# Initial Steps toward Automating the Fitting of DFTB $E_{rep}(r)$[†]

## J. M. Knaup,*,‡ B. Hourahine,§ and Th. Frauenheim‡

*Bremen Center for Computational Materials Science, Universität Bremen, P. O. Box 330440, D-28334 Bremen, Germany, and SUPA, Department of Physics, University of Strathclyde, John Anderson Building, 107 Rottenrow, Glasgow G4 0NG, United Kingdom*

The most time-consuming part of developing new parametrizations for the density functional based tight-binding (DFTB) method consists of producing accurate and transferable repulsive pair potentials. In the conventional approach to repulsive parametrization, every possible diatomic combination of the elements covered by the parametrization must be individually hand-constructed. We present an initial attempt to automate some of this time-consuming process. We consider a simple genetic algorithm-based approach to the fitting problem.

## I. Introduction

The density functional based tight-binding method (DFTB)[1−3] and its later self-consistent charge extension[4] (SCC−DFTB) are computationally very efficient approximations to fully self-consistent Kohn−Sham density functional theory. DFTB has been successfully applied to a wide range of problems in the fields of biomolecules, surfaces, and interfaces, as well as point and extended defects in solid-state systems. For several current examples of the general applicability of DFTB, see additional papers in this section. Depending on the reference systems used during parametrization, LDA, GGA, or hybrid-functional-like results have been obtained for energetics, structures, and vibrational modes.

**DFTB Energies and $E_{rep}(r)$.** The DFTB method is a two-center, minimal basis approximation to the Kohn−Sham problem.[5] For details, see refs 1, 3, and 6. As with empirical tight binding, a crucial component of the total energy is a pairwise repulsive interaction. The total energy of the original model[1] is given by

$$E_{DFTB} = \underbrace{Tr(\mathbf{H}^0[n_0] \cdot \rho)}_{E_{b.s.}} + \frac{1}{2} \sum_{ij}^{M} E_{rep}[n_0^i, n_0^j](|\, r_i - r_j \,|) \quad (1)$$

where the band-structure energy ($E_{b.s.}$) is given by the trace of the DFTB Hamiltonian for the reference system (the set of atoms with charge distributions of $n_0$) and the occupied single particle density matrix ($\rho$). The pairwise repulsive contribution depends on the separation and chemical species of the atoms involved.

The atomic reference systems which provide $n_0$ are chosen to be neutral, spin-unpolarized atoms in a confining potential. The matrix elements are calculated from diatomic pairs,[1,3,6] using a basis of the confined atomic Kohn−Sham wavefunctions in the chosen potential. When constructing $H^0$, DFT potentials for a chosen functional are evaluated either from superposition of the Kohn−Sham atomic potentials[1,6] or as a functional of the superposed densities.[4]

The total energy of the more recent self-consistent charge model (SCC−DFTB)[4] is given by

$$E_{SCC-DFTB} = Tr(\mathbf{H}^0[n_0] \cdot \rho) + \frac{1}{2} \sum_{ij}^{M} E_{rep}[n_0](|\, r_i - r_j \,|) + \underbrace{\frac{1}{2} \sum_{ab} \gamma_{ab}(U_a, U_b, r_{ab}) \Delta q_a \Delta q_b}_{E_{SCC}} \quad (2)$$

where $\gamma$ is a Coulombic-like interaction between sites (depending on the atomic Hubbard-$U$ parameters), and $\Delta q$ is the fluctuation of Mulliken charges compared to the reference system ($n_0$). $\gamma_{a,b}$ interpolates between two exact known limiting cases (correctly predicting the atomic chemical hardness as $r_{ab} \to 0$ and the Coulombic interaction as $r_{ab} \to \infty$). The SCC contributions are resolved either by atom or by individual l-shells of the atoms (thus both forms are rotationally invariant).

