


 Cite this: *RSC Adv.*, 2020, 10, 666

# Atomic partial charge predictions for furanoses by random forest regression with atom type symmetry function

 Xiacong Wang  and Jun Gao \*

Furanoses that are components for many important biomolecules have complicated conformational spaces due to the flexible ring and *exo*-cyclic moieties. Machine learning algorithms, which require descriptors as structural inputs, can be used to efficiently compute conformational adaptive (CA) charges to capture the electrostatic potential variations caused by the conformational changes in the molecular mechanics (MM) calculations. In the present study, we introduced atom type symmetry function (ATSF) developed based on atom centered symmetry function (ACSF) for describing conformations for furanoses, in which atoms were categorized by atom types defined by their properties or connectivity in classic molecular mechanics (MM) force field parameters to generate a suitable coordinate size. Random forest regression (RFR) models with ATSF showed improvements for predicting CA charges and dipole moments for furanoses compared to those with ACSF and atom name symmetry functions where atoms were categorized by their unique atom names. The CA charges predicted by RFR models with ATSF showed more comparable reproductions of the carbohydrate–water and carbohydrate–protein interactions computed with RESP charges individually derived from QM calculations than the ensemble-averaged atomic charge sets commonly employed in molecular mechanics force fields, suggesting that the predicted CA charges were capable of including electrostatic variations in their dynamic charge values. Improvements by ATSF showed that categorizing atoms by atom types introduced chemical structural perceptions to descriptors and produced a suitable coordinate size in ATSF to capture key structural features for furanoses. This categorizing scheme also allows ATSF to be readily adopted by other biomolecules thanks to the broad implementations of MM force fields.

 Received 10th November 2019  
 Accepted 18th December 2019

DOI: 10.1039/c9ra09337k

[rsc.li/rsc-advances](http://rsc.li/rsc-advances)

## Introduction

Furanoses are essential components for the backbones of nucleic acids and complex polysaccharides frequently found in organisms ranging from bacteria to protozoa, fungi to plants.<sup>1</sup> They have complicated conformational spaces as their five-membered ring can adopt multiple stable conformations in addition to the spinning of their abundant *exo*-cyclic groups in solution.<sup>2</sup> These conformational variations lead to heterogeneous intramolecular properties, such as electrostatic potentials, which affect their recognitions and interactions with other biomolecules.<sup>3</sup> These variations also made it difficult for classical molecular mechanics (MM) force fields to adequately represent the electrostatic properties for furanoses, as static atomic partial charge models are commonly employed.<sup>4</sup> These models, however computationally efficient, lack the accuracy to represent the intrinsic electrostatic potential variations. Efforts have been devoted to develop charge models that are capable of adapting conformational variations.<sup>4–10</sup> Approaches have been

proposed and developed, yet, adoptions of these approaches depend on their applicability and ease of use. It is also unfeasible to derive conformational adaptive (CA) charges during MM calculations directly from electrostatic potentials obtained from resource-hogging quantum mechanics (QM) calculations.<sup>5,11,12</sup> Therefore, it is desirable to efficiently compute QM-quality CA charges that can be used within the classical MM framework.

Machine learning algorithms have been implemented in calculating electrostatic potentials and provide a promising alternative approach for computing CA charges.<sup>13–15</sup> The accuracy and efficiency of machine learning algorithms critically depend on the descriptors that are used to represent molecular structures.<sup>16,17</sup> Descriptors for machine learning algorithms, unlike cartesian coordinates, are required to be invariant under permutations of atoms, as well as translations or rotations of the molecule, so as to represent any conformation in a unique set of coordinates.<sup>16</sup> These descriptors also need to describe the key structural features of molecules with a sufficient size of coordinates. The atom centered symmetry function<sup>18</sup> (ACSF) introduced by Behler and Parrinello in 2007 has become a prominent descriptor for machine learning algorithms and many successful implementations have been reported.<sup>19–24</sup> ACSF

Hubei Key Laboratory of Agricultural Bioinformatics, College of Informatics, Huazhong Agricultural University, Wuhan, China. E-mail: [gaojun@mail.hzau.edu.cn](mailto:gaojun@mail.hzau.edu.cn)



categorizes atoms by the number of atom types in the molecule, which determines the size of its coordinates. Thus, when describing furanoses and other biomolecules that usually possess complicated conformational spaces but limited types of chemical elements, improvements may be needed.

In the present study, we introduced atom type symmetry function (ATSF), that categorized atoms by their atom types defined in MM force fields<sup>3,25–27</sup> and provided more detailed structural descriptions, to predict CA charges with properly trained random forest regression (RFR) models.<sup>28</sup> Atom type is a well-established and crucial concept embedded in common MM force fields, in which atoms are categorized beyond chemical elements and by their properties or connectivity. In the furanose-specific GLYCAM force field,<sup>3</sup> atoms for furanoses (Fig. 1) that belong to three chemical elements were further categorized into eight different atom types and the size of coordinates for ATSF increased by more than three times comparing to that for ACSF. The Pearson correlation coefficients for predicted charge values by RFR models with ATSF with reference to RESP charges derived from QM calculations were all above 0.9 and increased averagely by 9% and 4%, respectively, comparing to ACSF and atom name symmetry functions (ANSF). In addition, the predictions of dipole moments were improved by 43% and 48% by charges predicted with ATSF comparing to those with ACSF and ANSF,

respectively. Furthermore, the electrostatic related interactions for furanosides, carbohydrate–water and carbohydrate–protein interactions, computed with the CA charges predicted by RFR models with ATSF reduced the average error by more than half for that calculated with the static ensemble-averaged charge models and individual RESP charges derived from QM calculations, which indicated that the predicted CA charges was capable of including electrostatic variations in their dynamic charge values.

