# Multi-View Polymer Representations for the Open Polymer Prediction

**Wonjin Jung**
SK Telecom
onejiin@gmail.com

**Yongseok Choi**
SK Telecom
yongseokchoi@sk.com

## Abstract

We address polymer property prediction with a multi-view design that exploits complementary representations. Our system integrates four families—(i) tabular RDKit/Morgan descriptors, (ii) graph neural networks, (iii) 3D-informed representations, and (iv) pretrained SMILES language models—and averages per-property predictions via a uniform ensemble. Models are trained with 10-fold splits and evaluated with SMILES test-time augmentation. The approach ranks **9th** of 2,241 teams in the *Open Polymer Prediction* challenge in NeurIPS 2025. The submitted ensemble achieves a Public MAE of 0.057 and a Private MAE of 0.082.

## 1 Introduction

Accurate polymer property prediction is critical to materials discovery and selection, yet laboratory measurement remains time- and cost-intensive. Recent advances in machine learning have made it increasingly feasible to predict key polymer properties directly from structure with both speed and accuracy. The *Open Polymer Prediction* (OPP) challenge Open Polymer Prediction Organizers [2025] evaluates such systems by providing polymer SMILES strings as inputs and requiring predictions for five targets: glass transition temperature ($Tg$), crystallization temperature ($Tc$), fractional free volume ($FFV$), density ($Density$), and radius of gyration ($Rg$).

The official training data contain 7,973 SMILES, with property availability varying by target ($Tg$: 511, $Tc$: 737, $FFV$: 7,030, $Density$: 613, $Rg$: 614). To strengthen supervision, we expand the training set to $Tg$: 13,415, $Tc$: 2,471, $FFV$: 7,892, $Density$: 1,399, and $Rg$: 614 by incorporating additional data. Beyond scale, our central hypothesis is that polymers admit multiple complementary representations, and that integrating these views yields more robust generalization than any single family alone.

We therefore propose a practical multi-view approach that combines four modeling families. **(1) Tabular:** RDKit-derived descriptors (e.g., Morgan fingerprints Rogers and Hahn [2010]) trained with boosted decision trees (XGBoost). **(2) Graph:** atom–bond graphs with node and edge attributes, modeled via graph neural networks (e.g., GNNs tailored to polymer structure). **(3) 3D-informed:** geometric information captured through a pretrained model to enhance structure-aware embeddings. **(4) Pretrained SMILES encoders:** language-model style representations ("BERT"-like) warm-started from models pretrained on large polymer corpora. These complementary inductive biases allow the system to capture local substructure, global connectivity, geometric cues, and data-driven sequence regularities.

To further improve robustness, we employ two training and inference practices. First, we train all base learners using 10-fold cross-validation to stabilize estimates and enable ensembling over out-of-fold (OOF) predictions. Second, at test time we apply SMILES-format test-time augmentation (TTA) by generating equivalent canonicalizations/reformulations and averaging predictions. Final outputs are produced by a simple, property-wise *uniform* ensemble that aggregates selected base models without overfitting to the public leaderboard.
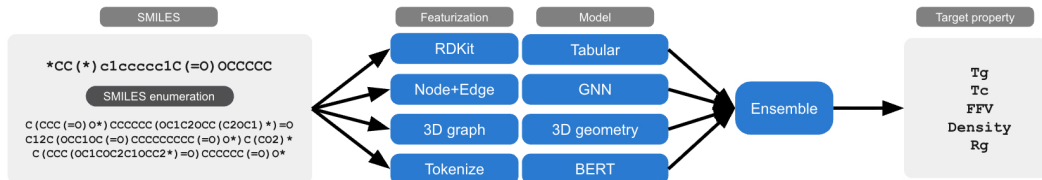
---

Preprint. Under review.

Figure 1: Pipeline overview

Empirically, this design delivers strong generalization in OPP: our system ranks **9th** out of 2,241 teams. Beyond this specific benchmark, we view the proposed generalized ensemble as a reusable pattern for materials AI, facilitating reliable property prediction across diverse manufacturing and application settings.