Further extensions have been built on top of these models. Due to the choice of neutral, unpolarized atoms as a reference, extensions such as spin polarization,[7,8] dispersion,[9,10] or orbitally dependent correlation[11] are additive in eq 2; hence, a previously parametrized DFTB or SCC−DFTB $E_{rep}(r)$ can also be used in these applications.

## II. Fitting $E_{rep}(r)$

Since each pairwise chemical combination requires an accurate $E_{rep}(r)$, a consistent set of $O((N^2 + N)/2)$ repulsive interactions needs to be constructed for the collection of atomic types present in the desired application. Once a highly transferable $E_{rep}(r)$ curve has been constructed, the DFTB method is very efficient. Much effort has been expended in adding new element combinations to existing sets of consistent repulsive sets.

$E_{rep}(r)$ is defined as the difference between $E_{DFT}$ and $E_{b.s.}$ (or $E_{b.s.} + E_{SCC}$ in the case of SCC−DFTB). From a purist point of view, the same functional should be used to construct both $H^0$ and $E_{rep}(r)$; however, within the DFTB approximations, the DFTB band structure energy is not strongly dependent on the functional used to generate the $H^0$ Hamiltonian, but the choice of functional for the repulsive reference has a stronger effect. Under this definition of $E_{rep}(r)$ as a difference to a DFT reference, $E_{rep}(r)$ is only guaranteed to be purely repulsive for systems consisting of only simple dimers. For more complicated

---
† Part of the "DFTB Special Section".
* Electronic address: Jan.Knaup@bccms.uni-bremen.de.
‡ Universität Bremen.
§ University of Strathclyde.

cases, such as bonds in larger molecules or extended solids, the difference definition allows $dE_{rep}(r)/dr$ to be positive in some regions.

**A. The Conventional Approach.** The manual fitting process typically proceeds by identifying a series of structures which possess examples of the chemical bonds which should be reproduced by the parameter set. The bond lengths are then systematically varied and the DFT total and DFTB band-structure energies are calculated. For example, the DFTB organic carbon set[4] used the single, double, and triple bonds in $C_2H_6$, $C_2H_4$, and $C_2H_2$, stretching each molecule over a range of C–C-distances, keeping the C–H bonds fixed, and fitting a smooth and short-ranged $E_{rep}(r)$ from these three piecewise sections. In this case, the electronic part of the energy was constructed with the Perdew, Burke, and Enzerhoff (PBE)[12,13] functional, while the $E_{rep}(r)$ was produced from a B3LYP[14,15] DFT reference. The $E_{rep}(r)$ curve is shifted so that at the cutoff distance it goes to 0 (ideally without discontinuity).

This process has also been performed with the intention of reproducing vibrational modes instead of energetics, by using an analogous method where the vibrational modes associated with chosen bonds as a function of length are instead fitted[16-18] with $E_{rep}(r)$ constructed by integrating the constructed derivatives.

Typically, homonuclear interactions are fitted first, then fixed while the heteronuclear cases are constructed.

**B. A First Attempt at Automation.** Reducing the amount of human intervention required in constructing $E_{rep}(r)$ is highly desirable. To achieve transferable results for a new system, perhaps up to 1 month of human time can be required to construct and verify a given pairwise interaction, while the computational cost of the parametrizing calculations, with current computers, is negligible by comparison with current computers. The human input into parametrization is rapidly becoming the time-determining point of calculations for a new material.

In some sense, the process of producing $E_{rep}(r)$ is similar to the "Learn on the Fly" method of Csanyi et al.,[19] where a series of environmentally dependent interatomic potentials is derived. This is typically done by constructing additional contributions to previously supplied potentials by fitting differences in forces when compared with quantum chemical calculations. In this case, in principle each pair of atoms requires a different correcting potential which may change during the calculation. In this, we are fortunate, since by construction $E_{rep}(r)$ is only required for each atomically distinct pair rather than for each chemical environment. We also have several further constraints on the form of $E_{rep}(r)$ that are present by definition, since $E_{rep}(r)$ should be continuous (ideally up to at least the third derivative), short-ranged, and often chosen to be monotonically decreasing. In this case, its first derivative would always be less than zero. Additionally, $E_{rep}(r)$ should be 0 at the end of the $E_{rep}(r)$ table.

