

# Robust Mixing for Ab-Initio Quantum Mechanical Calculations

L. D. Marks<sup>1</sup> and D. R. Luke<sup>2,\*</sup>

<sup>1</sup>*Department of Materials Science and Engineering, Northwestern University, Evanston, IL 60201.*

<sup>2</sup>*Department of Mathematical Sciences, University of Delaware, Newark DE 19716-2553, USA.*

(Dated: version 2.0—June 25, 2008)

We study the general problem of mixing for ab-initio quantum-mechanical problems. Guided by general mathematical principles and the underlying physics, we propose a multiseccant form of Broyden's second method for solving the self-consistent field equations of Kohn-Sham density functional theory. The algorithm is robust, requires relatively little fine-tuning and appears to outperform the current state of the art, converging for cases that defeat many other methods. We compare our technique to the conventional methods for problems ranging from simple to nearly pathological.

PACS numbers: PACS: 71.15.-m, 02.70.-c, 31.15.-p, 31.15.ec

## I. INTRODUCTION

We consider the problem of determining the electron density  $\rho$  that satisfies the self-consistent field equations according to the Kohn-Sham density functional theory<sup>1,2</sup>:

$$(H_0 + V_\rho)\phi_i = \epsilon_i\phi_i \quad (\text{I.1a})$$

$$\rho(x) = \sum_i \left(1 + e^{\beta(\epsilon_i - \mu)}\right)^{-1} |\phi_i(x)|^2. \quad (\text{I.1b})$$

Here  $H_0$  is the single-particle noninteracting Hamiltonian and  $V_\rho$  is an effective potential parameterized by the particle density  $\rho$ . The constant  $\beta$  is  $1/kT$  where  $k$  is Boltzmann's constant and  $T$  is temperature. The term  $(1 + e^{\beta(\epsilon_i - \mu)})^{-1}$  is the Fermi-Dirac occupation and the constant  $\mu$  is determined by  $\int \rho(x)dx = N$  for an  $N$ -body problem. Following<sup>3</sup> we let  $H_\rho := H_0 + \lambda V_\rho$ <sup>46</sup> denote the Kohn-Sham Hamiltonian and reformulate the above system of equations as a nonlinear fixed point problem: find  $\rho$  such that

$$F(\rho)(x) := \left(1 + e^{\beta(H_\rho - \mu)}\right)^{-1} (x, x) = \rho(x) \quad (\text{I.2})$$

where  $\mu$  is the unique solution to  $N = \text{trace}((1 + e^{\beta(H_\rho - \mu)})^{-1})$ . We refer to the operator  $F$  above as the self-consistent field (SCF) operator. We will not be concerned with the details of the SCF operator or its approximations since these tend to be specific to the application. Also, we will work with the discretized version of the SCF operator, which we will call the SCF *mapping* since it is a real vector-valued mapping of the discretized density. Throughout this work, however, we will point to instances where the form of this mapping can cause problems for numerical procedures for solving Eq.(I.2).

Numerical algorithms for solving Eq.(I.2) abound – the representative examples we focus on here are<sup>4–11</sup>. These are iterative procedures and the process of determining the desired density  $\rho$  from previous estimates has come to be known as “mixing” in the physical literature. For ab-initio methods there is frequently a user-provided mixing term which, if it is improperly chosen, will lead to divergence of the iterations. In many cases the user has to learn by failure what is the correct value to use, expending a fair amount of computer resources in the process. We will show that many of the methods found in the physical literature have counterparts in the mathematical literature where systematic approaches to the choice of algorithm parameters is well established. The goal of this work is the development of a method that does not require expert user input, is fast, and can handle many of the more complicated and poorly convergent problems such as metallic surfaces or heterostructures that can defeat a novice and sometimes an expert.

In the next two sections we discuss the leading methods in a novel analytical framework that clarifies similarities as well as fundamental differences. Our analysis sheds light on why the algorithms can fail which suggests strategies for design of an improved method. The class of algorithms we study are predicated upon mappings with a great deal of regularity – properties that the SCF mapping is not guaranteed to satisfy in all instances. Therefore, rather than viewing successive iterates as deterministic steps in a path to a solution, we treat the prior steps as random samples in a high-dimensional space. This viewpoint leads to a natural interpretation of the algorithm in terms of predicted

and unpredicted components, as well as the need for regularization and controls on the relative magnitude of the unpredicted step. In Section IV we present numerical results for both very easy problems as well as semi-pathological cases. The new approach outperforms existing algorithms in most cases, and does significantly better with poorly constructed Kohn-Sham mappings. The algorithm is also relatively insensitive to user input. We conclude with a discussion of some of the open issues.

## II. ITERATIVE METHODS FOR SOLVING NONLINEAR EQUATIONS

For fixed atom locations, we wish to determine the electron density  $\rho_*$ , a real-valued vector with  $k$  elements. With an estimated density  $\rho_n$  at the  $n$ -th step of an iterative procedure for determining  $\rho_*$ , we check whether our estimate satisfies the ab-initio self-consistent field (SCF) equations given by Eq.(I.2). Evaluation of the SCF mapping returns a modified density  $\rho'_n := F(\rho_n)$ , another real-valued vector with  $k$  elements. The density we seek is a fixed point of  $F$ , i.e., we solve the system of non-linear equations

$$F(\rho_*) - \rho_* = 0. \quad (\text{II.1})$$

This suggests the usual Newton algorithm as a possible numerical solution strategy.

### A. Newton-like Methods

Given a point  $\rho_n$ , Newton's method generates the next approximate solution to Eq.(II.1),  $\rho_{n+1}$ , by

$$\rho_{n+1} = \rho_n - (J(\rho_n) - I)^{-1}(F(\rho_n) - \rho_n). \quad (\text{II.2})$$

Under standard assumptions, this iteration can be shown to converge quadratically in a neighborhood of a local solution<sup>12</sup>. The computational cost of calculating the Jacobian is prohibitive for high-dimensional problems such as density functional calculations. Instead one can approximate the Jacobian via solutions to the matrix secant equation:  $B_n \approx (J(\rho_n) - I)$  where

$$B_n(\rho_n - \rho_{n-1}) = \left( (F(\rho_n) - \rho_n) - (F(\rho_{n-1}) - \rho_{n-1}) \right) \quad (\text{II.3})$$

Introducing new variables, this is represented as

$$B_n s_{n-1} = y_{n-1} \quad \text{or} \quad (\text{II.4})$$

$$H_n y_{n-1} = s_{n-1} \quad (\text{II.5})$$

where  $H_n = B_n^{-1}$  and

$$s_{n-1} = \rho_n - \rho_{n-1} \quad \text{and} \quad y_{n-1} = (F(\rho_n) - \rho_n) - (F(\rho_{n-1}) - \rho_{n-1}). \quad (\text{II.6})$$

The next density  $\rho_{n+1}$  is then generated either by the recursion

$$\rho_{n+1} = \rho_n - B_n^{-1}(F(\rho_n) - \rho_n) \quad (\text{II.7})$$

where  $B_n$  satisfies Eq.(II.4), or by

$$\rho_{n+1} = \rho_n - H_n(F(\rho_n) - \rho_n) \quad (\text{II.8})$$

where  $H_n$  satisfies Eq.(II.5). The variables in Eq.(II.4) and Eq.(II.5) are the matrices  $B_n$  and  $H_n$  respectively, and there are infinitely many possible solutions, each leading to a different numerical technique. Our focus in this study is on improvements of the Broyden family discussed next.

### B. Rank One Updates

A new matrix  $B_{n+1}$  is obtained by updating in some fashion  $B_n$  using the new data pair  $(s_n, y_n)$  combined with the prior information  $(s_0, y_0), (s_1, y_1), \dots, (s_{n-1}, y_{n-1})$  subject to the constraint that  $B_{n+1}$  satisfy Eq.(II.4). Broyden<sup>13</sup> looked at two approaches, given here in a multistep recursion.

The first (B1) is based on updates to the approximate Jacobian in Eq.(II.4), and is shown in<sup>14</sup> (Theorem 6.2) to be

$$B_{n+1}^{-1} = B_0^{-1} - (B_0^{-1}Y_n - \beta_n S_n) (L_n + S_n^T B_0^{-1}Y_n)^{-1} S_n^T B_0^{-1} \quad (\text{II.9})$$

where  $\beta_n$  is a scaling,

$$S_n := [s_0, s_1, s_2, \dots, s_n], \quad Y_n := [y_0, y_1, y_2, \dots, y_n] \quad (k\text{-by-}(n+1) \text{ matrices}) \quad (\text{II.10})$$

and  $(L_n)_{i,j} := \{-s_{i-1}^T s_{j-1} \text{ if } i > j; 0 \text{ otherwise}\}$ .

