

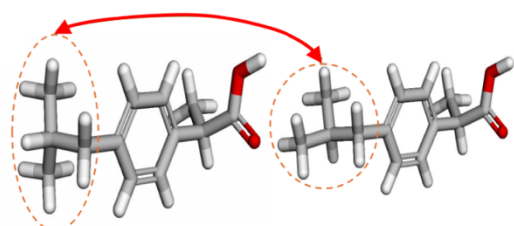
Structure-Aware Fusion with Progressive Injection for Multimodal Molecular Representation Learning

Authors: Zihao Jing, Yan Sun, Yan Yi Li, Sugitha Janarthanan, Alana Deng, Pingzhao Hu

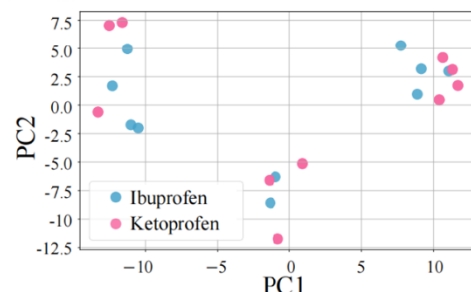
Presenter: Zihao Jing



Problem & Challenges



(a) Two conformers show local 3D variation



(b) Conformers reveals PCA embedding instability

Conformer instability and modality collapse motivate SFP + PI.

- **3D conformers are unstable**
→ **unreliable fusion**

Valid conformers vary by sampler/seed; early 3D-dependent fusion overfits these artifacts and yields inconsistent predictions.

- **Naive/symmetric fusion** → **modality collapse**

Equal, everywhere fusion lets high-variance 3D dominate, drowning SMILES semantics and hurting robustness.

Key Ideas / Contributions

Robust fusion of 2D+3D structure with SMILES, injected progressively into a sequence backbone—strong gains across 29 tasks.

- **Structured Fusion Pipeline (SFP)**: align/encode 2D topology + 3D geometry into a stable structural prior.
- **Progressive Injection (PI)**: asymmetrically inject that prior into the SMILES stream at later layers to avoid collapse.
- **Results**: avg +2.7% over best baseline across broad benchmarks; up to 27% gain on LD50.

