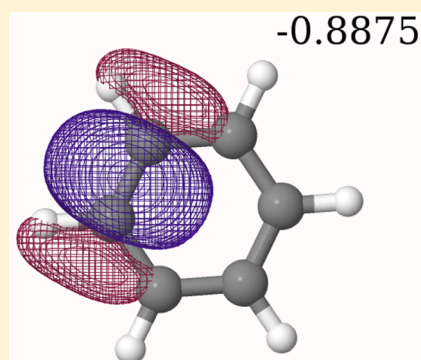


## Unitary Optimization of Localized Molecular Orbitals

Susi Lehtola<sup>\*,†</sup> and Hannes Jónsson<sup>†,‡</sup><sup>†</sup>COMP Centre of Excellence, Department of Applied Physics, School of Science, Aalto University, P.O. Box 11000, FI-00076 Aalto, Espoo, Finland<sup>‡</sup>Faculty of Physical Sciences, University of Iceland, 101 Reykjavík, Iceland

## S Supporting Information

**ABSTRACT:** A unified formalism and its implementation is presented for Foster–Boys, fourth moment, Pipek–Mezey, and Edmiston–Ruedenberg type localization schemes of molecular orbitals through unitary optimization of the localizing transform matrix using a recently proposed algorithm [Abrudan; et al. *Signal Processing* **2009**, 89, 1704]. A conjugate gradient algorithm is used with an efficient line search method. The option of using complex valued orbitals is included. Applications to fullerenes from C<sub>20</sub> to C<sub>100</sub>, as well as benzene and arachic acid are presented, showing the capability of the method, which has been implemented in ERKALE, an open source program for electronic structure calculations of atoms and molecules.



## ■ INTRODUCTION

The molecular orbitals commonly obtained from self-consistent field (SCF) calculations at the Hartree–Fock (HF) or Kohn–Sham density functional<sup>1,2</sup> (KS-DFT) level of theory are delocalized in space and typically do not correspond to the local orbitals chemical intuition is based on. This is because the orbitals obtained from the calculation are canonical molecular orbitals (CMOs) that diagonalize the Fock (or the Kohn–Sham–Fock) matrix, yielding well-defined orbital energies. A subsequent transformation from CMOs to localized, chemically intuitive orbitals is often used to better analyze the output from HF or KS-DFT calculations. This is the main motivation for the present work.

Another motivation for a localization transformation is the description of electron correlation on top of HF. When applied to the CMOs, such post-HF methods scale poorly with system size. But, correlation is mostly a local phenomenon, so orbitals can be used to construct computationally better scaling methods for treating correlation.<sup>3–5</sup> Localized orbitals can also be used at the HF level of theory for obtaining linear scaling.<sup>6,7</sup>

Because the total energy in both HF and KS-DFT is invariant under unitary transformations of the occupied–occupied and the virtual–virtual blocks, the solution to the SCF problem can also be described with the help of any set of orbitals connected to the CMOs via a unitary transformation:

$$|i\rangle \rightarrow \sum_j W_{ji} |j\rangle \quad (1)$$

where  $|i\rangle$  is the  $i$ th molecular orbital and  $\mathbf{W}^\dagger \mathbf{W} = \mathbf{1} = \mathbf{W} \mathbf{W}^\dagger$ . In particular, the unitary transformation may be constructed so that the orbitals become maximally localized. It is also possible

to sacrifice the orthogonality of the orbitals to make them more transferable between molecules<sup>9</sup> and to achieve even more locality of the individual orbitals,<sup>10,11</sup> but this complicates the physical interpretation because the density is then no more a simple sum of the individual orbital densities.

The most commonly used localization objective functions are the ones by Foster and Boys<sup>12</sup> (FB), Edmiston and Ruedenberg<sup>13</sup> (ER), and Pipek and Mezey<sup>14</sup> (PM). The FB scheme minimizes the orbital size as measured by its variance

$$\sum_i {}_2\Omega_i \quad (2)$$

where the sum over the molecular orbitals  $i$  runs over the block under treatment and the orbital variance is

$${}_n\Omega_i = \langle i | [\mathbf{r} - \langle i | \mathbf{r} | i \rangle]^n | i \rangle \quad (3)$$

The first problem with FB localization is that a small  $\sum_i {}_2\Omega_i$  does not imply that all  $\Omega_i$  are small. For this reason a generalization of the FB measure

$$\sum_i {}_2\Omega_i^p \quad p \geq 1 \quad (4)$$

has recently been suggested,<sup>15</sup> where the additional penalty for  $p > 1$  restricts the size of the largest orbital.

Second, because a small value of  ${}_2\Omega_i$  can occur even though an orbital has a significant tail, thus having significant nonlocal character, an alternative objective function was suggested<sup>16</sup>

**Received:** September 6, 2013

$$\sum_i 4\Omega_i^p \quad (5)$$

where the fourth moment (FM) adds more penalty to the tail parts of the localized orbitals. Again, for a truly local set of orbitals  $p > 1$  is necessary.<sup>16</sup>

The ER scheme maximizes the electronic self-repulsion energy

$$\sum_i (iilii) \quad (6)$$

or equivalently minimizes the interorbital repulsion energy

$$\sum_i \sum_{j \neq i} (iiljj) \quad (7)$$

where the electron repulsion integral (ERI) is

$$(ijkl) = \int \frac{\phi_i^*(\mathbf{r}) \phi_j(\mathbf{r}) \phi_k^*(\mathbf{r}') \phi_l(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} d^3r d^3r' \quad (8)$$

with  $\phi_i(\mathbf{r}) = \langle \mathbf{r} | i \rangle$ . A similar scheme has been suggested by von Niessen<sup>17</sup> (vN), with the difference that an overlap metric is used instead of the Coulomb metric:

$$\sum_i \sum_{j \neq i} \int \rho_i(\mathbf{r}) \rho_j(\mathbf{r}) d^3r \quad (9)$$

Its computational scaling in the linear combination of atomic orbitals (LCAO) treatment is similar to that of the ER procedure. The vN scheme has not become widely used and will not be discussed further here.

