

Advancing Force Fields Parameterization: A Directed Graph Attention Networks Approach

Gong Chen,* Théo Jaffrelot Inizan, Thomas Plé, Louis Lagardère, Jean-Philip Piquemal, and Yvon Maday*



Cite This: *J. Chem. Theory Comput.* 2024, 20, 5558–5569



Read Online

ACCESS |



Metrics & More

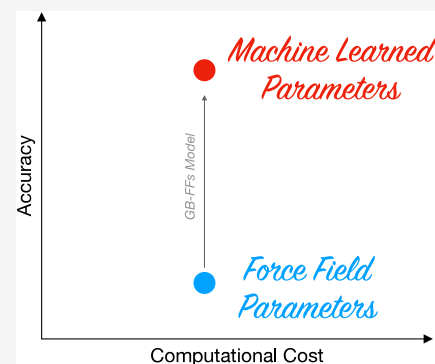


Article Recommendations



Supporting Information

ABSTRACT: Force fields (FFs) are an established tool for simulating large and complex molecular systems. However, parametrizing FFs is a challenging and time-consuming task that relies on empirical heuristics, experimental data, and computational data. Recent efforts aim to automate the assignment of FF parameters using pre-existing databases and on-the-fly *ab initio* data. In this study, we propose a graph-based force field (GB-FFs) model to directly derive parameters for the Generalized Amber Force Field (GAFF) from chemical environments and research into the influence of functional forms. Our end-to-end parametrization approach predicts parameters by aggregating the basic information in directed molecular graphs, eliminating the need for expert-defined procedures and enhances the accuracy and transferability of GAFF across a broader range of molecular complexes. Simulation results are compared to the original GAFF parametrization. In practice, our results demonstrate an improved transferability of the model, showcasing its improved accuracy in modeling intermolecular and torsional interactions, as well as improved solvation free energies. The optimization approach developed in this work is fully applicable to other nonpolarizable FFs as well as to polarizable ones.



INTRODUCTION

The high numerical cost of quantitative *ab initio* methods restricts their application to small systems, composed of a few hundred atoms,^{1,2} limiting the ability to study larger molecular systems. As an alternative, force fields (FFs) have emerged as valuable tools, employing physically motivated functional forms to model potential energy surfaces. These FFs are parametrized to match *ab initio* as well as experimental data, offering a computationally cheaper alternative for simulating diverse systems, ranging from biology to polymers and complex materials. Indeed, there is a wide range of FFs developed for different purposes and compounds.

FFs can be categorized into two main families. The most commonly used FFs are known as classical, or nonpolarizable, FFs, that include AMBER,^{3,4} CHARMM,^{5,6} and GAFF.^{7,8} These nonpolarizable FFs employ a combination of fixed-charge Coulomb potential and Lennard-Jones interactions to model intermolecular interactions. They are highly efficient numerically, enabling simulations of very large systems over long time scales.^{9,10} However, their simple functional form lacks polarization and many-body effects, which are crucial for accurately describing complex phenomena such as pi-stacking or allosteric effects.^{11,12}

On the other hand, there are polarizable force fields (PFFs) such as AMOEBA,^{13,14} AMOEBA+,^{15,16} CHARMM Drude,¹⁷ and SIBFA.^{18,19} These force fields have been specifically developed to incorporate polarization and many-body effects.

This enhanced flexibility and accuracy comes at a higher computational cost compared to nonpolarizable FFs. Nevertheless, PFFs provide a more comprehensive representation of intermolecular interactions and are particularly suitable for studying complex systems.^{20–23}

In recent years, significant attention and resources have been devoted to the development of Machine Learning Potentials (MLPs), aiming to bridge the accuracy and generality gap between FFs and *ab initio* methods.^{24–26} MLPs employ flexible functional forms from the field of Machine Learning (ML) to accurately fit *ab initio* energies or forces. With respect to the *ab initio* methods, they offer a favorable balance between computational efficiency and accuracy, circumventing the need for empirical functional forms used in FFs. MLPs possess the ability to capture complex interactions,^{27–31} that are challenging to model when using traditional nonpolarizable FFs.

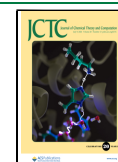
However, the accuracy of MLPs depends on the quality of the training data and the architecture of the ML model, which can limit their transferability. Moreover, MLP models are often

Received: December 29, 2023

Revised: May 23, 2024

Accepted: May 24, 2024

Published: June 14, 2024



difficult to interpret, posing challenges in identifying and understanding the underlying physics and chemistry of the systems under study. To address this, one strategy involves constructing hybrid models that combine MLPs and FFs, bridging the gap between the short-range accuracy of quantum mechanics and the computational efficiency of FFs.^{32–35}

However, FFs remain widely used across various areas due to their computational efficiency, and such hybrid MLP FF strategies still require an accurate FF parametrization. Parameterizing an FF is a challenging task since its accuracy and transferability heavily depend on the quality of its parameters. This process is time-consuming, often taking years, as it relies on empirical heuristics, experimental data, and computational data. FFs have established robust parametrization procedures, such as Antechamber for GAFF⁸ and polype2 for AMOEBA.^{36,37} Additionally, these FFs rely on empirically assigned atom types and atom classes to assign parameters (e.g., bonds, angles) that are specific to the local topology. To enhance the generalization and reliability of FFs, one tendency is to expand the atomic type space. However, this leads to an increasing number of possible valence compositions, introducing significant complexity into the parameter fitting process. Moreover, even with modern parameter optimization frameworks³⁸ and sufficient data, FF parameters defined by fixed atom types can sometimes suffer from low transferability.

Advances in the computational efficiency and scalability of *ab initio* methods have also provided new opportunities to enhance the transferability and accuracy of FFs by building larger and more accurate databases.^{39,40} A question arises: can ML be used to enhance the prediction of FF parameters by training on *ab initio* databases? Recently, the Espaloma model combined Graph Neural Networks (GNNs) and automatic differentiation to predict FF parameters.⁴¹ By focusing on intramolecular interactions, they demonstrated that GNNs can effectively predict FF parameters based on potential energies.

Here, we propose a Graph-Based Force Fields (GB-FFs) model for FF parametrization. GB-FFs automatically derive accurate FF parameters using atom features and bond features and aim to extend the generalization of FFs by using a novel Directed Graph multihead Attention Network architecture. It serves as a continuous alternative to traditional discrete atom typing schemes, eliminating the need for assigning atom types and obtaining FF parameters directly from atomic representations.

To assess the overall quality of our GB-FFs model, we compared its performance to the original GAFF parametrization using the highly parallel Tinker-HP GPU package.^{42,43} Our model is freely available (MIT License) on GitHub at <https://github.com/GongCHEN-1995/GB-FFs-Model>, and a tutorial is included for generating GAFF parameters from a given SMILES, as well as for fine-tuning the model on a newer database.

This work introduces several significant contributions. First, it improves the existing GAFF by refining its parameters, thereby improving its accuracy. Second, it treats molecules as directed molecular graphs and employs a self-attention mechanism to effectively aggregate information, enhancing the model's performance. Third, it incorporates an efficient charge transfer procedure to enhance the prediction of fixed atomic charges with an $O(N)$ complexity. Furthermore, the model's validation spans across various properties, including both quantum-mechanical and experimental properties, such as hydration free energies, across diverse molecular systems.

Additionally, the versatility and ease of use of the model allow its extension to other force fields. Lastly, the paper explores the limits of GAFF by modifying its functional forms, shedding light on its strengths and weaknesses.

METHODS

In the following sections, we provide a brief introduction to GAFF and present the GB-FFs model.

The General AMBER Force Field (GAFF). GAFF⁷ is among the most popular classical FFs to simulate organic molecules. It is an extension of the Amber force field.⁴ GAFF is specifically designed to be compatible with a broad range of organic molecules, including drug-like compounds, carbohydrates, and nucleic acids.

GAFF incorporates a comprehensive set of parameters for bond stretching, angle bending, torsional, and nonbonded interactions. These parameters allow for accurate modeling and simulation of the behavior of organic molecules under various conditions (e.g., high pressure, low temperature). Due to its computational efficiency, relative reliability and its relatively simple functional form, GAFF has been widely used and implemented in many popular molecular simulation software packages, such as AMBER,⁴ GROMACS,⁴⁴ CHARMM,⁶ and Tinker/Tinker-HP.^{42,43,45,46} Another advantage of GAFF is the public availability of its parameters. Thus, facilitating its widespread use within the scientific community.

In GAFF, the angle bending, bond stretching bonded interactions are modeled using a harmonic potential making it nonreactive thus greatly simplifying the parametrization process. The torsional potential is expressed as a Fourier series. For nonbonded interactions, the van der Waals (vdW) interactions are described by a 12–6 Lennard-Jones potential, and electrostatic interactions are modeled using atom-centered fixed charges.

$$E_{\text{potential}} = \sum_{\text{bonds}} K_r (r - r_{eq})^2 + \sum_{\text{angles}} K_\theta (\theta - \theta_{eq})^2 + \sum_{\text{dihedrals}} \frac{V_n}{2} [1 + \cos(n\phi - \gamma)] + \sum_{i < j} \left[\epsilon_{ij} \left(\frac{\sigma_{ij}^{12}}{R_{ij}^{12}} - 2 \frac{\sigma_{ij}^6}{R_{ij}^6} \right) + \frac{q_i q_j}{\epsilon R_{ij}} \right] \quad (1)$$

where r_{eq} and θ_{eq} are equilibrium structural parameters; K_r , K_θ , V_n are so-called “force constants”; n is multiplicity and γ is phase angle for torsional angle parameters. The ϵ , σ , and q parameters characterize the nonbonded vdW potential. ϵ , σ follow Lorentz–Berthelot combination rules.^{47,48} The GAFF parameters $\{K_r, r_{eq}, K_\theta, \theta_{eq}, V_n, \epsilon, \sigma\}$ are directly read from parameters table according to corresponding atom types while $n = 1, 2, 3, 4$ and $\gamma = 0$ or π .