Code Link



Paper Link

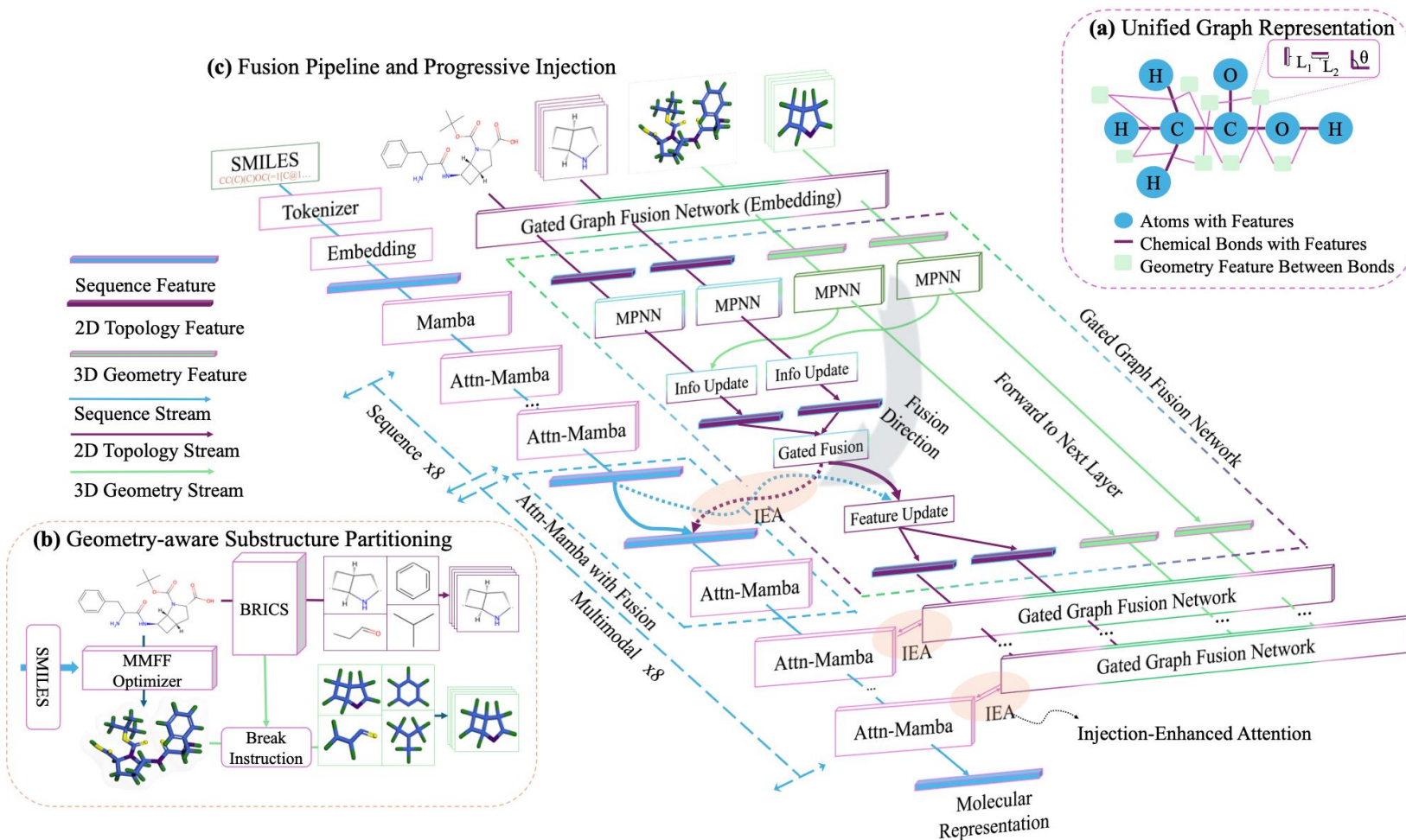


Contact Information

Zihao Jing, zjing29@uwo.ca

Pingzhao Hu, phu49@uwo.ca

Method Overview

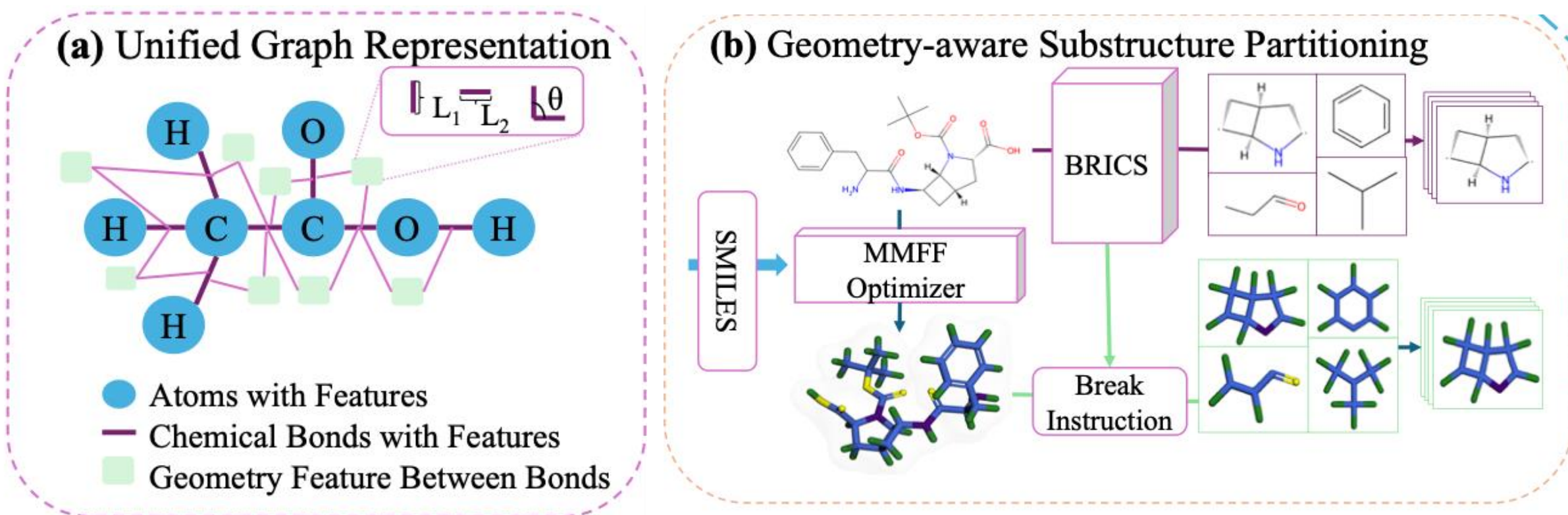


Two streams:

- Structural-fusion stream: 2D+3D fused → unified structural prior, propagated independently.
- Sequence stream: SMILES tokens modeled first; prior injected later via dedicated attention.