The second of Broyden's methods (B2) is based on updates to the approximate *inverse* Jacobian in Eq.(II.5) and is given by

$$H_{n+1} = H_0 \prod_{j=0}^n W_j + \sum_{j=0}^n \left( Z_j \prod_{i=j+1}^n W_i \right) \quad (\text{II.11})$$

where the products ascend from left to right with the empty product defined as 1, and

$$W_n := I - \frac{y_n y_n^T}{\|y_n\|^2} \quad \text{and} \quad Z_n := \beta_n \frac{s_n y_n^T}{\|y_n\|^2} \quad (n = 1, 2, \dots).$$

Here and throughout this work the norm  $\|y\| = \sqrt{y^T y}$  is the Euclidean norm and a vector (understood to be a *column* vector) or matrix raised to the power  $T$  indicates the transpose. Note that our sign convention is different to the sign in Broyden's paper where he takes  $H_{n+1} = -B_{n+1}^{-1}$ . Our recursion appears to be new, and for  $\beta_n = 1$  can be shown to be equivalent to a recursion proposed by Srivastava<sup>9</sup> with the same storage requirements.

Both Eq.(II.9) and Eq.(II.11) can be performed without storing or forming the matrix explicitly. In the recursions Eq.(II.9) and Eq.(II.11) the initial matrix,  $B_0$  and  $H_0$  respectively, is crucial; we explore scalings in greater detail in Subsection III C. Srivastava's formulation was initially implemented for LAPW code by<sup>11</sup> with  $H_0$  fixed, but a dynamic  $H_0$  yields substantially better performance.

Both of Broyden's methods are shown in<sup>15</sup> to converge locally superlinearly under the standard assumptions that the Jacobian exists, is nonsingular and Lipschitz continuous at the solution. Update Eq.(II.11), however, was not recommended by Broyden and subsequently became known as Broyden's "bad" method.

Broyden's updates are the nearest matrices to the previous matrix with respect to the Frobenius norm<sup>47</sup> that satisfy the current matrix secant equation Eq.(II.4) or Eq.(II.5). The main difference between B1 and B2 is the space in which the "nearest" criterion is applied<sup>16</sup>. For B1 the criterion is applied in the domain of the mapping, while B2 is applied in the range, where the domain of the mapping is the space of the density differences  $s_n$  and the range is the space of the residual differences  $y_n$ . We see from Eq.(II.7) that an ill-conditioned matrix update  $B_n$  will lead to a large and possibly unstable estimation of the step  $s_n$ . On the other hand, from Eq.(II.8) it is clear that a least change criterion in the space of the residual differences  $y_n$  will lead to smaller steps that could slow progress for well-conditioned problems; we return to this issue below.

### C. Multisecant Methods

To generate the  $n + 1$ -th Jacobian approximation the methods described above satisfy the matrix secant equation Eq.(II.4) or Eq.(II.5) for the current step  $s_n$  and residual difference  $y_n$ . Updating the Jacobian based only on the most recent sample and ignoring the other sample points imposes a bias toward the most recent step. Searching for the nearest matrix that satisfies the matrix secant equation only for the most recent sample point is a greedy strategy without recourse.

Multisecant techniques put the previous data on more equal footing with the most recent steps; that is, rather than satisfying the matrix secant equation for only the most recent step one satisfies *all* matrix secant equations *simultaneously*:

$$Y_n = B_n S_n \quad \text{or} \quad S_n = H_n Y_n \quad (\text{II.12})$$

where  $S_n = [s_{1,n}, s_{2,n}, \dots, s_{m,n}]$  and  $Y_n = [y_{1,n}, y_{2,n}, \dots, y_{m,n}]$  are  $k$ -by- $m$  ( $m \leq \min\{n, k\}$ ) matrices whose columns are previous steps and residual differences respectively. Multisecant techniques have been thoroughly studied in the mathematical literature<sup>14,17-25</sup>. Methods appearing independently in the physical sciences literature<sup>5-8,10</sup> are relaxations of more conventional multisecant methods. A very recent study of multisecant methods brought to our attention by an anonymous referee is<sup>26</sup>.

Many multiseccant methods are easily understood by formulating the underlying optimization problem each of the approximate Jacobians (implicitly) solves. We consider first the constrained optimization problem

$$\underset{X \in C}{\text{minimize}} \quad \frac{1}{2} \|A - X\|^2 \quad (\text{II.13})$$

where, throughout, the norm of a matrix is the *Frobenius* norm,  $A$  is a real  $k \times k$  matrix, and the set  $C := \{X \in \mathbb{R}^{k \times k} \text{ such that } XD = G\}$ . for  $D, G$  real  $k \times m$  matrices such that  $C$  is nonempty. If the columns of  $D$  are linearly independent, the solution  $X_*$  to the optimization problem Eq.(II.13) is the *orthogonal projection* of  $A$  onto  $C$ , written explicitly as

$$X_* = A + (G - AD) (D^T D)^{-1} D^T. \quad (\text{II.14})$$

Specializing to multiseccants, if  $A$  is an approximation to the Jacobian,  $D = S_n \in \mathbb{R}^{k \times m}$  and  $G = Y_n \in \mathbb{R}^{k \times m}$  ( $1 \leq m \leq n$ ), the columns of which are denoted  $y_j$  and  $s_j$  respectively ( $j \in [0, n]$ ), then we arrive at the multiseccant extension of the Broyden's first update (MSB1) as studied by<sup>18,19,21,23</sup>:

$$B_{n+1} = A + (Y_n - AS_n) (S_n^T S_n)^{-1} S_n^T. \quad (\text{II.15})$$

Elementary calculations using the Sherman-Morrison-Woodbury formula yield the multi-step recursion for  $B_{n+1}^{-1}$ , analogous to Eq.(II.9),

$$B_{n+1}^{-1} = A^{-1} + (S_n - A^{-1}Y_n) \left( (S_n^T S_n)^{-1} S_n^T A^{-1} Y_n \right)^{-1} (S_n^T S_n)^{-1} S_n^T A^{-1} \quad (\text{II.16})$$

Sequences based on update Eq.(II.15) are shown in<sup>23</sup> (Theorem 2.5) to be locally q-superlinearly convergent if, in addition to other standard assumptions, the approximate Jacobians,  $B_n$  stay close to the behavior of the true Jacobian, and if the columns of  $S_n$  are strongly linearly independent. Moreover, storage requirements for this formulation are no greater than those of Srivastava's implementation of Broyden's second method.

An alternative specialization of Eq.(II.14) leads to a multiseccant form of Broyden's second method (MSB2) if we let  $A$  be an approximation to the inverse of the Jacobian,  $D = Y_n$  and  $G = S_n$  so that

$$H_{n+1} = A + (S_n - AY_n) (Y_n^T Y_n)^{-1} Y_n^T. \quad (\text{II.17})$$

To our knowledge, there are no published numerical comparisons of Eq.(II.17) to alternatives, neither is there any published convergence theory, though we believe this is only a minor modification of the theory for Eq.(II.15). Again, the storage requirements for this recursion are equivalent to MSB1 and B2.

Independent studies appearing in the physics literature that parallel the mathematical literature have a different variational form. The various approaches can all be shown to be specializations of the optimization problem

$$\underset{X \in \mathbb{R}^{k \times k}}{\text{minimize}} \quad \frac{1}{2} \sum_{j=1}^n \alpha_j \text{dist}_{C_j}^2(X) + \frac{\alpha_0}{2} \|A - X\|^2 \quad (\text{II.18})$$

where each  $C_j := \{X \in \mathbb{R}^{k \times k} \text{ such that } XD_j = G_j\}$ ,  $A \in \mathbb{R}^{k \times k}$ , and  $\text{dist}_{C_j}(X)$  is the Euclidean distance of  $X$  to the set  $C_j$ . A short calculation yields the solution  $X_*$  to Eq.(II.18)

$$X_* = \sum_{j \in \mathbb{J}} \gamma_j A + \sum_{j=1}^n \gamma_j \left( (G_j - AD_j) (D_j^T D_j)^{-1} D_j^T \right) \quad \text{where} \quad \gamma_j := \frac{\alpha_j}{\sum_{j=0}^n \alpha_j}. \quad (\text{II.19})$$

Specializing to multiseccants, let  $A = B_n$ ,  $D_j = s_j$  and  $G_j = y_j$ , where  $s_j$  and  $y_j$  are defined by Eq.(II.6). Then the optimization problem Eq.(II.18) corresponds to the variational formulation of a method proposed by Vanderbilt and Louie<sup>10</sup>. A local convergence analysis, together with numerical tests are studied in<sup>5</sup>. Our derivation and formulation of the update, however, appears to be new and clarifies the connections between their method and Eq.(II.15) above:

$$B_{n+1} = \sum_{j=0}^n \gamma_j B_n + \sum_{j=1}^n \gamma_j \left( (y_j - B_n s_j) (s_j^T s_j)^{-1} s_j^T \right). \quad (\text{II.20})$$

If instead we let let  $A = H_n$ ,  $D_j = y_j$  and  $G_j = s_j$ , we get the update proposed by Johnson<sup>6</sup>:

$$H_{n+1} = \sum_{j=0}^n \gamma_j H_n + \sum_{j=1}^n \gamma_j \left( (s_j - H_n y_j) (y_j^T y_j)^{-1} y_j^T \right). \quad (\text{II.21})$$

Again, our derivation is different, and the new formulation makes the connection with Eq.(II.17) more transparent.

