

Improved prediction of solvation free energies by machine-learning polarizable continuum solvation model

Amin Alibakhshi ¹✉ & Bernd Hartke ¹

Theoretical estimation of solvation free energy by continuum solvation models, as a standard approach in computational chemistry, is extensively applied by a broad range of scientific disciplines. Nevertheless, the current widely accepted solvation models are either inaccurate in reproducing experimentally determined solvation free energies or require a number of macroscopic observables which are not always readily available. In the present study, we develop and introduce the Machine-Learning Polarizable Continuum solvation Model (ML-PCM) for a substantial improvement of the predictability of solvation free energy. The performance and reliability of the developed models are validated through a rigorous and demanding validation procedure. The ML-PCM models developed in the present study improve the accuracy of widely accepted continuum solvation models by almost one order of magnitude with almost no additional computational costs. A freely available software is developed and provided for a straightforward implementation of the new approach.

¹Theoretical Chemistry, Institute for Physical Chemistry, Christian-Albrechts-University, Olshausenstr. 40, Kiel, Germany. ✉email: alibakhshi@pctc.uni-kiel.de

Free energy of solvation is one of the key thermophysical properties in studying thermochemistry in solution, where the majority of real-life chemistry happens. In theoretical studies of solution chemistry, estimation of free energies allows evaluation of reaction rates and equilibrium constants of physical or chemical reactions of interest. Nevertheless, direct evaluation of free energies in solution can be quite challenging since it sometimes requires appropriate sampling of phase space^{1–3} and appropriate treatment of the non-covalent interactions between the solvent and solute, which can have a remarkable impact on electronic structures of both the solvent and solute and consequently on the microscopic and macroscopic observables^{4,5}.

Theoretical approaches for evaluating physical chemistry behind solvation free energy can be generally divided into two main categories, namely explicit solvent and implicit solvent approaches. In explicit solvent approaches, solvent molecules are treated explicitly, and the free energy is typically evaluated by analyzing the trajectory of time evolution of phase space obtained via molecular dynamics or Monte Carlo simulations. For that end, a number of efficient free energy estimators have been developed in the past decades such as thermodynamic integration, free-energy perturbation, and histogram analysis methods⁶.

Despite obvious advantages of applying the explicit solvent methods such as retaining the physically proper picture of discrete solvent molecules, they suffer by a number of limitations when applied to free-energy estimation. For example, in case of applying methods which evaluate the free energy through alchemical transformations (e.g., thermodynamic integration or free energy perturbation), defining intermediate states and pathways between the endpoints appropriately can be quite tricky⁷. Also, necessity of employing appropriate force fields, which for many solute-solvent mixtures requires to develop or reparametrize a force field, and running the simulations and trajectory analyses can be laborious and time-taking tasks.

To overcome the mentioned limitations, the implicit solvent approach has been developed and is widely applied as standard method for studying solvent effects in computational chemistry. In implicit solvent approaches, the solvent molecules are treated implicitly as a continuous medium and the solute is placed in a cavity of this implicitly defined solvent. The solute-solvent interactions are then evaluated via considering the solvent polarization due to the solute charge distribution and its resulting potential field acting on the solute, known as the reaction field⁵. For a moderate level of theory and medium-sized molecules, implicit solvent approaches can yield a reasonable estimation of the solvation free energy in few seconds to few minutes on a normal desktop PC, while for explicit solvent approaches it might take from hours to days.

The most widely applied implicit solvent approaches are those based on the so-called polarizable continuum model (PCM) proposed by Tomasi and co-workers⁸. In polarizable continuum models, the solvation free energy is constructed by summing the contributions of electrostatic interactions including electronic, nuclear, and polarization interactions (ΔG_{ENP}), changes in free energy by solvent cavity formation, dispersion energy and local solvent structure changes (G_{CDS}), and corrections for differences in molar densities in the two phases compared with the standard state (ΔG_{cons}°). The contributions of electrostatic interactions are evaluated by iteratively solving the following relationship:

$$\Delta G_{ENP} = \langle \Psi^{(1)} | H + \frac{1}{2} V | \Psi^{(1)} \rangle - \langle \Psi^{(0)} | H | \Psi^{(0)} \rangle \quad (1)$$

which is known as the self-consistent reaction-field (SCRF) calculations⁵. Here, superscripts (0) and (1) refer to the gas and solution phases, respectively, and V is the potential energy operator resulting from the reaction field. Various constructions

of the potential energy operator as well as G_{CDS} have resulted in different continuum solvation models. The parallel existence of several continuum solvation models is a good indicator that each of them has its own strengths and weaknesses, and choosing a single, optimal model is not trivial. It is totally impossible to provide a detailed overview here; a 2005 review of implicit solvation models⁹ covered 95 pages and cited 936 references. In the present study, we only consider the most widely used PCM-based models.

