \psi \left( M_{\text{pole}} + \frac{1}{2} + \frac{i}{\pi} x \right),$$

where  $\psi$  is the digamma function  $\psi(z) = \Gamma'(z)/\Gamma(z)$ .

In summary, we arrive at the following multipole representation for the Fermi operator [12]:

$$(22) \quad \rho = 1 - 4\text{Re} \sum_{n=1}^{N_G} \sum_{l=2^{n-1}}^{2^n-1} \frac{1}{\beta(\mathbf{H} - \mu) - (2l-1)\pi i} \sum_{\nu=0}^{P-1} \left( \frac{2(l-l_n)\pi i}{\beta(\mathbf{H} - \mu) - (2l-1)\pi i} \right)^\nu - \frac{2}{\pi} \text{Im} \psi \left( M_{\text{pole}} + \frac{1}{2} + \frac{i}{2\pi} \beta(\mathbf{H} - \mu) \right) + \mathcal{O}(N_G/3^P).$$

Here  $N_G$  is the number of groups in the multipole representation.  $M_{\text{pole}} = 2^{N_G} - 1$  is the number of poles that are effectively represented in the original Matsubara representation. In practice,  $N_G$  simple poles are first calculated, and then the  $N_G(P-1)$  multipoles can be constructed through matrix-matrix multiplication.

A disadvantage of (22) is that one needs to multiply simple poles together to get the multipoles before extracting the diagonal of Fermi operator. This prevents us from being able to directly apply the fast algorithms for extracting the diagonal of an inverse matrix, such as the one proposed in [13]. Therefore, it is useful to find an expansion similar to (22) that uses only simple poles. As we mentioned earlier, the key idea in deriving (22) is to combine the poles in each group together to form multipoles as the distance between them and the real axis is large. However, if instead we want an expansion that involves only simple poles, it is natural to revisit the variants of FMM that only use simple poles, for example, the version introduced in [15]. The basic idea there is to use a set of equivalent charges on a circle surrounding the poles in each group to reproduce the effective potential away from these poles.

Specifically, take the group of poles from  $l = 2^{n-1}$  to  $l = 2^n - 1$  for example. Consider a circle  $B_n$  with center  $c_n = (3 \cdot 2^{n-1} - 2)\pi i$  and radius  $r_n = 2^{n-1}\pi$ . It is clear that the circle  $B_n$  encloses the poles considered. Take  $P$  equally spaced points  $\{x_{n,k}\}_{k=1}^P$  on the circle  $B_n$ . Next, one needs to place equivalent charges  $\{\rho_{n,k}\}_{k=1}^P$  at these points such that the potential produced by these equivalent charges match with the potential produced by the poles inside  $B_n$  away from the circle. This can be done in several ways, for example, by matching the multipole expansion, by discretizing the potential on  $B_n$  generated by the poles, and so on. Here we follow the approach used in [15].

We simply take a bigger concentric circle  $\mathcal{B}_n$  outside  $B_n$  with radius  $R_n = 2^n\pi$  and match the potential generated on  $\mathcal{B}_n$  by the poles and by the equivalent charges on  $B_n$ . For this purpose, we solve for  $\rho_{n,k}$  the equations

$$(23) \quad \sum_{k=1}^P \frac{\rho_{n,k}}{y - x_{n,k}} = \sum_{l=2^{n-1}}^{2^n-1} \frac{1}{y - (2l-1)\pi i}, \quad y \in \mathcal{B}_n.$$

Regularization techniques such as Tikhonov regularization are required here since this is a first-kind Fredholm equation.

One can also prove that similar to the original version of the multipole representation, the error in the potential produced by the equivalent charges decay exponentially in  $P$ , the details can be found in [15]. Putting these altogether, we

can write down the following expansion of the Fermi-Dirac function

$$(24) \quad \rho = 1 - 4\text{Re} \sum_{n=1}^{N_G} \sum_{k=1}^P \frac{\rho_{n,k}}{\beta(\mathbf{H} - \mu) - x_{n,k}} - \frac{2}{\pi} \text{Im} \psi \left( M_{\text{pole}} + \frac{1}{2} + \frac{i}{2\pi} \beta(\mathbf{H} - \mu) \right) + \mathcal{O}(N_G/3^P).$$

The number of poles that are effectively represented in the original Matsubara representation is still  $M_{\text{pole}} = 2^{N_G} - 1$ .  $N_{\text{pole}} = N_G P$  simple poles are now to be calculated in practice.