The weighting scheme of<sup>6,10</sup> is similar to a technique proposed by Pulay<sup>8</sup>. A dynamic weighting scheme that optimizes the weights  $\gamma_j$  simultaneously with the determination of the matrix  $H_n$  or  $B_n$  is possible via the extended least squares techniques outlined in<sup>27</sup>. A variation of Eq.(II.21) due to Kawata et al<sup>7</sup> combines the method of Johnson with a construction of the columns of  $S_n$  and  $Y_n$  proposed by Pulay<sup>8</sup> and given in Eq.(III.3). We note that the methods summarized by Eq.(II.20)-(II.21) and their relatives solve single matrix secant equations *in parallel* and then average these solutions, while the methods summarized by Eq.(II.15)-(II.17) seek the single matrix that solves all the matrix secant equations *simultaneously*, which is more restrictive.

In the above analysis we are not specific about *how many* previous steps should be included in the matrices  $S_n$  and  $Y_n$ . Recall that these matrices are made up of  $m$  columns of previous step information where  $m \in [1, n]$ . If  $m < n$  then one is implicitly executing a *limited memory* technique<sup>14</sup>. If one constructs  $S_n$  and  $Y_n$  via Eq.(III.3), as we do in the following numerical experiments, then one would exclude points that are most distant from the current point  $\rho_n$ . This is a reasonable strategy for highly nonlinear problems, where the linear approximation that is at the heart of quasi-Newton methods is only valid on a local neighborhood of the current point. For extremely large problems such a strategy is also expedient since the matrix updates need not be explicitly stored as they can be constructed from a few stored vectors.

### III. SAFEGUARDED MULTI-SECANTS

Newton-like algorithms are not global techniques for solving equations and can behave wildly, even chaotically, far from a solution. For the practitioner who simply wants her software to converge for a particular example, unfortunately this means that the algorithms come only with extremely limited warranties that may not even be verifiable. The extent to which algorithms behave, or misbehave, depends on the functional properties of the SCF mapping. Consider the following simple algorithm

$$\rho_{n+1} = F(\rho_n). \quad (\text{III.1})$$

If  $F$  is a *contraction* on some closed subset of the space of densities (i.e. points move *closer* to one another under the mapping  $F$ ), then the sequence  $\rho_n$  converges to the *unique* fixed point  $\rho_*$  of  $F$  (Banach Contraction Theorem see, for instance<sup>28</sup>). If  $F$  is not a contraction, then Eq.(III.1) could continue forever without ever approaching a fixed point. Successive iterates might form a characteristic path, or they might behave chaotically. Less restrictive than contractive mappings are *nonexpansive* mappings (i.e. points do not move further apart under the mapping  $F$ ). If  $F$  is nonexpansive on a closed convex symmetric subset of Euclidean space  $\mathbb{X}$  and has fixed points (as it would if  $\mathbb{X}$  were *bounded*<sup>29,30</sup>), then for any  $\rho_0 \in \mathbb{X}$  the sequence of steps defined by the iteration

$$(n = 0, 1, 2, \dots) \quad \rho_{n+1} = \tilde{F}_\lambda(\rho_n) := \rho_n + \lambda(F(\rho_n) - \rho_n) \quad (\text{III.2})$$

converges to a fixed point of  $F$  (see Theorem 2.1 and Corollary 2.3 of<sup>31</sup>).

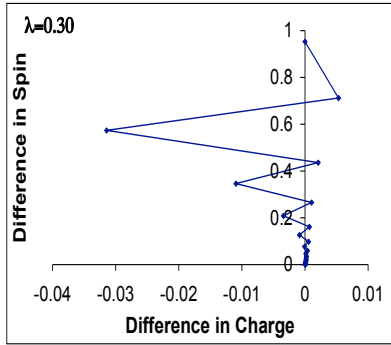
Most readers will recognize iteration Eq.(III.2) as the Pratt step<sup>32</sup>. When  $F(\rho_n) - \rho_n$  is an approximation to the gradient of the Kohn-Sham energy functional, then the Pratt step is simply an approximation to steepest descent with (fixed) step length  $\lambda$ . Though convergence is guaranteed for nonexpansive SCF mappings on compact convex regions it can be extremely slow. If  $F$  is not nonexpansive, then the numerical behavior of fixed point iterations like Eq.(III.2) and even Broyden's methods cannot be guaranteed.

Note that one can extend the above concepts to a subset of the density variables. For instance, the *sp*-electron states might converge quickly, while *d*-electron states might be very difficult to converge. Indeed, a frequent observation is that the density within the muffin-tins often behaves very differently to the density in the interstitial region.

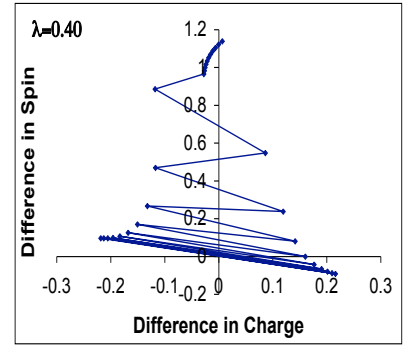
The problematic part of the Kohn-Sham mapping is the effective potential  $V_\rho$ . In general, there is no closed form for  $V_\rho$ . For certain approximations, denoted  $\tilde{V}$ , it is possible to prove the correspondence between the fixed points of the corresponding SCF mapping  $F_{\tilde{V}}$  and solutions to the Kohn-Sham equations<sup>3</sup>, and, moreover, that  $F_{\tilde{V}}$  is a contraction<sup>33</sup>. However, for exact  $V_\rho$  at finite temperatures existence and uniqueness of fixed points is an open question, further complicated by the occurrence of systems with multiple coexisting phases<sup>3</sup>.

With this in mind, and before we present the details of our algorithm, we describe in physical terms some of the features of ab-initio calculations that are problematic, together with common symptoms of poorly convergent problems.

- i. In many cases, for instance bulk MgO, the algorithms reach an acceptable solution in a surprisingly small number of iterations, e.g. 10 – 20 for  $10^4$  unknown density components. This implies that, at least for a substantial subset of the density parameters, the domain of attraction of the fixed point is large and the SCF mapping has “good” functional properties on this domain.



(a)



(b)

FIG. 1: (Color online) Iterates for an  $O_2$  molecule with atomic densities having a spin of +2 the other with 0. The figures show the difference between the spins (vertical axis) for the two atoms and the difference between the total charges (horizontal axis) within the muffin tins for iterates generated by the Pratt step (Eq.(III.2)) with different fixed step-length parameters  $\lambda$ . In frame (a) the Pratt step parameter is  $\lambda = 0.30$ , in frame (b) the Pratt step parameter is  $\lambda = 0.40$ .

- ii. In some cases there can be issues with the scaling of different parts of the density because they are represented in quite different fashions. For instance, with LAPW methods the plane wave components outside the muffin-tins are represented by the Fourier coefficients whereas the density inside the muffin-tins is expanded in terms of spherical harmonics.
- iii. The conventional wisdom for LAPW-based methods is that the muffin-tins should be as large as possible without overlapping. This implies that the basis set used for the muffin-tins is better suited for the physics or for the geometry of the atoms. This is manifested in more rapid convergence of the coefficients corresponding to these basis elements and indicates that the domain of attraction of the fixed point for these coefficients is large relative to the domain of attraction for the fixed point of the plane wave elements.
- iv. The most physically interesting problems are often harder to solve. A spin unpolarized DFT calculation of NiO, for example, may converge very slowly. The slow convergence of the mixing cycle is in part because spin unpolarized the system is metallic, but is also coincidental with an imperfect functional description of this system, in which case the Hamiltonian in Eq.(I.2) can be ill-posed. It is not uncommon to compromise on the physical model, particularly for large and complicated problems.
- v. In some cases, for instance when there are  $d$  or  $f$  electrons, charge carriers are in a large unit cell and for surfaces, mixing converges poorly and can easily diverge. In the literature this is called “charge sloshing” because one has oscillations of charge density between different spatial regions of the problem or between different local states such as  $d$ -electrons. Mathematically this sometimes corresponds to ill-conditioning when a small change in the density  $\rho$  can lead to large change in  $F(\rho)$ , with large eigenvalues of the matrix  $H$  (or small eigenvalues of  $B$ ). Alternatively, it may be that the higher-order terms in the Taylor series expansion of the Jacobian are large, so neglecting them is only appropriate for a very small change in the density. A third possibility in the case of charge sloshing is that the SCF mapping is not nonexpansive (and hence not contractive) along this trajectory. None of these possibilities is mutually exclusive.

