# Simplicial Message Passing for Chemical Property Prediction

**Hai Lan**
Fujian Institute of Research on the Structure of Matter
Chinese Academy of Sciences
lanhai09@mails.ucas.ac.cn

**Xian Wei**[*]
Software Engineering Department
East China Normal University
xian.wei@tum.de

## Abstract

Recently, message-passing Neural networks (MPNN) provide a promising tool for dealing with molecular graphs and have achieved remarkable success in facilitating the discovery and materials design with desired properties. However, the classical MPNN methods also suffer from a limitation in capturing the strong topological information hidden in molecular structures, such as nonisomorphic graphs. To address this problem, this work proposes a Simplicial Message Passing (SMP) framework to better capture the topological information from molecules, which can break through the limitation within the vanilla message-passing paradigm. In SMP, a generalized message-passing framework is established for aggregating the information from arbitrary-order simplicial complex, and a hierarchical structure is elaborated to allow information exchange between different order simplices. We apply the SMP framework within deep learning architectures for quantum-chemical properties prediction and achieve state-of-the-art results. The results show that compared to traditional MPNN, involving higher-order simplex can better capture the complex structure of molecules and substantially enhance the performance of tasks. The SMP-based model can provide a generalized framework for GNNs and aid in the discovery and design of materials with tailored properties for various applications.

## 1   Introduction

Capturing topological information is crucial for predicting the properties of functional materials[1]. The topological structure of materials, such as molecules or crystals, contains valuable information about their relationships, bond patterns, and atomic arrangements[2; 3]. For example, the topological arrangement of atoms in a crystal or molecule affects its electronic structure and energetics, which plays a vital role in understanding how materials interact with each other and undergo chemical reactions. Topological information determines the physical and chemical properties of materials. By capturing the topological information, the complex system can be represented in a way that preserves the essential geometric and connectivity features[4], which can more efficiently establish property-structure relationships, and facilitate the prediction and design of materials with desired properties.

To capture the topological information, one common approach is to represent molecules as graphs, where atoms are nodes and chemical bonds are edges[5; 6; 7; 8]. Each node can be associated with atom-specific features such as atomic number, hybridization state, or atomic mass. Edge features often include the type of bond, the length of the bond, or the angle of the bond. Compared to molecular fingerprints, such as Extended Connectivity Fingerprints (ECFP) and Simplified Molecular Input

---

[*]Corresponding Author

Line Entry System (SMILES), Graph representation is more general and intuitive[5]. Graph neural networks (GNNs) provide a promising tool for dealing with molecular graphs and have achieved remarkable success in molecular property prediction[9]. As one of the most widely used GNNs, Message Passing Neural Network (MPNN) operates by aggregating features between adjacent nodes[7] and provides a common framework for mainstream GNNs, such as Graph Convolutional Network (GCN)[10] and Graph Attention Network (GAT)[11]. More and more chemists have applied this technology to build end-to-end models for predicting the properties of molecules and designing new materials. However, it has been observed that its topological expressive power is limited by the Weisfeiler-Lehman (WL) isomorphism test[12]. As shown in Fig. 1, there exist two nonisomorphic graphs that have different topological structures, and the WL test generates identical coloring and fails to distinguish them. This reveals MPNN also encounters such insurmountable limitations and affects the performance on the molecular graph.
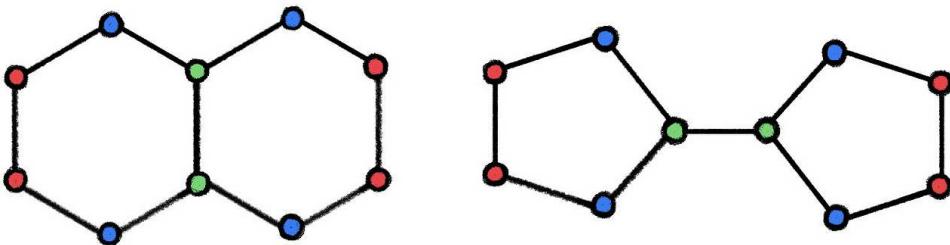


Figure 1: WL graph isomorphism test fails to distinguish two non-isomorphic graphs.

To solve this problem, in this paper, we extend the message-passing paradigm to a higher-order simplicial complex and propose a Simplicial Message Passing (SMP) to better capture topological information from molecular graphs. The main contributions of this paper are as below:

- We propose an SMP framework for aggregating the information from arbitrary order simplex, which can capture the complex interaction and intrinsic topological structure from low-order simplex to high-order simplex.

- We then develop a deep learning model following the SMP framework, called Simplicial Message Passing Neural Network (SMPNN), for quantum chemical properties prediction on both topological structure data and geometric data.

- We perform a series of experiments on quantum chemistry, including molecular properties prediction and molecular dynamics, and achieve state-of-the-art results.

In summary, involving message passing between higher-order simplices is crucial for capturing topological information hidden in complex systems. It provides insights into structural connectivity, chemical reactivity, electronic behavior, stability, symmetry, and property-structure relationships. Incorporating topological considerations allows for more accurate and comprehensive predictions, aiding in the discovery and design of materials with tailored properties for various applications. We believe this work not only equips chemists with a powerful tool for discovering new materials, but also provides data scientists and AI developers with a promising perspective and direction.

