

CHEMICAL-REACTION-AWARE MOLECULE REPRESENTATION LEARNING

Hongwei Wang¹, Weijiang Li¹, Xiaomeng Jin¹, Kyunghyun Cho^{2,3}, Heng Ji¹, Jiawei Han¹, Martin D. Burke¹

¹University of Illinois Urbana-Champaign, ²New York University, ³Genentech
{hongweiw, wl13, xjin17, hengji, hanj, mdburke}@illinois.edu, kyunghyun.cho@nyu.edu

1. Introduction and current work

- SMILES (Simplified Molecular-Input Line-Entry System) based LLMs.
- Delicately hand-crafted GNNs.
- Paper proposes using chemical reactions to learn *effective and generalizable* embeddings.
Eg. $\text{CH}_3\text{COOH} + \text{C}_2\text{H}_5\text{OH} \rightarrow \text{CH}_3\text{COOC}_2\text{H}_5 + \text{H}_2\text{O}$

$$\boxed{?} > h_{\text{CH}_3\text{COOH}} + h_{\text{C}_2\text{H}_5\text{OH}} = h_{\text{CH}_3\text{COOC}_2\text{H}_5} + h_{\text{H}_2\text{O}}$$

2. Brief results

- *Well-organized embeddings.*
- Can learn *reaction templates* (like esterification) using sumas READOUT.
- Beats many baselines. (Hit@1, AUC, RMSE across datasets and tasks)

3. Method

- Initial Encoding - “Each atom a_i has an initial feature vector x_i encoding its properties. In this work, we use four types of atom properties: **element type, charge, whether the atom is an aromatic ring, and the count of attached hydrogen atom(s)**. Each type of atom properties is represented as a one-hot vector, and we add an additional “**unknown**” entry for each one-hot vector to handle unknown values during inference. The four one-hot vectors are concatenated as the initial atom feature.”
- Bond-types are ignored! Graph becomes homogeneous.
- Preserving equivalence – (results in *equivalence classes*)

Proposition 1 Let M be the set of molecules, $R \subseteq M$ and $P \subseteq M$ be the reactant set and product set of a chemical reaction, respectively. If $R \rightarrow P \Leftrightarrow \sum_{r \in R} h_r = \sum_{p \in P} h_p$ for all chemical reactions, then “ \rightarrow ” is an equivalence relation on 2^M that satisfies the following three properties: (1) Reflexivity: $A \rightarrow A$, for all $A \in 2^M$; (2) Symmetry: $A \rightarrow B \Leftrightarrow B \rightarrow A$, for all $A, B \in 2^M$; (3) Transitivity: If $A \rightarrow B$ and $B \rightarrow C$, then $A \rightarrow C$, for all $A, B, C \in 2^M$.

3. Method

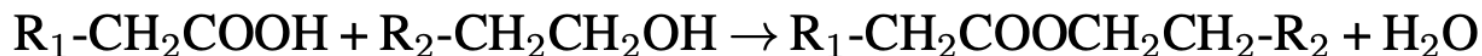
- The last condition is *very strong* leads to robust embeddings.
- *Eg.* in organic synthesis, a target compound \mathbf{t} may be made from three different sets of starting materials \mathbf{A} , \mathbf{B} , and \mathbf{C} . Then the sets \mathbf{A} , \mathbf{B} , \mathbf{C} as well as $\{\mathbf{t}\}$ belong to *one equivalence class*, and we have

$$\sum_{a \in A} h_a = \sum_{b \in B} h_b = \sum_{c \in C} h_c = h_t$$

- It also ***improves generalizability***.

Reaction Center – ‘The reaction center of $R \rightarrow P$ is defined as a subgraph of R consisting of atoms whose bonds have changed after reaction’.