To illustrate these features a simple model is an  $O_2$  molecule starting from atomic densities where the two atoms are deliberately treated differently, one starting with a spin of +2 the other with 0. Shown in Figure 1 is the difference between the spins (vertical axis) for the two atoms and the difference between the total charges (horizontal axis) within the muffin tins for iterates generated by the Pratt step (Eq.(III.2)) with different fixed step-length parameters  $\lambda$ . While the spins converge relatively smoothly to the final solution, the total charge oscillates or “sloshes”. The charge oscillations become unstable for a relatively small change in the Pratt step parameter.

## A. Mathematical Framework

There are two elements that distinguish our approach from previous work on matrix secant methods: first is the view of the matrix secant update as a *Jacobian simplex* of a vector-valued mapping, and second is the separation of the matrix update into *predicted* and *unpredicted* components. Both of these viewpoints are rooted in the observation that the dimension of the underlying problem is on the order of  $10^4$  or higher while the information used to model the fixed point mapping is at most dimension  $2n$  where  $n$  is the number of iterations (on the order of  $10^0$ ). The conventional view is that the  $n$  steps and residual differences generated in matrix secant methods are deterministic points on a path to the solution. Alternatively, we consider the  $n$  steps as random samples of a high-dimensional mapping.

The origins of many matrix secant methods are closely related to the conjugate gradient algorithm. According to this interpretation the construction of the matrices  $S_n$  and  $Y_n$  given in Eq.(II.6) is consistent with the columns of  $S_n$  being *conjugate directions*. Viewing the steps instead as *samples* on a small neighborhood of the current iterate leads us to the alternative centering

$$s_{j,n} = \rho_j - \rho_n \quad \text{and} \quad y_{j,n} = (F(\rho_j) - \rho_j) - (F(\rho_n) - \rho_n) \quad (j = 0, 1, \dots, n-1). \quad (\text{III.3})$$

The matrix secant update built from these step and residual differences is essentially a finite difference approximation to the Jacobian centered at  $\rho_n$  and, in the context of scalar-valued functions, yields what is known as the *gradient simplex*. The generalization in the present context is then appropriately called a Jacobian simplex. We therefore consider the vectors  $s_j$  and  $y_j$  given by Eq.(III.3) merely as data samples with no significance given to the order in which the samples were collected. This is a fundamentally different approach than the conventional matrix secant updates based on Eq.(II.6).

Independent of how one centers the step history is how one treats the components of the new step generated by the matrix secant update. Given the data samples, the algorithm *predicts* the behavior of the SCF mapping Eq.(I.2) at  $\rho_n$ . The multi-secant methods detailed in the previous sections can all be rewritten as

$$\rho_{n+1} = \rho_n - A_0 \left( I - Y_{n-1} A_n \right) g_n - S_{n-1} A_n g_n. \quad (\text{III.4})$$

where  $g_n = F(\rho_n) - \rho_n$ ,  $A_n$  is a matrix dependent on the method, and  $A_0$  is an *inverse* Jacobian estimate. Let us write this as

$$\rho_{n+1} - \rho_n = u_n + p_n \quad (\text{III.5})$$

where, according to Eq.(III.4),  $p_n = -S_{n-1} A_n g_n$  and  $u_n = -A_0 (I - Y_{n-1} A_n) g_n$ . We interpret  $p_n$  as the part of the vector  $g_n$  that can be explained by (is in the range of) the data at step  $n$ , and  $u_n$  is the component that is orthogonal to this information, and hence unpredicted. Intuitively, taking too large a step along the unpredicted component may be a bad idea and one of the sources of instabilities in these methods. We propose a strategy for controlling this step component in section III C. One of the main differences between (MS)B1 and (MS)B2 is the size of the unpredicted component. We show below that this component is inherently larger for (MS)B1 than for (MS)B2. Therefore if instabilities are due to large steps in the unpredicted direction it follows that B2 and MSB2 may perform better, which we will see later is the case for DFT problems.

This is made rigorous when we consider the multiseant formulation of Broyden's second method. Rewriting Eq.(II.8) with  $H_n$  replaced by Eq.(II.17) and rearranging components according to Eq.(III.4) yields

$$A_n := (Y_{n-1}^T Y_{n-1})^{-1} Y_{n-1}^T. \quad (\text{III.6})$$

Note that  $(Y_{n-1}^T Y_{n-1})^{-1} Y_{n-1}^T g_n$  is the solution to the least squares minimization problem

$$\underset{z \in \mathbb{R}^m}{\text{minimize}} \quad \frac{1}{2} \| Y_{n-1} z - g_n \|^2, \quad (\text{III.7})$$

where  $m \in [1, n-1]$  is the number of previous data points used in the update. Here,  $(Y_{n-1}^T Y_{n-1})^{-1} Y_{n-1}^T g_n$  is the element in the domain of  $Y_{n-1}$  that comes closest (in the least squares sense) to "predicting" the vector  $g_n$ . It follows, then, that

$$\left( I - Y_{n-1} A_n \right) g_n = \left( I - Y_{n-1} (Y_{n-1}^T Y_{n-1})^{-1} Y_{n-1}^T \right) g_n \quad (\text{III.8})$$

is the orthogonal projection of  $g_n$  onto the space orthogonal to the residual differences  $y_j$  defined by one of Eq.(II.6) or Eq.(III.3), our prior data.

The formalization for Broyden's first method is not as immediate. From Eq.(II.7) with  $B_n^{-1}$  replaced by Eq.(II.16), the modification of Eq.(III.4) for MSB1 amounts to

$$A_n := \left( (S_{n-1}^T S_{n-1})^{-1} S_{n-1}^T A_0 Y_{n-1} \right)^{-1} (S_{n-1}^T S_{n-1})^{-1} S_{n-1}^T A_0, \quad (\text{III.9})$$

where, again,  $A_0$  is an estimate of the *inverse* Jacobian. Again, we note that  $(S_{n-1}^T S_{n-1})^{-1} S_{n-1}^T A_0 w$  is the solution to the least squares problem  $\underset{z \in \mathbb{R}^{n-1}}{\text{minimize}} \frac{1}{2} \|S_{n-1} z - A_0 w\|^2$ . If, in addition,  $A_0 = \sigma I$ , then an elementary calculation yields the simplification to Eq.(III.9)

$$A_n = (S_{n-1}^T Y_{n-1})^{-1} S_{n-1}^T \quad (\text{III.10})$$

If  $(S_{n-1}^T Y_{n-1})^{-1}$  is well-defined, then the mapping  $I - Y_{n-1} A_n$  is a *nonorthogonal projection*<sup>48</sup> onto the nullspace of the columns of  $S_{n-1}$ , or in other words, a projection onto the space orthogonal to the range of the columns of  $S_{n-1}$ , our prior step data. Unlike Eq.(III.6) the projection is *not* to a nearest element in the range of  $S_{n-1}^\perp$ , hence, by definition, the resulting step will be larger than the orthogonal projection.