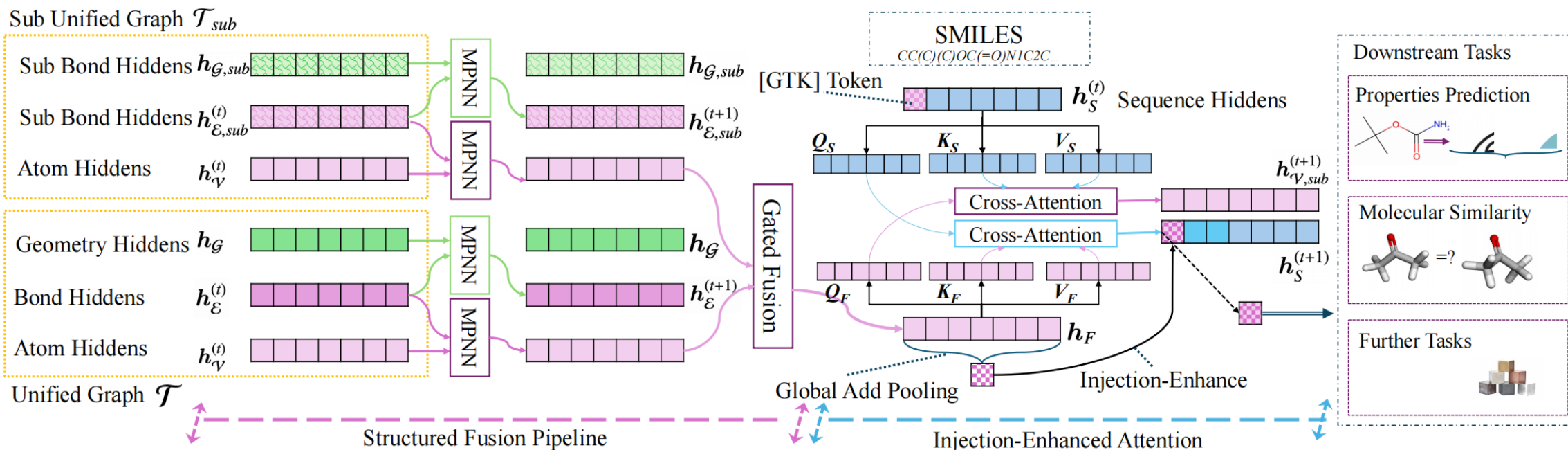
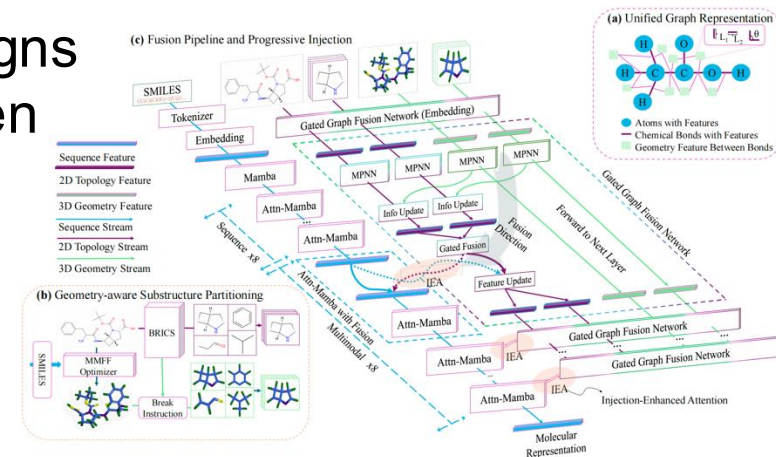
Novelty 1: Structured Fusion Pipeline (SFP)

- **Structural Unified Graph:** nodes (atoms), bonds, and auxiliary geometric linkages (lengths & angles) create a rotation-invariant joint structure.
- **Geometry-Aware Substructure Partitioning:** extend BRICS to 3D; fuse global + local features with gated fusion to form the **structural prior**.



