

## Parametrization of Atomic Energies to Improve Small Basis Set Density Functional Thermochemistry

Edward N. Brothers\* and Gustavo E. Scuseria

*Department of Chemistry, Mail Stop 60, Rice University, Houston, Texas 77005-1892*

Received March 23, 2006

**Abstract:** Enthalpies of formation predicted with density functional theory and small basis sets can be greatly improved by treating the atomic energies as empirical parameters. When a variety of functionals and small basis sets are used, the root-mean-square error in enthalpies of formation is reduced by a factor of approximately two for the least improved functional/basis set pair, with significantly larger reductions for other functionals, especially LSDA. When the 3-21G\* and 3-21+G\* basis sets are used with nonempirical functionals, it is possible to achieve accuracy greater than that of PM3, which was primarily designed to reproduce enthalpies of formation. In addition to decreasing statistical errors, our procedure can also remove qualitative errors in density functional/basis set pairs that fail for the prediction of enthalpies of formation.

### I. Introduction

This paper shows that combining the improvement in enthalpies of formation by optimizing atomic energies such as was done recently by Csonka et al.<sup>1</sup> with the ability of density functional theory (DFT) to predict enthalpies of formation with small basis sets<sup>2,3</sup> results in inexpensive density functional thermochemistry comparable to, or better than, semiempirical methods. This is notable in two regards; first, (very) small basis set density functional theory<sup>4</sup> has never (to our knowledge) outperformed PM3,<sup>5</sup> and second, it takes the corrective ability of parametrized atom energies beyond simply tightening error bars and actually removes qualitative errors.

The optimization of atomic energies to reduce thermochemical error was recently undertaken by Csonka et al.,<sup>1</sup> and they demonstrated that the errors in enthalpies of formation predicted by DFT can be partially attributed to problems with predicted atomic energies. When fixed geometries and tabulated frequency corrections are used with fairly large basis sets, this optimization substantially reduced the errors in enthalpies of formation for a variety of previously difficult molecules in the G3/99 set of compounds.<sup>6</sup>

In addition to improving enthalpies of formation, there is a second point that can be taken from that study which is

more subtle. It may be possible to create a set of atomic energies such that functionals heretofore considered unacceptable for thermochemistry because of large errors can in fact be useful. In other words, the atomic energy fitting procedure could correct qualitatively wrong results by removing the major impediment to calculating enthalpies of formation.

Concurrent with the atomic energy work cited above were two studies which demonstrated that DFT can predict enthalpies of formation accurately with some of the smallest common basis sets. In the first study, reasonable enthalpies of formation for many functionals were obtained, providing results comparable with semiempirical predictions while using geometries, energies, and frequencies all determined with the small basis sets.<sup>2</sup> In this case, all of the functionals that provided high-quality results were all based on both the density and the gradient of the density, that is, the generalized gradient approximation, or GGA. Meta-GGAs, which include terms based on the kinetic energy density, were not considered in that study, although they are included here.

The second study<sup>3</sup> developed fully an idea first examined in the course of analyzing small basis density functional thermochemistry.<sup>2</sup> LSDA,<sup>7</sup> which contains no gradient correction, was improved for a wide variety of properties through the use of an empirical parameter to scale the correlation, with special emphasis given to performance with small basis sets. This method was termed “cSVWN5”. It is

\* Corresponding author e-mail: enb@rice.edu.

useful for investigating large systems because small basis set methods are intrinsically fast, density-only functionals are slightly faster than GGAs and meta-GGAs, and more complicated theories are CPU-intensive.

In this paper, atomic energy optimization is applied to small basis set DFT, including the functional developed especially to work with small basis sets, to greatly improve thermochemical prediction.