## 2 Related Work

### 2.1 Deep Learning on Chemical Property Prediction

Machine learning algorithms are widely used to construct empirical potential and make increased efforts to improve computing accuracy and generalization capacity[1; 13; 14]. Among numerous machine learning models, neural networks (NN) become the first choice for Potential Energy Surface (PES) calculation due to their powerful nonlinear fitting ability. Most research works focused on designing an effective molecular descriptor for the input of NN. Behler and Parrinello proposed a generalized neural network framework in which a symmetric function was designed to handle the

input with different sizes[15], and this method achieved *ab initio* accuracy while it was several orders of magnitude faster. With further exploration of invariance regarding the particular molecule permutation group, Guo et al. proposed PIP-NN method[16] which used the Permutation Invariant Polynomial (PIP)[17] as the input of the neural network and accurately reproduced the analytical potential energy surfaces for $H + H_2$ and $Cl + H_2$ systems. To reduce the number of polynomial inputs in PIP-NN, Fundamental Invariant Neural Network (FI-NN)[18; 19] introduced an efficient representation of molecular permutation symmetry and broadened the application range for larger molecular systems. Deep Potential (DP)[20], a simple deep neural network representation of PES, eliminated ad hoc approximation of the above work by taking each atom into a single sub-network while respecting the symmetry of the system. Embedded Atom Neural Network (EANN)[21] replaced the scalar embedded atom density in empirical Embedded Atom Method (EAM)[22] with a Gaussian-type orbital based density vector which was sent into NN to obtain atomic energy. Gaussian Moment Neural Network[23] used an extendable invariant local molecular descriptor constructed from geometric moments as input of NN to get a high accuracy and efficiency model.

From the perspective of AI developers and data scientists, this hand-crafted input of NN is less intuitive and violates the end-to-end paradigm of deep learning. Therefore, Graph is adopted to encode the molecules and Graph Neural Network (GNN) is used as the backbone to follow the data-driven manner[24; 25; 26; 27]. Deep Tensor Neural Network (DTNN)[28] treated a molecule as a graph, and the molecular structure was encoded with an inter-atomic distance matrix and an atomic number vector so that the model could capture the interaction between atoms in a pairwise manner. A general framework named Message Passing Neural Networks (MPNN)[7] was distilled by reformulating existing models to exploit extra variations. Substantial refinement had been made on the basis of MPNN to improve the performance while reducing the model size and inference time, such as Polarizable Atom Interaction Neural Network (PAINN)[29] and SpookyNet[8]. On the other hand, following the success of the attention mechanism in Natural Language Processing (NLP) and Computer Vision (CV), Graph Attention Network[11] introduced the attention mechanism into GNN, which made it possible to distinguish the importance of neighbors and its own nodes during feature aggregation, rather than the averaging representation of adjacent nodes in MPNN. Many works had involved graph attention mechanism for molecular representation and demonstrated that the attention mechanism could effectively extract the nonlocal intramolecular interactions[30; 31; 32; 33].

## 2.2 Higher-Order Interaction in Complex System

Simplicial complexes and hypergraphs are natural candidates to equip a mathematical framework for describing group interactions in complex systems[34; 35]. Early research of random simplicial complexes, Linial–Meshulam model, is simply a higher-order refinement based on the Erdos Renyi (ER) model[36]. In this model, the connected graph of $n$ nodes is used as an initialization of m triangles which are formed by three nodes. Many variants based on this model focus on the topological data analysis[4]. However, with the blossoming of the research on geometric deep learning[37], the simplicial complex provides a novel theoretical perspective for GNNs. Motivated by Hodge-de Rham theory, Stefania Ebli et al. defined a simplicial convolutional operation to generalize the GNNs to process data living on simplicial complexes[38]. Corman et al. proposed the element-specific persistent homology (ESPH) method to represent 3D complex geometry by one-dimensional topological invariants[39]. The combination of ESPH and deep convolutional neural networks exhibited the favorable potential for retaining important biological information. Moreover, to improve the expressive power of GNNs, which is equivalent to the WL graph isomorphism test. Message-passing frameworks on simplicial complexes and cell complexes were proposed to capture the multi-level interactions presented in many complex systems[40; 41]. Specifically, the authors of [40] redefined four types of adjacent simplices to expand the perception of local structure and enhance expressive power. Different from these approaches, we just generalize the adjacent definition to higher-order simplex without increasing any adjacent type. Therefore, the vanilla message passing framework of the graph can be treated as a special case in our generalization.

## 3 Model Architecture

In this section, we present the SMP framework and the architecture of SMPNN. The later is referred to as a deep learning model in the SMP framework. Firstly, we give a brief introduction to simplicial

complex. Then, we review the message passing framework and describe the main idea of generalizing the concept to arbitrary order simlicial complex. Finally, we describe the SMPNN architecture in detail. The overall structure of the SMPNN is illustrated in Fig. 3.

## 3.1 Preliminary to Simplicial Complex

In this section, we briefly introduce the concept of simplicial complex, which is a set constructed by piecing simplices together. In discrete geometry, simplex is a fundamental concept used to generalize the notion of a triangle or tetrahedron to arbitrary dimensions. For example, the node of a graph can be treated as a 0-simplex, and the edge of a graph can be treated as a 1-simplex.

Moreover, for a $k$-simplex $\sigma^k$, the boundary $\partial\sigma^k$ of $\sigma^k$ is the closure of the set of all simplices $\sigma^{k-1}$. For example, the boundary of a 2-simplex (triangle) $\{v_1, v_2, v_3\}$ is the set that contains all three 1-simplex (edge) $\{v_1, v_2\}$, $\{v_2, v_3\}$, $\{v_3, v_1\}$ of this triangle, and the boundary of a 1-simplex (edge) is the set of two 0-simplex (vertex).

