

Improving the Accuracy of Composite Methods: A G4MP2 Method with G4-like Accuracy and Implications for Machine Learning

Naveen K. Dandu,* Rajeev S. Assary, Paul C. Redfern, Logan Ward, Ian Foster, and Larry A. Curtiss*



Cite This: *J. Phys. Chem. A* 2022, 126, 4528–4536



Read Online

ACCESS |



Metrics & More

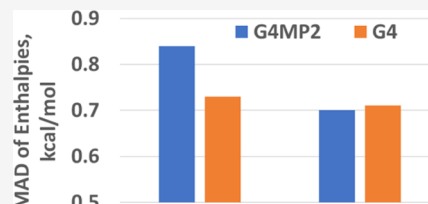


Article Recommendations



Supporting Information

ABSTRACT: G4MP2 theory has proven to be a reliable and accurate quantum chemical composite method for the calculation of molecular energies using an approximation based on second-order perturbation theory to lower computational costs compared to G4 theory. However, it has been found to have significantly increased errors when applied to larger organic molecules with 10 or more nonhydrogen atoms. We report here on an investigation of the cause of the failure of G4MP2 theory for such larger molecules. One source of error is found to be the “higher-level correction (HLC)”, which is meant to correct for deficiencies in correlation contributions to the calculated energies. This is because the HLC assumes that the contribution is independent of the element and the type of bonding involved, both of which become more important with larger molecules. We address this problem by adding an atom-specific correction, dependent on atom type but not bond type, to the higher-level correction. We find that a G4MP2 method that incorporates this modification of the higher-level correction, referred to as G4MP2A, becomes as accurate as G4 theory (for computing enthalpies of formation) for a test set of molecules with less than 10 nonhydrogen atoms as well as a set with 10–14 such atoms, the set of molecules considered here, with a much lower computational cost. The G4MP2A method is also found to significantly improve ionization potentials and electron affinities. Finally, we implemented the G4MP2A energies in a machine learning method to predict molecular energies.



INTRODUCTION

Gaussian-4 (G4)¹ theory, the fourth in a series of the Gn composite methods² developed to predict molecular energies to chemical accuracy, has been widely used since it was first developed. It has an accuracy of 0.83 kcal/mol (mean absolute deviation) for the 454 systems in the G3/05 test set. G4 theory uses a combination of single-point energies [CCSD(T)/6-31G(d), MP4/6-31+G(d), MP4/6-31G(2df,p), MP2(FU)/G3LargeXP, and HF limit] at the B3LYP/6-31G(2df,p) geometry and an empirical higher-level correction (HLC) to achieve this accuracy. A major limitation of G4 theory is the large increase in computational resources required as the molecules become larger, due to the MP4 calculations. To lower computational costs, the G4MP2 method was developed in which the fourth-order perturbation components of G4 theory³ are replaced with reduced perturbation theory levels. G4MP2 is approximately 6–8 times faster than G4.³ While this increase in speed comes with some loss of accuracy, with the mean absolute deviation for the G3/05 test set increasing to 1.03 kcal/mol, G4MP2 has been widely used.

G4 theory is one of a number of accurate quantum chemical methods that have been developed for predicting molecular energies to better than 1 kcal/mol accuracy. Other methods include the correlation consistent Composite Approach (ccCA) of Wilson and co-workers,⁴ Complete Basis Set (CBS) Method of Petersson and co-workers,⁵ coupled cluster-based methods of Feller, Peterson, and Dixon,⁶ the Weizmann methods of Martin and co-workers,^{7,8} and the Wuhan-

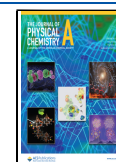
Minnesota scaling methods of Truhlar and co-workers.⁹ There have also been variations of G4MP2 theory, such as G4MP2-6X and G4MP2-XK, that have been meant to address certain types of problems and to extend to other elements.^{10–12}

In recent work, we have reported on investigations of quantum chemically informed machine learning for the prediction of G4MP2 energies of organic molecules. In one paper, we used quantum chemical energies of a set of 130,258 molecules with 1–9 nonhydrogen atoms¹³ to train and assess the performance of several machine learning (ML) methods.¹⁴ Two ML methods were used: a kernel-based ridge regression (FCHL)¹⁵ and a continuous filter convolutional neural network (SchNet),¹⁶ both combined with Δ -learning on the difference between the G4MP2 energies and B3LYP energies. The best-performing ML method was able to predict G4MP2 energies to within 0.1 kcal/mol.¹⁴ In a subsequent paper,¹⁷ we examined how the two best-performing ML methods from the first paper performed on molecules larger than the test set of 1–9 nonhydrogen atoms. In addition, in this study, we also

Received: February 24, 2022

Revised: June 15, 2022

Published: July 5, 2022



assessed the suitability of the ω B97XD density functional method for learning. To have accurate enthalpies of formation for molecules with 10 to 14 nonhydrogen atoms, we derived accurate experimental data for a set of 191 molecules. The better-performing ML method investigated in the second paper, FCHL- Δ , gave atomization energies for these 191 organic molecules within about 0.4 kcal/mol of accurate quantum chemical energies calculated by the G4MP2 method. Although this level of accuracy is less than that obtained with FCHL- Δ for a hold-out set of the smaller molecules (0.1 kcal/mol), it demonstrated that quantum chemically informed ML can be used to successfully predict energies of large organic molecules with sizes beyond those in the training set at a much lower cost in computer time.

Another result in ref 17 was that analysis of G4MP2 enthalpies of formation of the 191 molecules with 10–14 nonhydrogen atoms revealed that the G4MP2 method is significantly less accurate for these larger molecules than it is for the set of smaller organic molecules with 1–9 nonhydrogen atoms. Compared to accurate experimental data, the G4MP2 method had a mean absolute deviation of 1.41 kcal/mol for the molecules with 10–14 nonhydrogen atoms, which was dramatically larger than the 0.79 kcal/mol for 459 organic molecules having 1–9 nonhydrogen atoms. From an analysis of the molecules having large errors in the set of 191 molecules, it was found that G4MP2 does especially poorly for large aromatic molecules and for molecules with multiple nitrogen atoms. In addition, G4 theory was investigated for a subset of 54 aromatic molecules of the 10–14 nonhydrogen set and was found to do considerably better (mean absolute deviation of 1.15 kcal/mol compared to 1.84 kcal/mol for the same set with G4MP2).

