**FULL PAPER**

Journal of COMPUTATIONAL CHEMISTRY    WILEY

# ClassicalGSG: Prediction of log *P* using classical molecular force fields and geometric scattering for graphs

Nazanin Donyapour[1]    |    Matthew Hirn[1,2,3]    |    Alex Dickson[1,4]

[1]Department of Computational Mathematics, Science and Engineering, Michigan State University, East Lansing, Michigan, USA

[2]Department of Mathematics, Michigan State University, East Lansing, Michigan, USA

[3]Center for Quantum Computing, Science and Engineering, Michigan State University, East Lansing, Michigan, USA

[4]Department of Biochemistry and Molecular Biology, Michigan State University, East Lansing, Michigan, USA

**Correspondence**
Alex Dickson, Department of Computational Mathematics, Science and Engineering, Michigan State University, East Lansing, MI, USA.
Email: alexrd@msu.edu

## Abstract

This work examines methods for predicting the partition coefficient (log *P*) for a dataset of small molecules. Here, we use atomic attributes such as radius and partial charge, which are typically used as force field parameters in classical molecular dynamics simulations. These atomic attributes are transformed into index-invariant molecular features using a recently developed method called geometric scattering for graphs (GSG). We call this approach "ClassicalGSG" and examine its performance under a broad range of conditions and hyperparameters. We train ClassicalGSG log *P* predictors with neural networks using 10,722 molecules from the OpenChem dataset and apply them to predict the log *P* values from four independent test sets. The ClassicalGSG method's performance is compared to a baseline model that employs graph convolutional networks. Our results show that the best prediction accuracies are obtained using atomic attributes generated with the CHARMM generalized force field and 2D molecular structures.

**KEYWORDS**
geometric scattering for graphs, graph convolutional networks, log *P* prediction, partition coefficients

## 1 | INTRODUCTION

The partition coefficient (*P*) measures the relative solubility of a compound between two solvents. It is defined as the ratio of the concentration of an uncharged compound dissolved in an organic solvent (e.g., octanol) in comparison to water. The logarithm of this ratio (log *P*) is considered to be one of the main factors in determining the drug-likeness of a chemical compound. The log *P* determines the lipophilicity, which affects bioavailability, solubility, and membrane permeability of a drug compound in the body. Moreover, an orally active drug should have a log *P* value of less than five according to one of the criteria of the famous Lipinski's Rule of Five.[1] Thus, predicting log *P* plays an essential role in the drug discovery process and is our main focus in this study. The prediction of log *P* is also used as a stepping stone to calculate other molecular properties such as the distribution coefficient (log *D*),[2] drug solubility (log *S*),[3,4] and Lipophilic efficiency (LiPE).[5] All of these properties have been used in drug discovery and design: log *D* is an effective descriptor for lipophilicity

of an ionized compound,[6] LiPE combines the potency and lipophilicity of a drug compound to estimate its quality, and log *S* is important to model the solubility of a compound in the human body.

The partition coefficient is widely used in cheminformatics and generally there are diverse experimental methods to measure it.[7] However, these methods are time-consuming and expensive for large databases of compounds and not feasible for compounds that are not synthesized.[8] Therefore, a large number of computational methods have been developed to predict accurate values of log *P*. These methods have a long history and can be classified by both their input features (atomic/fragment, molecular and hybrid) and their models or algorithms (parameteric models vs. machine learning methods) (Table 1).

In atomic-based or fragment-based methods,[9–15] which is based on atomic or fragment contributions, a set of atomic or fragment based descriptors is used as input to the model, while "molecular" methods use descriptors of the whole molecule, such as topology or molecular fingerprints.[16–13] "Hybrid" models combine atomic and

**TABLE 1**  log *P* prediction model classifications

| | | Mathematical models | |
| --- | --- | --- | --- |
| | | Parametric models | Machine learning methods |
| Features | Atomic/Fragment | XlogP3,[9] AlogP,[10] ClogP,[11] KowWIN,[12] JPlogP[13] | László et al.,[14] Huuskonen et al.[15] |
| | Molecular | MlogP,[16] iLogP,[17,18] Manhold[19] | AlogPS,[20,21] S + logP,[22] CSLogP[23] |
| | Hybrid | Silicos-IT LogP[24] | TopP-S,[25] OpenChem[26] |

molecular descriptors as input to the model.[24–26] In general, there are challenges with both atomic and molecular descriptors. In atomic-based or fragment-based methods, the accuracy of the atomic contributions depends on the similarity of the atomic environments of the atoms in the group. Unfortunately, more training data is needed as the number of groups grows larger. Fragment-based methods can struggle with defining the optimal size of fragments that participate in predictions. On the other hand, the accuracy of property-based methods heavily depends on the choice of molecular descriptors. Common descriptors include: 3D molecular structure,[25] molecular volume and surface area,[27] solvation free energies,[17,18,28] number of carbon atoms and heteroatoms.[16] Furthermore, these molecular descriptors can be difficult or computationally costly to generate.

A second way of categorizing log *P* predictors is by the type of mathematical model used to process the input data. Parameteric models use methods such as least squares estimation or multiple linear regression to fit parameters that govern the relative contributions of the different input features. Machine learning based methods such as support vector machines (SVM),[29–31] neural networks (NNs),[25,29,30,32] and graph convolutional networks (GCN)[33] have been used to predict log *P* values.

Recently, some methods have been described that create their own custom molecular features from atomic attributes, which go beyond simple additive models. The TopP-S[25] predictor was developed by Wu et al. and uses the atomic positions to create topological descriptors called Betti barcodes. These Betti barcodes are used as molecular features that are input into deep NNs, along with atomic attributes such as atom type. Results were shown to further improve upon the addition of 633 "molecular fingerprints" calculated from ChemoPy.[34] TopP-S has shown success in predicting log *P* over other methods like XlogP3, ClogP, and KowWIN using a group of independent test sets.

Graph representations of molecules have also shown success in various applications including predicting molecular properties,,[33,35–39] virtual screening,[40] and molecular force field calculations.[41] In particular, OpenChem[26] uses a graph representation of the molecules. Each atom represents a node in a graph and has a vector of atomic attributes including element type, valence, charge, hybridization, and aromaticity; bonded atoms are connected by an edge in the graph. GCNs are then trained on the graph representations created using these atomic attributes and the 2D structure – or, "graph structure" – of the molecules.

Graph representations are beneficial in that they are invariant to translation, rotation, and reflection symmetries. Another molecular symmetry that should be respected is invariance to the re-indexation of atoms: changing the order in which atoms are input to the model should not affect the computed molecular features. Summation operations respect re-indexation symmetry but it is not straightforward how to capture more detailed information about molecular structure while maintaining re-indexation symmetry. A recently-described method, geometric scattering for graphs (GSG),[42] provides a solution to this problem. GSG, which is analogous to GCNs, creates molecular features by scattering atomic attributes across a graph using lazy random walks. GSG is fast in creating re-indexation invariant features and also its feature vectors have the same length allowing us to easily measure the similarity of molecules, even those with different numbers of atoms. It has shown promising results in the classification of social network data and predicting enzyme commission numbers.[42]

Given this abundance of algorithms for creating molecular features, we seek to compare some different methods based on molecular graph representations and their ability to predict log *P*. Here we use GSG in combination with a set of atomic descriptors that are generated for use with classical molecular dynamics force fields: partial charges, atom type, and Lennard-Jones interaction parameters. We call this method "ClassicalGSG" and examine its performance as a function of different atomic/molecular features. We compare the ClassicalGSG results with GCNs trained on the same data and using a variety of atomic attributes, including those from previous work.[33] We then evaluate the performance of ClassicalGSG on several independent test sets and study the properties of features generated in the pipeline of GSGNN models. In addition, we investigate the properties of molecules with high log *P* prediction error. We conclude with a discussion about the GSG method generated features, the computational cost of generating atomic attributes with CHARMM generalized force field (CGenFF) and general AMBER force field 2 (GAFF2), and the relative performance of 2D versus 3D structure in predicting log *P* values.

