

# Extrapolating DFT towards the complete basis set limit: Lessons from the PBE family of functionals

Peter Kraus\*

*School of Molecular and Life Sciences,*

*Curtin University,*

*GPO Box U1987, Perth 6845, WA, Australia*

E-mail: peter.kraus@curtin.edu.au

## Abstract

Extrapolation of density functional theory results from 2- and 3- $\zeta$  calculations is a promising method for extracting higher accuracy data from calculations of systems at the affordability limit. In this work, I present formulas for the determination of extrapolation parameters, that account for the make-up of the density functional approximation. The formulas are fitted to reproduce the complete basis set limit energies of PBE and related density functional approximations, using a set of 30 singlet diatomics. Their performance is extensively evaluated using standard benchmark datasets. The current formulas are shown to be transferable outside the PBE-family of functional approximations, with the resulting extrapolation parameters outperforming the previous, empirical values. A good performance of [2,3]- $\zeta$  extrapolations for interaction energies of systems with significant non-covalent character is confirmed, and holds even in systems of  $\sim 100$  atoms in size.

# 1 Introduction

Basis set extrapolation methods allow computational chemists to reliably approximate results of expensive calculations in larger basis sets by combining results obtained at the same level of theory from two or more cheaper calculations in smaller basis sets. This is particularly important in post-Hartree-Fock (HF) wavefunction theory, where the calculation of electronic correlation dominates the computational cost. It has long been known that the correlation energy ( $E^{\text{corr}}$ ) approaches its infinite basis set limit ( $E_{\infty}^{\text{corr}}$ ) as a function of a cubic power law of the basis set size,

$$E_{\infty}^{\text{corr}} = E_X^{\text{corr}} - AX^{-\alpha} \quad (1)$$

where the extrapolation parameter  $\alpha = 3$ ,  $X$  is the cardinal number of a finite correlation-consistent basis set, and  $A$  is a system-dependent parameter fitted to the  $E_X^{\text{corr}}$  data.<sup>1</sup> Notably, while this cubic power law is correct for all values of  $X$ s, other empirically determined values of  $\alpha$  may provide better results for  $X \leq 3$ .<sup>2</sup>

Extrapolation methods are not widely used in density functional theory (DFT) calculations. This may be due to several reasons: DFT calculations are computationally inexpensive compared to correlated WFT, reducing the applicability of such methods; there are many different empirical density functional approximations (DFAs) with each requiring its own  $\alpha$ , meaning such approaches are seen as less general; reference datasets of complete basis set limit energies from DFT are not as common as their WFT counterparts, hindering the development of such methods; and even the true basis set limit convergence of HF was until recently a matter of debate, putting any DFT extrapolation attempts on an uncertain theoretical basis.

In a previous work,<sup>3</sup> I have shown that the extrapolation parameters  $\alpha$  determined for a set of DFAs using numerical complete basis set calculations are surprisingly transferrable to other DFAs. However, the values of  $\alpha$  are strongly dependent on the basis set family. For

instance, extrapolated def2-[st]zvpd non-covalent interaction energies generally outperform their def2-qzvpd counterparts at 80% of the computational cost, which holds true even for double-hybrid density functional approximations (DHDFAs) as long as the correlated perturbation theory (PT) component is extrapolated appropriately. The best performing extrapolation function for DFT is the exponential-square root function (expsqrt), first proposed by Jensen:<sup>4</sup>

$$E_{\infty}^{\text{fctI}} = E_X^{\text{fctI}} - Ae^{-\alpha\sqrt{X}} \quad (2)$$

This function has since been proven to be the correct form for energy convergence in both HF and Kohn-Sham DFT with Gaussian basis sets.<sup>5</sup> However, several issues highlighted in the previous work remain. The determination of the values of the extrapolation parameters  $\alpha$ , as described in Ref. 3, is highly method-specific due to the introduced averaging. Additionally, the values of  $\alpha$  clearly scale with the inclusion of exact (HF) exchange in hybrid DFAs, but the introduction of formulas which correct the  $\alpha$  by the percentage of exact exchange did not improve the results. Finally, the influence of the reduced correlation component in double hybrid exchange–correlation functionals upon introduction of PT correlation has not been investigated.

In the current work, I address the three above issues by systematically investigating of PBE-like DFAs. Compared to the set of 9 diatomics used in previous work,<sup>3</sup> a larger dataset of complete basis set energies from finite element calculations is used, comprising 30 singlet diatomic species, and including cations and anions. This dataset is evaluated for PBE and related functionals with varying fractions of HF exchange and PT correlation. I investigate trends over a large number of basis set families, including most of the modern basis set families used with DFT and WFT. This leads to a formula which is specific for each basis set family, but universal for all DFAs, taking into account all components of the DFA recipe. The extrapolation parameters  $\alpha$  obtained from these formulas are then thoroughly examined using several benchmark databases.

## 2 Computational methods

Analogously to previous work, I use the following nomenclature for the total energy of a DFA ( $E^{\text{DFA}}$ ) and its components:

$$E^{\text{DFA}} = E^{\text{fctl}} + \Delta E^{\text{dh}} + \Delta E^{\text{disp}} + \Delta E^{\text{nl}} \quad (3)$$

The first component is the self-consistent-field (SCF) energy of the exchange-correlation functional (fctl), the second component is the double-hybrid (dh) correction, the third term corresponds to dispersion corrections, such as Grimme’s -D3,<sup>6</sup> and the fourth term is a non-local correlation term, such as in the VV10 DFA,<sup>7</sup> or generalised as -NL correction.<sup>8</sup>

The first two terms can be further split into the following components:

$$E^{\text{fctl}} = (1 - a_x)E^x + a_x E^{\text{HFx}} + (1 - a_c)E^c \quad (4)$$

$$\Delta E^{\text{dh}} = a_c E^{\text{PTc}} \quad (5)$$

Note that the coefficients  $a_x$  and  $a_c$  correspond to the fraction of HF exchange (HFx) and PT correlation (PTc) in the DFA recipe. While  $a_x$  is constrained to  $\leq 1$ ,<sup>9</sup> no such rigorous constraint on  $a_c$  exists.<sup>10</sup> In practice, most common DHDFAs obey both constraints, with the notable exception of the spin-component-scaled double hybrids (DSD family).<sup>11</sup> For simplicity, in the following discussion it is assumed both constraints hold. The fit is always performed using the exact values of the coefficients in front of  $E^x$  and  $E^c$  in Eq. (4), even if the latter is not equal to  $(1 - a_c)$ . Finally, while in principle the  $\Delta E^{\text{nl}}$  term contributes to the SCF cycle, it has been conclusively shown to produce nearly identical results when applied as a post-SCF correction.<sup>12</sup> The post-SCF treatment is applied throughout this work.

The  $\{a_x, a_c\}$  parameter space explored in this work consists of all 20 combinations of  $a_x \in$

{0.1, 0.3, 0.5, 0.7, 0.9} with  $a_c \in \{0.0, 0.2, 0.4, 0.6\}$  as well as the pure variant (i.e.  $a_x = a_c = 0.0$ ). The parameter sets are chosen to cover the range used by most of the common single and double hybrids. The exchange and correlation functionals of Perdew, Burke, and Ernzerhof (PBE)<sup>13</sup> are used, as they represent a “non-empirical” DFA, even though it could be argued parameter-counting is a questionable metric of empiricity.<sup>14</sup> The generalization of the results obtained from this family of PBE-like functionals is discussed below.