## II. Test Set and Computational Method

The G3/99 set of Curtiss and co-workers was selected here for use as a test bed because it has become a common standard for determining the accuracy of quantum chemical approaches. In total, there are 13 atom types and 223 compounds used to examine enthalpies of formation in G3/99; however, five of these atom types appear in three or less compounds, specifically lithium, beryllium, sodium, aluminum, and boron. These atom types and the compounds containing them were thus removed from the set to avoid creating biased parameters for those atom types, resulting in a total test set of 213 compounds consisting of nine atom types.

The basis sets chosen for this study were STO-3G,<sup>8</sup> 3-21G\*, and 3-21+G\*,<sup>9</sup> which are the smallest basis sets in common use.<sup>10</sup> It is important to note that the “\*” on 3-21G\* and 3-21+G\* denote placing a *d* function on atoms heavier than neon and not on all non-hydrogen atoms, as is the case with other Pople basis sets. Also, these basis sets use the default Cartesian basis functions; that is, they use 6*d* rather than 5*d* components. For molecules containing atoms larger than neon, any basis-set-specific parameters, such as the atomic energies in this study, would have to be adjusted to compensate for the change in basis if 5*d* was desired.

Several categories of functionals are represented in this study. The two density-only functionals used in this study are LSDA, which uses the local correlation functional of Vosko, Wilk, and Nusair,<sup>7</sup> and cSVWN5,<sup>3</sup> which is equivalent to the LSDA used in this study with the local correlation scaled by 0.3. cSVWN5 was optimized by Riley et al. for use with 3-21G\* and 3-21+G\* and, thus, is neglected for STO-3G in this paper. The GGA PBE<sup>11</sup> developed by Perdew, Burke, and Ernzerhof, and the meta-GGA TPSS<sup>12</sup> created by Tao, Perdew, Staroverov, and Scuseria, are also examined. PBEh,<sup>13,14</sup> which is PBE with 25% of the functional exchange replaced by exact exchange, and TPSSh,<sup>15</sup> which uses 10% exact exchange, are also evaluated. While not a density functional, HF<sup>16</sup> is included for comparison purposes. Note that, with the exception of cSVWN5, none of the functionals were developed by fitting internal parameters to enthalpies of formation or other experimental data; that is, they are nonempirical.

All calculations were performed in the Gaussian suite of programs.<sup>17</sup> For all of the methods used in this study, geometries were optimized and frequencies calculated at the method of interest; that is, the energies were functional/basis//functional/basis throughout. Gaussian defaults were used in all of the calculations, with the exception of the integration grid, which was a pruned (99,590), or “ultrafine”, grid.

After the data was collected, three separate parametrizations were attempted. First, a single multiplicative scaling parameter was used with all calculated atomic energies to see if the fit could be accomplished with one parameter. This fit was also attempted starting from the correct total atomic energies.<sup>18</sup> Finally, a full genetic algorithm<sup>19</sup> (GA) optimization was undertaken in which each atomic energy was treated as an empirical parameter and all parameters were simultaneously fit versus the entire set of compounds in this study. This optimization procedure was selected because we desired to optimize all of the parameters simultaneously. With nine parameters for each functional/basis pair, a grid search would be out of the question. Also, it was unknown at the beginning of the study how close to the final parameters the initial values were, and thus, any method consisting of multiple line searches would have to be restarted at a variety of different inertial points, which is a problem easily avoided by using a GA. Thus, a GA was selected for this problem. Note that all of the compounds were used, rather than a “jack-knife” set, in which fitting would be conducted versus some compounds and evaluated versus a larger set which includes the set parametrized against. For the purposes of optimization, the root-mean-square (RMS) error was treated as the error to be minimized, as this biases the parametrization to pull in the furthest outliers first.