## Methods

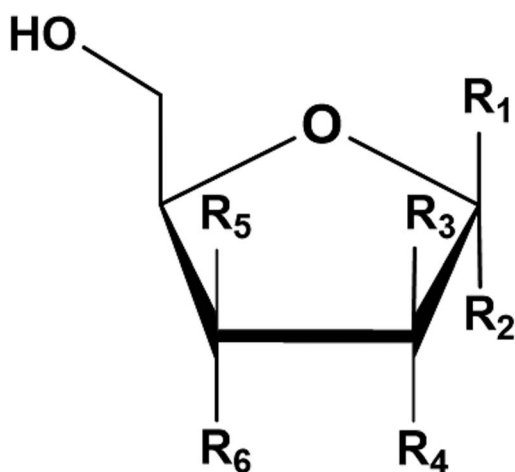
### Training and testing dataset formation

The training and testing datasets for RFR models were formed by combining the samplings for the *endo*- and *exo*-cyclic conformations for furanosides, in order to sufficiently cover their conformational spaces. The conformations for the five-membered ring are determined by the *endo*-cyclic rotations and can be described by the pseudorotational itinerary,<sup>29,30</sup> where the phase angle ( $P$ ) and puckering amplitude ( $\tau_m$ ) can be calculated by the five ring torsion angles (Fig. 2).

$$P = \tan^{-1} \frac{(\theta_2 + \theta_4) - (\theta_1 + \theta_3)}{2\theta_0(\sin 36^\circ + \sin 72^\circ)} \quad (1)$$

$$\tau_m = \frac{\theta_0}{\cos P} \quad (2)$$

In order to thoroughly sample the conformational spaces for the furanose ring,  $\tau_m$  was iterated from 3 to 45° at a 3° interval, and  $P$  was from 0 to 360° at a 6° interval. A total of 900 ring conformations were generated for each furanoside or furanose. Unlike the intertwining *endo*-cyclic torsions, all *exo*-cyclic



1 $\alpha$ : R<sub>1</sub> = H, R<sub>2</sub> = OCH<sub>3</sub>, R<sub>3</sub> = OH, R<sub>4</sub> = H, R<sub>5</sub> = H, R<sub>6</sub> = OH

1 $\beta$ : R<sub>1</sub> = OCH<sub>3</sub>, R<sub>2</sub> = H, R<sub>3</sub> = OH, R<sub>4</sub> = H, R<sub>5</sub> = H, R<sub>6</sub> = OH

2 $\alpha$ : R<sub>1</sub> = H, R<sub>2</sub> = OCH<sub>3</sub>, R<sub>3</sub> = OH, R<sub>4</sub> = H, R<sub>5</sub> = OH, R<sub>6</sub> = H

2 $\beta$ : R<sub>1</sub> = OCH<sub>3</sub>, R<sub>2</sub> = H, R<sub>3</sub> = OH, R<sub>4</sub> = H, R<sub>5</sub> = OH, R<sub>6</sub> = H

3 $\alpha$ : R<sub>1</sub> = H, R<sub>2</sub> = OCH<sub>3</sub>, R<sub>3</sub> = H, R<sub>4</sub> = OH, R<sub>5</sub> = H, R<sub>6</sub> = OH

3 $\beta$ : R<sub>1</sub> = OCH<sub>3</sub>, R<sub>2</sub> = H, R<sub>3</sub> = H, R<sub>4</sub> = OH, R<sub>5</sub> = H, R<sub>6</sub> = OH

4 $\alpha$ : R<sub>1</sub> = H, R<sub>2</sub> = OCH<sub>3</sub>, R<sub>3</sub> = H, R<sub>4</sub> = OH, R<sub>5</sub> = OH, R<sub>6</sub> = H

4 $\beta$ : R<sub>1</sub> = OCH<sub>3</sub>, R<sub>2</sub> = H, R<sub>3</sub> = H, R<sub>4</sub> = OH, R<sub>5</sub> = OH, R<sub>6</sub> = H

Fig. 1 Methyl furanosides for the present study: methyl *D*-arabino-furanoside (1), methyl *D*-lyxofuranoside (2), methyl *D*-ribofuranoside (3), and methyl *D*-xylofuranoside (4).

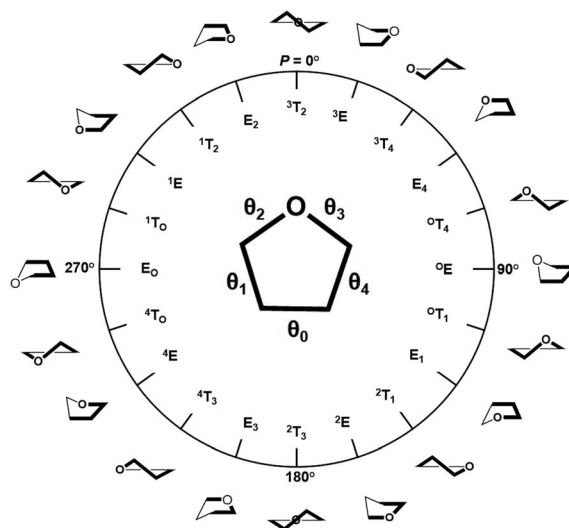


Fig. 2 Pseudorotational itinerary of furanoses depicting different Envelope (E) and Twist (T) ring conformations with associated conformational phase angle ( $P$ ) values in degrees. The inset demonstrated the definitions of five ring torsion angles:  $\theta_0 = \text{C1-C2-C3-C4}$ ,  $\theta_1 = \text{C2-C3-C4-O4}$ ,  $\theta_2 = \text{C3-C4-O4-C1}$ ,  $\theta_3 = \text{C4-O4-C1-C2}$ , and  $\theta_4 = \text{O4-C1-C2-C3}$ .



rotations were independent from each other, so it is hardly to systematically cover *exo*-cyclic conformational spaces. So, 30 combinations of randomly assigned values for all *exo*-cyclic dihedral angles were created for each ring conformation, in which 24 structures were randomly selected to construct the training dataset, and the rest were selected to construct the testing dataset.