The complete basis set energies of the 30 diatomic molecules listed in Table 1 are calculated using the diatomic component of the finite element code HelfEM,<sup>15</sup> using parameters determined with the `diatomic_cbasis` tool. Three anionic diatomics were difficult to converge using the following combinations of  $a_x$  and  $a_c$ :  $\text{SF}^-$  did not converge with  $a_x = 0.1$  and  $a_c = 0.6$ , and  $\text{OF}^-$  as well as  $\text{OH}^-$  failed to converge with the pure functional, as well as for all four values of  $a_c$  with  $a_x = 0.1$ . This is attributed to the diffuse nature of the anion exacerbated by the self-interaction error in PBE. Neither increasing the numerical integration radius from 40 Å to 120 Å, nor increasing the grid parameters led to convergence.

All other calculations were performed with Psi4 version 1.4a2.dev923,<sup>16</sup> or a development version of the program including basis set extrapolation routines for DFT implemented as part of previous<sup>3</sup> and current work. Calculations of the 30 diatomic molecules were carried out using the PK-supermatrix SCF algorithm to avoid issues with missing auxiliary basis sets, with a tightened energy convergence ( $10^{-10}$  Eh), and a large DFT quadrature with 150 radial and 974 angular points. The basis set families investigated in this work are listed in Table 2, significantly expanding on previous work where only four basis set families were compared. All basis sets were used as included with Psi4, or downloaded from the Basis Set Exchange.<sup>17,18</sup> It is notable that the energy obtained with 4- $\zeta$  variants of cc-pwcvXz and aug-cc-pwcvXz is in many cases lower than that with the 5- $\zeta$  variants; this non-convergence does not happen with other basis sets.

Results of the calculations of the diet100 variant<sup>37</sup> of the GMTKN55 database<sup>38</sup> are re-used from previous work, comprising 4 basis set families and 15 DFAs (see Table S1 for a detailed list); the WTMAD-2 criterion is used as figure of merit.<sup>3</sup> Additional calculations of

Table 1: Set of 30 diatomic molecules, including their bond lengths ( $R$ ), the  $\ell_{\max}$  of the basis sets, and the number of radial elements ( $n_{\text{elem}}$ ) used in the finite element calculations. The numerical integration radius was 40 Å in all cases, the number of nodes in each element was 15.

<b>Diatomic</b>	<b>R (Å)</b>	<b><math>\ell_{\max}</math></b>		<b><math>n_{\text{elem}}</math></b>
CH <sup>+</sup>	1.13085	15	11	3
OH <sup>+</sup>	1.02890	17	13	3
OH <sup>-</sup>	0.94246	15	13	3
FH	0.91696	17	13	5
C <sub>2</sub>	1.24780	17	11	3
CN <sup>+</sup>	1.17290	17	13	3
CN <sup>-</sup>	1.17160	17	13	3
N <sub>2</sub>	1.09434	17	11	3
NO <sup>+</sup>	1.06206	17	13	3
CO	1.12821	17	13	5
CF <sup>+</sup>	1.22875	19	15	5
OF <sup>-</sup>	1.49228	21	17	5
F <sub>2</sub>	1.41184	21	15	5
SiH <sup>+</sup>	1.50410	25	19	5
SH <sup>-</sup>	1.34993	27	21	5
HCl	1.29119	27	21	5
CP <sup>-</sup>	1.58753	27	21	5
CS	1.53442	27	21	7
SiN <sup>-</sup>	1.55578	27	21	5
NP	1.49085	27	21	5
SN <sup>+</sup>	1.44000	27	21	5
SiO	1.50974	25	19	5
PO <sup>+</sup>	1.41900	25	19	5
SF <sup>-</sup>	1.70395	29	23	7
ClF	1.66162	29	23	7
SiS	1.93000	31	23	7
P <sub>2</sub>	1.89340	31	23	7
PS <sup>+</sup>	1.87200	31	23	7
SCl <sup>-</sup>	2.14846	33	25	7
Cl <sub>2</sub>	2.04262	33	25	7

Table 2: Basis set families, the  $\zeta$  range used in extrapolation, and whether the basis set includes tighter core  $d$ -functions (Core) and/or diffuse augmentation (Diffuse).

Family	$\zeta$ range	Core?	Diffuse?	Reference
cc-pvXz <sup>†</sup>	2 – 6	×	×	Dunning <sup>19–22</sup>
cc-pwcvXz <sup>†</sup>	2 – 5	✓	×	Dunning <sup>19–21,23</sup>
aug-cc-pvXz <sup>†</sup>	2 – 6	×	✓	Dunning <sup>19–22</sup>
aug-cc-pwcvXz <sup>†</sup>	2 – 5	✓	✓	Dunning <sup>19–21,23</sup>
def2-Xzvp <sup>†</sup>	2 – 4	×	×	Ahlrichs <sup>24</sup>
def2-Xzvpp	2 – 4	×	×	Ahlrichs <sup>24</sup>
def2-Xzvpd <sup>†</sup>	2 – 4	×	✓	Ahlrichs <sup>24,25</sup>
def2-Xzvppd	2 – 4	×	✓	Ahlrichs <sup>24,25</sup>
pc- $N$	1 – 5*	×	×	Jensen <sup>26</sup>
pcseg- $N$ <sup>†</sup>	1 – 5*	×	×	Jensen <sup>27</sup>
aug-pc- $N$	1 – 5*	×	✓	Jensen <sup>28</sup>
aug-pcseg- $N$ <sup>†</sup>	1 – 5*	×	✓	Jensen <sup>27</sup>
Xzapa-nr	2 – 6	×	✓	Ranasinghe <sup>29</sup>
Xzapa-nr-cv	2 – 6	✓	✓	Ranasinghe <sup>29,30</sup>
jorge-Xzp	2 – 6	×	×	Jorge <sup>31–34</sup>
jorge-aXzp	2 – 5	×	✓	Jorge <sup>31,35,36</sup>

\* Note that for (aug-)pc- $N$  and (aug-)pcseg- $N$  families,  $N \approx X + 1$ .

<sup>†</sup> Basis set families used in the ASCDB benchmark.

the ASCDB database<sup>39</sup> were carried out using the standard density-fitted SCF algorithm, and the same (150, 974) grid, comprising 8 basis set families (marked with <sup>†</sup> in Table 2) and 3 PBE-like DFAs (PBE-D3(BJ),<sup>13,40</sup> PBE0-D3(BJ),<sup>40,41</sup> and PBE0DH-D3(BJ)<sup>42,43</sup>). Cases including elements for which basis sets are not defined within a given basis set family are excluded from the dataset of the corresponding basis set family. The ASCDB database allows for decomposition of the overall error metric into components of non-covalent interactions (NCI), thermochemistry (THCH), non-local effects (NLE), and unbiased calculations (UNBC),<sup>39</sup> allowing a more detailed insight into the performance of an extrapolation method than with the diet GMTKN55 subsets.<sup>37</sup>

The performance of basis set extrapolations for optimizing geometries of van der Waals complexes was evaluated using the NCDT database,<sup>44</sup> in its revised form,<sup>45</sup> using a development version of Psi4 with a tightened set of convergence criteria ( $\Delta E < 1 \mu\text{Eh}$ ,  $\max F < 15 \mu\text{Eh}/a_0$ ,  $\text{RMS } F < 10 \mu\text{Eh}/a_0$ ,  $\max d < 0.6 \text{ ma}_0$ ,  $\text{RMS } d < 0.4 \text{ ma}_0$ , i.e. interfrag-tight) and an unpruned (150, 974) quadrature grid. The geometries were optimized using the def2-