A fundamentally different scheme to the FB, FM, and ER schemes is the PM scheme, which contrary to the first three methods does not mix  $\sigma$  and  $\pi$  bonds and thus can correspond better to chemical intuition in (locally) planar systems. The conventional formulation of the scheme involves localization of the Mulliken charges<sup>18</sup> arising from the individual orbital densities by maximizing

$$\sum_i \sum_{\text{atoms } A} [\langle i | P_A | i \rangle]^2 \quad (10)$$

where  $P_A$  is an operator that projects on the basis functions centered on atom A, and  $\langle i | P_A | i \rangle = Q_A^i$  is the Mulliken charge on atom A arising from orbital  $i$ . Alternatively, Löwdin charges can be used instead of Mulliken charges. These have in fact been shown to yield better localized virtual orbitals.<sup>19</sup> Schemes equivalent to PM with Bader or Becke charges have also been suggested.<sup>20,21</sup> The discussion of various choices for the charges in the PM procedure will be presented elsewhere.<sup>22</sup> We note that relativistic versions of the FB, PM, and ER schemes have been developed as well.<sup>23,24</sup> Also, a localization scheme based on the electron localization function<sup>25</sup> has recently been suggested.<sup>26</sup>

In the current work, we examine the conventional scheme of first solving the SCF equations for the CMOs and then converting them to localized molecular orbitals (LMOs) with some localization procedure. An alternative approach of solving the SCF equations for the LMOs directly is also possible,<sup>6,7,13,29</sup> but this implies significant changes to the SCF procedure and will not be discussed further.

The conventional method for obtaining LMOs from the delocal CMOs, originally due to Edmiston and Ruedenberg<sup>13</sup> and later refined by Barr and Basch,<sup>8</sup> is to perform Jacobi sweeps (consecutive two-by-two rotations of orbital pairs) until

convergence to a properly localized set of orbitals is achieved. This method, however, has many shortcomings. Because only one orbital pair is treated at a time, separate rotations need to be performed for every orbital pair, and the optimal rotation of an orbital pair  $(i, j)$  may be destroyed by the optimal  $(i, j + 1)$  rotation. Due to this, the convergence may be slow, and the process can easily converge to a nonoptimal solution depending on the order in which the orbital pairs are treated. Furthermore, in the case of virtual orbitals, the Jacobi sweeps often do not converge at all. A workaround to this in the case of atom-centered basis sets has been proposed by Subotnik et al.,<sup>27</sup> where the full virtual space is divided into a molecular valence-like space and an atomic “hard virtual” space, out of which the former is easy to localize using standard methodology and the latter consists of functions that are largely unoccupied at the HF level of theory.

Progress has also been made on global optimization methods for orbital localization. This is of most importance for the ER objective function. Whereas in FB, FM, and PM localization all of the necessary matrices can be permanently stored in memory, in ER localization the most computational overhead comes from the calculation of the ERIs, which usually cannot be stored in memory due to their sheer number. Instead, the integrals are recomputed and contracted at every iteration, as in the direct SCF scheme.<sup>28</sup> In this case a global approach is much more efficient than Jacobi rotations, because in the former algorithm the integrals need only be computed once and contracted with all orbital densities, whereas in the latter they need to be recomputed for every orbital pair under treatment.

A steepest descent/steepest ascent (SDSA) approach, also presented by Edmiston and Ruedenberg,<sup>13</sup> was found to converge faster than the two-by-two orbital transformations.<sup>29</sup> Fletcher–Reeves<sup>30</sup> and Fletcher–Powell<sup>31</sup> conjugate gradient (CG) methods have also been investigated,<sup>32</sup> but these were not shown to be clearly superior to the SDSA method. Also, quasi-Newton methods such as the Broyden–Fletcher–Goldfarb–Shanno (BFGS) algorithm<sup>33</sup> have been investigated<sup>34</sup> for the FB and ER measures, but a comparison to CG or SDSA methods has not been presented.

Though gradient methods are powerful in the initial phase of an optimization, where large changes are made to the value of the objective function, the convergence near the optimum is often slow. Gradient methods focus on the best defined elements of the transformation, and the weakly defined elements are left undetermined.<sup>35</sup> Moreover, the transformation of one pair of orbitals may require adjustments to a third orbital. This is not taken into account in the gradient schemes and can lead to both premature apparent convergence and slow convergence.<sup>35</sup> To remedy this, a quadratically convergent Newton–Raphson (NR) procedure was suggested already more than 30 years ago by Leonard and Luken.<sup>35</sup>

The calculation of the Hessian for the second-order procedure adds only a small overhead for the FB and PM criteria, while reducing the necessary amount of iterations by 1–2 orders of magnitude.<sup>35</sup> However, in the case of the ER criterion the calculation of the off-diagonal elements of the Hessian is a nontrivial task because three-index integrals become necessary, requiring computationally expensive integral transforms and storage.<sup>35,36</sup>

A pure NR procedure has two general shortcomings. First, it only converges in a limited radius around the optimum, and second it converges to a saddle point if any eigenvalue of the Hessian is of the wrong sign. Leonard and Luken thus

advocated the use of a combined approach in which a first-order method is used until the radius of convergence is achieved, after which the full NR scheme is used.

More recently, a direct inversion in the iterative subspace method<sup>37</sup> has been proposed by Subotnik et al.<sup>38,39</sup> for performing orbital localization in the FB and the ER schemes, as the line searches needed for a steepest descent or conjugate gradient algorithm were found to be computationally too costly. However, the method proposed by Subotnik et al. requires the computationally challenging inverse matrix square root  $(\mathbf{R}^\dagger \mathbf{R})^{-1/2}$ ,  $\mathbf{R}$  being the orthogonal matrix that connects the CMO and LMO bases. This was found to be problematic in case of the FB objective function.<sup>38</sup> Also, the algorithm was found to be prone to converge on saddle points, for which a Hessian analysis was implemented.<sup>39</sup>

The use of a trust region (TR) method has also been recently suggested.<sup>40</sup> The TR method can be viewed as a smooth transition from a gradient method to a NR scheme, as far away from the optimum the search direction is given by the gradient, but near the optimum the steps resemble more NR. This makes the method unconditionally convergent and also avoids saddle point convergence.<sup>41</sup> The TR method has been shown to be able to localize the virtual spaces of large systems with diffuse basis sets in little time.<sup>42</sup> This or the Leonard–Luken scheme should be the method of choice for local correlation methods, for which a full set of local virtual orbitals are essential.

However, the occupied space is much easier to localize. As complex orbitals have recently been found to be indispensable<sup>43,44</sup> for the proper treatment of the Perdew–Zunger self-interaction correction<sup>45</sup> (PZ-SIC) of KS-DFT, which is mathematically similar<sup>a</sup> to ER localization,<sup>46</sup> the possible gain of complex degrees of freedom for conventional orbital localization methods needs to be examined as well. But, more often than not, unitarity is equated to orthogonality in the literature; indeed, all of the aforementioned localization algorithms have been restricted to real transformations.

The algorithms presented in the current work allow for complex-valued transformations that provide more freedom (and thus may provide better localization) than transformations restricted to real space. We present an implementation of a recently developed conjugate gradient unitary optimization algorithm<sup>47,48</sup> in the context of orbital localization. It has previously been used in the context of relativistic calculations<sup>23</sup> with an approximate line search; our implementation includes an efficient exact line search.