Our main contributions are as follows:

- **Multi-view modeling for generalization:** We diversify inductive biases by fusing four families (tabular RDKit/Morgan, GNNs on atom–bond graphs, 3D-informed representations, and pretrained SMILES LMs), which improves robustness over any single family.

- **Data-side generalization:** We stabilize training and evaluation via 10-fold splits and SMILES-based test-time augmentation (TTA), averaging predictions across equivalent reformulations.

- **Simple, reliable ensembling:** A per-property *uniform* ensemble ($1/n$ averaging) reduces variance and avoids leaderboard overfitting while preserving strong private-leaderboard performance.

## 2 Models

### 2.1 Overview

To explore complementary views of polymer SMILES, we evaluate a set of state-of-the-art models spanning four representation families: (i) *tabular* RDKit/Morgan descriptors, (ii) *graph neural networks* (GNNs) over atom–bond graphs, (iii) *3D-informed* representations via a pretrained model, and (iv) *pretrained SMILES language models*. Each family emphasizes a distinct inductive bias—local substructures and simple statistics (tabular), topological connectivity (graphs), geometric cues (3D-informed), and sequence-level regularities (LMs)—so that, when combined, they compensate for one another's weaknesses. An overview of the proposed pipeline is illustrated in Figure 1.

### 2.2 Tabular (RDKit/Morgan)

We extract fixed-length tabular features from SMILES using RDKit RDKit Development Team [2025], in particular circular fingerprints (Morgan). These features summarize atom environments into bit vectors or counts and are well-suited to tree-based learners. We use XGBoost as the primary tabular baseline, a strong choice in chemical domains for non-linear tabular data, and include Random Forest to further stabilize the tabular family under heterogeneous feature–target relations.

### 2.3 Graph Neural Networks (GNNs)

We convert each SMILES into a molecular graph with atoms as nodes, bonds as edges, and standard node/edge attributes. The graph datasets are trained with three complementary GNNs:

- **GINE (GIN with edge features).** A sum-aggregation MLP architecture with Weisfeiler–Lehman-level discriminative power; the "E" variant injects edge attributes into message updates to capture bond types and orders Hu et al. [2020].

- **GAT (Graph Attention Network).** Learns attention weights over neighbors so that more informative atoms/bonds receive higher contribution during message passing Veličković et al. [2018].

- **MPNN (Message Passing Neural Network).** A general message–update framework that parameterizes messages along edges and hidden-state updates at nodes, enabling flexible encoding of local chemistry Gilmer et al. [2017].

Using these three models provides diversity across aggregation (sum vs. attention) and parameterization (fixed aggregator vs. learned message functions), which improves coverage of graph-topological signals relevant to polymer properties.

### 2.4 3D-Informed Representation

Certain targets (e.g., radius of gyration, $Rg$) benefit from 3D cues. Directly generating reliable 3D conformations for long polymer SMILES is challenging, so we adopt a pretrained 3D-informed model to inject geometric priors with minimal overhead. In particular, we leverage GraphMVP as a backbone to obtain embeddings informed by 3D-aware pretraining Liu et al. [2021], while keeping the downstream head lightweight. Empirically, this preserves stable performance and supplies geometric signal where helpful.

### 2.5 Pretrained SMILES Language Models

To enhance generalization under limited labeled data, we warm-start from pretrained SMILES encoders and attach a regression head. We fine-tune PolyBERT Kuenneth and Ramprasad [2023], PolyCL Zhou et al. [2024], and TransPolymer Xu et al. [2023], which tokenize SMILES and exploit large-scale pretraining to capture grammar and frequent substructures. This family improves robustness to unseen SMILES variants and mitigates distribution shift by transferring sequence-level knowledge into the property predictor.