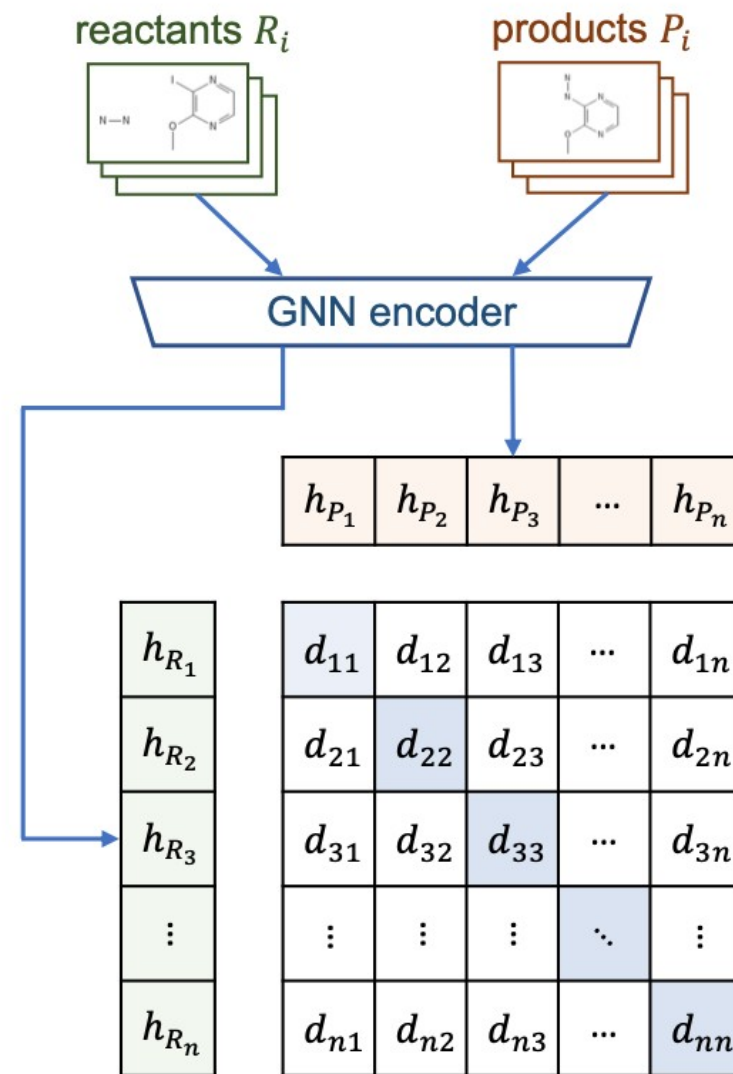
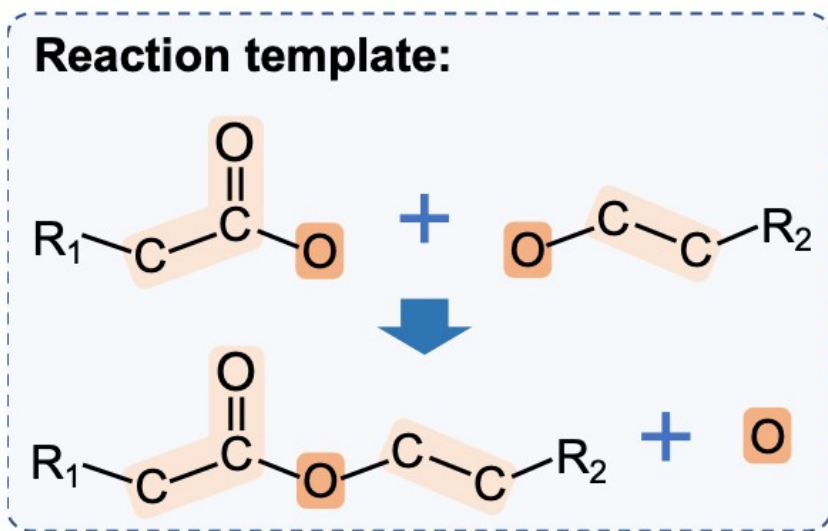
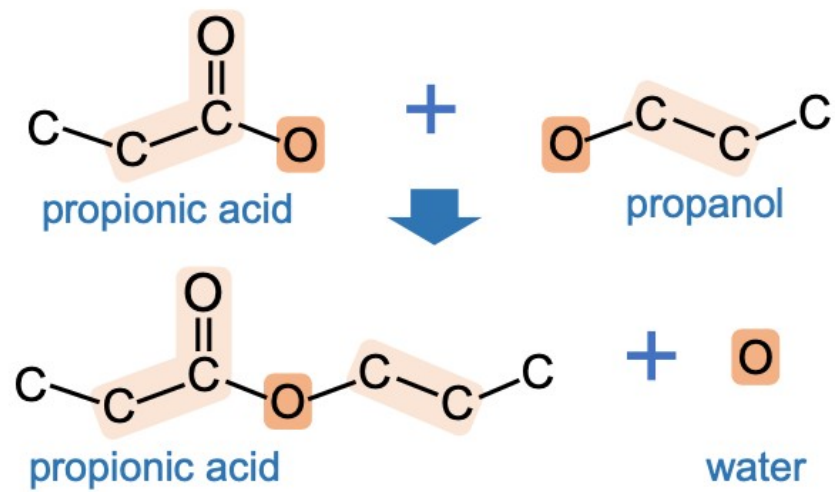
Proposition 2 *Let $R \rightarrow P$ be a chemical reaction where R is the reactant set and P is the product set, and C be its reaction center. Suppose that we use the GNN (whose number of layers is K) shown in Eqs. (1) and (2) as the molecule encoder, and set the READOUT function in Eq. (2) as summation. Then for an arbitrary atom a in one of the reactant whose final representation is h_a^K , the residual term $\sum_{r \in R} h_r - \sum_{p \in P} h_p$ is a function of h_a^K if and only if the distance between atom a and reaction center C is less than K .*