## 2 | METHODS

### 2.1 | Datasets and generation of atomic attributes

The dataset used in this work is generated by Korshunova et al.[26] from the public version of PHYSPROP database.[43] This dataset consists of 14,176 molecules in SMILES format and their corresponding log *P* values, which refer to as the OpenChem dataset.

**TABLE 2** Independent test sets used for evaluating ClassicalGSG models

| Test set name | Number of molecules |
|---|---|
| FDA[9] | 406 |
| Huuskonen[15] | 348 |
| Star[19] | 223 |
| NonStar[19] | 43 |

Abbreviation: FDA, food and drug administration.

**TABLE 3** Classifying atoms in 5 categories (AC5)

| Atom type | Category number |
|---|---|
| Hydrogen | 1 |
| Oxygen and Nitrogen | 2 |
| Carbon with hybridization value < 3 | 3 |
| Carbon with hybridization value = 3 | 4 |
| Others | 5 |

The molecules are converted from SMILES format to mol2 format and their 3D structure is created by OpenBabel (https://github.com/openbabel/openbabel). Then CGenFF[44,45] and GAFF2[46,47] parameter files are generated for each molecule. These force field parameter files are either created by CGenFF using the CGenFF tool of the SilcsBio software package (http://silcsbio.com) or by the Antichamber tool implemented in the Ambertools18[48] package. The process of generating CGenFF parameter files fails for 175 molecules, and GAFF2 for 681 molecules. These 774 molecules are removed from the OpenChem dataset, resulting in a dataset of 13,402 molecules. Then 80% of the molecules are used for training and the rest for testing. In addition, we evaluate our trained models on four independent test sets shown in Table 2. Star and NonStar[19] test sets are publicly available on https://ochem.eu/article/17434 and the Husskonen[15] test set can be found on https://ochem.eu/article/164. The food and drug administration (FDA) dataset contains drug molecules that are approved by the FDA of the United Sates and originally prepared by Chen et al.[9]

As mentioned above, we use atomic attributes including partial charges, atom type, and Lennard-Jones interaction parameters. Below, we explain how we generate these atomic attributes.

Atomic partial charges for each atom are extracted from the parameter files generated by either the CGenFF[44,45] or GAFF2[46,47] force field generator tools. To determine atom type classifications, we compared a number of different schemes. In one scheme, we classify atom types in one of five categories as shown in Table 3. This is referred to below as "AC5." Alternatively, we directly use the atom types as generated by either CGenFF or GAFF2; referred to below as "ACall." In the third classification scheme, we manually sorted CGenFF atom types into 36 groups (AC36; Table S1) and GAFF2 atom types into 31 groups (AC31; Table S2) based on chemical knowledge. Specifically, efforts were made to make new groups for atom types

with different elements and hybridization values and to separately identify atoms that are members of ring structures. Finally, a forth classification scheme simply uses a uniform atom type for all atoms, referred to as "AC1."

The two Lennard-Jones parameters – radius ($r$) and well-depth ($\varepsilon$) – are extracted from either CHARMM or AMBER parameter files for each atom type. In summary the atomic attributes are defined by both the force field generation tool (CGENFF or GAFF2) and the atom classification scheme AC1, AC5, AC36, AC31, or ACall and we generally refer to these as "FF" atomic attributes. The atomic charges and Lennard-Jones parameters are scalar values, but atomic types are represented in one-hot encoding format. For comparison, we also examine a different set of atomic attributes, following previous work for the log $P$ predictor from the OpenChem toolkit (https://github.com/Mariewelt/OpenChem.git).[26] These atomic attributes are defined as: atom element type, valence, charge, hybridization, and aromaticity. The 3D structure is generated from SMILES strings by RDkit (https://www.rdkit.org), then again using RDkit the atomic attributes are produced for each molecule and represented in one-hot encoding format as node features in OpenChem. These are referred to below as" OPENCHEM" atomic attributes.

## 2.2 | log *P* predictions using GSG

### 2.2.1 | Geometric scattering for graphs

The GSG method, which has been introduced in Reference 42 is a feature extraction method for graph data types and uses a geometric transform defined on the graph. In the GSG method, molecules are represented by a graph of atoms where each atom has a vector of attributes; a single attribute evaluated at each vertex is also referred to as a "graph signal". GSG combines the molecular structure (defined using an adjacency matrix) and atomic signal vectors to construct invariant, stable, and informative features as shown in Figure 1.

Now we briefly explain the mathematics behind this method. Let $G = (V, E, W)$ be a weighted graph and $\mathbf{x}(v_l)$, $1 < l \leq n$ be the signal function defined on each node where $n$ is the number of nodes in the graph and $v_l$ represents node $l$. A graph random walk operator is defined as $P = \frac{1}{2}(I + AD^{-1})$, where $A$ is the adjacency matrix and $D$ is the degree matrix. $P^t$ describes the probability distribution of a lazy random walk after $t$ steps and hence the term *lazy random walk* for $P$. In fact, graph random walks are low-pass filters that act like convolution filters to capture graph features at different scales. Graph wavelet operators – also known as "signal transform operators for graphs" – are defined at the scale $2^j$ as $\Psi_j = P^{2^{j-1}} - P^{2^j}$. Finally, a set of scattering transform operators $\mathbf{S}_0$, $\mathbf{S}_1$, and $\mathbf{S}_2$, are defined using the graph signals and wavelets to compute the molecular-level features. Unique $\mathbf{S}_0$, $\mathbf{S}_1$, and $\mathbf{S}_2$ vectors are created using different moments ($q$) of $\mathbf{x}$ and their wavelet coefficients $\psi$, as well as different wavelet scales indexed by $j$:
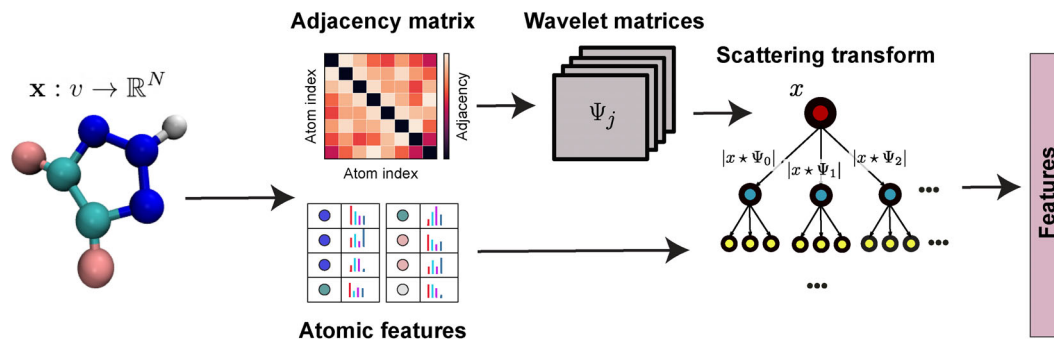
**FIGURE 1** Architecture of the geometric scattering for graphs method. The adjacency matrix describes the graph structure of the molecule. Each atom has a set of attributes that are shown as colored bars. Wavelet matrices ψ are built at different logarithmic scales, $j$, using the adjacency matrix as described in the text. Finally, the scattering transform is applied to get the graph features using both the wavelet matrices and the signal vectors. Modified from figure made by Feng et al.[42]

$$
\begin{aligned}
(\mathbf{S}_0)_q &= \sum_{l=1}^{n} \mathbf{x}(v_l)^q, & 1 \leq q \leq Q \\
(\mathbf{S}_1)_{j,q} &= \sum_{l=1}^{n} \left| \Psi_j \mathbf{x}(v_l) \right|^q, & 1 \leq j \leq J \quad 1 \leq q \leq Q \\
(\mathbf{S}_2)_{jj',q} &= \sum_{l=1}^{n} \left\| \Psi_{j'} |\Psi_j \mathbf{x}(v_l)| \right\|^q, & 1 \leq j < j' \leq J \quad 1 \leq q \leq Q
\end{aligned}
\tag{1}
$$

The scattering operators are named based on the number of times the wavelets ψ are used to transform the signal **x**. For example, $\mathbf{S}_0$ features do not use the wavelet operators, and are simply different moments of the atomic signals. These are called the "Zero order scattering moments." Accordingly, the $\mathbf{S}_1$ operators with one wavelet transformation are called "First order scattering moments" and $\mathbf{S}_2$ with two wavelet transformations are called the "Second order scattering moments." Note that these allow mixing between different wavelet scales as $j$ and $j'$ are set independently.