Xzvpd basis set family, used with revPBE-D3(BJ), M062X-D3, B97M-D3(BJ), B3LYP-D3(BJ), PBE0-D3, and  $\omega$ B97X-D3(BJ), selected based on their performance in GMTKN55,<sup>46</sup> and previous NCDT results.<sup>47</sup>

To evaluate the performance of the basis set extrapolations in their most likely use case, that is in large systems with a significant non-covalent character, the interaction energies of the geometries from the L7 dataset<sup>48</sup> were evaluated using M062X-D3/def2-Xzvpd and B97M-V/pcseg-*N*. This benchmark was designed to capture the dominant binding motifs in biological chemistry,<sup>48</sup> comprising 7 systems ranging in size between 48 and 112 atoms. The revised reference energies using correlation corrections obtained from DLPNO-CCSD( $T_0$ )/cc-pv[dt]z calculations<sup>49</sup> are used instead of the original values.<sup>48</sup> This revised set is in a good agreement with other recent reference data,<sup>50</sup> however even with significantly different reference values<sup>51</sup> the trends reported here are not affected. The calculations for the L7 dataset are performed with Psi4 version 1.4a2.dev923, using the default unpruned (75, 302) grid and the frozen core approximation for all double hybrid functionals.

The values of  $A$  and  $\alpha$  in Eq. (2) for each of the 30 diatomics are obtained using a non-linear least squares fit of Eq. (6) using the `curve_fit` function of the Python library `scipy.optimize`.<sup>52</sup>

$$\ln(E_X^{\text{fctl}} - E_\infty^{\text{fctl}}) = \ln(A) - \alpha\sqrt{X} \quad (6)$$

The  $E_X^{\text{fctl}}$  values correspond to the energies from individual basis sets in each family listed in Table 2, with  $X \approx \zeta$ . The complete basis set limit value  $E_\infty^{\text{fctl}}$  is the value obtained from HelFEM. The extrapolation parameters  $\alpha$  are then averaged based on the system charge, obtaining the means ( $\bar{\alpha}$ ) and medians ( $\tilde{\alpha}$ ) for the cationic ( $\alpha^+$ ), neutral ( $\alpha^0$ ), anionic ( $\alpha^-$ ), and all ( $\alpha$ ) diatomic species. Finally, a linear fit to the mean values  $\bar{\alpha}$  of all combinations of  $\{a_x, a_c\}$  is performed, using Eq. (7), obtaining the intercept  $\alpha_0$ , and two linear coefficients  $\alpha_x$  and  $\alpha_c$ .



$$\bar{\alpha} = \alpha_0 + \alpha_x a_x + \alpha_c a_c \quad (7)$$

## 3 Results

### 3.1 Extrapolation formulas

The results for each basis set family at all  $a_x$  with  $a_c = 0$  are presented in Fig. 1. The cationic outlier present in all basis sets is  $\text{OH}^+$ ; the anionic outlier most visible in def2-Xzvp is  $\text{OH}^-$ . With the exception of the Jensen basis set families, all  $\bar{\alpha}$ s show a positive correlation with  $a_x$ . The inclusion of HF exchange into the PBE recipe actually increases the overall self-interaction error (SIE) for hydrogenic atoms, as for the pure functional the exchange and correlation contributions to the SIE cancel each other out, and at non-zero  $a_x$  this error cancellation is lost.<sup>53</sup> However, for diatomics, the inclusion of HF exchange considerably improves the performance of the DFA in the SIE4x4 benchmark.<sup>53</sup> In any case, with a fraction of HF exchange, the electronic density is less delocalised. As the  $\exp(-\alpha\sqrt{X})$  term accounts for the error in the convergence of energy with respect to basis set size  $X$ ,<sup>54</sup> it is not surprising that for more compact densities the errors vanish quicker. This may also explain why no such correlation holds for  $a_c$ .

The augmentation of the basis sets by diffuse functions leads to a clear shift of the extrapolation parameters  $\alpha$  for anionic species (•) to values in line with the neutral and cationic species. This trend has been discussed in the context of basis set extrapolations previously.<sup>3,55</sup> Here, I wish to reiterate the need for the use of diffuse functions whenever anionic systems are involved. The data show no statistically significant difference between  $\bar{\alpha}^+$  and  $\bar{\alpha}^0$ .

The addition of tighter core  $d$ -electron functions pushes the extrapolation parameters  $\alpha$  for many of the diatomics upwards towards  $\sim 6$ . This trend is most apparent for aug-cc-pvXz  $\rightarrow$  aug-cc-pwcvXz augmentation, however it is also present in cc-pvXz  $\rightarrow$  cc-pwcvXz, as well as for Xzapa-nr  $\rightarrow$  Xzapa-nr-cv. It should be noted that the value of  $\alpha$  that is optimal for predicting

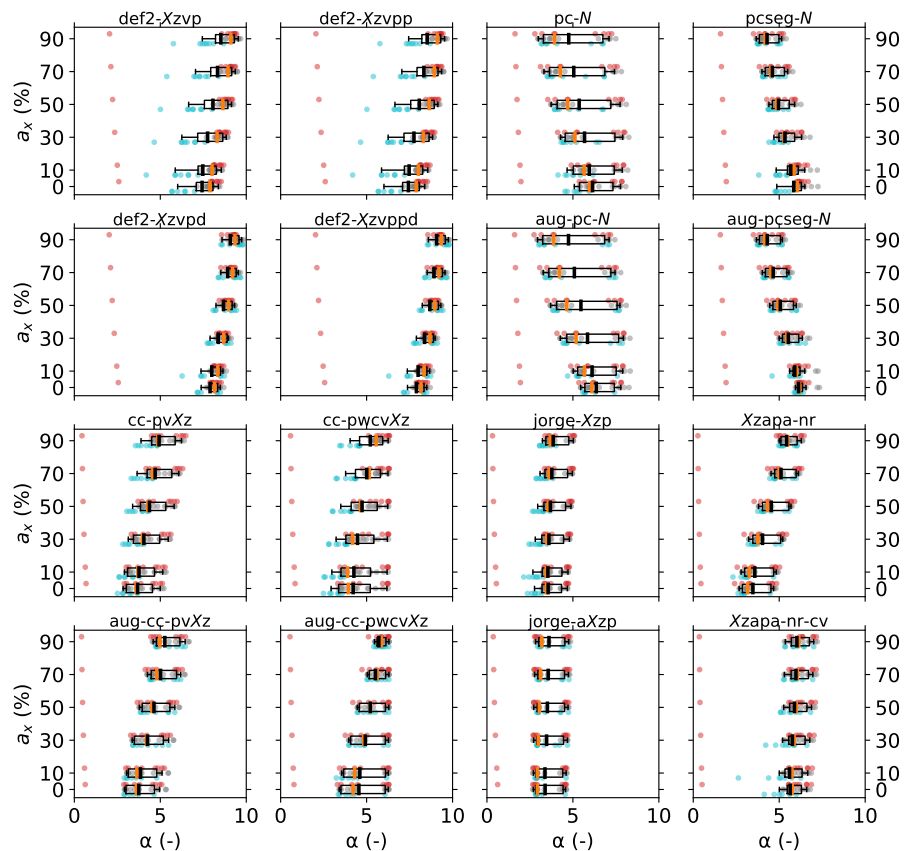


Figure 1: Extrapolation parameters  $\alpha$  obtained with various basis set families at  $a_c = 0$  and  $a_x$  as shown on the vertical axis. Colour coding corresponds to cationic ( $\bullet$ ), neutral ( $\bullet$ ), and anionic ( $\bullet$ ) species, with the box plots showing the overall median ( $\tilde{\alpha}$ ,  $\color{brown}{|}$ ), mean ( $\bar{\alpha}$ ,  $\color{black}{|}$ ), quartiles, and the 5th and 95th percentile.