The parametrization process of GAFF starts by assigning partial charges. In the early stages of GAFF, Hartree–Fock (HF) with the 6-31G* basis set were used to generate electrostatic potentials from which restrained electrostatic potential (RESP) charge^{49,50} fits were derived. This process proved to be expensive, especially for large molecules and led to the development of the AM1-BCC charge scheme that approximate HF/6-31G* RESP computation by first calculating charges using the AM1 semiempirical model and correct it via bond charge corrections.^{51,52}

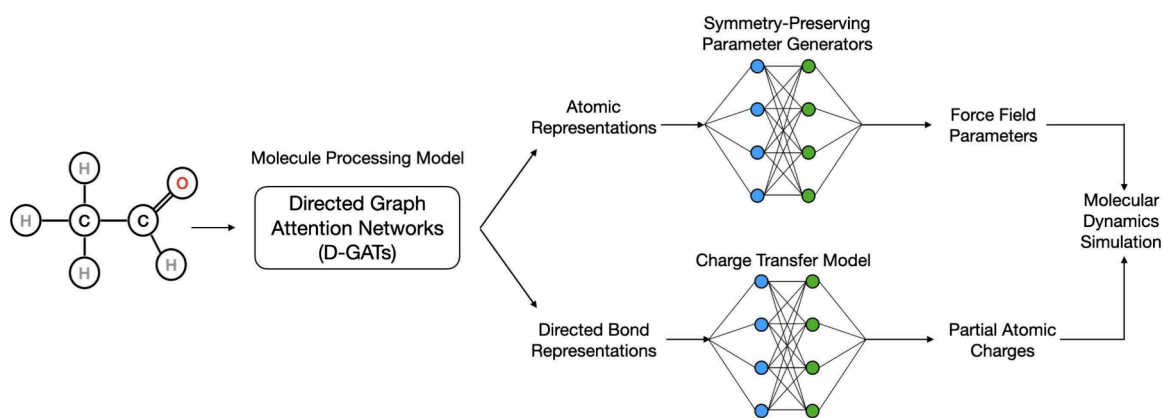


Figure 1. Framework of graph-based force fields (GB-FFs) model: It consists of molecule processing model, the symmetry-preserving parameter generator and charge transfer model.

In GAFF, the equilibrium bond length θ_{eq} is fitted through experimental data from X-ray and neutron diffraction, as well as MP2/6-31G* computations. On the other hand, bond angle parametrization uses reference from the Cambridge Structure Database, empirical rules and MP2/6-31G* computations. Finally, the strategy for developing torsional angle parameters involves performing torsional angle scanning and fitting the parameters to accurately reproduce the rotational profile obtained from MP2/6-31G* calculations. The vdW parameters are the same as those used by AMBER and thus extracted from a database.

While in this paper we specifically used GB-FFs on the second generation of GAFF parameters, this framework is general and can be used on other types of FFs including polarizable ones.

Graph-Based Force Fields Model, a Universal Parametrization Procedure. GNNs have been proved to be an efficient and powerful way to detect chemical environment and to extract molecular properties.^{53–56} In addition, GNNs also have shown potential in expressing atoms' representations and bonds' representations.^{53,54,57}

The model is composed of three modules: molecule processing model, symmetry-preserving parameter generator, and charge transfer model (see Figure 1). These components will be discussed in the following sections.

GB-FF's runtime complexity for a molecule with N atoms is $O(N)$ and processing a molecule with $N = 50$ atoms takes 0.04 s on a single GPU V100. In comparison, Antechamber takes 111 s to assign atom types and charges for the same molecule and its AM1-BCC charge model has a computational complexity of $O(N^2)$.

Molecules Processing Model. Assigning atom types and determining FF parameters are typical atom-level tasks. Building on the notion of directed bonds,⁵⁸ we used a Directed Graph Attention neTworks (D-GATs) model.⁵⁶ In contrast to other ML-based molecular processing models, D-GATs exhibit a remarkable ability to discern local chemical environments and eliminate unnecessary message flow. They have consistently outperformed state-of-the-art benchmarks in 13 out of 15 molecular property prediction tasks.⁵⁶

To enhance the robustness, we employ the Smooth Maximum Unit (SMU)⁵⁹ as an activation function. SMU offers a smooth approximation to the entire Maxout family, including ReLU, Leaky ReLU, and their variants.⁶⁰ Our goal is to predict a set of parameters that can bring molecular dynamic

simulation results closer to the *ab initio* data. The molecular potential energy surface is highly sensitive to these predicted parameters. Our test demonstrated that the discontinuity in the Maxout function can hinder the convergence of the models' loss.

To ensure compatibility with GAFF, we focused on compounds composed of C, N, O, H, S, P, F, Cl, Br, and I. Using the RDKit package,⁶¹ we extracted fundamental atomic and bond features (see Table 1). These features, in

Table 1. Input Features of the GB-FFs Model

Atom Features	Size(38)	Descriptions
atom symbol	11	[UNK,H,C,N,O,F,P,S,Cl,Br,I] (one-hot)
degree	6	number of covalent bonds [0, 1, 2, 3, 4, 5] (one-hot)
formal charge	7	[-3,-2,-1,-0,1,2,3] (one-hot)
hybridization	8	[unspecified, s, sp, sp2, sp3, sp3d, sp3d2, other] (one-hot)
chirality	4	[unspecified, tetrahedral_CW, tetrahedral_CCW, other] (one-hot)
ring	1	whether the atom is in ring [0/1] (one-hot)
aromaticity	1	whether the atom is part of an aromatic system [0/1] (one-hot)
Bond Features	Size(12)	Descriptions
bond type	4	[single, double, triple, aromatic] (one-hot)
conjugation	1	whether the bond is conjugated [0/1] (one-hot)
ring	1	whether the bond is in ring [0/1] (one-hot)
stereo type	6	[StereoNone, StereoAny, StereoZ, StereoE, Stereocis, Stereotrans] (one-hot)

conjunction with the molecular graph represented in Lewis structure,⁶² were then input into the GB-FFs model. For more information about the accuracy of the predicted atom types and features, refer to the section "Recovering atom types in GAFF" in the Supporting Information (SI). The model's outputs include atomic and bond representations. For this article, we specifically employ the directed bond representations to incorporate chemical information and the atomic representations for predicting FF parameters (as illustrated in Figure 2).

FFs are intricate and highly parameter-sensitive. To enhance the models' expressive capacity and expand their receptive field, we employ a hierarchical structure consisting of two stacked layers: Small and Larger Layers. Both layers share the same model architecture but operate in distinct dimensions.

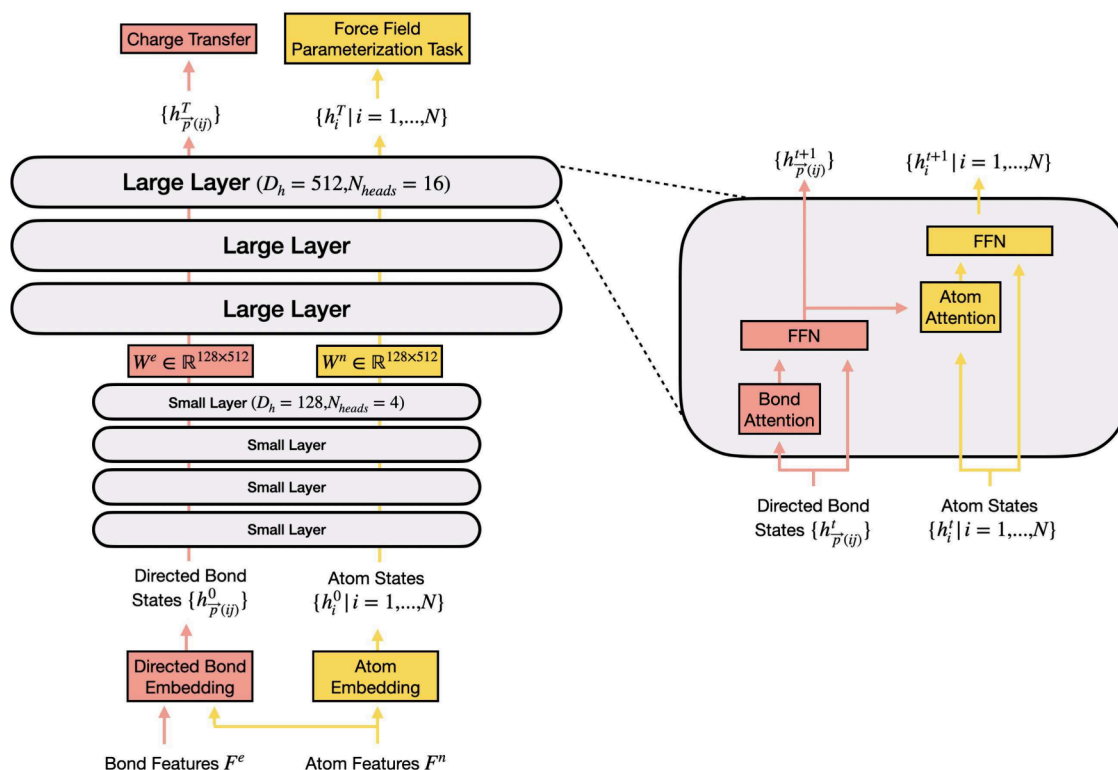


Figure 2. Molecule processing model: D-GATs model with D_h being the dimension of model and N_{heads} the number of heads in the multiattention mechanism. Between two stacked layers, W^e and W^n aim to convert the dimension of the embeddings. The stacked layers consist of several interaction layers.