We define the set **S** as the concatenation of all $\mathbf{S}_0, \mathbf{S}_1$, and $\mathbf{S}_2$ vectors using all possible values of $q$, $j$, and $j'$ as specified by the ranges in Equation (1). The set of graph features is generated deterministically from the atomic signals and adjacency matrix. The size of this set of features is equal to $QN_s(1 + J(J + 1)/2)$, where $N_s$ is the number of attributes per vertex, although this is lower if not all zeroth, first, and second order features are used.

As mentioned above, one of the inputs to GSG is the adjacency matrix of the graph and we consider two distinct ways to represent it. One is to use the 2D connectivity of the molecule, where $A_{ij} = 1$ indicates a bond between atoms $i$ and $j$, and $A_{ij} = 0$ otherwise. Alternatively, the adjacency matrix can be calculated using the 3D structure, where $A_{ij} = f(R_{ij})$, and $f(R)$ is a smooth and differentiable function that is equal to 1 at low $R$ and decreases to 0 as $R$ exceeds some cutoff value $R_c$[49]:

$$
f(R_{ij}) =
\begin{cases}
0.5\left( \cos\left( \dfrac{\pi R_{ij}}{R_c} \right) + 1 \right) & \text{for } R_{ij} \leq R_c \\
0.0 & \text{for } R_{ij} > R_c
\end{cases}
\tag{2}
$$

where $R_c$ is the radial cutoff and $R_{ij}$ is the Euclidean distance between atoms $i$ and $j$. If $f(R_{ij})$ is non-zero, nodes $i$ and $j$ are considered connected and disconnected otherwise. The wavelet operators in GSG can be defined using either discrete (2D) or continuous (3D) values in the adjacency matrix.

## 2.3 | Neural networks architecture and training

To predict the log $P$ values from the GSG features (**S**), we develop a feedforward NN using PyTorch[50] with a nonlinear activation function such as rectified linear unit (ReLU). To determine the best performing network, we consider three important hyperparamters including the number of hidden layers, the number of neurons in a hidden layer, and the dropout rate whose ranges are shown in Table 3. We perform cross validation using the PyTorch wrapper Skorch https://skorch.readthedocs.io/en/stable/ to tune the hyperparameters. The loss function of our NN model is MSELOSS (mean squared error) and we use the Adam (adaptive momentum estimation)[51] for optimizing the parameters. We use MultiStepLR, which has an adjustable learning rate set to the initial value of 0.005 and dynamically decreases during training every 15 steps by a factor of 0.5. In training our GSGNN models, we used two data regularization methods: L1 norm and standardization using the StandardScaler function from sci-kit learn.[52] For GSG features with maximum wavelet scales of $J = 4, 5$, and 6 we mostly benefit from using the standardization method. Other settings can be found in Table 4.

## 2.4 | log $P$ predictions using GCNs

### 2.4.1 | Graph convolutional neural network

Graph convolutional neural networks (GCNs) extend the application of convolutional NNs to graph data. The goal of GCNs is to take a graph and generate features according to node attributes and graph structure. The heart of the GCN method is the convolution operator, which aggregates the features of neighboring nodes

within the graph. Recently, various implementations of GCNs[38,51,53–57] have been developed to increase the speed and accuracy of the GCN models. In this paper, we employ the GCN model that was developed in OpenChem[26] based on the method introduced by Kipf and Welling.[55] Similar to the GSG model above, this model takes a graph $G = (V, E)$ as the input where each node has a vector of attributes $\mathbf{x}$. The model processes the graph by passing it through multiple hidden layers performing convolution operations (Figure 2).

The GCN method works by propagating a feature matrix $H$ – originally set to the value of the attribute vector $\mathbf{x}$ for all nodes – through a set of convolution operators as follows:

$$H^{(0)} = \mathbf{X},$$
$$H^{(l+1)} = \sigma\left(\tilde{A}H^{(l)}W^{(l)}\right), \tag{3}$$

where $\tilde{A} = A + I$, is the adjacency matrix of the input graph with self-connections, $I$ is the identity matrix, $W^{(l)}$ is the trainable weight matrix of layer $l$ and $\sigma$ is a non-linearity function such as ReLU. $H^{(l)}$ denotes the value of the feature matrix on layer $l$.

Each convolution layer is followed by a graph max-pooling layer introduced in Reference 57. Following these convolution and max-

polling operations, the final set of graph features is obtained by a graph gather layer, where the values of each feature are summed over nodes. This last layer gives GCNs their index-invariance property.

### 2.4.2 | GCN architecture and training

We predict log $P$ using GCNs using the model implemented in the OpenChem toolkit https://mariewelt.github.io/OpenChem/html/index.html. This model contains five layers of graph convolutions with a hidden layer size of 128. The GCN layers are followed by max-pooling and a graph gather layer. A 2-layer NN with ReLU as the activation function is added after the graph gather layer. The PyTorch Adam optimizer and the MSELOSS (mean squared error) are used as the parameter optimizer and training loss function, respectively. We use the MultiStepLR learning scheduler, implemented in PyTorch, with an initial value of 0.001, a step size of 15 and a learning coefficient of 0.5 for training the model. Parameters used for the GCN training are summarized in Table 5.

## 3 | RESULTS

### 3.1 | Evaluation of molecular representations

For our set of atomic attributes (FF), we use parameters from classical MD force fields: the partial charge $\sigma$, and the Lennard-Jones radii and well-depth. We first compare the performance of atomic attributes generated using two different algorithms: GAFF2 and CGenFF. Both algorithms are used to automatically generate force field parameters for small molecule ligands, by matching atomic environments with atoms that are part of existing force fields; GAFF2 uses the Amber[58] force field, while CGenFF uses the CHARMM[44,45] force field. Here we generate molecular features from atomic attributes using GSG (see Section 2). We examine the four different atom type classification

**TABLE 4** Neural network settings

| Parameter | Values |
| --- | --- |
| Number of hidden layers | [2, 3, 4, 5] |
| Size of hidden layers | [300, 400, 500] |
| Dropout rate | [0.2, 0.4] |
| Initial learning rate | 0.005 |
| Learning coefficient | 0.5 |
| Batch size | 256 |
| Max epoch size | 400 |

*Note:* Square brackets denote possible parameter values used in the grid search method.
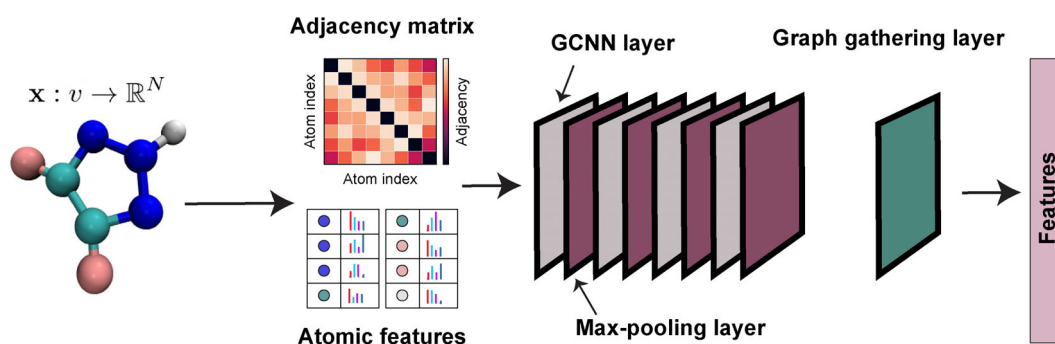


**FIGURE 2** Architecture of the graph convolutional network (GCN) method. The adjacency matrix describes the graph structure of the molecule. Each atom has a set of attributes and are shown as colored bars. GCN layers are shown by gray color and are followed a max-pooling layer which is shown in purple. The graph gathering layer is shown in green color adds features on all nodes to generate the molecular feature vector

**TABLE 5** Neural network settings

| Parameter | Values |
| --- | --- |
| Number of GCN hidden layers | 5 |
| Number of NN hidden layers | 2 |
| Size of hidden layers | 128 |
| Dropout rate | 0 |
| Initial learning rate | 0.01 |
| Learning coefficient | 0.5 |
| Batch size | 128 |
| Max epoch size | 200 |

Abbreviations: GCN, graph convolutional network; NN, neural network.