**Definition 1.** *If a $k$-simplex $\sigma^k$ is the boundary of a $(k+1)$-simplex $\tau^{k+1}$, we say $\sigma^k \prec \tau^{k+1}$.*

## 3.2 Message Passing Framework for k-simplex

A common undirected graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ contains the node features $h_v$ and edge features $e_vw$. The message passing function $M_t$ and the update function $U_t$ consist of a complete $t$-th message passing phase which totally runs for $T$ time steps. The features $h_v^t$ of node $v$ in the $t$-th iteration can be written as:

$$m_v^t = \sum_{w \in N(v)} M_{t-1}(h_v^{t-1}, h_w^{t-1}, e_{vw})$$
$$h_v^t = U_t(h_v^{t-1}, m_v^t)$$

(1)

We treat the graph as a simplicial complex and denote the node $\mathcal{V}$ as 0-simplex $\sigma^0$ and edge $\mathcal{E}$ as 1-simplex $\sigma^1$. Similar to how messages are passed along the edges of a graph, the message-passing mechanism in a higher-order simplicial complex involves passing messages along the "edges" between two simplices. For example, as shown in Fig. 2, the message passing between two 0-simplices is along a 1-simplex. The message passing between three 1-simplices is along the 2-simplices, that is the triangle constructed by them. Therefore, we generalize the concept of message passing between simplices as Definition 1.

**Definition 2.** *If a $k$-simplex $\sigma_i^k \prec \tau^{k+1}$, then message passing to $\sigma_i^k$ is along $\tau^{k+1}$ and involves all the $k$-simplex $\sigma^k \prec \tau^{k+1}$.*

According to Definition 2, we can generalize the vanilla message passing framework to SMP and rewrite the equation 1 as:

$$m_{\sigma^k}^t = \sum_{\sigma_i^k \prec \tau^{k+1}} M_{t-1}(h^{t-1}(\sigma_i^k), h^{t-1}(\tau^{k+1}))$$
$$h_{\sigma^k}^t = U_t(h_{\sigma^k}^{t-1}, m_{\sigma^k}^t)$$

(2)

## 3.3 Simplicial Message Passing Neural Network

According to the proposed SMP framework in Sec. 3.2, we elaborate SMP to process the downstream tasks of graphs, such as graph classification, node classification, and graph reconstruction. Instead of just passing messages between neighboring nodes like MPNN, we need to pass messages between neighboring simplices of different orders, referred to as an SMPNN network in the following. The whole SMPNN network can be divided into three steps.

Firstly, according to the specific tasks. Each simplex in the complex is assigned an initial representation or embedding. For example, in a molecular property prediction task, the node feature can be
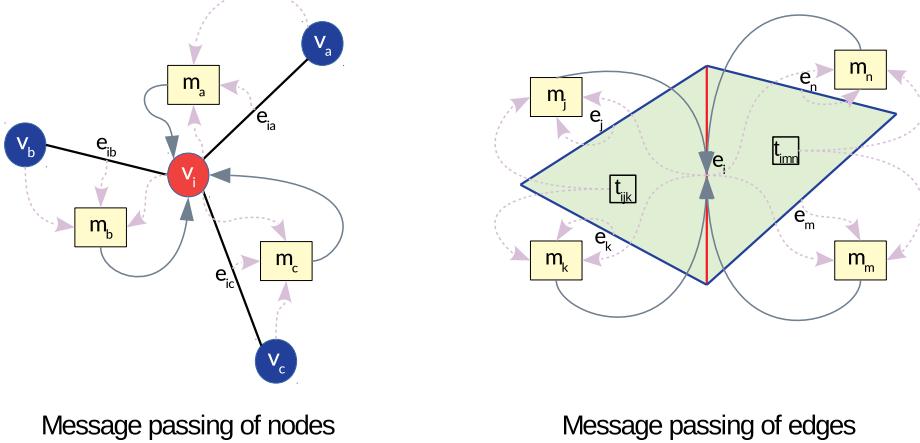
Figure 2: Message passing paradigm for simplices. The left demonstrates the message passing between nodes, and the right demonstrates the message passing between edges.

the atomic number, the edge feature can be the bonding length, the mesh feature can be the area or normal vector of the triangle, and the solid feature can be the volume of the tetrahedron.

Secondly, for each simplex, messages are computed based on the representations of its neighboring simplices. The message computation step involves applying a learnable function or neural network layer that combines the representations of the neighboring simplices in a meaningful way. Then the computed messages from neighboring simplices are aggregated for each simplex. The aggregation process combines the received messages to form a summary representation that captures information from the neighborhood. The aggregated messages are then used to update the representation of each simplex. This update function takes into account the simplex's current representation and the aggregated messages to generate an updated representation. The updated function can be a simple operation like concatenation or a more complex neural network layer. This step is repeated iteratively for a fixed number of steps or until convergence. This mechanism enables us to incorporate knowledge from the local and global neighborhoods of higher-order simplices and can be used for a variety of applications, such as shape analysis, topological data analysis, and protein structure prediction. The application in this paper refers specifically to the molecular property prediction task.

Finally, after iteratively passing and updating messages between simplices, we can capture and propagate information throughout the higher-order simplicial complex. An output layer is elaborated for different tasks. The details of the output layer are discussed in Sec. 3.3.1.

As shown in Fig. 3, we input simplicial complex into the proposed model instead of a graph. Since the molecule has strong 3D structures, we constrain the order of simplex to $2$. The features of simplex will be embedded into high dimensions by an MLP layer to enhance the expressivity. Then, the embedded features are sent into a message passing block. Notably, the $k$-simplicial message passing requires the features of $k + 1$ simplices, which forge a hierarchical structure. After several iterations, for prediction tasks of topological structure data, the features of $k$-simplices $x^k \in \mathbb{R}^{n \times d}$ will be sent into the pooling layer and concat layer to obtain the final representations $x^k \in \mathbb{R}^{d \times 3}$. And for prediction tasks of geometric data, only the features of 0-simplices are sent into the output block.