One of simplest and yet successful continuum solvation models is CPCM which implements the conductor-like screening solvation boundary condition within the PCM framework. In CPCM, the following correction of the polarization charge densities by the scaling factor x is employed¹⁰:

$$f(\epsilon) = \frac{\epsilon - 1}{\epsilon + x} \quad (2)$$

where ϵ is the solvent dielectric constant. One main advantage of CPCM is its much simpler defined boundary conditions. More importantly, unlike more advanced PCM-based models which require the normal component of the solute electric field as input, CPCM only requires the solute electrostatic potential; for this reason it is much less affected by outlying charge errors (OCE)^{11,12}. A more versatile model exploiting the conductor-like screening solvation boundary condition is COSMO-RS, developed by Klamt and co-workers^{13,14}, which although initially proposed in 1995, still is one of the most accurate available continuum solvation models. A more sophisticated treatment of the boundary condition is implemented in the integral equation formalism of PCM (IEF-PCM) taking into account apparent surface charge isotropic¹⁵ or anisotropic¹⁶ dielectric continuum solvation. Another extensively used continuum solvation model is the SMx family of methods which specifically focuses on more accurate estimation of the solvation free energy^{4,5}.

We already discussed the main advantages of continuum solvation models such as their efficiency in terms of computational cost. Nevertheless, it should be noted that all this has become possible for a considerable amount of assumptions and simplifications on the physics of the problem, such as overlooking the conformational entropy of solvent and solute which can have a significant contribution on the total free energy¹⁷, neglecting the site-specific solute-solvent interactions and decoupling the polar and nonpolar components of free energies and considering them independent, linear and additive^{18,19}. The inaccuracies resulting from such simplifications are commonly compensated for via incorporating additional macroscopic observables as well as adjustable parameters in the solvation models. In the CPCM model for example, this is achieved by implementing an ad hoc modification of the atomic radii via defining a number of adjustable parameters and empirical descriptors, such as the number of bonded hydrogens and the number of bonded active atoms¹⁰. In the COSMO-RS model, it is achieved by ad hoc modification of the interaction energies and effective contact area via some adjustable parameters¹⁴.

In contrast, in the SMx family of methods, to provide a more accurate estimation of the solvation free energy, an ad hoc modification of the G_{CDS} term in (1) has been proposed. For that end, employing additional macroscopic observables in the model has been considered⁴, including the refractive index, Abraham's hydrogen bond acidity and basicity of the solute, macroscopic surface tension of the solvent at the air/solvent interface at 298.15 K, the square of the fraction of solvent atoms that are aromatic carbon atoms, and the square of the fraction of solvent atoms that are F, Cl, or Br. Although these employed macroscopic observables indirectly introduce more physics into the model and hence provide the chance to make predictions of solvation free

energies more universal, except for the last two they are not readily available for many new compounds and their experimental or theoretical evaluation is not straightforward.

In a number of recent studies, Machine Learning (ML) has been exploited to map the highly complicated relationship between solvation free energy and potentially relevant macroscopic or microscopic observables.

Wang et al. employed a pool of 30 molecular representations which all are either per atom reaction field energies or partial charges, as the input of the learning-to-rank (LTR) machine learning algorithm, resulting in a root mean squared error (RMSE) of 1.05 kcal/mol¹⁸. Borhani et al. developed a QSPR model which requires 12 experimentally determined properties of solvent and 9 QM derived representations of solute as model input, yielding a Mean Unsigned Error (MUE) of 0.43 kcal/mol²⁰. Hutchinson and Kobayashi proposed a structure property relationship for prediction of hydration free energy which yields a RMSE of 1.65 kcal/mol²¹.

Another recent example is the kernel-based machine learning model of Rauer and Bereau which is developed to predict the free energy of solvating small organic molecules containing C, H, O, and N atoms in pure water via implicit-solvent molecular dynamics simulations²². For a 39-parameter model they reported a MUE of 1.06 kcal/mol.

The most recent example of employing machine learning for prediction of solvation free energy is the model developed by Vermeire and Green²³. Their model is developed based on the transfer of knowledge learned through one million data of QM evaluated free energies and fine tuning it to accurately reproduce the experimentally determined solvation free energies. They reported a MUE of 0.21 kcal/mol for their model which is currently the most accurate ever reported result for prediction of solvation free energy.

In the present study, we propose a machine-learning-based PCM model, which, similar to other conventional continuum solvation models, is based on considering the solvent as a continuous medium and calculating the solvation energy components of a solute placed in the cavity of this medium by the SCRF procedure. Nevertheless, unlike the conventional PCM models which propose simple and ad hoc expressions to integrate and modify those calculated energy components, we employ machine learning for this purpose and show its efficiency in substantial improvements of the predictability of solvation free energy.