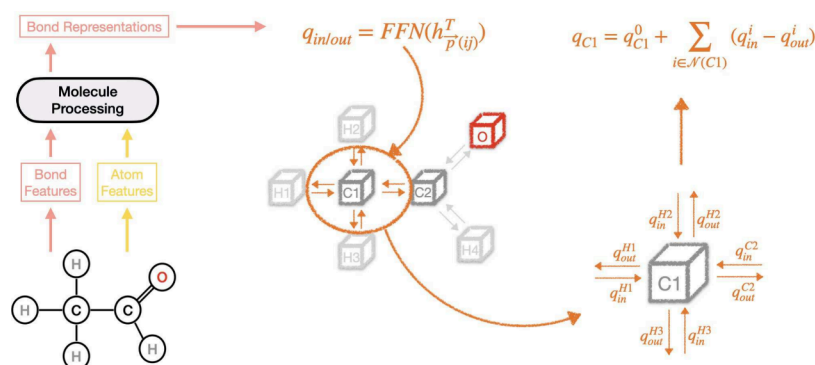


Figure 3. Charge transfer model: The charge is allowed to transfer between connected atoms and the charge in/out is directly calculated by the directed bond embeddings. The final partial charge of atom is the original formal charge plus charge flows in and minus the charge flows out.

Between the two stacked layers, we incorporate linear transformations (W^e and W^n) to convert the dimensions of atomic and bond representations.

As depicted in Figure 2, the Large Layers consist of 3 interaction layers with a model dimension (D_h) of 512. They play a central role in detecting chemical environments and forming atomic representations. The Small Layers, on the other hand, employ 4 interaction layers with a model dimension (D_h) of 128 and four attention heads. These layers are primarily used for embedding initialization and demand minimal computational resources.

We denote the output atomic representations as $h^T = h_i^T | i = 1, \dots, N$, while the output directed bond representations are denoted as $h_{p(ij)}^T$ for all connected atoms i and atoms j . It is important to emphasize that the

directionality in bond representation is critical, with $h_{p(ij)}^T$ indicating the bond from atom i to atom j .

Charge Transfer Model. To ensure the net charge of the molecule aligns with the actual scenario and to improve the physical meaning of charge distribution, we used directed bond states $\{h_{p(ij)}^T\}$ to predict the charge transfer between connected atoms. Our molecular processing model is based on directed graphs, eliminating the need for additional operations, and it can predict the charge transfer from one atom to its neighbors.

As illustrated in Figure 3, the directed bond features obtained from Figure 2 are fed into a feed-forward neural network (FFN) to determine the charge transfer in the corresponding bond direction. The final atomic charge is obtained by summing the original formal charge and the incoming charges while subtracting the outgoing charges.

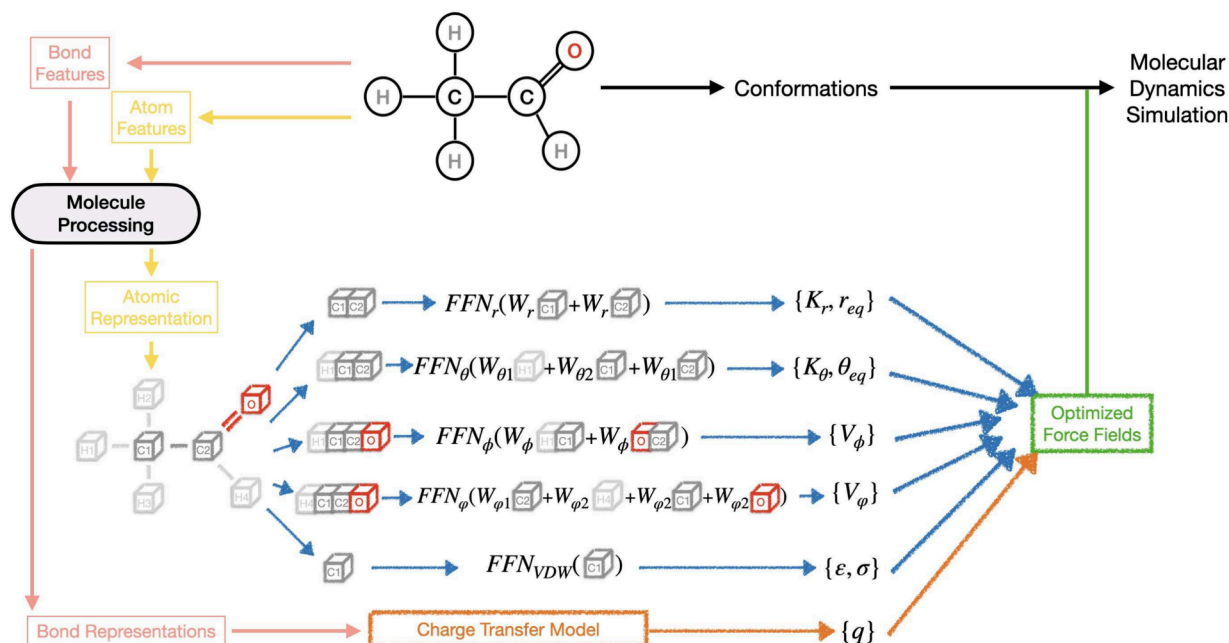


Figure 4. Symmetry-preserving parameter generator: For a specified molecule, we input atom and bond features to hierarchical D-GATs and obtain the atomic representations and directed bond representations. The symmetry-preserving parameter generator predicts all FF parameters, which can be used to do molecular dynamic simulation.

We used the AM1-BCC charge model^{49,50} to calculate the original partial atomic charges. Our predicted GB-FFs charges were compared with both AM1-BCC charges and the accurate wB97x/def2-TZVPP Minimal Basis Iterative Stockholder (MBIS) charges. We found that the charge transfer model performs comparably to the AM1-BCC charge.

Symmetry-Preserving Parameter Generator. The number of FF parameters depends on the molecule's geometry. According to the molecular geometry, we use RDKit to list all the combinations of bonds, angles, dihedrals, and nonbonded interaction pairs. We input atomic representations into the parameter generator based on the identified structure to predict all potential bonds, angles, dihedrals, and nonbonded parameters.

The parameter generator needs to ensure that atom ordering symmetries. For example, when bond parameters are predicted, if we exchange the order of two input atomic representations, the predicted parameters should be invariant. In the Espaloma model,⁴¹ the relevant equivalent atom permutations are enumerated, which has a computational cost. Here, we leverage this by splitting the input atom embeddings based on their intrinsic structure and applying linear transformations to ensure symmetry:

$$h_{r_i} = h_{r_j} = W_r h_i^T + W_r h_j^T \quad (2)$$

$$h_{\theta_{ip}} = h_{\theta_{pj}} = W_{\theta 1} h_i^T + W_{\theta 2} h_j^T + W_{\theta 1} h_p^T \quad (3)$$

$$h_{\phi_{ipq}} = h_{\phi_{qpi}} = W_{\phi} [h_i^T, h_j^T] + W_{\phi} [h_q^T, h_p^T] \quad (4)$$

$$\begin{aligned} h_{\phi_{ipq}} &= h_{\phi_{qpi}} = h_{\phi_{qpi}} = h_{\phi_{ipq}} = h_{\phi_{qpi}} = h_{\phi_{ipq}} \\ &= W_{\phi 1} h_i^T + W_{\phi 1} h_j^T + W_{\phi 2} h_p^T + W_{\phi 1} h_q^T \end{aligned} \quad (5)$$

$$h_{VdW_i} = W_{VdW} h_i^T \quad (6)$$

where $[.,.]$ denote concatenation and

$$h_r, h_{\theta}, h_{\phi}, h_{\phi} \in \mathbb{R}^{D_h}, W_r, W_{\theta 1}, W_{\theta 2}, W_{\phi 1}, W_{\phi 2},$$

$$W_{VdW} \in \mathbb{R}^{D_h \times D_h}, W_{\phi} \in \mathbb{R}^{2D_h \times D_h}$$

These embeddings for bond ($\{h_r\}$), angle ($\{h_{\theta}\}$), torsion ($\{h_{\phi}\}$), improper torsion term ($\{h_{\phi}\}$) and vdW interaction ($\{h_{VdW}\}$) are in the same dimension. Additionally, the number of parameters for each term is fixed (for example, one bond term needs two parameters, $\{K_r, r_{eq}\}$). According to the corresponding embeddings, we can use the fully connected NNs to predict the FF parameters (see Figure 4). Note that in our procedure all FF parameters are optimized together, as explained in the section **Training Strategies** below and presented in more detail in the **Supporting Information**.

We will now present as a perspective our attempts to modify GAFF's functional in order to improve its accuracy and transferability to a broader class of systems.

Perspective for GAFF Functional form. Here, we aim to enhance the accuracy of GAFF by modifying its functional form, especially by focusing on the bonded part. We modified the bond-stretching term of GAFF by using a Morse function and its angle bending energies with a Urey–Bradley term.

Morse Function for Stretching Energy. FFs typically adopt functional forms based on a balance between the computational cost and approximation effectiveness. The Morse function⁶³ provides a more accurate description of the bond potential, particularly for bonds that are stretched beyond their equilibrium values. We employ the following form of the Morse function (the same number of parameters as GAFF):

$$E_{bonds} = \sum_{bonds} \frac{K_r}{4} (e^{-2(r-r_{eq})} - 1)^2 \quad (7)$$

with $\{K_r, r_{eq}\}$ are the parameters directly from GAFF parameters.