### 3.3.1 Output Layer

We divide these tasks into two categories, the node classification tasks and the graph classification tasks. And two output blocks are designed for these two different downstream tasks. For graph

classification tasks of chemical property prediction, we input the topological structure of molecules, which is usually encoded by SMILES, and output a scalar result, such as solubility, and drug efficacy. As shown in Fig. 3, the output of message passing block, the features of $k$-simplices will be sent into the pooling layer, respectively. Then all the features will be concatenated and sent into an MLP layer for task prediction. On the other hand, for node classification tasks of chemical property prediction, the input often contains the Cartesian coordinates of every atom in the molecule. The output results, such as potential energy and force, can be represented as a sum of atomic contributions. Hence, we treat these tasks as node classification/regression tasks. The output block for these tasks is only to receive the features of 0-simplice and use an MLP layer to calculate the atomic contribution, then sum up them for the total output result.
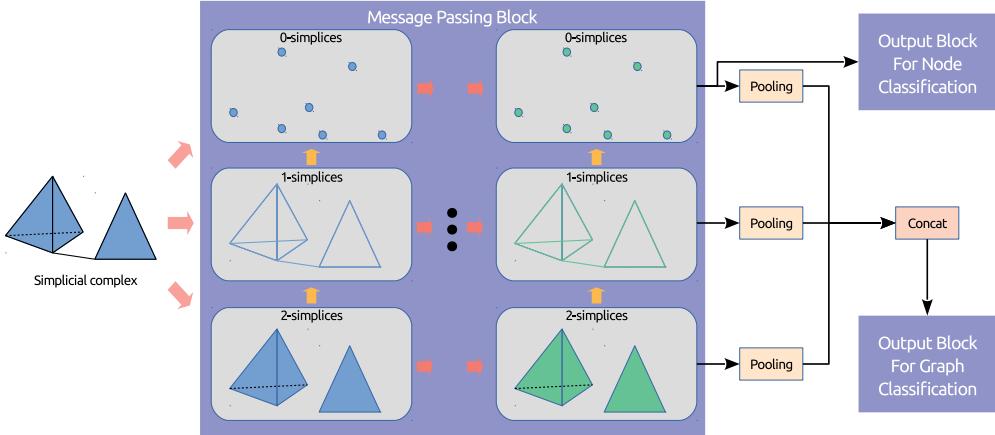


Figure 3: The architecture of proposed SMPNN.

# 4 Experiments

In this section, we evaluate the performance of SMPNN on several popular benchmarks of chemical properties prediction. First, we conduct the experiments on topological structure data in Sec. 4.2. Then, in Sec. 4.3, we evaluate the performances of SMPNN on geometric data, which is also called 3D molecules in [42], including quantum chemical properties prediction and molecular dynamics. Finally, an ablation study of SMPNN is conducted to investigate the contribution of higher-order SMP and the influence of model hyperparameters.

## 4.1 Experimental Configuration

The hardware environment used in the experiment is Core i7-6700 CPU@3.40 GHz, GeForce GTX1080Ti GPU. The algorithm is developed by Python3.8, and PyTorch1.10 is used as the backend. We use ADAM[43] optimizer for model training with an initial learning rate $1 \times 10^{-4}$, which is adjusted by cosine decay schedule[44]. The mini-batch size is set to 100, and the number of SMP blocks is set to 6, the embedding dimension is set to 128.

## 4.2 Chemical Property Prediction with Topological Data

To verify the effectiveness of SMPNN, we ran several experiments to compare the predictive performance on the solubility of our work and the popular benchmarks. These experiments used five public datasets, Delaney[45], Huuskonen[46], OCHEM[47], Tang[48], and Cui[49]. All the datasets were using SMILES as input, and we converted the SMILES into graphs in advance. The experimental results are shown in Tab. 1. The results show that our model achieves state-of-the-art performance in all the tasks. This confirms the better capacity of our model on capturing the topological structure.

6

Table 1: Mean predictive accuracy of solubility on five public datasets. The unit used in the experiment is log Mol/L.

| Dataset | Num of Sample | GCN | GAT | MPNN | D-MPNN | SMPNN |
|---------|---------------|-----|-----|------|--------|-------|
| Delaney[45] | 1144 | 0.878 | 1.116 | 0.903 | 0.725 | **0.709** |
| Huuskonen[46] | 1297 | 0.917 | 0.967 | 0.723 | 0.714 | **0.687** |
| OCHEM[47] | 1311 | 0.952 | 0.788 | 0.675 | 0.693 | **0.508** |
| Tang[48] | 4200 | 0.719 | 1.027 | 0.704 | 0.669 | **0.566** |
| Cui[49] | 9943 | 1.072 | 1.017 | 0.983 | 1.023 | **0.957** |

## 4.3 Chemical Property Prediction with Geometric Data

There is another type of molecular data in chemical properties prediction, which provides the 3D Cartesian coordinates of every atom in the molecule. Some researchers call it 3D molecular graph while we use the term "Geometric data" to distinguish it from chemical data that only have topological structure. In order to verify the effectiveness of the proposed algorithm for geometric data in this paper, we chose two public datasets, MD17[50] and QM9[51; 52] to conduct experiments. The molecular dynamics (MD17) datasets provide both the energy and atomic forces of eight small organic molecules as well as the corresponding atom coordinates of the thermalized system. The datasets range in size from 150k to nearly 1M conformational geometries. All trajectories are calculated at a temperature of 500 K and a resolution of 0.5 fs. The ground truth data are computed via molecular dynamics simulations using PBE+vdW-TS electronic structure method. In our experiment, all eight kinds of organic molecules were trained with the goal of providing highly accurate energy and force predictions. Another dataset QM9 contains $133,885$ stable small organic molecules with up to 7 heavy atoms (C, O, N, F) and up to 29 atoms in total including H. The ground truth is 12 quantum chemistry target properties, which are calculated by *ab initio* quantum chemistry methods at the $B3LYP/6 - 31G(2df, p)$ level. In the experiment of both datasets, the cutoff radius is set to 5Å.