### Quantum mechanics (QM) calculations

All of the QM calculations were performed under the same protocol as that in the development of furanose-specific GLYCAM force field (ref) with the Gaussian 16 software package<sup>31</sup> to maintain equal comparisons. Structural optimizations were performed at the HF/6-31G\* level of theory, with five ring torsion angles and all *exo*-cyclic torsion angles restrained (Fig. 1). The electrostatic potentials were calculated on these optimized structures at the B3LYP/cc-pVTZ level of theory. The atomic partial charges were derived by employing the restrained electrostatic potential (RESP) charge fitting methodology with a weak hyperbolic charge restraint weight of 0.0005.<sup>3,32</sup> The charge values for aliphatic hydrogen atoms were assigned to 0 by following GLYCAM force field parameter development philosophy.<sup>19</sup>

### Atom type symmetry function

ATSF was constructed under the framework of ACSF, whose coordinates were calculated from cartesian coordinates of atoms.<sup>18</sup> The cutoff function in ATSF was that same as that in ACSF:

$$f_c(R_{ij}) = \begin{cases} 0.5 \times \left[ \cos\left(\frac{\pi R_{ij}}{R_c}\right) + 1 \right] & \text{for } R_{ij} \leq R_c \\ 0 & \text{for } R_{ij} \geq R_c \end{cases} \quad (3)$$

with  $R_c$  set to 99 Å to include all atoms for molecules included in this study.

Radial components of atom  $i$  were calculated *via* a sum of Gaussians,

$$G_{i,J}^{\text{radial}} = \sum_{j \neq i}^{j \text{ in } J} e^{-\eta(R_{ij}-R_s)^2} f_c(R_{ij}) \quad (4)$$

in which  $R_s$  and  $\eta$  are both set to 1.0. Atom  $j$  is an atom in atom types  $J$ . Please note that the atom  $i$  could be in atom type  $J$ . The assembly of  $G_i^{\text{radial}}$  with different atom types constructs the radial components in ATSF for atom  $i$ .

Angular components for atom  $i$  were constructed as:

$$G_{i,J,K}^{\text{angular}} = 2^{1-\zeta} \sum_{j,k \neq i}^{j \text{ in } J \ \& \ k \text{ in } K} (1 + \lambda \cos \theta_{ijk})^\zeta \times e^{-\eta(R_{ij}^2 + R_{ik}^2 + R_{jk}^2)} f_c(R_{ij}) f_c(R_{ik}) f_c(R_{jk}) \quad (5)$$

with the parameters of  $\lambda = 1.0$ ,  $\zeta = 1.0$ . Atom  $j$  and  $k$  are atoms in atom types  $J$  and  $K$ , respectively. The atom  $i$  could also be in atom types  $J$  and  $K$ . The assembly of  $G_i^{\text{angular}}$  with different combinations of atom types constructs the angular components

in ATSF for atom  $i$ . For each atom  $i$  in the molecule, the ATSF coordinates can be assembled as:

$$X_i = \left\{ G_{i,AT1}^{\text{radial}}, \dots, G_{i,ATn}^{\text{radial}}, G_{i,AT1,AT1}^{\text{angular}}, \dots, G_{i,AT1,ATn}^{\text{angular}}, G_{i,AT2,AT2}^{\text{angular}}, \dots, G_{i,AT2,ATn}^{\text{angular}}, \dots, G_{i,ATn,ATn}^{\text{angular}} \right\} \quad (6)$$

where AT stands for atom type.

The training and testing sets were generated by  $S = \{(X_1, y_1), (X_2, y_2), \dots, (X_n, y_n)\}$ , where  $y_1, y_2, \dots, y_n$  are the RESP derived charges for the corresponding atoms.

The charge predictions were achieved *via* multiple random forest regression (RFR) models, the sum of the predicted charges for an individual molecule is not necessarily 0. The corrections were achieved by spreading the discrepancy based on the standard derivations for RFR predictions. This procedure was adopted from ref. 33.

### Random forest regression

RFR model was trained for atoms in each element using the scikit-learn library (version 0.18.1)<sup>34</sup> with the following parameters: number of trees = 200, maximum depth = 6, minimum number of samples to split = 6, and minimum number of samples in leaves = 6.

### Molecular dynamics (MD) simulations

The initial coordinates for furanosides 1–4 ( $\alpha$  and  $\beta$ ) were obtained from GLYCAM website (<http://www.glycam.org>). All systems were solvated with TIP3P water<sup>35</sup> using a 12 Å buffer in a cubic box, using the LEaP module in the AMBER16 software package.<sup>36</sup> Force field valence parameters were taken from furanose-specific parameters in GLYCAM.<sup>3</sup> The energy minimizations for these solvated furanoses were performed separately under NVT condition (500 steps steepest descent, followed by 24 500 steps of conjugate-gradient minimization). Subsequently, each system was heated to 300 K over a period of 50 ps, followed by equilibration at 300 K for a further 0.5 ns using NPT condition, with the Berendsen thermostat and barostat<sup>37</sup> for temperature and pressure control, respectively. SHAKE algorithm<sup>38</sup> was employed to constrain all covalent bonds involving hydrogen atoms, allowing a simulation time step of 2 fs throughout the simulations. After the equilibration, production simulations were carried out with the GPU implementation<sup>39</sup> of the PMEMD.MPI module and trajectory frames collected at every 1 ps from the total of 300 ns. A non-bonded cut-off of 8 Å was applied to van der Waals interactions, with long-range electrostatics treated with the particle mesh Ewald approximation.