In an attempt to remove some of the human effort in constructing $E_{rep}(r)$, we have first tried to generalize the fitting process such that automated comparisons between the DFT energy and the DFTB band structure are possible. Instead of restricting the comparisons to a set of example bonds, we instead define a path of distortions for a set of structures similar to (but with much smaller numbers of atoms than) the target system for which $E_{rep}(r)$ is needed. This path can be quite general; for example, it could come from a series of molecular dynamics or Monte Carlo steps, barrier crossing events, or simply structural relaxations. In addition, multiple paths may be described to capture other examples of bonding or other processes; for
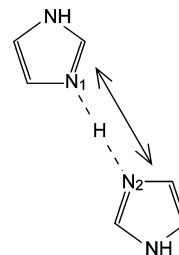


**Figure 1.** Schematic representation of the proton-transfer process. The dashed line marks the proton-transfer path.
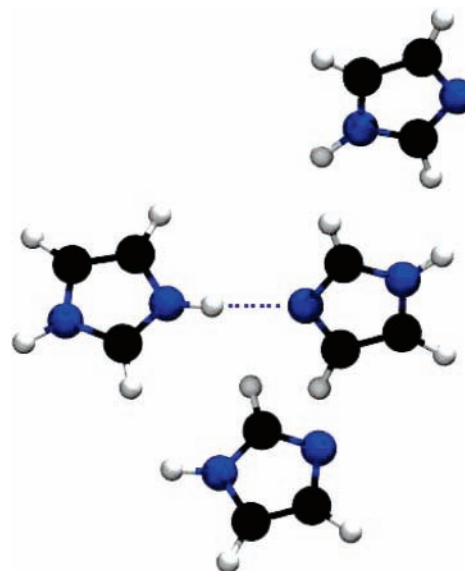


**Figure 2.** The proton-transfer path, marked by the dotted line, for a $N_1$–$N_2$ separation of 3.0 Å, including the two additional imidazole rings to simulate crystalline surroundings (at the top and bottom). C, black; N, blue; H, white.

example, including simple dimers or deformations of the unit cell for elemental solids makes sense. Over the course of exploring the paths, desired properties such as total energy, forces, and/or vibrational modes can then be monitored at the DFT and DFTB levels of theory. Fitting $E_{rep}(r)$ then becomes a general optimization process of minimizing the error in the DFTB properties, calculated to including the trial $E_{rep}(r)$, when compared against the DFT reference. Simultaneously, we must constrain the properties of $E_{rep}(r)$ to be short-ranged and continuous.

There are a wide range of techniques for such optimizations; one could consider least-squares fitting, for example or, as we do in this work, genetic algorithms.

## III. Genetic Optimization of $E_{rep}(r)$

Very good introductions into the field of genetic algorithms, GAs, can be found in refs 20–22. Here, we employ a simple scheme of genetic optimization, which is by no means fully developed but rather serves to test whether GA are generally useful for $E_{rep}(r)$ fitting. To do so, we represent the $E_{rep}(r)$ as a series of $E(r)$ points between which we interpolate using a natural cubic spline.[28] Since the first and second derivatives of natural cubic splines are continuous by definition, this representation automatically meets the continuity requirements on $E_{rep}(r)$.