Table 2: Mean absolute error per molecule of predictions for different target properties of the QM9 dataset using 110k training examples. The lowest error is emphasized in bold. The results from Provably powerful graph networks (PPGN)[24], SchNet[6] and enn-s2s[7] are compared.

| Target | Unit | PPGN | SchNet | enn-s2s | SMP |
|--------|------|------|--------|---------|-----|
| $\epsilon_{homo}$ | $meV$ | 40 | 41 | 43 | **33** |
| $\epsilon_{lumo}$ | $meV$ | 33 | 34 | 37 | **27** |
| $\Delta\epsilon$ | $meV$ | 60 | 63 | 69 | **54** |
| $ZPVE$ | $meV$ | 3.12 | 1.7 | **1.5** | 1.6 |
| $\mu$ | $Debye$ | 0.047 | 0.033 | **0.030** | 0.042 |
| $\alpha$ | $Bohr^3$ | 0.131 | 0.235 | 0.092 | **0.086** |
| $< R2 >$ | $Bohr^2$ | 0.592 | **0.073** | 0.180 | 0.241 |
| $U_0$ | $meV$ | 37 | 14 | 19 | **12** |
| $U$ | $meV$ | 37 | 19 | 19 | **15** |
| $H$ | $meV$ | 36 | **14** | 17 | 14 |
| $G$ | $meV$ | 36 | **14** | 19 | 14 |
| $C_v$ | $cal/molK$ | 0.055 | 0.033 | 0.040 | **0.030** |

We randomly choose 1000 and 50000 molecular configurations as the training set and the remaining data as the test set in the MD17 dataset. Mean Absolute Error (MAE) between predictions and ground truth per molecule is applied as evaluation metrics. Firstly, we compare our model with four benchmark methods (Schnet[6], EANN[21], GMNN[23], Physnet[53]). The SMPNN model was trained on $N = 1000$ and $N = 50000$ samples. The results of experiments, as shown in Tab. 3 and Tab. 4, demonstrate SMPNN performs best or at least equal to other models on 9 out of 16 targets with 50000 training samples and on half targets with 1000 training samples. Notably, for energy calculation, SMPNN achieves state-of-the-art performance on all the organic molecules except Benzene with 50000 training samples. When the number of training samples is decreased to 1000, SMPNN still outperforms others on most molecules datasets, which fully demonstrates the effectiveness of SMPNN in the prediction of chemical properties on geometric data.

Table 3: Mean Absolute Errors for energy and force prediction in kcal/mol and kcal/mol/Å on MD17 dataset, with $N = 1000$. The results provided by GM-sNN[23], EANN[21], SchNet[6] and PhysNet[53] are compared. EANN does not provide results on Benzene molecule.

| | | N=1000 | | | |
|---|---|---|---|---|---|
| | | Schnet | EANN | GM-sNN | SMPNN |
| Benzene | Energy | 0.08 | - | 0.08 | **0.06** |
| | Force | 0.31 | - | **0.21** | 0.34 |
| Toluene | Energy | 0.12 | **0.11** | 0.15 | **0.11** |
| | Force | 0.57 | 0.38 | **0.34** | 0.42 |
| Malonaldehyde | Energy | 0.13 | 0.14 | 0.12 | **0.10** |
| | Force | 0.66 | 0.62 | 0.45 | **0.40** |
| Salicylic acid | Energy | 0.20 | **0.14** | 0.19 | 0.16 |
| | Force | 0.85 | 0.51 | **0.49** | 0.52 |
| Aspirin | Energy | 0.37 | 0.33 | 0.38 | **0.32** |
| | Force | 1.35 | 0.99 | **0.69** | 0.99 |
| Ethanol | Energy | 0.08 | 0.10 | 0.10 | **0.07** |
| | Force | 0.39 | 0.47 | 0.33 | **0.22** |
| Uracil | Energy | 0.14 | **0.11** | 0.12 | **0.11** |
| | Force | 0.56 | 0.35 | **0.33** | 0.36 |
| Naphthalene | Energy | 0.16 | **0.12** | 0.17 | 0.15 |
| | Force | 0.58 | **0.27** | 0.36 | 0.36 |

Table 4: Mean Absolute Errors for energy and force prediction in kcal/mol and kcal/mol/Å on MD17 dataset, with $N = 50000$. The results provided by GM-sNN[23], EANN[21], SchNet[6] and PhysNet[53] are compared. EANN does not provide results on Benzene molecule.