**TABLE 6** Parameters for the geometric scattering for graphs algorithm

| Parameter | Values |
| --- | --- |
| Adjacency matrix ($A$) | [2D, 3D] |
| Wavelet maximum scale index ($J$) | [4, 5, 7, 8] |
| Scattering operators ($z, f, s$) | Four combinations $(z, f), (z, s), (f, s), (z, f, s)$ |
| Atom classification | [AC1, AC5, AC36/AC31, ACall] |

*Note:* The square brackets show all of the values examined for each parameter. For "scattering operators," "$z$" represents the zero order operator, "$f$" is first order, and "$s$" is second order.

schemes discussed above: AC1, AC5, AC36/AC31, and ACall. The GSG parameters used to construct the molecular features are shown in Table 6.

We trained 160 models for each of the CGenFF and GAFF2 atomic attributes sets, each trained using a 3-fold cross-validation method. After training, we ran tests on subsets of the full database using an 80:20 train:test split. We then calculated evaluation metrics such as the correlation coefficient ($r^2$), root mean squared errors (RMSE) and mean unsigned errors (MUE) between the predicted and experimental log $P$ values. In Figure 3, $r^2$ and RMSE values are shown for predictions using 2D and 3D molecular structures. Each $r^2$ and RMSE is averaged on 20 NN models that are created with all combinations of the four geometric scattering operator sets and the five wavelet step numbers as shown in Table 6.

We find that models with 2D structure universally have higher accuracy in prediction of log $P$ values. Additionally, it demonstrates that CGENFF atomic attributes on average are more accurate in predicting log $P$ values compared to GAFF2 atomic attributes, although this is not true for ACall. We find that good results on average are obtained by classifying atom types using AC36 categories; the best performing individual models are also found in this category. We thus use CGenFF, 2D structure and AC36 for all models going forward.

We next investigate different scattering moment sets and wavelet scales in Figure 4. We find that for each examined wavelet number

(4, 5, 6, 7, and 8) there is at least one model with an $r^2$ value of 0.9. The model also performs well for all combinations of scattering moments, making it difficult to determine an optimal set of parameters.

To further evaluate the accuracy of our method, we examine four different independent test sets: FDA, Huuskonen, Star and NonStar which are explained in Section 2. To avoid any accidental overlap, we identified shared molecules and removed them from our training dataset. The training dataset contains 10,722 randomly selected molecules from the OpenChem dataset that successfully are processed by CGenFF. The trained models are evaluated using these test sets and the $r^2$ between the actual and predicted log $P$ values are calculated. The results are shown in Figure 5. For the Huuskonen and FDA data sets we find that there is at least one GSGNN model with $r^2 = 0.92$ and $r^2 = 0.90$, respectively. The best performing models for the Star test set is a maximum wavelet scale of $J = 4$ and zero- and second-order scattering operators, while the best model for NonStar is a maximum wavelet scale of $J = 7$ and zero-, first- and second-order scattering operators.

The $r^2$, RMSE and MUE of best-performing GSGNN model for each test set are shown in Table 7. In paper[25] the log $P$ values for the FDA, Star, NonStar test sets are predicted using their topology-based models (TopP-S) and the different established methods such as ALOGPS, XLOGP3, and KowWIN. We find that our results do not outperform the TopP-S method but we do achieve higher accuracy compared to other known methods such as ALOGPS, KowWIN, and XLOGP3 for FDA test set (Table S3), methods like XLOGP3 and ALOGP for Star test set (Table S4) and XLOGP3 method for NonStar test set (Table S5).

## 3.2 | Geometric scattering versus graph convolution

To compare the performance of geometric scattering (GSGNN) models to the GCN models, we trained models using these two methods with two different sets of atomic attributes "FF" and "OPENCHEM" as described in Section 2. To rule out the influence of the training/test split, we trained five models for each method where the data is randomly divided into training and test sets with 80/20 ratio, respectively. Each of the five GSGNN models is trained using 5-fold cross validation. The RMSE and $r^2$ are determined for each model and are shown in Table 8. Note that the atomic attributes for these GSGNN models are generated by 2D molecular structure, CGenFF force fields, and AC36 atom type classification scheme. The GSG parameters that used to generate the molecular features are from one of the best performing models (a wavelet maximum scale of $J = 4$ and all three of the zero-, first- and second-order scattering operators) from Figure 4. Our results suggest that the ClassicalGSG method, which is a GSGNN model trained on FF atomic attributes, has the most accurate prediction of log $P$ values.
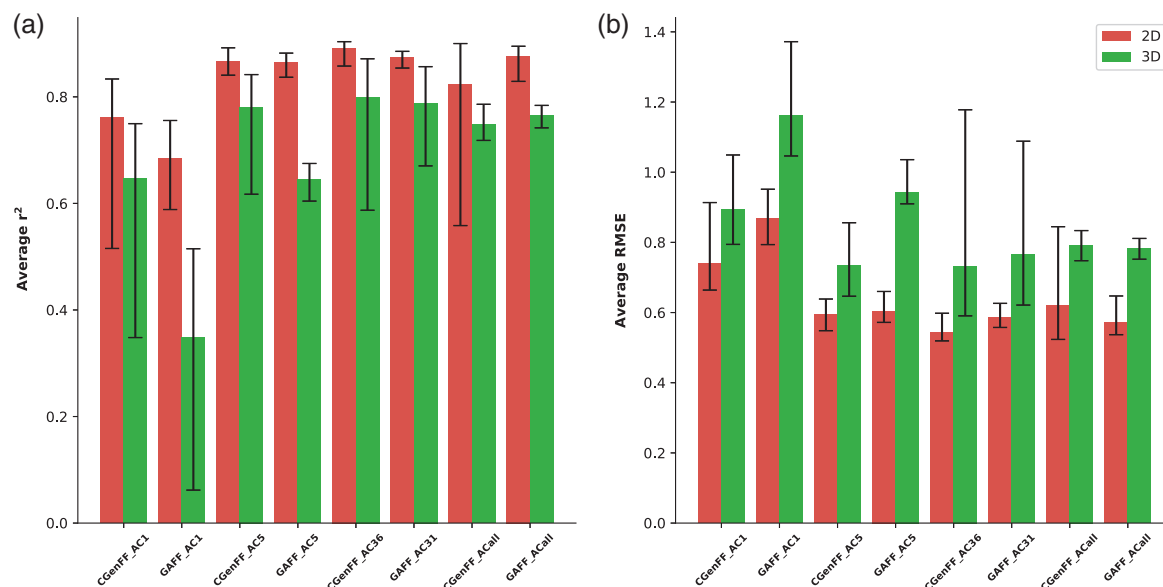
(a)



(b)



