



A Combined First-principles Calculation and Neural Networks Correction Approach for Evaluating Gibbs Energy of Formation

XiuJung Wang , LiHong Hu , LaiHo Wong & GuanHua Chen

To cite this article: XiuJung Wang , LiHong Hu , LaiHo Wong & GuanHua Chen (2004) A Combined First-principles Calculation and Neural Networks Correction Approach for Evaluating Gibbs Energy of Formation, Molecular Simulation, 30:1, 9-15, DOI: [10.1080/08927020310001631098](https://doi.org/10.1080/08927020310001631098)

To link to this article: <https://doi.org/10.1080/08927020310001631098>



Published online: 15 Aug 2006.



Submit your article to this journal [↗](#)



Article views: 38



View related articles [↗](#)



Citing articles: 17 View citing articles [↗](#)

A Combined First-principles Calculation and Neural Networks Correction Approach for Evaluating Gibbs Energy of Formation

XIUJUNG WANG, LIHONG HU, LAIHO WONG and GUANHUA CHEN*

Department of Chemistry, The University of Hong Kong, Pokfulam Road, Hong Kong, People's Republic of China

(Received September 2003; In final form October 2003)

Despite of their successes, the results of first-principles quantum mechanical calculations contain inherent numerical errors that are caused by inadequate treatment of electron correlation, incompleteness of basis sets, relativistic effects or approximated exchange-correlation functionals. In this work, we develop a combined density-functional theory and neural-network correction (DFT-NEURON) approach to reduce drastically these errors, and apply the resulting approach to determine the standard Gibbs energy of formation ΔG^0 at 298 K for small- and medium-sized organic molecules. The root mean square deviation of the calculated ΔG^0 for 180 molecules is reduced from $22.3 \text{ kcal} \cdot \text{mol}^{-1}$ to $3.0 \text{ kcal} \cdot \text{mol}^{-1}$ for B3LYP/6-311 + G(*d,p*). We examine further the selection of physical descriptors for the neural network.

Keywords: DFT; Neural network; Gibbs energy of formation; First-principles quantum mechanical methods

First-principles quantum mechanical methods have become indispensable research tools in chemistry, condensed matter physics, materials science and molecular biology [1,2]. Experimentalists increasingly rely on these methods to interpret their experimental findings. Despite their successes, first-principles quantum mechanical methods are often not quantitatively accurate enough to predict the experimental measurements, in particular, for large systems. This is caused by the inherent approximations adopted in the first-principles methods. Because of the computational costs, electron correlation has always been a difficult obstacle for first-principles calculations. For instance, highly accurate full configuration interaction

calculations have been limited to very small molecules [3]. Basis sets cannot cover entire physical space, and this introduces inherent computational errors [4]. In practice, limited by the computational resources we often adopt inadequate basis sets for large molecules. Effective core potential is frequently used to approximate the relativistic effects, which leads to inevitably the approximated calculation results for heavy element containing systems. Accuracy of a density-functional theory (DFT) calculation is mainly determined by the exchange-correlation (XC) functional that is employed [2]. No exact XC functional is known. All DFT calculations employ the approximated XC functionals, which lead to further calculation errors. Besides these inherent approximations, quantum mechanical calculations are often limited by finite computational resources and thus adopt less accurate methods. All these contribute to the discrepancies between the calculated and experimental results. An admirable objective in computational science is to predict the properties of matter prior to the experiments. To achieve this, we must eliminate the systematic deviations of the calculation results and reduce greatly the remaining numerical uncertainties. In addition, we need to quantify the accuracies of the numerical methods. Determining the calculation errors solely from the first principles has proven to be an extremely difficult task, and alternatives must be sought.

Despite the various approximations that the first-principles quantum mechanical calculations adopt, the calculated results capture the essence of the physics. For instance, although their absolute values

*Corresponding author. Fax: +852-2857-1586. E-mail: ghc@yangtze.hku.hk

may not agree with the experimental data, the calculated results of different molecules often have the same relative tendency as the corresponding experimental results. To predict a physical property of a material, it may thus be sufficient to correct its raw result from the first-principles calculation. The discrepancies between the calculated and experimental results depend on the physical properties of the material. These physical properties include predominantly the property of interest and, to a lesser degree, other related properties of the materials which can again be evaluated via first-principles calculations. In other words, a quantitative relationship exists between the experimental results and calculated properties. Although it is exceedingly difficult to be determined from the first-principles, the quantitative relationship can be obtained empirically. We propose here a neural-network [5–8] based approach, DFT-NEURON, to determine this quantitative relationship which can be subsequently used to reduce drastically the numerical errors of first-principles calculations. As a demonstration, the standard Gibbs energies of formation ΔG^0 s at 298.15 K of 180 small- or medium-sized organic molecules are evaluated via B3LYP calculations and subsequently corrected by the DFT-NEURON approach. Different physical descriptors are adopted to examine the validity of the DFT-NEURON approach.