### Hydration free energy and protein–carbohydrate interaction energy calculations

Hydration free energies for 1–4 ( $\alpha$  and  $\beta$ ) and protein–carbohydrate interaction energies for 3 ATP-binding cassette (ABC) transporters with furanoses as ligands (PDB ID: 2VK2,  $\alpha$ -D-Galf-OH and  $\alpha$ -D-Galf-OH as the ligands; PDB ID: 3KSM,  $\beta$ -D-Ribf-OH as the ligand) were calculated with molecular mechanics-



generalized born surface area (MM-GBSA) method using single trajectory approach. Monosaccharides in all systems were taken as the ligand in the calculation. Solvent molecules and proteins were taken as the receptors in hydration free energy and protein-carbohydrate interaction calculations, respectively. Each MM-GBSA calculation was performed on the 10 000 evenly extracted structures of each solvated system with 3 different charge sets: RESP charges individually derived from QM calculations, predicted charges by RFR models with ATSF, and the ensemble-averaged atomic charges from GLYCAM force field.

## Results and discussion

### Different atom categorizing schemes in symmetry functions

ATSF in the present study employed atom types defined in furanose-specific GLYCAM force field<sup>3</sup> (Fig. 3) to divide atoms into more groups beyond chemical elements. When categorizing atoms by only three chemical elements, ACSF contains 9 coordinates. The size of ATSF increased to 41 coordinates when categorizing atoms by a total of eight different atom types. The descriptor needs to be sufficiently large to ensure an unambiguous distinction of different conformations.<sup>16</sup> In addition, categorizing atoms by their atom types introduced structural perceptions to the descriptor by adding information of property or connectivity for atoms. For complete comparisons, the most subtle categorizing scheme was also employed to generate atom name symmetry function (ANSF), where atoms were divided by their atom names (Fig. 3) and each single atom was in a unique category (Fig. 3). The size of ANSF for furanoside was dramatically increased to 276 coordinates and undoubtedly made this scheme unpractical for efficient calculations. Remarkably, ANSF abolished all chemical or structural information from the descriptor.

### Performances for atom type symmetry function

The performances of ATSF, ACSF, and ANSF were, firstly, evaluated individually by comparing the predicted charge values to the corresponding expected values, aka RESP values.<sup>32</sup>

Data in panel C of Fig. 4 appeared to be less scattered than those in panel B and C, which suggested a better correlation achieved by ATSF than ACSF and ANSF. The Pearson correlation coefficients for the CA charges predicted with ATSF reference to RESP charges derived from QM calculations increased averagely by 9% and 4%, respectively, comparing to ACSF and ANSF. The Pearson coefficients for all atoms were larger than 0.9 for those predicted with ATSF and systematically higher than the other two descriptors. This stronger correlation indicated RFR models with ATSF were able to produce more accurate CA charges than ACSF and ANSF. The results of linear fittings between predicted and RESP-fit charges were shown in Table 1. The slopes and intercepts to Y-axis for predictions with ATSF were close to 1 and 0, respectively, implied charges predicted values were not overestimated or underestimated. It is worth noting that RFR models with ANSF, which has a substantially larger size, did not produce higher quality predictions comparing to ATSF. This suggested that the chemical perceptions in categorizing scheme is crucial and increasing the size of coordinates without chemical perceptions would not necessarily guarantee more accurate predictions.

RFR models with ATSF have shown great potentials in predicting atomic partial charge values from different conformations of furanosides. It is also essential to demonstrate their advances in representing the electrostatic potentials of furanosides by reproducing molecular dipole moments,<sup>40–44</sup> which is the first order of multipole expansions of electrostatic potentials and strongly depend on the conformation of the molecule. The average absolute differences of dipole moments between QM calculated values and the corresponding MM calculated values with different charge models were shown in Table 2. The differences of dipole moments computed from atomic partial charges predicted with ATSF were 43% and 48% smaller than those computed with ACSF and ANSF, respectively. It is worth noting that RESP charge models has the lowest dipole moment differences, which suggested that this charge model is appropriate to represent the electrostatic potential variations for