In the following, the model is denoted as “GB-FFs Morse”.

Urey–Bradley Terms. The Urey–Bradley (UB) terms⁶⁴ serve as a cross-term addressing 1–3 nonbonded interactions that are not adequately covered by the bond and angle terms. We make use the following function:

$$E_{UB} = \sum_{\text{angles}} K_{UB} \left(\left(\frac{r_{UB_{eq}}}{r_{UB}} \right)^2 - 1 \right)^2 \quad (8)$$

with $\{K_{UB}, r_{UB_{eq}}\}$ are the new FF parameters. The details to calculate $r_{UB_{eq}}$ are presented in the [Supporting Information](#) (Urey–Bradley Terms).

In the following, the model will be denoted as “GB-FFs UB”.

Training Strategies. Our training strategy consists of several stages aimed at enhancing the robustness and performance of our model. In the first stage, the model is trained on the ANI-1 database, one of the largest database of Density Functional Theory (DFT) computations for small organic molecules.⁶⁵ However, due to the relatively low *ab initio* accuracy of the ANI-1 database, we fine-tuned the model on two more chemically accurate databases: SPICE⁶⁶ and DES370K.⁶⁷

The SPICE data set is used to ensure the accuracy on intramolecular interactions of biomolecular systems. SPICE is a collection of DFT data mainly built to train the MLP for simulating drug molecules and proteins. The computations are performed at the ω B97M-D3(BJ) functional^{68,69} with the def2-TZVPPD basis set.^{70,71}

In contrast, the DES370K data set is used to ensures the accuracy of intermolecular interactions. In DES370K, the reference interaction energies for these systems are computed using the highly accurate coupled-cluster singles and doubles with perturbative triples (CCSD(T))⁷² level of theory with a complete basis set (CBS).⁷³ The complexes in this database represent most of the molecular interactions that could occur in chemistry, including electrostatic-dominated (hydrogen bonding), dispersion-dominated, and mixed (electrostatic/dispersion) interactions.

Given that GAFF includes only the elements H, C, N, O, F, P, S, Cl, Br, and I, we exclude the molecules containing elements Li, Na, Mg, K, Ca both from SPICE and DES370K data sets. This left us with a total of 29,389 compounds. The compounds were randomly divided into training/validation/test sets following an 8:1:1 ratio.

Fitting the potential energy and atomic forces simultaneously poses a significant challenge. First, because it is unclear whether a unique global best fit exists and, second, if there exist local ones. To address this challenge, we proposed a multistage training strategy. The first stage consists of pretraining on the large ANI-1 Database: initially, training on FF parameters only (bond, angle, dihedral angle, vdW, AM1-BCC charges) to match GAFF parameters and ensure the physical relevance of the generated parameters; subsequently, training on FF parameters and GAFF energy components; and finally, training on FF parameters as well as *ab initio* energies and forces computed with the ω B97x functional with the 6-31G(d) basis set. Once this pretraining is completed, the second stage involves fine-tuning: training on FF parameters along with energies and forces from the SPICE database and interaction energies from the DES370K database. More information can be found in the [SI](#) (“Pretraining on ANI-1 Database” and “Fine-Tuning Strategy on SPICE and DES370K Databases”). This second stage allows us to capture not only a much higher

level of theory but also more accurate chemical diversity and environments. In each training stage, we kept training on FF parameters to maintain physical awareness and prevent a significant deviation from the GAFF parameters.

RESULTS AND DISCUSSION

The results of the fine-tuning process are shown in [Table 2](#). For the SPICE database, compared to the original GAFF, our

Table 2. RMSE Comparison on SPICE Dataset and DES370K Database’s *Ab Initio* Data for Potential Energy, Atomic Forces, and AM1-BCC Charges, between GAFF and GB-FFs GAFF Models

	SPICE			DES370K	
	Energy (kcal/mol)	Force (kcal/mol/Å)	Charge (C)	Energy (kcal/mol)	Charge (C)
GAFF	5.7804	13.4398		1.1470	
GB-FFs GAFF	2.9706	5.9232	0.0500	1.4146	0.0713

GB-FFs GAFF model significantly reduces the Root Mean Square Error (RMSE) for the energies from 5.8 kcal/mol to less than 3.0 kcal/mol and for the forces from 13.4 kcal/mol/Å to 6.0 kcal/mol/Å. However, for the DES370K database, the RMSE for interaction energy has increased from 1.1 to 1.4 kcal/mol due to the sensitivity of vdW parameters.

Additionally, as stated before, GB-FFs is 3–4 orders of magnitude faster than AM1-BCC calculations using *Antechamber* (the command in AMBER), from 111 to 0.04 s for 50 atoms (more details are in [SI](#) “Computational Resources”). Furthermore, the GB-FFs charges provided by the charge transfer model closely approximate AM1-BCC charges ([Table 2](#)).

Finally, the results presented in [Table 3](#) confirm the improvements in all quantities on these two databases obtained from an improved functional form.

Table 3. RMSE Comparison on SPICE Dataset and DES370K Database’s *Ab Initio* data for Potential Energy, Atomic Forces, and AM1-BCC Charges, with the Improved Functional Forms GB-FFs Morse and GB-FFs UB

	SPICE			DES370K	
	Energy (kcal/mol)	Force (kcal/mol/Å)	Charge (C)	Energy (kcal/mol)	Charge (C)
GB-FFs Morse	2.8812	5.3050	0.0492	1.3884	0.0698
GB-FFs UB	2.5723	4.1416	0.0491	1.0941	0.0526

Intermolecular Interaction Accuracy: S66×8 Benchmark. In previous databases, models are trained and tested on the same databases. In this subsection, we aim to assess the models’ ability to generalize to unseen molecular systems and conformations.

The S66×8 database⁷⁴ comprises 66 dimers positioned at 8 distinct intermolecular distances, resulting in a total of 528 unique structures. The S66×8 database is a widely known reference database for assessing the accuracy of intermolecular interactions.

The minimum distance between two monomers ranges from 0.9 to 2.0 times the equilibrium value. When the intermolecular distance varies, the monomers have fixed

geometries, meaning that the deformation energies of monomers are not considered.

The Mean Absolute Error (MAE) and RMSE on the overall data set are shown in Table 4. Further details regarding the

Table 4. MAE and RMSE of the Interaction Energy on the S66×8 Database, as Well as the MAE and RMSE of the Potential Energy on the Torsion Scan Database for GAFF and the GB-FFs GAFF Models

	S66×8		Torsion Scan	
	MAE (kcal/mol)	RMSE (kcal/mol)	MAE (kcal/mol)	RMSE (kcal/mol)
GAFF	0.9368	1.8388	1.9694	3.5351
GB-FFs GAFF	0.5087	0.8766	0.9892	1.4843

accuracy of predicted forces, including Mean Absolute Percentage Error (MAPE) and atomwise performance, are available in SI “Accuracy of predicted Atomic Forces”. Compared to GAFF, the GB-FFs model reduced by more than half of the RMSE. This demonstrates that the training process significantly benefits the approximation of long-range interactions in fitting potential energy and atomic forces. Further improvements are presented in the SI with the improved functional forms of GB-FFs Morse and GB-FFs UB.

Figure 5 depicts the results for the four dimers. The remaining dimers can be found in SI “Full results on S66×8 database”. The GB-FFs GAFF models almost perfectly reproduce the intermolecular energy surface of the water dimer, which is a critical aspect for simulating solvated biomolecules. In other instances, GB-FFs models are often comparable or outperform GAFF.

Torsion Profiles of 62 Drug-like Fragments. After evaluating intermolecular interactions between molecules, we

now turn our attention to assessing the accuracy of predicting intramolecular interactions.

Torsion energies play a crucial role in biology and in small molecular systems. However, accurately assessing torsional parameters in FF is challenging, as they necessitate computationally expensive calculations and complicated fitting procedure. Additionally, these parameters are highly sensitive to the local chemical environment, making them difficult to transfer across different molecular systems. Consequently, they often rely on simplistic transferability rules, which can lead to inaccuracies.

Thus, achieving accurate torsion profiles while avoiding the need for extensive torsion fitting is of great importance in FF parametrization. In this context, the performance of the GB-FFs parametrization is also evaluated on a highly accurate torsion scan database.⁷⁵ It comprises 62 fragments with drug-like functional groups and their CCSD(T) /CBS single point energies calculated on optimized geometries using MP2^{76,77}/6-311+G**.^{78,79}

For each molecule of the 62 fragments, a specific dihedral angle is varied from -170° to 170° in increments of 10° (the chosen dihedral angle for modification is indicated in SI “Full results on Torsion Scan database”).

The overall performance is recorded in Table 4. Compared with the original GAFF, the GB-FFs models provide FF parameters that better fit the potential energy changes caused by dihedral angle variations. In some cases, although there may be a gap between GB-FFs model’s predicted results and the reference energies (see Figure 6c and 6d), the observed trends in these changes correspond with the actual scenarios.

This assessment aims to highlight the capabilities of the GB-FFs model in accurately capturing torsional energies.