Although a TR scheme has good convergence properties, we are not aware of a complex unitary TR algorithm. A complex unitary NR scheme has, though, been published.<sup>49</sup> A complex unitary trust region approach should be relatively straightforward to implement following, e.g., the procedure of ref 40.<sup>b</sup> However, in the context of ER localization scheme and the PZ-SIC method, a scheme involving the Hessian would be computationally very costly and we are not aware of any work on second-order methods for these applications. The gradient scheme presented here is easy to implement, is applicable to complex transformations, and is significantly better than Jacobi rotations.

The organization of the manuscript is the following. In the next section, we present a brief outline of the unitary optimization algorithm. In the Implementation section we give the formulation of the FB, FM, PM, and ER type objective functions and their derivatives and discuss possible starting points for the unitary optimization. In the Results section we

present the application of the methods, and we conclude with a brief summary.

## METHOD

The unitary optimization method we base our implementation on has been presented in refs 47 and 48, and only a brief outline is given here. Given an objective function  $\mathcal{J}(\mathbf{W})$ , where  $\mathbf{W}$  is a unitary matrix  $\mathbf{W}^\dagger \mathbf{W} = \mathbf{1} = \mathbf{W} \mathbf{W}^\dagger$  and a starting point  $\mathbf{W}_0$ , the Euclidean derivative

$$\mathbf{\Gamma}_k = \left. \frac{\partial \mathcal{J}}{\partial \mathbf{W}^*} \right|_{\mathbf{W}_k} \quad (11)$$

is computed and transformed into the unitary manifold to obtain the Riemannian derivative

$$\mathbf{G}_k = \mathbf{\Gamma}_k \mathbf{W}_k^\dagger - \mathbf{W}_k \mathbf{\Gamma}_k^\dagger \quad (12)$$

The search ascent direction is set to

$$\mathbf{H}_k = \mathbf{G}_k + \gamma_k \mathbf{H}_{k-1} \quad (13)$$

where  $\gamma_0 = 0$  and for  $k \geq 1$  the update factor is

$$\gamma_k^{\text{SDSA}} = 0 \quad (14)$$

$$\gamma_k^{\text{CGFR}} = \frac{\langle \mathbf{G}_k, \mathbf{G}_k \rangle}{\langle \mathbf{G}_{k-1}, \mathbf{G}_{k-1} \rangle} \quad (15)$$

$$\gamma_k^{\text{CGPR}} = \frac{\langle \mathbf{G}_k, \mathbf{G}_k - \mathbf{G}_{k-1} \rangle}{\langle \mathbf{G}_{k-1}, \mathbf{G}_{k-1} \rangle} \quad (16)$$

in which  $\langle \mathbf{X}, \mathbf{Y} \rangle = \text{Re Tr } \mathbf{X} \mathbf{Y}^\dagger / 2$  and CGFR and CGPR correspond to the Fletcher–Reeves<sup>30</sup> and Polak–Ribière–Polyak<sup>50,51</sup> choices for the conjugate gradient update factor. If  $\langle \mathbf{H}_k, \mathbf{G}_k \rangle < 0$  after the update, the search ascent direction is reset to  $\mathbf{H}_k = \mathbf{G}_k$ . Additionally, because in an  $N$ -dimensional vector space only  $N$  vectors can be orthogonal, the search direction is also reset if  $(k-1) \bmod N = 0$ .

Having the search direction  $\pm \mathbf{H}_k$ , where the positive sign corresponds to maximization and the negative to minimization, a line optimization of the objective function is performed

$$\left[ \frac{\partial}{\partial \mu} \mathcal{J}(e^{\pm \mu \mathbf{H}_k} \mathbf{W}_k) \right]_{\mu_{\text{opt}}} = 0 \quad (17)$$

yielding an optimal step size  $\mu_{\text{opt}}$ . The transformation matrix is then updated,

$$\mathbf{W}_{k+1} = e^{\pm \mu_{\text{opt}} \mathbf{H}_k} \mathbf{W}_k \quad (18)$$

$k$  is incremented and the procedure restarts from the computation of the derivative in eq 11. The iteration is continued until the norm of the gradient  $G_k = \langle \mathbf{G}_k, \mathbf{G}_k \rangle$  decreases below the user-adjustable threshold. Following Høyvik et al.,<sup>40</sup> a convergence threshold of  $G_k \leq 10^{-5}$  is used in the current work.

The line search in eq 17 is performed on a grid with a length scale determined by the highest frequency component of the search direction  $\mathbf{H}_k$  as

$$T_\mu = 2\pi / q \omega_{\text{max}} \quad (19)$$

where  $\omega_{\text{max}} = \max_i |\omega_i|$  and  $q$  is the order in which  $\mathbf{W}$  appears in the Taylor series of  $\mathcal{J}(\mathbf{W})$ :  $q = 4$  for PM and ER,  $q = 4p$  for FB, and  $q = 8p$  for FM.

## IMPLEMENTATION

The method with line optimizations based on an Armijo line search,<sup>47</sup> or the exact solution of eq 17 based on polynomial interpolation or Fourier transform<sup>48</sup> has been implemented in the ERKALE program.<sup>52,53</sup> The matrix exponentials in eqs 17 and 18 are computed by diagonalization<sup>c</sup>. For large problems, other methods<sup>54</sup> may be pursued.<sup>d</sup>

We have implemented the FB, FM, and ER schemes, as well as PM with Mulliken and Löwdin charges. Below, we give the objective functions and their derivatives for the various localization schemes.  $i$  and  $j$  denote LMO indices, and  $r, s, t$ , and  $u$  denote CMO indices.

Because self-consistent field (SCF) theory is rarely performed with complex coefficients,<sup>55,56</sup> we assume real basis functions and canonical orbitals, without imposing any additional limitations on the method.

**Foster–Boys.** The objective function to be minimized in the FB scheme is given by eqs 3 and 4

$$\mathcal{J}^{\text{FB}}(\mathbf{w}) = \sum_i [r_{ii}^2 - \sum_{\alpha \in x,y,z} (r_{ii}^\alpha)^2]^p \quad (20)$$

where the matrix elements in the LMO basis can be written as

$$r_{ii}^2 = \sum_{rs} W_{ri}^{*2} W_{si} \quad (21)$$

$r_{rs}^2 = \langle r | r^2 | s \rangle$  being the matrix element in the CMO basis. The derivative is found to be

$$\begin{aligned} \frac{\partial \mathcal{J}^{\text{FB}}}{\partial W_{tj}^*} &= p [r_{jj}^2 - \sum_{\alpha \in x,y,z} (r_{jj}^\alpha)^2]^{p-1} \\ &\times \left[ \frac{\partial r_{jj}^2}{\partial W_{tj}^*} - \sum_{\alpha \in x,y,z} 2r_{jj}^\alpha \frac{\partial r_{jj}^\alpha}{\partial W_{tj}^*} \right] \end{aligned} \quad (22)$$

where the derivative of the matrix element is

$$\frac{\partial r_{ii}^2}{\partial W_{tj}^*} = \sum_s \delta_{ij} r_{ts}^2 W_{sj} \quad (23)$$