Novelty 2: Progressive Injection (PI)

- **Injection Enhanced Attention:** bidirectional cross-attention aligns structure \leftrightarrow sequence, then **pooled prior** updates a global token [GTK]; prior continues evolving in state-space layers (Mamba) across depth.
- **Why progressive?** Early symmetric fusion distorts semantics; delayed injection yields better stability and convergence.



Benchmarks

MODELS	BBB	HIA	PGP	BIOAV.	Tox-Avg.	CYP-Avg.	Top2Cnt/10
TDC DATASETS - CLASSIFICATION - AUROC \uparrow							
ATTENTIVEFP	0.855 _{0.011}	0.974 _{0.007}	0.892 _{0.012}	0.632 _{0.039}	0.842 _{0.010}	0.749 _{0.008}	0
FPGNN	0.888 _{0.018}	0.958 _{0.012}	0.930 _{0.007}	0.666 _{0.035}	0.860 _{0.017}	0.866 _{0.004}	4
DMPNN	0.864 _{0.010}	0.976 _{0.004}	0.889 _{0.005}	0.617 _{0.050}	0.821 _{0.019}	0.819 _{0.004}	2
ATTRMASKING	0.892 _{0.012}	0.978 _{0.006}	0.929 _{0.006}	0.577 _{0.087}	0.846 _{0.021}	0.817 _{0.005}	4
CONTEXTPred	0.897 _{0.004}	0.975 _{0.004}	0.923 _{0.005}	0.671 _{0.026}	0.818 _{0.017}	0.827 _{0.003}	1
TRANFOXmol	0.868 _{0.019}	0.951 _{0.036}	0.875 _{0.011}	0.619 _{0.019}	0.837 _{0.017}	0.860 _{0.006}	0
DEEPMOL	0.774 _{0.023}	0.880 _{0.012}	0.821 _{0.007}	0.509 _{0.026}	0.735 _{0.015}	0.770 _{0.008}	0
MuMo	0.899 _{0.014}	0.979 _{0.013}	0.942 _{0.019}	0.714 _{0.021}	0.840 _{0.015}	0.880 _{0.017}	7
MODELS	BACE-R	BACE-S	BBBP-R	BBBP-S	CLINTOX	SIDER	TOX21
MOLECULENET - CLASSIFICATION - AUROC \uparrow							
FPGNN	0.831 _{0.011}	0.831 _{0.011}	0.904 _{0.020}	0.892 _{0.019}	0.732 _{0.068}	0.661 _{0.014}	0.833 _{0.004}
TRANSFOXmol	0.780 _{0.032}	0.780 _{0.032}	0.907 _{0.024}	0.881 _{0.015}	0.830 _{0.047}	0.636 _{0.022}	0.816 _{0.011}
CHEMBERTA-2	0.848 _{0.037}	0.848 _{0.037}	0.932 _{0.037}	0.892 _{0.019}	0.933 _{0.054}	0.708 _{0.090}	0.809 _{0.029}
MOLFORMER	0.873 _{0.009}	0.833 _{0.009}	0.889 _{0.028}	0.868 _{0.013}	0.888 _{0.044}	0.651 _{0.016}	0.804 _{0.013}
MOLBERT	0.882 _{0.015}	0.832 _{0.015}	0.955 _{0.008}	0.949 _{0.013}	0.875 _{0.041}	-	-
GROVER	0.779 _{0.059}	0.779 _{0.059}	0.849 _{0.008}	0.823 _{0.020}	0.685 _{0.066}	0.635 _{0.034}	0.808 _{0.014}
UNI-MOL	0.840 _{0.031}	0.840 _{0.031}	0.889 _{0.025}	0.886 _{0.016}	0.818 _{0.065}	0.666 _{0.021}	0.812 _{0.007}
MuMo	0.878 _{0.046}	0.849 _{0.014}	0.962 _{0.007}	0.957 _{0.011}	0.985 _{0.011}	0.677 _{0.009}	0.834 _{0.009}
MODELS	LD50	CACO-2	PPBR	LIPO	MODELS	ESOL	FREESOLV
TDC DATASETS - REGRESSION - MAE \downarrow				MOLECULENET - REGRESSION - RMSE \downarrow			
ATTENTIVEFP	0.678 _{0.012}	0.401 _{0.032}	9.373 _{0.335}	0.572 _{0.007}	CHEMBERTA-2	0.633 _{0.132}	1.219 _{0.206}
FPGNN	0.638 _{0.024}	0.326 _{0.040}	8.465 _{1.709}	0.544 _{0.011}	FPGNN	0.658 _{0.006}	1.106 _{0.195}
DMPNN	0.607 _{0.022}	0.388 _{0.077}	8.158 _{0.314}	0.448 _{0.014}	GROVER	0.617 _{0.077}	1.901 _{0.459}
ATTRMASKING	0.685 _{0.025}	0.546 _{0.052}	10.075 _{0.202}	0.547 _{0.024}	MOLFORMER	0.653 _{0.029}	1.190 _{0.046}
CONTEXTPred	0.669 _{0.030}	0.502 _{0.036}	9.445 _{0.224}	0.535 _{0.012}	MOLBERT	0.617 _{0.091}	1.311 _{0.257}
TRANFOXmol	0.645 _{0.036}	0.487 _{0.068}	9.055 _{0.523}	0.525 _{0.024}	TRANFOXmol	0.930 _{0.261}	1.225 _{0.155}
DEEPMOL	0.589 _{0.006}	0.327 _{0.012}	9.533 _{0.162}	0.660 _{0.004}	UNI-MOL	0.769 _{0.153}	1.598 _{0.153}
MuMo	0.426 _{0.031}	0.315 _{0.055}	7.324 _{0.323}	0.448 _{0.007}	MuMo	0.536 _{0.061}	1.082 _{0.088}

