

# Density Functional Theory Augmented with an Empirical Dispersion Term. Interaction Energies and Geometries of 80 Noncovalent Complexes Compared with *Ab Initio* Quantum Mechanics Calculations\*

PETR JUREČKA,<sup>1,2\*</sup> JIŘÍ ČERNÝ,<sup>2</sup> PAVEL HOBZA,<sup>2</sup> DENNIS R. SALAHUB<sup>1</sup>

<sup>1</sup>Department of Chemistry, University of Calgary, 2500 University Drive NW Calgary, Alberta, Canada T2N 1N4

<sup>2</sup>Institute of Organic Chemistry and Biochemistry, Academy of Sciences of the Czech Republic and Center for Biomolecules and Complex Molecular Systems, Flemingovo nám. 2, 166 10 Prague, Czech Republic

Received 27 February 2006; Revised 5 May 2006; Accepted 6 May 2006

DOI 10.1002/jcc.20570

Published online 21 December 2006 in Wiley InterScience (www.interscience.wiley.com).

**Abstract:** Standard density functional theory (DFT) is augmented with a damped empirical dispersion term. The damping function is optimized on a small, well balanced set of 22 van der Waals (vdW) complexes and verified on a validation set of 58 vdW complexes. Both sets contain biologically relevant molecules such as nucleic acid bases. Results are in remarkable agreement with reference high-level wave function data based on the CCSD(T) method. The geometries obtained by full gradient optimization are in very good agreement with the best available theoretical reference. In terms of the standard deviation and average errors, results including the empirical dispersion term are clearly superior to all pure density functionals investigated—B-LYP, B3-LYP, PBE, TPSS, TPSSh, and BH-LYP—and even surpass the MP2/cc-pVTZ method. The combination of empirical dispersion with the TPSS functional performs remarkably well. The most critical part of the empirical dispersion approach is the damping function. The damping parameters should be optimized for each density functional/basis set combination separately. To keep the method simple, we optimized mainly a single factor,  $s_R$ , scaling globally the vdW radii. For good results, a basis set of at least triple- $\zeta$  quality is required and diffuse functions are recommended, since the basis set superposition error seriously deteriorates the results. On average, the dispersion contribution to the interaction energy missing in the DFT functionals examined here is about 15 and 100% for the hydrogen-bonded and stacked complexes considered, respectively.

© 2006 Wiley Periodicals, Inc. J Comput Chem 28: 555–569, 2007

**Key words:** dispersion interaction; density functional theory; empirical corrections; van der Waals complexes

## Introduction

The dispersion force, also called the London force, is an attractive interaction arising from the electrostatic interaction of the fluctuating charge distributions—instantaneous dipoles and higher multipoles. Regarding its magnitude, dispersion stabilization in biomolecules is usually (but not always) weaker than the stabilization from hydrogen bonding—for illustration, the dispersion contribution to the intrinsic stability of the DNA double helix is about 30 and 50% in the CG and AT rich sequences, respectively.<sup>1</sup> Dispersion energy also stabilizes the DNA–intercalator complex<sup>2</sup> and it contributes to the surprising stability of DNA modified by the nonpolar base analogues.<sup>3</sup> In proteins, a substantial dispersion contribution can be expected whenever

\*Present address: Department of Physical Chemistry, Palacký University, Tr. Svobody 26, 771 46 Olomouc, Czech Republic

**Correspondence to:** P. Jurečka; e-mail: jurecka@marge.uochb.cas.cz

Contract/grant sponsor: NSERC of Canada; contract/grant number: RGP-262037

Contract/grant sponsor: Grant Agency of the Czech Republic; contract/grant number: 203/05/0009

Contract/grant sponsor: Grant Agency of the Academy of Sciences of the Czech Republic; contract/grant number: LC 512

Contract/grant sponsor: MŠMT of the Czech Republic; contract/grant number: A400550510

Contract/grant sponsor: Ministry of Education, Youth and Sports, Czech Republic; contract/grant number: MSM 6198959216

\*This article contains supplementary material available via the Internet at <http://www.interscience.wiley.com/jpages/0192-8651/suppmat>.

nonpolar, especially, aromatic residues are present, but also the peptide bond itself contributes remarkably.<sup>4</sup> If dispersion is not included in molecular dynamics (MD) simulations of proteins, divergence from the X-ray geometry may occur, however, this may be overlooked on a short time scale.<sup>5</sup> Neglect of dispersion can result in wrong binding energies for the ligands<sup>6–8</sup> and in wrong geometries.<sup>9</sup> Also, dispersion is supposed to be an important driving force for protein folding<sup>10</sup> and a main stabilization factor of amyloid fibrils.<sup>11</sup> Therefore, an accurate description of the dispersion forces is vital if reliable simulations of DNA, lipids, peptides, and proteins are to be obtained.

The present paper addresses the question of the intermolecular dispersion interaction in density functional theory (DFT). We will focus on the intermolecular nonlocal correlation contribution at distances close to and longer than typical vdW distances. It is not completely accurate to call these forces dispersion forces, because the dispersion interaction is strictly defined in the asymptotic region only, and the terminology becomes ambiguous when molecular densities overlap, especially at the equilibrium distances where the overlap may be substantial. Nevertheless, for the sake of simplicity (and because we have no better expression) we will use the term “dispersion” for these forces a few times in the following text, keeping the subtlety of its definition in mind.

Although the number of DFT-based calculations and simulations on biologically relevant molecules exceeds the number of wave function theory (WFT)-based applications, the role of the dispersion interaction in DFT is rarely discussed. In the WFT-correlated treatments dispersion is naturally included in terms of the simultaneous electron excitations into distant virtual orbitals. On the other hand, current density functionals are based on the local density expansion and are therefore inherently incapable of describing long-range nonlocal interactions such as dispersion or long-range exchange repulsion. The fact that dispersion energy is missing in the local DFT energy was not easily accepted by the DFT community and the first papers showing the drawbacks appeared in the nineties of the last century.<sup>12,13</sup> This applies to all LDA, and GGA functionals, and also to the most advanced meta-GGA functionals (for the long-range exchange repulsion it was recently shown by Boese and Handy<sup>14</sup>). Despite this, DFT results for hydrogen bonds are often in fair agreement with reference data, both experimental and WFT based.<sup>15,16</sup> However, it seems that good performance in these cases stems from the error cancellation. The missing dispersive attraction is in some density functionals partly simulated by the (erroneous) attraction of the exchange functional.<sup>17</sup> A large portion of the missing dispersion attraction may be compensated by the basis set superposition error (BSSE), which amounts to several kcal mol<sup>−1</sup> for the medium-sized dimers and double- $\zeta$  quality basis sets (see later). Although these factors are able to produce enough additional unphysical attraction for the H-bonded complexes, they are not large enough in the case of the dispersion-bonded complexes, which results in unacceptable errors. It should be noted that this applies also to single molecules, if their dimensions are comparable to or larger than sums of the vdW radii of constituent atoms, especially if highly polarizable groups like phenyl rings or peptide bonds are present. This renders DFT-based methods highly unreliable, since they often fail to reproduce proper geometries of molecules with nonpolar moieties, like peptides.<sup>5,18</sup>

Numerous attempts to include the dispersion interaction into DFT can be divided into three groups: (i) theoretically based approaches with as little empiricism as possible—nonempirical approaches. (ii) Attempts to reparametrize the existing density functionals so that they describe dispersion properly. Note that there are very few dispersion-bonded complexes in widely used training sets, therefore good performance of the current functionals can not be anticipated. (iii) Empirical approaches, based on the force-field-like terms. Dispersion energy is simply added to the total DFT energy.

Nonempirical ways to include dispersion into DFT will be mentioned just briefly. Usually, the molecular complex is divided into two subsystems and the dispersion interaction between them is calculated either from intermolecular perturbation theory,<sup>19,20</sup> from the dynamic polarizabilities,<sup>21,22</sup> from the ground state densities only,<sup>23–26</sup> or in some other way.<sup>27–29</sup> Note that division of the complex into subsystems is in fact one of the major limitations of these methods, since *intramolecular* dispersion in larger molecules (peptide chains, proteins, nucleic acids) can not be accounted for. Methods relying on molecular fragmentation as a principle make an exception and a promising solution has been presented recently.<sup>22</sup> Density functionals not requiring splitting of the system (seamless methods) were suggested by Kohn et al.<sup>30</sup> and Dion et al.<sup>31</sup> While the density functional of Kohn et al. is prohibitively computationally demanding, the method of Dion et al. is only moderately expensive and may appear useful after the necessary testing is done. However, evaluation of the dispersion energy by these methods is either very tedious or it is based on simplifications leading to some inaccuracies, typically over 10% and sometimes over 20%. Also, in most cases relatively large basis sets are required to describe the dynamic polarizabilities of the molecules correctly. Therefore, the errors in computed dispersion energies are usually not much smaller than the errors which appear in the empirical dispersion calculations due to approximate dispersion coefficients and damping. Construction of a computationally inexpensive density functional accounting for dispersion in a general and seamless way thus remains a great challenge.

Regarding attempts to modify the existing density functionals, the mainstream is represented by modification of the exchange functional.<sup>32–35</sup> The exchange functionals certainly need to be improved since the long-range exchange which co-determines the vdW bonding is often predicted only poorly.<sup>33</sup> However, the dispersion interaction should be described rather by the correlation functional, since it is purely a correlation effect. Also, while dispersion requires the known  $1/r^6$ ,  $1/r^8$ , ... , behavior of the functional, the exchange functional should exhibit the  $1/r$  behavior. Nevertheless, these contradictions do not preclude that some partial success be achieved also by the exchange functional modification. In the work of Xu and Goddard,<sup>34</sup> the exchange functional was optimized considering among others the rare gas dimers; the results on the real systems however turned out to be poor.<sup>9</sup> Zhao and Truhlar<sup>36</sup> have tested their kinetics-optimized functionals on several larger dispersion-bonded complexes with reasonable success; these studies, however, suffer from small basis sets and more thorough testing would be desirable. Kurita et al.'s modified PW91 functional<sup>35</sup> performs relatively well for the stacked cytosine dimers, however, at the expense of the

hydrogen-bonded complexes. This feature, i.e., that certain functionals work either for the hydrogen-bonded or for the stacked complexes, but not for both, seems to be shared by all the approaches based on modifying the exchange functional.

In the empirical dispersion approach, the dispersion energy is represented by the well known  $C_6/R^6$  formula (see Methods).  $E_{\text{Disp}}$  is calculated separately from the DFT calculation, damped by an appropriate distance-dependent damping function to correct for the overlap effects, and simply added to the DFT energy. It is assumed that since dispersion acts mostly on the long-range, it has only a small effect on the electron densities, and its separate (uncoupled) calculation should not be problematic. Thus, one relies on the dispersion influencing the chemistry of the system mainly through the geometries, which is a great simplification. This may not be true in some specific situations such as dipole bound states<sup>37</sup> or excited states, which are, however, beyond the scope of this study.