The objectives of the work reported in this paper were to determine the cause of the failure of the G4MP2 theory for the larger molecules identified in ref 17 and to investigate remedies for the problem. One source of error investigated was the higher-level correction, which is meant to correct for deficiencies in correlation contributions to the molecular energies. One problem for this approximation is that it assumes that the contribution is independent of the element and the type of bonding involved. This approximation seems to be a problem that becomes worse in larger molecules, which may have types of molecules not present in the sets of smaller molecules on which G4MP2 performs well. To address this problem, we have added an atom-specific correction to the HLC. Importantly, this correction is dependent only on the type of atom, not on the bond type, as the latter would make implementation difficult. With this simple addition to the HLC, we find that the G4MP2 method, referred to as G4MP2A, becomes as accurate as G4 theory, with a much smaller cost. We also added this type of modification to the G4 method, referred to as G4A, but found little improvement. The G4MP2A method is also found to do better than G4MP2 for ionization potentials and electron affinities. The G4MP2A method has been implemented in the ML learning method previously reported for G4MP2. This implementation improves the ML-predicted energies considerably because they are now trained on data with G4-like rather than G4MP2 accuracy.

In the **Methods** section, we describe the theoretical methods and test sets used in this paper. In the **Results and Discussion** section, we report on how the new parameterization performs for enthalpies, ionization potentials, and electron affinities, as

well as an analysis of the results. In this section, we also describe the new ML results for the G4MP2A method. Conclusions are drawn in the **Conclusions** section.

METHODS

The energies of the molecules and ions considered in this work were calculated using the composite G4MP2 method.³ Some results for the G4 method¹ are included for comparison, and a similar procedure is used. In the G4MP2 method, a higher-level correction (HLC) is added on to the series of calculations to account for remaining deficiencies in the energies. The HLC has six parameters. The HLC includes corrections for the number of pairs of electrons in closed-shell molecules (or ions), the number of pairs of electrons and unpaired electrons in open-shell molecules (or ions), and the number of pairs of electrons and unpaired electrons in atoms and atomic ions. There is also a correction for the energy of a pair of electrons in molecular and atomic nonhydrogen species having only one pair of valence electrons.

In the current work, we have added an “atom-specific correction”, i.e., a correction that is specific to the number of each type of atom in the molecule or ion, that increases or decreases the atomic energies to minimize the mean absolute deviation with respect to experiment for the test sets as described below. The atom-specific correction is similar to a method proposed by Perdew et al.¹⁸ for improving several density functional methods. The new G4MP2 method with atom-specific corrections is referred to as G4MP2A. Details of the atom-specific corrections and results using the new method are given in the next section. The G4MP2 and G4 methods as implemented within the Gaussian 16 package are used with the standard settings.¹⁹ In addition, we determined atom-specific parameters for B3LYP^{20,21} and ω B97XD²² for enthalpies in a similar manner. The B3LYP study was done with the 6-31G(2df,p) basis, and the ω B97XD study was done with the 6-311+G(3df,2p) basis as in our previous study.¹⁷

We used a set of 681 neutral enthalpies of formation to optimize the atom-specific parameters. This set is composed of the 222 first- and second-row molecules from the G3/99 test set²³ and the 459 first-row molecules of the Pedley test set¹³ (referred to as the PDS(9)) composed of 1–9 nonhydrogen atoms. The PDS(9) test set includes experimental data for 175 hydrocarbons and 284 substituted hydrocarbons from experimental data in the Pedley compilation²⁴ that was critically evaluated.¹³

We also used a set of larger molecules with 10–14 nonhydrogen atoms, referred to as the PDS(10–14) test set,¹⁷ to assess how well the new methods based on atom-specific correction perform on larger molecules than in the original set. The larger molecules include 71 hydrocarbons and 120 substituted hydrocarbons. The PDS(10–14) test set¹⁷ was derived from experimental data for molecules with 10–14 nonhydrogen atoms in the Pedley compilation.²⁴ The experimental data was checked by an isodesmic scheme²⁵ that eliminated 34 of 225 molecules to get to the 191 molecules.²⁴ We have computed standard enthalpies of formation at 298 K using methods described previously.^{26,27} Some G4MP2A energies still had large errors for the PDS(10–14) set even after application of the isodesmic scheme and these were investigated for possible incorrect conformers as described in the next section. Six molecules were further eliminated after this process (see below).

The electron affinities and ionization potentials used to derive atom-specific corrections come from the G3/99 test set. This test set has 88 ionization potentials and 58 electron affinities of first- and second-row molecules that were used in this study. The ionization potentials and electron affinities for this study were calculated as described previously.²³

RESULTS AND DISCUSSION

Atom-Specific Corrections. The atom-specific corrections for G4MP2 are defined as follows for any element X

$$E_e(\text{G4MP2A}, X) = E_e(\text{G4MP2}, X) - E_{\text{ASC}}(\text{G4MP2A}, X) \quad (1)$$

where $E_{\text{ASC}}(\text{G4MP2A}, X)$ is the atomic-specific correction (ASC) to an atomic energy of element X and $E_e(\text{G4MP2}, X)$ is the G4MP2 energy of element X calculated as described in ref 3. The $E_{\text{ASC}}(\text{G4MP2A}, X)$ corrections have been optimized to give the smallest mean absolute deviation for G4MP2A enthalpies of formation compared to the experiment for the 681 molecules in the G3/99 and PDS(9) test sets, which are described in the Methods section. The G4MP2A enthalpies are calculated from G4MP4 energies³ of the molecules and atomic energies corrected as in eq 1. The values of the corrections are given in Table 1 for the elements H, B–F, and Al–Cl. (The ASC parameters are not used for Li, Be, Na, and Mg due to the lack of data for molecules containing these elements in the test sets.) For example, for the ethane atomization energy, the corrections would be as follows

C_2H_6 : 0 mH

2C: $2^*(-0.444 \text{ mH})$

6H: $6^*(0.194 \text{ mH})$

where the atomic values are from Table 1. Alternatively, the corrections could be subtracted from the molecule energy. The higher-level corrections used for G4MP2A are the same as those used in G4MP4; they are also given in Table 1. The atom-specific corrections for G4A theory are defined in a similar way to that for G4MP2A. For any element X

$$E_e(\text{G4A}, X) = E_e(\text{G4}, X) - E_{\text{ASC}}(\text{G4A}, X) \quad (2)$$

where $E_{\text{ASC}}(\text{G4A}, X)$ is the atomic-specific correction and $E_e(\text{G4}, X)$ is the G4 energy of element X calculated as described in ref 1. The $E_{\text{ASC}}(\text{G4A}, X)$ parameters were optimized to give the smallest mean absolute deviation for G4A enthalpies of formation compared to the experiment for the 681 molecules in the G3/99 and PDS(9) test sets, which are described in the Methods section. The G4A enthalpies are calculated from G4 energies¹ of the molecules and atomic energies corrected as in eq 2. The values of these corrections are given in Table 1. The higher-level corrections used for G4A are the same as G4 and are also given in Table 1.