### B. Regularization, and preconditioning: the matrix $A_n$

The discussion in the previous subsection of MSB2 assumes that  $Y_{n-1}$  is full-rank. If the columns of  $Y_{n-1}$  are nearly linearly dependent, then the inverse  $(Y_{n-1}^T Y_{n-1})^{-1}$  can be numerically unstable. More fundamentally, we are implicitly assuming that the approximation to the Jacobian in Eq.(II.2) is, first of all, valid on the neighborhood of  $\rho_n$  defined by the other data points and, second of all, that the Newton step is the *right* step to take. If either one of these assumptions does not hold, as would be the case when we are far from the solution and our sample points are far apart, conventional optimization strategies link local and global techniques by allowing steps to rotate between the steepest descent direction (in the present setting, the direction of the vector  $g_n$ ) and a Newton-like direction. One well-known strategy of this kind is the Levenberg-Marquardt algorithm<sup>34,35</sup>. We propose a different technique that is an unusual use of a classical *regularization* technique usually attributed to Tikhonov<sup>36-38</sup> and rediscovered in the statistics community under the name of ridge regression<sup>39</sup>, though the more general notion of proximal mappings due to Moreau<sup>40</sup> predates both of these. In particular we regularize Eq.(III.7) in the usual way:

$$\underset{z \in \mathbb{R}^m}{\text{minimize}} \frac{1}{2} \|Y_{n-1} z - g_n\|^2 + \frac{\alpha}{2} \|z\|^2, \quad (\alpha > 0). \quad (\text{III.11})$$

The solution to Eq.(III.11) is

$$z_n = (Y_{n-1}^T Y_{n-1} + \alpha I)^{-1} Y_{n-1}^T g_n \quad (\text{III.12})$$

which yields the following regularization of  $A_n$  given by Eq.(III.6):

$$A_n^\alpha := (Y_{n-1}^T Y_{n-1} + \alpha I)^{-1} Y_{n-1}^T. \quad (\text{III.13})$$

Note that as  $\alpha \rightarrow \infty$ ,  $A_n^\alpha \rightarrow 0$ , and the step generated by Eq.(III.4) rotates to the direction  $A_0 g_n$ . We thus interpret the regularization parameter in both the conventional way, stabilizing  $(Y_{n-1}^T Y_{n-1})^{-1}$ , and as an estimation of the uncertainty of the approximate Newton step. Given our understanding of the previous step data as pseudo-random samples from an unknown process, the latter interpretation has a very natural explanation in terms of the Wiener filter for a signal with normally distributed zero-mean white noise. The size of the regularization parameter corresponds to the energy of the noise, or uncertainty in our model.

Similarly, the discussion of MSB1 in the previous subsection also assumes that  $(S_{n-1}^T Y_{n-1})^{-1}$  is well-defined, but this says nothing of whether or not  $S_{n-1}^T Y_{n-1}$  is well-conditioned. Regularization of  $(S_{n-1}^T Y_{n-1})^{-1}$  in Eq.(III.10) gives

$$A_n^\alpha := (S_{n-1}^T Y_{n-1} + \alpha I)^{-1} S_{n-1}^T \quad (\text{III.14})$$

which shifts the eigenvalues to the right. Since  $S_{n-1}^T Y_{n-1}$  could have negative eigenvalues, unless  $\alpha$  is chosen large enough, this regularization could result in an even *more* ill-conditioned matrix. Our numerical experience is that  $\alpha > 10^{-6}$  is sufficiently large to avoid this possibility for the applications of interest to us.



Johnson<sup>6</sup> proposes a normalization of the columns of the matrices of  $Y_n$  and  $S_n$  for numerical reasons, though this can easily be shown to have no formal impact on the algorithm. Such a normalization can, however, have a significant effect on the choice of the regularization parameter. This is equivalent to multiplication of the matrices  $Y_n$  and  $S_n$  on the right by the diagonal matrix  $\Psi_n$ . We show the formalism for MSB2 – MSB1 is handled similarly. The least squares problem analogous to Eq.(III.11) under such a renormalization is

$$\underset{z \in \mathbb{R}^m}{\text{minimize}} \quad \frac{1}{2} \|Y_{n-1} \Psi_n z - g_n\|^2 + \frac{\alpha}{2} \|z\|^2, \quad (\alpha > 0) \quad (\text{III.15})$$

with the solution  $(\Psi_n Y_{n-1}^T Y_{n-1} \Psi_n + \alpha I)^{-1} \Psi_n Y_{n-1}^T g_n$ . It follows immediately from this that if we normalize the columns of  $Y_{n-1}$ ,

$$\Psi_n = \begin{pmatrix} 1/\|y_1^{(n-1)}\| & 0 & \dots & 0 \\ 0 & 1/\|y_2^{(n-1)}\| & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & 1/\|y_m^{(n-1)}\| \end{pmatrix} \quad (\text{III.16})$$

where  $y_j^{(n-1)}$  is the  $j$ -th column of  $Y_{n-1}$ , then our regularization parameter will be *independent* of multiple scales between the columns of the matrix  $Y_{n-1}$ . Viewing the regularization as a Wiener filter applied to the approximate Newton step, the normalization reduces the effect of outliers on the regularization parameter in the least squares estimation, these outliers coming from steps that are relatively far away from the current point. We denote the matrix corresponding to this renormalization, together with the regularization  $\alpha$  by  $A_n^{\alpha, \Psi}$  where

$$A_n^{\alpha, \Psi} := \begin{cases} \Psi_n (\Psi_n S_{n-1}^T Y_{n-1} \Psi_n + \alpha I)^{-1} \Psi_n S_{n-1}^T & (\text{MSB1}), \text{ or} \\ \Psi_n (\Psi_n Y_{n-1}^T Y_{n-1} \Psi_n + \alpha I)^{-1} \Psi_n Y_{n-1}^T & (\text{MSB2}). \end{cases} \quad (\text{III.17})$$

We turn next to preconditioning. We propose rescaling elements of the density  $\rho_n$  to account for multiple scales between the interstitial electrons and the muffin tin electrons. Such a preconditioning is generically represented by multiplying the density  $\rho_n$  at each iteration  $n$  *on the left* by an arbitrary invertible diagonal matrix  $\Omega_n$ . One need not change any of the formalism above; specifically, one replaces  $Y_n$ ,  $S_n$ , and  $A_n$  in Eq.(III.4) with  $\widehat{Y}_n := \Omega_n Y_n$ ,  $\widehat{S}_n := \Omega_n S_n$ , and,

$$A_n^{\alpha, \Psi_n, \Omega_n} := \begin{cases} \Psi_n (\Psi_n \widehat{S}_{n-1}^T \widehat{Y}_{n-1} \Psi_n + \alpha I)^{-1} \Psi_n \widehat{S}_{n-1}^T \Omega_n & (\text{MSB1}), \text{ or} \\ \Psi_n (\Psi_n \widehat{Y}_{n-1}^T \widehat{Y}_{n-1} \Psi_n + \alpha I)^{-1} \Psi_n \widehat{Y}_{n-1}^T \Omega_n & (\text{MSB2}). \end{cases} \quad (\text{III.18})$$

The preconditioner used in the numerical experiments in Section IV rescales the change in the interstitial electrons relative to that in the muffin-tin electrons. We represent the interstitial and muffin-tin portions of the residual  $g_n = F(\rho_n) - \rho_n$  by  $g_n^{(I)}$  and  $g_n^{(M)}$  respectively where  $g_n = \left( g_n^{(I)T}, g_n^{(M)T} \right)^T$ . The averages of the residuals of these components separately are

$$\bar{g}_n^{(I)} = \sum_{j=0}^n \|g_j^{(I)}\| / \|g_j\|, \quad \text{and} \quad \bar{g}_n^{(M)} = \sum_{j=0}^n \|g_j^{(M)}\| / \|g_j\|, \quad (\text{III.19})$$

Our preconditioner  $\Omega_n$  is defined by

$$\Omega_n = \begin{pmatrix} \omega_n I_1 & 0 \\ 0 & I_2 \end{pmatrix} \quad \text{where} \quad \omega_n = \sqrt{\frac{\bar{g}_n^{(M)}}{\bar{g}_n^{(I)}}} \quad (\text{III.20})$$

and  $I_j$  is the  $l_j \times l_j$  identity matrix where  $l_j$  is the dimension of the interstitial/muffin tin electrons respectively. We note that the  $\omega_n$  term enters the multiseant form squared, hence our use of a square root. Removing this square root is also reasonable, and in some cases is better in numerical tests, but it can be less stable and lead to runaway behavior where the interstitial regions converge too rapidly. More sophisticated preconditioning are also plausible, for instance a dielectric term for the plane waves<sup>41,42</sup>, though we found this simple form to be very effective.

Before concluding this section we note that, by construction, the vectors  $s_n$ ,  $y_n$  conserve charge, as does the residual  $g_n$ , and the preconditioners and normalizations do not have any effect on the charge. The result will then conserve charge automatically within numerical accuracy, so no explicit charge constraint is necessary.

### C. Step Control and the generating matrix $A_0$

In Broyden's original numerical experiments he constructed the initial matrix  $A_0$  from a finite difference approximation to the true Jacobian (see<sup>13</sup> (Section 7)). This is not a practical approach for DFT calculations. The convention for the initial estimate is a scaling of the identity; that is, at each iteration  $n$  we choose  $A_{0,n} = \sigma_n I$ . From the previous analysis, the magnitude of the unpredicted step depends upon  $\sigma_n$ , increasing linearly for both B2 and MSB2 and in general increasing for both B1 and MSB1 as well although in a more complicated fashion. Therefore, by controlling  $\sigma_n$  we control the size of the step in the unpredicted direction. The choice of the scaling is critical – if it is poorly chosen iterations can stagnate or diverge. A more technical discussion of strategies for choosing  $\sigma_n$  are intimately connected to a convergence analysis of the algorithms, which is the topic of subsequent work.