The tail part can be approximated using a Chebyshev polynomial expansion. Similar to the analysis in [12], it can be shown that the complexity of the expansion is  $\mathcal{O}(\log \beta \Delta E)$ . As we pointed out earlier, the advantage of (24) over (22) is that only simple poles are involved in the formula. This is useful when combined with fast algorithms for extracting the diagonal of an inverse matrix [13].

Note that in (22) and (24), for  $2^{n-1} < P$  there would be no savings if we use  $P$  terms in the expansion. They are written in this form just for simplicity. In practice the first  $P$  simple poles will be calculated separately and the multipole expansion will be used starting from the  $(P+1)$ -th term and the starting level is  $n = \log_2 P + 1$ . We show in Fig. 5 a typical configuration of the set of poles in the multipole representation type algorithm.

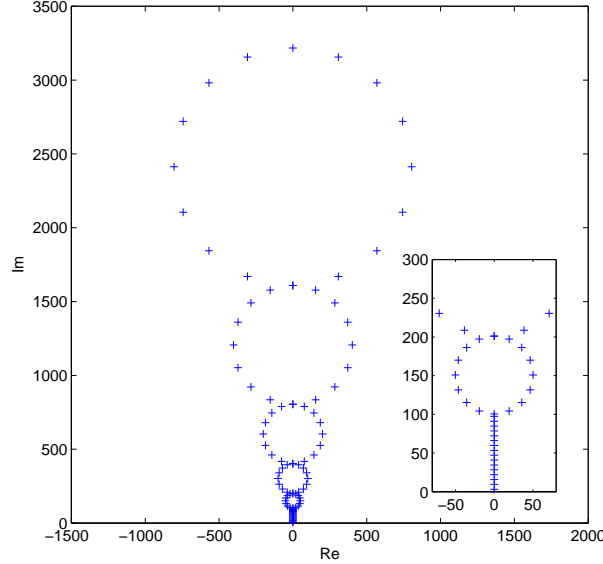


FIGURE 5. A typical configuration of the poles in the multipole representation type algorithm.  $M_{\text{pole}} = 512$  and  $P = 16$  is used in this figure. The poles with negative imaginary parts are not explicitly shown. The inset shows the first few poles. The first 16 poles are calculated separately and the starting level is  $n = 5$ .

$\beta\Delta E$	$N_{\text{pole}}$	$\Delta\rho_{\text{rel}}$
4, 208	40	$5.68 \times 10^{-7}$
8, 416	44	$3.86 \times 10^{-7}$
16, 832	44	$3.60 \times 10^{-7}$
33, 664	44	$3.55 \times 10^{-7}$
67, 328	44	$3.57 \times 10^{-7}$
134, 656	44	$3.47 \times 10^{-7}$
269, 312	44	$3.55 \times 10^{-7}$

TABLE 1.  $N_{\text{pole}}$  and  $L^1$  error of electronic density per electron with respect to various  $\beta\Delta E$ . The energy gap  $E_g \approx 0.01$ . The contour integral representation for gapped system at finite temperature is used for the calculation. The performance of the algorithm depends weakly on  $\beta\Delta E$ .

#### 4. NUMERICAL RESULTS

We test the algorithms described above using a two dimensional nearest neighbor tight binding model for the Hamiltonian. The matrix components of the Hamiltonian can be written as (in atomic units),

$$(25) \quad H_{i'j';ij} = \begin{cases} 2 + V_{ij}, & i' = i, j' = j, \\ -1/2 + V_{ij}, & i' = i \pm 1, j' = j \text{ or } i' = i, j' = j \pm 1. \end{cases}$$

The on-site potential energy  $V_{ij}$  is chosen to be a uniform random number between 0 and  $10^{-3}$ . The domain size is  $32 \times 32$  with periodic boundary condition. The chemical potential will be specified later. The accuracy is measured by the  $L^1$  error of the electronic density profile per electron

$$(26) \quad \Delta\rho_{\text{rel}} = \frac{\text{Tr} |\hat{P} - P|}{N_{\text{Electron}}}.$$