The results for the modified G4MP2 and G4 methods with atom-specific corrections added, G4MP2A and G4A, are given in Table 2. These results include the performance on the G3/99 test set of 222 enthalpies of formation and the PDS(9) test set of 459 enthalpies of formation. We note that the addition of the ASC parameter to G4 theory has a modest effect, only slightly decreasing the overall mean absolute deviation (MAD) for both test sets with respect to experiment from 0.73 to 0.71 kcal/mol. The decrease for the individual test sets is similar: 0.73 to 0.71 kcal/mol for the G3/99 test set and 0.74 to 0.72 kcal/mol for the PDS(9) test set. In contrast, a much larger improvement is seen when the ASC parameters are added to the G4MP2 method. The addition of the ASC parameter to G4MP2 decreases the overall MAD with respect to the experiment from 0.84 to 0.70 kcal/mol. The decrease is from 0.93 to 0.71 kcal/mol for the G3/99 test set and from 0.79 to 0.69 kcal/mol for the PDS(9) test set. The breakdown of the results in the types of molecules is also given in Table 2. For the G3/99 test set, the G4MP2A improvement occurs for all categories of molecules. For the PDS(9) test set, the largest improvement occurs for hydrocarbons (0.68 to 0.52 kcal/mol), although substituted hydrocarbons (0.86 to 0.81 kcal/mol) are also improved. The key result for the ASC parameterization is that its inclusion in the G4MP2 energies gives an overall accuracy for enthalpies of formation that is about the same as G4 theory for the combined test set of 681 molecules, i.e., 0.70 kcal/mol, while requiring much less computer time.

The optimization of the ASC parameters was also done without the HLC parameterization. For G4, the MAD increases to 0.77 kcal/mol (from 0.71 kcal/mol), and for G4MP2, it increases only to 0.72 kcal/mol (from 0.70 kcal/mol), both of these results being for the combined test set of 681 molecules. Thus, G4MP2A could be formulated without an HLC with essentially the same results. For this study, we have included the HLC in G4MP2A, to avoid any change in the basic method.

While the G4MP2A method also gave some improvement for the larger molecules of the PDS(10–14) set compared to G4MP2 from our previous study,³ we observed that some of

Table 1. Higher-Level Correction (HLC) and Atom-Specific Corrections (ASC) for G4 and G4MP2 Methods

	HLC values (mhartrees) ^a			
	G4MP2	G4MP2A	G4	G4A
A	9.472	9.472	6.947	6.947
B	3.102	3.102	2.441	2.441
C	9.741	9.741	7.116	7.116
D	2.115	2.115	1.414	1.414
A'	9.769	9.769	7.128	7.128
E	2.379	2.379	2.745	2.745
	ASC (mhartrees) ^{b,c}			
H		0.194		0.002
B		4.614		0.015
C		−0.444		0.021
N		0.508		−0.227
O		0.461		−0.361
F		0.247		−0.011
Al		−1.879		−1.450
Si		−1.933		−0.783
P		−1.037		0.220
S		−1.924		−0.373
Cl		−0.614	−	0.274

^aThe HLC for the G4A and G4MP2A methods are not optimized and are taken to be the same as those for G4 and G4MP2. ^bThe ASC parameters are optimized for the 691 neutral enthalpies of formation in the G3/99 + PDS(9) test sets; the ASC parameter is not used for Li, Be, Na, and Mg due to the lack of data for molecules containing these elements in the test sets. ^cThese parameters are added to the atomic energies (Note: they are subtracted from neutral molecular energies for calculation of ionization potentials and electron affinities where the cations and anions are adjusted by parameters optimized for ionization potentials and electron affinities; see Table S3).

Table 2. G4A and G4MP2A Mean Absolute Deviations with Experimental Enthalpies of Formation for Three Test Sets G3/99, PDS(9), and PDS(10–14)^a

		mean absolute deviations, kcal/mol			
	test set	G4MP2	G4MP2A ^b	G4	G4A ^b
G3/99	all (222)	0.93	0.71	0.73	0.71
	nonhydrogens (47)	1.42	0.96	1.08	1.08
	hydrocarbons (38)	0.63	0.48	0.48	0.46
	substituted hydrocarbons (91)	0.81	0.68	0.67	0.66
	inorganic hydrides (15)	0.94	0.45	0.80	0.73
	radicals (31)	0.86	0.83	0.66	0.61
PDS(9)	all (459)	0.79	0.69	0.74	0.72
	hydrocarbons (175)	0.68	0.52	0.54	0.55
	substituted hydrocarbons (284)	0.86	0.81	0.86	0.82
G3/99 + PDS(9)	all (681)	0.84	0.70 ^c	0.73	0.71 ^d
PDS(10–14)	all (185) ^e	1.43(4.8)	0.96(3.3)		
	hydrocarbons (68)	1.81	0.93		
	substituted hydrocarbons (117)	1.20	0.98		
	aromatics (54) ^f	1.93	1.07	1.15	1.16
	all aromatics (87)	1.92	1.06		

^aG4MP2 results for G3/99 are from ref 2, for PDS(9) are from ref 13, and for PDS(10–14) are from ref 17 (but updated for the elimination of 6 molecules and new geometries, see text). G4 results for G3/99 are from ref 2. The G4 results for PDS(9) are from this work. The G4 PDS(10–14) aromatics are from ref 17 (but updated for new geometries, see text). The COF₂ molecules were not included in the G3/99 results due to a large experimental error.²³ ^bAtom-specific correction parameters optimized for the combined set of 681 enthalpies of neutral molecules in the G3/99 and PDS(9) test sets as described in the text. ^cOptimization of the atom-specific correction parameters with no higher-level correction (HLC) gives 0.72 kcal/mol. ^dOptimization of the atom-specific correction parameters with no higher-level correction (HLC) gives 0.77 kcal/mol. ^eNumbers in parentheses are for max deviations. ^fFor 54 of the 87 aromatic molecules in the PDS(10–14) test set.

the molecules still performed poorly. We identified one possible cause of this to be that the structure that we used in the previous study was not the lowest-energy conformer. Due to the large number of potential conformers of the larger molecules in this test set, locating the lowest-energy conformer is more difficult using traditional geometry optimization techniques. As a result, we used an advanced sampling technique to search for other more stable conformers of these molecules based on G4MP2 energies. For this purpose, we used the CREST (Conformer–Rotamer Ensemble Sampling Tool) utility program²⁸ to perform conformational searches on the previously published PDS(10–14) set of geometries as a starting point.¹⁷ When performing these conformational searches, referred to as Method-1, we optimized the resultant conformers using the GFN2-xTB method^{29,30} and then selected the lowest-energy conformer for carrying out a new G4MP2 energy calculation. Although it provided lower-energy conformers for some of the molecules, we still had cases where some previously published¹⁷ conformers were more stable than those found by the CREST program. Thus, we also performed CREST starting from SMILES strings, referred to as Method-2. However, even this method did not always provide the lowest-energy conformers. Overall, we updated our previously published PDS(10–14) set of molecules containing 191 species with new geometries when Method-1 or Method-2 gave the lowest energies based on comparing G4MP2 energies. With this improved set of geometries, we then eliminated six molecules out of 191, based on isodesmic reactions performed on the resultant geometries. The isodesmic evaluation was done similarly to our previous study.²⁵ The new geometries of the 185 are given in an Excel spreadsheet in the [Supplemental Information](#).