**FIGURE 3** Average $r^2$ (A) and root mean squared errors (B) for the OpenChem test set using GSGNN models. Each average is calculated over 20 individual parameter values and the error bars show the best and worst performing models. The atomic attributes are generated with either CHARMM generalized force field or general AMBER force field 2 force fields and using one of three atom type classification schemes ("AC1," "AC5," "AC36/AC31," or "ACall")
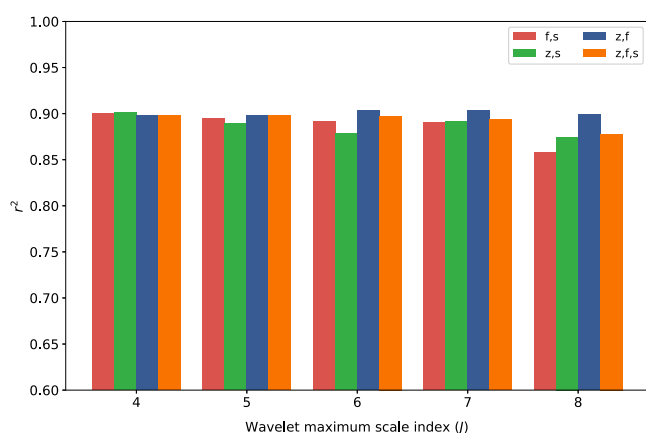


**FIGURE 4** The $r^2$ for the OpenChem test set using GSGNN models. The atomic attributes are all generated with CHARMM generalized force field force fields, AC36 atom type classification scheme, and 2D molecular structure

## 3.3 | Features visualization

To examine the molecular features that are generated in the pipeline of the ClassicalGSG method, we visualize the features generated by the GSG method and the last layer of NN model using $t$-distributed stochastic neighbor embedding ($t$-SNE)[59] plots. The $t$-SNE plots are intended for projecting high dimensional data into the low dimensional space so it can be visualized readily. Figure 6 shows the GSG and NN features for molecules in the OpenChem test set visualized in a 2D space. Here each point shows a molecule in the OpenChem test set and is colored by its log $P$ value. The initial dimension of GSG features is 1716 while NN features

have size of 400. These GSG features are generated from CGenFF atomic charges, 2D molecular structure, AC36 type classification scheme and all three scattering moment operators with wavelet step number of 4.

We can observe from Figure 6 that features extracted from NN are more discriminative with respect to log $P$ values. More specifically, in the last layer of the NN model, molecules with similar log $P$ value tend to be near each other and form distinct clusters. In contrast, GSG features are less discriminative where nearby molecules have different log $P$ values. To further verify this, five nearest neighbors are determined for each point in the GSG and NN reduced feature space. Then, the difference between the actual log $P$ value of each point and its five neighbors are calculated and averaged. This value averaged over all the molecules and is shown by $\langle \Delta \log P \rangle_N$ inside the figure. The $\langle \Delta \log P \rangle_N$ in the reduced GSG features space is 0.81 while this value for the reduced NN features space is only 0.17. The average log $P$ difference between two random points is 2.00. This shows that molecules with similar log $P$ values are closer in the NN features space.

This is noteworthy, as the GSG features are task independent and therefore are not adapted to any particular prediction task, including log $P$ prediction. On the other hand, the NN model, which takes as input the GSG features, is a supervised model that is trained for the specific task of log $P$ prediction. The results in Figure 6, in addition to the results in Tables 7 and 8, illustrate that the GSG features provide not only a translation, rotation, and permutation invariant representation of the molecules, but one that is also sufficiently rich so that, when combined with a downstream supervised NN, the resulting model provides accurate estimates of log $P$ values for molecules.
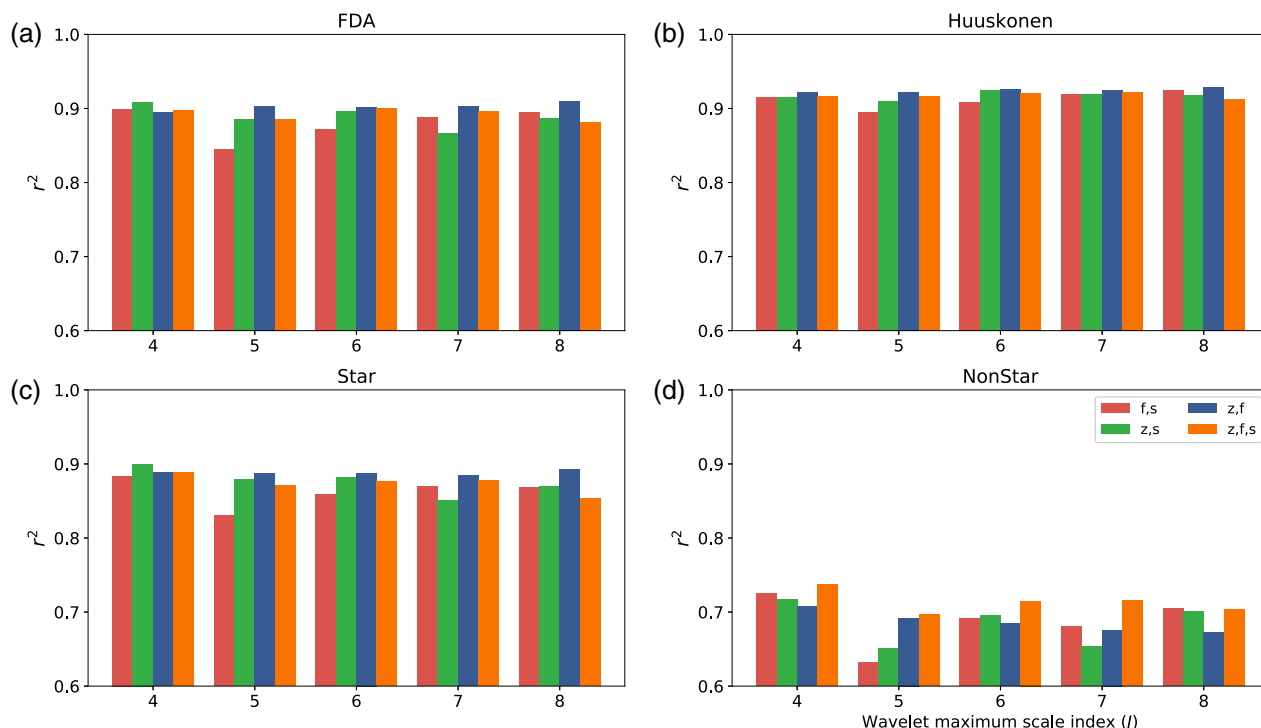
**FIGURE 5** The $r^2$ for different test sets using GSGNN models. (A) Shows $r^2$ for the food and drug administration test set. (B) Represent $r^2$ for the Huuskonen test set. (C) and (D) show $r^2$ for the Star and NonStar test sets, respectively. The horizontal axis indicates the maximum wavelet scale $J$. The atomic attributes are generated with 2D molecular structure, CHARMM generalized force field force fields and using AC36 atom type classification scheme

**TABLE 7** The log $P$ prediction performance results for four independent test sets

| Dataset | $r^2$ | RMSE | MUE |
|---|---|---|---|
| FDA | 0.91 | 0.56 | 0.35 |
| Star | 0.90 | 0.49 | 0.35 |
| NonStar | 0.72 | 1.02 | 0.81 |
| Huuskonen | 0.93 | 0.37 | 0.23 |

Abbreviations: FDA, food and drug administration; MUE, mean unsigned errors; RMSE, root mean squared errors.

## 3.4 | Features contribution analysis

In order to determine important atomic attributes that contributed to the prediction of log $P$ we performed ablation analysis on the test sets used in Section 3. The molecular features for these models were created by atomic attributes, including atom partial charge, Lennard-Jones radius and epsilon parameters and atomic classification scheme AC36. Feature ablation analysis identifies important features using the change in the model's prediction upon removal of each feature or group of features. The ablation importance values for 39 atomic attributes are determined using the Captum library,[60] and are averaged over the five models created by 5-fold cross validation in Section 3. The top 10 most important attributes are listed in Table 9. Intuitively, these attributes are associated with polar groups that could have large impacts on molecular hydrophobicity. While none of these atomic

types is extremely rare – the least common is N3, which occurs on average 0.33 times per molecule (Table S6) – the most commonly observed atomic types (e.g., non-polar hydrogen, sp3 carbon) are not placed highly in this list. The ablation analysis thus reveals some important information on which features have outsized importance for log $P$ prediction. We also note that the scalar molecular dynamics parameters (LJ epsilon, atomic charge, LJ radius) are found to be significant for log $P$ prediction, implying that they contain useful information that is not covered by atomic types alone.