**4.1. Contour integral representation: gapped case.** The error of the contour integral representation is determined by  $N_{\text{pole}}$ . At finite temperature  $N_{\text{pole}} = 2Q$ , while at zero temperature  $N_{\text{pole}} = Q$ , with  $Q$  being the quadrature points on one loop of the contour. The performance of the algorithm is studied by the minimum number of  $N_{\text{pole}}$  such that  $\Delta\rho_{\text{rel}}$  (the  $L^1$  error in the electronic density per electron) is smaller than  $10^{-6}$ . For a given temperature, the chemical potential  $\mu$  is set to satisfy

$$(27) \quad \text{Tr} P = N_{\text{Electron}}.$$

In our setup the energy gap  $E_g \approx 0.01$  Hartree = 0.27 eV and  $E_M \approx 4$  Hartree. Therefore, this system can be regarded as a crude model for semiconductor with a small energy gap. The number of  $N_{\text{pole}}$  and the error  $\Delta\rho_{\text{rel}}$  are shown in Table 1 with respect to  $\beta\Delta E$  ranging between 4,000 and up to 270,000. Because of the existence of the finite energy gap, the performance is essentially independent of  $\beta\Delta E$ , as is clearly shown in Table 1.

When the temperature is low and therefore when  $\beta$  is large, as discussed before the finite temperature result is well approximated by the zero temperature Fermi operator, *i.e.*, the matrix sign function. In such case the quadrature formula is

given by (10). Only the contour that encircles the spectrum lower than chemical potential is calculated, and  $N_{\text{pole}} = Q$ .

In order to study the dependence of  $\Delta\rho_{\text{rel}}$  on the number of poles  $N_{\text{pole}}$ , we tune artificially the chemical potential to reduce the energy gap to  $10^{-6}$  Hartree. Fig. 6 shows the exponential decay of  $\Delta\rho_{\text{rel}}$  with respect to  $N_{\text{pole}}$ . For example, in order to reach the  $10^{-6}$  error criterion,  $N_{\text{pole}} \approx 50$  is sufficient. The increase in  $N_{\text{pole}}$  is very small compared to the large decrease of energy gap and this is consistent the logarithmic dependence of  $N_{\text{pole}}$  on  $E_g$  given by (8).

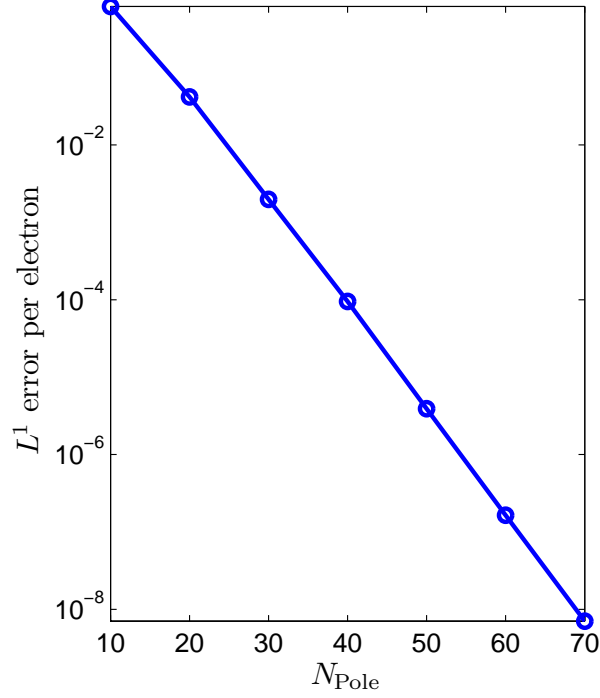


FIGURE 6. The lin-log plot of the  $L^1$  error of electronic density per electron with respect to  $N_{\text{pole}}$ . The energy gap  $E_g \approx 10^{-6}$ . The contour integral representation for gapped system at zero-temperature is used for calculation.

**4.2. Contour integral representation: gapless case.** For gapless systems such as metallic systems, our quadrature formula in (11) exploits the effective gap on the imaginary axis due to finite temperature. In the following results the chemical potential is set artificially so that  $E_g = 0$ .  $E_M \approx 4$  Hartree and the error criterion is still  $10^{-6}$  as in the gapped case. Table 2 reports the number of poles  $N_{\text{pole}}$  and the error  $\Delta\rho_{\text{rel}}$  with respect to  $\beta\Delta E$  ranging from 4,000 up to 4 million. These results are further summarized in Fig. 7 to show the logarithmic dependence of  $N_{\text{pole}}$  on  $\beta\Delta E$ , as predicted in the analysis of (15).