Results and discussions

After setting up and training the neural networks and screening the appropriately trained models via the post-validation strategy discussed in the previous section, the best results with MUE of 0.52526 and 0.40011 kcal/mol were observed for the computations at B3LYP/6–31 G* and DSD-PBEP86-D3/def2TZVP levels of theory, respectively. The two models employed SCRF energy components and solvation free energy computed via CPCM_{x=0.5} solvation model in both cases and 100 and 130 hidden layer neurons, respectively. These two models are denoted by ML-PCM(B3LYP) and ML-PCM(DSD-PBEP86) hereafter, respectively. Details of the selected input variables and implementation instructions for all selected models are provided in Supplementary Software 1. These results show a substantial improvement compared to the original continuum solvation model CPCM_{x=0.5}, which for the same dataset yielded MUE of 3.1611 and 2.9130 kcal/mol, respectively.

In comparison to the SMD model, which for the same dataset and solvation free energy computations at B3LYP/6–31 G* and DSD-PBEP86-D3/def2TZVP levels yields MUE of 0.78623 and 0.85396 kcal/mol, respectively, the obtained results still show a higher accuracy, without requiring additional solvent parameters

needed in the SMD approach. In comparison to the MUE of 0.4214 kcal/mol reported by Klamt and Diedenhofen²⁴ for employing one of the recent versions of the COSMO-RS model for the same dataset, the ML-PCM(DSD-PBEP86) provides a slightly higher accuracy. Also, in terms of maximum unsigned error, the two ML-PCM models which yield maximum unsigned error of 6.2252 and 3.8799 kcal/mol, respectively, are more accurate than that of COSMO-RS for which this value is 6.8701 kcal/mol. For other continuum solvation models studied for the same dataset, the maximum unsigned error of the SMD, PCM, CPCM and CPCM_{x=0.5} were 11.311, 12.75, 12.2, 12.6 kcal/mol for B3LYP/6–31 G* and 11.311, 12.83, 12.31, 12.68 kcal/mol for DSD-PBEP86-D3/def2TZVP levels of theory, which are all substantially higher than those achievable by the ML based models.

The higher accuracy of the predicted solvation free energies by the COSMO-RS model compared to the other conventional solvation models also motivated us to study neural networks which take SCRF energy components computed via PCM or CPCM models in addition to the solvation free energies predicted via COSMO-RS as neural network feeds. For these updates, the best results with MUEs of 0.26057 and 0.24387 kcal/mol and maximum unsigned errors of 7.1349 and 2.9154 kcal/mol were obtained for energy components calculated via CPCM_{x=0.5} and CPCM solvation models, 130 and 120 hidden layer neurons, and computations at B3LYP/6–31 G* and DSD-PBEP86-D3/def2TZVP levels of theory, respectively. These two models, which are denoted by ML-PCM/COSMO-RS(B3LYP) and ML-PCM/COSMO-RS(DSD-PBEP86) hereafter, respectively, show a remarkable improvement in predicted solvation free energy compared to those obtained via the original implementation of COSMO-RS reported by Klamt and Diedenhofen²⁴. This implies considerable flexibility of the proposed approach in improving accuracy of various solvation models. Nevertheless, it should be noted that the solvation free energies evaluated by COSMO-RS which were used as additional model inputs in the present study were evaluated using the 2015 version of that method. Using free energies evaluated by more recent versions of COSMO-RS and also the energy terms computed with this method, will probably result in more accurate predictions of the solvation free energy by the presented ML-PCM.

As the most important parameter in developing ANN models, we studied the impact of the selected number of hidden layer neurons on the performance of the developed machine learning models. As can be seen in Fig. 1, by increasing the number of hidden layer neurons, the predictability of the solvation free

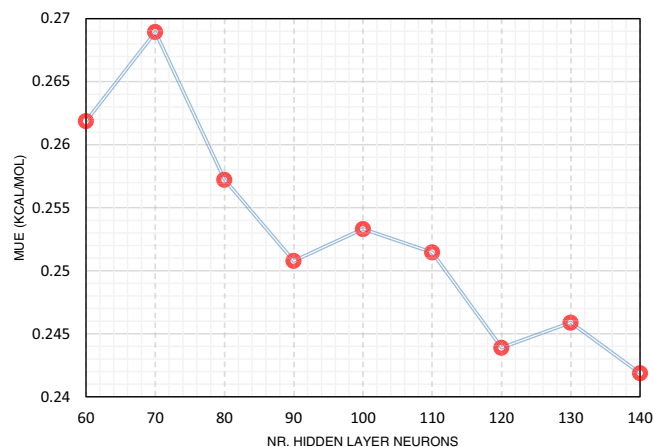


Fig. 1 MUE of developed ML-PCM/COSMO-RS(B3LYP) models versus the number of hidden layer neurons. The general trend shows the reducing pattern in MUE with increasing the size of the neural network model.

energy is generally improved. This is due to the larger number of adjustable parameters of the resulting models and their consequently higher flexibility to map complicated functionalities. However, at the same time this may reduce the extrapolation capability of the model, i.e., it may reduce performance when applied to samples remarkably different from those already examined in developing the models.