## 3.5 | Distinguishing features of failed molecules

To investigate the common characteristics of failed molecules during the prediction of their log $P$ value, we created molecular fingerprints, which have been used extensively in cheminformatics, QSAR/QSPR predictions, and drug design.[25] Here we employ ChemoPy[34] to create a set of constitutional fingerprints. To specify the failed molecules we define a failure cutoff value (here, 0.5) where if the difference between the actual and predicted log $P$ values is larger than the cutoff, we consider the molecule as a failed prediction. We determine the failed molecules from ClassicalGSG models we described in Section 3. These models are trained on features constructed by CGenFF force fields, 2D structure and AC36 atom type classification method. The probability distribution of 30 constitutional fingerprints are calculated for all molecules in the OpenChem dataset and for failed

| Method name | Atomic attributes | Model | $r^2$ (STD) | RMSE (STD) |
|---|---|---|---|---|
| ClassicalGSG | FF | GSGNN | 0.91 (0.003) | 0.52 (0.009) |
| - | OPENCHEM | GSGNN | 0.89 (0.003) | 0.57 (0.005) |
| - | FF | GCN | 0.75 (0.091) | 0.91 (0.18) |
| OpenChem | OPENCHEM | GCN | 0.79 (0.052) | 0.83 (0.122) |

**TABLE 8** The log $P$ prediction results using different set of features and models

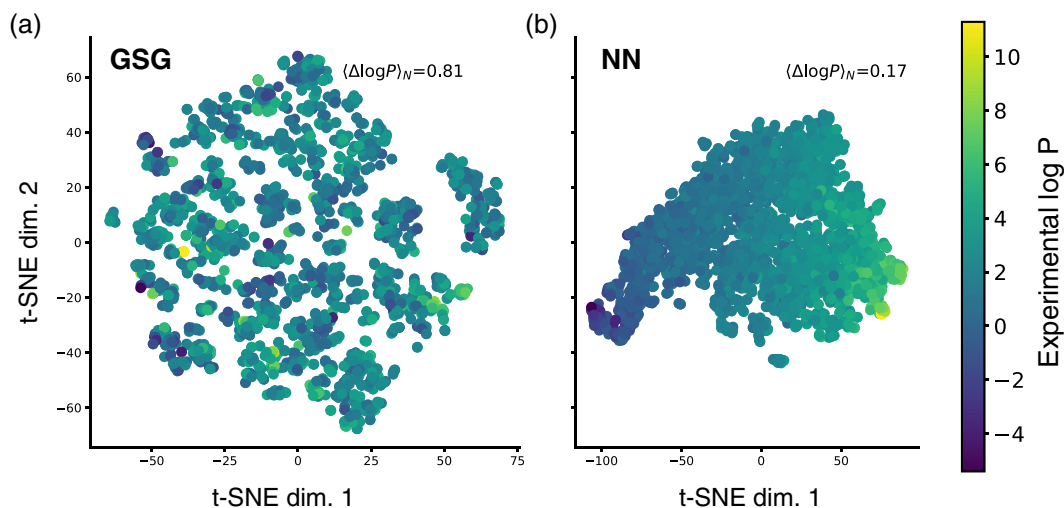Abbreviations: GCN, graph convolutional network; RMSE, root mean squared errors.



**FIGURE 6** The $t$-distributed stochastic neighbor embedding ($t$-SNE) plots with geometric scattering for graphs and neural network (NN) features of the OpenChem test set molecules. Each represents a molecule and is colored by its actual log $P$ value. $\langle \Delta \log P \rangle_N$ shows the mean log $P$ difference value calculated over the nearest neighbors in the $t$-SNE plot. (A) the geometric scattering for graphs features of size 1716 are projected into 2-dimensional space. (B) The NN features from the last hidden layer with size of 400 are projected into 2-dimensional space

| Rank | Atomic parameter | Average importance | STE |
|---|---|---|---|
| 1 | AC36 O2 (sp2 oxygen) | 0.0716 | 0.0078 |
| 2 | AC36 N3 (sp3 nitrogen) | 0.0635 | 0.0054 |
| 3 | AC36 N2R (sp2 nitrogen in ring) | 0.0560 | 0.0049 |
| 4 | LJ epsilon | 0.0516 | 0.0048 |
| 5 | AC36 H2 (polar hydrogen) | 0.0499 | 0.0017 |
| 6 | Atomic charge | 0.0475 | 0.0029 |
| 7 | AC36 C2 (sp2 carbon) | 0.0448 | 0.0033 |
| 8 | LJ radius | 0.0416 | 0.0058 |
| 9 | AC36 O3 (sp3 oxygen) | 0.0392 | 0.0030 |
| 10 | AC36 C2R (sp2 carbon in ring) | 0.0373 | 0.0086 |

**TABLE 9** The most important atomic attributes as predicted by ablation analysis

*Note:* Atom types are groups of CGenFF atom types. See Table S1 for more detail.

molecules in each of the 5 GSGNN models. Kullback–Leibler divergence (KL-divergence)[61] values are determined by comparing probability distributions of all data with distributions of the failed molecules from each model. The KL-divergence values are averaged over five models and their standard error (STE) is shown in Table 10. We note that the PCX descriptors count the number of shortest paths of length X. The attributes with the highest KL-divergence values are: PC counts with lengths of 2, 1, 3, and 4;

molecular weight (weight); the number of carbon atoms (ncarb); and the number of heavy atoms (nhev). In other words, the distributions of failed molecules are enriched in particular values of these attributes.

The probability distributions of fingerprints with the highest KL-divergence values are shown in Figure 7. The distributions for the failed molecules largely follow the complete dataset distribution, but are enriched toward higher values, suggesting that the

**TABLE 10** The averaged KL-divergence between fingerprint distributions of all data versus failed molecules averaged over 5 GCGNN models

| Fingerprint | Average_KL | STE | Fingerprint | Average_KL | STE |
|---|---|---|---|---|---|
| PC2 | 0.048 | 0.005 | nhet | 0.028 | 0.003 |
| PC1 | 0.047 | 0.006 | noxy | 0.026 | 0.003 |
| PC3 | 0.046 | 0.007 | nrot | 0.026 | 0.004 |
| Weight | 0.045 | 0.004 | nsulph | 0.026 | 0.003 |
| PC4 | 0.044 | 0.004 | ndb | 0.017 | 0.002 |
| ncarb | 0.043 | 0.005 | ndonr | 0.012 | 0.004 |
| nhev | 0.043 | 0.007 | ncof | 0.011 | 0.001 |
| nta | 0.043 | 0.007 | AWeight | 0.011 | 0.002 |
| PC5 | 0.04 | 0.005 | nnitro | 0.01 | 0.001 |
| naro | 0.04 | 0.003 | ncocl | 0.009 | 0.002 |
| PC6 | 0.039 | 0.005 | nhal | 0.008 | 0.001 |
| nsb | 0.038 | 0.004 | nphos | 0.007 | 0.003 |
| naccr | 0.036 | 0.002 | ncobr | 0.006 | 0 |
| nring | 0.035 | 0.004 | ncoi | 0.005 | 0.002 |
| nhyd | 0.031 | 0.003 | ntb | 0.001 | 0 |