Among the first, Scoles and coworkers<sup>38,39</sup> combined empirical dispersion, damped by a sophisticated many-parameter damping function, with Hartree–Fock calculations. This method was fairly successful since it dealt with particularly weak rare gas atom complexes, where the intermolecular repulsion is fairly well described by the HF nonlocal treatment. A similar approach was successfully used later by Hobza et al. for the description of larger complexes including DNA base pairs:  $(\text{H}_2\text{O})_2$ ,  $\text{H}_2\text{O} \cdots \text{HCX}_3$  ( $\text{X} = \text{F}, \text{Cl}$ ),<sup>40</sup> formamide dimer, formamide  $\cdots \text{HCX}_3$  ( $\text{X} = \text{F}, \text{Cl}$ ),<sup>41</sup>  $(\text{H}_2\text{O})_2$ ,  $\text{H}_2\text{O} \cdots \text{CH}_4$ ,  $\text{CH}_3\text{Cl}$ ,  $\text{CH}_2\text{Cl}_2$ ,  $\text{CCl}_4$ ,<sup>42</sup> and H-bonded DNA pairs.<sup>43</sup> Since then, this initial success energized several attempts to apply this idea to DFT. Works of Meijer and Sprik,<sup>44</sup> Mooij et al.,<sup>45</sup> and Wu et al.<sup>46</sup> showed significant improvements to pure DFT but also pointed to several questions regarding the reliability of DFT for the long-range exchange repulsion. Wu and Yang<sup>47</sup> and Zimmerli et al.<sup>48</sup> touched the question of the damping function. Unfortunately, all these studies share a drawback of only one or very few molecules tested, which casts some doubt on the transferability of the results. Significant progress in this direction was made in a study by Grimme<sup>49</sup> in which various density functionals were tested on a set of molecules. The basis set dependence was examined and basis sets of at least triple- $\zeta$  quality were recommended. Notably, Grimme recognized the need for the dispersion to be adjusted for a given functional form and introduced a simple scaling factor optimized for each particular density functional. An empirical dispersion term was successfully applied also in connection with the approximate SCC-DFTB method by Elstner et al.<sup>50</sup> and Zhechkov et al.<sup>51</sup> Among other works we mention studies of Lilienfeld et al.<sup>52</sup> in which dispersion is simulated by reparametrization of the nonlocal part of the pseudopotential in the plane wave DFT.

The aim of this work is to suggest a computationally efficient empirical model with good transferability and reasonably small error. We assume that this model is used mainly for the DFT-based calculations and simulations of peptides and nucleic acids. Such large scale calculations are usually performed with a non-hybrid XC functional because the density fitting approximation can be exploited. The density fitting approximation brings significant speed-ups but can not be applied to the hybrid functionals which contain explicit HF exchange. Unfortunately, this sacrificed explicit Hartree–Fock exchange seems to be crucial for the accuracy of the intermolecular interactions.<sup>29,53</sup> Therefore, we

will not focus on the accuracy of the dispersion in the first place, but rather on the transferability and simplicity of the model.

Results will be judged on the basis of standard deviation and average (signed) errors. The signed average carries information about average over- or underestimation of a given interaction by a given method and can be controlled by the global dispersion scaling parameter  $s_R$ . The standard deviation indicates errors in the XC functional, variable double counting of dispersion and errors due to the averaged dispersion coefficients. As a reference, we use a training set of 22 accurate intermolecular interaction energies and geometries and a validation set of 58 molecules adopted from our previous papers.<sup>54–58</sup> All geometries and interaction energies used in this paper are available in ref. 54.

## Methods

In the atomic dispersion scheme, the total dispersion energy is calculated as a sum of all possible pairwise atomic contributions [eq. (1)].

$$E_{\text{dis}} = - \sum_{ij} f_{\text{damp}}(r_{ij}, R_{ij}^0) C_{6ij} r_{ij}^{-6} \quad (1)$$

This sum is simply added to the total DFT energy, and the gradient of the dispersion energy is added to the QM gradient during optimization. In eq. (1),  $r_{ij}$  is the interatomic distance and  $R_{ij}^0$  is the equilibrium van der Waals (vdW) separation derived from the atomic vdW radii. A particular dispersion scheme is defined by a set of the atomic dispersion coefficients,  $C_6$ , a set of the atomic vdW radii,  $R_{ij}^0$ , and by a damping function,  $f_{\text{damp}}$ . Here we do not consider higher terms in the dispersion energy expansion such as  $C_8$ ,  $C_{10}$ ,  $\dots$ , for the sake of simplicity. The error due to this simplification should be small (although the  $C^8$  contribution is generally not small<sup>22</sup>), since the damping optimization procedure (see later) partially accounts for the missing portion of the attractive interaction. We used a set of the  $C_6$  coefficients by Grimme<sup>49</sup> (0.16, 1.65, 1.11, and 0.7 J nm<sup>6</sup> mol<sup>−1</sup> for H, C, N, and O). The main advantage of this choice is that in Grimme's work each element is assigned only one dispersion coefficient—i.e., there are no different atom types for the same element as in MM force fields. This allows us to treat situations in which atom types would change during an MD simulation, such as certain transition states. Also the input preparation is simplified as no atom type assignment algorithm is needed. Loss of accuracy due to this approximation is not negligible. In the case of the carbon atom it can be as large as 20% of the total dispersion contribution if a particular atom is in sp<sup>3</sup> or sp hybridization because Grimme's coefficient corresponds to the average of all possible states, i.e., it is close to the  $C_6$  of an sp<sup>2</sup> hybridized atom. However, as will be shown later, errors coming from the XC functional inaccuracies are of the same order and sometimes larger; therefore, averaging over atom types is not a major issue (and, moreover, can be easily implemented if needed). Also note that small inaccuracies in the absolute values of the  $C_6$  coefficients can be partially compensated by the damping optimization procedure (see later) and thus only relative values need to

be accurate. Nevertheless, while averaging may work reasonably for H, C, N and O, it is obvious that it can not be applied in many cases such as S ( $\text{SO}_4^{2-}$ ,  $\text{S}^{2-}$ ), P ( $\text{PO}_4^{3-}$ ,  $\text{PH}_3$ ), etc., which we did not introduce into our study for the sake of simplicity at this stage.

The need for the damping function stems from the fact that the  $1/r^6$  form is only an asymptotic expansion, i.e., it is not valid at short distances. What is more, some short-range correlation effects are already contained in the density functional and the damping function has to take this into account. Many different damping functions have been suggested in the literature.<sup>38,39,44–51</sup> We considered a damping function tested by Grimme<sup>49</sup> [eq. (2)] who used  $d = 23$  in the exponent:

$$f_{\text{damp}} = 1/(1 + \exp(-d(r_{ij}/R_{ij}^0 - 1))). \quad (2)$$

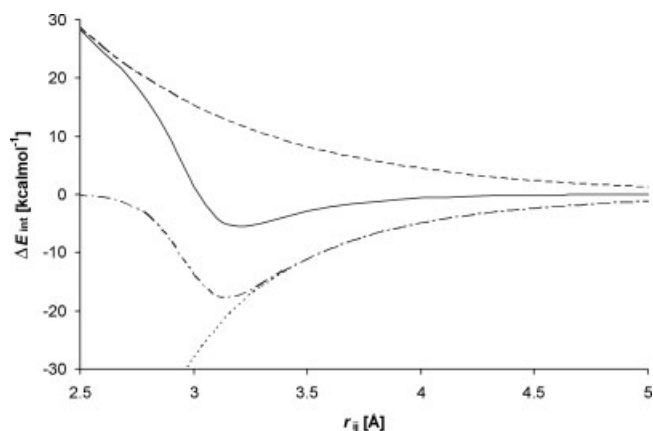
According to our experience, damping is a critical component of the empirical dispersion approach as the errors introduced by incorrect damping can be larger than the errors due to inaccurate  $C_6$  coefficients. The major change to Grimme's<sup>49</sup> scheme is that we discarded the  $s_6$  coefficient introduced to adapt the dispersion strength to various density functionals [ $s_6$  would multiply the RHS (right hand side) of eq. (1)]. The main reason is that  $s_6$  scales the dispersion strength also far from the overlap region where the dispersion interaction is not affected by the choice of a particular density functional. This behavior is unphysical and it introduces some error. Nevertheless, it is obvious that the strength of the added empirical dispersion must be adapted to a functional form given, since the long-range performance of individual XC-functionals varies substantially. Here we will scale the vdW radii by a certain coefficient ( $s_R$ ) rather than scaling the strength of the interaction [eq. (3)].

$$f_{\text{damp}} = 1/(1 + \exp(-d(r_{ij}/(s_R R_{ij}^0) - 1))) \quad (3)$$

In this way the long distance nonoverlapping interactions remain unchanged while the short-range dispersion is damped and extensive double counting is avoided. Our approach is similar to the one adopted in ref. 59. The scaling coefficient  $s_R$  reflects the range of the exchange and correlation covered by a particular XC-functional in terms of distances, which is very natural for the given problem. It also allows for correction of the inaccuracy (or, better, fitness) of the absolute values of vdW radii, and thus only relative values of the radii need to be correct. Figure 1 shows the effect of the damped dispersion (dotted-dashed line) on the otherwise repulsive DFT potential (dashed line).

In a few cases it appeared fruitful to vary also another two parameters of the damped dispersion formula. First, in the case of the B-LYP functional,<sup>60,61</sup> we reintroduced Grimme's  $s_6$  global scaling parameter as a correction for the overrepulsive B88<sup>60</sup> functional. Although the average strength of the interaction could have been corrected solely by the  $s_R$  factor, simultaneous optimization of  $s_6$  and  $s_R$  led to significantly lower errors for B-LYP (in contrary to other functionals). Second, somewhat steeper damping (exponent  $d$  larger than 23) significantly lowered errors for the B3-LYP and TPSS functionals in combination with larger basis sets.

Bondi vdW radii<sup>62</sup> were used throughout this study (i.e., 1.20, 1.70, 1.55, and 1.52 Å for H, C, N, and O, respectively). For a pair of unlike atoms, a combination rule must be applied



**Figure 1.** Dependence of  $\Delta E_{\text{int}}$  of a typical dispersion-bonded complex on the intermolecular separation  $r_{ij}$ . The dashed line represents a purely repulsive potential with no dispersion contribution, the dotted and dotted-dashed lines stand for nondamped and damped dispersion, and the solid line is the resulting potential (sum of the dashed and dotted-dashed lines).

to calculate  $R_{ij}^0$  from  $R_{ii}^0$  and  $R_{jj}^0$ . We have tested both the simple average as used by Grimme<sup>49</sup> and the cubic mean suggested by Halgren<sup>63</sup> and since the cubic mean [eq. (4)] yielded lower errors we will use it throughout this study.

$$R_{ij}^0 = (R_{ii}^{03} + R_{jj}^{03})/(R_{ii}^{02} + R_{jj}^{02}) \quad (4)$$

For  $C_6$  coefficients, we have chosen the combination rule suggested by Wu and Yang<sup>47</sup> [eq. (5)], which performed better than the simple average. The Slater–Kirkwood effective number of electrons ( $N_{\text{eff}}$ ) required can be easily found in the literature for the most common elements ( $N_{\text{eff}} = 0.80, 2.5, 2.82$ , and  $3.15$  for H, C, N, and O, respectively).<sup>63</sup>

$$C_{6ij} = 2(C_{6ii}^2 C_{6jj}^2 N_{\text{eff}i} N_{\text{eff}j})^{1/3} / ((C_{6ii} N_{\text{eff}i}^2)^{1/3} + (C_{6jj} N_{\text{eff}j}^2)^{1/3}) \quad (5)$$

### Notation

In the following text the optimized dispersion parameters will be abbreviated as  $B$ - $s_R$ - $d$ - $s_6$ , where  $B$  is for Bondi vdW radii and  $d$ ,  $s_R$ , and  $s_6$  are the exponents [see eq. (3)] and the respective global scaling factors. For instance,  $B$ -0.90-23-1.25 means that the Bondi vdW radii are multiplied by a factor of 0.90 ( $s_R = 0.90$ ), exponent  $d = 23$ , and the whole dispersion is multiplied by a factor of 1.25 ( $s_6 = 1.25$ ). In most cases  $s_6 = 1.00$  and is omitted. Thus, combination of the TPSS functional, TZVP basis set, and its optimized dispersion reads TPSS/TZVP/ $B$ -0.98-35.