For our purposes it suffices to give a number of effective controls. Our strategy for implementing a dynamic step length  $\sigma_n$  has three parts. First we constrain  $\sigma_n$  so that the step in the direction of the unpredicted component has an upper bound that is proportional to the size of the predicted component:

$$\sigma_n \leq R|p_n|/|g_n| \quad (\text{III.21})$$

where  $R$  is a fixed parameter. One has to take some step along this component, as otherwise no new information is generated; however if too large a step is taken the algorithm can diverge. As a second level of control, we bound the total variation between successive scalings:

$$\tilde{\sigma}_n = \sigma_{n-1} * \max(0.5, \min(2.0, \|g_{n-1}\|/\|g_n\|)). \quad (\text{III.22})$$

Note that we do not reject steps that yield a larger residual  $g_n$ , but rather reduce the size of the step in the unpredicted direction. In almost all cases a large improvement is achieved in the next step by retaining the bad step. As a third level of control, we include an upper bound on the absolute value of the scaling,  $\bar{\sigma}$ .

Of these controls on  $\sigma_n$  our numerical experience is that the parameter  $R$  is the most important. For hard problems we have found that a value of  $R$  from 0.05 to 0.15 works well. The upper bound on  $\sigma_n$  that we have found to be effective is  $\bar{\sigma} \approx 0.1 - 0.2$ . These values are problem specific, however, and may fail for examples we have not considered. An automatic dynamic choice for  $\sigma_n$  in conjunction with standard trust region strategies is the subject of future research

Before concluding we note that for the very first cycle we take a small step with

$$\sigma_0 = \bar{\sigma} * (0.1 + \exp(-2.0 * \max(dQ, dPW/3.5, dRMT))), \quad (\text{III.23})$$

where  $dQ$  is the change in the charge within the muffin tins,  $dPW$  is the change in the rescaled plane waves and  $dRMT$  is the change of the density within the muffin tins. This form is based upon numerical experience with WIEN2k, and is somewhat conservative.

### D. Summary

#### Algorithm III.1 (Regularized, preconditioned, limited-memory multiseant method)

0. Choose an initial  $\rho_0$ ,  $\sigma_0$  according to Eq.(III.23), generate  $\rho_1 = \rho_0 + \lambda(F(\rho_0) - \rho_0)$  for  $\lambda > 0$  some appropriately chosen step length (this is the Pratt step Eq.(III.2)), set  $n = 1$  and fix  $\alpha > 0$  ( $10^{-6}$  to  $10^{-4}$ ).
1. If the convergence criterion is met, terminate. Otherwise, given  $S_{n-1}$  and  $Y_{n-1}$ , whose columns are steps  $s_j$  and residual differences  $y_j$  respectively ( $j = n - m, n - (m - 1), \dots, n - 1$  for some appropriate number of prior steps, e.g.  $m = \min\{n, 8\}$ ) centered on the current point  $\rho_n$  as in Eq.(III.3), calculate  $A_n^{\alpha, \Psi_n, \Omega_n}$  via Eq.(III.18) for either MSB2 or MSB1 with the scaling  $\Psi_n$  given by Eq.(III.16) and the preconditioner  $\Omega_n$  given by Eq.(III.19)-(III.20). Determine the value of  $\sigma_n$  according to

$$\sigma_n = \min\{\tilde{\sigma}_n, R|p_n|/|g_n|, \bar{\sigma}\} \quad (\text{III.24})$$

where  $\tilde{\sigma}_n$  is given by Eq.(III.22) and  $\bar{\sigma}$  is some appropriately chosen upper bound (0.1 to 0.2). Calculate the next step  $\rho_{n+1}$  according to Eq.(III.4) with  $A_n$  replaced by  $A_n^{\alpha, \Psi_n, \Omega_n}$ .

2. Evaluate  $F(\rho_{n+1})$ , set  $n = n + 1$  and repeat Step 1.

#### IV. RESULTS

We test the performance of the algorithm on five examples of increasing physical difficulty, all run using the WIEN2k code<sup>4</sup> and the PBE functional<sup>43</sup>; we provide the details below with technical information so they can be reproduced as well as reasons for their choice.

**Model 1** Simple bulk MgO, spin-unpolarized with RMT's of 1.8 a.u., an RKMAX of 7 and a  $5 \times 5 \times 5$   $k$ -point mesh and a Mermin-functional<sup>2</sup> (i.e. Fermi-Dirac distribution) with a temperature of 0.0068eV. This is a very easy to solve problem.

**Model 2** Bulk Pd, spin-unpolarized with RMT's of 2.0 a.u., an RKMAX of 7.5, a  $5 \times 5 \times 5$   $k$ -point mesh and a Mermin-functional with a temperature of 0.0068eV. This is slightly harder because of the possibility of sloshing between the  $d$ -electron states and the fact that one should use a larger sampling of reciprocal space.

**Model 3** A bulk silicon cell with an RMT of 2.16 a.u., an RKMAX of 7.0, a  $6 \times 6 \times 6$   $k$ -point mesh and a Mermin-functional with a temperature of 0.0013eV.

**Model 4** A  $2 \times 2 \times 2$  Pd supercell with a vacancy at the origin, RMT's of 2.5 a.u., an RKMAX of 6.5, a  $k$ -point mesh of  $3 \times 3 \times 3$  and a Mermin-functional with a temperature of 0.0068eV. Here, in addition to sloshing between  $d$ -electron states one can have longer-range dielectric sloshing. In addition, this is a poorly constructed problem because the RKMAX is too small as is the  $k$ -point mesh.

**Model 5** A  $4.757 \times 4.757 \times 34.957$  a.u., spin-polarized (111) fcc nickel surface with seven atoms in the range  $-1/3 \leq z \leq 1/3$ . Technical parameters were RMTs of 2.13, an RKMAX of 7 and a  $11 \times 11 \times 1$   $k$ -point mesh, also with a Mermin-function temperature of 0.0068eV. It should be noted that the two surfaces are sufficiently close together, so there is real electron density in the vacuum. In this case one can have spin sloshing,  $d$ -electron sloshing as well as long-range Coulomb sloshing of electrons in the vacuum.

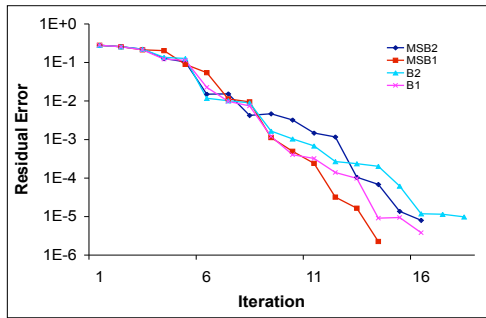
In all cases we started from densities calculated as a sum of independent atoms, and the calculations were run with both forms of Broyden multiseccants given by Eq.(II.17) and Eq.(II.15), as well as the more conventional Broyden first Eq.(II.9) and second Eq.(II.11) methods. Convergence criteria were an energy change of  $10^{-5}$  Rydbergs and an RMS convergence of the charge within the muffin tins of  $10^{-5}$  electrons. For the multiseccant implementations eight prior memory steps were used. To simplify the results, unless noted otherwise we used fixed values of the regularization parameter  $\alpha$  of  $10^{-4}$  and  $R = 0.1$ . In almost all cases, Figure 2 shows that the convergence appears to be linear, although the precision of the calculations does not allow one to observe the final asymptotic behavior, including rates of convergence, of the algorithms.

TABLE I: Iterations to convergence as a function of  $\sigma$  for models 1 – 5 with fixed  $\alpha = 10^{-4}$  and  $R = 0.1$ . The mean and standard deviation are for  $\bar{\sigma}$  between 0.05 and 0.8 for Models 1 – 3, 0.05 to 0.5 for Models 4 and 5.