Our GA recombination operation simply cuts two parent $E_{rep}(r)$ at the same randomly chosen data point and exchanges the sections. Since the same parent can be chosen twice for

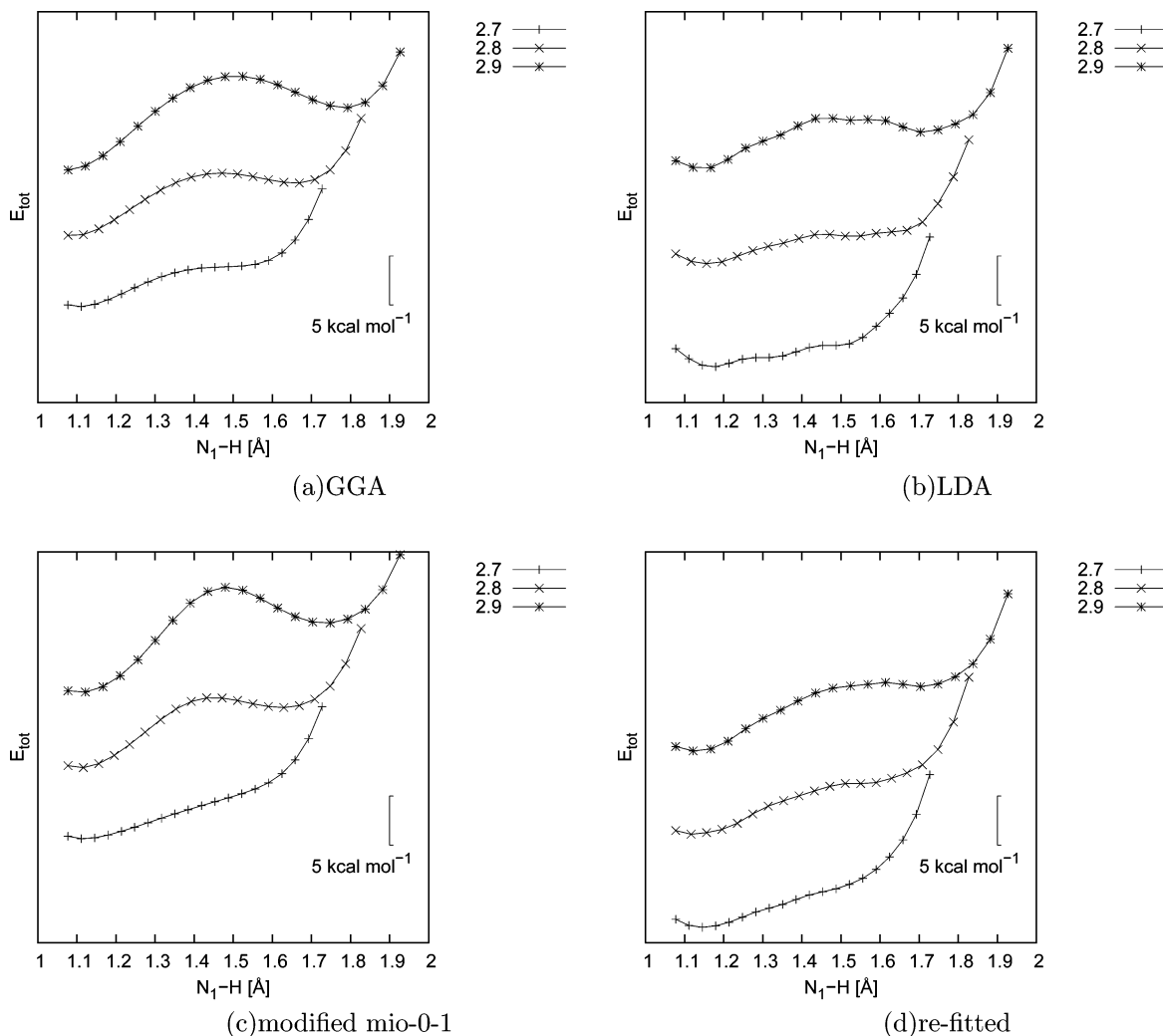Automated Fitting of DFTB $E_{rep}(r)$

*J. Phys. Chem. A, Vol. 111, No. 26, 2007* **5639**



(a)GGA



(b)LDA



(c)modified mio-0-1



(d)re-fitted

**Figure 3.** Proton-transfer energy profiles along the $N_1-N_2$ axis for different $N_1-N_2$ distances, calculated using GGA, LDA, and two versions of the $N-H$ $E_{rep}(r)$. Curves are shifted along the $y$-axis for clarity; a 5 kcal mol$^{-1}$ interval is indicated in each plot for reference.