Table 2 shows the mean absolute deviations (MADs) of the standard enthalpies of formation at 298 K for the 185

PDS(10–14) molecules computed with the G4MP2 and G4MP2A methods. The standard enthalpies of formation for the 185 molecules at the G4MP2A level of theory using the updated molecules along with their experimental values are given in Table S1. This table also shows the deviations for both the G4MP2A and G4MP2 methods from experiment for the 185 molecules. Additionally, deviations in standard enthalpies of the six eliminated molecules are given in Table S2. In Table 2, the overall MAD of G4MP2 of 1.43 kcal/mol, relative to experimental values, decreases to 0.96 kcal/mol for the G4MP2A method. Among the 185 molecules, there are 68 hydrocarbons and 117 substituted hydrocarbons. The MAD of the 71 hydrocarbons decreases dramatically from 1.80 kcal/mol for G4MP2 to 0.93 kcal/mol for the new G4MP2A method. The MADs of the substituted hydrocarbons also decrease, to 1.20 kcal/mol for G4MP2 and to 0.98 kcal/mol and G4MP2A.

The G4MP2 and G4MP2A MADs for the molecules in the PDS(9) and PD(10–14) datasets as a function of size are given in Table 3. The dramatic improvement in the G4MP2A energies compared to G4MP2 for molecules with 10–14 nonhydrogen atoms is illustrated in Figure 1. When the comparison is made on a per electron pair basis as a function of nonhydrogen atoms in Figure 2, G4MP2A also shows a significant improvement over G4MP2 for the molecule set with 10–14 nonhydrogen atoms. In addition, we have performed computations with G4 theory for a limited set of the larger molecules in the 185 PDS(10–14) set, to enable comparison with G4MP2A. This set includes the 54 aromatic molecules in the PDS(10–14) set. The results in Table 2 show that for G4MP2, the error is reduced from 1.93 to 1.07 kcal/mol for this set of molecules, which is actually better than the G4 result (1.16 kcal/mol). We note that the set of data used to optimize the atomic correction parameters did not include the

Table 3. Mean Absolute Deviation (MAD) from Experiment of Enthalpies of Formation at 298 K for the PDS(9) and PDS(10–14) Datasets for the G4MP2 and G4MP2A Methods (in kcal/mol)

number of nonhydrogen atoms	number of molecules	MAD of enthalpy of formation, kcal/mol		
		G4MP2	G4MP2A	Per heavy atom for G4MP2A
1	1	0.13	0.10	0.10
2	6	0.53	0.36	0.18
3	17	0.52	0.38	0.13
4	42	0.64	0.59	0.15
5	63	0.55	0.51	0.10
6	91	0.68	0.62	0.10
7	92	0.81	0.76	0.11
8	105	0.99	0.88	0.11
9	42	1.24	0.88	0.10
10	80	1.08	0.79	0.08
11	43	1.28	0.93	0.10
12	31	1.67	1.19	0.09
13	24	1.77	1.19	0.10
14	13	2.87	1.20	0.07

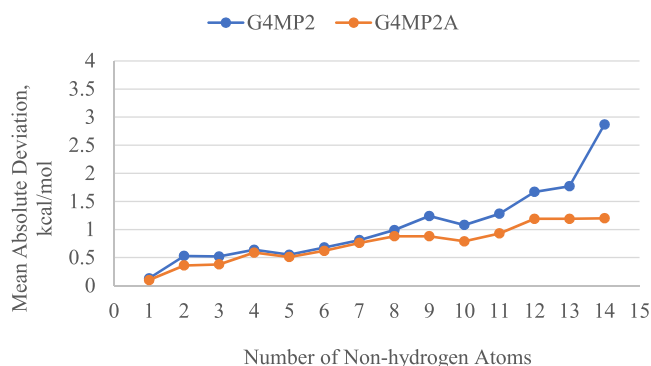


Figure 1. Comparison of G4MP2 vs G4MP2A mean absolute deviations for the PDS(9) and PDS(10–14) test sets as a function of nonhydrogen atoms.

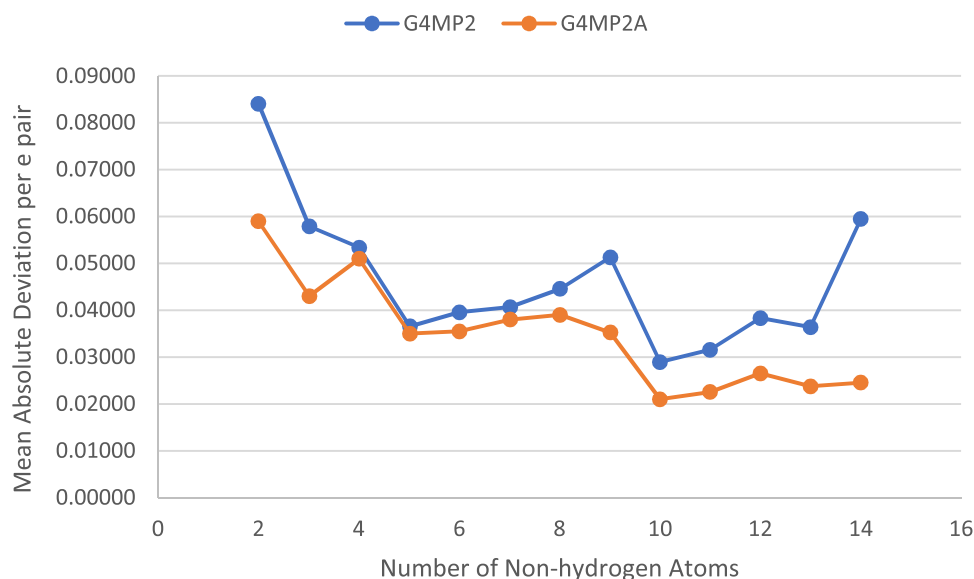


Figure 2. Comparison of G4MP2 vs G4MP2A mean absolute deviations per electron pair as a function of nonhydrogen atoms for the PDS(9) and PDS(10–14) test sets.

PDS(10–14) set of 185 molecules, indicating that the parameters should be valid for larger molecules.