To investigate the impact of the number of hidden layer neurons on extrapolation performance of the models developed in the present study, we re-examined the trained models for out-of-sample predictions, following the approach proposed by Vermeire and Green²³. For that end, we compared the results of models for which a group of samples with either a specific element or a specific solvent were included in the training dataset with the same models trained with a dataset excluding that specific group of samples. We studied out-of-sample prediction performance for 20 solvents and 6 solute elements most frequently encountered in our studied dataset. The obtained results are reported in Tables 1 and 2. According to the results, the developed models show an excellent extrapolation capability for out of sample predictions of solvent splits, while for the element splits, the extrapolation is slightly less accurate. Furthermore, except for the element Br, the out-of-sample predictions tested for ML-PCM/COSMO-RS(B3LYP) are within chemical accuracy.

A comparison of predicted and experimentally determined free energies is depicted in Fig. 2. As can be seen, the linear correlation

between the predicted and reference data is more evident for the newly derived models, compared to the conventionally accepted ones.

The overall results obtained via newly developed ML models are compared with various other models proposed in the literature in Table 3. Although a more informative comparison would be possible if different models were compared for the same dataset and, if applicable, the same level of theory, the larger size of the benchmark dataset used in the present study compared to most of the other works confirms the superior accuracy of the newly proposed method compared to the majority of the widely accepted ones. In comparison to the model developed by Vermeire and Green²³ which yields MUE of 0.21 kcal/mol, our results are slightly less accurate, but it should be noted that our results are obtained for a much lower number of neurons and model parameters.

Furthermore, it should be noted that the inaccuracies inherent in the reference data of solvation free energies (Aleatoric uncertainty) can also impact both the training efficiency and inferences about model performances, as pointed out by Vermeire and Green²³.

To summarize, we have demonstrated substantial improvements of continuum solvation models in evaluating solvation free energy with the help of machine learning. For that end, we proposed a more versatile machine learning assisted integration of the continuum solvation energy components calculated in SCRFF computations which can be used to modify the predicted solvation free energy by various solvation models. It allowed us to achieve

Table 1 Out-of-sample predictions for solvent splits.

Solvent	Nr. Samples	ML-PCM/COSMO-RS(B3LYP)		ML-PCM/COSMO-RS(DSD-PBEP86)	
		MUE (solvent included)	MUE (solvent excluded)	MUE (solvent included)	MUE (solvent excluded)
Water	261	0.13921	0.53856	0.12724	0.52107
n-Octanol	199	0.21116	0.40528	0.19416	0.34079
n-Hexadecane	184	0.47931	0.63312	0.42914	0.38652
Chloroform	102	0.2962	0.33126	0.27975	0.28894
CycloHexane	88	0.27941	0.30729	0.30521	0.35877
CarbonTetraChloride	73	0.37704	0.38958	0.30407	0.31146
Benzene	71	0.21953	0.24581	0.37323	0.52627
DiethylEther	66	0.23975	0.29187	0.22181	0.22156
Heptane	64	0.41215	0.4795	0.2033	0.19233
n-Hexane	57	0.19548	0.19648	0.28332	0.3775
Toluene	49	0.22219	0.20023	0.31435	0.33986
Xylene-mixture	46	0.25694	0.22309	0.27209	0.27789
DiChloroEthane	37	0.38469	0.49085	0.22075	0.28748
n-Decane	37	0.21171	0.25761	0.15559	0.17148
ChloroBenzene	36	0.2183	0.2374	0.20119	0.22425
n-Octane	35	0.13265	0.13455	0.17431	0.19447
2,2,4-TriMethylPentane	32	0.2097	0.20388	0.23426	0.23618
EthylBenzene	27	0.20878	0.23166	0.20471	0.24331
BromoBenzene	24	0.16054	0.22188	0.13648	0.18478
Decalin-mixture	24	0.39408	0.44484	0.31164	0.33004

Table 2 Out-of-sample predictions for element splits.

Element	Nr. Samples	ML-PCM/COSMO-RS(B3LYP)		ML-PCM/COSMO-RS(DSD-PBEP86)	
		MUE (element included)	MUE (element excluded)	MUE (element included)	MUE (element excluded)
N	611	0.25549	0.40139	0.25486	0.37139
F	81	0.29188	0.38812	0.32345	0.48087
P	62	0.12927	0.64773	0.20648	0.95936
S	91	0.26592	0.50868	0.29104	0.53079
Cl	174	0.25295	0.5194	0.17956	0.47383
Br	102	0.25005	1.4559	0.26268	0.91972