To adjust the dispersion parameters to a given XC functional/basis set combination, we varied the  $s_R$  parameter so that the difference between the reference single point energy and the DFT-D (we use the abbreviation from ref. 49) energy on the same geometry is minimized (i.e., the DFT energy was calculated only once and the dispersion was adjusted). Ideally, the DFT-D method should perform equally well for all sorts of complexes regardless



of whether they are stabilized by induction, dispersion, electrostatics, or any combination of these forces. Unfortunately, it appears that there are systematic differences in the relative stabilities of the hydrogen-bonded and dispersion-bonded complexes. By addition of the empirical dispersion to the DFT energy, these differences are reduced, but not completely removed (see discussion on  $\Delta\Delta_{\text{H-S}}$  term later). Therefore, optimal dispersion parameters are always chosen as a compromise between quality of the description of the complexes of different nature. We have considered the following three criteria (in this order):

1. Minimal standard deviation for the hydrogen-bonded and dispersion-bonded molecules separately. Note that due to some shift in average errors of both groups this does not mean that also the standard deviation for the whole set is minimal.
2. Minimal average error (or stabilization energy shift) for the H-bonded complexes. In most cases the standard deviation of dispersion complexes was sacrificed to maintain good performance for the hydrogen bonds.
3. Minimal difference between average errors of H-bonded and dispersion-bonded complexes. This leads to minimal standard deviation for the whole set of molecules.

Mostly only the  $s_R$  scaling parameter was varied, but in some cases, especially with large basis sets, we changed also the exponent  $d$  [see eq. (3)]. The damping function is thus steeper and smaller  $\Delta\Delta_{\text{H-S}}$  and consequently smaller overall standard deviation is achieved.

The parameters obtained in this way were used to optimize the geometries of the complexes without any restraints and derivatives were evaluated analytically. The interaction energies for optimized geometries were calculated at the same level as the optimization, without counterpoise correction.

All S-VWN,<sup>64,65</sup> B-VWN,<sup>60,65</sup> B-LYP,<sup>60,61</sup> B3-LYP,<sup>66</sup> BH-LYP,<sup>67</sup> PBE,<sup>68</sup> TPSS,<sup>69</sup> and TPSSH<sup>70</sup> calculations were performed with TurboMole 5.7 program.<sup>71</sup> The dispersion energy was calculated by a simple Fortran code coupled to TurboMole by a modified “jobex” script. All MP2 calculations and the DFT calculations without explicit HF exchange were done with the RI approximation<sup>72–76</sup> and the default grid, after we made sure that neither RI approximation nor the grid size introduced any significant error. The following basis sets were used: Pople 6-31G\*\*,<sup>77</sup> 6-311++G(3df,3pd)<sup>78</sup> (abbreviated as LP in this paper), Ahlrichs TZVP,<sup>79</sup> Dunning’s aug-cc-pVQZ<sup>80</sup> with both  $g$  functions and the most diffuse  $f$  function removed from the heavy atoms, and analogically reduced for the hydrogen atom (abbreviated as aQZ') and cc-pVTZ.<sup>83a</sup>

### Training Set of Molecules

The quality of the results of any semiempirical density functional or an empirical method necessarily depends on the quality and balance of the training set. Common training sets for DFT usually contain only very few hydrogen-bonded complexes and even fewer dispersion-bonded ones. What is more, in our experience, popular small dispersion-bonded complexes like rare gas dimers do not represent larger molecules well. They exhibit different error cancellation behavior, show different basis set dependencies

of the interaction energy components, and usually have quite small interaction energy, implying that even large (relative) errors have small weights in the fitting procedure. As a result, inclusion of such small molecules does not guarantee good results for stronger molecular complexes (see ref. 8). Therefore, we believe that inclusion of several larger complexes in the training set is unavoidable if it is to serve its purpose well.

A balanced training set for parametrization of weak interactions should contain all sorts of complexes with interactions varying from weak to strong. Also, the contribution of different interactions should correspond to their relative strength and abundance in nature. Here we used a set published in one of our previous paper,<sup>54</sup> named S22. It consists of 22 complexes, out of which seven are hydrogen-bonded with interaction energies between  $-3.2$  and  $-20.7$  kcal mol<sup>-1</sup>, eight are predominantly dispersion-bonded (between  $-0.5$  and  $-12.2$  kcal mol<sup>-1</sup>), and seven are mixed complexes in which dispersion and electrostatic contributions are comparable (from  $-1.2$  to  $-7.1$  kcal mol<sup>-1</sup>). Total and average interaction energies of these three groups are  $-98$  and  $-14$  kcal mol<sup>-1</sup> for hydrogen-bonded,  $-38$  and  $-4.8$  kcal mol<sup>-1</sup> for dispersion-bonded, and  $-27$  and  $-3.9$  kcal mol<sup>-1</sup> for mixed complexes. We believe that these values represent a typical situation for proteins, DNA, and RNA (see also refs. 1 and 54).

Reference geometries were obtained either by CCSD(T)/cc-pVTZ or cc-pVQZ numerical gradient optimization in the case of small complexes, or by MP2/cc-pVTZ counterpoise corrected optimization with analytical gradients (see Table 1). Both levels of theory should provide very good geometries, although partly on account of some error cancellation.<sup>81</sup> Interaction energies were calculated as a sum of the MP2 complete basis set energy (estimated by an extrapolation procedure<sup>82,83</sup>) and a correction for higher order correlation effects  $\Delta\text{CCSD(T)}$  (for more details see Table 1 and ref. 55). All geometries and interaction energies used in this paper are available in ref. 54.

All energy values in this paper are interaction energies not corrected for the deformation energy of the monomers. This greatly simplifies the calculations and also circumvents additional complications with potentially inaccurate deformation energies of monomers.

## Results and Discussion

Organization of the discussion is as follows. First, we will discuss the main errors in the DFT calculations and their consequences for the combined DFT-D method. Then we will focus on the single point calculations on reference geometries and comment on the limitations of the DFT-D approach and on the limitations coming from the current XC-functionals. After that, results of the geometry optimizations and single point calculations on the optimized geometries will be reviewed. Finally, transferability of the dispersion parameters obtained for the S22 training set will be assessed using a larger validation set.

### Basis Set Size and BSSE

BSSE is certainly one of the most serious errors in the WFT but also in the DFT calculations of intermolecular interaction ener-

**Table 1.** Reference WFT Interaction Energies for the Training Set of Molecules S22, Adapted from ref. 54.

Complex (symmetry)	$\Delta E_{\text{int}}^{\text{a}}$ (kcal mol <sup>-1</sup> )	Energy <sup>b</sup> MP2/CCSD(T)	Geometry <sup>c</sup>
H-bonded (7)			
(NH <sub>3</sub> ) <sub>2</sub> ( <i>C</i> <sub>2h</sub> )	-3.17	QZ-5Z/TZ	CCSD(T)/QZ
(H <sub>2</sub> O) <sub>2</sub> ( <i>C</i> <sub>s</sub> )	-5.02	QZ-5Z/TZ	CCSD(T)/QZ
Formic acid dimer ( <i>C</i> <sub>2h</sub> )	-18.61	QZ-5Z/TZ	CCSD(T)/TZ
Formamide dimer ( <i>C</i> <sub>2h</sub> )	-15.96	QZ-5Z/TZ	CCSD(T)/TZ
Uracil dimer ( <i>C</i> <sub>2h</sub> )	-20.65	TZ-QZ/TZ'	MP2/TZ-CP
2-Pyridoxine...2-aminopyridine ( <i>C</i> <sub>1</sub> )	-16.71	TZ-QZ/TZ'	MP2/TZ-CP
A...T Watson-Creek ( <i>C</i> <sub>1</sub> )	-16.37	TZ-QZ/DZ	MP2/TZ-CP
Dispersion-bonded (8)			
(CH <sub>4</sub> ) <sub>2</sub> ( <i>D</i> <sub>3d</sub> )	-0.53	QZ-5Z/TZ	CCSD(T)/TZ
(C <sub>2</sub> H <sub>4</sub> ) <sub>2</sub> ( <i>D</i> <sub>2d</sub> )	-1.51	QZ-5Z/TZ	CCSD(T)/TZ
Benzene...CH <sub>4</sub> ( <i>C</i> <sub>3</sub> )	-1.50	QZ-5Z/TZ'	MP2/TZ-CP
Benzene dimer ( <i>C</i> <sub>2h</sub> )	-2.73	aDZ-aTZ/aDZ	MP2/TZ-CP
Pyrazine dimer ( <i>C</i> <sub>s</sub> )	-4.42	aTZ-aQZ/TZ'	MP2/TZ-CP
Uracil dimer ( <i>C</i> <sub>1</sub> )	-10.12	TZ-QZ/TZ'	MP2/TZ-CP
Indole...benzene ( <i>C</i> <sub>1</sub> )	-5.22	aDZ-aTZ/DZ	MP2/TZ-CP
A...T stack ( <i>C</i> <sub>1</sub> )	-12.23	aDZ-aTZ/DZ	MP2/TZ-CP
Mixed (7)			
Ethene...ethyne ( <i>C</i> <sub>2v</sub> )	-1.53	aQZ-a5Z/TZ	CCSD(T)/TZ
Benzene...H <sub>2</sub> O ( <i>C</i> <sub>2</sub> )	-3.28	QZ-5Z/TZ'	MP2/TZ-CP
Benzene...NH <sub>3</sub> ( <i>C</i> <sub>s</sub> )	-2.35	QZ-5Z/TZ'	MP2/TZ-CP
Benzene...HCN ( <i>C</i> <sub>s</sub> )	-4.46	aTZ-aQZ/TZ'	MP2/TZ-CP
Benzene dimer ( <i>C</i> <sub>2v</sub> )	-2.74	aDZ-aTZ/aDZ	MP2/TZ-CP
Indole...benzene T-shape ( <i>C</i> <sub>1</sub> )	-5.73	aDZ-aTZ/DZ	MP2/TZ-CP
Phenol dimer ( <i>C</i> <sub>1</sub> )	-7.05	TZ-QZ/TZ'	MP2/TZ-CP

Deformation energy of monomers is not included.

<sup>a</sup>Deformation energy of monomers and ZPE corrections are not included.

<sup>b</sup>Basis sets used for MP2 extrapolation/ $\Delta$ CCSD(T) correction. XZ and aXZ stands for cc-pVXZ, and aug-cc-pVXZ, respectively. TZ-QZ means a value extrapolated using these two basis sets. TZ' is TZ basis sets with modified polarization coefficients—see ref. 54.