## III. Results

Before discussing the results of the parametrization, it is necessary to briefly mention what the optimized atomic energies represent. Optimization does not necessarily move the atomic energies closer to the exact values. (A list of the difference between the exact energies<sup>18</sup> and the calculated and parametrized energies using carbon as an example is available in the Supporting Information.) The difference between exact and calculated energies ranges up to 0.8 au, and the difference between the optimized and original atomic energies is small relative to the difference between the exact and calculated values. The parameters also compensate for issues with the basis set. There are several functionals which when used with large basis sets already produce very good enthalpies of formation,<sup>6</sup> but by using small basis sets such as the ones considered in this study, the parameters are forced to deal with both functional shortcomings and the paucity of the basis set. Therefore, because the parameters do not represent an improvement in atomic energies versus exact values and because they are basis-set-specific, it would be erroneous to assign them any physical interpretation. They are simply empirical parameters whose strength is their efficacy.

Attempts were made to scale the atomic energies, meaning that all atomic energies were multiplied by a single parameter. Scaling exact atomic energies did not improve accuracy, and in fact, the results were worse than those using the original atomic energies. Scaling the calculated energies with a single parameter does improve the enthalpies of formation slightly, but the improvement is marginal at best. Thus, this data is not presented; it is mentioned to explain why using one parameter for each atomic energy was undertaken. The balance of this paper will deal with the outcome of the GA optimization.

**Table 1.** Mean Error (ME), Mean Absolute Error (MAE), Root-Mean-Squared Error (RMS), and Standard Deviation (SD) for Enthalpies of Formation in the G3/99 Set Using Original and Optimized Atomic Energies<sup>a</sup>

basis	method	original				GA optimized			
		ME	MAE	RMS	SD	ME	MAE	RMS	SD
STO-3G	LSDA	-268.0	268.8	335.3	202.0	2.5	20.1	29.1	29.1
	PBE	-145.5	151.6	190.9	123.8	0.7	17.5	27.3	27.3
	PBEh	-111.7	126.7	165.1	121.9	0.5	18.7	29.1	29.1
	TPSS	-118.4	128.0	164.0	113.8	0.1	17.3	27.4	27.5
	TPSSh	-108.5	121.9	158.3	115.5	0.1	17.7	28.1	28.1
	HF	140.4	142.2	171.3	98.5	-1.7	22.0	36.1	36.2
3-21G*	LSDA	-113.5	113.5	136.1	75.2	2.8	7.3	9.5	9.1
	PBE	-13.2	15.5	19.8	14.9	0.9	4.3	6.1	6.1
	PBEh	15.0	15.2	19.8	12.9	1.0	4.4	5.8	5.7
	TPSS	11.5	13.5	17.2	12.8	0.3	4.1	6.1	6.1
	TPSSh	19.4	20.5	25.2	16.1	0.4	4.0	5.7	5.7
	HF	261.5	261.5	298.9	145.1	-1.5	7.8	11.7	11.7
3-21+G*	cSVWN5	-7.4	19.1	28.0	27.1	1.8	6.0	7.9	7.7
	LSDA	-100.8	100.8	121.7	68.4	2.1	6.3	8.1	7.8
	PBE	-1.0	7.5	9.9	9.9	0.3	3.6	5.7	5.7
	PBEh	24.2	24.2	29.0	16.0	0.5	3.8	5.1	5.1
	TPSS	21.5	22.5	27.2	16.8	-0.3	3.7	5.8	5.9
	TPSSh	28.5	29.2	35.0	20.5	-0.1	3.5	5.3	5.3
	HF	266.4	266.4	304.7	148.2	-1.8	7.9	11.7	11.6
	cSVWN5	7.9	16.3	19.3	17.7	1.1	4.9	6.6	6.5
	PM3	-1.0	6.9	9.5	9.4				

<sup>a</sup> All values are in kcal/mol. PM3 is included on the last line for comparison purposes.