The 180 small- or medium-sized organic molecules listed in Table I, are selected for the DFT-NEURON approach that combines first-principles calculation and neural network based correction. These molecules can be found in all three Refs. [9–11]. These molecules contain elements such as H, C, N, O, F, Si, S, Cl and Br. The heaviest molecule contains 14 heavy atoms, and the largest contains 32 atoms. We divide these molecules randomly into the training set containing 150 molecules and the testing set containing 30 molecules. The geometries of 180 molecules are optimized via B3LYP/6-311 + G(*d,p*) [12] calculations and zero-point energies (ZPE) are calculated subsequently at the same level. The Gibbs energy of formation of each molecule is calculated at both B3LYP/6-311 + G(*d,p*) and B3LYP/6-311 + G(3*df,2p*) levels. To calculate ΔG^0 of a molecule, we calculate first its ΔH^0 and its entropy S^0 . For instance, for a molecule A_xB_y , its ΔG^0 at 298 K can be evaluated as:

$$\begin{aligned} \Delta G^0(A_xB_y, 298\text{ K}) &= \Delta H^0(A_xB_y, 298\text{ K}) - 298.15 \\ &\times \left(S^0(A_xB_y, 298\text{ K})_{\text{st}} - \left(xS^0(A, 298\text{ K})_{\text{st}} \right. \right. \\ &\left. \left. + yS^0(B, 298\text{ K})_{\text{st}} \right) \right) \end{aligned} \quad (1)$$

where $\Delta H^0(A_xB_y, 298\text{ K})$ is the standard enthalpy of formation for A_xB_y at 298 K. $S^0(A_xB_y, 298\text{ K})$ is

the entropy of A_xB_y at 298 K, $S^0(A, 298\text{ K})_{\text{st}}$ and $S^0(B, 298\text{ K})_{\text{st}}$ are the standard entropies of species A and B relative to the elements in their reference states, respectively. $S^0(A, 298\text{ K})_{\text{st}}$ and $S^0(B, 298\text{ K})_{\text{st}}$ are taken from Ref. [13]. The strategies in Ref. [14] are adopted to calculate ΔH^0 . B3LYP/6-311 + G(3*df,2p*) employs a larger basis set than B3LYP/6-311 + G(*d,p*). The unscaled B3LYP/6-311 + G(*d,p*) ZPE is employed in the ΔG^0 calculation. The ΔG^0 values for both B3LYP/6-311 + G(*d,p*) and B3LYP/6-311 + G(3*df,2p*) are compared to their experimental data in Fig. 1a,b, respectively. The horizontal coordinates are the experimental ΔG^0 , and the vertical coordinates are the raw calculated values. The dashed line is where the vertical and horizontal coordinates are equal, i.e. where the B3LYP calculations and experiments would have the perfect match. The raw calculated values are mostly above the dashed line, i.e. most calculated ΔG^0 are larger than the experimental data. Compared to the experimental measurements, the root mean square (RMS) deviations are 22.3 and 12.9 kcal·mol⁻¹ for the raw B3LYP/6-311 + G(*d,p*) and B3LYP/6-311 + G(3*df,2p*) results, respectively. In Table II, we list the experimented and calculated ΔG^0 for 16 of 180 molecules. Overall, B3LYP/6-311 + G(3*df,2p*) calculations yield better agreements with the experiments than B3LYP/6-311 + G(*d,p*). In particular, for small molecules with few heavy elements B3LYP/6-311 + G(3*df,2p*) calculations result in very small deviations from the experiments. For instance, the deviations for CH₂F₂ and CS₂ are only 0.3 and 0.5 kcal·mol⁻¹, respectively. For large molecules, both B3LYP/6-311 + G(*d,p*) and B3LYP/6-311 + G(3*df,2p*) calculations yield quite large deviations from their experimental data.

Next we construct our neural network [15]. There is an input layer consisted of input from physical descriptors and a bias, a hidden layer, and an output layer [16]. The most important issue is to select the proper physical descriptors of our molecules, which are to be used as the input for our neural network. The calculated Gibbs energy of formation ΔG^0 contains the essential value of the experimental ΔG^0 , and is thus an obvious choice for the primary descriptor. We observe that the size of a molecule affects the accuracy of the calculated ΔG^0 . The more atoms a molecule has the worse calculated ΔG^0 is. This is consistent with the general observations in the field [14]. The total number of atoms N_t in a molecule is thus chosen as the second descriptor for the molecule. ZPE is an important component of ΔG^0 . Its calculated value is often scaled in evaluating ΔH^0 [14,17], which has the dominant contribution for ΔG^0 . It is thus taken as the third physical descriptor. Hydrogen atom is much lighter than other organic elements, and has the most quantum characteristics among all the elements. Its contributions to energy