160 the error of the energy expectation value approaches  $\pi\sqrt{3} \approx 5.44$  from below as basis set size  
 161 approaches completeness.<sup>54</sup>

Table 3: Exchange and correlation dependence of  $\bar{\alpha}$  for each basis set family studied. Coefficients are as per Eq. (7).

Family	$\alpha_0$	$\alpha_x$	$\alpha_c$
cc-pvXz	3.622	1.511	0.005
cc-pwcvXz	4.157	1.192	-0.048
aug-cc-pvXz	3.676	1.887	0.139
aug-cc-pwcvXz	4.485	1.445	0.085
pc- <i>N</i>	6.172	-1.623	0.183
pcseg- <i>N</i>	5.883	-1.825	0.227
aug-pc- <i>N</i>	6.390	-1.874	0.260
aug-pcseg- <i>N</i>	6.166	-2.137	0.296
def2-Xzvp	7.406	1.266	-0.046
def2-Xzvpp	7.408	1.267	-0.046
def2-Xzvdp	7.925	1.370	0.101
def2-Xzvppd	7.927	1.371	0.101
jorge-Xzp	3.531	0.338	-0.011
jorge-aXzp	3.386	0.245	0.033
Xzapa-nr	3.306	2.525	-0.040
Xzapa-nr-cv	5.618	0.490	-0.016

## 162 3.2 ASCDB database

163 The mean unsigned errors (MUE) in the ASCDB database are shown for three widely used  
 164 PBE-like functionals in Fig. 2. As expected, the pure functional PBE-D3(BJ) is outperformed by  
 165 the single hybrid PBE0-D3(BJ). However, the double hybrid PBE0DH-D3(BJ) does not provide a  
 166 systematic improvement over the single hybrid. Diffusely augmented basis set families outperform  
 167 the corresponding non-augmented families in all cases, which is not surprising given the inclusion  
 168 of non-covalent interactions in the ASCDB database. The Jensen basis set families (pcseg-*N*  
 169 and aug-pcseg-*N*) systematically outperform both the original Dunning (cc-pvXz and aug-cc-  
 170 pvXz) as well as the Ahlrichs (def2-Xzvp and def2-Xzvdp) sets. The weighted core-valence  
 171 correlation consistent Dunning sets including tight *d*-functions (cc-pwcvXz and aug-cc-pwcvXz)  
 172 offer comparable performance to the Jensen sets, albeit at a higher computational cost.

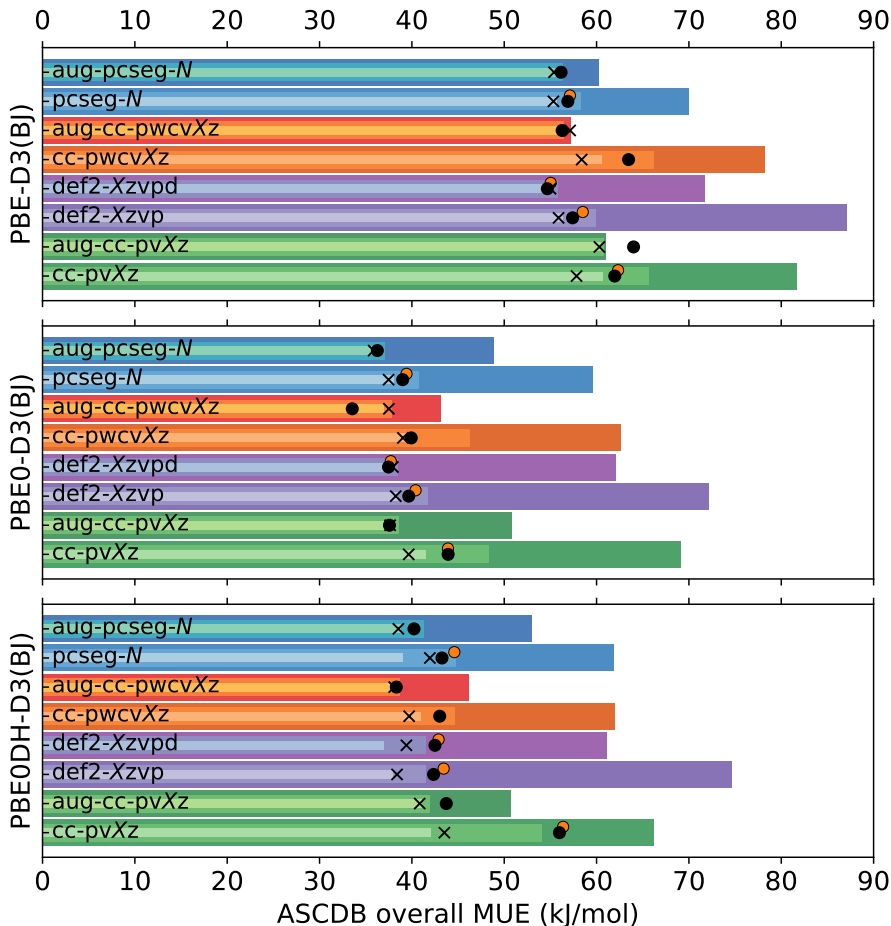


Figure 2: The mean unsigned errors of various functional and basis set combinations in the ASCDB database. Calculations with 2-, 3-, and 4- $\zeta$  basis sets shown as bars. Results from the [2,3]- $\zeta$  extrapolation (●) and [3,4]- $\zeta$  extrapolation (×) from current work compared to previous [2,3]- $\zeta$  results (●), where available.

The performance of the [2,3]- $\zeta$  extrapolation (●) using the formulas from Table 3 is encouraging. The extrapolations generally outperform the 3- $\zeta$  calculations, with the following exceptions: The extrapolation using aug-cc-pvXz basis sets with PBE-D3(BJ) performs worse than the 2- $\zeta$  calculations; and the extrapolations using def2-Xzvp, def2-Xzvpd, cc-pvXz, or aug-cc-pvXz basis sets with PBE0DH-D3(BJ) perform worse than 3- $\zeta$  results. These exceptions will be analysed further below. Notably, the previously determined [2,3]- $\zeta$  extrapolation parameters (●)<sup>3</sup> perform worse than the current formulas in all cases. The results are especially encouraging as the revised formulas perform significantly better with the Jensen basis sets. Finally, the [3,4]- $\zeta$  extrapolation (×) rarely outperforms the 4- $\zeta$  results, and as such cannot be recommended.