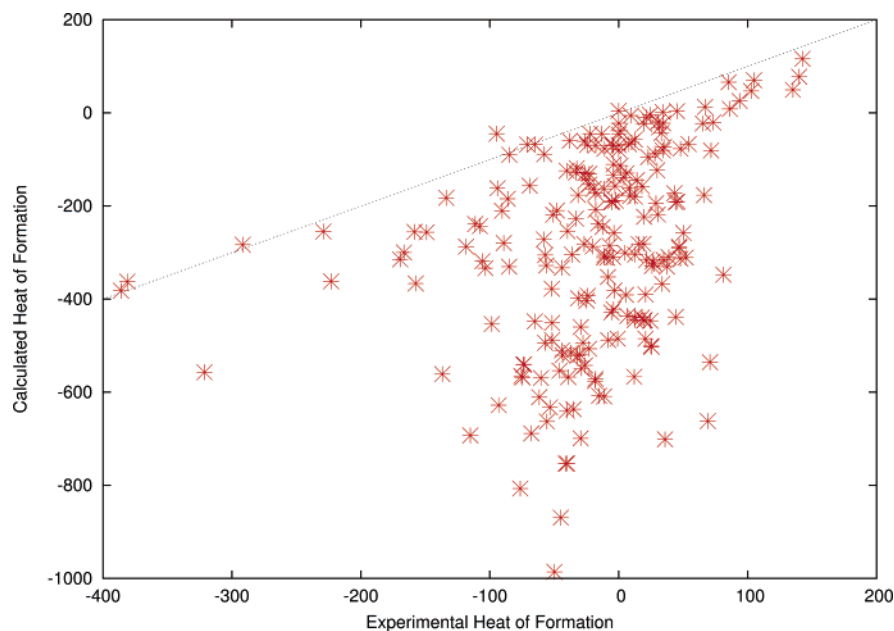
The results of the GA optimization listed in Table 1 demonstrate the effectiveness of the optimization, especially in light of how poor many of the original results are. By any reasonable criteria, the enthalpies of formation calculated for STO-3G are a failure, with enormous systematic errors due to overbinding for all functionals and underbinding for HF. Also, at any of the three basis sets considered in this study, LSDA fails. The best functional/basis set pair with regular energies is PBE/3-21+G\*. This occurs because of a cancellation of errors, as the 3-21G\* results show that PBE underbinds while PBEh, TPSS, and TPSSh overbind, and the addition of the diffuse function increases binding for all functionals, resulting in the good performance of PBE/3-21+G\*.

In contrast, after the parametrization, the errors that result from systematic overbinding or underbinding (exhibited through nonzero mean errors) are nearly removed. The performance of LSDA is also brought much closer to the performance of the GGA and meta-GGA functionals, especially when the comparison is done with STO-3G. In fact, the errors after parametrization are small enough that no functional can be said to fail for enthalpies of formation. This is not to say that all basis sets are equally well suited for thermochemistry as long as parameters are present but rather that the parametrized methods at STO-3G are no longer qualitatively incorrect. This improvement from qualitatively incorrect to qualitatively correct can be most easily seen by examining Figures 1 and 2. With the original atomic energies, the enthalpies of formation are predicted to be hundreds of kilocalories per mole low and are shown in the graph to correspond very poorly to the experimental values, appearing almost randomly scattered. The parametrized results qualitatively correspond to the experimental values, albeit still with significant errors, similar in size to TPSS with 3-21+G\* and no parameters.