| | | N=50000 | | | |
|---|---|---|---|---|---|
| | | Schnet | PhysNet | GM-sNN | SMP |
| Benzene | Energy | **0.07** | **0.07** | **0.07** | 0.10 |
| | Force | 0.17 | 0.15 | **0.14** | 0.17 |
| Toluene | Energy | 0.09 | 0.10 | 0.14 | **0.07** |
| | Force | 0.09 | **0.03** | 0.10 | 0.05 |
| Malonaldehyde | Energy | 0.08 | **0.07** | 0.12 | **0.07** |
| | Force | 0.08 | **0.04** | 0.08 | **0.04** |
| Salicylic acid | Energy | 0.10 | 0.11 | 0.19 | **0.09** |
| | Force | 0.19 | **0.04** | 0.14 | 0.09 |
| Aspirin | Energy | 0.12 | 0.12 | 0.19 | **0.11** |
| | Force | 0.33 | **0.06** | 0.26 | 0.14 |
| Ethanol | Energy | 0.05 | 0.05 | 0.05 | **0.04** |
| | Force | 0.05 | 0.03 | 0.06 | **0.02** |
| Uracil | Energy | **0.10** | **0.10** | **0.10** | **0.10** |
| | Force | 0.11 | **0.03** | 0.07 | 0.04 |
| Naphthalene | Energy | 0.11 | 0.12 | 0.13 | **0.08** |
| | Force | 0.11 | **0.04** | 0.13 | 0.05 |

## 4.4 Ablation Study

To further investigate the influence of higher-order simplex on the expressive power of our model, we removed the 2-simplicial and 1-simplicial message-passing block and show the results in Tab. 5. We carried out ablation experiments and selected two organic compounds with the largest and smallest number of samples from MD17 dataset, Malonaldehyde with $993,237$ data samples and Uracil with $133,770$ data samples. We also carried out the experiments on the whole MD17 dataset and QM9 dataset, due to the page limit, the results can be found in the appendix.

As shown in Tab. 5, compared to 0-SMP, which is identical to the traditional MPNN, adding 1-SMP can substantially enhance the performance on force prediction but provides a trivial improvement

on energy prediction, while involving 2-SMP can remarkably boost the accuracy on both tasks. The experimental results indicated that ignoring higher-order interactions, such as the edge and mesh, will restrain the performance of the model. Therefore, we speculated that the intrinsic relationship captured by our model is more proximate to the first principles.

Furthermore, we varied the dimension of the embedding features and the number of layers in the message-passing block to investigate the influence of the width and depth of the model. As shown in Fig. 4(a), we enlarged the model width by increasing the dimension of the embedding features. For a large dataset (Malonaldehyde), the performance improved and then soon saturated, while there was no substantial improvement for a small dataset (Uracil). Meanwhile, Fig. 4(b) shows that increasing the number of message-passing layers can give rise to better performance on both datasets. However, a trade-off should be considered because stacking too many layers will consume more computational resources and training time.

Table 5: Results of removing the higher-order simplex message-passing block. Scores are given by MAE of energy (kcal/mol) and force (kcal/mol/Å) prediction. 0-SMP denotes 0-simplicial message passing, and the same for other cases.

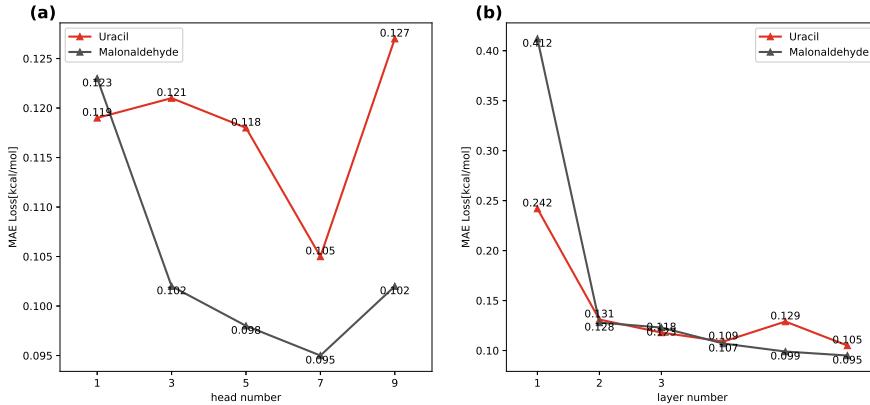| 0-SMP | 1-SMP | 2-SMP | Malonaldehyde | | Uracil | |
|---|---|---|---|---|---|---|
| | | | energy | force | energy | force |
| ✓ | ✗ | ✗ | 0.13 | 0.66 | 0.14 | 0.56 |
| ✓ | ✓ | ✗ | 0.14 | 0.49 | 0.12 | 0.44 |
| ✓ | ✓ | ✓ | 0.10 | 0.44 | 0.11 | 0.36 |



Figure 4: The comparison with width(left) and depth(right) of message-passing block.

## 5   Conclusion

In this study, we established a framework for message-passing in simplicial complex and demonstrated that our work can provide a comprehensive representation of the molecular topology, allowing for a more nuanced analysis of its structural features. The results of experiments showed that, by involving higher-order simplices, our model encompassed both the fine-grained details and the larger-scale patterns, thereby capturing crucial aspects of molecular connectivity and arrangement. However, the introduction of high-order simplex often consumes more computing resources, and the success of our work is initially verified by the empirical results of experiments. In future work, we develop more efficient methods for optimizing the algorithm and delve into the theoretical analysis. This will provide new explorations for the development of geometric deep learning in chemistry and other related fields.

# References

[1] Sergei Manzhos and Tucker Carrington Jr. Neural network potential energy surfaces for small molecules and reactions. *Chemical Reviews*, 121(16):10187–10217, 2020.

[2] Kenneth Atz, Francesca Grisoni, and Gisbert Schneider. Geometric deep learning on molecular representations. *Nature Machine Intelligence*, 3(12):1023–1032, 2021.