this operation, then several copies of the more successful individuals can pass to the mutation stage. We define the mutation scheme to change the data point at index $c$, with the scaling factor $s$ to limit the amount of mutation and a random number $p = \text{rnd}(-1, 1)$ in the range $[-1, ..., 1]$ as follows:

$$E_{rep,c}^{new} = E_{rep,c}^{old} + p \cdot s \cdot (E_{rep,c-1}^{old} - E_{rep,c}^{old}) \quad \text{if } p < 0 \quad (3)$$

$$E_{rep,c}^{new} = E_{rep,c}^{old} + p \cdot s \cdot (E_{rep,c}^{old} - E_{rep,c+1}^{old}) \quad \text{if } p \geq 0 \quad (4)$$

This shifts $E_{rep,c}^{old}$ between the neighboring data points centered around $E_{rep,c}^{old}$, thus providing scaling of the mutation. Since, in this work, we are trying to refine an existing $E_{rep}(r)$, we set $s$ to 0.1, to stabilize the system. We also append two additional data points with zero $E_{rep}(r)$, spaced at 0.1 au at the end of the data set, and keep these fixed, to ensure that $E_{rep}(r)$ and its first two derivatives vanish at the desired cutoff range. Additionally, we keep the first data point fixed.

In this work, we define the fitness $F$ of $E_{rep}(r)$ to be the sum of squares of force differences between the sampled points of our PES

$$F = \sum_{i}^{\text{all atoms}} (F_i^{E_{rep}} + F_i^{DFTB} - F_i^{ref})^2 \quad (5)$$

where the superscripts denote the forces calculated from the repulsive potential, DFTB electronics, and the reference method, respectively.

We generated the initial population of 20 $N-H$ $E_{rep}(r)$'s by taking an initial $E_{rep}(r)$ from a modified mio-0−1 DFTB parameter set[4,23] and creating 19 mutant versions.

We then start the iterative optimization scheme in which we first calculate the fitness for each $E_{rep}(r)$ and sort the $E_{rep}(r)$'s by fitness. Subsequently, we eliminate all but the five fittest $E_{rep}(r)$'s and refill the population by combining random pairs from the surviving $E_{rep}(r)$'s and mutating each child (but leaving the five parents unchanged).

We currently do not define an automatic convergence criterion, but periodically analyze the PES by hand.

## IV. Example: The N−H $E_{rep}(r)$

At first glance, the interaction of an element with hydrogen appears to be comparatively easy to parametrize, even with the traditional approach, as hydrogen is single-valent and thus no combination of different bond types must be regarded. Yet, the description of proton-transfer processes, especially their barriers, depends delicately upon a correct $E_{rep}(r)$.

One example for ths dependency is the onset of the proton-transfer barrier in an imidazole crystal, dependent on the distance between imidazole rings, as shown in Figures 1 and 2. Here, a

proton-transfer barrier appears for a $N_1-N_2$ distance greater than ~2.8 Å. The separation at which this barrier is present is very important for molecular dynamics simulations to determine the diffusion properties of protons in imidazole crystals.

When using density functional theory (DFT), the calculated onset of the proton-transfer barrier depends on the level of theory which is employed. Using the local density approximation (LDA) functional by Perdew and Zunger,[24] with the parametrization by Ceperley and Alder,[25] the barrier appears at a $N_1-N_2$ distance of ~2.8–2.9 Å, while the PBE[12,13] gives an onset of ~2.7 Å (cf. Figure 3).