<sup>c</sup>Method and basis set for gradient optimization. CP stands for counterpoise correction applied in the course of optimization.

gies. Table 2 shows the basis set dependence of the average interaction energies and the BSSE for the TPSS functional and S22 training set. Regarding the BSSE, other functionals tested behave similarly. For comparison, also the MP2/cc-pVTZ results are shown. Table 2 illustrates a well known fact that BSSE is much larger for the WFT-correlated calculations than for DFT. The average TPSS/6-31G\*\* BSSE is similar to the MP2/cc-pVTZ one and the BSSE of the DFT/TZVP is about four times smaller than the MP2/cc-pVTZ one (although the TZVP basis is smaller than the cc-pVTZ basis). Interestingly, no significant improvement is achieved by passing from TZVP to 6-311++G(3df,3pd), where the average BSSE is still close to 0.5 kcal mol<sup>-1</sup>. For our largest basis set, aQZ', BSSE is negligibly small. Basis sets of the TZVP quality are often used to parametrize empirical XC functionals, which is probably justified in the case of thermochemical data.<sup>84</sup> However, Table 2 shows that triple- $\zeta$  quality basis sets are incapable of describing the noncovalent interactions with sufficient accuracy. If the counterpoise correction<sup>85</sup> can be applied, at least a triple- $\zeta$  basis set augmented with diffuse functions should be used for the intermolecular interactions. If the CP correction can not be used for practical or

principal reasons (like in the case of intramolecular hydrogen bonds), a still larger basis set such as aQZ' is necessary.

### Dispersion Energy and BSSE

As mentioned in the introduction, in the case of the hydrogen-bonded complexes, BSSE provides some artificial attraction, which may compensate for the missing dispersive interaction in DFT. Indeed, the average BSSE value of 3.17 kcal mol<sup>-1</sup> for commonly used 6-31G\*\* basis set (see Table 2) more than compensates for the dispersion complementing the aQZ' DFT calculation (-2.01 kcal mol<sup>-1</sup>). As a result, the counterpoise uncorrected TPSS/6-31G\*\* intermolecular interaction energies of hydrogen-bonded complexes are overestimated. This applies to all functionals tested if a small basis set is used (see Table 3). The overestimated interaction energy is then reflected in too short intermolecular distances for these functionals (see later). The only exception is the B-LYP/6-31G\*\* combination which predicts slightly too short intermolecular separations for small hydrogen-bonded complexes but slightly too long for larger complexes (for more detailed discussion see Geometries section). When enlarging the basis set to

**Table 2.** Average Interaction Energies and BSSE for S22 Training Set of Molecules and Selected DFT and WFT Methods and Basis Sets.

Complexes	Methods					WFT <sup>c</sup>
	DFT 6-31G**	DFT TZVP	DFT LP <sup>a</sup>	DFT aQZ <sup>b</sup>	MP2 cc-pVTZ	
H-bonded (7)						−13.78
TPSS	−15.77	−12.46	−12.33	−11.94		
TPSS CP	−12.60	−11.84	−11.87	−11.90		
BSSE	3.17	0.62	0.46	0.03	2.57	
D <sup>d</sup>	−1.78	−2.19	−2.12	−2.01		
TPSS/CP+D	−14.38	−14.03	−13.99	−13.91	−12.74	
Disp. bonded (8)						−4.78
TPSS	−0.22	1.56	0.54	1.17		
TPSS CP	1.73	2.08	1.30	1.21		
BSSE	1.94	0.52	0.76	0.04	2.11	
D <sup>d</sup>	−5.08	−5.85	−5.71	−5.72		
TPSS/CP+D	−3.14	−3.77	−4.41	−4.51	−5.19	
Mixed (7)						−3.88
TPSS	−2.51	−1.11	−1.72	−1.18		
TPSS CP	−1.40	−0.72	−1.16	−1.16		
BSSE	1.11	0.39	0.56	0.02	1.35	
D <sup>d</sup>	−2.07	−2.43	−2.37	−2.27		
TPSS/CP+D	−3.47	−3.15	−3.53	−3.43	−3.83	

Values in kcal mol<sup>−1</sup>; CP for counterpoise correction. Averages over 7 hydrogen-bonded, 8 dispersion-bonded, and 7 mixed complexes.

<sup>a</sup>LP: 6-311++G(3df,3pd).

<sup>b</sup>Modified aug-cc-pVQZ (see Methods).

<sup>c</sup>See Table 1.

<sup>d</sup>Dispersion optimized for a given TPSS/CP/basis set combination, see below.

TZVP or bigger, interaction energies for reference geometries become underestimated for all functionals except LDA, most likely due to the missing dispersion interaction. Regarding the dispersion-bonded complexes, interaction energies calculated by pure DFT on reference geometries are on average positive (repulsive), except for S-VWN. Ordering of the functionals according to the average error in interaction energy is S-VWN(−0.15) < PBE(4.90) < BH-LYP(5.31) < TPSSh(5.81) < B3-LYP(6.56) < TPSS(6.75) < B-LYP(7.74) < B-VWN(10.35 kcal mol<sup>−1</sup>) for the TZVP/CP calculation. The BSSE error is not sufficient to supply enough spurious attraction and the average error of the GGA-based functionals is thus close to 100% for the dispersion-bonded complexes.

#### Dispersion Adjustment to an XC-Functional/Basis Set Combination. Transferability of the Dispersion Parameters

The subdivision of the noncovalent complexes into hydrogen-bonded and dispersion-bonded subgroups in Table 1 is not purposeless. It appears that the two categories of complexes show significantly different behavior when the DFT energy is combined with the empirical dispersion. In general, the amount of the dispersion energy, which is necessary for the dispersion-bonded complexes, leads to overestimation of the hydrogen-bonded complexes. A (necessarily very arbitrary) measure of this effect is a difference between the average errors of the interaction energy of the hydrogen-bonded complexes ( $\Delta_H$ , Avg for the H-bonded

complexes in Tables 3 and 4), and the dispersion-bonded complexes ( $\Delta_D$ , Avg for the dispersion-bonded complexes in Tables 3 and 4),  $\Delta\Delta_{D-H} = \Delta_D - \Delta_H$  (last columns of these tables). For the single point DFT-D calculations on the reference geometries it amounts to 0.3–1.3 kcal mol<sup>−1</sup> with the 6-311++G(3df,3pd) basis set, depending on the functional. In fact, this discrepancy between the estimated interaction energies of different types of complexes represents the most serious problem of our DFT-D scheme. Analysis of Table 4 in ref. 49 indicates that this behavior could be an inherent problem of the DFT-D approach. Nevertheless, improvement over pure DFT (see Table 3) is remarkable (3–12 times smaller  $\Delta\Delta_{D-H}$  shift, depending on the XC functional used), since the major part of the pure DFT  $\Delta\Delta_{D-H}$  value comes from the missing dispersion energy. Note also that the MP2 method encounters a similar problem, except that  $\Delta\Delta_{D-H}$  is negative (see the last rows of Tables 3 and 4). In the MP2 method, the  $\Delta\Delta_{D-H}$  shift is a consequence of the known overestimation of the dispersion energy and is reflected by a positive value of  $\Delta\text{CCSD(T)}$  for the stacked base pairs.<sup>55</sup>

What is the origin of the remaining (nondispersion)  $\Delta\Delta_{D-H}$  shift? A possible explanation is, that while monomers in the dispersion-only complexes are almost unperturbed by the interaction, the hydrogen-bonded monomers are significantly perturbed by mutual orbital interaction (induction). The correlation effects can thus reach further, being mediated by the shared electron density. Can we correct for those effects in DFT-D? More detailed analysis shows that if hydrogen-bonded or dispersion-bonded complexes are treated separately, very small DFT-D errors (both average and standard deviation) can be achieved for both groups, but with different  $s_R$  scaling. This indicates that double counting is treated effectively by our (simple) damping function, but different kinds of complexes require different scaling. Experiments with steeper or differently shaped damping functions (the steeper the function the smaller the  $\Delta\Delta_{D-H}$  shift) revealed that similarly good performance for both groups at the same time can not be achieved with a single  $s_R$  for any damping function. (Note that zeroing of the  $\Delta\Delta_{D-H}$  shift through simultaneous  $s_6$  and  $s_R$  scaling in the case of the B-LYP functional should be understood rather as a coincidence. Simultaneous  $s_6$  and  $s_R$  scaling brings only small or no improvement for other XC-functionals. Another (and in our opinion more likely) explanation is that the  $\Delta\Delta_{D-H}$  shift is a consequence of the known erroneous long-range behavior of the XC functionals.<sup>32</sup> This idea is supported by the relatively wide variation of the  $\Delta\Delta_{D-H}$  shift among different XC functionals (see Table 3). Moreover, more advanced functionals (hybrid B3-LYP or meta-GGA TPSS) exhibit much smaller  $\Delta\Delta_{D-H}$  shift than simple LDA or GGA functionals. In our opinion, regarding intermolecular interactions, current GGA-based functionals have not reached their limits yet. We believe that significant improvements can be achieved through appropriate optimization of their large-gradient–small-density parts, irrespective of whether this step is physically sound or just a practically motivated remedy. Work on this topic is under way.

#### Overall Performance of the DFT-D Method

In the Tables 3 and 4 pure DFT and DFT-D results of the single point calculations on the reference geometries are shown. For the hydrogen-bonded complexes the results of most GGA functionals

**Table 3.** Single-Point Calculations on Reference Geometries, S22 Set: Comparison of the DFT Results with the Reference WFT Data.

X-C funct.	Basis set <sup>a</sup>	S22		H-bonded		Dispersion bonded		Mixed		$\Delta\Delta_{D-H}$
		SD	Avg	SD	Avg	SD	Avg	SD	Avg	
PBE	6-31G**	3.17	0.34	0.97	−3.17	2.56	3.24	1.00	0.55	6.41
	6-31G** CP	3.29	2.66	0.99	0.40	3.98	5.40	1.25	1.80	5.01
	TZVP	2.70	2.06	1.00	0.05	3.07	4.36	1.16	1.45	4.31
	TZVP CP	2.87	2.63	1.04	0.77	3.56	4.90	1.16	1.88	4.14
	LP	2.50	2.20	0.84	0.75	3.23	4.08	1.13	1.51	3.33
	LP CP	2.76	2.69	0.86	1.11	3.64	4.69	1.27	1.98	3.58
	aQZ'	2.71	2.56	0.79	0.95	3.54	4.56	1.23	1.90	3.61
B-LYP	6-31G**	3.83	2.15	0.99	−1.56	3.76	5.56	1.45	1.95	7.13
	6-31G** CP	4.23	4.81	1.64	2.51	5.40	8.07	1.77	3.37	5.56
	TZVP	3.72	4.29	1.69	2.11	4.55	7.19	1.78	3.16	5.08
	TZVP CP	3.91	4.87	1.74	2.85	5.04	7.74	1.80	3.61	4.89
	LP	3.61	4.62	1.50	2.89	4.72	7.23	1.73	3.38	4.33
	LP CP	3.82	5.03	1.54	3.21	5.05	7.74	1.87	3.76	4.53
	aQZ'	3.74	4.81	1.47	3.06	5.07	7.33	1.79	3.66	4.27
TPSS	6-31G**	3.37	1.40	1.00	−2.01	3.14	4.45	1.10	1.33	6.47
	6-31G** CP	3.60	3.47	1.18	1.16	4.45	6.39	1.41	2.45	5.24
	TZVP	3.24	3.55	1.21	1.30	3.81	6.23	1.49	2.73	4.93
	TZVP CP	3.43	4.06	1.25	1.92	4.28	6.75	1.52	3.12	4.83
	LP	2.86	3.03	1.08	1.43	3.64	5.21	1.32	2.12	3.78
	LP CP	3.17	3.63	1.09	1.89	4.12	5.97	1.50	2.68	4.08
	aQZ'	3.21	3.55	1.05	1.82	4.27	5.84	1.48	2.66	4.02
B3-LYP	TZVP	3.07	2.86	1.10	0.83	3.75	5.32	1.35	2.08	4.48
	TZVP CP	3.25	3.36	1.15	1.45	4.19	5.81	1.39	2.46	4.36
	6-31G**	3.51	1.55	0.97	−1.98	3.28	4.73	1.21	1.45	6.72
	6-31G** CP	3.82	3.85	1.27	1.44	4.72	6.97	1.49	2.68	5.52
	TZVP	3.28	3.30	1.24	1.06	3.85	6.01	1.48	2.45	4.95
	TZVP CP	3.47	3.84	1.28	1.70	4.34	6.56	1.52	2.86	4.86
	LP	3.12	3.50	1.17	1.73	3.94	5.91	1.41	2.52	4.18
S-VWN	LP CP									
	6-31G**	4.17	−4.14	3.06	−9.45	0.86	−1.19	1.29	−2.19	8.26
	6-31G** CP	3.18	−1.79	1.96	−5.76	1.30	0.91	0.76	−0.92	6.66
	TZVP	2.70	−2.73	1.82	−6.19	0.59	−0.80	0.99	−1.48	5.39
B-VWN	TZVP CP	2.61	−2.13	1.70	−5.50	0.59	−0.15	0.77	−1.02	5.35
BH-LYP	TZVP CP	5.23	6.59	2.81	5.47	6.64	9.88	3.50	3.94	4.41
MP2	TZVP CP	5.13	7.39	2.88	6.10	7.08	10.35	2.56	5.28	4.25
MP2	TZVP	2.96	2.31	0.93	−0.01	3.36	4.81	1.26	1.76	4.82
MP2	TZVP CP	3.13	2.77	0.94	0.53	3.79	5.31	1.29	2.10	4.78
MP2	cc-pVTZ	1.28	−1.95	0.46	−1.56	1.91	−2.61	0.60	−1.35	−1.06
	cc-pVTZ CP	0.82	0.15	0.57	1.02	0.56	−0.50	0.01	0.31	−1.51