$\beta\Delta E$	$N_{\text{pole}}$	$\Delta\rho_{\text{rel}}$
4,208	58	$1.90 \times 10^{-7}$
8,416	62	$5.32 \times 10^{-7}$
16,832	66	$8.28 \times 10^{-7}$
33,664	72	$3.55 \times 10^{-7}$
67,328	76	$3.46 \times 10^{-7}$
134,656	80	$1.69 \times 10^{-7}$
269,312	84	$8.89 \times 10^{-8}$
538,624	88	$7.09 \times 10^{-8}$
1,077,248	88	$8.94 \times 10^{-7}$
2,154,496	88	$4.25 \times 10^{-7}$
4,308,992	92	$3.43 \times 10^{-7}$

TABLE 2.  $N_{\text{pole}}$  and  $L^1$  error of electronic density per electron with respect to various  $\beta\Delta E$ .  $E_g = 0$ . The contour integral representation for gapless system is used for the calculation.

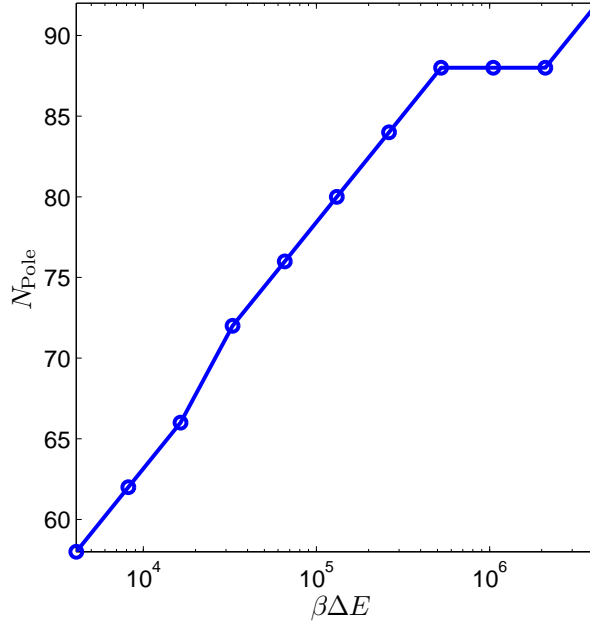


FIGURE 7. Log-lin plot of  $N_{\text{pole}}$  with respect to  $\beta\Delta E$ . The contour integral representation for gapless system is used for the calculation.

**4.3. Multipole representation.** The approach (24) based on the multipole representation has three parts of error: the finite-term multipole expansion, the finite-term Chebyshev expansion for the tail part, and the truncated matrix-matrix multiplication in the Chebyshev expansion.

$\beta\Delta E$	$M_{\text{pole}}$	$N_{\text{pole}}$	$N_{\text{Cheb}}$	$\Delta\rho_{\text{rel}}$
4,208	512	96	22	$4.61 \times 10^{-7}$
8,416	1,024	112	22	$4.76 \times 10^{-7}$
16,832	2,048	128	22	$4.84 \times 10^{-7}$
33,664	4,096	144	22	$4.88 \times 10^{-7}$
67,328	8,192	160	22	$4.90 \times 10^{-7}$
134,656	16,384	176	22	$4.90 \times 10^{-7}$
269,312	32,768	192	22	$6.98 \times 10^{-7}$
538,624	65,536	208	22	$3.20 \times 10^{-6}$
1,077,248	131,072	224	22	$7.60 \times 10^{-6}$

TABLE 3. The number of poles calculated  $N_{\text{pole}}$ , the order of Chebyshev expansion for the tail part  $N_{\text{Cheb}}$ , and the  $L^1$  error of electronic density per electron with respect to various  $\beta\Delta E$ . The number of poles excluded in the tail part  $M_{\text{pole}}$  is chosen to be proportional to  $\beta\Delta E$ .

The error from the multipole expansion is well controlled by  $P$  in (24). When  $P = 16$ ,  $1/3^P \sim \mathcal{O}(10^{-8})$ . The number of groups  $N_G$  is usually no more than 20, and therefore the error introduced by multipole expansion is around  $\mathcal{O}(10^{-7})$ , which is much less than the error criterion  $10^{-6}$ .