To compare the descriptions of proton transfer between two imidazole molecules in a crystal-like setting using the different methods and parameter sets, we calculate the potential energy surface (PES) using the $N_1-N_2$ distance as the first axis and the $N_1-H$ distance as the second axis. We vary the $N_1-N_2$ distance between 2.4 and 3.0 Å.

For LDA, we calculate a reference PES using SIESTA,[26] choosing a double-$\zeta$ basis with polarization functions (dzp), determining the basis cutoff by a shift of $2 \times 10^{-3}$ Ry. We set the cutoff for the auxiliary $k$-space grid to 250 Ry. We employ Troullier-Martins pseudopotentials.[27] As can be seen in Figure 3B, the PES is not completely smooth. This is probably due to the still rather small basis set. Since a larger basis set would have to be generated by hand, we chose to use the largest automatically generated set available, as the exact form of the PES is less relevant for this study of DFTB parametrization.

For PBE, we calculate the PES using an all-electron calculation in a 6-31G* basis set. The resulting DFT PES's are shown in Figure 3a,b.

The N−H interaction from the reference DFTB parameter set[4,23] reproduces the GGA results very well, as can be seen in Figure 3c. The question arises regarding whether the details of the PES are determined mostly by the electronic part or the $E_{rep}(r)$ in the parameter set.

To determine this, we start from the existing N−H $E_{rep}(r)$ from the modified mio-0−1 set and refine the repulsive potential using a genetic algorithm, using a LDA-derived PES as a target. (It should be noted that the GGA results describe the proton transfer more accurately than the LDA ones. Thus, modifying the parameter sets to reproduce LDA rather than GGA provides no practical improvement.)

After about 365 000 mutations of the $E_{rep}(r)$ (~2.4% of which led to an improvement of the maximum fitness) performed in a little less than 72 h on a single PC workstation, we find that the DFTB-PES resembles the LDA results closely (cf. Figure 3d). Figure 3 shows that the details of the proton-transfer PES between imidazole molecules in a crystalline-like environment can be fitted to different levels of DFT theory, modifying the $E_{rep}(r)$ only. The rather low rate of improving mutations indicates that either the scale of mutation is too large or our mutation/recombination scheme tends to easily generate individual unfavorable $E_{rep}(r)$'s during the optimization.

**LDA and B3LYP Repulsives.** As shown in Figure 4, both the original B3LYP and the refitted LDA repulsives decay over the length scale of the N−H transfer process. The newly fitted $E_{rep}(r)$ is close to a monoexponential, while the B3LYP shows a much flatter (very slightly attractive) region between 1.2 and 1.5 Å. The deviation from exponential decay for this $E_{rep}(r)$ is strongest at 1.3 Å. On first inspection, this does not appear to correspond to a feature in the GGA/DFTB reaction barrier (Figure 3), when plotted in the $N_1-H$ reaction coordinate. However, if the $N_2-H$ separation is instead considered, this matches the peak of the barrier position for the 2.8 and 2.9 Å
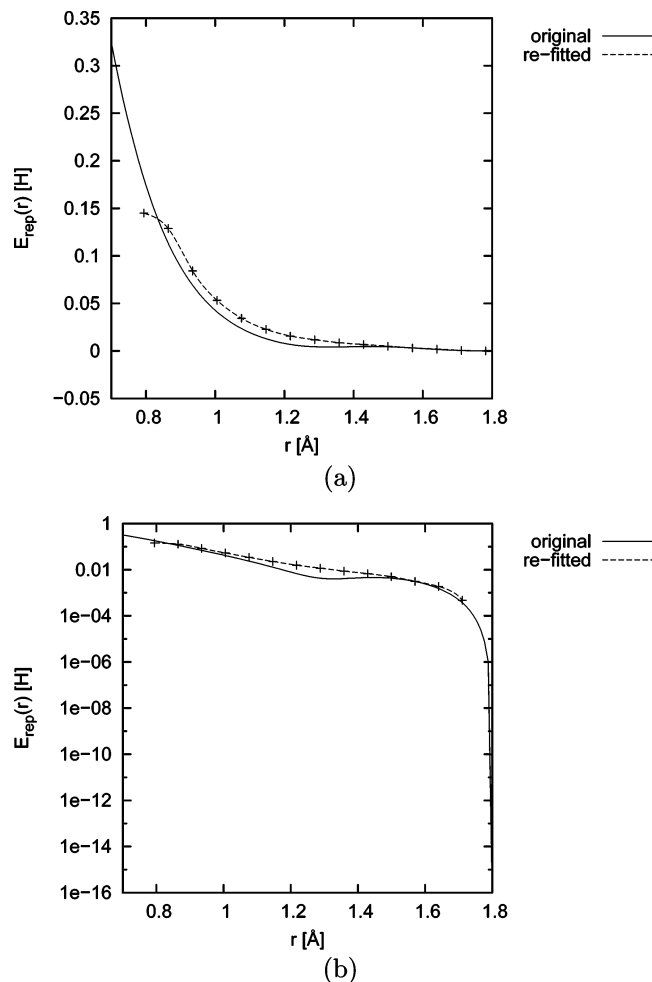