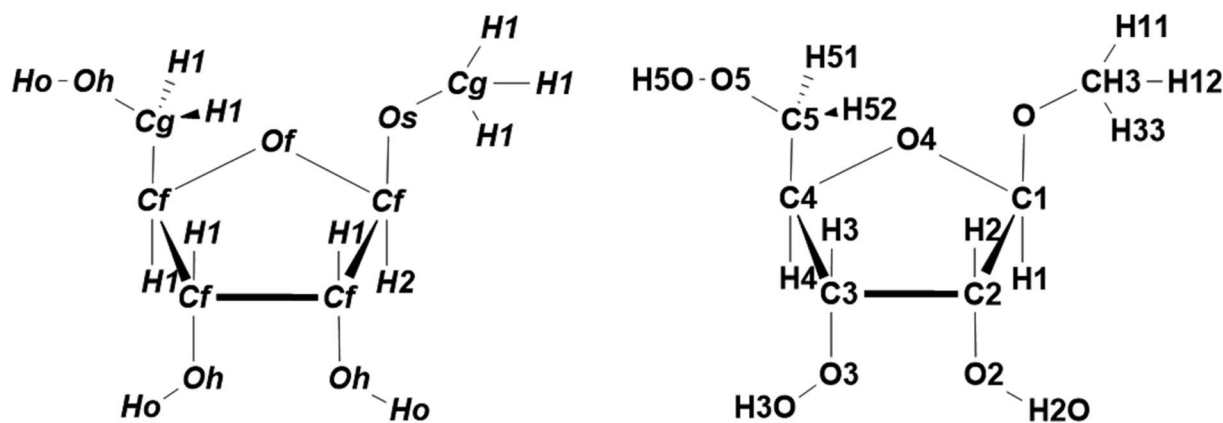


Fig. 3 Atom types (left) and atom names (right) for furanosides in GLYCAM force field. In GLYCAM force field,<sup>3</sup> "Cf" and "Cg" are for *endo*- and *exo*-cyclic carbon atoms, respectively. "Of" stands for the *endo*-cyclic oxygen atoms; "Os" and "Oh" are for *exo*-cyclic oxygen atoms in ether and hydroxyl groups, respectively. "H1" and "H2" stand for hydrogen atoms attached to a carbon atom that is bonded with one and two electron-withdrawing atoms, respectively; "Ho" is for the hydrogen atoms in hydroxyl groups.



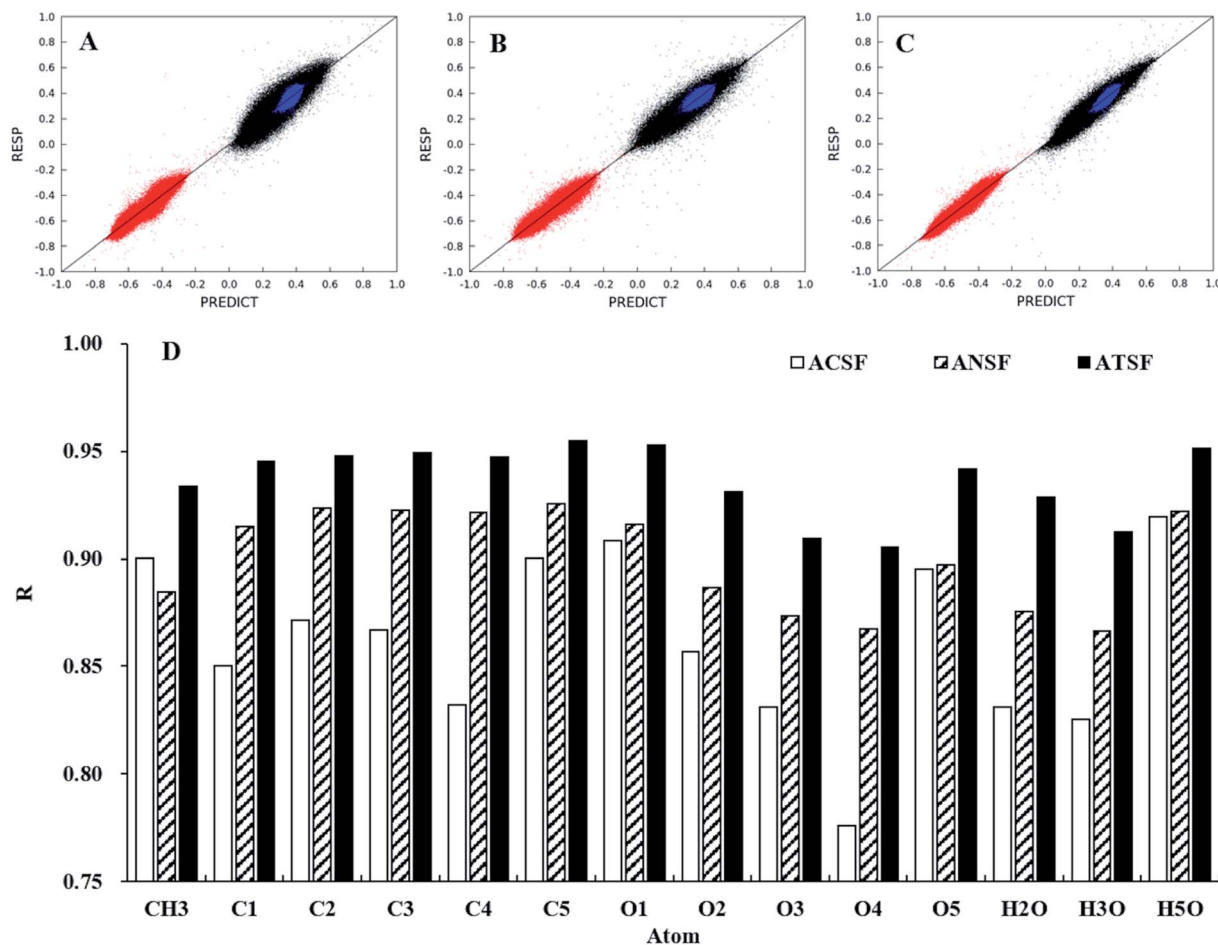


Fig. 4 Comparisons for predicted and the corresponding RESP charges in 1–4 ( $\alpha$  and  $\beta$ ). Predictions of the CA charges for carbon (black), oxygen (red), and hydrogen (blue) atoms were performed under trained RFR models with ACSF (A), ANSF (B), and ATSF (C). The predictions were evaluated with their Pearson coefficients for different atoms (D).

furanosides and be utilized as the references for RFR model training.

So far, RFR models with ATSF demonstrated their capabilities of predicting atomic partial charges with QM quality. To further confirm the validities of ATSF and its CA charges, the electrostatic-related interactions, carbohydrate–water and carbohydrate–protein interactions, for furanoses computed with predicted charges were compared to those with RESP charges.