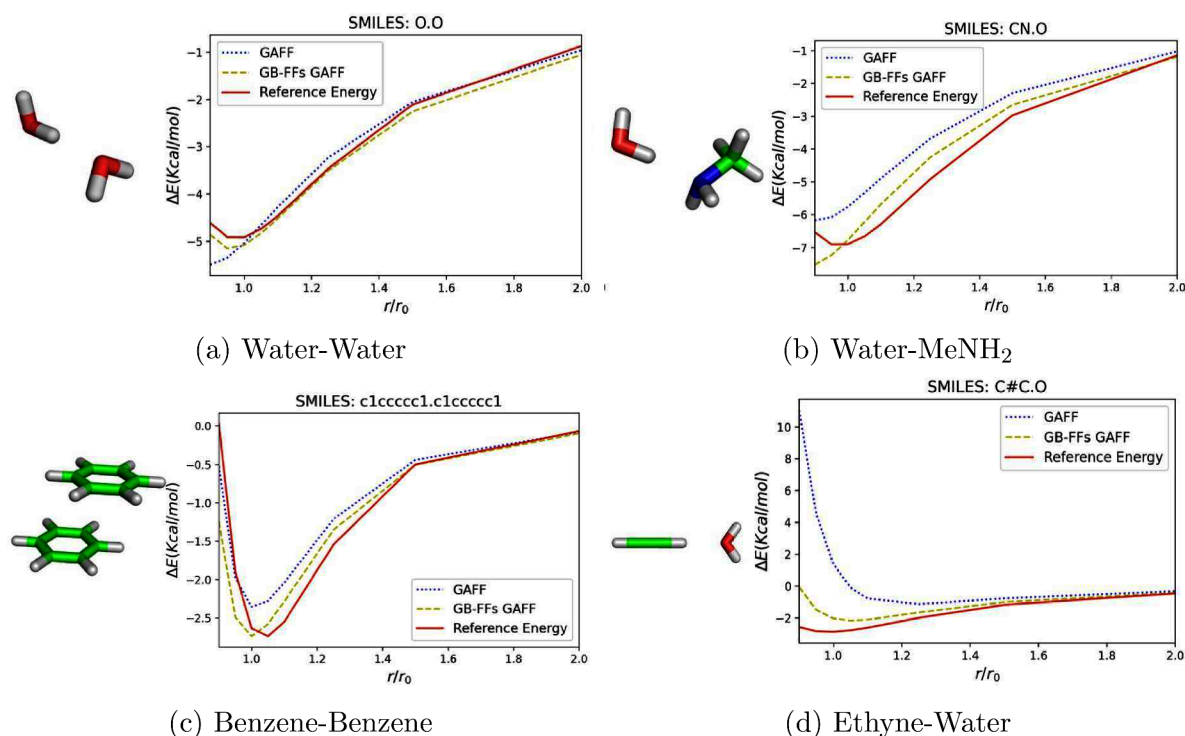


Figure 5. Results of four example on S66×8 database. (a) Water–Water, (b) Water–MeNH₂, (c) Benzene–Benzene, and (d) Ethyne–Water.

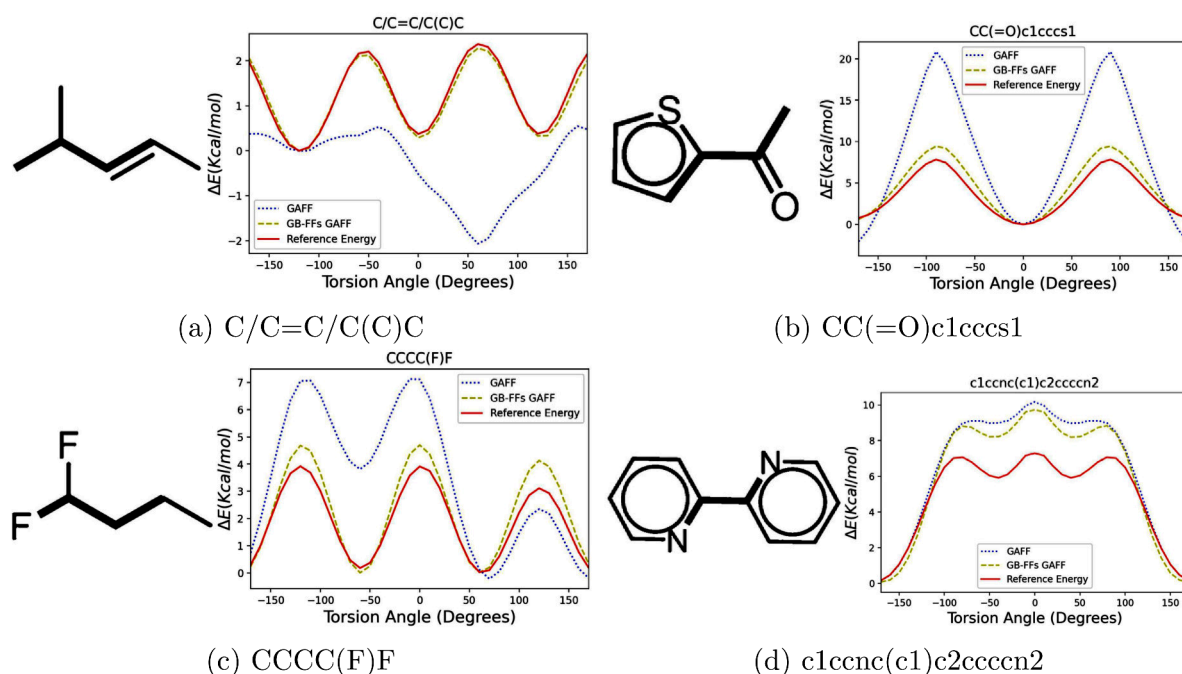


Figure 6. Results of four examples on 1D torsion scan database. (a) C/C=C/C(C)C, (b) CC(=O)c1cccs1, (c) CCCC(F)F, and (d) c1ccnc(c1)c2cccn2.

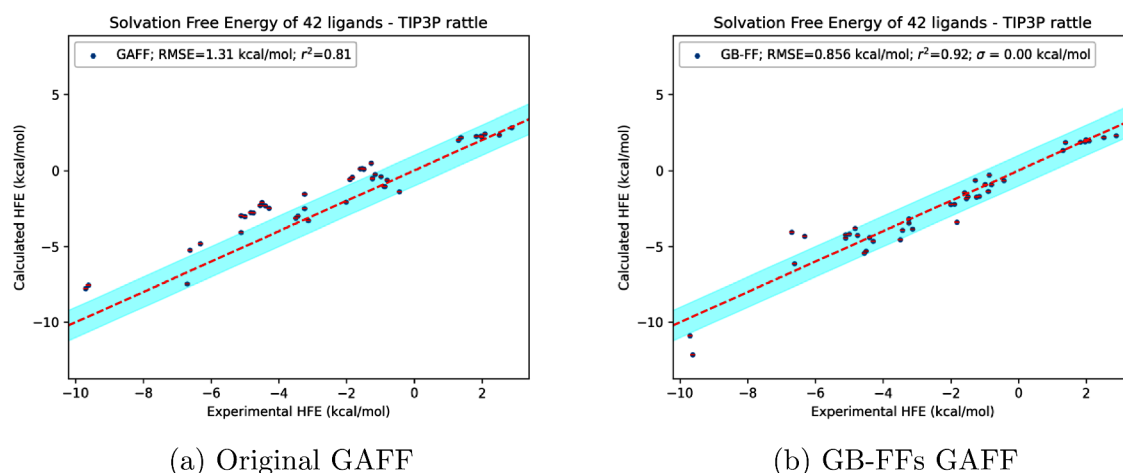


Figure 7. Comparison of hydration free energies: Computed models versus experimental values for 42 molecules in TIP3P water. The computations used λ -ABF over a total of 10 ns simulation time in the NPT ensemble with Berendsen barostat and the velocity Verlet integrator, using a 1 fs time step. (a) Hydration free energy calculated by original GAFF. (b) Hydration free energy calculated by GB-FFs GAFF.

Hydration Free Energies. While our previous focus was on potential energy and atomic forces, this subsection tackles a more challenging property: hydration free energy.

We computed the hydration free energies for a set of 42 small molecules,³⁷ previously employed to evaluate AMOEBA's performance and the recently introduced ANI/AMOEBA model.³² This set covers common chemistry examples, including benzene, acetic acid, and ethane. The experimental values are sourced from the Guthrie solvation database.⁸⁰

In our experiments, the FF parameters for the solvated molecules are taken from either the original GAFF or GB-FFs models, while the surrounding water molecules are modeled using the TIP3P water model,⁸¹ bonds and angles are constrained using the rattle algorithm.⁸² The calculations were performed using the newly introduced λ -ABF method

within Tinker-HP, over a total of 10 ns of simulation time in the NPT ensemble with the Berendsen Barostat and the velocity Verlet integrator.⁸³ Based on the original paper describing the method and our tests (see Figure S6 in the SI), for most systems, 10 ns of simulation time is sufficient to reach an error below 0.2 kcal/mol.

The overall performance is depicted in Figure 7, and the comprehensive data details are documented in Tables S5 and S6.

The original GAFF resulted in an RMSE of 1.310 kcal/mol and an r^2 value of 0.81 (see Figure 7a). In comparison, our best model achieved an RMSE of 0.856 kcal/mol and an r^2 of 0.92. Across all metrics, our model demonstrated superior performance, notably by dividing the RMSE by a factor of almost 2.

CONCLUSION

In summary, the results that have been presented here have demonstrated the efficiency and viability of using Directed Graph Attention neTworks (D-GATs) to predict parameters for the General AMBER Force Field (GAFF). The proposed parametrization approach offered several advantages. First, thanks to its low computational cost and $O(N)$ complexity, GB-FFs can assign parameters of molecules encompassing hundreds of atoms, within a few hundredths of a second. This enables the simultaneous and self-consistent parametrization of small molecules and biomolecular systems, eliminating the need for multiple distinct methodologies. Additionally, the automated workflow can leverage large databases, which could further enhance the development of FFs.