3. Method

- Important implications –
 1. No complicated networks required unlike previous work + *few-shot* adaptation.
 2. Can learn strong properties even though reactions can be unbalanced, eg. H_2O is omitted from all Esterification reactions.
 3. Because of Proposition-2, number of GNN layers is a ***very important factor***.
- Training –

$$L_{\mathcal{B}} = \frac{1}{|\mathcal{B}|} \sum_i \left\| \sum_{r \in R_i} h_r - \sum_{p \in P_i} h_p \right\|_2 + \frac{1}{|\mathcal{B}|(|\mathcal{B}| - 1)} \sum_{i \neq j} \max \left(\gamma - \left\| \sum_{r \in R_i} h_r - \sum_{p \in P_j} h_p \right\|_2, 0 \right),$$



4. Experiments

- Reaction Product Prediction - USPTO dataset

Metrics	MRR	MR	Hit@1	Hit@3	Hit@5	Hit@10
Mol2vec	0.681	483.7	0.614	0.725	0.759	0.798
Mol2vec-FT1	0.688 ± 0.000	417.6 ± 0.1	0.620 ± 0.000	0.734 ± 0.000	0.767 ± 0.000	0.806 ± 0.000
MolBERT	0.708	460.7	0.623	0.768	0.811	0.858
MolBERT-FT1	0.731 ± 0.000	457.9 ± 0.0	0.649 ± 0.000	0.790 ± 0.000	0.831 ± 0.000	0.873 ± 0.000
MolBERT-FT2	0.776 ± 0.000	459.6 ± 0.2	0.708 ± 0.000	0.827 ± 0.000	0.859 ± 0.000	0.891 ± 0.000
MolR-GCN	0.905 ± 0.001	34.5 ± 2.4	0.867 ± 0.001	0.938 ± 0.001	0.950 ± 0.001	0.961 ± 0.002
MolR-GAT	0.903 ± 0.002	35.3 ± 2.8	0.864 ± 0.002	0.935 ± 0.003	0.948 ± 0.003	0.961 ± 0.003
MolR-SAGE	0.903 ± 0.004	53.0 ± 4.6	0.865 ± 0.005	0.935 ± 0.004	0.948 ± 0.004	0.961 ± 0.002
MolR-TAG	0.918 ± 0.000	27.4 ± 0.4	0.882 ± 0.000	0.949 ± 0.001	0.960 ± 0.001	0.970 ± 0.000
MolR-TAG (1% training data)	0.904 ± 0.002	33.0 ± 3.7	0.865 ± 0.003	0.937 ± 0.003	0.951 ± 0.002	0.963 ± 0.002

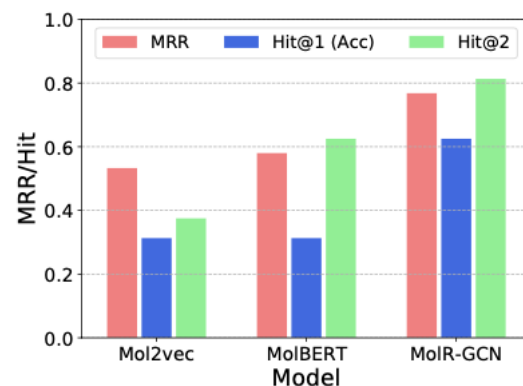


Figure 2: Result of answering real multi-choice questions on product prediction.