#### Performances on carbohydrate–water interaction energy calculations

Carbohydrate–water interactions, quantified by their hydration free energies, substantially depend on their electrostatic interactions. Thus, the quality of the atomic charges can be evaluated by their computed hydration free energies.<sup>33,45</sup> In terms of conformation adaptive charges, the validity can be tested by comparing their computed hydration free energies to those computed with individual RESP-fit charges. Moreover, the hydration free energy from a single solute conformation could introduce errors.<sup>46</sup> So, the averaged hydration free energies

Table 1 Linear fits<sup>a</sup> results between atomic charges predicted with different descriptors and their corresponding RESP charges derived from QM calculations

Atom	ACSF		ANSF		ATSF	
	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>
CH3	1.038	−0.008	0.902	0.021	1.042	−0.009
C1	1.054	−0.020	0.931	0.026	1.047	−0.017
C2	1.068	−0.018	0.935	0.016	1.043	−0.011
C3	1.071	−0.017	0.930	0.017	1.055	−0.013
C4	1.077	−0.017	0.938	0.014	1.081	−0.018
C5	1.077	−0.017	0.934	0.015	1.034	−0.008
O1	1.036	0.014	0.924	−0.030	1.040	0.016
O2	1.063	0.038	0.907	−0.057	1.066	0.040
O3	1.096	0.059	0.886	−0.070	1.087	0.053
O4	1.151	0.065	0.886	−0.049	1.119	0.052
O5	1.051	0.030	0.911	−0.053	1.053	0.031
H2O	1.081	−0.030	0.900	0.038	1.062	−0.023
H3O	1.075	−0.027	0.882	0.045	1.087	−0.033
H5O	1.051	−0.019	0.927	0.027	1.035	−0.013

<sup>a</sup> Linear fit was achieved by  $y = a \times x + b$ .



**Table 2** Average absolute differences between QM calculated dipole moments and the corresponding MM calculated values with different charge models for all conformations of furanosides in the testing and training data set

Charges	RFR models with			
	ACSF	ANSF	ATSF	RESP
$\langle  \text{Dipole difference}  \rangle^a$	$0.28 \pm 0.23$	$0.31 \pm 0.32$	$0.16 \pm 0.16$	$0.04 \pm 0.03$

<sup>a</sup> In Debye (D).

computed with conformational adaptive, RESP-fit, and ensemble-averaged charge sets for 100 000 structures of each furanoside in 1–4 ( $\alpha$  and  $\beta$ ) extracted from explicit solvent MD simulations were compared (Table 3). The hydration free energies calculated with CA charges predicted with ATSF are comparable to those computed with RESP-fit charges. The difference is  $0.4 \text{ kcal mol}^{-1}$ , which is significantly less than that computed from ensemble-averaged charge sets ( $1.0 \text{ kcal mol}^{-1}$ ). The differences among these calculated hydration free energies are significant ( $p$ -value  $< 0.0001$ ), because of the large amount of structures employed in hydration free energy calculations, although the standard deviations are mostly over  $3.0 \text{ kcal mol}^{-1}$ .

It is worth noting that the hydration free energies do not include the entropic penalties, therefore, the values may be more negative than the experimental measured values.

**Table 3** Hydration free energies for 1–4 (both  $\alpha$  and  $\beta$ ) computed with different atomic partial charge sets

	Ensemble-averaged		RFR model with ATSF		QM derived RESP	
	Average	Stdev	Average	Stdev	Average	Stdev
1 $\alpha$	-6.8	3.2	-8.2	3.2	-7.8	3.4
1 $\beta$	-6.9	3.2	-7.4	3.2	-8.0	3.3
2 $\alpha$	-6.8	3.1	-7.9	3.1	-7.8	3.2
2 $\beta$	-6.5	3.0	-8.2	2.9	-7.6	3.1
3 $\alpha$	-6.7	3.1	-7.0	3.1	-7.9	3.3
3 $\beta$	-6.7	3.1	-7.2	3.1	-7.4	3.2
4 $\alpha$	-7.1	3.2	-8.5	3.2	-8.3	3.3
4 $\beta$	-6.9	3.2	-7.3	3.1	-7.9	3.4
$\langle  \text{Difference}  \rangle$	1.0		0.4			

**Table 4** MM-GBSA energies for three ABC proteins computed with different atomic partial charge sets

Ligands	QM derived RESP		RFR model with ATSF		Ensemble-averaged	
	Average	Stdev	Average	Stdev	Average	Stdev
$\alpha$ -D-Galf-OH	-31.7	1.9	-31.3	2.2	-31.9	2.3
$\beta$ -D-Galf-OH	-26.8	2.0	-27.0	2.4	-27.9	2.1
$\beta$ -D-Ribf-OH	-19.9	2.1	-19.6	2.3	-21.2	2.4
$\langle  \text{Error}  \rangle$			0.3		0.9	

### Performance on carbohydrate–protein interaction energy calculations

Hydrogen bonding interaction is one of most popular and crucial hydrophilic interactions between carbohydrate molecules and proteins,<sup>47–49</sup> due to the richness of hydroxyl groups presence in the *exo*-cyclic moieties. Accurately representing electrostatic potentials of carbohydrate molecules is crucial for correctly modeling the strength of hydrogen bonds between carbohydrate molecules and proteins. Yet, the static charge model lacks the accuracy for representing the electrostatic variations due to the changes from both *endo*- and *exo*-cyclic conformations of furanoses while interacting with proteins. Thus, the CA charge sets could improve the accuracy for carbohydrate–protein interaction energy calculations.