TABLE I 180 molecules employed in the calculation

Formula	Name	Formula	Name	Formula	Name
CBrF ₃	Bromotrifluoromethane	CClF ₃	Chlorotrifluoromethane	CCIN	Cyanogens chloride
CCl ₂ O	Phosgene	CF ₂ O	Carbonyl fluoride	CF ₄	Carbon tetrafluoride
CHCl ₃	Chloroform	CH ₂ Cl ₂	Dichloromethane	CH ₂ F ₂	Difluoromethane
CH ₂ O ₂	Formic acid	CH ₃ Br	Methyl bromide	CH ₃ NO ₂	Nitromethane
CH ₃ NO ₂	Methyl nitrite	CH ₄	Methane	CH ₄ O	Methanol
CH ₄ S	Methyl mercaptan	CH ₅ N	Methylamine	COS	Carbonyl sulfide
CS ₂	Carbon disulfide	C ₂ H ₂	Acetylene	C ₂ H ₂ Cl ₂	1,1-Dichloroethylene
C ₂ H ₂ O ₄	Oxalic acid	C ₂ H ₃ Br	Vinyl bromide	C ₂ H ₃ ClO	Acetyl chloride
C ₂ H ₃ F	Vinyl fluoride	C ₂ H ₄	Ethylene	C ₂ H ₄ Br ₂	1,2-Dibromoethane
C ₂ H ₄ Cl ₂	1,1-Dichloroethane	C ₂ H ₄ O	Ethylene oxide	C ₂ H ₄ O ₂	Acetic acid
C ₂ H ₄ S	Thiacyclopropane	C ₂ H ₅ Br	Bromoethane	C ₂ H ₅ N	Ethyleneimine
C ₂ H ₅ NO	Acetamide	C ₂ H ₅ NO ₂	Nitroethane	C ₂ H ₆	Ethane
C ₂ H ₆ O	Dimethyl ether	C ₂ H ₆ S	Dimethyl sulfide	C ₂ H ₇ N	Ethylamine
C ₂ H ₈ N ₂	Ethylenediamine	C ₂ N ₂	Cyanogen	C ₃ H ₃ NO	Oxazole
C ₃ H ₄	Methylacetylene	C ₃ H ₄	Propadiene	C ₃ H ₄ O ₃	Ethylene carbonate
C ₃ H ₅ Cl ₃	1,2,3-trichloropropane	C ₃ H ₆	Cyclopropane	C ₃ H ₆ Cl ₂	1,2-Dichloropropane
C ₃ H ₆ O	Acetone	C ₃ H ₆ O ₂	Methyl acetate	C ₃ H ₆ O ₂	Propionic acid
C ₃ H ₆ S	Thiacyclobutane	C ₃ H ₇ Br	1-Bromopropane	C ₃ H ₇ Br	2-Bromopropane
C ₃ H ₇ F	1-Fluoropropane	C ₃ H ₇ NO	<i>N,N</i> -dimethylformamide	C ₃ H ₇ NO ₂	2-Nitropropane
C ₃ H ₇ NO ₃	Propyl nitrate	C ₃ H ₇ NO ₃	Isopropyl nitrate	C ₃ H ₈	Propane
C ₃ H ₈ O	Methyl ethyl ether	C ₃ H ₈ S	<i>n</i> -Propyl mercaptan	C ₃ H ₈ S	Isopropyl mercaptan
C ₃ H ₈ S	Ethyl methyl sulfide	C ₃ H ₉ N	<i>n</i> -Propylamine	C ₃ H ₉ N	Isopropylamine
C ₃ H ₉ N	Trimethylamine	C ₃ H ₁₀ N ₂	1,2-Propanediamine	C ₄ H ₄ N ₂	Succinonitrile
C ₄ H ₆	1,2-Butadiene	C ₄ H ₈	1-Butene	C ₄ H ₈ O	Isobutyraldehyde
C ₄ H ₈ O ₂	Ethyl acetate	C ₄ H ₉ Cl	tert-Butyl chloride	C ₄ H ₁₀ O	sec-Butanol
C ₄ H ₁₀ O ₂	1,4-Butanediol	C ₄ H ₁₀ S	Isobutyl mercaptan	C ₄ H ₁₀ S	Methyl propyl sulfide
C ₄ H ₁₁ N	tert-Butylamine	C ₅ H ₅ N	Pyridine	C ₅ H ₆ S	2-Methylthiophene
C ₅ H ₈ O ₂	Acetylacetone	C ₅ H ₁₀	Cyclopentane	C ₅ H ₁₀	2-Methyl-1-butene
C ₅ H ₁₀	2-Methyl-2-butene	C ₅ H ₁₀	3-Methyl-1-butene	C ₅ H ₁₀	<i>trans</i> -2-Pentene
C ₅ H ₁₀ O	2-Pentanone	C ₅ H ₁₀ O	Valeraldehyde	C ₅ H ₁₀ O ₂	Valeric acid
C ₅ H ₁₀ S	Thiacyclohexane	C ₅ H ₁₀ S	cyclopentanethiol	C ₅ H ₁₁ Br	1-Bromopentane