## 3 Data and Training Strategy

We train on the official OPP split and augment supervision with additional labels per target (see Introduction for counts). Our goal is to improve robustness under limited labeled data by (i) exploiting SMILES-level augmentation at inference and (ii) maximizing data usage via $K$-fold training while avoiding overfitting.

### 3.1 SMILES Augmentation

We apply SMILES-format augmentation for a given molecule, we generate multiple *equivalent* SMILES strings and average predictions across them Bjerrum [2017]. Concretely, we (1) randomize the traversal order when writing SMILES, (2) preserve stereochemistry and aromaticity, and (3) optionally toggle canonical/non-canonical rendering to diversify token sequences while keeping molecular identity unchanged. This yields a set of semantically identical inputs whose predictions are averaged, improving stability to sequence-level perturbations.

### 3.2 $K$-Fold Training

To make effective use of limited data per target, we adopt $K$-fold training for all base learners. In each fold, a disjoint validation subset is held out for model selection, and out-of-fold predictions are reserved for downstream ensembling. Rotating the validation set across folds increases data utilization while preserving an unbiased estimate. Although this increases the number of trained models by a factor of $K$, it materially improves generalization and variance reduction in the final ensemble.

## 4 Preprocessing and Inference Protocol

### 4.1 Chemical Standardization and Target Policy

We parse SMILES with RDKit and standardize core chemistry conventions (valence sanitization and, where applicable, kekulization) while preserving stereochemistry and aromaticity. Targets are used as

provided for each property; no target smoothing or external relabeling is applied. All models optimize mean absolute error (MAE), and evaluation follows the OPP public/private leaderboard protocol.

## 4.2 Inference-Time Aggregation and Ensembling

At inference, each base model produces per-property predictions across a set of augmented, *semantically equivalent* SMILES. Predictions are averaged per model and then fused via a property-wise *uniform* ensemble Wolpert [1992] (see Section 2). We consider model ensembles, including an 8-model set comprising *Tabular* (2), *GNN* (2), *BERT* (3), and *3D* (1). For leaderboard generalization, the submission uses uniform $1/n$ averaging per property. The combination of multi-view modeling, SMILES test-time augmentation, and $K$-fold training yields a robust predictor under sparse supervision.

# 5 Experiments

We report mean absolute error (MAE) on the public and private OPP leaderboards. Unless noted otherwise, results are aggregated over the five target properties. To promote generalization, we vary key hyperparameters across families, train all base models with 10-fold splits, and evaluate with SMILES test-time augmentation (TTA). We additionally enable stochastic inference via dropout at test time (MC dropout) and select the number of TTA variants (up to 10) based on validation performance.

## 5.1 Single-Model Baselines

Table 1 summarizes Public/Private MAE for each base model. These results reflect per-family configurations tuned within modest ranges under the shared 10-fold protocol and TTA.

Table 1: Single models: Public/Private MAE.

| Model | Public ↓ | Private ↓ |
|---|---|---|
| XGBoost | 0.067 | 0.092 |
| Random Forest | 0.067 | 0.091 |
| PolyBERT | 0.059 | 0.092 |
| PolyCL | 0.062 | 0.090 |
| TransPolymer | 0.059 | 0.088 |
| GraphMVP | 0.069 | 0.096 |
| MPNN | 0.061 | 0.088 |
| GAT | 0.063 | 0.077 |

## 5.2 Ensembles

Table 2 reports uniform-averaging ensembles built from the above families. The first row aggregates 16 models (2 Tabular + 6 GNN + 7 BERT + 1 3D) whose heads differ in capacity and regularization; the second row uses a compact 8-model subset (2 Tabular + 2 GNN + 3 BERT + 1 3D). For leaderboard robustness, the submitted system employs property-wise uniform $1/n$ averaging.