[3] Kevin Ryan, Jeff Lengyel, and Michael Shatruk. Crystal structure prediction via deep learning. *Journal of the American Chemical Society*, 140(32):10158–10168, 2018.

[4] Matthew Kahle et al. Topology of random simplicial complexes: a survey. *AMS Contemp. Math*, 620:201–222, 2014.

[5] David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P Adams. Convolutional networks on graphs for learning molecular fingerprints. *Advances in neural information processing systems (NeurIPS)*, 28, 2015.

[6] Kristof Schütt, Pieter-Jan Kindermans, Huziel Enoc Sauceda Felix, Stefan Chmiela, Alexandre Tkatchenko, and Klaus-Robert Müller. Schnet: A continuous-filter convolutional neural network for modeling quantum interactions. *Advances in neural information processing systems (NeurIPS)*, 30, 2017.

[7] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International conference on machine learning (ICML)*, pages 1263–1272. PMLR, 2017.

[8] Oliver T Unke, Stefan Chmiela, Michael Gastegger, Kristof T Schütt, Huziel E Sauceda, and Klaus-Robert Müller. Spookynet: Learning force fields with electronic degrees of freedom and nonlocal effects. *Nature communications*, 12(1):7273, 2021.

[9] Kristof T Schütt, Huziel E Sauceda, P-J Kindermans, Alexandre Tkatchenko, and K-R Müller. Schnet–a deep learning architecture for molecules and materials. *The Journal of Chemical Physics*, 148(24):241722, 2018.

[10] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. *Advances in neural information processing systems*, 29, 2016.

[11] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations (ICLR)*, 2018.

[12] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.

[13] Bin Jiang, Jun Li, and Hua Guo. High-fidelity potential energy surfaces for gas-phase and gas–surface scattering processes from machine learning. *The Journal of Physical Chemistry Letters*, 11(13):5120–5131, 2020.

[14] Pei-Lin Kang, Cheng Shang, and Zhi-Pan Liu. Large-scale atomic simulation via machine learning potentials constructed by global potential energy surface exploration. *Accounts of Chemical Research*, 53(10):2119–2129, 2020.

[15] Jörg Behler and Michele Parrinello. Generalized neural-network representation of high-dimensional potential-energy surfaces. *Physical review letters*, 98(14):146401, 2007.

[16] Bin Jiang and Hua Guo. Permutation invariant polynomial neural network approach to fitting potential energy surfaces. *The Journal of chemical physics*, 139(5):054112, 2013.

[17] Zhen Xie and Joel M Bowman. Permutationally invariant polynomial basis for molecular energy surface fitting via monomial symmetrization. *Journal of Chemical Theory and Computation*, 6(1):26–34, 2010.

[18] Rongjun Chen, Kejie Shao, Bina Fu, and Dong H. Zhang. Fitting potential energy surfaces with fundamental invariant neural network. ii. generating fundamental invariants for molecular systems with up to ten atoms. *The Journal of Chemical Physics*, 152(20):204307, 2020.

[19] Lu Xiaoxiao, Chenyao Shang, Lulu Li, Rongjun Chen, Bina Fu, Xin Xu, and Dong Zhang. Unexpected steric hindrance failure in the gas phase f- + (ch3)3ci sn2 reaction. *Nature Communications*, 13:4427, 07 2022.

[20] Jiequn Han, Linfeng Zhang, Roberto Car, et al. Deep potential: A general representation of a many-body potential energy surface. *Communications in Computational Physics*, 23(3), 2018.

[21] Yaolong Zhang, Ce Hu, and Bin Jiang. Embedded atom neural network potentials: Efficient and accurate machine learning with a physically inspired representation. *The journal of physical chemistry letters*, 10(17):4962–4967, 2019.

[22] SM Foiles, MI Baskes, and Murray S Daw. Embedded-atom-method functions for the fcc metals cu, ag, au, ni, pd, pt, and their alloys. *Physical review B*, 33(12):7983, 1986.

[23] V. Zaverkin and J. Kästner. Gaussian moments as physically inspired molecular descriptors for accurate and scalable machine learning potentials. *Journal of Chemical Theory and Computation*, 16(8):5410–5421, 2020.

[24] Haggai Maron, Heli Ben-Hamu, Hadar Serviansky, and Yaron Lipman. Provably powerful graph networks. *Advances in neural information processing systems (NeurIPS)*, 32, 2019.

[25] Ke Liu, Xiangyan Sun, Lei Jia, Jun Ma, Haoming Xing, Junqiu Wu, Hua Gao, Yax Sun, Florian Boulnois, and Jie Fan. Chemi-net: a molecular graph convolutional network for accurate drug property prediction. *International journal of molecular sciences*, 20(14):3389, 2019.

[26] Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Junction tree variational autoencoder for molecular graph generation. In *International conference on machine learning (ICML)*, pages 2323–2332. PMLR, 2018.

[27] Thin Nguyen, Hang Le, Thomas P Quinn, Tri Nguyen, Thuc Duy Le, and Svetha Venkatesh. Graphdta: Predicting drug–target binding affinity with graph neural networks. *Bioinformatics*, 37(8):1140–1147, 2021.

[28] Kristof T Schütt, Farhad Arbabzadah, Stefan Chmiela, Klaus R Müller, and Alexandre Tkatchenko. Quantum-chemical insights from deep tensor neural networks. *Nature communications*, 8(1):13890, 2017.

[29] Kristof Schütt, Oliver Unke, and Michael Gastegger. Equivariant message passing for the prediction of tensorial properties and molecular spectra. In *International Conference on Machine Learning (ICML)*, pages 9377–9388. PMLR, 2021.