C ₅ H ₁₁ Cl	1-Chloropentane	C ₅ H ₁₁ N	Piperidine	C ₅ H ₁₂	Isopentane
C ₅ H ₁₂	<i>n</i> -Pentane	C ₅ H ₁₂ O	3-Methyl-1-butanol	C ₅ H ₁₂ O	3-Methyl-2-butanol
C ₅ H ₁₂ O	2-Pentanol	C ₅ H ₁₂ O	3-Pentanol	C ₅ H ₁₂ O	Ethyl propyl ether
C ₅ H ₁₂ O ₄	Pentaerythritol	C ₅ H ₁₂ S	<i>n</i> -Pentyl mercaptan	C ₆ F ₆	Hexafluorobenzene
C ₆ H ₄ C ₁₂	<i>m</i> -Dichlorobenzene	C ₆ H ₄ F ₂	<i>p</i> -Difluorobenzene	C ₆ H ₅ Cl	Monochlorobenzene
C ₆ H ₅ F	Fluorobenzene	C ₆ H ₅ NO ₂	Nitrobenzene	C ₆ H ₆	Benzene
C ₆ H ₆ N ₂ O ₂	<i>m</i> -Nitroaniline	C ₆ H ₆ O	Phenol	C ₆ H ₆ O ₂	1,3-Benzenediol
C ₆ H ₇ N	2-Methylpyridine	C ₆ H ₁₀ O ₃	Propionic anhydride	C ₆ H ₁₁ NO	<i>e</i> -Caprolactam
C ₆ H ₁₂	<i>trans</i> -3-Hexene	C ₆ H ₁₂ O	Butyl vinyl ether	C ₆ H ₁₂ O	3-Hexanone
C ₆ H ₁₄	3-Methylpentane	C ₆ H ₁₄ S	Methyl pentyl sulfide	C ₇ H ₅ N	Benzonitrile
C ₇ H ₆ O	Benzaldehyde	C ₇ H ₈	Toluene	C ₇ H ₈ O	<i>o</i> -Cresol
C ₇ H ₉ N	2,6-Dimethylpyridine	C ₇ H ₁₄	<i>cis</i> -1,2-Dimethylcyclopentane	C ₇ H ₁₅ Br	1-Bromoheptane
C ₇ H ₁₆	3,3-Dimethylpentane	C ₇ H ₁₆	2,2,3-Trimethylbutane	C ₇ H ₁₆ S	<i>n</i> -Heptyl mercaptan
C ₈ H ₆ O ₄	Terephthalic acid	C ₈ H ₈ O	Acetophenone	C ₈ H ₁₆	<i>cis</i> -1,2-Dimethylcyclohexane
C ₈ H ₁₆	<i>trans</i> -1,4-Dimethylcyclohexane	C ₈ H ₁₈	2,3-Dimethylhexane	C ₈ H ₁₈	3-Ethylhexane
C ₈ H ₁₈	4-Methylheptane	C ₈ H ₁₈ S ₂	Dibutyl disulfide	C ₉ H ₁₂	<i>m</i> -Ethyltoluene
C ₉ H ₁₂	1,2,3-Trimethylbenzene	C ₉ H ₁₈ O	Diisobutyl ketone	C ₉ H ₂₀	3,3-Diethylpentane
C ₉ H ₂₀	2,2,3,4-Tetramethylpentane	C ₁₀ H ₁₄	sec-Butylbenzene	C ₁₀ H ₁₄	Isobutylbenzene
C ₁₀ H ₁₈ O ₄	Sebacic acid	C ₁₀ H ₂₀ O ₂	<i>n</i> -Decanoic acid	C ₁₂ H ₁₀	Acenaphthene
CBrC ₁₃	Bromotrichloromethane	CHF ₃	Trifluoromethane	C ₂ H ₂ F ₂	1,1-Difluoroethylene
C ₂ H ₃ ClO ₂	Chloroacetic acid	C ₂ H ₃ Cl ₃	1,1,1-Trichloroethane	C ₂ H ₄ Cl ₂	1,2-Dichloroethane
C ₂ H ₄ F ₂	1,1-Difluoroethane	C ₂ H ₅ Cl	Ethyl chloride	C ₂ H ₅ NO ₃	Ethyl nitrate
C ₂ H ₇ N	Dimethylamine	C ₃ H ₆	Propylene	C ₃ H ₆ Br ₂	1,2-Dibromopropane
C ₃ H ₇ Cl	Isopropyl chloride	C ₃ H ₇ Cl	<i>n</i> -Propyl chloride	C ₃ H ₇ NO ₂	1-Nitropropane
C ₄ H ₆ O	Divinyl ether	C ₄ H ₉ Br	1-Bromobutane	C ₅ H ₈	<i>trans</i> -1,3-Pentadiene
C ₅ H ₁₀	1-Pentene	C ₅ H ₁₀	<i>cis</i> -2-Pentene	C ₅ H ₁₂ O	2-Methyl-1-butanol
C ₅ H ₁₂ S	Butyl methyl sulfide	C ₆ H ₈ N ₂	Adiponitrile	C ₆ H ₁₀	1-Methylcyclopentene
C ₆ H ₁₀	1,5-Hexadiene	C ₈ H ₁₀	<i>o</i> -Xylene	C ₈ H ₁₀ O	3,4-Xylenol
C ₈ H ₁₆	2,4,4-Trimethyl-2-pentene	C ₈ H ₁₈	2,3,4-Trimethylpentane	C ₈ H ₁₈ O	2-Ethyl-1-hexanol