Detailed results for the various subsets of the ASCDB database are included in the Supplementary Information (Figs. S1–S4). In the non-covalent interactions subset, [2,3]- $\zeta$  extrapolation performs significantly better than the 3- $\zeta$  calculations using all basis set families and functionals. Generally results comparable to 4- $\zeta$  quality can be expected. In the thermochemistry subset, the performance of the [2,3]- $\zeta$  extrapolation is more mixed. In the most extreme cases, the MUE increases with basis set size (e.g. PBE-D3(BJ) with aug-cc-pvXz or aug-cc-pwcvXz); the resulting poorer performance of extrapolated data can be attributed to the functional as opposed to the basis sets. That said, when 3- $\zeta$  and 4- $\zeta$  results are compared, the MUE in the thermochemistry data does not improve as significantly as for non-covalent interactions, and the [2,3]- $\zeta$  extrapolation is unlikely to increase the errors. In systems with significant non-local character, extrapolation methods offer no benefit; in fact it is this class of problems which causes the poor performance of PBE0DH-D3(BJ) when coupled with a [2,3]- $\zeta$  extrapolation using def2-Xzvp, def2-Xzvdp, aug-cc-pvXz, or aug-cc-pwcvXz basis sets. This is of course a more general issue with DFT as opposed to a basis set completeness error. Notably, the non-local effects subset is the only subset where the double-hybrid PBE0DH-D3(BJ) consistently outperforms the single hybrid PBE0-D3(BJ) when used with 4- $\zeta$  basis sets; the [2,3]- $\zeta$  extrapolation applied to cc-pwcvXz and aug-cc-pwcvXz yields results comparable to the 4- $\zeta$  data. Finally, in the unbiased benchmarks subset of ASCDB, the performance of the [2,3]- $\zeta$  extrapolation is also underwhelming. This might be correlated with the significant increase in performance of the 3- $\zeta$  basis sets compared to the 2- $\zeta$  results, and the comparably minor improvement upon a further increase in the basis set size to 4- $\zeta$ . The 2- $\zeta$  results are likely too inaccurate to improve the 3- $\zeta$  results by extrapolation.

### 3.3 GMTKN55 database

The results of basis set extrapolation with four representative DFAs for the diet100 subset of the GMTKN55 database is shown in Fig. 3. More detailed results, calculated with a wider range of functionals and the cc-pvXz-pp, def2-Xzvp, def2-Xzvdp, and pcseg- $N$  basis set families, are shown in the Supplementary Information (Figs. S5). In the vast majority of cases, the current

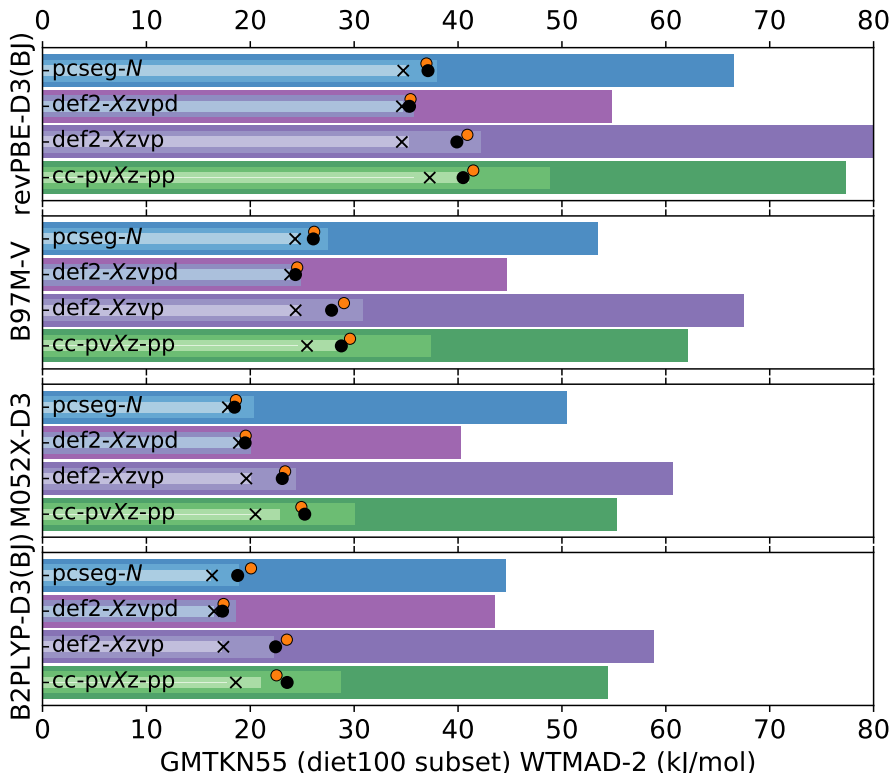


Figure 3: The weighted total mean absolute deviation (WTMAD-2) of various functional and basis set combinations in the diet100 subset of the GMTKN55 database. Calculations with 2-, 3-, and 4- $\zeta$  basis sets shown as bars. Results from the [2,3]- $\zeta$  extrapolation ( $\bullet$ ) and [3,4]- $\zeta$  extrapolation ( $\times$ ) from current work compared to previous [2,3]- $\zeta$  results ( $\bullet$ ), where available.

[2,3]- $\zeta$  extrapolation ( $\bullet$ ) outperforms or matches the previous extrapolation results ( $\bullet$ ), with the only exception being the cc-pvXz-pp results (see e.g. M052X-D3 and B2PLYP-D3(BJ) data). The four DFAs shown in Fig. 3 are all significantly different from the PBE family of functionals (the closest being revPBE, where the Lieb-Oxford bound is relaxed<sup>56</sup>), therefore this improvement in performance confirms the transferrability of the extrapolation formulas, which are based on PBE-like DFAs. Additionally, in all but two cases (DSD-BLYP-D3(BJ) with def2-Xzvp and pcseg-N basis sets), the current [2,3]- $\zeta$  extrapolation performs at least as well as 3- $\zeta$  calculations. As also shown in previous work,<sup>3</sup> the extrapolation using def2-[s,t]zvpd basis sets matches def2-qzvpd performance in this GMTKN55 subset. This is encouraging, as diffusely-augmented basis sets are necessary for proper description of anionic systems, and the def2-Xzvpd sets are defined for the whole periodic table. The performance of the [2,3]- $\zeta$  extrapolation with pcseg-N basis

sets is generally also improved, especially with DHDFAs. However, as also noted previously,<sup>3</sup> the pcseg-2 results (3- $\zeta$  quality) generally outperform the [2,3]- $\zeta$  extrapolation with def2-Xzvp and cc-pvXz-pp and are often comparable to the cc-pvqz-pp results. The choice of an appropriate basis set family therefore remains crucial and cannot be remedied by basis set extrapolation.

### 3.4 NCDT database

Given the good performance of extrapolation schemes for interaction energies of systems with significant non-covalent character, a similarly good performance may be expected for obtaining geometries of non-covalent complexes with such schemes. The results for the NCDT database, presented as the mean absolute deviations of selected bonds in each complex, summed over all 17 complexes, are shown in Fig. 4 for six density functional approximations, three of which have been recommended based on the GMTKN55 benchmarks,<sup>46</sup> and two (B3LYP-D3(BJ) and PBE0-D3) have performed well in previous NCDT benchmarks.<sup>47</sup>

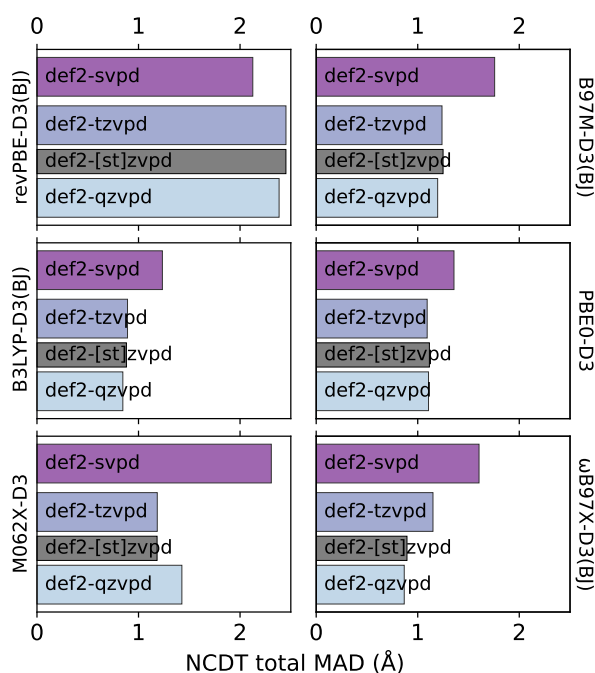


Figure 4: The total of the mean absolute deviations in the optimized structures of the 17 complexes in the NCDT database. The performance of the def2-[st]zvpd extrapolation is indicated by the gray bar.