Standard deviations (SD) and average signed errors (Avg) in kcal mol<sup>−1</sup>.<sup>a</sup>LP is 6-311++G(3df,3pd), aQZ' is a modified aug-cc-pVQZ (see Methods).

are relatively acceptable, depending on the size of the basis set. In general, when enlarging the basis set, the interaction energy decreases and becomes smaller than the reference WFT value. As expected, the PBE functional gives the largest and the B-LYP functional the smallest interaction energies. Regarding the dispersion-bonded complexes, performance of the GGA functionals is traditionally very poor with errors over 100%. Note a remarkably good performance of the S-VWN functional for stacked complexes and at the same time bad results for PBE. According to Zhang

et al.,<sup>17</sup> S-VWN greatly overestimates bonding in the rare gas dimers, while PBE behaves well. Such contradictory results make us emphasize again the need for larger dispersion-bonded complexes in the training sets. Obviously, small dispersion-bonded complexes with interaction energies lower than ~1 kcal mol<sup>−1</sup> misrepresent the dispersion interaction in biological molecules and their inclusion into training sets is, from this perspective, pointless.

The empirical dispersion term improves the results significantly, if damped properly. In most cases it is enough to opti-



**Table 4.** Single-Point Calculations on the Reference Geometries, S22 Set: Comparison of the DFT-D Results with the Reference WFT Data.

X-C funct.	Basis set <sup>a</sup>	Dispersion type <sup>b</sup>	S22		H-bonded		Dispersion bonded		Mixed		$\Delta\Delta_{D-H}$
			SD	Avg	SD	Avg	SD	Avg	SD	Avg	
PBE	6-31G**	B-1.15-23	2.25	−1.15	1.01	−4.01	0.96	0.87	0.60	−0.61	4.88
	6-31G** CP	B-1.00-23	1.02	−0.14	0.58	−1.20	0.81	0.76	0.34	−0.11	1.96
	TZVP	B-1.05-23	0.91	−0.28	0.56	−1.27	0.55	0.53	0.48	−0.20	1.80
	TZVP CP	B-1.02-33	0.71	−0.01	0.55	−0.70	0.64	0.52	0.26	0.08	1.22
	LP	B-1.06-23	0.66	−0.04	0.38	−0.52	0.80	0.41	0.31	−0.08	0.93
	LP CP	B-1.02-33	0.56	0.05	0.38	−0.36	0.68	0.30	0.33	0.19	0.66
	aQZ'	B-1.05-23	0.78	0.22	0.38	−0.38	0.96	0.72	0.37	0.25	1.10
B-LYP	6-31G**	B-1.05-23	2.17	−0.19	0.73	−2.89	1.17	1.73	0.84	0.31	4.62
	6-31G** CP	B-0.87-23	1.40	0.98	0.65	−0.01	1.63	2.15	0.49	0.64	2.17
	TZVP	B-0.87-23	1.00	0.47	0.72	−0.42	0.81	1.27	0.66	0.44	1.69
	TZVP CP	B-0.82-23	1.10	0.81	0.70	−0.01	1.21	1.72	0.35	0.59	1.73
	LP	B-0.79-23	0.81	0.48	0.35	−0.10	0.90	1.19	0.29	0.26	1.29
	LP CP	B-0.75-23	1.04	0.81	0.38	0.04	1.21	1.69	0.30	0.56	1.65
	LP	B-0.90-23-1.25	0.32	0.09	0.36	0.03	0.30	0.00	0.30	0.23	−0.03
	aQZ'	B-0.78-23	1.00	0.65	0.34	0.03	1.39	1.30	0.29	0.52	1.27
TPSS	6-31G**	B-1.05-23	1.85	−0.94	0.98	−3.34	0.56	0.62	0.54	−0.31	3.96
	6-31G** CP	B-0.97-23	1.07	0.40	0.56	−0.62	1.03	1.31	0.40	0.38	1.93
	TZVP	B-0.98-35	0.84	0.53	0.56	−0.38	0.40	1.15	0.64	0.74	1.53
	TZVP CP	B-0.91-33	0.73	0.46	0.58	−0.28	0.54	0.90	0.46	0.69	1.18
	LP	B-0.96-27	0.38	−0.15	0.48	−0.39	0.30	−0.07	0.25	0.01	0.32
	LP CP	B-0.88-23	0.45	0.12	0.41	−0.23	0.42	0.26	0.35	0.32	0.49
	aQZ'	B-0.93-35	0.56	0.11	0.44	−0.19	0.69	0.12	0.39	0.39	0.31
	TZVP	B-1.00-33	0.86	0.23	0.44	−0.64	0.73	0.93	0.48	0.29	1.56
	TZVP CP	B-0.92-27	0.69	0.34	0.44	−0.23	0.76	0.74	0.41	0.47	0.98
B3-LYP	6-31G**	B-1.10-23	2.21	−0.34	0.78	−3.05	1.24	1.69	0.68	0.06	4.74
	6-31G** CP	B-0.95-33	1.10	0.57	0.47	−0.42	1.10	1.48	0.55	0.53	1.91
	TZVP	B-0.95-27	0.80	0.05	0.42	−0.83	0.42	0.61	0.70	0.28	1.43
	TZVP CP	B-0.93-33	0.75	0.40	0.44	−0.31	0.65	0.87	0.59	0.58	1.18
	LP	B-0.93-35	0.43	0.23	0.35	−0.14	0.32	0.43	0.41	0.37	0.57
MP2	cc-pVTZ		1.28	−1.95	0.46	−1.56	1.91	−2.61	0.60	−1.35	−1.06
	cc-pVTZ CP		0.82	0.15	0.57	1.02	0.56	−0.50	0.31	0.01	−1.51

Standard deviations (SD) and average signed errors (Avg) in kcal mol<sup>−1</sup>.<sup>a</sup>LP is 6-311++G(3df,3pd), aQZ' is a modified aug-cc-pVQZ (see Methods).<sup>b</sup>For the optimized dispersion notation see Methods.

mize the global scaling factor  $s_R$ , which accounts for the reach of correlation and exchange effects in a particular XC functional. In a few cases we optimized also the steepness of the damping function [exponent  $d$  in the eq. (2)] because it decreased the  $\Delta\Delta_{D-H}$  shift and led to smaller overall errors. In the case of the B-LYP functional we found that reintroduction of Grimme's<sup>49</sup> global  $s_6$  scaling in combination with our  $s_R$  scaling brings a significant improvement for large basis sets (see Table 4). Because the plane wave basis sets are equivalent to very large standard basis sets, we can recommend the B-0.90-23-1.25 combination also for the plane wave-based B-LYP calculations. This combination, together with TPSS/LP/B-0.96-27, performs remarkably well and (statistically) significantly surpasses the quality of the MP2/cc-pVTZ results (see the last line of Table 4). To a lesser extent this holds also for several other

combinations even with the TZVP basis set. In most cases, increase in the quality of the basis set leads to improvement of the DFT-D results, contrary to pure DFT results. This behavior is most desirable and indicates that the dispersion correction to the GGA and meta-GGA DFT theory is physically sound.

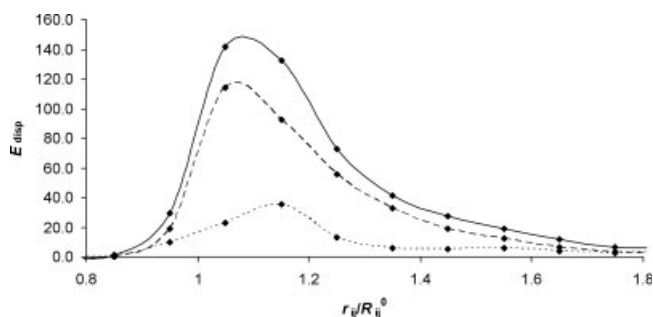
When the dispersion-bonded complexes were fitted separately from the others, a standard deviation of about 0.3 kcal mol<sup>−1</sup> could be typically achieved along with a small average error, despite a relatively crude dispersion model. Note that dispersion-bonded complexes are very sensitive to the choice of the  $C^6$  coefficients as the dispersion contribution is large when compared with the total interaction energy. However, using the same damping parameters for the rest of the complexes would lead to unacceptable standard deviation close to 1 kcal mol<sup>−1</sup>. The influence of the damping can therefore be considered more criti-

cal than the errors due to inaccurate atomic  $C^6$  coefficients for the complexes considered.