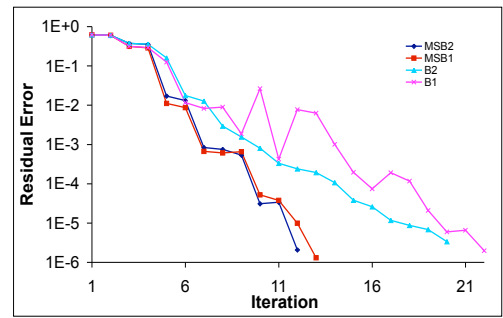
	MSB2		MSB1		B2		B1	
	mean	stdev	mean	stdev	mean	stdev	mean	stdev
Model 1	16.22	0.44	14.56	1.01	18.67	1.66	22.44	12.71
Model 2	12	0	12.89	0.33	20.11	0.78	39.78	9.58
Model 3	15.44	2.35	16.78	0.67	25.22	2.95	57.56	7.76
Model 4	24.17	2.04	29.17	4.92	–	–	–	–
Model 5	54.60	3.51	–	–	–	–	–	–

For the very simple Model 1 all the methods converge quickly and the parameter  $\bar{\sigma}$  has no significant impact on performance. The MSB1 method is slightly faster, but as the latter results indicate this is an exception. If  $\bar{\sigma}$  is too small (below 0.025) convergence is slower. Interestingly, even for this very simple case the multiseccant methods are significantly faster.

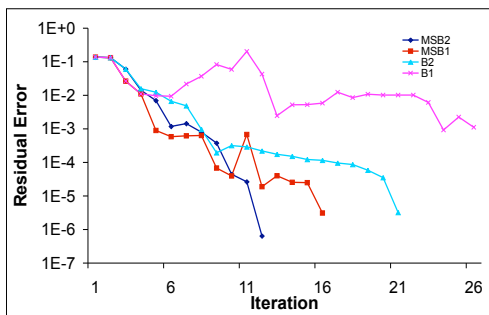
For the slightly more complicated Model 2, both multiseccant methods converge rapidly, whereas the B2 method converges more slowly and the B1 method is worst by a significant margin. The principal difference between Models 1 and 2 is that in Model 1 there are large changes during the iterations both within the muffin-tins, as well as for the plane waves, whereas in Model 2 almost all the changes are in the plane waves. This supports the rule-of-thumb discussed earlier that one should make the muffin tins as large as possible without overlapping.



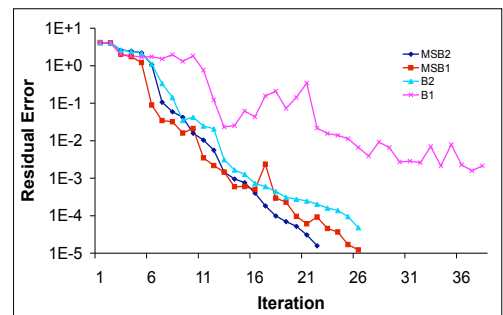
(a)



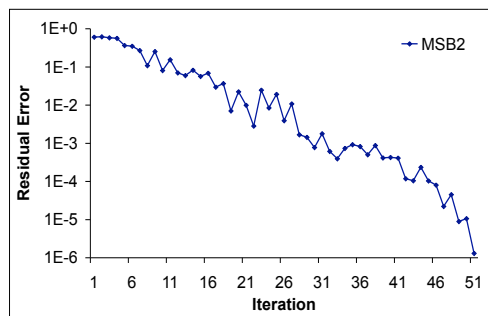
(b)



(c)



(d)



(e)

FIG. 2: (Color online) Plot of the convergence for models 1-5 (frames (a)-(e) respectively). using the multiseant update based on Broyden’s first method (MSB1) and second method (MSB2), compared to Broyden’s second method (B2) and Broyden’s first method (B1). In model 5 the only algorithm to converge is MSB2.

With Model 3 the multiseant methods significantly outperform the classical secant methods. For bulk silicon much of the covalent bonding lies in the interstitial region. We conjecture, therefore, that the improvement is due to the improved step direction and size for the multiseant methods that allow these methods to handle the greater variations of the Kohn -Sham mapping for this basis set.

The same trend continues with both Model 4 and Model 5 to the extent that B2 and B1 only converge for “good” values of  $\bar{\sigma}$  (which have to be found by trial and error) and in many cases diverge. For the hardest problem we report here, Model 5, only the MSB2 method converged. If one added a line search the other methods would probably converge albeit less rapidly and with a many more SCF evaluations.

The  $\sigma$  parameter in the MSB2 update gives one direct control over the size of the steps, which is an important feature for models with strong variations. The control of steps is less immediate for the MSB1 update and involves

a more sensitive coupling of the regularization parameter  $\alpha$  and the step size parameter  $\sigma$ . This is illustrated by the greater variance in performance of the MSB1 update versus MSB2 for models 1, 2, 4, and 5 shown in Table I.

## V. DISCUSSION

To summarize the main points of this work:

- We argue that for DFT problems, where many physically interesting models result in noncontractive SCF mappings, one should consider the information from previous points of the SCF cycle more as samples of a higher-dimensional space than as part of a deterministic path. As a consequence multiseant methods are better than sequential secant updates, as born out in the results.
- There is a fundamental difference between methods based upon Broyden’s first (B1) and second (B2) methods in terms of the space they operate in. The second method is more robust and handles poorly constructed, (nearly) ill-posed problems better – in general these are the more interesting physical problems.
- Scaling, regularization and preconditioning have a significant impact on algorithm performance. Moreover, regularization acts simultaneously to reduce instabilities due both to linear dependencies as well as to deficiencies in the model.
- Controlling the step size  $\sigma_n$  along the direction about which no information is available is critical. For difficult problems, this step should in general be *smaller* than for easy problems.
- The multiseant method based upon Broyden’s second formulation (MSB2) with appropriate safeguards simply and quickly solves problems which may defeat a novice, sometimes even an expert.

The method we have detailed (MSB2) is robust and has been part of the main WIEN2k distribution since August 2007 without any apparent problems. Even in the hands of an experienced user for complicated problems such as LDA+U we have been told of cases where the MSB2 version is three times faster than the earlier B2 code. The default values of  $\alpha = 10^{-4}$  and  $R = 0.1$  will be approximately correct for a pseudopotential code where preconditioning the variables is not necessary though there are strong variations Kohn-Sham mapping. We have not attempted to implement the MSB2 algorithm for a pseudopotential code but see no reason why it should not work at least as well. One can of course adjust these parameters to improve a single problem, but we recommend values that perhaps are slightly slower in a few cases, but more robust for a wide variety of problems. There may also be ways to stabilize MSB1 so that it could possibly work better for pseudopotential codes where preconditioning is easier.

We acknowledge that we have only considered relatively small problems here, but experience indicates that the convergence depends only very weakly (if at all) upon the size of the problem either in terms of the size of the basis set or the number of atoms. For instance, for a h-BN/Rh(111) nanomesh slab of 1108 atoms<sup>44</sup> with the earlier Broyden mixing algorithm it did not converge even after 200 iterations, but did in 30-40 with the new algorithm. For other large structures, for instance a Si (111) 7x7 surface with 498 atoms, starting from neutral atoms the convergence is only slightly slower than it is for Model 3. We emphasize once again the link between convergence of the mixing process and the functional properties of the underlying Kohn-Sham mapping. A poorly constructed problem will in most cases converge much more slowly than a well constructed one; a single atom may converge slower than  $10^4$  atoms. This may be a consequence of short-cuts in the DFT calculation, e.g. too few  $k$ -points or numerical errors in an iterative diagonalization, or it can be due to a poorly constructed Hamiltonian or perhaps density functional. For the general user poor convergence should be taken as a suggestion that the model of the physics may not have been properly constructed.

Some additional comments are appropriate about the role of the term in the regularization. As mentioned earlier, we are using this *simultaneously* in three ways, firstly as a standard regularization technique to avoid ill-conditioning associated with near linear dependence of the columns of  $Y_n$ , secondly as a Levenberg-Marquardt-type strategy to rotate the step and thirdly in a standard Wiener filter sense to account for model uncertainty. The regularization parameter can be considered to scale proportionally to the noise or uncertainty in the secant equations. Far from the solution the quasi-Newton step may not be appropriate, suggesting that one should use a larger regularization. Similarly, near the solution if the quasi-Newton step is accurate, it will yield faster rates of convergence, in which case one would choose a smaller regularization. While one could dynamically adjust the regularization parameter, for our numerical experiments we choose a relatively large fixed value of  $\alpha$  ( $10^{-4}$ ). This, in our experience, yields adequate overall convergence and better convergence in the “dangerous” early stages of the iterations.

The fact that we that we obtain improvement under the assumption that most models of physical interest do not lead to contractive, or more generally monotone SCF mappings raises some questions. It is well established that

current density functionals are inexact descriptions of the physics, but the exact analytic properties of many physical systems are unknown. In particular, for many systems it is not known whether the SCF operator is monotone, let alone that it has fixed points, although it is hard to conceive of an experimentally observable equilibrium structure that does not have fixed points. An interesting question to raise is whether the SCF operator is monotonic with the "true" density functional that correctly describes the physics. Since in many cases the effective potential  $V_\rho$  has no closed form, it is not known whether many of these theoretical properties are verifiable. It is tempting to infer analytic properties from numerical experiments – and we have made numerical progress by doing just this – but one cannot on numerical evidence alone determine the extent to which numerical behavior is indicative of the true nature of the physical system. As a final speculation, we raise the question of whether the character of the SCF mapping can be experimentally measured, or whether this type of behavior is a mathematical anomaly resulting from being much further away from equilibrium than any feasible experimental system will ever be.