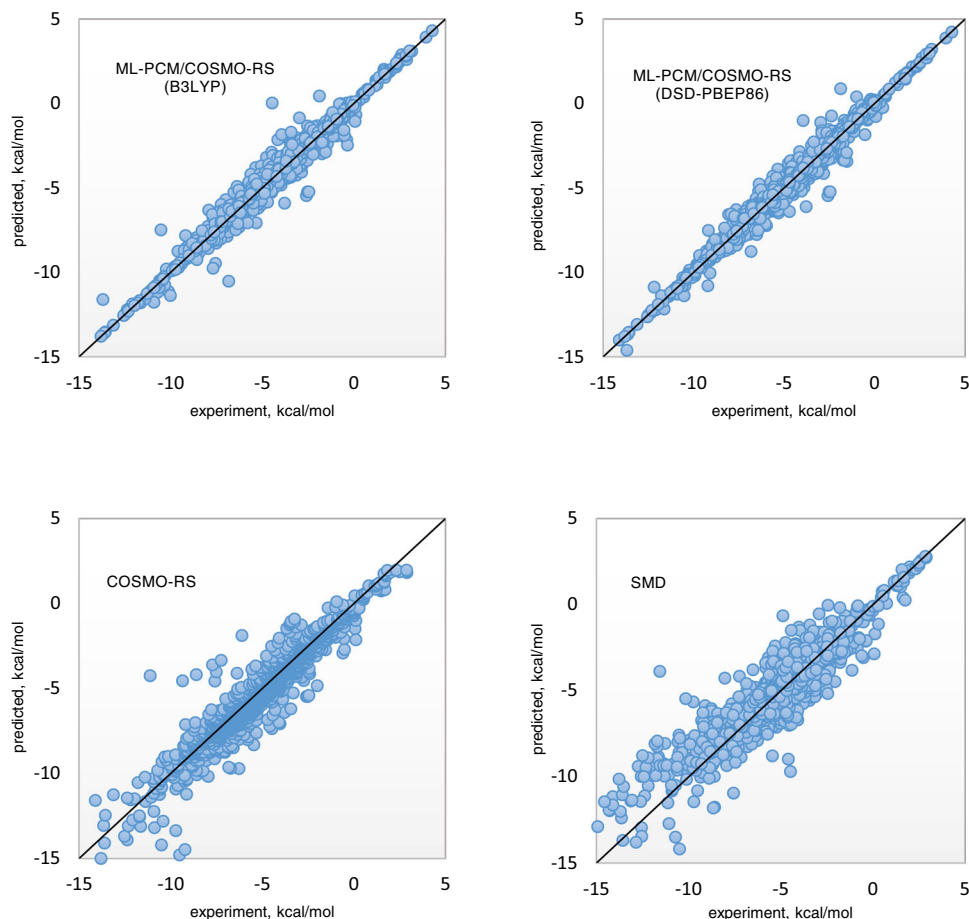


Fig. 2 Comparison of experimentally determined and predicted solvation free energies for various solvation models. The results show a higher correlation between the experimentally determined and predicted data for the proposed machine learning solvation models compared to the SMD or COSMO-RS models.

accurate predictions of solvation free energy with MUE as low as 0.2439 kcal/mol for a large dataset of 2493 binary mixtures of 435 neutral solutes and 91 solvents from diverse chemical families.

Methods

Dataset. To benchmark our results, we used the solvation free energy data of 2493 binary mixtures of 435 neutral solutes and 91 solvents from diverse chemical families available in the Minnesota solvation database⁴. The full list of the studied samples can be found as Supplementary Data 1.

Computational details. The performance of models is reported as mean unsigned error (MUE) and root mean squared error (RMSE) defined as:

$$MUE = \frac{1}{N} \sum (|y_i^{\text{exp}} - y_i^{\text{pred}}|) \quad (3)$$

$$RMSE = \left(\frac{1}{N} \sum \left((y_i^{\text{exp}} - y_i^{\text{pred}})^2 \right) \right)^{\frac{1}{2}} \quad (4)$$

where y_i^{exp} and y_i^{pred} are experimentally determined and predicted solvation free energies, respectively.

Prior to SCRF computations, all solute geometries were optimized in vacuo at the B3LYP/6-31 G* level of theory. Using the optimized structures, the SCRF principal energy components listed in Table 4 were computed for each compound at the B3LYP/6-31 G* and DSD-PBEP86-D3/def2TZVP levels of theory. The latter method as a double hybrid has been shown to yield more precise charge distributions and energy estimations compared to lower-rung DFT or MP2 methods, for a cost comparable to that of the MP2 calculation²⁵.