Table 2: Ensembles (uniform averaging).

| Ensemble | Public ↓ | Private ↓ |
|---|---|---|
| 16 models (2 Tabular + 6 GNN + 7 BERT + 1 3D) | 0.056 | 0.084 |
| 8 models (2 Tabular + 2 GNN + 3 BERT + 1 3D) | 0.057 | 0.082 |

**Observations.** The uniform ensemble is competitive and stable. Notably, the 8-model subset slightly improves Private MAE, suggesting that moderate diversity with controlled variance can generalize better than a larger pool. Across runs, increasing TTA (up to 10 SMILES variants) and

enabling test-time dropout consistently reduced variance, while 10-fold training provided reliable out-of-fold signals for ensembling without overfitting to the public leaderboard.

## 6 Conclusion

We presented a practical multi-view system for polymer property prediction that integrates four complementary families—tabular RDKit/Morgan descriptors, graph neural networks over atom–bond graphs, 3D-informed representations via a pretrained model, and pretrained SMILES language models—and aggregates per-property outputs with a simple, uniform $1/n$ ensemble. Combined with $K$-fold training and SMILES test-time augmentation, this design yielded strong and stable performance across targets.

While ensembling heterogeneous inductive biases improves robustness, our uniform averaging does not explicitly account for per-model reliability across properties or chemical subspaces, and the benefit of SMILES-based TTA depends on the diversity of valid reformulations. Future work includes training a meta-ensemble model to learn effective property-specific or data-dependent ensemble weights, developing uncertainty-aware stacking and calibration, and broadening augmentation beyond string reformulations (e.g., graph/substructure or 3D-aware variants).

## References

E. J. Bjerrum. Smiles enumeration as data augmentation for neural network modeling of molecules. *arXiv preprint arXiv:1703.07076*, 2017. URL https://arxiv.org/abs/1703.07076.

J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing for quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, volume 70 of *Proceedings of Machine Learning Research*, pages 1263–1272. PMLR, 2017.

W. Hu, B. Liu, J. Gomes, M. Zitnik, P. Liang, V. Pande, and J. Leskovec. Strategies for pre-training graph neural networks. In *International Conference on Learning Representations (ICLR)*, 2020. URL https://arxiv.org/abs/1905.12265. Introduces GINE (GIN with edge features).

C. Kuenneth and R. Ramprasad. polybert: a chemical language model to enable fully machine-driven ultrafast polymer informatics. *Nature Communications*, 14:4099, 2023. doi: 10.1038/s41467-023-39868-6.

S. Liu, H. Wang, W. Liu, J. Lasenby, H. Guo, and J. Tang. Pre-training molecular graph representation with 3d geometry. *arXiv preprint arXiv:2110.07728*, 2021. doi: 10.48550/arXiv.2110.07728. URL https://arxiv.org/abs/2110.07728. GraphMVP: multi-view pre-training leveraging 2D–3D consistency.

Open Polymer Prediction Organizers. Open polymer prediction challenge. https://open-polymer-challenge.github.io/, 2025. Accessed: 2025-09-21.

RDKit Development Team. RDKit: Open-source cheminformatics. https://www.rdkit.org, 2025. Accessed: 2025-09-21.

D. Rogers and M. Hahn. Extended-connectivity fingerprints. *Journal of Chemical Information and Modeling*, 50(5):742–754, 2010. doi: 10.1021/ci100050t.

P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. Graph attention networks. In *International Conference on Learning Representations (ICLR)*, 2018. URL http://arxiv.org/abs/1710.10903.

D. H. Wolpert. Stacked generalization. *Neural Networks*, 5(2):241–259, 1992. doi: 10.1016/S0893-6080(05)80023-1.

C. Xu, Y. Wang, and A. Barati Farimani. Transpolymer: a transformer-based language model for polymer property predictions. *npj Computational Materials*, 9:64, 2023. doi: 10.1038/s41524-023-01016-5.

J. Zhou, Y. Yang, A. M. Mroz, and K. E. Jelfs. Polycl: Contrastive learning for polymer representation learning via explicit and implicit augmentations, 2024.