There are several directions of research with regard to algorithms. Firstly, the heuristics for adjusting the step size  $\sigma_n$  need to be put on firm mathematical footing. This would accompany a study of the asymptotic behavior of the algorithm and is the subject of future research. While the analysis of Eq.(III.6) has attractive interpretations in terms of nearest points in the range and space orthogonal to the prior data, the notion of "nearest" is with respect to the usual Euclidean ( $L^2$ ) norm, which is biased towards outliers. One could consider the development of algorithms based on weighted norms, or even non-Euclidean prox mappings as opposed to those detailed in Subsection II C. The  $\Omega_n$  considered by<sup>5,6,10</sup> is in the spirit of weighted norms. Other areas for improvement could be found in the initialization of the iterations. We used the Pratt step, however one could use information from a previous SCF iteration.

Finally, while we have used some physics in helping to design the algorithm, there may be more that could be exploited. We find particularly appealing the observation discussed at the beginning of Section III that the density appears to be separable into distinct subsets. One might envision tailoring algorithms to exploit this property. For instance, one could iterate on the components of the density associated with the muffin-tins, while holding the interstitial electron density fixed. Alternatively, one could iterate on the  $sp$ -electron density holding the  $d$ -electron density fixed, or one could iterate on other observables such as the spin associated with a particular atom. Such an approach might allow one to isolate irregular variables within the SCF mapping and design algorithms accordingly. This general approach is known as *operator splitting* about which there is a vast literature. (see, for instance<sup>45</sup> and references therein). This would allow one to isolate the analytical properties of the SCF operator and work more directly with specific physical quantities.

### Acknowledgments

This work was funded by NSF under Grants #DMR-0455371/001 (LDM) and #DMS-0712796 (DRL).

---

\* Electronic address: [rluke@math.udel.edu](mailto:rluke@math.udel.edu)

<sup>1</sup> W. Kohn and S. L. J., Phys. Rev. **140**, A 1133 (1965).

<sup>2</sup> N. D. Mermin, Phys. Rev. **137**, A 1441 (1965).

<sup>3</sup> E. Prodan, J. Phys. A. **38**, 5647 (2005).

<sup>4</sup> P. Blaha, K. Schwarz, G. Madsen, D. Kvasnicka, and J. Luitz, *WIEN2k, An Augmented Plane Wave + Local Orbitals Program for Calculating Crystal Properties* (Institute for Materials Chemistry, TU Vienna, <http://www.wien2k.at/>, 2006).

<sup>5</sup> F. Crittin and M. Bierlaire, in *Proceedings of the 3rd Swiss Transport Research Conference* (STRC, 2003).

<sup>6</sup> D. D. Johnson, Phys. Rev. B **38**, 12807 (1988).

<sup>7</sup> M. Kawata, C. M. Cortis, and R. A. Friesner, J. Chem. Phys. **108**, 4426 (1998).

<sup>8</sup> P. Pulay, Chem. Phys. Lett. **73**, 393 (1980).

<sup>9</sup> G. P. Srivastava, J. Phys. A: Math. Gen. **17**, L317 (1984).

<sup>10</sup> D. Vanderbilt and S. G. Louie, Phys. Rev. B **30**, 6118 (1984).

<sup>11</sup> D. Singh, H. Krakauer, and C. S. Wang, Phys. Rev. B **34**, 8391 (1986).

<sup>12</sup> S. Boyd and L. Vandenberghe, *Convex Optimization* (Oxford University Press, New York, 2003).

<sup>13</sup> C. G. Broyden, Mathematics of Computation **19**, 577 (1965).

<sup>14</sup> R. H. Byrd, J. Nocedal, and R. B. Schnabel, Math. Prog. **63**, 129 (1994).

<sup>15</sup> J. E. Dennis and J. J. Moré, SIAM Rev. **19**, 46 (1977).

<sup>16</sup> C. M. IP and M. J. Todd, SIAM J. Numer. Anal. **25**, 206 (1988).

<sup>17</sup> P. Wolfe, Comm. ACM **12**, 12 (1959).

<sup>18</sup> J. G. P. Barnes, Comput. J. **8**, 66 (1965).

<sup>19</sup> J. M. Ortega and W. C. Rheinboldt, *Iterative Solution of Nonlinear Equations in Several Variables* (Academic Press, New York, 1970).

- <sup>20</sup> W. Gragg and G. Stewart, *SIAM J. Numer. Anal.* **13**, 889 (1976).
- <sup>21</sup> D. M. Gay and R. B. Schnabel, in *Nonlinear Programming*, edited by O. L. Mangasarian, R. R. Meyer, and S. M. Robinson (Academic Press, New York, 1978), vol. 3, pp. 245–281.
- <sup>22</sup> J. M. Martínez, *BIT* **19**, 236 (1979).
- <sup>23</sup> R. B. Schnabel, Tech. Rep. CU-CS-247-83, University of Colorado, Boulder (1983).
- <sup>24</sup> J. E. Dennis and R. Schnabel, *Numerical Methods for Unconstrained Optimization and Nonlinear Equations* (Prentice Hall, Englewood Cliffs, NJ, 1996).
- <sup>25</sup> J. Ford and I. Moghrabi, *J. Comp. Appl. Math.* **82**, 105 (1997).
- <sup>26</sup> H. Fang and Y. Saad, Tech. Rep., Department of Computer Science and Engineering; University of Minnesota (2007).
- <sup>27</sup> D. R. Luke, J. V. Burke, and R. G. Lyon, *SIAM Rev.* **44**, 169 (2002).
- <sup>28</sup> J. M. Borwein and A. S. Lewis, *Convex analysis and nonlinear optimization : theory and examples* (Springer Verlag, New York, 2006), 2nd ed.
- <sup>29</sup> F. E. Browder, *Proc. Nat. Acad. Sci. U.S.A.* **54**, 1041 (1965).
- <sup>30</sup> W. A. Kirk, *Amer. Math. Monthly* **72**, 1004 (1965).
- <sup>31</sup> J. B. Baillon, R. E. Bruck, and S. Reich, *Houston J. Math.* **4**, 1 (1978).
- <sup>32</sup> G. Pratt, *Phys. Rev.* **88**, 1217 (1952).
- <sup>33</sup> E. Prodan and N. P., *J. Stat. Phys.* **111**, 967 (2003).
- <sup>34</sup> K. Levenberg, *The Quarterly of Applied Mathematics* **2**, 164 (1944).
- <sup>35</sup> D. Marquardt, *SIAM J. Appl. Math.* **11**, 431 (1963).
- <sup>36</sup> A. N. Tihonov, *Dokl. Akad. Nauk SSSR* **151**, 501 (1963).
- <sup>37</sup> A. N. Tihonov, *Dokl. Akad. Nauk SSSR* **153**, 49 (1963).
- <sup>38</sup> P. C. Hansen, *Rank-Deficient and Discrete Ill-Posed Problems* (SIAM, 1998).
- <sup>39</sup> A. E. Hoerl and R. W. Kennard, *Technometrics* **12**, 55 (1970).
- <sup>40</sup> J. J. Moreau, *Comptes Rendus de l'Académie des Sciences de Paris* **255**, 2897 (1962).
- <sup>41</sup> D. Raczkowski, A. Canning, and L. W. Wang, *Phys. Rev. B* **64**, 121101 (R) (2001).
- <sup>42</sup> K. M. Ho, J. Ihm, and J. D. Joannopoulos, *Phys. Rev. B* **25**, 4260 (1982).
- <sup>43</sup> J. P. Perdew, K. Burke, and M. Ernzerhof, *Phys. Rev. Lett.* **77**, 3865 (1996).
- <sup>44</sup> R. Laskowski and P. Blaha, *Journal of Physics: Condensed Matter* **20**, 064207 (2008).
- <sup>45</sup> D. R. Luke, *SIAM J. Optim.* ((to appear)).
- <sup>46</sup> “:=” distinguishes a *definition* from an equation.
- <sup>47</sup> The Frobenious norm of a matrix is the square root of the sum of squares of the matrix entries.
- <sup>48</sup> A projection is defined as any mapping  $P$  such that  $P^2 = P$ . An orthogonal projection onto a set  $C$  is the point in  $C$  that is nearest, with respect to the norm, to the point being projected.