The SCRF energy components listed in Table 4 were computed for two widely accepted polarizable continuum models, namely the IEF-PCM and CPCM, as implemented in Gaussian 16 (ref. 26). For CPCM, the default value of zero is considered as the scaling factor x in relationship (2). However, a value of 0.5 has been shown to be a more reasonable choice for this scaling factor^{11,27}. Therefore, in

addition to the default implementation of CPCM in Gaussian 16, we also employed a CPCM model with a scaling factor of $x=0.5$ and denote it by CPCM _{$x=0.5$} . For that, we replaced the original dielectric constant of the solvent with an effective dielectric constant $\tilde{\epsilon}(\epsilon, x)$ calculated via:

$$\tilde{\epsilon}(\epsilon, x) = \frac{\epsilon + x}{x + 1} \quad (5)$$

as suggested by Klamt et al.¹¹. For comparison purposes, we also calculated the solvation free energy via the SMD approach.

We employed feed-forward neural networks to map the relationship between the solvation free energy and the calculated SCRF energy components, which in addition to the solvation free energy estimated by the applied continuum solvation model and to the dielectric constant of the solvent, comprised our model inputs.

The obtained pool of model inputs was further screened using the Minimum Redundancy and Maximum Relevance (MRMR) algorithm²⁸ resulting in various 8–16 membered combinations of those variables. MRMR is a highly efficient algorithm for selecting most effective sets of variables for developing robust machine-learning-based models²⁹. For each number of selected variables, 25 different settings of the MRMR algorithm were applied, distinguished by the employed quantization level, level of dependency, forward or backward variable selection and considering pseudo-samples based on Bayesian statistics or not²⁸. In many cases, this resulted in diversely selected set of variables, even for the same applied level of theory and continuum solvation model.

In the next step, various configurations of neural network models were set up and their reliability were examined with a demanding procedure based on the guidelines presented in a previous study³⁰. Accordingly, we assigned large parts of the dataset for test (25%) and validation (15%), and only 60% of the dataset compounds were used for training the models.

To improve the transferability of the developed models for out-of-sample predictions, validation and test sets were selected in a way to include either solvent or solute elements not available in the training set.

We employed Levenberg-Marquardt backpropagation and Gradient descent backpropagation training algorithms, and hidden layer transfer functions of the logarithm-sigmoid and tangent-sigmoid types³¹. We only employed neural

Table 3 Comparison of the results of the new method with other models.

Method	Source	Nr. Samples	Nr. Solvents	Nr. Solutes	Deviation measure	Deviation (kcal/mol)
ML-PCM/COSMO-RS(DSD-PBEP86)	Present study	2224	88	300	MUE	0.24387
ML-PCM/COSMO-RS(B3LYP)	Present study	2224	88	300	RMSE	0.37252
ML-PCM (DSD-PBEP86)	Present study	2488	91	435	MUE	0.26057
ML-PCM (B3LYP)	Present study	2493	91	435	RMSE	0.43623
Other models found in the literature:						
Machine learning	Vermeire and Green ²³	10145	291	1368	MUE	0.40011
					RMSE	0.56014
COSMO-RS	Klamt and Diedenhofen ²⁴	2346	91	318	MUE	0.52526
					RMSE	0.75112
SM12	Marenich et al. ³²	2403	91	352	MUE	0.21
QSPR	Borhani et al. ²⁰	1777	210	295	RMSE	0.44
					MUE	0.42145
					RMSE	0.69644
DCOSMO-RS	Klamt and Diedenhofen ²⁴	2346	91	318	MUE	0.5457-0.6717
					RMSE	0.43
SMD (B3LYP)	Present study	2493	91	435	MUE	0.52
					RMSE	0.6584
SMD (DSD-PBEP86)	Present study	2488	91	435	RMSE	0.99724
					MUE	0.78623
					RMSE	1.1633
Feature Functional Theory	Wang et al. ¹⁸	668	1 (water)	668	MUE	0.85396
kernel-based machine learning	Rauer and Bereau ²²	355	1 (water)	355	RMSE	1.3362
atoms-in-molecules neural network	Zubatyuk et al. ³³	-	-	414	MUE	1.05
Structure-Property Relationship	Hutchinson and Kobayashi ²¹	-	1 (water)	-	MUE	1.06
CPCM(B3LYP)	Present study	2493	91	435	RMSE	1.1
					MUE	1.65
PCM(B3LYP)	Present study	2493	91	435	RMSE	2.6942
					MUE	3.1733
CPCM _{x=0.5} (B3LYP)	Present study	2493	91	435	MUE	2.9054
					RMSE	3.3948
CPCM(DSD-PBEP86)	Present study	2488	91	435	MUE	2.9130
					RMSE	3.3985
PCM (DSD-PBEP86)	Present study	2488	91	435	MUE	2.9651
					RMSE	3.4426
CPCM _{x=0.5} (DSD-PBEP86)	Present study	2488	91	435	MUE	3.1569
					RMSE	3.6445
					MUE	3.1611
					RMSE	3.6466

Table 4 The components of the continuum solvation model.