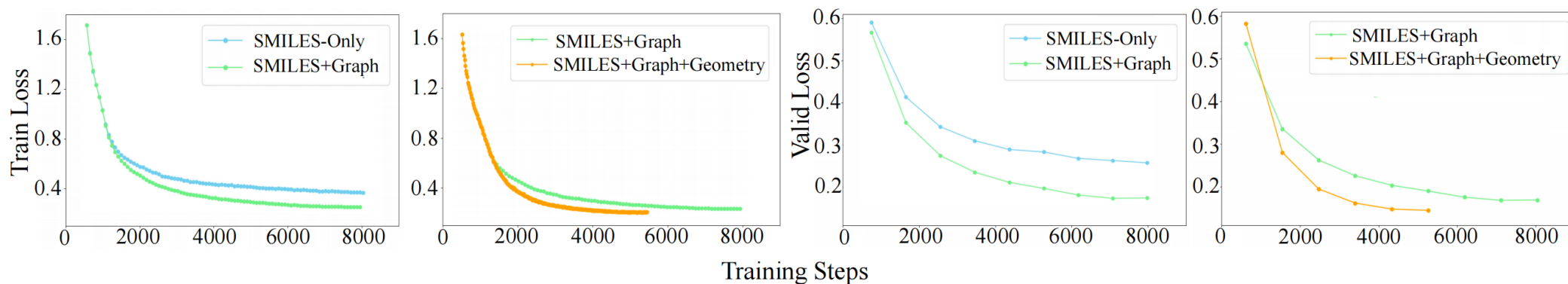
- Benchmarks: 29 tasks (TDC, MoleculeNet, Reaxtica).
- Overall: MuMo averages +2.7% over best baseline; rank-1 in 17/21 (TDC + MoleculeNet) and +27% on LD50 from TDC.
- QM datasets: wins 7/10 tasks, robust to conformer sensitivity.
- Broader chemical reaction benchmarks.

MODEL	HOMO/LUMO/GAP \downarrow	α \downarrow	C_v \downarrow	μ \downarrow	R^2 \downarrow	ZPVE \downarrow
GROVER-BASE	0.0079 _{3E-04}	2.365 _{0.302}	1.103 _{0.339}	0.618 _{0.002}	113.01 _{4.206}	0.0035 _{3E-04}
GROVER-LARGE	0.0083 _{6E-04}	2.240 _{0.385}	0.853 _{0.186}	0.623 _{0.006}	85.85 _{6.816}	0.0038 _{5E-04}
GEM	0.0067 _{4E-05}	0.589 _{0.0042}	0.237 _{0.0137}	0.444 _{0.0015}	25.67 _{0.743}	0.0011 _{2E-05}
UNI-MOL	0.0043 _{2E-05}	0.363 _{0.009}	0.183 _{0.002}	0.155 _{0.0015}	4.805 _{0.055}	0.0011 _{3E-05}
UNI-MOL2 310M	0.0036 _{1E-05}	0.315 _{0.003}	0.143 _{0.002}	0.092 _{0.0013}	4.672 _{0.245}	0.0005 _{1E-05}
UNI-MOL2 570M	0.0036 _{2E-05}	0.315 _{0.004}	0.147 _{0.007}	0.089 _{0.0015}	4.523 _{0.080}	0.0005 _{3E-05}
UNI-MOL2 1.1B	0.0035 _{1E-05}	0.305 _{0.003}	0.144 _{0.002}	0.089 _{0.0004}	4.265 _{0.067}	0.0005 _{8E-05}
MuMo 505M	0.0030 _{1E-05}	0.283 _{0.003}	0.126 _{0.003}	0.400 _{0.0018}	18.08 _{0.533}	0.0005 _{1E-05}