Thus we evaluated the accuracy and performance of the FF parameters from the GB-FFs model through extensive assessments on different databases.

First, we fine-tuned GB-FFs models on the SPICE and DES370K databases, which explore a wide range of chemical space. The resulting Root Mean Square Error (RMSE) of the model for SPICE energies is 3.06 kcal/mol, improving that of the original GAFF, which exhibits a 6.03 kcal/mol error. The GB-FFs model also improves the GAFF performance on the DES370K database.

To further assess the precision of our model in capturing intermolecular interactions, such as vdW's and Coulomb's, we tested its accuracy on the S66 \times 8 database. Our results showcase a reduction in RMSE by nearly half compared to the original GAFF, highlighting the improved accuracy of our approach in modeling intermolecular interactions. Moreover, we have evaluated the model's transferability and accuracy in capturing torsional interactions by computing one-dimensional torsion profiles. The GB-FFs parametrization exhibits excellent performance in capturing torsional properties (energy RMSE from 3.53 to 1.34 kcal/mol). Lastly, we examined the model's capability in predicting hydration free energies for various systems. Our model achieved lower RMSE errors 0.72 kcal/mol compared to the original GAFF parametrization 1.31 kcal/mol, showcasing its high transferability to chemically relevant systems.

The open-source code (MIT license) is available on GitHub at <https://github.com/GongCHEN-1995/GB-FFs-Model> and can be used on multiple GPUs, enabling accelerated calculations and ensuring efficient processing of multiple molecules simultaneously. Its flexibility allows for its easy integration into popular molecular dynamics workflows.

While in this article our focus was specifically on optimizing GAFF, it is important to note that our model can be extended to other nonpolarizable FFs without the need for scheme modifications, including with alternate functional forms that leads to improvements on the evaluations of energies and forces as demonstrated above. However, extending it to polarizable FFs presents a more complex challenge that demands further research, given the more complex nature of the associated parameters.

Another improvement could be done regarding the assignment of partial charges to refine the charge transfer model, either by directly accounting for polarization effects or through the electronegativity equalization approach proposed by Gilson et al.⁸⁴

Additionally, a way to enhance the accuracy of the model in simulating condensed-phase systems would be to add addi-

tional quantities, such as binding free energies into the training process, thus helping the model in capturing complex molecular behavior in condensed phase.

Finally, another avenue for improvement involves integrating polarization effects (references: 18, 20–22) and/or neural network components (references: 32, 34) into their functional form, along with accounting for nuclear quantum effects within their dynamics (references: 85, 86).

ASSOCIATED CONTENT

Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/acs.jctc.3c01421>.

Strategy to recover atom types in GAFF; details on the pretraining; details on the computations resources requirements and on the accuracy of predicted forces; full results on the S66 \times 8 and on the torsion scan databases as full hydration free energies evaluations; energy results for alternative GAFF functional forms (PDF)

AUTHOR INFORMATION

Corresponding Authors

Gong Chen – Sorbonne Université, CNRS, Université Paris Cité, Laboratoire Jacques-Louis Lions (LJLL), UMR 7598 CNRS, 75005 Paris, France; orcid.org/0000-0002-2632-3890; Email: gong.chen@sorbonne-universite.fr

Yvon Maday – Sorbonne Université, CNRS, Université Paris Cité, Laboratoire Jacques-Louis Lions (LJLL), UMR 7598 CNRS, 75005 Paris, France; orcid.org/0000-0002-0443-6544; Email: yvon.maday@sorbonne-universite.fr

Authors

Théo Jaffrelo Inizan – Sorbonne Université, Laboratoire de Chimie Théorique (LCT), UMR 7616 CNRS, 75005 Paris, France; Present Address: University of California Berkeley, Bakar Institute of Digital Materials for the Planet, College of Computing, Data Science, and Society, Berkeley 94720, USA

Thomas Plé – Sorbonne Université, Laboratoire de Chimie Théorique (LCT), UMR 7616 CNRS, 75005 Paris, France

Louis Lagardère – Sorbonne Université, Laboratoire de Chimie Théorique (LCT), UMR 7616 CNRS, 75005 Paris, France; orcid.org/0000-0002-7251-0910

Jean-Philip Piquemal – Sorbonne Université, Laboratoire de Chimie Théorique (LCT), UMR 7616 CNRS, 75005 Paris, France; orcid.org/0000-0001-6615-9426

Complete contact information is available at: <https://pubs.acs.org/doi/10.1021/acs.jctc.3c01421>

Notes

The authors declare no competing financial interest.

ACKNOWLEDGMENTS

This work has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation program (grant agreement No 810367), project EMC2 (JPP,YM). Simulations have been performed at GENCI on the Jean Zay machine (IDRIS, Orsay, France) on grant no A0070707671.