The number of terms in the Chebyshev expansion for the tail part  $N_{\text{Cheb}}$  is  $\mathcal{O}(\frac{\beta\Delta E}{M_{\text{pole}}})$ , with  $M_{\text{pole}}$  being the number of poles that are excluded in the tail part in the pole expansion. The truncation radius for the tail part is  $\mathcal{O}(\exp(-C\frac{\beta\Delta E}{M_{\text{pole}}}))$ . In order to reach a fixed target accuracy, we set  $M_{\text{pole}}$  to be proportional to  $\beta\Delta E$ . Due to the fact that  $M_{\text{pole}} \approx 2^{N_G} \approx 2^{N_{\text{pole}}/P}$ , this requires  $N_{\text{pole}}$  to grow logarithmically with respect to  $\beta\Delta E$ .

The target accuracy for the Chebyshev expansion is set to be  $10^{-7}$  and the truncation radius for the tail is set to be 4 for the metallic system under consideration. For  $\beta\Delta E = 4208$ ,  $M_{\text{pole}}$  is set to be 512 so that the error is smaller than  $10^{-6}$ . For other cases,  $M_{\text{pole}}$  scales linearly with  $\beta\Delta E$ . The lin-log plot in Fig. 8 shows the logarithmic dependence of  $N_{\text{pole}}$  with respect to  $\beta\Delta E$ . For more detailed results, Table 3 measures  $M_{\text{pole}}$ ,  $N_{\text{pole}}$ ,  $N_{\text{Cheb}}$ , and  $\Delta\rho_{\text{rel}}$  for  $\beta\Delta E$  ranging from 4000 up to 1 million. For all cases,  $N_{\text{Cheb}}$  is kept as a small constant. Note that the truncation radius is always set to be a small number 4, and this indicates the tail part is extremely localized in the multipole representation due to the effectively raised temperature.

Table 3 indicates that the error exhibits some slight growth. We believe that it comes from the growth of the number of groups in the multipole representation (24) and also the extra log log dependence on  $\beta\Delta E$  (see [12] for details). When compared with the results reported in Table 2, we see that for the current application to electronic structure, the contour integral representation outperforms the multipole representation in terms of both the accuracy and the number of poles used.

## 5. CONCLUSION

We propose two approaches for the expansion of Fermi operator: a rational approximation based on the contour integral idea introduced in [8] and a variant of

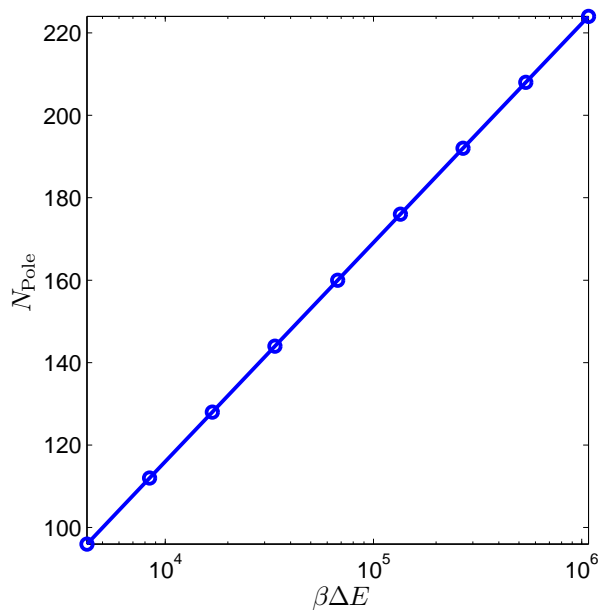


FIGURE 8. log-lin plot of  $N_{\text{pole}}$  with respect to  $\beta\Delta E$ . The multipole representation is used for the calculation.

the multipole representation in [12] using only simple poles. Both approximations result in logarithmic scaling complexity with respect to  $\beta\Delta\epsilon$  with small prefactor. Fast algorithms for electronic structure calculations can be obtained by combining these approaches with the algorithm introduced in [13] for extracting the diagonal of the inverse of a matrix.