To bring perspective to the extrapolation results (indicated by an the gray bar in Fig. 4), let us start with a few comments on the general performance of the DFAs and trends in the dataset. Firstly, with the exception of M062X-D3, all DFAs struggle with optimization of the interfragment angle of the  $\text{HCN} \cdots \text{H}_2\text{CO}$  structure; in the case of B3LYP-D3(BJ) this complex alone accounts for  $\sim 25\%$  of the total MAD. Secondly, the three functionals chosen based on their GMTKN55 performance perform significantly worse than both B3LYP-D3(BJ) and PBE0-D3, which remain good methods for structure optimization. However, the comparison is not quite fair, as the range-separated hybrids  $\omega\text{B97M-V}$  and  $\omega\text{B97X-V}$ , which were previously found to offer excellent performance with both GMTKN55 and NCDT, are here excluded from the study as analytical gradients for the non-local correction are not available in Psi4. Accordingly, the  $\omega\text{B97X-D3(BJ)}$  variant, for which analytic gradients are available, is also evaluated and with a 4- $\zeta$  basis it shows results on par with B3LYP-D3(BJ). Finally, in most cases, the improvement from def2-svpd to def2-tzvpd is significant, while the difference between def2-tzvpd and def2-qzvpd is marginal, if any. With revPBE-D3(BJ), increasing the basis set size decreases the agreement with the reference data, which is likely an issue with the DFA. The M062X-D3 functional was trained using a 3- $\zeta$  basis set, therefore a decrease in performance with a 4- $\zeta$  basis is not unexpected. The comparably poor performance of the basis set extrapolation can then be attributed to the small potential for improvement in geometries past 3- $\zeta$  basis sets, and the fact that more accurate interaction energies do not necessarily translate to more accurate geometries. Basis set extrapolations therefore cannot be recommended to obtain more accurate structures with DFT.

### 3.5 L7 database

The most likely use-case of extrapolation methods in DFT would be for systems where a calculation of 3- $\zeta$  quality would approach the limit of affordability. The systems in the L7 database are comparably large, between 48 and 112 atoms in size,<sup>48</sup> yet reference interaction energies of near-CCSD(T) quality are now available.<sup>49</sup> This dataset is therefore an appropriate proxy for the performance of methods in such applications.



257 The results with a variety of combinations of DFAs and basis set families are shown in Fig. 5.  
 258 The three double hybrids (DSD-BLYP-D3(BJ), DSD-PBEP86-NL, and B2PLYP-D3(BJ)) per-  
 259 form poorly compared to the less expensive methods, with a 3- $\zeta$  MAE above 20 kJ/mol and a  
 260 systematic overbinding regardless of the basis set family. This may be due to the exceptionally  
 261 slow convergence of the double hybrid correction  $\Delta E^{\text{dh}}$ , especially in the two spin-component-  
 262 scaled double hybrids. Indeed, the reported MAE of MP2 extrapolated towards the complete  
 263 basis set limit exceeds 30 kJ/mol.<sup>49</sup> On the contrary, a comparably good performance of double  
 264 hybrid DFAs with -NL correction for the L7 database has been reported by Calbo et al., who  
 265 compared B2PLYP-NL and revPBE0-DH-NL to their single hybrid and non-hybrid counterparts,  
 266 albeit with different reference energies.<sup>51</sup> The use of the reference energies of Calbo et al.<sup>51</sup> or  
 267 Al-Hamdani et al.,<sup>50</sup> as opposed to the data of Ballesteros et al.,<sup>49</sup> has no impact on the trends  
 268 shown in Fig. 5. A more detailed analysis of the poor performance of double hybrid DFAs is  
 269 beyond the scope of this work.

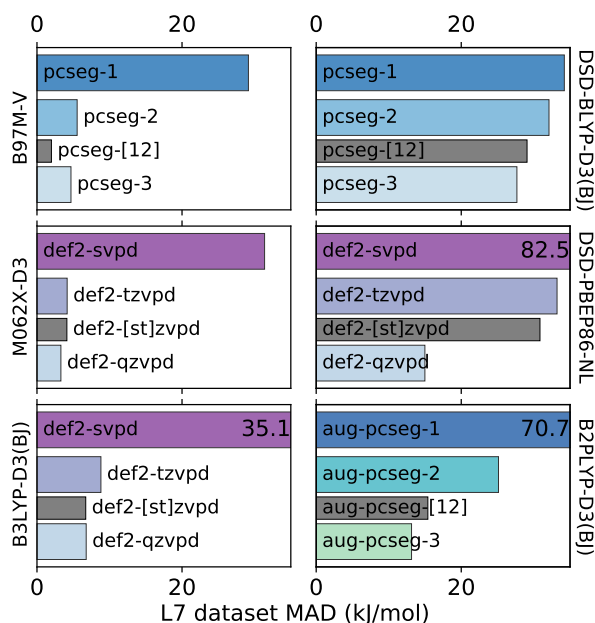


Figure 5: The mean absolute deviations in the interaction energies of the 7 structures in the L7 database. The performance of the [2,3]- $\zeta$  extrapolation methods is indicated by the gray bar.

270 As also shown in Fig. 5, the [2,3]- $\zeta$  extrapolations outperform the 3- $\zeta$  results in all cases,  
 271 with significant improvements over 3- $\zeta$  results using B97M-V/pcseg-[12], B3LYP-D3(BJ)/def2-

[st]zvpd, as well as the three double hybrids. Notably, none of the six methods presented in Fig. 5 are members of the PBE-family of functionals, confirming the good transferability of the extrapolation formulas developed in the current work.

## 4 Conclusion

A systematic way of extrapolating density functional theory results towards the complete basis set limit can be derived from data computed with the PBE family of density functional approximations. The proposed formulas adjust the extrapolation  $\alpha$  for each DFA based on the admixture of Hartree-Fock exchange and second-order perturbation theory correlation into the functional recipe. The extrapolation parameters  $\alpha$  are also dependent on the basis set family used in the extrapolation. The use of diffusely-augmented basis sets is strongly recommended for anionic species.

A set of complete basis set energies for 30 singlet diatomic molecules has been calculated with 21 PBE-related DFAs using the finite element code HelFEM. The dataset is complemented by energies of the same 30 diatomic molecules, calculated with 16 basis set families and the same 21 PBE-related DFAs. This comprehensive set was used to fit linear formulas for scaling the extrapolation parameter  $\alpha$  based on the fraction of HF exchange  $a_x$  and PT correlation  $a_c$  in each of the 21 DFAs. The formulas for calculating the parameter  $\alpha$  are determined individually for each basis set family, while the extrapolation parameters  $\alpha$  derived from such formulas are DFA as well as basis set family specific.

The current, systematically derived, extrapolation parameters outperform the previous, empirically averaged values. For single point energies in the ASCDB database, the performance of [2,3]- $\zeta$  extrapolations exceeds 3- $\zeta$  and generally approaches 4- $\zeta$  results in systems with significant non-covalent character. For thermochemistry calculations the [2,3]- $\zeta$  results at worst match and often outperform 3- $\zeta$  results. The improvement over 3- $\zeta$  results in interaction energies for the L7 database of large systems is also significant. On the contrary, for geometry optimizations of

non-covalent complexes, the [2,3]- $\zeta$  extrapolation barely outperforms 3- $\zeta$  results and therefore cannot be recommended. The basis set extrapolation methods are implemented in a development version of Psi4, with extrapolated analytic gradient calculations as well as extrapolated numerical Hessians available.