The gradual increase in error with molecule size in Figure 1 for G4MP2A is probably due to the accumulation of error with the number of atoms in the molecule, not any inherent problem with the methodology. We found that the MAD for the enthalpies of formation for the 87 aromatic molecules in the PDS(10–14) set is 1.92 kcal/mol. In contrast, the MAD for the remaining 98 (nonaromatic) molecules is 0.99 kcal/mol (almost one-half compared to the MAD of aromatic molecules). The signed error per carbon atom increases from 1.08 kcal/mol for 6 carbon atoms to 3.04 kcal/mol for 14 carbon atoms in the aromatic molecules in PDS(10–14). The results for all 87 molecules as a function of the number of carbon atoms in the molecule are shown in Figure 3, and the results are in Table 4. The results show a significant increase in errors on G4MP2 energies with an increase in the number of carbon atoms for the 87 of PDS(10–14) aromatic molecules. However, there is no such effect in the case of the G4MP2A enthalpies. Thus, it is apparent that the higher-level corrections (HLC) used in G4MP2 do not do very well for aromatics. However, in the case of G4, for similarly obtained HLCs, the MADs are much better because G4 has less approximations than G4MP2. The average error per nonhydrogen atom in the molecule, shown in Table 3, is quite constant with increasing molecule size, a result consistent with an accumulation of error. For G4MP2A, only 16 of the 185 molecules in Table S1 have errors of 2 kcal/mol or more (only two of which are slightly greater than 3 kcal/mol). Of the 16, these are distributed as follows: 11 are substituted or multiple aromatic rings, three are substituted alkyl or ether chain molecules, and two are multiple aliphatic rings. In comparison, the G4MP2 method does much worse, with 40 out of the 185 having errors of 2 kcal/mol or more, and one-half of those greater than 3 kcal/mol.

Table 5 gives the G4MP2A results, i.e., results for the inclusion of atom-specific corrections in the G4MP2 method, for the calculation of ionization potentials and electron affinities in the G33 test set. The atomic-specific corrections for anions and cations are optimized in the same way as for the

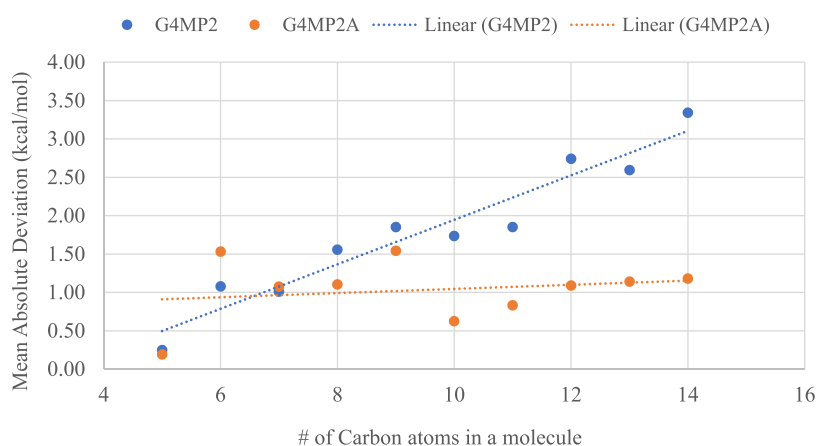


Figure 3. Mean absolute deviation (MAD) from experiment of enthalpies of formation at 298 K for the 87 aromatic molecules of the PDS(10–14) dataset grouped based on the number of carbon atoms in a molecule for the G4MP2 and G4MP2A methods. Data from Table 4.

Table 4. Mean Absolute Deviation (MAD) from Experiment of Enthalpies of Formation at 298 K for the 87 Aromatic Molecules of the PDS(10–14) Dataset Grouped Based on the Number of Carbon Atoms in a Molecule for the G4MP2 and G4MP2A Methods

number of carbon atoms in a molecule	number of molecules in the group	MAD of enthalpy of formation, kcal/mol	
		G4MP2	G4MP2A
5	1	0.25	0.19
6	7	1.08	1.53
7	10	1.01	1.08
8	11	1.56	1.10
9	8	1.85	1.54
10	14	1.74	0.63
11	11	1.85	0.83
12	14	2.74	1.09
13	3	2.59	1.14
14	8	3.34	1.18

Table 5. G4A and G4MP2A Mean Absolute Deviations (kcal/mol) with respect to Experiment for the G3/99 Test Set²³ of Ionization Energies and Electron Affinities^a

test set		G4MP2	G4MP2A	G4
IPs	All (69)	0.93	0.79	0.87
EAs	All (47)	1.05	0.83	0.80

^aNo atoms included in comparison to experiment. B2F4 excluded based on recently identified problems with experiment.³³

neutrals for this G33 set and are used with those for the neutral molecules. The results of the optimized atom-specific parameters are given in Table S4. For example, for the NH₃ ionization energy, corrections would be calculated as follows

$$\text{NH}_3: (-0.508 \text{ mH}) + 3 \times (0.194 \text{ mH})$$

$$\text{NH}_3^+: (-0.70 \text{ mH}) + 3 \times (-0.05 \text{ mH})$$

where the neutral atomic values are from Table 1. As shown in Table 5, the G4MP2A MAD for the 69 ionization potentials in the G33 test set is 0.79 eV, a significant improvement over 0.93 eV for G4MP2. Similarly, the G4MP2A MAD for the 47 electron affinities in the G33 test set is 0.83 eV, a significant improvement over 1.04 eV for G4MP2. Similar to the case for

enthalpies of neutrals, the results in Table 5 indicate that the G4MP2A method has an accuracy comparable to full G4 theory for ionization potentials and electron affinities.

We also investigated the use of atom-specific corrections with two density functional methods to see if they provide significant improvement. The two methods used are the B3LYP and ω B97xd functionals, as described in the Methods section; results are given in Table 6. For both functionals, the atom-specific corrections were optimized to give the smallest mean absolute deviation for enthalpies of formation compared to experiment for the 681 molecules in the G3/99 and PDS(9) test sets as described in the Methods section for G4MP2A. The results in Table 6 indicate that the use of the atomic-specific corrections significantly improves the MADs for both functionals. For B3LYP, the improvement is from 4.16 to 2.36 kcal/mol, and for ω B97XD, the improvement is from 1.92 to 1.48 kcal/mol. The improvement for B3LYP is better than that for ω B97XD, possibly because there is more room for improvement. The ω B97XD functional is still significantly better than B3LYP for enthalpies. The use of atom-specific corrections for DFT calculations of ionization potentials and electron affinities will be reported on in another publication.

Machine Learning Based on Methods with Atom-Specific Corrections. We have previously investigated various ML methods for predicting enthalpies of formation of neutral molecules.^{14,17} The best-performing method was based on an FCHL- Δ ML approach. The FCHL method is based on the Kernel-Ridge regression model.¹⁵ In this model, molecular properties are learned using a many-body expansion distribution of the local environment. Each atom of a molecule is represented as a weighted sum of Gaussian distributions. The approach measures the distance between molecules by summing the distances between each atom in the two molecules. The distance between two atomic environments is measured using a weighted sum of the distances between each many-body expansion, based on geometries and local chemical environment. The FCHL- Δ approach is a variant of FCHL in that it learns the difference between properties computed with high- and low-fidelity quantum chemical methods on a training set of molecules.