The necessary moment matrices are generated via recurrence relations.<sup>57</sup>

**Fourth Moment.** The FM objective function can be written as

$$\begin{aligned} \mathcal{J}^{\text{FM}} &= \sum_i 4\Omega_i^p = \sum_i [\langle \varphi_i | [\mathbf{r} - \langle \varphi_i | \mathbf{r} | \varphi_i \rangle]^4 | \varphi_i \rangle]^p \\ &= \sum_i [\langle \varphi_i | \sum_{\alpha\beta} [r^\alpha - \langle \varphi_i | r^\alpha | \varphi_i \rangle]^2 [r^\beta - \langle \varphi_i | r^\beta | \varphi_i \rangle]^2 | \varphi_i \rangle]^p \end{aligned} \quad (24)$$

where  $\alpha, \beta \in x, y, z$ . Compacting the notation, it is easy to see that the expectation value (taken over the  $i$ th orbital) can be written as

$$\begin{aligned} &\langle (r^\alpha - \langle r^\alpha \rangle)^2 (r^\beta - \langle r^\beta \rangle)^2 \rangle \\ &= \langle \sum_{\alpha\beta} (r^\alpha r^\beta)^2 \rangle - 4 \sum_{\alpha} \langle r^\alpha \mathbf{r}^2 \rangle \langle r^\alpha \rangle + 2 \langle \mathbf{r}^2 \rangle \langle \sum_{\alpha} \langle r^\alpha \rangle^2 \rangle \\ &\quad + 4 \sum_{\alpha\beta} \langle r^\alpha \rangle \langle r^\alpha r^\beta \rangle \langle r^\beta \rangle - 3 \langle \sum_{\alpha} \langle r^\alpha \rangle^2 \rangle^2 \end{aligned} \quad (25)$$

where  $\mathbf{r}^2 = x^2 + y^2 + z^2$ . The derivative is somewhat more involved

$$\begin{aligned} \frac{\partial \mathcal{J}^{\text{FM}}}{\partial W_{tj}^*} &= p 4\Omega_i^{p-1} \left[ \frac{\partial \langle \sum_{\alpha\beta} (r^\alpha r^\beta)^2 \rangle}{\partial W_{tj}^*} \right. \\ &\quad - 4 \sum_{\alpha} \left( \frac{\partial \langle r^\alpha \mathbf{r}^2 \rangle}{\partial W_{tj}^*} \langle r^\alpha \rangle + \langle r^\alpha \mathbf{r}^2 \rangle \frac{\partial \langle r^\alpha \rangle}{\partial W_{tj}^*} \right) \\ &\quad + 2 \frac{\partial \langle \mathbf{r}^2 \rangle}{\partial W_{tj}^*} \langle \sum_{\alpha} \langle r^\alpha \rangle^2 \rangle + 4 \langle \mathbf{r}^2 \rangle \left( \sum_{\alpha} \langle r^\alpha \rangle \frac{\partial \langle r^\alpha \rangle}{\partial W_{tj}^*} \right) \\ &\quad + 8 \sum_{\alpha\beta} \langle r^\alpha \rangle \langle r^\alpha r^\beta \rangle \frac{\partial \langle r^\beta \rangle}{\partial W_{tj}^*} \\ &\quad + 4 \sum_{\alpha\beta} \langle r^\alpha \rangle \frac{\partial \langle r^\alpha r^\beta \rangle}{\partial W_{tj}^*} \langle r^\beta \rangle \\ &\quad \left. - 12 \left( \sum_{\alpha} \langle r^\alpha \rangle^2 \right) \left( \sum_{\beta} \langle r^\beta \rangle \frac{\partial \langle r^\beta \rangle}{\partial W_{tj}^*} \right) \right] \end{aligned} \quad (26)$$

where the expectation values are now taken over the  $j$ th orbital.

**Pipek–Mezey.** The objective function to maximize in the PM scheme is

$$\mathcal{J}^{\text{PM}}(\mathbf{w}) = \sum_{i=1}^N \sum_{A=1}^n [Q_{ii}^A]^2 \quad (27)$$

and its derivative is

$$\frac{\partial \mathcal{J}^{\text{PM}}}{\partial W_{tj}^*} = \sum_{A=1}^n 2Q_{jj}^A \frac{\partial Q_{jj}^A}{\partial W_{tj}^*} \quad (28)$$

The atomic charges in the CMO basis are

$$Q_{rs}^A = \sum_{\mu \in A} \frac{1}{2} [c_{\mu r}(\mathbf{S}\mathbf{c})_{\mu s} + c_{\mu s}(\mathbf{S}\mathbf{c})_{\mu r}] \quad (29)$$

in the Mulliken case and

$$Q_{rs}^A = \sum_{\mu \in A} (\mathbf{S}^{1/2}\mathbf{c})_{\mu r} (\mathbf{S}^{1/2}\mathbf{c})_{\mu s} \quad (30)$$

in the Löwdin case. Here  $c_{\mu r}$  is the coefficient of the  $\mu$ th basis function in the  $r$ th canonical orbital, and  $\mathbf{S}$  is the overlap matrix.

**Edmiston–Ruedenberg.** The objective function is

$$\mathcal{J}^{\text{ER}}(\mathbf{w}) = \sum_{i=1}^N (\text{ilii}) \quad (31)$$

and the derivative is

$$\frac{\partial \mathcal{J}^{\text{ER}}}{\partial W_{tj}^*} = 2(\text{tjljj}) \quad (32)$$

In the LCAO formulation we can recast eqs 31 and 32 to

$$\mathcal{J}^{\text{ER}} = \sum_i \text{Tr} \mathbf{P}^i \mathbf{J}^i \quad (33)$$

$$\frac{\partial \mathcal{J}^{\text{ER}}}{\partial W_{tj}^*} = 2\mathbf{c}_t^T \mathbf{J}^j \mathbf{c}_j \quad (34)$$

where  $\mathbf{P}^i$  is the orbital density matrix

$$P_{\mu\nu}^i = \text{Re} c_{\mu i}^* c_{\nu i} \quad (35)$$



and  $J^i$  is the corresponding Coulomb matrix

$$J_{\mu\nu}^i = \sum_{\sigma\rho} (\mu\nu|\sigma\rho) P_{\sigma\rho}^i \quad (36)$$

The orbital Coulomb matrices are computed using density fitting<sup>58,59</sup> with an automatically generated fitting basis set.<sup>60</sup> Alternatively, an exchange matrix  $K^i$  can be used in eqs 33 and 34 as in ref 38, because famously  $J^i = K^i$  as only a single orbital is involved.