**Acknowledgement:** This is a continuation of the work that was done jointly with Roberto Car, to whom we are very grateful for many stimulating discussions. This work was partially supported by DOE under Contract No. DE-FG02-03ER25587 and by ONR under Contract No. N00014-01-1-0674 (L. L., J. L. and W. E), and by an Alfred P. Sloan Research Fellowship and a startup grant from University of Texas at Austin (L. Y.).

#### REFERENCES

- [1] S. Baroni and P. Giannozzi, *Towards very large-scale electronic-structure calculations*, Europhys. Lett. **17** (1992), no. 6, 547–552.
- [2] M. Ceriotti, T.D. Kühne, and M. Parrinello, *A hybrid approach to Fermi operator expansion*, arXiv:0809.2232v1 (2008).
- [3] ———, *An efficient and accurate decomposition of the Fermi operator.*, J. Chem. Phys **129** (2008), no. 2, 024707.
- [4] S. Goedecker and L. Colombo, *Efficient linear scaling algorithm for tight-binding molecular dynamics*, Phys. Rev. Lett. **73** (1994Jul), no. 1, 122–125.
- [5] S. Goedecker and M. Teter, *Tight-binding electronic-structure calculations and tight-binding molecular dynamics with localized orbitals*, Phys. Rev. B **51** (1995Apr), no. 15, 9455–9464.
- [6] S. Goedecker, *Linear scaling electronic structure methods*, Rev. Mod. Phys. **71** (1999), no. 4, 1085–1123.

- [7] L. Greengard and V. Rokhlin, *A fast algorithm for particle simulations*, J. Comput. Phys. **73** (1987), no. 2, 325–348.
- [8] N. Hale, N. J. Higham, and L. N. Trefethen, *Computing  $A^\alpha$ ,  $\log(A)$ , and related matrix functions by contour integrals*, SIAM J. Numer. Anal. **46** (2008), no. 5, 2505–2523.
- [9] F.R. Krajewski and M. Parrinello, *Stochastic linear scaling for metals and nonmetals*, Phys. Rev. B **71** (2005), no. 23, 233105.
- [10] W. Liang, R. Baer, C. Saravanan, Y. Shao, A. T. Bell, and M. Head-Gordon, *Fast methods for resumming matrix polynomials and chebyshev matrix polynomials*, J. Comput. Phys. **194** (2004), no. 2, 575–587.
- [11] W. Liang, C. Saravanan, Y. Shao, R. Baer, A. T. Bell, and M. Head-Gordon, *Improved fermi operator expansion methods for fast electronic structure calculations*, J. Chem. Phys. **119** (2003), no. 8, 4117–4125.
- [12] L. Lin, J. Lu, R. Car, and W. E, *Multipole representation of the fermi operator with application to the electronic structure analysis of metallic systems*, Phys. Rev. B (2009), 115133.
- [13] L. Lin, J. Lu, L. Ying, R. Car, and W. E, *Fast algorithm for extracting the diagonal of the inverse matrix with application to the electronic structure analysis of metallic systems*, submitted to Comm. Math. Sci. (2009).
- [14] T. Ozaki, *Continued fraction representation of the Fermi-Dirac function for large-scale electronic structure calculations*, Phys. Rev. B **75** (2007), no. 3, 035123.
- [15] L. Ying, G. Biros, and D. Zorin, *A kernel-independent adaptive fast multipole algorithm in two and three dimensions*, J. Comput. Phys. **196** (2004), no. 2, 591–626.

PROGRAM IN APPLIED AND COMPUTATIONAL MATHEMATICS, PRINCETON UNIVERSITY, PRINCETON, NJ 08544. EMAIL: LINLIN@MATH.PRINCETON.EDU

PROGRAM IN APPLIED AND COMPUTATIONAL MATHEMATICS, PRINCETON UNIVERSITY, PRINCETON, NJ 08544. EMAIL: JIANFENG@MATH.PRINCETON.EDU

DEPARTMENT OF MATHEMATICS AND ICES, UNIVERSITY OF TEXAS AT AUSTIN, 1 UNIVERSITY STATION/C1200, AUSTIN, TX 78712. EMAIL: LEXING@MATH.UTEXAS.EDU

DEPARTMENT OF MATHEMATICS AND PACM, PRINCETON UNIVERSITY, PRINCETON, NJ 08544. EMAIL: WEINAN@MATH.PRINCETON.EDU