In this paper, we trained the FCHL- Δ model using the difference between G4MP2A and two density functional methods (B3LYP and ω B97XD) for atomization energies. The training sets were derived from the previously published

Table 6. Mean Absolute Deviations (kcal/mol) of B3LYP and ω B97XD with and without Atom-Specific Corrections^a

		mean absolute deviations (kcal/mol) with respect to experiment			
	test set	B3LYP/6-31(2df,p)	B3LYPA/6-31(2df,p) ^b	ω B97XD/6-311+G(3df,2p)	ω B97XD-A/6-311+G(3df,2p) ^b
G3/99	all (222)	4.52	2.92	2.05	1.88
	nonhydrogens (47)	7.58	4.98	3.40	3.76
	hydrocarbons (38)	3.29	2.11	1.45	1.28
	substituted hydrocarbons (91)	3.78	2.07	1.57	1.19
	inorganic hydrides (15)	2.70	3.09	1.60	1.91
PDS(9)	radicals (31)	4.43	3.22	2.39	1.79
	all (459)	3.99	2.09	1.85	1.28
	substituted hydrocarbons	4.74	2.16	2.16	1.43
	hydrocarbons (175)	2.77	1.35	1.35	1.05
	all (681)	4.16	2.36	1.92	1.48
G3/99 + PDS(9)	all (185)	4.77 (25.2)	4.01 (−15.37)	2.96 (13.75)	2.02 (−8.55)
	hydrocarbons (68)	3.77	5.03	3.01	1.96
	substituted hydrocarbons (117)	5.35	3.42	2.93	2.05
	all aromatics (87)	4.08	2.61	3.17	2.13

^aThe B3LYP and ω B97XD results are from the references in the footnote a of Table 2. ^bThe atom-specific correction is optimized on both G33/99 and PDS(9) sets combined. The ω B97XD-A/6-311+G(3df,2p) method is based on B3LYP/6-31G(2df,p) geometries.

energies for 133 K molecules at the G4MP2 and B3LYP/6-31G(2df) levels of theory.¹³ For G4MP2A, these energies were modified by the atom-specific corrections derived from the set of 681 molecules as described above. The FCHL- Δ model was used with these G4MP2A and B3LYP energies and is referred to as ML-B3LYP-A. We expect this use of the improved G4MP2A energies to result in the FCHL- Δ trained ML models achieving energy accuracies close to those of “G4”, rather than “G4MP2”, since G4MP2A has approximately the same accuracy as G4 theory, as shown in the last section. In addition, we have also trained the FCHL- Δ model using the difference between G4MP2A and ω B97XD atomization energies for comparison with results using B3LYP. This method is referred to as ML- ω B97XD-A. The FCHL- Δ trained ML models used in this work are available on DLHub³¹ and the scripts used to perform parameter optimization and their outputs are posted to GitHub.³²

The ML results for calculations of enthalpies of formation are summarized in Table 7. In the FCHL- Δ trained ML model for which results are shown in Table 7, we used a training set of 10,000 molecules taken from 13,026 molecules selected randomly from the GDB-9 set of 133 K molecules. The ML-B3LYP-A model gave a small MAD of 0.11 kcal/mol between the ML-predicted energies and the G4MP2A energies. This was based on the hold-out set (10%) of 1302 molecules from the pool of 13,026. The ML- ω B97XD-A model gave a slightly smaller MAD of 0.09 kcal/mol between the ML-predicted energies and the G4MP2A energies. For both the ML-B3LYP-A and ML- ω B97XD-A methods, the use of G4MP2 without atom-specific corrections (ML-B3LYP and ML- ω B97XD in Table 7) gave the same MADs as G4MP2A with atom-specific corrections. In addition, the use of B3LYP or ω B97XD with atomic-specific corrections gave no improvement (Table 7). Overall, we find that the predicted energies from ML-B3LYP-A and ML- ω B97XD-A based on G4MP2A energies are more accurate than ML-B3LYP and ML- ω B97XD since the MADs of the former are with respect to G4MP2A energies instead of G4MP2.

Table 7 also gives the results for the 185 larger molecules of the PDS(10–14) dataset using the same methods described

Table 7. Prediction of G4MP2A Atomization Energies Using Δ -FCHL Machine Learning Method Developed with 13,026 Molecules Randomly Selected from the 133 K Molecules of GDB-9 (ref 14)^{a,h}

method ^a	energy data for ML training	test set	MAD of Δ FCHL, kcal/mol ^b
ML-B3LYP-A	B3LYP/G4MP2A	1302 from GDB-9 pool	0.11 (0.11) ^c
ML-B3LYP	B3LYP/G4MP2		0.11 ^{d,f}
ML- ω B97XD-A	ω B97XD/G4MP2A		0.09
ML- ω B97XD	ω B97XD/G4MP2		0.09 ^{e,g}
ML-B3LYP-A	B3LYP/G4MP2A	185 from PDS(10–14)	0.43 (0.43) ^c
ML-B3LYP	B3LYP/G4MP2		0.44 ^d
ML- ω B97XD-A	ω B97XD/G4MP2A		0.34
ML- ω B97XD	ω B97XD/G4MP2		0.34 ^e

^aMethod is abbreviated as ML-DFT, where ML = Δ FCHL machine learning, DFT is the specified DFT method (see Table 5 for basis set), and -A means that the method is based on G4MP2A with atom-specific corrections. A method without -A in its name is based on G4MP2 energies. ^bMean absolute deviation (MAD) of Δ FCHL prediction from the G4MP2A (or G4MP2) result for the test set. ^cValue in parentheses is based on B3LYPA energies instead of B3LYP, showing that there is no difference. ^dBased on G4MP2 and B3LYP energies with no atomic-specific corrections. ^eBased on G4MP2 and ω B97XD energies with no atomic-specific corrections. ^fThis was incorrectly reported as 0.18 kcal/mol in ref 17. ^gThis was incorrectly reported as 0.12 kcal/mol in ref 17. ^hTraining set size = 10,000 from the 13,026 selected sample.

above for the GDB-9 set to determine the effect of use of the ML models on molecules larger than in the training set. The performance is not as good when applied to molecules larger than in the training set. The ML-B3LYP-A model gave a MAD of 0.44 kcal/mol between the ML-predicted energies and the G4MP2A energies for the 185 molecules. The ML- ω B97XD-A model gave a somewhat smaller MAD of 0.34 kcal/mol between the ML-predicted energies and the G4MP2A energies. The errors for the 185 molecules for both models are given in Table S1.