The upper panel is for the training set, and the lower panel is for the testing set.

and entropy are different from heavy elements, and need to be distinguished from other types of atoms. We thus set the number of hydrogen atoms of a molecule N_h as the fourth descriptor.

To ensure the quality of our Neural Network, a cross-validation procedure is adopted to determine

our neural network [7,18–20]. We divide further randomly 150 training molecules into five subsets of equal size. Four of them are used to train the neural network, the fifth to validate its predictions. This procedure is repeated 5 times in rotation. The number of neurons in the hidden layer is

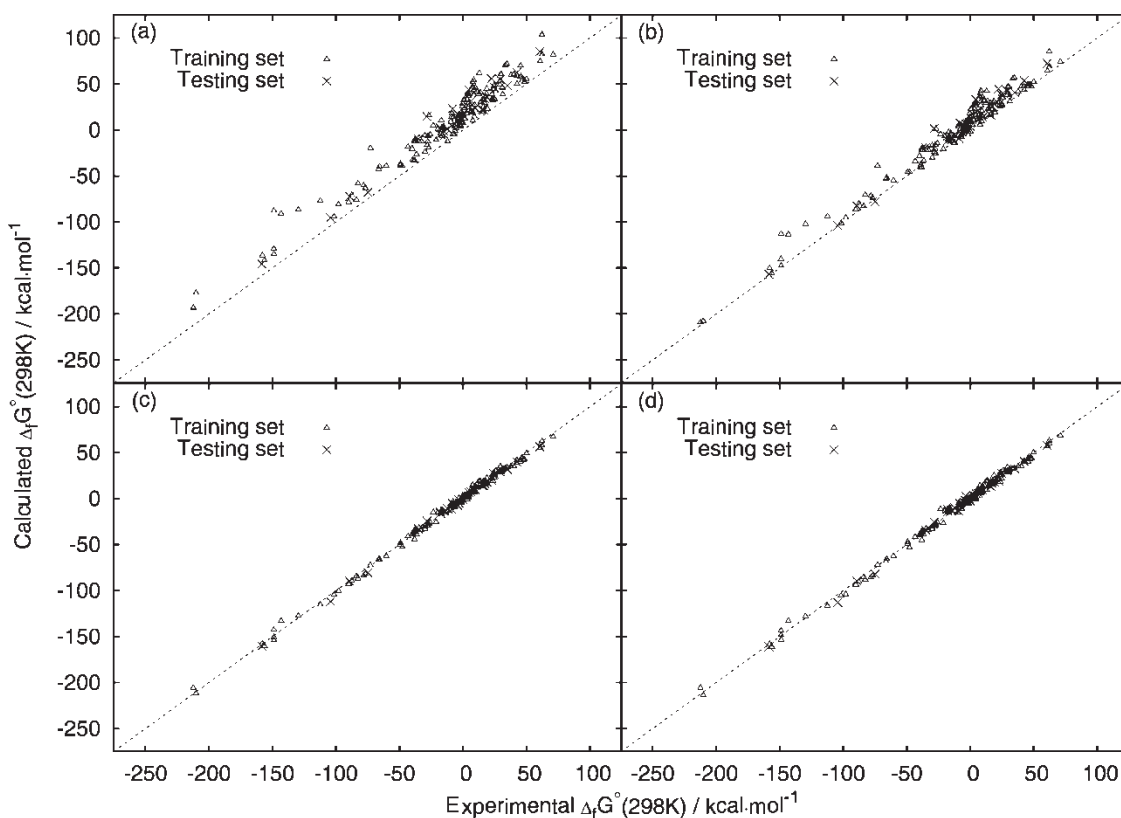


FIGURE 1 Calculated ΔG^0 versus experimental ΔG^0 for all 180 compounds. (a) and (b) are for raw B3LYP/6-311 + G(d,p) and B3LYP/6-311 + G(3df,2p) results, respectively. (c) and (d) are for neural-network corrected B3LYP/6-311 + G(d,p) and B3LYP/6-311 + G(3df,2p) ΔG^0 's, respectively. Triangles are for the training set and crosses for the testing set. The dashed lines are where the calculations and experiments have perfect fits. The correlation coefficient of the linear fits are 0.996, 0.995 in (c) and (d), respectively. All values are in the units of kcal · mol⁻¹.

varied from 2 to 10 and the number of descriptors in the input layer is switched from 2 to 4 in order to determine the optimal structure of our neural network. We find that an input layer containing

four descriptors (ΔG^0 , N_v , ZPE and Nh) and a hidden layer containing two neurons yield the best overall results. The 5-2-1 structure is thus adopted for our neural network (See Fig. 2). It is found that

TABLE II Experimental and calculated Gibbs energies of formation for sixteen selected compounds (all data are in the units of kcal · mol⁻¹)

Formula	Name	Experiment ΔG^0 (298.15K)*	Deviation (Theory – Experiment)			
			DFT1 [†]	DFT1-NN [‡]	DFT2 [§]	DFT2-NN [§]
CBrF ₃	Bromotrifluoromethane	-148.82 ± 0.7	14.0	-5.2	1.0	-4.9
CClF ₃	Chlorotrifluoromethane	-156.30 ± 0.8	14.4	-4.2	0.3	-4.8
CF ₂ O	Carbonyl fluoride	-148.99 ± 0.4	19.8	6.4	8.5	5.4
CF ₄	Carbon tetrafluoride	-212.34 ± 0.3	19.0	6.5	3.3	6.6
CHCl ₃	Chloroform	-16.38 ± 0.3	16.1	4.4	8.0	6.4
CHF ₃	Trifluoromethane	-158.48 ± 0.8	12.8	-1.8	1.5	-2.7
CH ₂ Cl ₂	Dichloromethane	-16.46 ± 0.3	10.2	3.0	4.6	4.3
CH ₂ F ₂	Difluoromethane	-101.66 ± 0.4	7.8	-2.8	0.3	-4.0
CH ₄	Methane	-12.15 ± 0.1	-0.2	-0.7	-1.9	-1.4
CS ₂	Carbon disulfide	15.99 ± 0.2	8.6	2.2	0.5	1.5
C ₅ H ₁₂	<i>n</i> -Pentane	-2.00 ± 0.4	17.4	-1.6	10.6	-1.9
C ₅ H ₁₂ O	3-Methyl-2-butanol	-37.37 ± 0.7	25.1	1.3	15.8	0.6
C ₅ H ₁₂ O	2-Pentanol	-38.07 ± 0.6	25.7	2.3	16.5	1.4
C ₆ H ₁₄	3-Methylpentane	-0.51 ± 0.4	26.2	1.5	18.1	1.1
C ₈ H ₁₀	<i>o</i> -Xylene	29.18 ± 0.6	25.8	-0.2	13.4	0.6
C ₉ H ₁₂	1,2,3-Trimethylbenzene	29.77 ± 0.6	31.1	-0.9	17.3	-0.7

*The experimental values are taken from Ref. [9] and the uncertainties are estimated from Refs. [13,21,22]. [†]The deviations of calculated ΔG^0 by using B3LYP/6-311 + G(d,p) geometries, zero point energies. [‡]The deviations of calculated ΔG^0 by B3LYP/6-311 + G(d,p)—Neural Networks approach. [§]The deviations of calculated ΔG^0 by using the 6-311 + G(d,p) geometries, zero point energies and recalculated total energies by 6-311 + G(3df,2p) basis. [§]The deviations of calculated ΔG^0 by B3LYP/6-311 + G(3df,2p)—Neural Networks approach.