(a)


(b)

**Figure 4.** Comparison of $E_{rep}(r)$ before and after the genetic refitting process in (a) linear and (b) logarithmic energy scales. Symbols in the refitted potential mark the data points used in the fitting process; curves are spline-interpolated in the same manner as in the DFTB+ implementation.

$N_1-N_2$ separation closely (to within 0.1 Å). This suggests that, in order to reproduce a GGA-like surface, the requirement of monotonicity must be lifted for the $E_{rep}(r)$. However, since the new $E_{rep}(r)$ reproduces the LDA PES well, in this case a monotonic, near exponentially decaying, $E_{rep}(r)$ suffices.

## V. Summary

From our experience with proton-transfer paths in imidazole crystals, we can conclude that genetic algorithms are well-suited to perform at least the fine fitting of $E_{rep}(r)$. This allows for a significant simplification of the $E_{rep}(r)$ fitting process, even for new materials: Starting from a preliminary $E_{rep}(r)$, which could be derived from just one simple run using target paths similar to the traditional fitting process, the repulsive potential can then be simultaneously optimized for a number of fit systems. In the case of N−H, we have shown that even fine details of the potential energy surfaces of heteronuclear interactions can be reproduced by fitting the $E_{rep}(r)$ via a genetic algorithm. To achieve a semiautomatic generation scheme for DFTB repulsive potentials, two further steps now will have to be developed.

1. A strategy as to which general target properties are the most crucial to reproduce, e.g., forces, reaction barriers, atomization energies, bulk moduli, and so forth.

2. Improvements in the genetic fitting itself. In this work, we use a very simple algorithm in a rather crude implementation

Automated Fitting of DFTB $E_{rep}(r)$

*J. Phys. Chem. A, Vol. 111, No. 26, 2007* **5641**

which has by no means been tuned for efficiency. It can be expected that taking advantage of the vast experience, which the bioinformatics community has accumulated over the past decades of research on genetic and Monte Carlo algorithms, will further improve the quality $E_{rep}(r)$ fits can attain as well as the efficiency of the fitting process.

This new fitting procedure will not only facilitate the generation of parameter sets for new elements or interactions, but it also allows for the generation of specifically tuned repulsive potentials for use in applications where different levels of theory are applied in a tiered manner; e.g., when using DFTB calculations to generate geometries for computationally expensive DFT or even higher-level calculations, it would be possible to generate parameter sets which reproduce the equilibrium configurations of the higher-level method even better than DFTB already manages to do, thus further improving the precision of such multilevel approaches.

## References and Notes

(1) Eschrig, H.; Seifert, G. *Zeitschrift fur Physikalische Chemie Neue Folge* **1986**, *267*, 529.

(2) Porezag, D.; Frauenheim, T.; Köhler, T.; Seifert, G.; Kaschner, R. *Phys. Rev. B* **1995**, *51*, 12947−12957.