**Starting point.** The speed of the localization procedure depends on its proper initialization. However, from eq 12 for the Riemannian derivative and the derivatives of the FB, PM, and ER objective functions (eqs 22, 28 and 32, respectively), it is clear that if the optimization is started from a real  $W_0$ , the optimization will be restricted to the domain of real matrices. Obtaining a complex localizing matrix  $W$  thus requires the optimization to be started from a complex  $W_0$ .

There are three obvious ways to start the optimization:

1. from the canonical orbitals ( $W = I$ , restriction to real  $W$ ),
2. from a random orthogonal matrix (real  $W$ ), or
3. from a random unitary matrix (complex  $W$ ).

The use of other starting guesses, such as one based on a Cholesky decomposition of the density matrix<sup>63</sup> or a natural orbital analysis,<sup>64</sup> is also possible. In our test calculations, we did not see significant performance differences between these starting guesses, as the current method is very powerful in the initial phase of the optimization, far away from the optimum. However, more elaborate schemes may be of importance in very large systems, e.g., when sparse matrix techniques become necessary.

Starting the initialization from the CMOs is usually a suboptimal choice due to their delocalized nature. Also, because the first derivatives between CMOs belonging to different symmetries vanish,<sup>32</sup> starting the optimization from the CMOs may result in spurious saddle point convergence. The symmetries of the CMOs must first be broken to allow them to mix to form the LMOs, which usually belong to different point group symmetries than the CMOs. In the current method, initialization with a random unitary matrix breaks the orbital symmetries and thus is less prone to this problem. Second-order methods such as the NR and the TR algorithms are not affected by this problem, because the Hessian information is taken into account when the search direction is chosen.

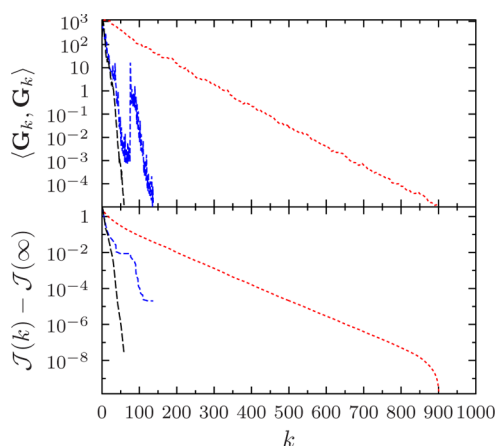
The complex degree of freedom in the unitary optimization was not found to improve on any of the localization measures, compared to an orthogonal optimization. The results shown in the following sections were obtained by starting from a random orthogonal guess.

A large penalty exponent  $p$  for the FB and FM objective functions may result in large derivatives when a bad starting guess  $W_0$  is used, causing grave problems for the optimization. For this reason we perform the localization with  $p > 1$  by recursive initialization with a localizing matrix optimized for  $p' = p - 1$ . This procedure carries little extra cost, as it allows a larger step size to be used in the initial phase where the bulk of the orbitals become localized. Increasing  $p$  then pulls in the tails of the distribution. For a similar reason we initialize FM localization with FB localized orbitals.

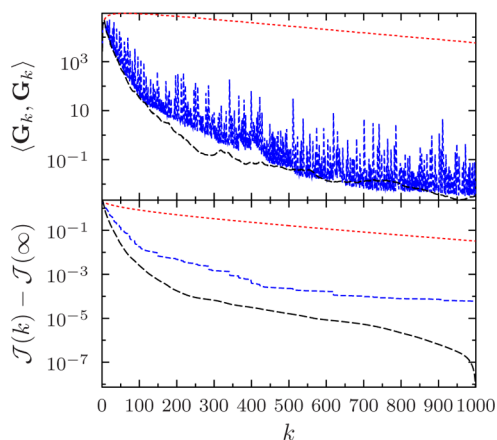
## RESULTS

**Search Direction and Line Search.** First, we examine the choice of the line search and of the search direction in an FB localization for the benzene molecule in the aug-cc-pVDZ basis set<sup>61,62</sup> (21 occupied and 171 virtual orbitals, the used geometry is available in the Supporting Information). All the calculations in the current manuscript were performed at the spin-restricted Hartree–Fock level of theory.

First, the optimal search direction is obtained. A comparison of the SDSA, CGFR, and CGPR schemes using the polynomial line search is shown in Figure 1. From the evolution of the



(a) Occupied space

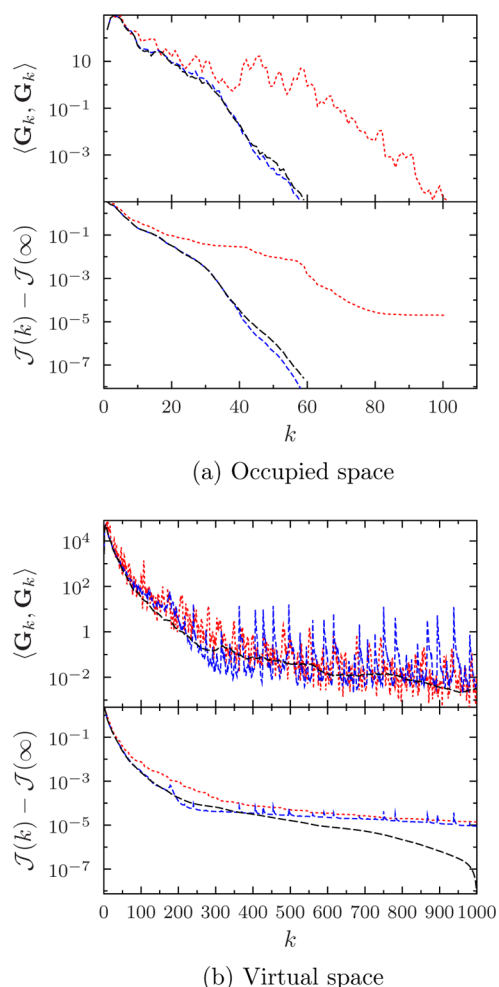


(b) Virtual space

**Figure 1.** Performance of search direction schemes in the FB localization of the MOs in benzene: CGFR in red, SDSA in blue, and CGPR in black.

SDSA gradients, it is apparent that the objective function landscape is highly structured, with abrupt changes in the function and its derivative. Though both CGFR and CGPR even out these fluctuations in the norm of the derivative, CGFR fails to accelerate the convergence compared to SDSA and actually slows it down considerably. The CGPR choice for the update factor is by far the best.

A comparison of the Armijo, Fourier transform, and polynomial line search within the steepest descent scheme is shown in Figure 2. Due to its consistent performance and low computational effort, the polynomial line search will be used in the remainder of the work.