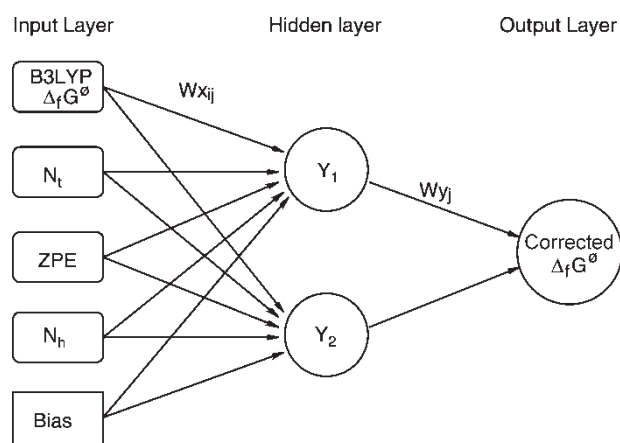


FIGURE 2 The structure of our neural network. The input for the bias is always 1.

our neural network can be tuned to fit the trained data and the corresponding experimental data very well for the training set. Moreover, for the testing set our neural network predicts the values that agree very well with experiments. The quantitative relationship between the calculated and the experimental data is thus established. In Fig. 1c,d, the triangles and crosses belong to the training and testing set, respectively. Compared to the raw calculated results, the neural-network corrected values are much closer to the experimental values for both training and testing sets. More importantly, the systematic deviations in Fig. 1a,b are eliminated, and the remaining numerical errors are reduced

substantially. This can be further demonstrated by the error analysis performed for B3LYP/6-311 + G(*d,p*) ΔG^0 of all 180 molecules. In Fig. 3, we plot the histograms for the deviations (from the experiments) of the raw B3LYP/6-311 + G(*d,p*) ΔG^0 s and their neural-network corrected values. Figure 3a,c shows the raw ΔG^0 s and neural-network corrected ΔG^0 s of training set, respectively. Figure 3b,d shows the raw ΔG^0 s and neural-network corrected ΔG^0 s of testing set, respectively. For the training set, the RMS errors before and after the neural-network correction are 22.6 and 3.1 kcal·mol⁻¹, respectively; while for the testing set they are 20.7 and 2.8 kcal·mol⁻¹, respectively. Overall, the raw calculated ΔG^0 s have a large systematic deviations of 22.3 kcal·mol⁻¹ while the neural-network corrected ΔG^0 s have virtually no systematic deviation at all. Moreover, the remaining numerical deviations are much smaller for the neural-network corrected ΔG^0 s, i.e. the RMS deviations is reduced from 22.3 to 3.0 kcal·mol⁻¹ upon neural-network correction. Note that the error distribution after the neural-network correction is of an approximate Gaussian type (see Fig. 3c,d). We have performed the same error analysis for B3LYP/6-311 + G(3*df*,2*p*) ΔG^0 s and have reached a similar conclusion. The RMS deviation for B3LYP/6-311 + G(3*df*,2*p*) ΔG^0 s is reduced from 13.2 to 3.4 kcal·mol⁻¹ upon neural-network correction. In Table II, we list the neural-network corrected ΔG^0 s. The deviations of large molecules are of the same magnitude as those of small molecules. Our neural-network correction

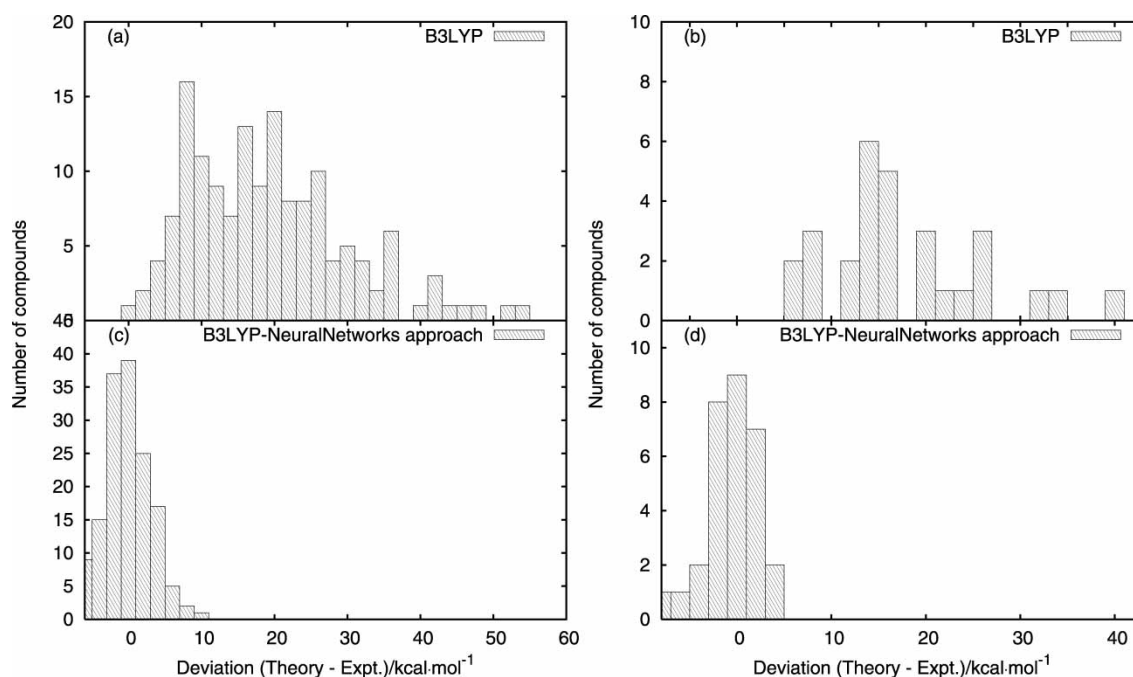


FIGURE 3 Histograms for the deviations of B3LYP/6-311 + G(*d,p*) calculated ΔG^0 for all 180 compounds. (a) and (b) are for the raw calculated ΔG^0 of training and testing set, respectively, and (c) and (d) are for the neural network corrected ΔG^0 of training and testing set, respectively. ΔG^0 is in the units of kcal·mol⁻¹.