### Dispersion Contribution to the Total Interaction Energy

The optimized empirical dispersion contribution to the DFT interaction energy varies for different functionals in a relatively wide range. On average (the S22 average), it amounts to  $-2.34$  kcal mol $^{-1}$  for the PBE/aQZ',  $-3.44$  kcal mol $^{-1}$  for the TPSS/aQZ', and  $-4.78$  kcal mol $^{-1}$  for the B-LYP/aQZ' combination. If we choose the TPSS XC-functional as the best compromise (PBE exhibits too attractive exchange repulsion and B-LYP is overrepulsive), the dispersion contribution to the interaction energy of the hydrogen-bonded complexes is  $-2.01$  out of  $-13.76$  kcal mol $^{-1}$ , i.e., about 15%, and it is close to 100% for the dispersion-bonded complexes (see Table 2). What is the relationship of these numbers to other available theoretical results? We have compared our empirical dispersion with the published SAPT- or TD-DFT-based data available for the water dimer, ammonia dimer, and benzene dimer<sup>19,20,22</sup> and also to our own SAPT calculations (unpublished). For both hydrogen-bonded and dispersion-bonded dimers our *non-damped* empirical values are relatively close to the SAPT data or to the TD-DFT polarizability based data. This is not surprising, since the empirical dispersion estimates based on the polarization expansion are widely believed to be fairly accurate if proper  $C^6$  coefficients are used. Interestingly, when damping is applied, its effect on the hydrogen-bonded complexes is much more profound than its effect on the dispersion-bonded ones. While the former lose about 80–90% of the dispersion contribution by damping, the latter lose only about 10–15%. Obviously, combination of the non-damped dispersion, like the SAPT dispersion, with the pure DFT results would lead to enormous errors (overestimation) in the case of the hydrogen-bonded complexes. A simple explanation could be that DFT already covers some intermolecular correlation interaction at relatively short distances found in the hydrogen-bonded complexes characterized by significant density overlap. This question also relates to the  $\Delta\Delta_{H-S}$  shift discussed earlier.

Figure 2 shows the distance dependence of the damped dispersion energy contribution for the TPSS/LP/B-0.96-27 method. It is calculated as a sum of the dispersion contributions for 80 molecules of the training and validation sets; the dotted, dashed, and



**Figure 2.** Damped dispersion energy as a function of  $r_{ij}/R_{ij}^0$  (interatomic distance divided by the vdW equilibrium distance for the respective atom pair). Average over 80 molecules of the training and validation sets (solid), and over all hydrogen-bonded (dotted) and dispersion-bonded complexes (dashed).

full lines correspond to the hydrogen-bonded, dispersion-bonded, and all complexes, respectively. On the y axis  $r_{ij}/R_{ij}^0$  is an interatomic distance divided by the vdW equilibrium distance for the respective atom pair. Because most of the complexes in our set are optimized *in vacuo*, the distribution of the  $r_{ij}/R_{ij}^0$  corresponds to the complexes in their optimal geometries *in vacuo*. For instance, the dispersion contribution of all the atom pairs which are closer than their equilibrium vdW distance (according to Bondi vdW radii) is the integral under the full curve from  $r_{ij}/R_{ij}^0 = 0$  to  $r_{ij}/R_{ij}^0 = 1$ . Apparently, most of the dispersion comes from the  $r_{ij}/R_{ij}^0$  between 1.0 and 1.4, i.e., from the interatomic distances between about 3.3 and 4.6 Å. At these distances most of the widely used basis sets still have some overlap, although this overlap varies over a very wide range, depending on the basis set size and the Gaussian exponents. This means that frequent attempts to construct a local or semilocal XC functional, which gives an additional attractive contribution originating from the small-density-large-gradient regions, similar in magnitude to the missing long-range correlation, may be partly successful. Such a functional would not, of course, represent dispersion in the strict sense of the nonlocal interactions at large distances. It might, however, represent a reasonable, or at least better, approximation for correlation in the small overlap regime which is often associated with the term “dispersion.” Although such a method can not be accurate asymptotically, it could turn out to be useful until theoretically sound and practical nonlocal XC functionals are designed.

### Geometry Optimization with DFT-D

For the B-LYP, PBE, and TPSS functionals, the dispersion parameters adjusted to fit the reference WFT interaction energies were used to optimize the geometries of the complexes. In this DFT-D optimization we let the whole complex relax and the gradients were calculated analytically (the dispersion gradient was added to the DFT gradient). Table 5 shows a comparison of the DFT-D geometries with the reference WFT geometries for the 22 complexes of the S22 set (for more details see Table S1 in Supplementary Material which contains also the additional 21 optimized complexes of the validation set). We focused on the distance of the centers of mass of the two interacting molecules ( $d_{AB}$ ) and RMS fit of the optimized and reference geometries. The DFT-D results are shown in the second half of the table, while the first half displays  $d_{AB}$  and RMS for the full gradient optimizations by pure DFT for comparison.

For H-bonded complexes pure DFT performs relatively well. With small basis sets like 6-31G\*\* the PBE, TPSS, and B3-LYP functionals give slightly too short hydrogen bonds, partly due to large BSSE. On the other hand, B-LYP hydrogen bonds are on average about right for the 6-31G\*\* basis set, while they are getting too long for TZVP and 6-311++G(3df,3pd) basis sets. Good performance of the B-LYP and B3-LYP functionals with small basis sets is highly valued in applications. When increasing the basis set to TZVP or 6-311++G(3df,3pd), intermolecular distances in most cases prolong, partly due to reduced BSSE. As expected, the plain DFT results for dispersion-bonded complexes are unsatisfying. Some monomers, like the AT stack, form different complexes during the optimization (hydrogen-bonded, if possible). In most cases, the intermolecular distance increases until the gradient criterion is met, but the interaction may be

**Table 5.** Comparison of the DFT and DFT-D Geometries with Reference WFT Geometries.

X-C func.	Basis set <sup>a</sup>	Dispersion type <sup>b</sup>	$\Delta d_{AB}$				RMS			
			All	H-bonded	Disp. bonded	Mixed	All	H-bonded	Disp. bonded	Mixed
PBE	6-31G**	None	0.36	−0.056	0.93	0.12	7.64	0.47	5.74	1.43
	TZVP	None	0.40	−0.015	0.89	0.24	6.56	0.23	5.14	1.19
	LP	None	0.27	−0.017	0.63	0.14	4.52	0.15	3.77	0.59
B-LYP	6-31G**	None	0.62	0.011	1.43	0.31	9.72	0.32	7.60	1.80
	TZVP	None	0.68	0.058	1.42	0.47	10.16	0.30	7.95	1.91
	LP	None	0.59	0.063	1.16	0.48	8.28	0.29	6.13	1.86
TPSS	6-31G**	None	0.47	−0.032	1.06	0.30	8.04	0.27	6.29	1.48
	TZVP	None	0.56	0.003	1.13	0.47	7.57	0.26	5.39	1.92
	LP	None	0.34	−0.002	0.69	0.27	4.73	0.15	3.49	1.09
B-3LYP	6-31G**	None	0.36	−0.007	0.83	0.19	6.87	0.18	5.33	1.36
S-VWN	6-31G**	None	−0.25	−0.248	−0.21	−0.28	4.20	1.55	1.20	1.44
	TZVP	None	−0.18	−0.200	−0.13	−0.22	2.82	0.94	0.78	1.11
PBE	6-31G**	B-1.15-23	0.16	−0.051	0.41	0.07	3.60	0.31	2.72	0.57
	TZVP	B-1.05-23	0.09	−0.023	0.20	0.08	1.77	0.23	0.84	0.70
	LP	B-1.06-23	0.08	−0.021	0.17	0.07	1.18	0.17	0.72	0.30
B-LYP	6-31G**	B-1.05-23	0.13	0.005	0.22	0.16	2.12	0.24	1.02	0.86
	TZVP	B-0.87-23	0.07	0.020	0.11	0.08	1.36	0.21	0.65	0.50
	LP	B-0.79-23	0.01	−0.019	0.07	−0.02	1.07	0.19	0.64	0.25
	LP	B-0.90-23-1.25	0.04	0.033	0.00	0.08	0.77	0.20	0.24	0.33
TPSS	6-31G**	B-1.05-23	0.10	−0.035	0.24	0.07	2.16	0.24	1.36	0.55
	TZVP	B-0.98-35	0.09	−0.015	0.10	0.19	1.56	0.15	0.51	0.90
	LP	B-0.96-27	0.02	−0.005	0.02	0.05	0.60	0.21	0.20	0.19

$\Delta d_{AB}$  is average signed error of the distances between the molecular centers of mass ( $d_{AB}$ , Å). RMS is root mean square error of the optimized geometries with respect to the WFT reference (Å). Averages over all S22 complexes and over the hydrogen-bonded, dispersion-bonded, and mixed complexes, respectively.

<sup>a</sup>LP is 6-311++G(3df,3pd).

<sup>b</sup>For the optimized dispersion notation see Methods.

slightly repulsive (either the minimum on a very flat PES was not reached, or there is no minimum). As a result, RMS errors for mixed, but mainly for dispersion-bonded complexes, are unacceptably large (see Table 5). This applies also to the PBE functional, which is otherwise believed to perform relatively well for the vdW interactions. Apparently, the error cancellation that leads to very good performance for hydrogen-bonded complexes does not work either for complexes with mixed H-bond/dispersion character or for stacked complexes.

The empirical dispersion correction removes most of the earlier mentioned errors and remarkably improves the intermolecular geometries if at least a TZVP basis set is used. This applies especially to the dispersion-bonded complexes, where the RMS error is reduced up to 25 times (B-LYP/LP/B-0.90-23-1.25). In the case of the hydrogen-bonded complexes, the situation is less clear and we will discuss it in more detail in the following text. The issue is that additional dispersion energy causes overestimation of the hydrogen bonds and consequently their shortening. For those functionals, for which intermolecular distances are already too short even without dispersion (like PBE or TPSS), only further shortening occurs. On the other hand, for the pure

B-LYP functional the intermolecular distances are too long and addition of the dispersion energy brings significant improvement. We would like to emphasize that (in this respect) the behavior of the B-LYP functional is correct (in trend, but not in magnitude) and the behavior of the PBE and TPSS functionals is obviously wrong. The resulting error can not be attributed to the dispersion correction, and only partly to BSSE. A physically correct local or semilocal GGA XC functional should predict slightly underestimated hydrogen bonds and reflect in this way the missing dispersion interaction. For small hydrogen-bonded complexes this underestimation should be only small, whereas for larger complexes like DNA base pairs it is expected to be fairly large as the dispersion interaction plays a more important role here (large polarizable monomers). According to our calculations this trend is more or less noticeable for all the functionals tested (see also Table S1 in Supplementary Material).

Interestingly, in the case of the PBE/6-31G\*\*/B-1.15-23 combination the additional (always stabilizing) dispersion eventually lengthens the hydrogen bonds (on average). This is not an error. Due to large (scaled) vdW radii the minimum of the dispersion term alone is further than the minimum of the pure PBE intermolecular

Table 6. Single-Point Calculations on the Geometries Optimized Using DFT and DFT-D Methods.