REFERENCES

- (1) Riplinger, C.; Sandhoefer, B.; Hansen, A.; Neese, F. Natural triple excitations in local coupled cluster calculations with pair natural orbitals. *J. Chem. Phys.* **2013**, *139*, 134101.
- (2) Liakos, D. G.; Guo, Y.; Neese, F. Comprehensive Benchmark Results for the Domain Based Local Pair Natural Orbital Coupled Cluster Method (DLPNO-CCSD(T)) for Closed- and Open-Shell Systems. *J. Phys. Chem. A* **2020**, *124*, 90–100.
- (3) Pearlman, D. A.; Case, D. A.; Caldwell, J. W.; Ross, W. S.; Cheatham III, T. E.; DeBolt, S.; Ferguson, D.; Seibel, G.; Kollman, P. AMBER, a package of computer programs for applying molecular mechanics, normal mode analysis, molecular dynamics and free energy calculations to simulate the structural and energetic properties of molecules. *Comput. Phys. Commun.* **1995**, *91*, 1–41.
- (4) Case, D. A.; Cheatham III, T. E.; Darden, T.; Gohlke, H.; Luo, R.; Merz Jr, K. M.; Onufriev, A.; Simmerling, C.; Wang, B.; Woods, R. J. The Amber biomolecular simulation programs. *J. Comput. Chem.* **2005**, *26*, 1668–1688.
- (5) Brooks, B. R.; Brucoleri, R. E.; Olafson, B. D.; States, D. J.; Swaminathan, S. a.; Karplus, M. CHARMM: a program for macromolecular energy, minimization, and dynamics calculations. *J. Comput. Chem.* **1983**, *4*, 187–217.
- (6) Brooks, B. R.; Brooks III, C. L.; Mackerell Jr, A. D.; Nilsson, L.; Petrella, R. J.; Roux, B.; Won, Y.; Archontis, G.; Bartels, C.; Boresch, S.; et al. CHARMM: the biomolecular simulation program. *J. Comput. Chem.* **2009**, *30*, 1545–1614.
- (7) Wang, J.; Wolf, R. M.; Caldwell, J. W.; Kollman, P. A.; Case, D. A. Development and testing of a general amber force field. *J. Comput. Chem.* **2004**, *25*, 1157–1174.
- (8) Wang, J.; Wang, W.; Kollman, P. A.; Case, D. A. Automatic atom type and bond type perception in molecular mechanical calculations. *Journal of Molecular Graphics and Modelling* **2006**, *25*, 247–260.
- (9) Stevens, J. A.; Grünewald, F.; van Tilburg, P. A. M.; König, M.; Gilbert, B. R.; Brier, T. A.; Thornburg, Z. R.; Luthey-Schulten, Z.; Marrink, S. J. Molecular dynamics simulation of an entire cell. *Frontiers in Chemistry* **2023**, *11*, 1106495.
- (10) Phillips, J. C.; Hardy, D. J.; Maia, J. D.; Stone, J. E.; Ribeiro, J. V.; Bernardi, R. C.; Buch, R.; Fiorin, G.; Hénin, J.; Jiang, W.; et al. Scalable molecular dynamics on CPU and GPU architectures with NAMD. *J. Chem. Phys.* **2020**, *153*, 044130.
- (11) Jaffrelot Inizan, T.; Célerse, F.; Adjoua, O.; El Ahdab, D.; Jolly, L.-H.; Liu, C.; Ren, P.; Montes, M.; Lagarde, N.; Lagardère, L.; et al. High-resolution mining of the SARS-CoV-2 main protease conformational space: supercomputer-driven unsupervised adaptive sampling. *Chemical Science* **2021**, *12*, 4889–4907.
- (12) El Ahdab, D.; Lagardère, L.; Inizan, T. J.; Célerse, F.; Liu, C.; Adjoua, O.; Jolly, L.-H.; Gresh, N.; Hobaika, Z.; Ren, P.; et al. Interfacial Water Many-Body Effects Drive Structural Dynamics and Allosteric Interactions in SARS-CoV-2 Main Protease Dimerization Interface. *J. Phys. Chem. Lett.* **2021**, *12*, 6218–6226.
- (13) Ponder, J. W.; Wu, C.; Ren, P.; Pande, V. S.; Chodera, J. D.; Schnieders, M. J.; Haque, I.; Mobley, D. L.; Lambrecht, D. S.; DiStasio Jr, R. A.; et al. Current status of the AMOEBA polarizable force field. *J. Phys. Chem. B* **2010**, *114*, 2549–2564.
- (14) Zhang, C.; Lu, C.; Jing, Z.; Wu, C.; Piquemal, J.-P.; Ponder, J. W.; Ren, P. AMOEBA Polarizable Atomic Multipole Force Field for Nucleic Acids. *J. Chem. Theory Comput.* **2018**, *14*, 2084–2108.
- (15) Liu, C.; Piquemal, J.-P.; Ren, P. AMOEBA+ Classical Potential for Modeling Molecular Interactions. *J. Chem. Theory Comput.* **2019**, *15*, 4122–4139.
- (16) Liu, C.; Piquemal, J.-P.; Ren, P. Implementation of Geometry-Dependent Charge Flux into the Polarizable AMOEBA+ Potential. *J. Phys. Chem. Lett.* **2020**, *11*, 419–426.
- (17) Lemkul, J. A.; Huang, J.; Roux, B.; MacKerell, A. D. J. An Empirical Polarizable Force Field Based on the Classical Drude Oscillator Model: Development History and Recent Applications. *Chem. Rev.* **2016**, *116*, 4983–5013.
- (18) Gresh, N.; Cisneros, G. A.; Darden, T. A.; Piquemal, J.-P. Anisotropic Polarizable Molecular Mechanics Studies of Inter- and Intramolecular Interactions and Ligand- Macromolecule Complexes. A Bottom-up Strategy. *J. Chem. Theory Comput.* **2007**, *3*, 1960–1986.
- (19) Naseem-Khan, S.; Lagardère, L.; Narth, C.; Cisneros, G. A.; Ren, P.; Gresh, N.; Piquemal, J.-P. Development of the Quantum-Inspired SIBFA Many-Body Polarizable Force Field: Enabling Condensed-Phase Molecular Dynamics Simulations. *J. Chem. Theory Comput.* **2022**, *18*, 3607–3621.
- (20) Jing, Z.; Liu, C.; Cheng, S. Y.; Qi, R.; Walker, B. D.; Piquemal, J.-P.; Ren, P. Polarizable Force Fields for Biomolecular Simulations: Recent Advances and Applications. *Annual Review of Biophysics* **2019**, *48*, 371–394.
- (21) Shi, Y.; Ren, P.; Schnieders, M.; Piquemal, J.-P. Polarizable force fields for biomolecular modeling. *Reviews in Computational Chemistry* **2015**, *28*, 51–86.
- (22) Melcr, J.; Piquemal, J.-P. Accurate Biomolecular Simulations Account for Electronic Polarization. *Frontiers in Molecular Biosciences* **2019**, *6*, 143.
- (23) El Khoury, L.; Célerse, F. a.; Lagardère, L.; Jolly, L.-H.; Derat, E.; Hobaika, Z.; Maroun, R. G.; Ren, P.; Bouaziz, S.; Gresh, N.; et al. Reconciling NMR Structures of the HIV-1 Nucleocapsid Protein NCp7 Using Extensive Polarizable Force Field Free-Energy Simulations. *J. Chem. Theory Comput.* **2020**, *16*, 2013–2020.
- (24) Behler, J. Four generations of high-dimensional neural network potentials. *Chem. Rev.* **2021**, *121*, 10037–10072.
- (25) Unke, O. T.; Chmiela, S.; Sauceda, H. E.; Gastegger, M.; Poltavsky, I.; Schütt, K. T.; Tkatchenko, A.; Müller, K.-R. Machine learning force fields. *Chem. Rev.* **2021**, *121*, 10142–10186.
- (26) Kocer, E.; Ko, T. W.; Behler, J. Neural network potentials: A concise overview of methods. *Annu. Rev. Phys. Chem.* **2022**, *73*, 163–186.
- (27) Artrith, N.; Behler, J. High-dimensional neural network potentials for metal surfaces: A prototype study for copper. *Phys. Rev. B* **2012**, *85*, 045439.
- (28) Sun, G.; Sautet, P. Toward fast and reliable potential energy surfaces for metallic Pt clusters by hierarchical delta neural networks. *J. Chem. Theory Comput.* **2019**, *15*, 5614–5627.
- (29) Raimbault, N.; Grisafi, A.; Ceriotti, M.; Rossi, M. Using Gaussian process regression to simulate the vibrational Raman spectra of molecular crystals. *New J. Phys.* **2019**, *21*, 105001.
- (30) Sommers, G. M.; Calegari Andrade, M. F.; Zhang, L.; Wang, H.; Car, R. Raman spectrum and polarizability of liquid water from deep neural networks. *Phys. Chem. Chem. Phys.* **2020**, *22*, 10592–10602.
- (31) Poier, P. P.; Jaffrelot Inizan, T.; Adjoua, O.; Lagardère, L.; Piquemal, J.-P. Accurate Deep Learning-Aided Density-Free Strategy for Many-Body Dispersion-Corrected Density Functional Theory. *J. Phys. Chem. Lett.* **2022**, *13*, 4381–4388.
- (32) Jaffrelot Inizan, T.; Plé, T.; Adjoua, O.; Ren, P.; Gökcen, H.; Isayev, O.; Lagardère, L.; Piquemal, J.-P. Scalable hybrid deep neural networks/polarizable potentials biomolecular simulations including long-range effects. *Chemical Science* **2023**, *14*, 5438–5452.
- (33) Wang, Y.; Inizan, T. J.; Liu, C.; Piquemal, J.-P.; Ren, P. Incorporating Neural Networks into the AMOEBA Polarizable Force Field. *J. Phys. Chem. B* **2024**, *128*, 2381–2388.
- (34) Plé, T.; Lagardère, L.; Piquemal, J.-P. Force-field-enhanced neural network interactions: from local equivariant embedding to atom-in-molecule properties and long-range effects. *Chemical Science* **2023**, *14*, 12554–12569.
- (35) Illarionov, A.; Sakipov, S.; Pereyaslavets, L.; Kurnikov, I. V.; Kamath, G.; Butin, O.; Voronina, E.; Ivahnenko, I.; Leontyev, I.; Nawrocki, G.; et al. Combining Force Fields and Neural Networks for an Accurate Representation of Chemically Diverse Molecular Interactions. *J. Am. Chem. Soc.* **2023**, *145*, 23620–23629.
- (36) Wu, J. C.; Chattree, G.; Ren, P. Automation of AMOEBA polarizable force field parameterization for small molecules. *Theor. Chem. Acc.* **2012**, *131*, 1138.
- (37) Walker, B.; Liu, C.; Wait, E.; Ren, P. Automation of AMOEBA polarizable force field for small molecules: Poltype 2. *J. Comput. Chem.* **2022**, *43*, 1530–1542.