## Acknowledgement

I would like to thank Susi Lehtola for his help with HelFEM; Jan Gerit Brandenburg, Lars Goerigk, Amir Karton, Julian Gale, and Laura McKemmish for their comments and feedback; and the Forrest Research Foundation for funding. This work was supported by resources provided by the Pawsey Supercomputing Centre (project f97) and the National Computational Infrastructure (project f97), with funding from the Australian Government and the Government of Western Australia.

## Supporting Information Available

Supporting information available: List of methods used with the diet100 variant of GMTKN55; Figures of results of the subsets of the ASCDB database; Figures of results of the diet100 subset of the GMTKN55 database. Additional supporting information, including the scripts generating the inputs, parsing the outputs, and creating figures are available on Zenodo under DOI: 10.5281/zenodo.4783007.

## Table of contents graphic

## References

- (1) Helgaker, T.; Klopper, W.; Koch, H.; Noga, J. Basis-set convergence of correlated calculations on water. *The Journal of Chemical Physics* **1997**, *106*, 9639–9646.

- (2) Truhlar, D. G. Basis-set extrapolation. *Chemical Physics Letters* **1998**, *294*, 45–48.
- (3) Kraus, P. Basis set extrapolations for density functional theory. *J. Chem. Theory Comput.* **2020**, *16*, 5712–5722.
- (4) Jensen, F. Estimating the Hartree–Fock limit from finite basis set calculations. *Theor Chem Acc* **2005**, *113*, 267–273.
- (5) Shaw, R. A. The completeness properties of Gaussian-type orbitals in quantum chemistry. *Int J Quantum Chem* **2020**, *120*.
- (6) Grimme, S.; Antony, J.; Ehrlich, S.; Krieg, H. A consistent and accurate *ab initio* parametrization of density functional dispersion correction (DFT-D) for the 94 elements H–Pu. *The Journal of Chemical Physics* **2010**, *132*, 154104.
- (7) Vydrov, O. A.; Van Voorhis, T. Nonlocal van der Waals density functional: The simpler the better. *Journal of Chemical Physics* **2010**, *133*.
- (8) Hujo, W.; Grimme, S. Performance of the van der waals density functional VV10 and (hybrid)GGA variants for thermochemistry and noncovalent interactions. *Journal of Chemical Theory and Computation* **2011**, *7*, 3866–3871.
- (9) Becke, A. D. A new mixing of Hartree–Fock and local density-functional theories. *The Journal of Chemical Physics* **1993**, *98*, 1372–1377.
- (10) Grimme, S. Semiempirical hybrid density functional with perturbative second-order correlation. *Journal of Chemical Physics* **2006**, *124*, 034108–034108.
- (11) Kozuch, S.; Martin, J. M. L. Spin-component-scaled double hybrids: An extensive search for the best fifth-rung functionals blending DFT and perturbation theory. *J. Comput. Chem.* **2013**, 2327–2344.

- (12) Najibi, A.; Goerigk, L. The nonlocal kernel in van der Waals density functionals as an additive correction: An extensive analysis with special emphasis on the B97M-V and  $\omega$ B97M-V approaches. *J. Chem. Theory Comput.* **2018**, *14*, 5725–5738.
- (13) Perdew, J. P.; Burke, K.; Ernzerhof, M. Generalized gradient approximation made simple. *Phys. Rev. Lett.* **1996**, *77*, 3865–3868.
- (14) Peverati, R. Fitting elephants in the density functionals zoo: Statistical criteria for the evaluation of density functional theory methods as a suitable replacement for counting parameters. *Int J Quantum Chem* **2020**,
- (15) Lehtola, S. Fully numerical Hartree-Fock and density functional calculations. II. Diatomic molecules. *Int J Quantum Chem* **2019**, *119*.
- (16) Smith, D. G. A.; Burns, L. A.; Simmonett, A. C.; Parrish, R. M.; Schieber, M. C.; Galvelis, R.; Kraus, P.; Kruse, H.; Di Remigio, R.; Alenaizan, A.; James, A. M.; Lehtola, S.; Misiewicz, J. P.; Scheurer, M.; Shaw, R. A.; Schriber, J. B.; Xie, Y.; Glick, Z. L.; Sirianni, D. A.; O'Brien, J. S.; Waldrop, J. M.; Kumar, A.; Hohenstein, E. G.; Pritchard, B. P.; Brooks, B. R.; Schaefer, H. F.; Sokolov, A. Y.; Patkowski, K.; DePrince, A. E.; Bozkaya, U.; King, R. A.; Evangelista, F. A.; Turney, J. M.; Crawford, T. D.; Sherrill, C. D. Psi4 1.4: Open-source software for high-throughput quantum chemistry. *J. Chem. Phys.* **2020**, *152*, 184108.
- (17) Schuchardt, K. L.; Didier, B. T.; Elsethagen, T.; Sun, L.; Gurumoorthi, V.; Chase, J.; Li, J.; Windus, T. L. Basis Set Exchange: A Community Database for Computational Sciences. *J. Chem. Inf. Model.* **2007**, *47*, 1045–1052.
- (18) Pritchard, B. P.; Altarawy, D.; Didier, B.; Gibson, T. D.; Windus, T. L. New Basis Set Exchange: An Open, Up-to-Date Resource for the Molecular Sciences Community. *J. Chem. Inf. Model.* **2019**, *59*, 4814–4820.

- (19) Dunning, T. H. Gaussian basis sets for use in correlated molecular calculations. I. The atoms boron through neon and hydrogen. *J. Chem. Phys.* **1989**, *90*, 1007–1022.
- (20) Kendall, R. A.; Dunning, T. H.; Harrison, R. J. Electron affinities of the first-row atoms revisited. Systematic basis sets and wave functions. *The Journal of Chemical Physics* **1992**, *96*, 6796–6806.
- (21) Woon, D. E.; Dunning, T. H. Gaussian basis sets for use in correlated molecular calculations. III. The atoms aluminum through argon. *The Journal of Chemical Physics* **1993**, *98*, 1358–1371.
- (22) Mourik, T. V.; Dunning, T. H. Gaussian basis sets for use in correlated molecular calculations. VIII. Standard and augmented sextuple zeta correlation consistent basis sets for aluminum through argon. *International Journal of Quantum Chemistry* **2000**, *76*, 205–221.
- (23) Peterson, K. A.; Dunning, T. H. Accurate correlation consistent basis sets for molecular core–valence correlation effects: The second row atoms Al–Ar, and the first row atoms B–Ne revisited. *The Journal of Chemical Physics* **2002**, *117*, 10548–10560.
- (24) Weigend, F.; Ahlrichs, R. Balanced basis sets of split valence, triple zeta valence and quadruple zeta valence quality for H to Rn: Design and assessment of accuracy. *Phys. Chem. Chem. Phys.* **2005**, *7*, 3297.
- (25) Rappoport, D.; Furche, F. Property-optimized Gaussian basis sets for molecular response calculations. *Journal of Chemical Physics* **2010**, *133*, 0–11.
- (26) Jensen, F. Polarization consistent basis sets. II. Estimating the Kohn–Sham basis set limit. *The Journal of Chemical Physics* **2002**, *116*, 7372–7379.
- (27) Jensen, F. Unifying General and Segmented Contracted Basis Sets. Segmented Polarization Consistent Basis Sets. *J. Chem. Theory Comput.* **2014**, *10*, 1074–1085.