4. Experiments

- Case study on first 20 reactions

No.	Reactant(s)	Ground-truth product	Predicted product by MolR-GCN	Predicted product by Mol2vec	Predicted product by MolBERT
5					
			(No. 32353)	(No. 32353)	(No. 32353)
6			Same as ground-truth		
				(No. 39181)	(No. 24126)
8			Same as ground-truth		
				(No. 11233)	(No. 17526)

10			(No. 18889)	(No. 18889)	(No. 18889)
13			Same as ground-truth	Same as ground-truth	
					(No. 37024)
16			Same as ground-truth		Same as ground-truth
					(No. 21947)
17			Same as ground-truth		
					(No. 2029)
					(No. 22247)

4. Experiments

- Molecule Property Prediction - 5 datasets *BBBP*, *HIV*, *BACE*, *Tox21*, and *ClinTox*

Datasets	BBBP	HIV	BACE	Tox21	ClinTox
SMILES-Transformers	0.704	0.729	0.701	0.802	0.954
ECFP4	0.729	<u>0.792</u>	<u>0.867</u>	0.822	0.799
GraphConv	0.690	0.763	0.783	<u>0.829</u>	0.807
Weave	0.671	0.703	0.806	0.820	0.832
ChemBERTa	0.643	0.622	-	0.728	0.733
D-MPNN	0.708	0.752	-	0.688	0.906
CDDD	0.761 \pm 0.00	0.753 \pm 0.00	0.833 \pm 0.00	-	-
MolBERT	0.762 \pm 0.00	0.783 \pm 0.00	0.866 \pm 0.00	-	-
GraphCL	0.695 \pm 0.005	0.776 \pm 0.009	0.782 \pm 0.012	0.754 \pm 0.009	0.701 \pm 0.019
GraphLoG	0.725 \pm 0.008	0.778 \pm 0.008	0.835 \pm 0.012	0.757 \pm 0.005	0.767 \pm 0.033
Mol2vec	<u>0.872</u> \pm 0.021	0.769 \pm 0.021	0.862 \pm 0.027	0.803 \pm 0.041	0.841 \pm 0.062
MolR-GCN	0.890 \pm 0.032	0.802 \pm 0.024	0.882 \pm 0.019	0.818 \pm 0.023	0.916 \pm 0.039
MolR-GAT	0.887 \pm 0.026	0.794 \pm 0.022	0.863 \pm 0.026	0.839 \pm 0.039	0.908 \pm 0.039
MolR-SAGE	0.879 \pm 0.032	0.793 \pm 0.026	0.859 \pm 0.029	0.811 \pm 0.039	0.890 \pm 0.058
MolR-TAG	0.895 \pm 0.031	0.801 \pm 0.023	0.875 \pm 0.023	0.820 \pm 0.028	0.913 \pm 0.043

4. Experiments

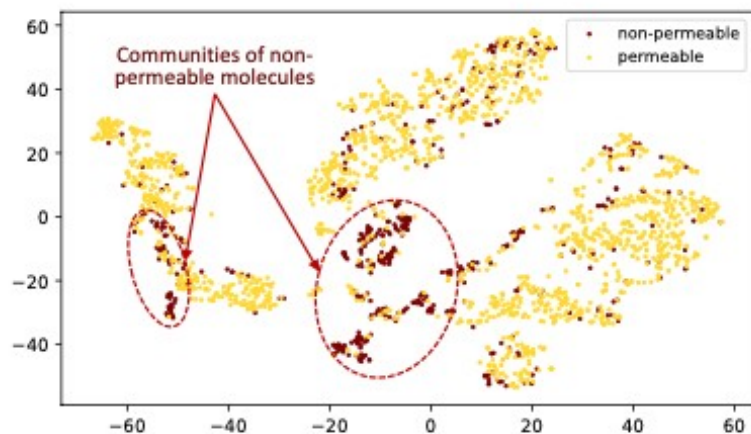
- Chemical Reaction Classification – USPTO-1k-TPL dataset