An analysis of the molecules with larger errors contributing to the larger MAD indicates that ML-B3LYP-A had 12

molecules with a deviation of 1 kcal/mol or more. Of these, five were over 2 kcal/mol: 1,4-dinitrosopiperazine, di-tert-butyl diazene, *N,N,N',N'*-tetrafluoro-1,1-heptanediamine, tert-butylperoxide, and trinitromethane. The ML- ω B97XD-A model had eight with errors over 1 kcal/mol, of which only two were over 2 kcal/mol: 1,4-dinitropiperazine and pentafluoropropanoic acid methyl ester. A likely reason for the larger MADs seen for the listed molecules is the inclusion of new bonding types that are not present in the set of molecules with nine nonhydrogen atoms. For example, the molecules with larger errors noted above include ones with multiple nitro groups or many fluorines, which are molecular types that are not included in the smaller molecule set. In future work, we will investigate how to improve the ML models to make them less dependent on molecule size.

CONCLUSIONS

In this paper, we investigated the cause of the increase in error in G4MP2 theory when applied to larger organic molecules with 10 or more nonhydrogen atoms. The following conclusions can be drawn from this study:

- (1) A key source of error is found to be the “higher-level correction”, which is meant to correct for deficiencies in correlation contributions to the energies. The problem is that the higher-level correction is assumed to be independent of the element and the type of bonding involved.
- (2) To address this problem, we have added an atom-specific correction to the higher-level correction, which is only dependent on the type of atom, not on the bond type. With this simple addition to the higher-level correction, we find that a modified G4MP2 method, referred to as G4MP2A, is as accurate as G4 theory for the set of molecules considered here, with a much smaller computational cost. Upon addition of this correction, the G4MP2 MAD is reduced from 0.84 to 0.70 kcal/mol for the 681 molecules with up to nine nonhydrogen molecules and from 1.43 to 0.96 kcal/mol for the 185 molecules with 10–14 nonhydrogen atoms.
- (3) The G4MP2A method is also more accurate in predicting ionization potentials and electron affinities of molecules. The G4MP2 MAD is improved from 0.93 to 0.79 eV for ionization potentials and from 1.05 to 0.83 for electron affinities.
- (4) Finally, we have implemented the G4MP2A energies in a Δ -learning ML method that, when trained on density functional and G4MP2A energies, is shown to predict molecular energies with improved accuracy with respect to experiment than comparable ML methods based on G4MP2 reported previously.¹⁷ However, it is found that extension to larger molecules than in the training set results in an increase in error, especially for larger molecules with new bonding types.

ASSOCIATED CONTENT

Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/acs.jpca.2c01327>.

Comparison of methods for all PDS(10–14) molecules, molecules with large errors in PDS(10–14), and atom-specific corrections for cations and anions (PDF)

AUTHOR INFORMATION

Corresponding Authors

Naveen K. Dandu – Materials Science Division, Argonne National Laboratory, Lemont, Illinois 60439, United States; Joint Center for Energy Storage Research, Argonne National Laboratory, Lemont, Illinois 60439, United States; Chemical Engineering Department, University of Illinois-Chicago, Chicago, Illinois 60607, United States; orcid.org/0000-0001-7122-8537; Email: ndandu@anl.gov

Larry A. Curtiss – Materials Science Division, Argonne National Laboratory, Lemont, Illinois 60439, United States; Joint Center for Energy Storage Research, Argonne National Laboratory, Lemont, Illinois 60439, United States; orcid.org/0000-0001-8855-8006; Email: curtiss@anl.gov

Authors

Rajeev S. Assary – Materials Science Division, Argonne National Laboratory, Lemont, Illinois 60439, United States; Joint Center for Energy Storage Research, Argonne National Laboratory, Lemont, Illinois 60439, United States; orcid.org/0000-0002-9571-3307

Paul C. Redfern – Materials Science Division, Argonne National Laboratory, Lemont, Illinois 60439, United States
Logan Ward – Joint Center for Energy Storage Research and Data Science and Learning Division, Argonne National Laboratory, Lemont, Illinois 60439, United States; orcid.org/0000-0002-1323-5939

Ian Foster – Data Science and Learning Division, Argonne National Laboratory, Lemont, Illinois 60439, United States; Department of Computer Science, University of Chicago, Chicago, Illinois 60637, United States

Complete contact information is available at: <https://pubs.acs.org/10.1021/acs.jpca.2c01327>

Notes

The authors declare no competing financial interest.

ACKNOWLEDGMENTS

This work was supported by the Joint Center for Energy Storage Research (JCESR), an Energy Innovation Hub funded by the U.S. Department of Energy, Office of Science, Basic Energy Sciences. The authors acknowledge a generous grant of computer time from the ANL Laboratory Computing Resource Center (Bebop).

REFERENCES

- (1) Curtiss, L. A.; Redfern, P. C.; Raghavachari, K. Gaussian-4 Theory. *J. Chem. Phys.* **2007**, *126*, No. 084108.
- (2) Curtiss, L. A.; Redfern, P. C.; Raghavachari, K. G *n* Theory. *WIREs Comput. Mol. Sci.* **2011**, *1*, 810–825.
- (3) Curtiss, L. A.; Redfern, P. C.; Raghavachari, K. Gaussian-4 Theory Using Reduced Order Perturbation Theory. *J. Chem. Phys.* **2007**, *127*, No. 124105.
- (4) DeYonker, N. J.; Cundari, T. R.; Wilson, A. K. The Correlation Consistent Composite Approach (Cc CA): An Alternative to the Gaussian-*n* Methods. *J. Chem. Phys.* **2006**, *124*, No. 114104.
- (5) Montgomery, J. A.; Frisch, M. J.; Ochterski, J. W.; Petersson, G. A. A Complete Basis Set Model Chemistry. VI. Use of Density Functional Geometries and Frequencies. *J. Chem. Phys.* **1999**, *110*, 2822–2827.
- (6) Feller, D.; Dixon, D. A. Extended Benchmark Studies of Coupled Cluster Theory through Triple Excitations. *J. Chem. Phys.* **2001**, *115*, 3484–3496.