- (38) Wang, L.-P.; Chen, J.; Van Voorhis, T. Systematic parametrization of polarizable force fields from quantum chemistry data. *J. Chem. Theory Comput.* **2013**, *9*, 452–460.
- (39) Kumar, A.; Pandey, P.; Chatterjee, P.; MacKerell, A. D. J. Deep Neural Network Model to Predict the Electrostatic Parameters in the Polarizable Classical Drude Oscillator Force Field. *J. Chem. Theory Comput.* **2022**, *18*, 1711–1725.
- (40) Thürlmann, M.; Bösel, L.; Riniker, S. Regularized by Physics: Graph Neural Network Parametrized Potentials for the Description of Intermolecular Interactions. *J. Chem. Theory Comput.* **2023**, *19*, 562–579.
- (41) Wang, Y.; Fass, J.; Kaminow, B.; Herr, J. E.; Rufa, D.; Zhang, L.; Pulido, I.; Henry, M.; Macdonald, H. E. B.; Takaba, K.; et al. End-to-end differentiable construction of molecular mechanics force fields. *Chemical Science* **2022**, *13*, 12016–12033.
- (42) Lagardère, L.; Jolly, L.-H.; Lipparini, F.; Aviat, F.; Stamm, B.; Jing, Z. F.; Harger, M.; Torabifard, H.; Cisneros, G. A.; Schnieders, M. J.; et al. Tinker-HP: a massively parallel molecular dynamics package for multiscale simulations of large complex systems with advanced point dipole polarizable force fields. *Chemical Science* **2018**, *9*, 956–972.
- (43) Adjoua, O.; Lagardère, L.; Jolly, L.-H.; Durocher, A.; Very, T.; Dupays, I.; Wang, Z.; Inizan, T. J.; Célerse, F.; Ren, P.; et al. Tinker-HP: Accelerating Molecular Dynamics Simulations of Large Complex Systems with Advanced Point Dipole Polarizable Force Fields Using GPUs and Multi-GPU Systems. *J. Chem. Theory Comput.* **2021**, *17*, 2034–2053.
- (44) Van Der Spoel, D.; Lindahl, E.; Hess, B.; Groenhof, G.; Mark, A. E.; Berendsen, H. J. GROMACS: fast, flexible, and free. *J. Comput. Chem.* **2005**, *26*, 1701–1718.
- (45) Ponder, J. W. *TINKER: Software tools for molecular design*; Washington University School of Medicine, Saint Louis, MO, 2004, 3, 116.
- (46) Rackers, J. A.; Wang, Z.; Lu, C.; Laury, M. L.; Lagardère, L.; Schnieders, M. J.; Piquemal, J.-P.; Ren, P.; Ponder, J. W. Tinker 8: Software Tools for Molecular Design. *J. Chem. Theory Comput.* **2018**, *14*, 5273–5289.
- (47) Lorentz, H. A. Ueber die Anwendung des Satzes vom Virial in der kinetischen Theorie der Gase. *Annalen der Physik* **1881**, *248*, 127–136.
- (48) Berthelot, D. Sur le mélange des gaz. *Compt. Rendus* **1898**, *126*, 15.
- (49) Bayly, C. I.; Cieplak, P.; Cornell, W.; Kollman, P. A. A well-behaved electrostatic potential based method using charge restraints for deriving atomic charges: the RESP model. *J. Phys. Chem.* **1993**, *97*, 10269–10280.
- (50) Cornell, W. D.; Cieplak, P.; Bayly, C. I.; Kollman, P. A. Application of RESP charges to calculate conformational energies, hydrogen bond energies, and free energies of solvation. *J. Am. Chem. Soc.* **1993**, *115*, 9620–9631.
- (51) Jakalian, A.; Bush, B. L.; Jack, D. B.; Bayly, C. I. Fast, efficient generation of high-quality atomic charges. AM1-BCC model: I. Method. *J. Comput. Chem.* **2000**, *21*, 132–146.
- (52) Jakalian, A.; Jack, D. B.; Bayly, C. I. Fast, efficient generation of high-quality atomic charges. AM1-BCC model: II. Parameterization and validation. *J. Comput. Chem.* **2002**, *23*, 1623–1641.
- (53) Hu, W.; Liu, B.; Gomes, J.; Zitnik, M.; Liang, P.; Pande, V.; Leskovec, J. Strategies for pre-training graph neural networks. *arXiv*, May 29, 2019, 1905.12265.
- (54) Rong, Y.; Bian, Y.; Xu, T.; Xie, W.; Wei, Y.; Huang, W.; Huang, J. Self-supervised graph transformer on large-scale molecular data. *Advances in Neural Information Processing Systems* **2020**, *33*, 12559–12571.
- (55) Fang, X.; Liu, L.; Lei, J.; He, D.; Zhang, S.; Zhou, J.; Wang, F.; Wu, H.; Wang, H. Geometry-enhanced molecular representation learning for property prediction. *Nature Machine Intelligence* **2022**, *4*, 127–134.
- (56) Chen, G.; Maday, Y. Directed message passing based on attention for prediction of molecular properties. *Comput. Mater. Sci.* **2023**, *229*, 112443.
- (57) Liu, S.; Wang, H.; Liu, W.; Lasenby, J.; Guo, H.; Tang, J. Pre-training molecular graph representation with 3d geometry. *arXiv*, October 7, 2021, 2110.07728.
- (58) Yang, K.; Swanson, K.; Jin, W.; Coley, C.; Eiden, P.; Gao, H.; Guzman-Perez, A.; Hopper, T.; Kelley, B.; Mathea, M.; et al. Analyzing learned molecular representations for property prediction. *J. Chem. Inf. Model.* **2019**, *59*, 3370–3388.
- (59) Biswas, K.; Kumar, S.; Banerjee, S.; Pandey, A. K. SMU: smooth activation function for deep networks using smoothing maximum technique. *arXiv*, November 8, 2021, 2111.04682.
- (60) Goodfellow, I.; Warde-Farley, D.; Mirza, M.; Courville, A.; Bengio, Y. Maxout networks. *Proceedings of Machine Learning Research* **2013**, *28*, 1319–1327.
- (61) Landrum, G. *RDKit: Open-source cheminformatics*, 2020. <https://www.rdkit.org>.
- (62) Muller, P. Glossary of terms used in physical organic chemistry (IUPAC Recommendations 1994). *Pure Appl. Chem.* **1994**, *66*, 1077–1184.
- (63) Morse, P. M. Diatomic molecules according to the wave mechanics. II. Vibrational levels. *Phys. Rev.* **1929**, *34*, 57.
- (64) Devlin, J. P. Urey-Bradley “Nonbonded” Forces. *J. Chem. Phys.* **1963**, *39*, 2385–2385.
- (65) Smith, J. S.; Isayev, O.; Roitberg, A. E. ANI-1: an extensible neural network potential with DFT accuracy at force field computational cost. *Chemical Science* **2017**, *8*, 3192–3203.
- (66) Eastman, P.; Behara, P. K.; Dotson, D. L.; Galvelis, R.; Herr, J. E.; Horton, J. T.; Mao, Y.; Chodera, J. D.; Pritchard, B. P.; Wang, Y.; et al. SPICE, A Dataset of Drug-like Molecules and Peptides for Training Machine Learning Potentials. *Scientific Data* **2023**, *10*, 11.
- (67) Donchev, A. G.; Taube, A. G.; Decolvenaere, E.; Hargus, C.; McGibbon, R. T.; Law, K.-H.; Gregersen, B. A.; Li, J.-L.; Palmo, K.; Siva, K.; et al. Quantum chemical benchmark databases of gold-standard dimer interaction energies. *Scientific Data* **2021**, *8*, 55.
- (68) Najibi, A.; Goerigk, L. The nonlocal kernel in van der Waals density functionals as an additive correction: An extensive analysis with special emphasis on the B97M-V and ω B97M-V approaches. *J. Chem. Theory Comput.* **2018**, *14*, 5725–5738.
- (69) Mardirossian, N.; Head-Gordon, M. ω B97M-V: A combinatorially optimized, range-separated hybrid, meta-GGA density functional with VV10 nonlocal correlation. *J. Chem. Phys.* **2016**, *144*, 214110.
- (70) Weigend, F.; Ahlrichs, R. Balanced basis sets of split valence, triple zeta valence and quadruple zeta valence quality for H to Rn: Design and assessment of accuracy. *Phys. Chem. Chem. Phys.* **2005**, *7*, 3297–3305.
- (71) Rappoport, D.; Furche, F. Property-optimized Gaussian basis sets for molecular response calculations. *J. Chem. Phys.* **2010**, *133*, 134105.
- (72) Raghavachari, K.; Trucks, G. W.; Pople, J. A.; Head-Gordon, M. A fifth-order perturbation comparison of electron correlation theories. *Chem. Phys. Lett.* **1989**, *157*, 479–483.
- (73) Montgomery, Jr. J. A.; Frisch, M. J.; Ochterski, J. W.; Petersson, G. A. A complete basis set model chemistry. VI. Use of density functional geometries and frequencies. *J. Chem. Phys.* **1999**, *110*, 2822–2827.
- (74) Rezac, J.; Riley, K. E.; Hobza, P. S66: A well-balanced database of benchmark interaction energies relevant to biomolecular structures. *J. Chem. Theory Comput.* **2011**, *7*, 2427–2438.
- (75) Sellers, B. D.; James, N. C.; Gobbi, A. A comparison of quantum and molecular mechanical methods to estimate strain energy in druglike fragments. *J. Chem. Inf. Model.* **2017**, *57*, 1265–1275.
- (76) Frisch, M. J.; Head-Gordon, M.; Pople, J. A. A direct MP2 gradient method. *Chem. Phys. Lett.* **1990**, *166*, 275–280.
- (77) Head-Gordon, M.; Head-Gordon, T. Analytic MP2 frequencies without fifth-order storage. Theory and application to bifurcated

hydrogen bonds in the water hexamer. *Chem. Phys. Lett.* **1994**, *220*, 122–128.

(78) Hehre, W. J.; Ditchfield, R.; Pople, J. A. Self-consistent molecular orbital methods. XII. Further extensions of Gaussian-type basis sets for use in molecular orbital studies of organic molecules. *J. Chem. Phys.* **1972**, *56*, 2257–2261.

(79) McLean, A.; Chandler, G. Contracted Gaussian basis sets for molecular calculations. I. Second row atoms, $Z = 11–18$. *J. Chem. Phys.* **1980**, *72*, 5639–5648.

(80) Guthrie, J. P.; Mobley, D. L. *The Guthrie Hydration Free Energy Database of Experimental Small Molecule Hydration Free Energies*. UC Irvine, 2018. Department of Pharmaceutical Sciences, UCI. Retrieved from <https://escholarship.org/uc/item/53n2h10t>.

(81) Jorgensen, W. L.; Chandrasekhar, J.; Madura, J. D.; Impey, R. W.; Klein, M. L. Comparison of simple potential functions for simulating liquid water. *J. Chem. Phys.* **1983**, *79*, 926–935.

(82) Andersen, H. C. Rattle: A “velocity” version of the shake algorithm for molecular dynamics calculations. *J. Comput. Phys.* **1983**, *52*, 24–34.

(83) Lagardère, L.; Maurin, L.; Adjoua, O.; Hage, K. E.; Monmarché, P.; Piquemal, J.-P.; Hénin, J. Lambda-ABF: Simplified, Accurate and Cost-effective Alchemical Free Energy Computations. *J. Chem. Theory Comput.* **2024**, DOI: 10.1021/acs.jctc.3c01249.

(84) Gilson, M. K.; Gilson, H. S.; Potter, M. J. Fast assignment of accurate partial atomic charges: an electronegativity equalization method that accounts for alternate resonance forms. *Journal of chemical information and computer sciences* **2003**, *43*, 1982–1997.

(85) Mauger, N.; Plé, T.; Lagardère, L.; Bonella, S.; Mangaud, É.; Piquemal, J.-P.; Huppert, S. Nuclear quantum effects in liquid water at near classical computational cost using the adaptive quantum thermal bath. *J. Phys. Chem. Lett.* **2021**, *12*, 8285–8291.

(86) Plé, T.; Mauger, N.; Adjoua, O.; Inizan, T. J.; Lagardère, L.; Huppert, S.; Piquemal, J.-P. Routine Molecular Dynamics Simulations Including Nuclear Quantum Effects: From Force Fields to Machine Learning Potentials. *J. Chem. Theory Comput.* **2023**, *19*, 1432–1445.



CAS BIOFINDER DISCOVERY PLATFORM™

ELIMINATE DATA SILOS. FIND WHAT YOU NEED, WHEN YOU NEED IT.

A single platform for relevant, high-quality biological and toxicology research

Streamline your R&D

CAS
A division of the American Chemical Society