TABLE III The DFT-neural network RMS with different descriptors (all data are in the units of kcal · mol⁻¹)

Method	B3LYP/6-311+G(d,p)-NN			B3LYP/6-311+G(3df,2p)-NN		
	Training set	Testing set	Overall	Training set	Testing set	Overall
A	4.3	5.0	4.4	3.7	4.3	3.8
B	3.1	2.7	3.1	3.7	4.0	3.7
C	3.1	2.8	3.1	3.5	3.6	3.5
D	3.1	2.8	3.0	3.4	3.4	3.4

Descriptors used in: Method A, ΔG^0 and N_t ; Method B, ΔG^0 , N_t and ZPE; Method C, ΔG^0 , N_t , ZPE and N_{db} ; Method D, ΔG^0 , N_t , ZPE and N_H .

strategy does not discriminate against large molecules, unlike most calculations, which yield worse results for large molecules than for small ones. Although the raw B3LYP/6-311 + G(*d,p*) results have much larger deviations than those of B3LYP/6-311 + G(3df,2p), the neural-network corrected values of both calculations have deviations of the same magnitude. This implies that it is sufficient to employ the smaller basis set B3LYP/6-311 + G(*d,p*) in our DFT-NEURON approach. The neural-network step can easily correct the deficiency of a small basis set. Therefore, our approach can potentially be applied to much larger systems.

Analysis of our neural network reveals that the weights connecting the input for ΔG^0 have the dominant contribution in all cases. This confirms our fundamental assumption that the calculated ΔG^0 captures the essential values of exact ΔG^0 . The input for the second physical descriptor, N_t , has quite large weights in all cases. In particular, when the smaller basis set 6-311 + G(*d,p*) is adopted in the B3LYP calculations, N_t has the second largest weights. It is found that the raw ΔG^0 deviations are roughly proportional to N_t , which confirms the importance of N_t as a significant descriptor of our neural network. We use the raw calculated ΔG^0 and N_t as the only two descriptors for the input layer, and find that the RMS deviations of the training and testing sets are corrected to 4.3, 5.0 kcal · mol⁻¹ for B3LYP/6-311 + G(*d,p*), and 3.7, 4.3 kcal · mol⁻¹ for B3LYP/6-311 + G(3df,2p) (See Table III), respectively. We thus conclude that raw ΔG^0 and N_t are the two main physical properties that determine the experimental values of ΔG^0 . This is the main reason why our neural network approach is successful here. The bias contributes to the correction of systematic deviations in the raw calculated data, and has thus significant weights. When the larger basis set 6-311 + G(3df,2p) is used, the bias has the second largest weights for all cases. ZPE has been often scaled to account for the discrepancies of ΔH^0 between calculations and experiments [14,17], and it is thus expected to have large weights. This is indeed the case, especially when the smaller basis set 6-311 + G(*d,p*) is adopted in calculations. In all cases the number of hydrogen atoms, N_H , has the smallest but non-negligible weights. The number of double bonds, N_{db} , is

the measurement of chemical structure of a molecule, and may be used as another physical descriptor. When N_H are replaced by N_{db} , the similar results are obtained as shown in Table III. This similarity can be rationalized. Both descriptors, i.e. the number of hydrogen N_H and the number of double bond N_{db} , are reflecting the structural information of individual molecules in different aspects.

The neural-network based approach adopted in this work reduces substantially the raw first-principles calculation errors, and can be used to quantify the uncertainty of the resulting ΔG^0 . Our DFT-NEURON approach has a RMS error of ~ 3 kcal · mol⁻¹ for the 180 small- to medium-sized organic molecules. This is larger than the experimental error bars of our 180 molecules [9–11,13]. G2 method [14] results more accurate energies for small molecules. However, our approach is more efficient and can be applied to larger systems. The physical descriptors adopted in our neural network, the raw ΔG^0 , N_t , ZPE and N_H , are quite general, and are not limited to the special motifs of our molecules. Our DFT-NEURON approach is therefore expected to yield the RMS deviations of ~ 3 kcal · mol⁻¹ for any small to medium sized organic molecules. To improve the accuracy for our DFT-NEURON approach, we may need more and better physical descriptors for the molecules, and possibly, more and better experimental data. The larger the experimental database is, the more accurately the DFT-NEURON approach predicts. In principle, the accuracy of our approach is limited only by the precision and size of the experimental database.

Besides the Gibbs energy of formation, our DFT-NEURON approach can be generalized to calculate other properties such as heat of formation, ionization energy, dissociation energy, absorption frequency and etc. The success of our approach is determined by the selection of the physical descriptors. The raw first-principles calculation property of interest contains its essential value, and is thus always the primary descriptor. The raw calculation error accumulates as the molecular size increases. The number of atom N_t should thus be selected for any Neural Networks correction procedure.

Additional physical descriptors should be chosen according to their relations to the property of interest and the physical and chemical structures of the compounds. Others have used Neural Networks to determine the quantitative structure relationship between the experimental measured thermodynamic properties and the structure parameters of the molecules [18]. We distinguish our work from others by utilizing specifically the first-principles methods. Since the first-principles calculations captures more readily the essences of the property of interest, our approach is more reliable and covers a much wider range of molecules or compounds.

To summarize, we have applied the DFT-NEURON approach to calculate the first-principles Gibbs energy of formation for small- and medium-sized organic molecules. The accuracy of the combined first-principle calculations and Neural Networks correction approach can be systematically increased, as more and better experimental data are available. Since the requirements on first-principles methods are modest, our approach is very efficient compared to the sophisticated first-principles methods of similar accuracy. More importantly, it can be easily extended to very large molecular systems.