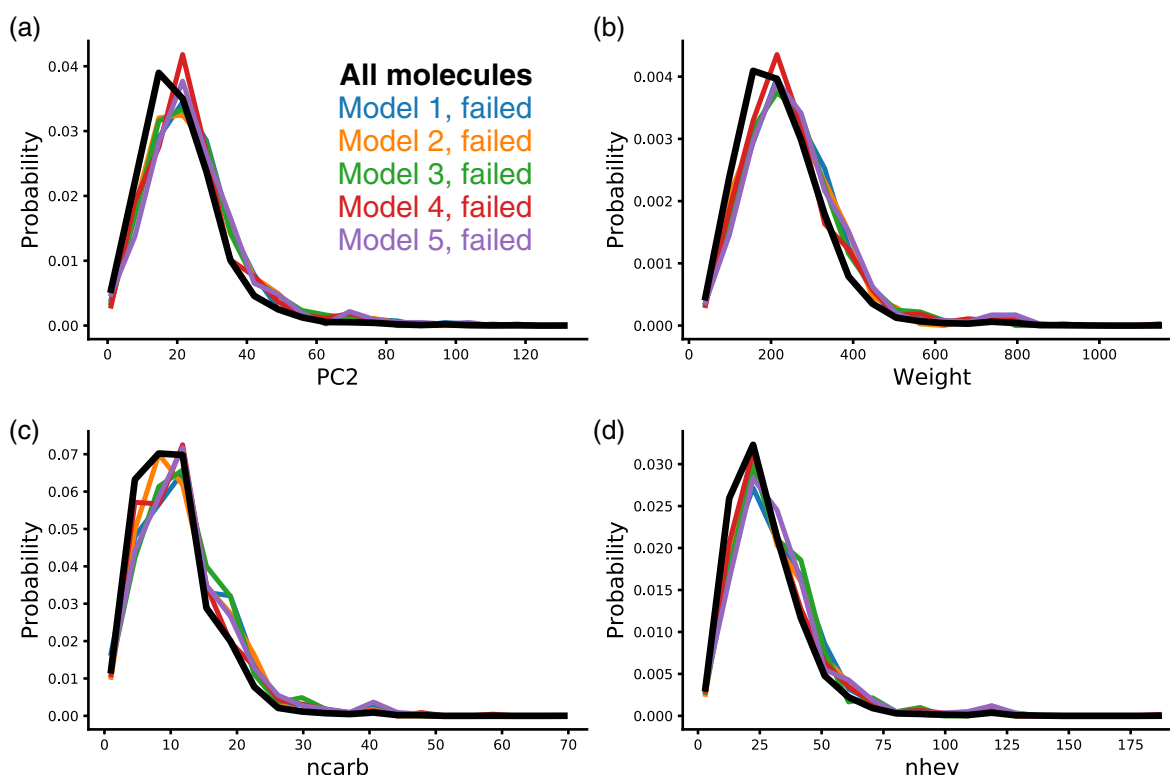Abbreviation: CGenFF, CHARMM generalized force field.



**FIGURE 7** Probability distributions of molecular fingerprints. The histograms show the distribution of fingerprints of all data and failed molecules of 5 GCGNN models. The distribution of all data is shown in thick black line. (A) The number of shortest paths of length 2, (B) the atomic weight, (C) the number of carbon atoms (ncarb), and (D) the number of heavy atoms

log $P$ predictions are more likely to fail for larger molecules. Interestingly the distributions for attributes that count rare element types (e.g., number of phosphorus [nphos], number of iodine atoms [ncoi], and number of bromine atoms [ncobr]) do not show large KL-divergence values, although this might have been expected given that they are poorly represented in the training set. The atom element types and their count in our dataset is shown in Figure S1.

# 4 | DISCUSSION AND CONCLUSIONS

In this paper, we introduced a method called "ClassicalGSG" for predicting the partition coefficient. Our method uses atomic attributes that are usually utilized as parameters for classical MD simulations. These parameters include partial charges, Lennard-Jones well depth, Lennard-Jones radius and atomic type. The pipeline for generating these parameters includes: creating 3D structure from SMILES, creating PDB and Mol2 formatted files, and generating atomic parameters using either Antechamber (GAFF2) or CGenFF tools. We note that, in our implementation, the Antechamber tool is about 200 times slower than CGenFF, requiring a couple of days to process 10,000 molecules.

We employ the GSG method to transform the atomic attributes into molecular features that satisfy re-indexation symmetry. Our results indicate that GSG is powerful enough to capture the universal molecular features, as the average log $P$ difference for all pairs of adjacent molecules in a $t$-SNE plot ($\langle \Delta \log P \rangle_N$) is 0.81, which is a low value compared to the random pair log $P$ difference (2.00). As these features are general to the molecule and not specific to log $P$, this suggests that they can be used in multi-task NNs to predict other molecular properties, such as solubility (log $S$), melting point, pKa, and intestinal permeability (e.g., Caco2).[62]

Our results show that employing a 2D molecular structure in ClassicalGSG yields accurate log $P$ predictions compared to 3D structures, and this confirms the same conclusion achieved previously.[63,64] This could be due to difficulties in generating appropriate 3D structures, or that a single 3D structure is insufficient to capture the high probability conformations of a given molecule. Additionally, our 3D adjacency matrix did not explicitly distinguish between bonded and non-bonded interactions, where the former are much more important to determine molecular properties.

The results reported here for four independent external test sets show that our log $P$ ClassicalGSG method is generalizable to new molecules. However, we do not expect this model to perform well for molecules with new elements or functional groups that are not covered in the training set. Like other empirical methods, we expect the accuracy will improve as the availability of training data grows.

## ORCID

*Nazanin Donyapour* https://orcid.org/0000-0002-2269-722X
*Matthew Hirn* https://orcid.org/0000-0003-0290-4292
*Alex Dickson* https://orcid.org/0000-0002-9640-1380

## REFERENCES

[1] C. A. Lipinski, F. Lombardo, B. W. Dominy, P. J. Feeney, *Adv. Drug Delivery Rev.* **1997**, *23*, 3.
[2] Y. Kwon, *Handbook of Essential Pharmacokinetics, Pharmacodynamics and Drug Metabolism for Industrial Scientists*, Springer Science & Business Media, New York **2001**.
[3] Y. Ran, S. H. Yalkowsky, *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 354.
[4] S. H. Yalkowsky, S. C. Valvani, *J. Pharm. Sci.* **1980**, *69*, 912.
[5] T. Ryckmans, M. P. Edwards, V. A. Horne, A. M. Correia, D. R. Owen, L. R. Thompson, I. Tran, M. F. Tutt, T. Young, *Bioorg. Med. Chem. Lett.* **2009**, *19*, 4406.
[6] S. K. Bhal, K. Kassam, I. G. Peirson, G. M. Pearl, *Mol. Pharmaceutics* **2007**, *4*, 556.
[7] Z. Chen, S. G. Weber, *Anal. Chem.* **2007**, *79*, 1043.
[8] J. Sangster, *Octanol-water Partition Coefficients: Fundamentals and Physical Chemistry*, Vol. 1, John Wiley & Sons, Chichester **1997**.
[9] T. Cheng, Y. Zhao, X. Li, F. Lin, Y. Xu, X. Zhang, Y. Li, R. Wang, L. Lai, *J. Chem. Inf. Model.* **2007**, *47*, 2140.
[10] A. K. Ghose, G. M. Crippen, *J. Comput. Chem.* **1986**, *7*, 565.
[11] A. J. Leo, *Chem. Rev.* **1993**, *93*, 1281.
[12] W. M. Meylan, P. H. Howard, *J. Pharm. Sci.* **1995**, *84*, 83.
[13] J. Plante, S. Werner, *J. Cheminf.* **2018**, *10*, 61.
[14] L. Molnár, G. M. Keserű, Á. Papp, Z. Gulyás, F. Darvas, *Bioorg. Med. Chem. Lett.* **2004**, *14*, 851.
[15] J. J. Huuskonen, D. J. Livingstone, I. V. Tetko, *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 947.
[16] I. Moriguchi, S. Hirono, Q. Liu, I. Nakagome, Y. Matsushita, *Chem. Pharm. Bull.* **1992**, *40*, 127.
[17] A. Daina, O. Michielin, V. Zoete, *J. Chem. Inf. Model.* **2014**, *54*, 3284.
[18] D. Chen, Q. Wang, Y. Li, Y. Li, H. Zhou, Y. Fan, *Chemosphere* **2020**, *247*, 125869.
[19] R. Mannhold, G. I. Poda, C. Ostermann, I. V. Tetko, *J. Pharm. Sci.* **2009**, *98*, 861.
[20] I. V. Tetko, V. Y. Tanchuk, A. E. Villa, *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 1407.
[21] I. V. Tetko, V. Y. Tanchuk, *J. Chem. Inf. Comput. Sci.* **2002**, *42*, 1136.
[22] A. Predictor, Inc., Lancaster, CA. **2009**.
[23] T. ChemSilico LLC, Cslogp program, http://www.chemsilico.com/CS_prLogP/LPhome.html
[24] Silicos-it, Filter-it software, http://silicos-it.be.s3-website-eu-west-1.amazonaws.com/software/filter-it/1.0.2/filter-it.html
[25] K. Wu, Z. Zhao, R. Wang, G.-W. Wei, *J. Comput. Chem.* **2018**, *39*, 1444.
[26] M. Korshunova, B. Ginsburg, A. Tropsha, O. Isayev, *J. Chem. Inf. Model.* **2021**, *61*, 7.
[27] J.-W. Zou, W.-N. Zhao, Z.-C. Shang, M.-L. Huang, M. Guo, Q.-S. Yu, *J. Phys. Chem. A* **2002**, *106*, 11550.
[28] M. Riquelme, E. Vöhringer-Martinez, *J. Comput.-Aided Mol. Des.* **2020**, *34*, 327.
[29] H.-F. Chen, *Chem. Biol. Drug Des.* **2009**, *74*, 142.
[30] E. W. Lowe, M. Butkiewicz, M. Spellings, A. Omlor, J. Meiler, Paper presented at: 2011 IEEE Symposium on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB), IEEE, **2011**, p. 1.
[31] Q. Liao, J. Yao, S. Yuan, *Mol. Diversity* **2006**, *10*, 301.
[32] A. Breindl, B. Beck, T. Clark, R. C. Glen, *Mol. Model. Annu.* **1997**, *3*, 142.
[33] M. Popova, O. Isayev, A. Tropsha, *Sci. Adv.* **2018**, *4*, eaap7885.
[34] D.-S. Cao, Q.-S. Xu, Q.-N. Hu, Y.-Z. Liang, *Bioinformatics* **2013**, *29*, 1092.
[35] A. Lusci, G. Pollastri, P. Baldi, *J. Chem. Inf. Model.* **2013**, *53*, 1563.
[36] E. N. Feinberg, D. Sur, Z. Wu, B. E. Husic, H. Mai, Y. Li, S. Sun, J. Yang, B. Ramsundar, V. S. Pande, *ACS Cent. Sci.* **2018**, *4*, 1520.
[37] P. Gao, J. Zhang, J. Sun, J. Yu, *Phys. Chem. Chem. Phys.* **2020**, *22*, 23766.