(3) Seifert, G.; Porezag, D.; Frauenheim, T. *Int. J. Quantum Chem.* **1996**, *58*, 185−192.

(4) Elstner, M.; Porezag, D.; Jungnickel, G.; Elsner, J.; Haugk, M.; Frauenheim, T.; Shuhai, S.; Seifert, G. *Phys. Rev. B* **1998**, *58*, 7260−7268.

(5) Kohn, W.; Sham, L. J. *Phys. Rev.* **1965**, *140*, A1133−A1138.

(6) Seifert, G. *J. Phys. Chem. A* **2007**, *111*, 5609−5613.

(7) Köhler, C.; Seifert, G.; Frauenheim, T. *Chem. Phys.* **2005**, *309*, 23.

(8) Köhler, C. *J. Phys. Chem. A* **2007**, *111*, 5622−5629.

(9) Elstner, M.; Hobza, P.; Frauenheim, T.; Suhai, S.; Kaxiras, E. *J. Chem. Phys.* **2001**, *114*, 5149.

(10) Zhechkov, L.; Heine, T.; Patchkovskii, S.; Seifert, G.; Duarte, H. A. *J. Chem. Theor. Comput.* **2005**, *1*, 841−847.

(11) Sanna, S.; Hourahine, B.; Gallauner, Th.; Frauenheim, Th. *J. Phys. Chem. A* **2007**, *111*, 5665−5670.

(12) Perdew, J. P.; Burke, K.; Ernzerhof, M. *Phys. Rev. Lett.* **1997**, *78*, 1396.

(13) Perdew, J. P.; Burke, K.; Ernzerhof, M. *Phys. Rev. Lett.* **1996**, *77*, 3865.

(14) Becke, A. D. *J. Chem. Phys.* **1993**, *98*, 5648.

(15) Lee, C.; Yang, W.; Parr, R. G. *Phys. Rev. B* **1988**, *37*, 785.

(16) Witek, H. A.; Irle, S.; Morokuma, K. *J. Chem. Phys.* **2004**, *121*, 5163−5170.

(17) Małolepsza, E.; Witek, H. A.; Morokuma, K. *Chem. Phys. Lett.* **2005**, *412*, 237.

(18) Zheng, G.; Irle, S.; Frisch, M.; Morokuma, K. *J. Phys. Chem. A*, submitted.

(19) Csanyi, G.; Albaret, T.; Payne, M. C.; De Vita, A. *Phys. Rev. Lett.* **2004**, *93*, 175503.

(20) Goldberg, D. E. *Genetic Algorithms in Search, Optimization and Machine Learning*; Addison-Wesley Publishing Company, Inc.: Reading, MA, 1989.

(21) Fogel, D. B. *Evolutionay Computation: toward a new philosophy of machine intelligence*; IEEE Press: New York, 1995.

(22) Bäck, T. *Evolutionary Algorithms in Theory and Practice*; Oxford University Press: New York, 1996.

(23) Elstner, M. Optimized version of the mio-0-1 N−H parameter.

(24) Perdew, J. P.; Zunger, A. *Phys. Rev. B* **1981**, *23*, 5048−5079.

(25) Ceperley, D. M.; Alder, B. J. *Phys. Rev. Lett.* **1980**, *45*, 566−569.

(26) Soler, J. M.; Artacho, E.; Gale, J. D.; García, A.; Junquera, J.; Ordejón, P.; Sanchéz-Portal, D. *J. Phys.: Condens. Matter* **2002**, *14*, 2745.

(27) Troullier, N.; Martins, J. L. *Phys. Rev. B* **1991**, *43*, 1996.

(28) The current implementations of DFTB also represent the $E_{rep}(r)$ as c-splines but store the spline coefficients rather than data points. To convert between the two formats, we use the same tool that is used to generate the spline coefficients for the $E_{rep}(r)$'s fitted in the traditional way.