The computed MM-GBSA energies for three ABC transporter complexes with furanoses as ligands were listed in Table 4. The  $\langle |\text{error}| \rangle$  for MM-GBSA energies with CA charges predicted by RFR models with ATSF, comparing to that calculated with RESP charges derived from QM calculations, was only  $0.3 \text{ kcal mol}^{-1}$ , while that with the ensemble-averaged charges was  $0.9 \text{ kcal mol}^{-1}$ .

The CA charge set predicted by RFR models with ATSF, comparing to the static charge model, showed improvements for including the electrostatic potential variations seen by individually derived charge values from QM calculations in their dynamic charge values. Similar to hydration free energies, these values do not include the entropic penalties, therefore, do not reflect the results measured from experiments.

## Conclusions

Atom type symmetry function (ATSF) categorized atoms by their atom types defined by the properties and connectivity of atoms in MM force field, beyond chemical elements in ACSF, and formed a more detailed structural description for furanoses that have complicated conformational spaces but limited chemical elements. Hence, the RFR models with ATSF produced more accurate predictions of CA charges and dipole moments than those with ACSF, which suggested ATSF obtained improvements in representing structural information for furanoses. The CA charge predicted by RFR models with ATSF, comparing to the static ensemble-averaged charges, employed in computing carbohydrate–water and carbohydrate–protein interaction energies showed a better agreement to those computed with individually derived RESP charges, which also demonstrated



that CA charges were able to include the electrostatic potentials variations into the dynamic charge values.

Improvements achieved by ATSF in representing structural information for furanoses suggested that introducing structural perceptions to the descriptor and increasing the size of the coordinates could improve the performance of ACSF in describing furanoses. Furthermore, ATSF outperforming ANSF that had a significant larger size of coordinates but removed all chemical or structural perceptions of atoms suggested that categorizing atoms by atom types generated a suitable size of coordinates that represented the key structural features for furanoses. Additionally, this categorizing scheme endowed ATSF with the exceeding potent transferability to other biomolecules thanks to the broad implementations of MM force fields for biomolecules.

## Conflicts of interest

There are no conflicts to declare.

## Acknowledgements

This work was supported by the National Key R&D Program of China [Grant No. 2017YFB0203405], National Natural Science Foundation of China [No. 21873034], Fundamental Research Funds for the Central Universities (Project 2662018JC027) and Huazhong Agricultural University Scientific and Technological Self-innovation Foundation [Program No. 2015RC008].

## References

- H. A. Taha, M. R. Richards and T. L. Lowary, *Chem. Rev.*, 2013, **113**, 1851.
- K. N. Kirschner and R. J. Woods, *Proc. Natl. Acad. Sci. U. S. A.*, 2001, **98**, 10541.
- X. C. Wang and R. J. Woods, *J. Biomol. NMR*, 2016, **64**, 291.
- M. Jana and A. D. Mackerell, *J. Phys. Chem. B*, 2015, **119**, 7846.
- L. K. Du, J. Gao, F. Z. Bi, L. L. Wang and C. B. Liu, *J. Comput. Chem.*, 2013, **34**, 2032.
- J. A. Lemkul, J. Huang, B. Roux and A. D. MacKerell, *Chem. Rev.*, 2016, **116**, 4983.
- P. Cieplak, F. Y. Dupradeau, Y. Duan and J. M. Wang, *J. Phys. Condens. Matter*, 2009, **21**, 21.
- T. L. Fletcher and P. L. A. Popelier, *J. Comput. Chem.*, 2017, **38**, 1005.
- A. Warshel, M. Kato and A. V. Pisiakov, *J. Chem. Theory Comput.*, 2007, **3**, 2034.
- C. M. Baker, *Wiley Interdiscip. Rev.: Comput. Mol. Sci.*, 2015, **5**, 241.
- L. Guan, W. L. Wang, R. Shao, F. Y. Liu and S. W. Yin, *J. Mol. Model.*, 2015, **21**, 9.
- F. Liu, L. K. Du, J. Gao, L. L. Wang, B. Song and C. B. Liu, *J. Comput. Chem.*, 2015, **36**, 441.
- A. E. Sifain, N. Lubbers, B. T. Nebgen, J. S. Smith, A. Y. Lokhov, O. Isayev, A. E. Roitberg, K. Barros and S. Tretiak, *J. Phys. Chem. Lett.*, 2018, **9**, 4495.
- P. Bleiziffer, K. Schaller and S. Riniker, *J. Chem. Inf. Model.*, 2018, **58**, 579.
- J. C. Snyder, M. Rupp, K. Hansen, K. R. Muller and K. Burke, *Phys. Rev. Lett.*, 2012, **108**, 253002.
- J. Behler, *J. Chem. Phys.*, 2011, **134**, 13.
- J. Behler, *J. Chem. Phys.*, 2016, **145**, 9.
- J. Behler and M. Parrinello, *Phys. Rev. Lett.*, 2007, **98**, 146401.
- A. M. Goryaeva, J. B. Maillet and M. C. Marinica, *Comput. Mater. Sci.*, 2019, **166**, 200.
- K. Lee, D. Yoo, W. Jeong and S. Han, *Comput. Phys. Commun.*, 2019, **242**, 95.
- S. Rostami, M. Amsler and S. A. Ghasemi, *J. Chem. Phys.*, 2018, **149**, 8.
- M. Gastegger, L. Schwiedrzik, M. Bittermann, F. Berzsenyi and P. Marquetand, *J. Chem. Phys.*, 2018, **148**, 11.
- J. S. Smith, O. Isayev and A. E. Roitberg, *Chem. Sci.*, 2017, **8**, 3192.
- J. R. Boes and J. R. Kitchin, *Mol. Simul.*, 2017, **43**, 346.
- K. N. Kirschner, A. B. Yongye, S. M. Tschampel, J. Gonzalez-Outeirino, C. R. Daniels, B. L. Foley and R. J. Woods, *J. Comput. Chem.*, 2008, **29**, 622.
- E. Hatcher, O. Guvench and A. D. MacKerell, *J. Phys. Chem. B*, 2009, **113**, 12466.
- K. Nester, K. Gaweda and W. Plazinski, *J. Chem. Theory Comput.*, 2019, **15**, 1168.
- L. Breiman, *Mach. Learn.*, 2001, **45**, 5.
- C. Altona and M. Sundaralingam, *J. Am. Chem. Soc.*, 1973, **95**, 2333.
- C. Altona and M. Sundaral, *J. Am. Chem. Soc.*, 1972, **94**, 8205.
- M. J. Frisch, G. W. Trucks, H. B. Schlegel, G. E. Scuseria, M. A. Robb, J. R. Cheeseman, G. Scalmani, V. Barone, G. A. Petersson, H. Nakatsuji, X. Li, M. Caricato, A. V. Marenich, J. Bloino, B. G. Janesko, R. Gomperts, B. Mennucci, H. P. Hratchian, J. V. Ortiz, A. F. Izmaylov, J. L. Sonnenberg, D. Williams-Young, F. Ding, F. Lipparini, F. Egidi, J. Goings, B. Peng, A. Petrone, T. Henderson, D. Ranasinghe, V. G. Zakrzewski, J. Gao, N. Rega, G. Zheng, W. Liang, M. Hada, M. Ehara, K. Toyota, R. Fukuda, J. Hasegawa, M. Ishida, T. Nakajima, Y. Honda, O. Kitao, H. Nakai, T. Vreven, K. Throssell, J. A. Montgomery Jr, J. E. Peralta, F. Ogliaro, M. J. Bearpark, J. J. Heyd, E. N. Brothers, K. N. Kudin, V. N. Staroverov, T. A. Keith, R. Kobayashi, J. Normand, K. Raghavachari, A. P. Rendell, J. C. Burant, S. S. Iyengar, J. Tomasi, M. Cossi, J. M. Millam, M. Klene, C. Adamo, R. Cammi, J. W. Ochterski, R. L. Martin, K. Morokuma, O. Farkas, J. B. Foresman and D. J. Fox, *Gaussian 16 Rev. C.01*, Wallingford, CT, 2016.
- C. I. Bayly, P. Cieplak, W. D. Cornell and P. A. Kollman, *J. Phys. Chem.*, 1993, **97**, 10269.
- B. K. Rai and G. A. Bakken, *J. Comput. Chem.*, 2013, **34**, 1661.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot and E. Duchesnay, *J. Mach. Learn. Res.*, 2011, **12**, 2825.