Acknowledgements

Support from the Hong Kong Research Grant Council (RGC) and the Committee for Research and Conference Grants (CRCG) of the University of Hong Kong is gratefully acknowledged.

References

- [1] Cramer, C.J. (2002) *Essentials of Computational Chemistry: Theories and Models* (John Wiley, West Sussex, England).
- [2] Parr, R.G. and Yang, W. (1989) *Density-Functional Theory of Atoms and Molecules* (Oxford University Press, New York).
- [3] Foresman, J.B. and Frisch, A. (1996) *Exploring Chemistry with Electronic Structure Method*, 2nd Ed. (Gaussian Inc., Pittsburgh, PA).
- [4] Irikura, K.K. and Frurip, D.J. (1998) *Computational Thermochemistry: Prediction and Estimation of Molecular Thermodynamics* (American Chemical Society, Washington, D.C.).
- [5] Harvery, R.L. (1994) *Neural Network Principles* (Prentice-Hall, Englewood Cliffs, NJ).
- [6] Ripley, B.D. (1996) *Pattern Recognition and Neural Networks* (Cambridge, New York).
- [7] Haykin, S. (1999) *Neural Networks: a Comprehensive Foundation* (Prentice Hall, Upper Saddle River).
- [8] Rumelhart, D.E., Hinton, G.E. and Williams, R.J. (1986) "Learning representations by back-propagating errors", *Nature* **323**, 533.
- [9] Yaws, C.L. (1999) *Chemical Properties Handbook* (McGraw-Hill, New York).
- [10] Lide, D.R. (2000) *CRC Handbook of Chemistry and Physics*, 3rd Electronic ed. (CRC Press, Boca Raton, FL).
- [11] Dean, J.A. (1999) *Lange's Handbook of Chemistry*, 15th Ed. (McGraw-Hill, New York).
- [12] Frisch, M.J., Trucks, G.W., Schlegel, H.B., Scuseria, G.E., Robb, M.A., Cheeseman, J.R., Zakrzewski, V.G., Montgomery, J.A., Stratmann, R.E., Burant, J.C., Dapprich, S., Millam, J.M., Daniels, A.D., Kudin, K.N., Strain, M.C., Farkas, O., Tomasi, J., Barone, V., Cossi, M., Cammi, R., Mennucci, B., Pomelli, C., Adamo, C., Clifford, S., Ochterski, J., Petersson, G.A., Ayala, P.Y., Cui, Q., Morokuma, K., Salvador, P., Dannenberg, J.J., Malick, D.K., Rabuck, A.D., Raghavachari, K., Foresman, J.B., Cioslowski, J., Ortiz, J.V., Baboul, A.G., Stefanov, B.B., Liu, G., Liashenko, A., Piskorz, P., Komaromi, I., Gomperts, R., 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., Andres, J.L., Gonzalez, C., Head-Gordon, M., Replogle, E.S. and Pople, J.A. (2002) *Gaussian 98*, Revision A.11.3. Gaussian, Inc., Pittsburgh, PA.
- [13] Chase, M.W., Davies, C.A., Downey, J.R., Frurip, D.J., McDonald, R.A. and Syverud, A.N. (1985), *J. Phys. Chem. Ref. Data* **14**(Suppl. No. 1).
- [14] Curtiss, L.A., Raghavachari, K., Redfern, P.C. and Pople, J.A. (1997) "Assessment of gaussian-2 and density functional theories for the computation of enthalpies of formation", *J. Chem. Phys.* **106**, 1063.
- [15] Ivanciuc, O., Rabine, J.-P., Cabrol-Bass, D., Panaye, A. and Doucet, J.P. (1997) "¹³C NMR chemical shift prediction of the sp³ carbon atoms in the α position relative to the double bond in acyclic alkenes", *J. Chem. Inf. Comput. Sci.* **37**, 587.
- [16] Nohair, M. and Zakarya, D. (2002) "Autocorrelation method adapted to generate new atomic environments: application for the prediction of 13-C chemical shifts of alkanes", *J. Chem. Inf. Comput. Sci.* **42**, 586.
- [17] Curtiss, L.A., Raghavachari, K., Redfern, P.C. and Pople, J.A. (1997) "Investigation of the use of B3LYP zero-point energies and geometries in the calculation of enthalpies of formation", *Chem. Phys. Lett.* **270**, 419.
- [18] Yao, X., Zhang, X., Zhang, R., Liu, M., Hu, Z. and Fan, B. (2001) "Prediction of enthalpy of alkanes by the use of radial basis function neural networks", *Comput. Chem.* **25**, 475.
- [19] Pompe, M. and Novič, M. (1999) "Prediction of gas-chromatographic retention indices using topological descriptors", *J. Chem. Inf. Comput. Sci.* **39**, 59.
- [20] Dong, L., Yan, A., Chen, X., Xu, H. and Hu, Z. (2001) "Research and prediction of coordination reactions between CPA-mA and some metal ions using artificial neural networks", *Comput. Chem.* **25**, 551.
- [21] Afeefy, H. Y., Liebman, J. F. and Stein, S.E. (2003) "Neutral Thermochemical Data", In: Linstrom, P.J. and Mallard, W.G., eds, *NIST Chemistry WebBook NIST Standard Reference Database Number 69*, (National Institute of Standards and Technology, Gaithersburg MD, 20899), (<http://webbook.nist.gov>).
- [22] Taylor, J.R. (1982) *An Introduction to Error analysis: The Study of Uncertainties in Physical Measurements* (University Science Books, Mill Valley, CA).