1	Solvation free energy calculated by the continuum solvation model
2	$\langle \Psi^{(0)} H \Psi^{(0)} \rangle$
3	$\langle \Psi^{(0)} H + V^{(0)} / 2 \Psi^{(0)} \rangle$
4	$\langle \Psi^{(0)} H + V^{(1)} / 2 \Psi^{(0)} \rangle$
5	$\langle \Psi^{(1)} H \Psi^{(1)} \rangle$
6	$\langle \Psi^{(1)} H + V^{(1)} / 2 \Psi^{(1)} \rangle$
7	Interaction energy of unpolarized solute and polarized solvent
8	Interaction energy of polarized solute and polarized solvent
9	Solute polarization energy
10	Total electrostatic interaction energy
11	Cavity surface area
12	Cavity volume
13	Total kinetic energy
14	Total potential energy
15	Sum of kinetic and potential energy

networks with one hidden layer and 1 to 140 neurons in the hidden layer, with intervals of 10 neurons for ANNs with more than 50 neurons in the hidden layer. For each neural network configuration, training was carried out for 60 randomly selected training, validation and test sets, and for each one 40 different initializations of weight and bias constants of the neural networks were made. Above all, to avoid getting misleading data affected by favorable or unfavorable division of dataset into training, validation and test sets, the post validation strategy proposed in a previous study³⁰ was carried out. Accordingly, during the initial training of the neural networks, for the models which yielded mean absolute percentage errors lower than 22%, the final optimized weights and bias constants of the neural network models were recorded. These recorded constants were used as the initial guess to train, validate and test the same neural network configurations but under 100 different randomly selected training, validation and test sets. The models for which in at least 80 out of 100 iterations their test and training sets

errors had the same means and variances as evaluated by the two sample t-test method with 5% significance level were considered as reliably trained models. For them, the average of the ANN-predicted results in all repeats were reported as the performance of that model. Setting up and running the neural network models were implemented in Matlab software. A freely available C++ code for practical use of our proposed ML-PCM models, with detailed user instructions, is provided in Supplementary Software 1.

All the computations were carried out on the High Performance Computing center clusters of the Christian-Albrechts-University of Kiel.

Data availability

All data produced in this study are available and can be provided by contacting the corresponding author.

Code availability

The source file of the C++ code developed for implementing the proposed method with detailed used instructions are available in Supplementary Software 1 or can be provided by contacting the corresponding author.

Received: 11 August 2020; Accepted: 12 May 2021;

Published online: 18 June 2021

References

- Dittner, M. & Hartke, B. Globally optimal catalytic fields—inverse design of abstract embeddings for maximum reaction rate acceleration. *J. Chem. theory Comput.* **14**, 3547–3564 (2018).
- Gauthier, J. A., Dickens, C. F., Chen, L. D., Doyle, A. D. & Nørskov, J. K. Solvation effects for oxygen evolution reaction catalysis on IrO₂ (110). *The. J. Phys. Chem. C.* **121**, 11455–11463 (2017).
- Sakong, S. & Groß, A. The importance of the electrochemical environment in the electro-oxidation of methanol on Pt (111). *ACS Catal.* **6**, 5575–5586 (2016).