[38] D. K. Duvenaud, D. Maclaurin, J. Iparraguirre, R. Bombarell, T. Hirzel, A. Aspuru-Guzik, R. P. Adams, *Advances in Neural Information Processing Systems*. 28, Curran Associates, Inc., Red Hook, NY **2015**, p. 2224.

[39] K. Yang, K. Swanson, W. Jin, C. Coley, P. Eiden, H. Gao, A. Guzman-Perez, T. Hopper, B. Kelley, M. Mathea, A. Palmer, V. Settels, T. Jaakkola, K. Jensen, R. Barzilay, *J. Chem. Inf. Model.* **2019**, *59*, 3370.

[40] J. Ma, R. P. Sheridan, A. Liaw, G. E. Dahl, V. Svetnik, *J. Chem. Inf. Model.* **2015**, *55*, 263.

[41] J. S. Smith, O. Isayev, A. E. Roitberg, *Chem. Sci.* **2017**, *8*, 3192.

[42] F. Gao, G. Wolf, M. Hirn, International Conference on Machine Learning, **2019**, p. 2122.

[43] P. Howard, W. Meylan, *Physical/Chemical Property Database (Physprop)*. New York **1999**.

[44] K. Vanommeslaeghe, A. D. MacKerell Jr.., *J. Chem. Inf. Model.* **2012**, *52*, 3144.

[45] K. Vanommeslaeghe, E. P. Raman, A. D. MacKerell Jr.., *J. Chem. Inf. Model.* **2012**, *52*, 3155.

[46] J. A. Maier, C. Martinez, K. Kasavajhala, L. Wickstrom, K. E. Hauser, C. Simmerling, *J. Chem. Theory Comput.* **2015**, *11*, 3696.

[47] D. Vassetti, M. Pagliai, P. Procacci, *J. Chem. Theory Comput.* **2019**, *15*, 1983.

[48] D. A. Case, K. Belfon, I. Y. Ben-Shalom, S. R. Brozell, D. S. Cerutti, T. E. Cheatham, III, V. W. D. Cruzeiro, T. A. Darden, R. E. Duke, G. Giambasu, M. K. Gilson, H. Gohlke, A. W. Goetz, R. Harris, S. Izadi, S. A. Izmailov, K. Kasavajhala, A. Kovalenko, R. Krasny, T. Kurtzman, T. S. Lee, S. LeGrand, P. Li, C. Lin, J. Liu, T. Luchko, R. Luo, V. Man, K. M. Merz, Y. Miao, O. Mikhailovskii, G. Monard, H. Nguyen, A. Onufriev, F. Pan, S. Pantano, R. Qi, D. R. Roe, A. Roitberg, C. Sagui, S. Schott-Verdugo, J. Shen, C. Simmerling, N. R. Skrynnikov, J. Smith, J. Swails, R. C. Walker, J. Wang, L. Wilson, R. M. Wolf, X. Wu, Y. Xiong, Y. Xue, D. M. York and P. A. Kollman, AMBER 2018, University of California, San Francisco, **2018**.

[49] J. Behler, M. Parrinello, *Phys. Rev. Lett.* **2007**, *98*, 146401.

[50] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai and S. Chintala, *Advances in Neural Information Processing Systems*. 32, Curran Associates, Inc., Red Hook, NY **2019**, p. 8024.

[51] W. Hamilton, Z. Ying, J. Leskovec, *Advances in Neural Information Processing Systems*. 30, Curran Associates, Inc., Red Hook, NY **2017**, p. 1024.

[52] 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 É. Duchesnay, *J. Mach. Learn. Res.* **2011**, *12*, 2825.

[53] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, arXiv preprint arXiv: 1312.6203, **2013**.

[54] M. Henaff, J. Bruna, Y. LeCun, arXiv preprint arXiv:1506.05163, **2015**.

[55] T. N. Kipf, M. Welling, 5th International Conference on Learning Representations, **2017**.

[56] H. Altae-Tran, B. Ramsundar, A. S. Pappu, V. Pande, *ACS Cent. Sci.* **2017**, *3*, 283.

[57] J. Wang, R. M. Wolf, J. W. Caldwell, P. A. Kollman, D. A. Case, *J. Comput. Chem.* **2004**, *25*, 1157.

[58] G. E. Hinton, S. T. Roweis, *Advances in Neural Information Processing Systems*. 15, MIT Press, Cambridge, MA **2002**, p. 857.

[59] N. Kokhlikyan, V. Miglani, M. Martin, E. Wang, B. Alsallakh, J. Reynolds, A. Melnikov, N. Kliushkina, C. Araya, S. Yan and O. Reblitz-Richardson, arXiv preprint arXiv:2009.07896, **2020**.

[60] S. Kullback, R. A. Leibler, *Ann. Math. Stat.* **1951**, *22*, 79.

[61] I. J. Hidalgo, T. J. Raub, R. T. Borchardt, *Gastroenterology* **1989**, *96*, 736.

[62] R. D. Brown, Y. C. Martin, *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 1.

[63] G. B. McGaughey, R. P. Sheridan, C. I. Bayly, J. C. Culberson, C. Kreatsoulas, S. Lindsley, V. Maiorov, J.-F. Truchon, W. D. Cornell, *J. Chem. Inf. Model.* **2007**, *47*, 1504.

[64] D. P. Kingma, J. Ba, 3rd International Conference on Learning Representations, **2015**.

## SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of this article.