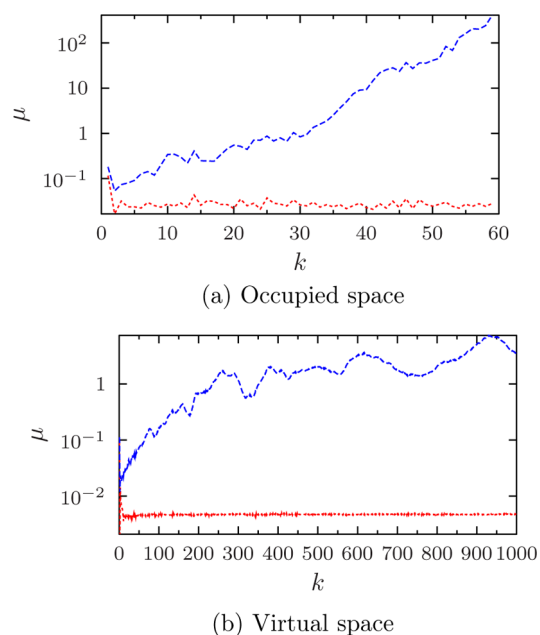


**Figure 2.** Performance of line search schemes in the FB localization of the MOs in benzene: Armijo in red, Fourier in blue, and polynomial in black.

Finally, we examine the evolution of the step size chosen by the line search algorithm and the period of the highest frequency  $T_\mu$  (defined in eq 19) in Figure 3. Surprisingly, the optimal step size is the same for a large variety of iterations and gradient norms.

**Fullerenes.** To study the performance of the optimization algorithms for a range of similar molecules spanning a wide range in size, we performed calculations on the occupied orbitals of the fullerene series  $C_{20}$ ,  $C_{40}$ ,  $C_{60}$ ,  $C_{80}$ , and  $C_{100}$ , with 60–300 occupied orbitals. The cc-pVDZ and cc-pVTZ basis sets, as well as the much more diffuse pcemd-3 and pcemd-4 basis sets<sup>65</sup> were used. The full set of results is available in the Supporting Information. A summary of the results is presented here.

Calculations were performed from five different random complex unitary and real orthogonal starting point matrices, as well as from the canonical orbitals. This demonstrates the flexibility of the optimization algorithms in optimizing from bad starting points and probes the existence of multiple local optima<sup>32,40</sup> in the objective functions. Indeed, different final values for the objective functions are found. It is possible that part of the variation in the results is caused by convergence to saddle points, but due to the line search procedure, CG methods in general are not prone to convergence on saddle points.



**Figure 3.** Evolution of the step size for FB optimization of benzene: the chosen step size  $\mu_{\text{opt}}$  in red and  $T_\mu$  in blue.

The PM objective function is the easiest to optimize, the optimization usually converging in less than 100 iterations. Furthermore, the necessary number of iterations does not seem to be size dependent. Because the present implementation performs all the possible rotations at once, this is an expected result for simple localization criteria. Quite the opposite is seen for the FB criterion, which often takes 2–5 times more iterations than PM and for which there is a clear increase in the necessary number of iterations as the system grows. This can tentatively be attributed to the following. Though the PM criterion is dominated by contributions from a few atomic regions, the FB criterion involves contributions from the tail parts of the orbitals, which are interdependent and harder to optimize. The increase of the penalty exponent  $p$  makes the optimization more difficult by requiring the tails to be more similar in size. An extreme case is the FM criterion, which places even more weight on the tail parts, nearly making the applicability of the current method questionable.

The ER criterion is computationally much more demanding due to the evaluation of the ERIs, and thus optimization was only feasible for small systems. A size dependence of the number of iterations is clearly present also in the case of ER localization.

**Benchmark.** We report here performance numbers for the unitary optimization algorithm applied to PM, FB, and ER localization of orbitals in the arachic acid molecule,  $C_{20}H_{40}O_2$  (the used geometry is available in the Supporting Information). This molecule has 88 occupied orbitals. Various basis sets have been used (Table 1), the number of resulting virtual orbitals ranging from 420 to 2322. The aug-cc-pVXZ basis<sup>61,62</sup> with diffuse functions deleted from the hydrogen atoms (plain cc-pVXZ basis) is used, as timings have previously been reported for this system in ref 42.

The PM localization of the occupied orbitals converges quickly, in fewer than 100 iterations and only a few seconds; see Table 1. Similarly to PM, the ER localization converges in a relatively small number of iterations but takes much longer computer time due to the direct evaluation of the ERIs. The ER localization of the occupied space in the cc-pVQZ basis set did

Table 1. Timings for the Unitary Optimization for Arachic Acid

	$N_{\text{virt}}$	PM				FB				ER	
		occupied		virtual		occupied		virtual		occupied	
		$n_{\text{iter}}$	$t$ , s	$n_{\text{iter}}$	$t$	$n_{\text{iter}}$	$t$ , s	$n_{\text{iter}}$	$t$	$n_{\text{iter}}$	$t$
cc-pVDZ	420	90	5.5	304	25 min 32 s	882	36	10175	9 h 22 min	66	4 h 45 min
aug-cc-pVDZ	618	44	2.7	374	1 h 38 min	1105	46	4946	14 h 3 min	61	9 h 49 min
cc-pVTZ	1132	47	2.8	544	30 h 41 min	962	40	4522	3 days 20 h	69	2 days 8 h 12 min
cc-pVQZ	2322	56	3.5	no convergence		1368	57	no convergence		no convergence	

not reach convergence within the time limit of the runs, 7 days of wall clock time.

In contrast, the convergence of the FB localization is much slower, requiring an order of magnitude more iterations. The FM optimization did not converge even for the occupied space and the smallest basis set (cc-pVDZ basis) with  $p = 1$ . As demonstrated by Figure 4, the objective function becomes

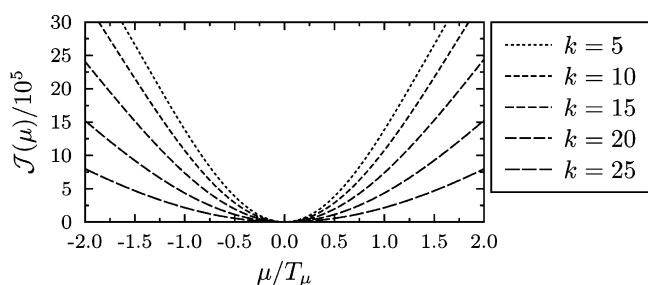


Figure 4. Line search in FM optimization for five iterations  $k$ , started from FB orbitals.

flatter as more iterations are performed. In our experience, the problems become worse with increasing penalty exponent  $p$ . This can clearly be interpreted by the importance of second-order effects as discussed in the Introduction: the optimization of an orbital pair requires adjustments to a third orbital, which fundamentally cannot be described by a gradient method. Instead, the use of a second-order method such as NR or TR is necessary to describe the coupled relaxation.

The localization of the virtual orbitals in the PM and FB schemes required an order of magnitude more iterations than that of the occupied orbitals. The reason is that the objective functions have more gradual variations for the virtual space. For the largest basis set, with 2322 virtual orbitals, the calculations did not complete for any of the methods.