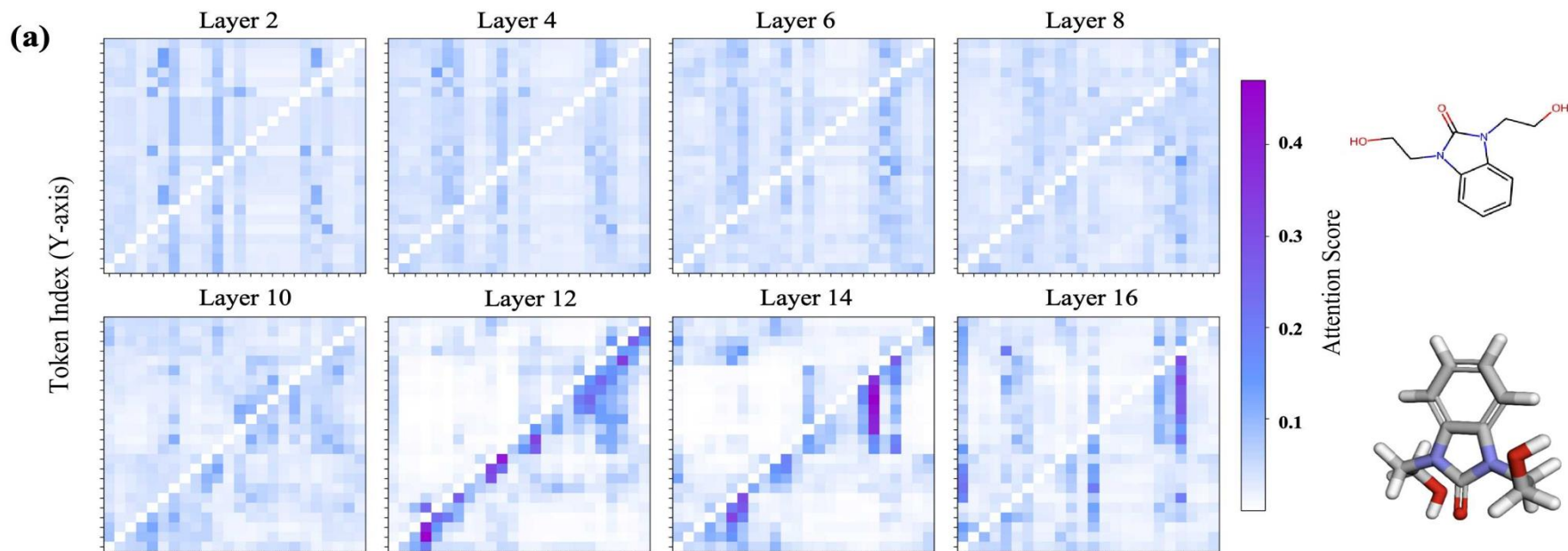
BHC (R^2 \uparrow , REACTION YIELD)		CPA (MAE \downarrow , CATALYTIC ACTIVITY)		HTE (R^2 \uparrow , REACTION YIELD)	
MODELS	VALUE	MODELS	VALUE	MODELS	VALUE
REAXTICA	0.94	REAXTICA	0.144	REAXTICA	0.87
MFF	0.92	MFF	0.144	RXNFP	0.81
RXNFP	0.95	DENMARK ET AL.	0.152	DRFP	0.85
MuMo	0.952 _{0.002}	MuMo	0.144 _{0.000}	MuMo	0.873 _{0.002}

Insights

- Pretraining losses show the effectiveness of modality fusion.



- Attention visualization. From Layer 9, the graph and geometry feature start interacting with the SMILES stream.



Welcome to connect and collaborate!