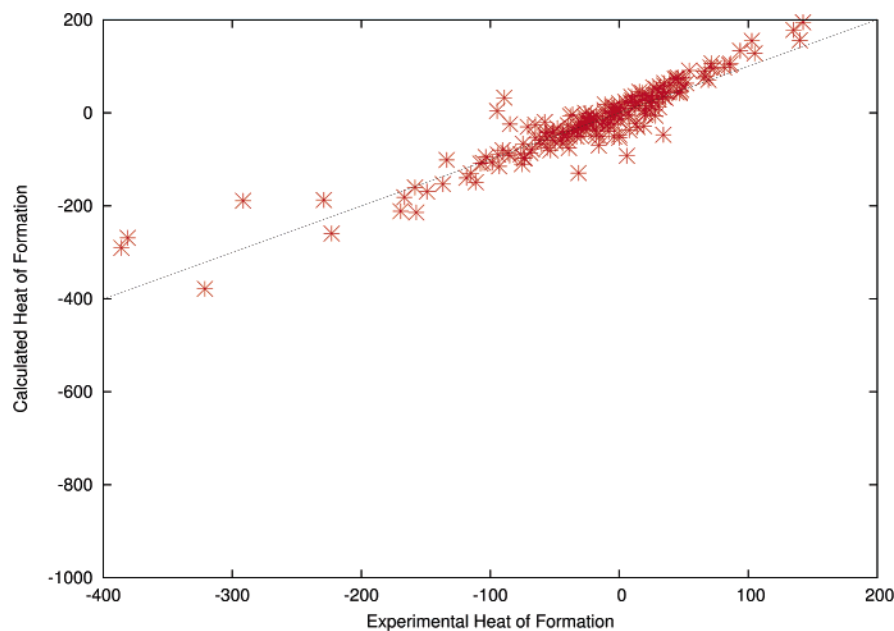
Another measure of the success of the parametrization is the comparison with PM3 results.<sup>5</sup> PM3 was parametrized primarily versus enthalpies of formation and is consistently more accurate than DFT with small basis sets and standard atomic energies.<sup>2</sup> PM3 still outperforms all of the STO-3G results, but functionals with split-valence basis sets are more accurate than PM3 after atomic energy optimization. Thus, using nonempirical functionals and small basis sets and correcting the errors in atomic energies allows thermochemical accuracy as high as that of semiempirical methods, which use far more parameters.

Performance is still determined by the basis set and functional, with the larger basis sets and functionals which consider more density-related quantities providing higher accuracy both before and after parametrization. For any optimized result in this study, the RMS for STO-3G is larger than that for 3-21G\* and larger for 3-21G\* than for 3-21+G\*. The trend for mean absolute error follows the same trend for all methods except HF, and the deviation from the trend for HF is negligible. (More information on parametrization to improve HF results can be found in the work of Ruzsinszky and Csonka.<sup>20</sup>) Thus, the new atomic energies do not alter the fact that bigger basis sets work better. Also, while the difference between density-only methods and the more modern GGAs and meta-GGAs is not particularly noticeable at the STO-3G level, with all methods providing errors of around 20 kcal/mol, there is a clear advantage to PBE, TPSS, and their hybrids with the split-valence basis sets.

It should be noted that the inclusion of exact exchange with small basis sets does not improve the thermochemistry without the atomic energy parametrization, and with parametrization, the pure functionals are outperformed by the hybrids only with the split-valence basis sets. Also, HF even after parametrization does worse than any of the DFT



**Figure 1.** Enthalpies of formation calculated using LSDA/STO-3G. A line with a slope of unity is added to guide the eye.



**Figure 2.** Enthalpies of formation calculated using LSDA/STO-3G and the optimized atomic energies. A line with a slope of unity is added to guide the eye.

methods. This can be attributed to the fact that, with small basis sets, exact exchange is a hindrance,<sup>2</sup> unlike in large basis sets, where it can improve enthalpies of formation greatly.<sup>6</sup>

The only empirical functional presented in this study (cSVWN5) has errors approximately halfway between the best functionals and LSDA after atomic energy optimization. It thus performs better than LSDA, which it was developed from, as well as PM3, and is slightly cheaper than functionals that include terms other than the density. This makes it ideal for investigating large systems.