X-C	Basis set <sup>a</sup>	Dispersion type <sup>b</sup>	Full set		H-bonded		Dispersion bonded		Mixed		$\Delta\Delta_{D-H}$
			SD	Avg	SD	Avg	SD	Avg	SD	Avg	
PBE	6-31G**	None	4.87	−0.26	2.46	−5.24	4.66	3.73	1.10	0.15	5.25
	TZVP	None	2.99	1.22	1.14	−1.28	3.50	3.56	1.02	1.05	3.32
	LP	None	2.85	1.55	1.34	−0.60	3.51	3.67	0.89	1.27	2.82
B-LYP	6-31G**	None	4.10	1.11	1.63	−2.82	4.31	4.42	1.27	1.27	3.94
	TZVP	None	2.39	2.70	1.45	1.36	3.07	4.23	1.23	2.28	1.62
	LP	None	2.61	3.23	1.35	2.21	3.70	4.72	1.18	2.55	1.33
TPSS	6-31G**	None	4.34	0.76	2.02	−3.46	4.44	4.31	1.11	0.92	4.13
	TZVP	None	2.87	2.01	1.21	0.02	3.70	4.02	1.08	1.70	3.08
	LP	None	1.92	1.86	1.29	0.47	2.19	3.19	0.99	1.74	2.72
B-3LYP	6-31G**	None	4.02	−1.14	1.38	−2.60	4.29	4.54	1.24	0.99	3.65
S-VWN	6-31G**	None	9.33	−8.01	10.07	−18.49	3.15	−2.97	1.87	−3.30	15.51
	TZVP	None	6.50	−5.37	7.00	−12.75	1.66	−1.63	1.44	−2.26	11.12
PBE	6-31G**	B-1.15-23	3.04	−2.44	2.47	−6.00	1.87	−0.73	0.74	−0.85	5.26
	TZVP	B-1.05-23	1.34	−1.04	1.16	−2.62	0.35	−0.14	0.64	−0.48	2.48
	LP	B-1.06-23	1.10	−0.70	1.30	−1.84	0.50	−0.14	0.25	−0.20	1.70
B-LYP	6-31G**	B-1.05-23	2.37	−1.09	1.54	−4.10	0.65	0.67	0.99	−0.09	4.77
	TZVP	B-0.87-23	1.20	0.08	0.72	−1.22	0.77	1.11	0.66	0.20	2.34
	LP	B-0.79-23	1.14	0.07	0.80	−1.12	0.87	1.06	0.31	0.14	2.17
	LP	B-0.90-23-1.25	0.55	−0.24	0.64	−0.71	0.40	−0.06	0.25	0.03	0.65
TPSS	6-31G**	B-1.05-23	2.44	−1.75	2.07	−4.79	0.38	−0.11	0.71	−0.59	4.67
	TZVP	B-0.98-35	1.09	−0.49	0.94	−1.73	0.40	0.18	0.68	−0.02	1.91
	LP	B-0.96-27	0.92	−0.59	1.22	−1.45	0.30	−0.26	0.28	−0.10	1.19

Single-point calculations are performed at the same level as optimization, without counterpoise correction. Comparison of the DFT and DFT-D results with high-level WFT data. Standard deviations (SD) and average signed errors (Avg) in kcal mol<sup>−1</sup>.

<sup>a</sup>LP is 6-311++G(3df,3pd).

<sup>b</sup>For the optimized dispersion notation see Methods.

potential alone, and the dispersion contribution to the gradient is therefore repulsive close to the PBE minimum (see also ref. 51 where similar effect is achieved by an additional repulsion potential). In this particular case, the dispersion interaction causes (desirable) lengthening of the intermolecular distances and lowers the RMS geometry error. In general, if the pure DFT interaction energy was underestimated and the intermolecular distance too short at the same time (this is not impossible), both the energy and the geometry could be improved in DFT-D. However, this is not the case, and although the PBE+D and TPSS+D geometries are improved compared with pure DFT, the interaction energies change for the worse (increase).

For all three functionals tested (B-LYP, PBE, TPSS), optimized geometries improve with the size of the basis set used. In the case of our smallest basis set, 6-31G\*\*, best geometries can be obtained with the B-LYP functional, mainly because PBE and TPSS predict too short hydrogen bonds (see earlier). Compared with pure B-LYP, the B-LYP/6-31G\*\*/B-1.05-23 method decreases the RMS for the dispersion-bonded complexes from 7.6 to 1.0 and improves slightly also the hydrogen-bonded com-

plexes. With the 6-31G\*\* optimized parameters similar results are expected also for other double- $\zeta$  basis sets like cc-pVDZ or SVP. The TZVP basis set brings very significant improvements and more balanced results for both geometries and interaction energies and we recommend using it instead of double- $\zeta$  basis sets whenever possible. The TPSS/TZVP/B-0.98-35 and also the B-LYP/TZVP/0.87-23 combinations appear to be better than the PBE-based DFT-D. Overall best results were achieved with large 6-311++G(3df,3pd) basis set (LP) containing diffuse functions (the B-LYP/LP/B-0.90-23-1.25 and TPSS/LP/B-0.98-35 combinations). The B-0.90-23-1.25 dispersion could be probably recommended also for the plane wave-based B-LYP calculations. Note that, regarding single point calculations (see Table 4), the two earlier mentioned combinations perform on average better than the much more demanding MP2 method with cc-pVTZ basis set.

Interaction energies (errors) calculated for the optimized DFT and DFT-D geometries are shown in Table 6. All energies are evaluated by the same method which was used for optimization. Note that results are not corrected for BSSE (for practical reasons



Table 7. Comparison of Errors for the Training Set S22 and the Validation Set of 58 Molecules.

Method <sup>a</sup>	Set	Full set		H-bonded		Dispersion bonded		$\Delta\Delta_{D-H}$
		SD	Avg	SD	Avg	SD	Avg	
PBE/TZVP/ B-1.05-23	S22	0.91	−0.28	0.56	−1.27	0.55	0.53	1.80
	S58	1.05	−0.29	0.63	−1.19	0.81	0.45	1.64
B-LYP/TZVP/ B-0.87-23	S22	1.00	0.47	0.72	−0.42	0.81	1.27	1.69
	S58	1.05	1.03	0.64	0.05	0.39	1.93	1.87
TPSS/6-31G**/ B-1.05-23	S22	1.85	−0.94	0.98	−3.34	0.56	0.62	3.96
	S58	2.22	−1.05	0.68	−3.51	0.74	0.86	4.38
TPSS/TZVP/ B-0.98-35	S22	0.84	0.53	0.56	−0.38	0.40	1.15	1.53
	S58	0.92	0.78	0.62	−0.07	0.49	1.49	1.57
TPSS/LP/ B-0.96-27	S22	0.38	−0.15	0.48	−0.39	0.30	−0.07	0.32
	S58	0.59	0.21	0.60	−0.02	0.55	0.54	0.45
B3-LYP/TZVP/ B-0.95-27	S22	0.80	0.05	0.42	−0.83	0.42	0.61	1.43
	S58	0.70	0.16	0.40	−0.52	0.37	0.71	1.23

Standard deviations (SD) and average signed errors (Avg) in kcal mol<sup>−1</sup>.

<sup>a</sup>LP is 6-311++G(3df,3pd). For the optimized dispersion notation see Methods.

our DFT+D parameters are intended to be used without the CP correction). Compared with Tables 3 and 4 performance of DFT-D for the dispersion-bonded complexes is in most cases (except the smallest basis sets) about the same or even slightly improved. This indicates that a parametrization procedure based on the interaction energies only (not considering gradients) can yield satisfying results and that the dispersion description used is sound close to the potential minima. However, results for hydrogen-bonded complexes are less satisfying as there is some increase in overbinding, especially, for small basis sets. As explained earlier, a large portion of this overbinding stems from the uncorrected BSSE, but there is also some contribution from the XC functional (the errors are larger for PBE and TPSS functionals than for B-LYP). Likely, further improvement could be achieved by using a still larger basis set, such as aQZ' for optimization (see Tables 3 and 4), provided that known incorrect behavior of the exchange functionals in the large reduced gradient regions does not introduce unexpectedly large additional errors.

Along with the hydrogen bonds overbinding there is also an increase in the  $\Delta\Delta_{D-H}$  shift due to optimization. Nevertheless, the final  $\Delta\Delta_{D-H}$  value is still about two times smaller than for the respective pure DFT method if a large enough basis set is used. What is more, lowering of the pure DFT  $\Delta\Delta_{D-H}$  after optimization is mostly due to incorrect geometries of dispersion-bonded complexes. Because the  $\Delta\Delta_{D-H}$  shift is most important when comparing relative stabilities of complexes with mixed dispersion and hydrogen-bonded character like different peptide conformations, only large basis sets with diffuse functions should be used for these purposes.

#### Validation on a Set of 58 vdW Complexes

The dispersion parameters optimized on the S22 training set were tested on a validation set of 58 molecular complexes. The validation data were collected from our previous papers,<sup>54–58</sup> and contain stacked and hydrogen-bonded nucleic acid base pairs in their

optimal (MP2/cc-pVTZ) geometries,<sup>54–56</sup> in the general geometries generated to map the whole potential energy surface<sup>57</sup> or in the geometries close to those found in natural DNA.<sup>58</sup> Reference energies were calculated similarly to the procedure described in Methods, but with smaller basis sets. For a list of the molecules included in the validation set S58 see Table S1 in Supplementary Material. Almost all XC functional/basis set/dispersion combinations parametrized on the training set were tested on the validation set. For the sake of simplicity we show just a small part of these results in Table 7 (remaining data are very similar). Comparison of the errors for the training set S22 and the validation set S58 shows no significant trend. Standard deviations for hydrogen-bonded and stacked systems (separately) are sometimes somewhat larger for the validation set, which could be explained also by the larger error expected in the validation reference WFT data (lower level of theory). In general, the errors for both training and validation sets are quite similar, despite the rather different composition of those two sets. This suggests that the S22 training set is quite representative regarding intermolecular interactions in DNA and dispersion parameters obtained are well transferable in this case. We believe that this will hold also for other kinds of molecules, like peptides (work on a similar comparison for peptides is under way).

#### Conclusions

Because of its empirical character, the dispersion correction can not challenge the accuracy and reliability of the high-level WFT QM data. Nevertheless, it is possible to improve the overall performance of the DFT-based calculations significantly while keeping the DFT-D treatment very simple and fast. This is certainly vital for the DFT-based calculations at least until more rigorous and rapid nonempirical approaches become available. Since we are primarily interested in the applications of DFT to large biological systems or molecular dynamics, the question of

computational efficiency becomes crucial. Therefore, we have focused mainly on the pure, nonhybrid, GGA, or meta-GGA XC-functionals for which significant speedups can be achieved by making use of the density fitting (or RI) approximation.

1. A simple empirical dispersion correction can significantly improve the description of the intermolecular interactions. DFT-D interaction energies and geometries are in very good agreement with high quality reference data for a large set of vdW complexes.
2. Damping of the empirical dispersion term is apparently the most critical point in DFT-D, since it concerns questions of double counting of correlation energy and range of a given correlation functional. The error introduced by inaccurate dispersion  $C^6$  coefficients is probably secondary.
3. As the basis set quality increases, the importance of the empirical correction becomes more and more obvious. At the same time, the results including the empirical dispersion term converge to the reference WFT data. This basis set convergence behavior supports our opinion that the dispersion correction to the current DFT functionals is physically sound and desirable.
4. In the case of hydrogen-bonded systems the BSSE for the basis sets of DZ quality is of the same direction and of almost the same size as the dispersion interaction. Therefore, double- $\zeta$  BSSE uncorrected DFT calculations give fairly reasonable accordance with experiment. The role of the dispersion energy in hydrogen bonding is thus frequently overlooked.
5. With respect to the previous two points, we recommend that the empirical dispersion correction be added to sufficiently large basis sets only. At least a TZVP basis set is recommended for reasonable results and triple- $\zeta$  quality basis set augmented with diffuse functions is recommended for the dispersion interaction in  $\pi$ -interacting systems, which are particularly sensitive to the basis set size.
6. The quality of the results is considerably limited by the quality of the XC functional used. The relatively large variation between individual XC functionals indicates that, regarding the intermolecular interactions, current GGA-based functionals have not yet reached their limits. It is encouraging that, in general, the more advanced functionals (hybrid functionals, meta-GGA functionals) give better results than LDA or simple GGA approximations when combined with empirical dispersion. When modifying just a global scaling factor  $s_R$  and the steepness of the damping function (exponent), the best results were obtained with the TPSS meta-GGA functional.

## Acknowledgment

This work was part of the research project Z4 055 0506.