[30] Zhaoping Xiong, Dingyan Wang, Xiaohong Liu, Feisheng Zhong, Xiaozhe Wan, Xutong Li, Zhaojun Li, Xiaomin Luo, Kaixian Chen, Hualiang Jiang, and Mingyue Zheng. Pushing the boundaries of molecular representation for drug discovery with the graph attention mechanism. *Journal of medicinal chemistry*, 63(16):8749–8760, 2020.

[31] Connor W Coley, Wengong Jin, Luke Rogers, Timothy F Jamison, Tommi S Jaakkola, William H Green, Regina Barzilay, and Klavs F Jensen. A graph-convolutional neural network model for the prediction of chemical reactivity. *Chemical science*, 10(2):370–377, 2019.

[32] Anthony Yu-Tung Wang, Steven K Kauwe, Ryan J Murdock, and Taylor D Sparks. Compositionally restricted attention-based network for materials property predictions. *Npj Computational Materials*, 7(1):1–10, 2021.

[33] Steph-Yves Louis, Yong Zhao, Alireza Nasiri, Xiran Wang, Yuqi Song, Fei Liu, and Jianjun Hu. Graph convolutional neural networks with global attention for improved materials property prediction. *Physical Chemistry Chemical Physics*, 22(32):18141–18148, 2020.

[34] Federico Battiston, Giulia Cencetti, Iacopo Iacopini, Vito Latora, Maxime Lucas, Alice Patania, Jean-Gabriel Young, and Giovanni Petri. Networks beyond pairwise interactions: structure and dynamics. *Physics Reports*, 874:1–92, 2020.

[35] Matthew Kahle. Random geometric complexes. *Discrete & Computational Geometry*, 45:553–573, 2011.

[36] Nathan Linial* and Roy Meshulam*. Homological connectivity of random 2-complexes. *Combinatorica*, 26(4):475–487, 2006.

[37] Michael M Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017.

[38] Stefania Ebli, Michaël Defferrard, and Gard Spreemann. Simplicial neural networks. *arXiv preprint arXiv:2010.03633*, 2020.

[39] Zixuan Cang and Guo-Wei Wei. Topologynet: Topology based deep convolutional and multi-task neural networks for biomolecular property predictions. *PLoS computational biology*, 13(7):e1005690, 2017.

[40] Cristian Bodnar, Fabrizio Frasca, Yuguang Wang, Nina Otter, Guido F Montufar, Pietro Lio, and Michael Bronstein. Weisfeiler and lehman go topological: Message passing simplicial networks. In *International Conference on Machine Learning (ICML)*, pages 1026–1037. PMLR, 2021.

[41] Cristian Bodnar, Fabrizio Frasca, Nina Otter, Yuguang Wang, Pietro Lio, Guido F Montufar, and Michael Bronstein. Weisfeiler and lehman go cellular: Cw networks. *Advances in Neural Information Processing Systems*, 34:2625–2640, 2021.

[42] Yi Liu, Limei Wang, Meng Liu, Yuchao Lin, Xuan Zhang, Bora Oztekin, and Shuiwang Ji. Spherical message passing for 3d molecular graphs. In *International Conference on Learning Representations (ICLR)*, 2022.

[43] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.

[44] Ilya Loshchilov and Frank Hutter. SGDR: stochastic gradient descent with warm restarts. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.

[45] John S. Delaney. Esol: Estimating aqueous solubility directly from molecular structure. *Journal of Chemical Information and Computer Sciences*, 44(3):1000–1005, 2004.

[46] Jarmo Huuskonen. Estimation of aqueous solubility for a diverse set of organic compounds based on molecular topology. *Journal of Chemical Information and Computer Sciences*, 40(3):773–777, 2000.

[47] I. Sushko, S. Novotarskyi, R. Körner, and et al. Online chemical modeling environment (ochem): web platform for data storage, model development and publishing of chemical information. *Journal of Computer-Aided Molecular Design*, 25:533–554, 2011.

[48] Bowen Tang, Skyler T Kramer, Meijuan Fang, Yingkun Qiu, Zhen Wu, and Dong Xu. A self-attention based message passing neural network for predicting molecular lipophilicity and aqueous solubility. *Journal of cheminformatics*, 12(1):1–9, 2020.

[49] Qiuji Cui, Shuai Lu, Bingwei Ni, Xian Zeng, Ying Tan, Ya Dong Chen, and Hongping Zhao. Improved prediction of aqueous solubility of novel compounds by going deeper with deep learning. *Frontiers in oncology*, 10:121, 2020.

[50] Stefan Chmiela, Alexandre Tkatchenko, Huziel E Sauceda, Igor Poltavsky, Kristof T Schütt, and Klaus-Robert Müller. Machine learning of accurate energy-conserving molecular force fields. *Science advances*, 3(5):e1603015, 2017.

[51] Raghunathan Ramakrishnan, Pavlo O Dral, Matthias Rupp, and O Anatole Von Lilienfeld. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific data*, 1(1):1–7, 2014.

[52] Lars Ruddigkeit, Ruud Van Deursen, Lorenz C Blum, and Jean-Louis Reymond. Enumeration of 166 billion organic small molecules in the chemical universe database gdb-17. *Journal of chemical information and modeling*, 52(11):2864–2875, 2012.

[53] Oliver T. Unke and Markus Meuwly. Physnet: A neural network for predicting energies, forces, dipole moments, and partial charges. *Journal of Chemical Theory and Computation*, 15(6):3678–3693, 2019.