- (7) Chan, B.; Radom, L. W3X: A Cost-Effective Post-CCSD(T) Composite Procedure. *J. Chem. Theory Comput.* **2013**, *9*, 4769–4778.
- (8) Karton, A.; Martin, J. M. L. Explicitly Correlated W_n Theory: W1-F12 and W2-F12. *J. Chem. Phys.* **2012**, *136*, No. 124114.
- (9) Zhao, Y.; Xia, L.; Liao, X.; He, Q.; Zhao, M. X.; Truhlar, D. G. Extrapolation of High-Order Correlation Energies: The WMS Model. *Phys. Chem. Chem. Phys.* **2018**, *20*, 27375–27384.
- (10) Chan, B.; Deng, J.; Radom, L. G4(MP2)-6X: A Cost-Effective Improvement to G4(MP2). *J. Chem. Theory Comput.* **2011**, *7*, 112–120.
- (11) Chan, B.; Karton, A.; Raghavachari, K. G4(MP2)-XK: A Variant of the G4(MP2)-6X Composite Method with Expanded Applicability for Main-Group Elements up to Radon. *J. Chem. Theory Comput.* **2019**, *15*, 4478–4484.
- (12) Semidalas, E.; Martin, J. M. L. Canonical and DLPNO-Based G4(MP2)XK-Inspired Composite Wave Function Methods Parametrized against Large and Chemically Diverse Training Sets: Are They More Accurate and/or Robust than Double-Hybrid DFT? *J. Chem. Theory Comput.* **2020**, *16*, 4238–4255.
- (13) Narayanan, B.; Redfern, P. C.; Assary, R. S.; Curtiss, L. A. Accurate Quantum Chemical Energies for 133 000 Organic Molecules. *Chem. Sci.* **2019**, *10*, 7449–7455.
- (14) Ward, L.; Blaiszik, B.; Foster, I.; Assary, R. S.; Narayanan, B.; Curtiss, L. Machine Learning Prediction of Accurate Atomization Energies of Organic Molecules from Low-Fidelity Quantum Chemical Calculations. *MRS Commun.* **2019**, *9*, 891–899.
- (15) Faber, F. A.; Christensen, A. S.; Huang, B.; von Lilienfeld, O. A. Alchemical and Structural Distribution Based Representation for Universal Quantum Machine Learning. *J. Chem. Phys.* **2018**, *148*, No. 241717.
- (16) Schütt, K. T.; Sauceda, H. E.; Kindermans, P.-J.; Tkatchenko, A.; Müller, K.-R. SchNet – A Deep Learning Architecture for Molecules and Materials. *J. Chem. Phys.* **2018**, *148*, No. 241722.
- (17) Dandu, N.; Ward, L.; Assary, R. S.; Redfern, P. C.; Narayanan, B.; Foster, I. T.; Curtiss, L. A. Quantum-Chemically Informed Machine Learning: Prediction of Energies of Organic Molecules with 10 to 14 Non-Hydrogen Atoms. *J. Phys. Chem. A* **2020**, *124*, 5804–5811.
- (18) Csonka, G. I.; Ruzsinszky, A.; Tao, J.; Perdew, J. P. Energies of Organic Molecules and Atoms in Density Functional Theory. *Int. J. Quantum Chem.* **2005**, *101*, 506–511.
- (19) Frisch, M. J.; Trucks, G. W.; Cheeseman, J. R.; Scalmani, G.; Caricato, M.; Hratchian, H. P.; Li, X.; Barone, V.; Bloino, J.; Zheng, G. et al. *Gaussian 16*; Gaussian, Inc.: Wallingford CT, 2016.
- (20) Becke, A. D. Density-functional Thermochemistry. III. The Role of Exact Exchange. *J. Chem. Phys.* **1993**, *98*, 5648–5652.
- (21) Becke, A. D. A New Mixing of Hartree–Fock and Local Density-functional Theories. *J. Chem. Phys.* **1993**, *98*, 1372–1377.
- (22) Mardirossian, N.; Head-Gordon, M. Thirty Years of Density Functional Theory in Computational Chemistry: An Overview and Extensive Assessment of 200 Density Functionals. *Mol. Phys.* **2017**, *115*, 2315–2372.
- (23) Curtiss, L. A.; Raghavachari, K.; Redfern, P. C.; Pople, J. A. Assessment of Gaussian-3 and Density Functional Theories for a Larger Experimental Test Set. *J. Chem. Phys.* **2005**, *123*, 7374–7383.
- (24) Pedley, J. B.; Naylor, R. D.; Kirby, S. P. *Thermochemical Data and Structures of Organic Compounds*; Springer: Netherlands, 2012.
- (25) Raghavachari, K.; Stefanov, B. B.; Curtiss, L. A. Accurate Thermochemistry for Larger Molecules: Gaussian-2 Theory with Bond Separation Energies. *J. Chem. Phys.* **1997**, *106*, 6764–6767.
- (26) Curtiss, L. A.; Raghavachari, K.; Redfern, P. C.; Pople, J. A. Assessment of Gaussian-2 and Density Functional Theories for the Computation of Enthalpies of Formation. *J. Chem. Phys.* **1997**, *106*, 1063–1079.
- (27) Curtiss, L. A.; Raghavachari, K.; Trucks, G. W.; Pople, J. A. Gaussian-2 Theory for Molecular Energies of First- and Second-row Compounds. *J. Chem. Phys.* **1991**, *94*, 7221–7230.
- (28) Pracht, P.; Bohle, F.; Grimme, S. Automated Exploration of the Low-Energy Chemical Space with Fast Quantum Chemical Methods. *Phys. Chem. Chem. Phys.* **2020**, *22*, 7169–7192.
- (29) Bannwarth, C.; Ehlert, S.; Grimme, S. GFN2-XTB—An Accurate and Broadly Parametrized Self-Consistent Tight-Binding Quantum Chemical Method with Multipole Electrostatics and Density-Dependent Dispersion Contributions. *J. Chem. Theory Comput.* **2019**, *15*, 1652–1671.
- (30) Pracht, P.; Caldeweyher, E.; Ehlert, S.; Grimme, S. *A Robust Non-Self-Consistent Tight-Binding Quantum Chemistry Method for Large Molecules*, preprint; Chemistry, 2019 DOI: 10.26434/chemrxiv.8326202.v1.
- (31) <https://Github.Com/Globus-Labs/G4mp2-Atomization-Energy/Tree/Master/Dlhub>.
- (32) <https://Github.Com/Globus-Labs/G4mp2-Atomization-Energy>.
- (33) Chan, B.; Trevitt, A. J.; Blanksby, S. J.; Radom, L. Comment on the Ionization Energy of B_2F_4 . *J. Phys. Chem. A* **2012**, *116*, 9214–9215.

Recommended by ACS

Benchmark of GW Methods for Core-Level Binding Energies

Jiachen Li, Dorothea Golze, et al.

NOVEMBER 02, 2022

JOURNAL OF CHEMICAL THEORY AND COMPUTATION

READ 

Recommendation of Orbitals for G_0W_0 Calculations on Molecules and Crystals

Linyao Zhang, Donald G. Truhlar, et al.

MAY 17, 2022

JOURNAL OF CHEMICAL THEORY AND COMPUTATION

READ 

Basis Set Selection for Molecular Core-Level GW Calculations

Daniel Mejia-Rodriguez, Niranjana Govind, et al.

JULY 11, 2022

JOURNAL OF CHEMICAL THEORY AND COMPUTATION

READ 

Large-Scale Benchmarking of Multireference Vertical-Excitation Calculations via Automated Active-Space Selection

Daniel S. King, Laura Gagliardi, et al.

SEPTEMBER 16, 2022

JOURNAL OF CHEMICAL THEORY AND COMPUTATION

READ 

Get More Suggestions >