To examine the size dependence of the errors, the error in the predicted enthalpy of formation was plotted versus the number of carbon atoms for the first eight alkanes. This curve was then fit to a line, and the slope of the line is

presented in Table 2. In this case, the larger the magnitude of the slope, the greater the size dependence of the error. With the split-valence basis sets and GGA, meta-GGA, hybrid functionals, and the empirical functional cSVWN5, the original size dependence of the error is small to begin with, and in most cases, the optimization reduces it further. The exceptions to this are cSVWN5/3-21G\* and PBEh/3-21G\*, and in both of those cases, the optimized values are still relatively small. For these functionals with STO-3G, the size dependence of the error is large and is reduced greatly by the parametrization. LSDA and HF also have extremely large initial slopes and are compensated for by the parametrization. For these methods, this improvement is probably due to the removal of large errors throughout the total test set. Finally, the large initial dependence of HF on size is



**Table 2.** Slope of the Line Created by Plotting the Number of Carbons versus the Difference between Experimental and Calculated Enthalpies of Formation for the First Eight Alkanes

method	original			GA optimized		
	STO-3G	3-21G*	3-21+G*	STO-3G	3-21G*	3-21+G*
LSDA	-111.6	-41.3	-37.8	-3.5	-2.0	-1.4
PBE	-66.0	-5.8	-2.3	-2.5	-0.9	-0.2
PBEh	-60.5	-0.4	2.4	-2.7	-0.9	-0.3
TPSS	-58.7	0.6	3.7	-2.0	-0.3	0.4
TPSSh	-57.7	1.6	4.5	-2.1	-0.3	0.3
HF	16.4	78.9	81.1	-2.0	0.8	1.5
cSVWN5		-0.8	3.3		-1.6	-1.0

due to the incompleteness of the these basis sets relative to the size necessary to converge exact exchange.

The differences between optimized and regular atomic energies are listed in the Supporting Information. Several of the series of corrections are all negative, in which case they are correcting overbinding, and the HF parameters are almost all positive, to correct underbinding. Methods without large mean errors, implying no large systematic errors, have a mixture of positive and negative values. As the values themselves are only interesting for implementation, they are omitted here.

#### IV. Conclusion

The use of optimized atomic energies to calculate enthalpies of formation calculated with small basis sets results in significant improvements when compared to experimental values. The improvements are large enough to allow small basis set density functional methods to achieve higher accuracy than PM3 for the first time. The improvement is even greater for STO-3G calculations, as using this basis previously gave enthalpies of formation that were off by hundreds of kilocalories per mole, and with atomic energy parameters, STO-3G now provides qualitatively correct values.

**Acknowledgment.** This work was supported by NSF-CHE-0457030 and the Welch Foundation. E.N.B. would like to acknowledge helpful comments from Nicole Brothers.

**Supporting Information Available:** The Supporting Information for this paper consists of the optimized parameters (three tables) and the difference of carbon atomic energies from the exact values both before and after parametrization (one table). This information is available free of charge via the Internet at <http://pubs.acs.org>.

#### References

- (1) Csonka, G. I.; Ruzsinszky, A.; Tao, J.; Perdew, J. P. *Int. J. Quantum Chem.* **2005**, *101*, 506.
- (2) Brothers, E. N.; Merz, K. M., Jr. *J. Phys. Chem. A* **2004**, *108*, 2904.
- (3) Riley, K. E.; Brothers, E. N.; Ayers, K. B.; Merz, K. M., Jr. *J. Chem. Theory Comput.* **2005**, *1*, 546.
- (4) Hohenberg, P.; Kohn, W. *Phys. Rev. B* **1964**, *136*, 864. Kohn, W.; Sham, L. J. *Phys. Rev. A* **1965**, *140*, 1133.