## SUMMARY

We have presented the application of a unitary optimization algorithm and conjugate gradient minimization to a variety of orbital localization methods, including FB, FM, PM, and ER localization objective functions, and implementation of the method in the ERKALE<sup>52,53</sup> program. The unitary optimization algorithm performs significantly better than the commonly used Jacobi sweeps where consecutive pairwise rotations are performed, because it simultaneously optimizes all of the orbitals. The Polak–Ribière–Polyak CG algorithm was found to yield significant improvement over a steepest descent approach, and the polynomial line search to provide the best combination of performance and computational cost.

In the sample calculations, the conjugate gradient method worked well for orbital localization using the PM and ER measures but did not work well for the FB and FM schemes

due to the importance of tail effects. Convergence problems with the FB cost function have also been previously encountered by Subotnik et al.<sup>38</sup> The use of a second-order method, such as the ones proposed by Leonard and Luken<sup>35</sup> or Høyvik et al.,<sup>40</sup> seems to be necessary in general for a proper optimization of the FB and FM cost functions, or localization of the virtual space.

Because the current method efficiently localizes the occupied space within the PM and ER schemes, it is promising for the implementation of a self-consistent PZ-SIC scheme for which complex orbitals are essential, and for which a second-order scheme would be very costly due to the necessity of three-index ERIs.<sup>35,36</sup>

## APPENDIX

The algorithms were implemented using the Armadillo C++ linear algebra library,<sup>66</sup> the underlying linear algebra operations being performed with version 2.8.0 of the OpenBLAS library.<sup>67</sup> The programs were compiled with GNU G++ 4.4.7 at the -O2 optimization level. The calculations were performed on nodes with 2 Intel Xeon X5650 processors, totalling 12 cores per calculation.

## ASSOCIATED CONTENT

### Supporting Information

The resulting objective function values for Pipek–Mezey localization with Löwdin and Mulliken charges, Foster–Boys localization, fourth moment localization and Edmiston–Ruedenberg localization for fullerenes. The used geometries for all the molecules studied in the current work. This material is available free of charge via the Internet at <http://pubs.acs.org>.

## AUTHOR INFORMATION

### Corresponding Author

\*S. Lehtola: e-mail, [susi.lehtola@alumni.helsinki.fi](mailto:susi.lehtola@alumni.helsinki.fi).

### Notes

The authors declare no competing financial interest.

## ACKNOWLEDGMENTS

S.L. thanks Traian Abrudan for discussions. We gratefully acknowledge computational resources provided by CSC–IT Center for Science Ltd. (Espoo, Finland). This work has been supported by the Academy of Finland through its Centers of Excellence (grant 251748) and FiDiPro program (grant 263294).

## ADDITIONAL NOTES

<sup>a</sup>The objectives differ only by the apparition of the exchange–correlation contribution to the orbital self-interaction energy in PZ-SIC.

<sup>b</sup>Given a starting point  $\mathbf{W}_0$  for the optimization, and allowing for complex rotations  $\exp[-(\boldsymbol{\kappa} + i\boldsymbol{\lambda})]$ , where  $\boldsymbol{\kappa}$  and  $\boldsymbol{\lambda}$

correspond to real and imaginary rotations, respectively, the conditions on the matrices become  $\kappa^T = -\kappa$  for the real part and  $\lambda^T = \lambda$  for the imaginary part.

$\mathbf{H}_k$  is antihermitian, thus having purely imaginary eigenvalues of the form  $i\omega_j$ . The matrix can be decomposed as  $\mathbf{H}_k = \mathbf{U}(i\omega)\mathbf{U}^\dagger$ , where the unitary matrix  $\mathbf{U}$  holds the eigenvectors. The exponential is then obtained as  $\exp(\pm i\mu\mathbf{H}) = \mathbf{U}\exp(\pm i\mu\omega)\mathbf{U}^\dagger$ . The eigenvectors  $\mathbf{U}$  and eigenvalues  $\omega$  can be found by diagonalizing the hermitian matrix  $-i\mathbf{H}$ .

<sup>d</sup>The line minimization only requires knowledge of the largest absolute eigenvalue  $|\omega|$  and the computation of the matrix exponential.<sup>48</sup> Although obtaining the canonical orbitals entails a diagonalization of the  $N \times N$  Fock matrix, the localization of the occupied (virtual) orbitals only needs the  $N_{\text{occ}} \times N_{\text{occ}}$  ( $N_{\text{virt}} \times N_{\text{virt}}$ ) matrix exponential, where  $N = N_{\text{occ}} + N_{\text{virt}}$ . However, the unitary optimization may require orders of magnitude more steps than the solution of the SCF equations.