- 35 W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey and M. L. Klein, *J. Chem. Phys.*, 1983, **79**, 926.
- 36 D. A. Case, R. M. Betz, D. S. Cerutti, T. E. Cheatham III, T. A. Darden, R. E. Duke, T. J. Giese, H. Gohlke, A. W. Goetz, N. Homeyer, S. Izadi, P. Janowski, J. Kaus, A. Kovalenko, T. S. Lee, S. LeGrand, P. Li, C. Lin, T. Luchko, R. Luo, B. Madej, D. J. Mermelstein, K. M. Merz, G. Monard, H. Nguyen, H. T. Nguyen, I. Omelyan, A. Onufriev, D. R. Roe, A. Roitberg, C. Sagui, C. L. Simmerling, W. M. Botello-Smith, J. Swails, R. C. Walker, J. Wang, R. M. Wolf, X. Wu, L. Xiao and P. A. Kollman, *AMBER*, 2016, University of California, San Francisco.
- 37 H. J. C. Berendsen, J. P. M. Postma, W. F. Vangunsteren, A. Dinola and J. R. Haak, *J. Chem. Phys.*, 1984, **81**, 3684.
- 38 W. F. Vangunsteren and H. J. C. Berendsen, *Mol. Phys.*, 1977, **34**, 1311.
- 39 A. W. Gotz, M. J. Williamson, D. Xu, D. Poole, S. Le Grand and R. C. Walker, *J. Chem. Theory Comput.*, 2012, **8**, 1542.
- 40 P. Zhang, P. Bao and J. Gao, *J. Comput. Chem.*, 2011, **32**, 2127–2139.
- 41 B. T. Thole, *Chem. Phys.*, 1981, **59**, 341.
- 42 J. W. Storer, D. J. Giesen, C. J. Cramer and D. G. Truhlar, *J. Comput.-Aided Mol. Des.*, 1995, **9**, 87.
- 43 M. Swart, P. T. Van Duijnen and J. G. Snijders, *J. Comput. Chem.*, 2001, **22**, 79.
- 44 J. B. Li, T. H. Zhu, C. J. Cramer and D. G. Truhlar, *J. Phys. Chem. A*, 1998, **102**, 1820.
- 45 C. Oostenbrink, A. Villa, A. E. Mark and W. F. Van Gunsteren, *J. Comput. Chem.*, 2004, **25**, 1656.
- 46 D. L. Mobley, K. A. Dill and J. D. Chodera, *J. Phys. Chem. B*, 2008, **112**, 938.
- 47 W. I. Weis and K. Drickamer, *Annu. Rev. Biochem.*, 1996, **65**, 441.
- 48 V. Spiwok, *Molecules*, 2017, **22**, 1038.
- 49 F. A. Quirocho, *Annu. Rev. Biochem.*, 1986, **55**, 287.