## References

1. Dąbkowska, I.; Gonzalez, H. V.; Jurečka, P.; Hobza, P. *J Phys Chem A* 2005, 109, 1131.
2. Kubař, T.; Hanus, M.; Ryjáček, F.; Hobza, P. *Chem Eur J* 2006, 12, 280.
3. Řeha, D.; Hocek, M.; Hobza, P. *Chem Eur J* 2006, 12, 3587.
4. Vondrášek, J.; Bendová, L.; Klusák, V.; Hobza, P. *J Am Chem Soc* 2005, 127, 2615.
5. Liu, H.; Elstner, M.; Kaxiras, E.; Frauenheim, T.; Hermans, J.; Yang, W. *Proteins* 2001, 44, 484.
6. Guallar, V.; Borrelli, K. W. *Proc Natl Acad Sci USA* 2005, 102, 3954.
7. Alber, F.; Kuonen, O.; Scapozza, L.; Folkers, G.; Carloni, P. *Proteins* 1998, 31, 453.
8. Klusák, V.; Havlas, Z.; Rulíšek, L.; Vondrášek, J.; Svatoš, A. *Chem Biol* 2003, 10, 331.
9. Černý, J.; Hobza, P. *Phys Chem Chem Phys* 2005, 7, 1624.
10. Bashford, D.; Chothia, C.; Lesk, A. M. *J Mol Biol* 1987, 196, 199.
11. Gazit, E. *FASEB J* 2002, 16, 77.
12. Kristyán, S.; Pulay, P. *Chem Phys Lett* 1994, 229, 175.
13. Hobza, P.; Šponer, J.; Reschel, T. *J Comput Chem* 1995, 16, 1315.
14. Boese, A. D.; Handy, N. C. *J Chem Phys* 2002, 116, 9559.
15. Sim, F.; St-Amant, A.; Papai, I.; Salahub, D. R. *J Am Chem Soc* 1992, 114, 4391.
16. Sirois, S.; Proynov, E. I.; Nguyen, D. T.; Salahub, D. R. *J Chem Phys* 1997, 107, 6770.
17. Zhang, Y.; Pan, W.; Yang, W. *J Chem Phys* 1997, 107, 7921.
18. Řeha, D.; Valdés, H.; Vondrášek, J.; Hobza, P.; Abu-Riziq, A.; Crews, B.; de Vries, M. S. *Chem Eur J* 2005, 11, 6803.
19. Hesselmann, A.; Jansen, G. *Chem Phys Lett* 2003, 367, 778.
20. Misquitta, A. J.; Jeziorski, B.; Szalewicz, K. *Phys Rev Lett* 2003, 91, 033201.
21. Osinga, V. P.; vanGisbergen, S. J. A.; Snijders, J. G.; Baerends, E. J. *J Chem Phys* 1997, 106, 5091.
22. Adamovic, I.; Gordon, M. S. *Mol Phys* 2005, 103, 379.
23. Andersson, Y.; Langreth, D. C.; Lundqvist, B. I. *Phys Rev Lett* 1996, 76, 102.
24. Dobson, J. F. *Int J Quantum Chem* 1998, 69, 615.
25. Sato, T.; Tsuneda, T.; Hirao, K. *J Chem Phys* 2005, 123, 104307.
26. Rapcewicz, K.; Ashcroft, N. W. *Phys Rev B* 1991, 44, 4032.
27. Wesolowski, T. A.; Tran, F. *J Chem Phys* 2003, 118, 2072.
28. Becke, A. D.; Johnson, E. R. *J Chem Phys* 2005, 122, 154104.
29. Walsh, T. R. *Phys Chem Chem Phys* 2005, 7, 443.
30. Kohn, W.; Meir, Y.; Makarov, D. E. *Phys Rev Lett* 1998, 80, 4153.
31. Dion, M.; Rydberg, H.; Schröder, E.; Langreth, D. C.; Lundqvist, B. I. *Phys Rev Lett* 2004, 92, 246401.
32. Lacks, D. J.; Gordon, R. G. *Phys Rev A* 1993, 47, 4681.
33. Adamo, C.; Barone, V. *J Chem Phys* 1998, 108, 664.
34. Xu, X.; Goddard, W. A., III. *Proc Natl Acad Sci USA* 2004, 101, 2673.
35. Kurita, N.; Inoue, H.; Sekino, H. *Chem Phys Lett* 2003, 370, 161.
36. Zhao, Y.; Truhlar, D. G. *Phys Chem Chem Phys* 2005, 7, 2701.
37. (a) Gutowski, M.; Jordan, K. D.; Skurski, P. *J Phys Chem A* 1998, 102, 2624; (b) Svozil, D.; Frigato, T.; Havlas, Z.; Jungwirth, P. *Phys Chem Chem Phys* 2005, 7, 840.
38. Hepburn, J.; Scoles, G. *Chem Phys Lett* 1975, 36, 451.
39. Ahlrichs, R.; Penco, R.; Scoles, G. *Chem Phys* 1977, 19, 119.
40. Hobza, P.; Mulder, F.; Sandorfy, C. *J Am Chem Soc* 1981, 103, 1360.
41. Hobza, P.; Mulder, F.; Sandorfy, C. *J Am Chem Soc* 1982, 104, 925.
42. Hobza, P.; Sandorfy, C. *Can J Chem* 1984, 62, 606.
43. Hobza, P.; Sandorfy, C. *J Am Chem Soc* 1987, 109, 1302.
44. Meijer, E. J.; Sprik, M. *J Chem Phys* 1996, 105, 8684.
45. Mooij, W. T. M.; van Duijneveldt, F. B.; van Duijneveldt-van de Rijdt, J. G. C. M.; van Eijck, B. P. *J Phys Chem A* 1999, 103, 9872.
46. Wu, X.; Vargas, M. C.; Nayak, S.; Lotrich, V.; Scoles, G. *J Chem Phys* 2001, 115, 8748.
47. Wu, Q.; Yang, W. *J Chem Phys* 2002, 116, 515.

48. Zimmerli, U.; Parrinello, M.; Koumotsakos, P. *J Chem Phys* 2004, 120, 2693.
49. Grimme, S. *J Comput Chem* 2004, 25, 1463.
50. Elstner, M.; Hobza, P.; Frauenheim, T.; Suhai, S.; Kaxiras, E. *J Chem Phys* 2001, 114, 5149.
51. Zhechkov, L.; Heine, T.; Patchkowskii, S.; Seifert, G.; Duarte, H. A. *J Chem Theory Comput* 2005, 1, 841.
52. Lilienfeld, O. A.; Tavernelli, I.; Rothlisberger, U. *Phys Rev Lett* 2004, 93, 153004.
53. Pérez-Jordá, J. M.; Becke, A. D. *Chem Phys Lett* 1995, 233, 134.
54. Jurečka, P.; Šponer, J.; Černý, J.; Hobza, P. *Phys Chem Chem Phys* 2006, 8, 1985.
55. Jurečka, P.; Hobza, P. *J Am Chem Soc* 2003, 125, 15608.
56. Šponer, J.; Jurečka, P.; Hobza, P. *J Am Chem Soc* 2004, 126, 10142.
57. Jurečka, P.; Šponer, J.; Hobza, P. *J Phys Chem B* 2004, 108, 5466.
58. Šponer, J.; Jurečka, P.; Marchan, I.; Luque, F. J.; Orozco, M.; Hobza, P. *Chem Eur J* 2006, 12, 2854.
59. Serra, S.; Iarlori, S.; Tosatti, S.; Scandolo, S.; Santoro, G. *Chem Phys Lett* 2000, 331, 339.
60. Becke, A. D. *Phys Rev A* 1988, 38, 3098.
61. Lee, C.; Yang, W.; Parr, R. G. *Phys Rev B* 1988, 37, 785.
62. Bondi, A. *J Chem Phys* 1964, 68, 441.
63. Halgren, T. A. *J Am Chem Soc* 1992, 114, 7827.
64. Slater, J. C. *Phys Rev* 1951, 81, 385.
65. Vosko, S. H.; Wilk, L.; Nusair, M. *Can J Phys* 1980, 58, 1200.
66. Becke, A. D. *J Chem Phys* 1993, 98, 5648.
67. Becke, A. D. *J Chem Phys* 1993, 98, 1372.
68. Perdew, J. P.; Burke, K.; Ernzerhof, M. *Phys Rev Lett* 1996, 77, 3865.
69. Tao, J.; Perdew, J. P.; Staroverov, V. N.; Scuseria, G. E. *Phys Rev Lett* 2003, 91, 146401.
70. Staroverov, V. N.; Scuseria, G. E.; Tao, J.; Perdew, J. P. *J Chem Phys* 2003, 119, 12129.
71. Ahlrichs, R.; Bär, M.; Häser, M.; Horn, H.; Kölmel, C. *Chem Phys Lett* 1989, 162, 165.
72. (a) Eichkorn, K.; Treutler, O.; Öhm, H.; Häser, M.; Ahlrichs, R. *Chem Phys Lett* 1995, 240, 283; (b) Eichkorn, K.; Treutler, O.; Öhm, H.; Häser, M.; Ahlrichs, R. *Chem Phys Lett* 1995, 242, 652.
73. Eichkorn, K.; Weigend, F.; Treutler, O.; Ahlrichs, R. *Theor Chem Acc* 1997, 97, 119.
74. Weigend, F.; Häser, M.; Patzelt, H.; Ahlrichs, R. *Chem Phys Lett* 1998, 294, 143.
75. Weigend, F.; Köhn, A.; Hättig, C. *J Chem Phys* 2002, 116, 3175.
76. Weigend, F.; Häser, M. *Theor Chem Acc* 1997, 97, 331.
77. Hariharan, P. C.; Pople, J. A. *Theor Chim Acta* 1973, 28, 213.
78. (a) Krishnan, R.; Binkley, J. S.; Seeger, R.; Pople, J. A. *J Chem Phys* 1980, 72, 650; (b) Clark, T.; Chandrasekhar, J.; Spitznagel, G. W.; von R Schleyer, P. *J Comput Chem* 1983, 4, 294; (c) Gill, P. M. W.; Johnson, B. G.; Pople, J. A.; Frisch, M. J. *Chem Phys Lett* 1992, 197, 499; (d) Frisch, M. J.; Pople, J. A.; Binkley, J. S. *J Chem Phys* 1984, 80, 3265.
79. Schäfer, A.; Huber, C.; Ahlrichs, R. *J Chem Phys* 1994, 100, 5829.
80. (a) Dunning, T. H., Jr. *J Chem Phys* 1989, 90, 1007; (b) Kendall, R. A.; Dunning, T. H., Jr.; Harrison, R. J. *J Chem Phys* 1992, 96, 6796.
81. Dąbkowska, I.; Jurečka, P.; Hobza, P. *J Chem Phys* 2005, 122, 204322.
82. Halkier, A.; Helgaker, T.; Jørgensen, P.; Klopper, W.; Koch, H.; Olsen, J. *Chem Phys Lett* 1999, 302, 437.
83. Halkier, A.; Helgaker, T.; Jørgensen, P.; Klopper, W.; Koch, H.; Olsen, J.; Wilson, A. K. *Chem Phys Lett* 1998, 286, 243.
84. Boese, A. D.; Martin, J. M.; Handy, N. C. *J Chem Phys* 2003, 119, 3005.
85. Boys, S. F.; Bernardi, F. *Mol Phys* 1970, 19, 553.