## REFERENCES

- (1) Hohenberg, P.; Kohn, W. *Phys. Rev.* **1964**, *136*, B864.
- (2) Kohn, W.; Sham, L. J. *Phys. Rev.* **1965**, *140*, A1133.
- (3) Saebo, S.; Pulay, P. *Annu. Rev. Phys. Chem.* **1993**, *44*, 213.
- (4) Hampel, C.; Werner, H.-J. *J. Chem. Phys.* **1996**, *104*, 6286.
- (5) Werner, H.-J.; Pflüger, K. *Annu. Rep. Comput. Chem.* **2006**, *2*, 53.
- (6) Polly, R.; Werner, H.-J.; Manby, F. R.; Knowles, P. J. *Mol. Phys.* **2004**, *102*, 2311.
- (7) Peng, L.; Gu, F. L.; Yang, W. *Phys. Chem. Chem. Phys.* **2013**, *15*, 15518.
- (8) Barr, R.; Basch, H. *Chem. Phys. Lett.* **1975**, *32*, 537.
- (9) Cioslowski, J. *Int. J. Quantum Chem. Symp.* **1990**, *24*, 015.
- (10) Liu, S.; Pérez-Jordá, J. M.; Yang, W. *J. Chem. Phys.* **2000**, *112*, 1634.
- (11) Feng, H.; Jiang, B.; Li, L.; Yang, W. *J. Chem. Phys.* **2004**, *120*, 9458.
- (12) Foster, J.; Boys, S. *Rev. Mod. Phys.* **1960**, *32*, 300.
- (13) Edmiston, C.; Ruedenberg, K. *Rev. Mod. Phys.* **1963**, *35*, 457.
- (14) Pipek, J.; Mezey, P. G. *J. Chem. Phys.* **1989**, *90*, 4916.
- (15) Jansik, B.; Host, S.; Kristensen, K.; Jørgensen, P. *J. Chem. Phys.* **2011**, *134*, 194104.
- (16) Høyvik, I.-M.; Jansik, B.; Jørgensen, P. *J. Chem. Phys.* **2012**, *137*, 224114.
- (17) von Niessen, W. *J. Chem. Phys.* **1972**, *56*, 4290.
- (18) Mulliken, R. S. *J. Chem. Phys.* **1955**, *23*, 1833.
- (19) Høyvik, I.-M.; Jansik, B.; Jørgensen, P. *J. Comput. Chem.* **2013**, *34*, 1456.
- (20) Cioslowski, J. *J. Math. Chem.* **1991**, *8*, 169.
- (21) Alcoba, D. R.; Lain, L.; Torre, A.; Bochicchio, R. C. *J. Comput. Chem.* **2006**, *27*, 596.
- (22) Lehtola, S.; Jónsson, H. Manuscript in preparation.
- (23) Ciupka, J.; Hanrath, M.; Dolg, M. *J. Chem. Phys.* **2011**, *135*, 244101.
- (24) Dubillard, S.; Rota, J.-B.; Saue, T.; Faegri, K. *J. Chem. Phys.* **2006**, *124*, 154307.
- (25) Becke, A. D.; Edgecombe, K. E. *J. Chem. Phys.* **1990**, *92*, 5397.
- (26) Oña, O. B.; Alcoba, D. R.; Tiznado, W.; Torre, A.; Lain, L. *Int. J. Quantum Chem.* **2013**, *113*, 1401.
- (27) Subotnik, J. E.; Dutoi, A. D.; Head-Gordon, M. *J. Chem. Phys.* **2005**, *123*, 114108.
- (28) Almlöf, J. E.; Faegri, K., Jr.; Korsell, K. *J. Comput. Chem.* **1982**, *3*, 385.
- (29) Edmiston, C.; Ruedenberg, K. *J. Chem. Phys.* **1965**, *43*, S97.
- (30) Fletcher, R.; Reeves, C. M. *The Computer Journal* **1964**, *7*, 149.
- (31) Fletcher, R.; Powell, M. J. D. *The Computer Journal* **1963**, *6*, 2.
- (32) Ryback, W.; Poirier, R.; Kari, R. *Int. J. Quantum Chem.* **1978**, *13*, 1.
- (33) Broyden, C. G. *J. Inst. Math. Appl.* **1975**, *6*, 222. Fletcher, R. *The Computer Journal* **1970**, *13*, 317. Goldfarb, D. *Mathematics of Computation* **1970**, *24*, 22. Shanno, D. F. *Mathematics of Computation* **1970**, *24*, 647.
- (34) Kari, R. *Int. J. Quantum Chem.* **1984**, *25*, 321.
- (35) Leonard, J. M.; Luken, W. L. *Theor. Chem. Acc.* **1982**, *62*, 107.
- (36) Leonard, J. M.; Luken, W. L. *Int. J. Quantum Chem.* **1984**, *25*, 355.
- (37) Pulay, P. *Chem. Phys. Lett.* **1980**, *73*, 393.
- (38) Subotnik, J. E.; Shao, Y.; Liang, W.; Head-Gordon, M. *J. Chem. Phys.* **2004**, *121*, 9220.
- (39) Subotnik, J. E.; Sodt, A.; Head-Gordon, M. *Phys. Chem. Chem. Phys.* **2007**, *9*, 5522.
- (40) Høyvik, I.-M.; Jansik, B.; Jørgensen, P. *J. Chem. Theor. Comput.* **2012**, *8*, 3137.
- (41) Absil, P. A.; Mahony, R.; Sepulchre, R. *Optimization algorithms on matrix manifolds*; Princeton University Press: Princeton, NJ, 2008; p 136.
- (42) Høyvik, I.-M.; Jørgensen, P. *J. Chem. Phys.* **2013**, *138*, 204104.
- (43) Klüpfel, S.; Klüpfel, P.; Jónsson, H. *Phys. Rev. A* **2011**, *84*, 050501.
- (44) Klüpfel, S.; Klüpfel, P.; Jónsson, H. *J. Chem. Phys.* **2012**, *137*, 124102.
- (45) Perdew, J. P.; Zunger, A. *Phys. Rev. B* **1981**, *23*, 5048.
- (46) Klüpfel, S.; Klüpfel, P.; Tsemekhman, K.; Jónsson, H. *Lecture Notes in Computer Science* **2012**, *7134*, 23.
- (47) Abrudan, T. E.; Eriksson, J.; Koivunen, V. *IEEE Transactions on Signal Processing* **2008**, *56*, 1134.
- (48) Abrudan, T.; Eriksson, J.; Koivunen, V. *Signal Processing* **2009**, *89*, 1704.
- (49) Manton, J. H. *IEEE Transactions on Signal Processing* **2002**, *50*, 635.
- (50) Polak, E.; Ribière, G. *Revue française d'informatique et de recherche opérationnelle* **1969**, *3*, 35.
- (51) Polyak, B. *USSR Computational Mathematics and Mathematical Physics* **1969**, *9*, 94.
- (52) Lehtola, S. ERKALE – HF/DFT from Hel. 2013; <http://erkale.googlecode.com>.
- (53) Lehtola, J.; Hakala, M.; Sakko, A.; Hämäläinen, K. *J. Comput. Chem.* **2012**, *33*, 1572.
- (54) Moler, C.; Van Loan, C. *SIAM Rev.* **1978**, *20*, 801.
- (55) Ostlund, N. S. *J. Chem. Phys.* **1972**, *57*, 2994.
- (56) Edwards, W. D. *Int. J. Quantum Chem.* **1988**, *34*, 549.
- (57) Obara, S.; Saika, A. *J. Chem. Phys.* **1986**, *84*, 3963.
- (58) Baerends, E. J. *J. Chem. Phys.* **1973**, *2*, 41.
- (59) Sambe, H.; Felton, R. H. *J. Chem. Phys.* **1975**, *62*, 1122.
- (60) Yang, R.; Rendell, A. P.; Frisch, M. J. *J. Chem. Phys.* **2007**, *127*, 074102.
- (61) Dunning, T. H. *J. Chem. Phys.* **1989**, *90*, 1007.
- (62) Kendall, R. A.; Dunning, T. H.; Harrison, R. J. *J. Chem. Phys.* **1992**, *96*, 6796.
- (63) Aquilante, F.; Pedersen, T. B.; Sánchez de Merás, A.; Koch, H. *J. Chem. Phys.* **2006**, *125*, 174101.
- (64) Reed, A. E.; Weinhold, F. *J. Chem. Phys.* **1985**, *83*, 1736.
- (65) Lehtola, S.; Manninen, P.; Hakala, M.; Hämäläinen, K. *J. Chem. Phys.* **2013**, *138*, 044109.
- (66) Sanderson C. Armadillo: An Open Source C++ Linear Algebra Library for Fast Prototyping and Computationally Intensive Experiments; Technical Report; NICTA: St. Lucia, QLD, Australia, 2010
- (67) OpenBLAS, an optimized BLAS library. <http://xianyi.github.com/OpenBLAS>.