- Marenich, A. V., Cramer, C. J. & Truhlar, D. G. Universal solvation model based on solute electron density and on a continuum model of the solvent defined by the bulk dielectric constant and atomic surface tensions. *J. Phys. Chem. B* **113**, 6378–6396 (2009).
- Cramer, C. J. & Truhlar, D. G. A universal approach to solvation modeling. *Acc. Chem. Res.* **41**, 760–768 (2008).
- Chipot, C. & Pohorille, A. *Free energy calculations*. (Springer, 2007).
- Pohorille, A., Jarzynski, C. & Chipot, C. Good practices in free-energy calculations. *J. Phys. Chem. B* **114**, 10235–10253 (2010).
- Miertuś, S., Scrocco, E. & Tomasi, J. Electrostatic interaction of a solute with a continuum. A direct utilization of AB initio molecular potentials for the prevision of solvent effects. *Chem. Phys.* **55**, 117–129 (1981).
- Tomasi, J., Mennucci, B. & Cammi, R. Quantum mechanical continuum solvation models. *Chem. Rev.* **105**, 2999 (2005).
- Barone, V. & Cossi, M. Quantum calculation of molecular energies and energy gradients in solution by a conductor solvent model. *J. Phys. Chem. A* **102**, 1995–2001 (1998).
- Klamt, A., Moya, C. & Palomar, J. A comprehensive comparison of the IEFPCM and SS (V) PE continuum solvation methods with the COSMO approach. *J. Chem. Theory Comput.* **11**, 4220–4225 (2015).
- Klamt, A. & Jonas, V. Treatment of the outlying charge in continuum solvation models. *J. Chem. Phys.* **105**, 9972–9981 (1996).
- Klamt, A. Conductor-like screening model for real solvents: a new approach to the quantitative calculation of solvation phenomena. *J. Phys. Chem.* **99**, 2224–2235 (1995).
- Klamt, A., Jonas, V., Bürger, T. & Lohrenz, J. C. Refinement and parametrization of COSMO-RS. *The J. Phys. Chem. A* **102**, 5074–5085 (1998).
- Mennucci, B., Cammi, R. & Tomasi, J. Excited states and solvatochromic shifts within a nonequilibrium solvation approach: A new formulation of the integral equation formalism method at the self-consistent field, configuration interaction, and multiconfiguration self-consistent field level. *J. Chem. Phys.* **109**, 2798–2807 (1998).
- Cances, E., Mennucci, B. & Tomasi, J. A new integral equation formalism for the polarizable continuum model: theoretical background and applications to isotropic and anisotropic dielectrics. *J. Chem. Phys.* **107**, 3032–3041 (1997).
- Suárez, E., Díaz, N. & Suárez, D. Entropy calculations of single molecules by combining the rigid-rotor and harmonic-oscillator approximations with conformational entropy estimations from molecular dynamics simulations. *J. Chem. Theory Comput.* **7**, 2638–2653 (2011).
- Wang, B., Wang, C., Wu, K. & Wei, G. W. Breaking the polar-nonpolar division in solvation free energy prediction. *J. Comput. Chem.* **39**, 217–233 (2018).
- Dzubiella, J., Swanson, J. M. & McCammon, J. Coupling hydrophobicity, dispersion, and electrostatics in continuum solvent models. *Phys. Rev. Lett.* **96**, 087802 (2006).
- Borhani, T. N., García-Muñoz, S., Luciani, C. V., Galindo, A. & Adjiman, C. S. Hybrid QSPR models for the prediction of the free energy of solvation of organic solute/solvent pairs. *Phys. Chem. Chem. Phys.* **21**, 13706–13720 (2019).
- Hutchinson, S. T. & Kobayashi, R. Solvent-specific featurization for predicting free energies of solvation through machine learning. *J. Chem. Inf. Modeling* **59**, 1338–1346 (2019).
- Rauer, C. & Bereau, T. Hydration free energies from kernel-based machine learning: compound-database bias. *J. Chem. Phys.* **153**, 014101 (2020).
- Vermeire, F. H. & Green, W. H. Transfer learning for solvation free energies: from quantum chemistry to experiments. *Chem. Eng. J.* **418**, 129307 (2021).
- Klamt, A. & Diedenhofen, M. Calculation of solvation free energies with DCOSMO-RS. *J. Phys. Chem. A* **119**, 5439–5445 (2015).
- Kozuch, S. & Martin, J. M. DSD-PBEP86: in search of the best double-hybrid DFT with spin-component scaled MP2 and dispersion corrections. *Phys. Chem. Chem. Phys.* **13**, 20104–20107 (2011).
- Frisch, M. et al. (Gaussian, Inc. Wallingford, CT, 2016).
- Cossi, M., Rega, N., Scalmani, G. & Barone, V. Energies, structures, and electronic properties of molecules in solution with the C-PCM solvation model. *J. Comput. Chem.* **24**, 669–681 (2003).
- Peng, H., Long, F. & Ding, C. Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy. *IEEE Trans. Pattern Anal. Mach. Intell.* **27**, 1226–1238 (2005).
- Brown, Gavin. A new perspective for information theoretic feature selection. *Artificial intelligence and statistics*, PMLR, pp. 49–56 (2009).
- Alibakshi, A. Strategies to develop robust neural network models: prediction of flash point as a case study. *Anal. Chim. Acta* **1026**, 69–76 (2018).
- Demuth, H. & Beale, M. *Neural Network Toolbox For Use with Matlab--User'S Guide Verion 3.0*. (1993).
- Marenich, A. V., Cramer, C. J. & Truhlar, D. G. Generalized born solvation model SM12. *J. Chem. Theory Comput.* **9**, 609–620 (2013).
- Zubatyyuk, R., Smith, J. S., Leszczynski, J. & Isayev, O. Accurate and transferable multitask prediction of chemical properties with an atoms-in-molecules neural network. *Sci. Adv.* **5**, eaav6490 (2019).

Acknowledgements

The authors wish to thank Karsten Balzer at the high performance computing center of Kiel University for his support and assistance in running the computations there. The Authors wish to thank the referees for their careful review of our work and fruitful discussions and comments.

Author contributions

A.A. has contributed to method development, carried out the computations and contributed to writing the manuscript. B.H. supervised the project and contributed to method development and writing the manuscript.

Funding

Open Access funding enabled and organized by Projekt DEAL.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41467-021-23724-6>.

Correspondence and requests for materials should be addressed to A.A.

Peer review information *Nature Communications* thanks John Keith and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. Peer reviewer reports are available.

Reprints and permission information is available at <http://www.nature.com/reprints>

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021