Methods	Accuracy
RXNFP	0.989
AP3-256-5NN	0.295
AP3-256-MLP	0.809
DRFP-5NN	0.917
DRFP-MLP	0.977
MolR-GCN	0.931 ± 0.022
MolR-GAT	0.930 ± 0.017
MolR-SAGE	0.936 ± 0.025
MolR-TAG	0.962 ± 0.028

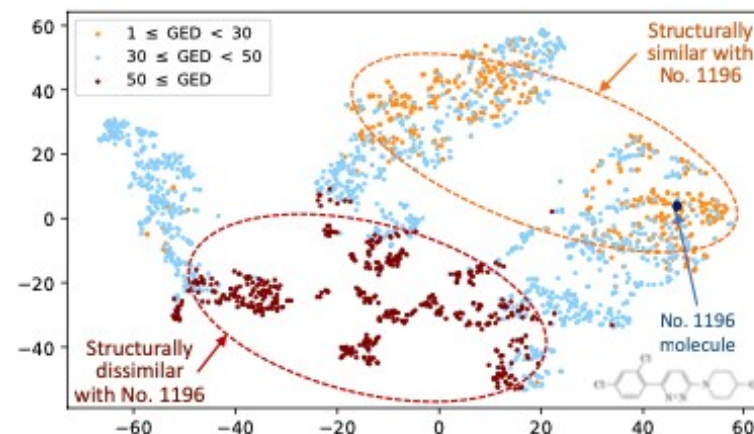
- Graph Edit Distance Prediction – subset of the QM9 dataset.

Feature mode	Concat	Subtract
Mol2vec	1.140 ± 0.041	0.995 ± 0.034
MolBERT	1.127 ± 0.042	0.937 ± 0.029
MolR-GCN	0.976 ± 0.026	0.922 ± 0.019
MolR-GAT	1.007 ± 0.021	0.943 ± 0.016
MolR-SAGE	0.918 ± 0.028	0.817 ± 0.013
MolR-TAG	0.960 ± 0.027	0.911 ± 0.027

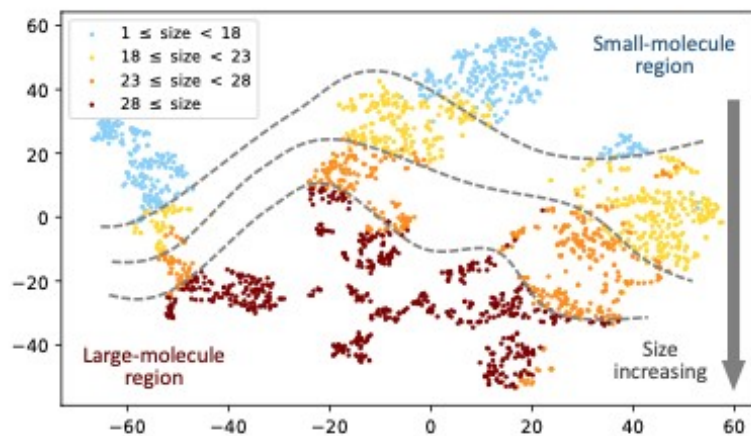
5. Visualization of embeddings (using t-SNE)



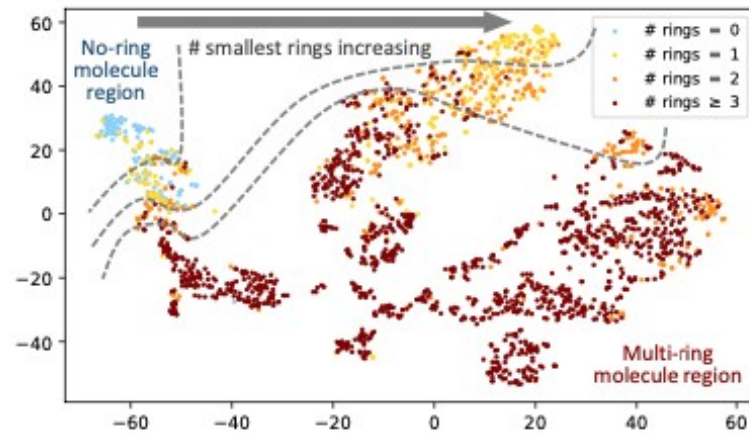
(a) Molecule property



(b) GED w.r.t. No. 1196 molecule



(c) Molecule size



(d) # smallest rings