- (5) Stewart, J. J. P. *J. Comput. Chem.* **1989**, *10*, 209.
- (6) Curtiss, L. A.; Raghavachari, K.; Redfern, P. C.; Pople, J. A. *J. Chem. Phys.* **2000**, *112*, 7374.
- (7) Vosko, S. H.; Wilk, L.; Nusair, M. *Can. J. Phys.* **1980**, *58*, 1200. LSDA is equivalent to the Gaussian keyword "SVWN5" and refers to using the fifth formula for local correlation proposed in the paper.
- (8) Hehre, W. J.; Stewart, R. F.; Pople, J. A. *J. Chem. Phys.* **1969**, *51*, 2657.
- (9) Binkley, J. S.; Pople, J. A.; Hehre, W. J. *J. Am. Chem. Soc.* **1980**, *102*, 939.
- (10) Note that MIDI! (Easton, R. E.; Giesen, D. J.; Welch, A.; Cramer, C. J.; Truhlar, D. G. *Theor. Chim. Acta* **1996**, *93*, 281) could also have been selected as it is of comparable size and performance to these sets, although slightly larger than 3-21G\* and 3-21+G\*, and has approximately the same performance.
- (11) Perdew, J. P.; Burke, K.; Ernzerhof, M. *Phys. Rev. Lett.* **1996**, *77*, 3865.
- (12) Tao, J.; Perdew, J. P.; Staroverov, V. N.; Scuseria, G. E. *Phys. Rev. Lett.* **2003**, *91*, 146401.
- (13) Ernzerhof, M.; Scuseria, G. E. *J. Chem. Phys.* **1999**, *110*, 5029.
- (14) Adamo, C.; Barone, V. *J. Chem. Phys.* **1999**, *110*, 6158.
- (15) Staroverov, V. N.; Scuseria, G. E.; Tao, J.; Perdew, J. P. *J. Chem. Phys.* **2003**, *119*, 12129.
- (16) Roothaan, C. C. J. *Rev. Mod. Phys.* **1951**, *23*, 69.
- (17) Frisch, M. J.; Trucks, G. W.; Schlegel, H. B.; Scuseria, G. E.; Robb, M. A.; Cheeseman, J. R.; Montgomery, J. A., Jr.; Vreven, T.; Kudin, K. N.; Burant, J. C.; Millam, J. M.; Iyengar, S. S.; Tomasi, J.; Barone, V.; Mennucci, B.; Cossi, M.; Scalmani, G.; Rega, N.; Petersson, G. A.; Nakatsuji, H.; Hada, M.; Ehara, M.; Toyota, K.; Fukuda, R.; Hasegawa, J.; Ishida, M.; Nakajima, T.; Honda, Y.; Kitao, O.; Nakai, H.; Klene, M.; Li, X.; Knox, J. E.; Hratchian, H. P.; Cross, J. B.; Bakken, V.; Adamo, C.; Jaramillo, J.; Gomperts, R.; Stratmann, R. E.; Yazyev, O.; Austin, A. J.; Cammi, R.; Pomelli, C.; Ochterski, J. W.; Ayala, P. Y.; Morokuma, K.; Voth, G. A.; Salvador, P.; Dannenberg, J. J.; Zakrzewski, V. G.; Dapprich, S.; Daniels, A. D.; Strain, M. C.; Farkas, O.; Malick, D. K.; Rabuck, A. D.; Raghavachari, K.; Foresman, J. B.; Ortiz, J. V.; Cui, Q.; Baboul, A. G.; Clifford, S.; Cioslowski, J.; Stefanov, B. B.; Liu, G.; Liashenko, A.; Piskorz, P.; Komaromi, I.; Martin, R. L.; Fox, D. J.; Keith, T.; Al-Laham, M. A.; Peng, C. Y.; Nanayakkara, A.; Challacombe, M.; Gill, P. M. W.; Johnson, B.; Chen, W.; Wong, M. W.; Gonzalez, C.; Pople, J. A. *Gaussian 03*, revision C.01; Gaussian, Inc.: Wallingford, CT, 2004.
- (18) Chakravorty, S. J.; Gwaltney, S. R.; Davidson, E. R.; Parpia, F. A.; Fischer, C. F. P. *Phys. Rev. A* **1993**, *47*, 3649.
- (19) Goldberg, D. *Genetic Algorithms in Search, Optimization and Machine Learning*; Addison-Wesley: San Mateo, CA, 1989.
- (20) Ruzsinszky, A.; Csonka, G. I. *J. Phys. Chem. A* **2003**, *107*, 8687, and references therein.

CT600109X