- (28) Jensen, F. Polarization consistent basis sets. III. The importance of diffuse functions. *The Journal of Chemical Physics* **2002**, *117*, 9234–9240.
- (29) Ranasinghe, D. S.; Petersson, G. A. CCSD(T)/CBS atomic and molecular benchmarks for H through Ar. *The Journal of Chemical Physics* **2013**, *138*, 144104.
- (30) Ranasinghe, D. S.; Frisch, M. J.; Petersson, G. A. Core-core and core-valence correlation energy atomic and molecular benchmarks for Li through Ar. *The Journal of Chemical Physics* **2015**, *143*, 214110.
- (31) Canal Neto, A.; Muniz, E.; Centoducatte, R.; Jorge, F. Gaussian basis sets for correlated wave functions. Hydrogen, helium, first- and second-row atoms. *Journal of Molecular Structure: THEOCHEM* **2005**, *718*, 219–224.
- (32) Barbieri, P. L.; Fantin, P. A.; Jorge, F. E. Gaussian basis sets of triple and quadruple zeta valence quality for correlated wave functions. *Molecular Physics* **2006**, *104*, 2945–2954.
- (33) Jorge, F.; Sagrillo, P.; de Oliveira, A. Gaussian basis sets of 5 zeta valence quality for correlated wave functions. *Chemical Physics Letters* **2006**, *432*, 558–563.
- (34) Campos, C.; Ceolin, G.; Canal Neto, A.; Jorge, F.; Pansini, F. Gaussian basis set of sextuple zeta quality for hydrogen through argon. *Chemical Physics Letters* **2011**, *516*, 125–130.
- (35) Fantin, P.; Barbieri, P.; Canal Neto, A.; Jorge, F. Augmented Gaussian basis sets of triple and quadruple zeta valence quality for the atoms H and from Li to Ar: Applications in HF, MP2, and DFT calculations of molecular dipole moment and dipole (hyper)polarizability. *Journal of Molecular Structure: THEOCHEM* **2007**, *810*, 103–111.
- (36) de Oliveira, P.; Jorge, F. Augmented Gaussian basis set of quintuple zeta valence quality for H and from Li to Ar: Applications in DFT calculations of molecular electric properties. *Chemical Physics Letters* **2008**, *463*, 235–239.

- (37) Gould, T. 'Diet GMTKN55' offers accelerated benchmarking through a representative subset approach. *Phys. Chem. Chem. Phys.* **2018**, *20*, 27735–27739.
- (38) Goerigk, L.; Hansen, A.; Bauer, C.; Ehrlich, S.; Najibi, A.; Grimme, S. A look at the density functional theory zoo with the advanced GMTKN55 database for general main group thermochemistry, kinetics and noncovalent interactions. *Phys. Chem. Chem. Phys.* **2017**, *19*, 32184–32215.
- (39) Morgante, P.; Peverati, R. Statistically representative databases for density functional theory via data science. *Phys. Chem. Chem. Phys.* **2019**, *21*, 19092–19103.
- (40) Grimme, S.; Ehrlich, S.; Goerigk, L. Effect of the damping function in dispersion corrected density functional theory. *J. Comput. Chem.* **2011**, *32*, 1456–1465.
- (41) Adamo, C.; Barone, V. Toward reliable density functional methods without adjustable parameters: The PBE0 model. *The Journal of Chemical Physics* **1999**, *110*, 6158–6158.
- (42) Brémond, E.; Adamo, C. Seeking for parameter-free double-hybrid functionals: The PBE0-DH model. *The Journal of Chemical Physics* **2011**, *135*, 024106.
- (43) Bousquet, D.; Brémond, E.; Sancho-García, J. C.; Ciofini, I.; Adamo, C. Non-parametrized functionals with empirical dispersion corrections: A happy match? *Theor Chem Acc* **2015**, *134*, 1602.
- (44) Kraus, P.; Obenchain, D. A.; Frank, I. Benchmark-quality semiexperimental structural parameters of van der Waals complexes. *J. Phys. Chem. A* **2018**, *122*, 1077–1087.
- (45) Kraus, P. Non-covalent dimers and trimers (NCDDT) database version 2.1. 2020.
- (46) Goerigk, L.; Mehta, N. A trip to the density functional theory zoo: Warnings and recommendations for the user. *Aust. J. Chem.* **2019**, *72*, 563.
- (47) Kraus, P.; Frank, I. Density functional theory for microwave spectroscopy of noncovalent complexes: A benchmark study. *J. Phys. Chem. A* **2018**, *122*, 4894–4901.



- (48) Sedlak, R.; Janowski, T.; Pitoňák, M.; Řezáč, J.; Pulay, P.; Hobza, P. Accuracy of Quantum Chemical Methods for Large Noncovalent Complexes. *J. Chem. Theory Comput.* **2013**, *9*, 3364–3374.
- (49) Ballesteros, F.; Dunivan, S.; Lao, K. U. Coupled cluster benchmarks of large noncovalent complexes: The L7 dataset as well as DNA–ellipticine and buckycatcher–fullerene. *J. Chem. Phys.* **2021**, *154*, 154104.
- (50) Al-Hamdani, Y. S.; Nagy, P. R.; Barton, D.; Kállay, M.; Brandenburg, J. G.; Tkatchenko, A. *Interactions between Large Molecules: Puzzle for Reference Quantum-Mechanical Methods*; arXiv [physics.chem-ph] 2009.08927, 2020.
- (51) Calbo, J.; Ortí, E.; Sancho-García, J. C.; Aragó, J. Accurate treatment of large supramolecular complexes by double-hybrid density functionals coupled with nonlocal van der Waals corrections. *J. Chem. Theory Comput.* **2015**, *11*, 932–939.
- (52) SciPy 1.0 Contributors; Virtanen, P.; Gommers, R.; Oliphant, T. E.; Haberland, M.; Reddy, T.; Cournapeau, D.; Burovski, E.; Peterson, P.; Weckesser, W.; Bright, J.; van der Walt, S. J.; Brett, M.; Wilson, J.; Millman, K. J.; Mayorov, N.; Nelson, A. R. J.; Jones, E.; Kern, R.; Larson, E.; Carey, C. J.; Polat, İ.; Feng, Y.; Moore, E. W.; VanderPlas, J.; Laxalde, D.; Perktold, J.; Cimrman, R.; Henriksen, I.; Quintero, E. A.; Harris, C. R.; Archibald, A. M.; Ribeiro, A. H.; Pedregosa, F.; van Mulbregt, P. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat Methods* **2020**, *17*, 261–272.
- (53) Lonsdale, D. R.; Goerigk, L. The one-electron self-interaction error in 74 density functional approximations: A case study on hydrogenic mono- and dinuclear systems. *Phys. Chem. Chem. Phys.* **2020**,
- (54) Kutzelnigg, W. Expansion of a wave function in a Gaussian basis. I. Local versus global approximation. *Int. J. Quantum Chem.* **2013**, *113*, 203–217.

- 458 (55) Varandas, A. J. Straightening the hierarchical staircase for basis set extrapolations: A low-  
459 cost approach to high-accuracy computational chemistry. *Annu. Rev. Phys. Chem.* **2018**,  
460 *69*, 177–203.
- 461 (56) Zhang, Y.; Yang, W. Comment on “Generalized gradient approximation made simple”.  
462 *Phys. Rev. Lett.* **1998**, *80*, 890–890.