

---

# Expressivity and Generalization: Fragment-Biases for Molecular GNNs

---

**Tom Wollschläger\***

Technical University of Munich  
t.wollschlaeger@tum.de

**Niklas Kemper\***

Technical University of Munich  
niklas.kemper@tum.de

**Leon Hetzel**

Technical University of Munich  
leon.hetzel@in.tum.de

**Johanna Sommer**

Technical University of Munich  
sommer@in.tum.de

**Stephan Günnemann**

Technical University of Munich  
s.guennemann@tum.de

## Abstract

Although recent advances in higher-order Graph Neural Networks (GNNs) improve the theoretical expressivity and molecular property predictive performance, they often fall short of the empirical performance of models that explicitly use fragment information as inductive bias. However, for these approaches, there exists no theoretic expressivity study. In this work, we propose the *Fragment-WL* test, an extension to the well-known Weisfeiler & Leman (WL) test, which enables the theoretic analysis of these fragment-biased GNNs. Building on the insights gained from the Fragment-WL test, we develop a new GNN architecture and a fragmentation with infinite vocabulary that significantly boosts expressiveness. We show the effectiveness of our model on synthetic and real-world data where we outperform all GNNs on Peptides and have 12% lower error than all GNNs on ZINC and 34% lower error than other fragment-biased models. Furthermore, we show that our model exhibits superior generalization capabilities compared to the latest transformer-based architectures, positioning it as a robust solution for a range of molecular modeling tasks.

## 1 Introduction

A common issue with Graph Neural Networks (GNNs) is their lack of expressiveness, including their inability to recognize substructures, which could limit their empirical performance [40]. In machine learning for chemistry, frequently occurring substructures, or fragments, are key predictors of molecular properties [39]. These fragments become even more crucial in larger systems like proteins [44].

To address this, recent methods enhance GNNs by improving their ability to distinguish non-isomorphic graphs, with the Weisfeiler & Leman (WL) test used to measure expressiveness [52; 40]. However, these approaches often emphasize theoretical expressiveness over practical performance and suffer from poor generalization to data that does not perfectly fit the training distribution [6]. Fragment-biased GNNs which incorporate fragment information directly as inductive biases tend to perform better but typically lack theoretical analysis [17; 51].

In this work, we bridge the gap between theory and practical performance by introducing the Fragment-WL test, which extends the standard WL test to analyze fragment-biased models. We also propose a new model that integrates a general fragmentation within the message-passing framework, improving generalization to out-of-distribution data. Our model allows for a novel fragmentation of molecular graphs with an infinite vocabulary composed of basic building blocks. We demonstrate state-of-the-art performance across various molecular datasets and tasks, including outperforming transformer-based architectures in some cases.

---

\*Equal contribution.

Our core contributions are as follows:

- We provide a more fine-grained hierarchy on the expressivity for a multitude of models that incorporate substructures as inductive bias.
- We propose a new general architecture that performs message passing along substructures together with a new fragmentation for molecules, achieving highest expressivity.
- We evaluate predictive power, long-range performance, and generalization across extensive experiments.

## 2 Weisfeiler & Lehman Go Fragments

Existing fragment-biased GNNs vary not only in their vocabulary (what substructures are considered) but also in how fragmentation information is integrated into the model (see Related Work in Appendix B). Approaches range from including fragment information as *node features* (NF) [4], learning an explicit *representation for each fragment* (FR) that exchanges messages with the underlying nodes [51; 54] or, building a *higher-level graph* (HLG) on which neighboring fragments can exchange additional messages [17]. This variability makes a direct comparison of the expressiveness of these models challenging. To address this, we develop in Appendix C new variants of the Weisfeiler & Lehman (WL) test (NF-WL, FR-WL, HLG-WL) that model these different approaches to incorporate fragment information. Our Fragment-WL test framework allows us to bound the expressivity of most existing fragment-biased models (see Table 2 in Appendix C) and proves that the expressivity strictly increases from NF to FR to HLG approaches (see Theorems C.6 to C.8 in Appendix C). Formally, we establish the following hierarchy:

$$2\text{-WL} < \text{NF-WL} < \text{FR-WL} < \text{HLG-WL}. \quad (1)$$

This hierarchy demonstrates that the method of incorporating fragment information significantly impacts model expressivity, with higher-level abstractions yielding more powerful models.

## 3 Fragment Graph Neural Network

Based on the insights of the previous section, we propose our new model architecture and a new fragmentation scheme with infinite vocabulary consisting of only basic building blocks. Given the higher expressiveness, our model can differentiate complex substructures given only these basic building blocks, as it is able to learn the dependencies on the higher-level graph. We empirically confirm this in Section 4.

### 3.1 Model

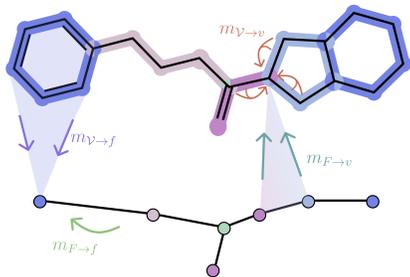
Building on the theoretical findings that a higher-level graph (see Theorem C.8) improves expressiveness, we propose FragNet, a general model for any fragmentation  $F$  that performs message-passing on the original graph *and* a higher-level graph of fragments, i.e., the HLG approach. Correspondingly, we have learned representations for every node  $v \in \mathcal{V}$  and every fragment  $f \in F$ . Conceptually, each node  $v$  receives messages from all its neighbors, and all fragment of which it is part of. Similarly, each fragment  $f$  receives messages from neighboring fragments and all nodes that are part of it. The aggregated messages together with the previous representation are then used to update the representation. The message passing scheme is illustrated in Figure 1 and precisely defined in Appendix E. The final graph-level output is computed by aggregating all representations after  $T$  layers. Note that the complexity of our FragNet model is linear in the number of nodes and fragments (assuming that each node is only part of a constant number of fragments).

Additionally, our FragNet model achieves the highest expressiveness in our Fragment-WL hierarchy and also compared to other fragment-biased GNNs.

**Theorem 3.1.** *FragNets are at most as powerful as HLG-WL. Additionally, when using injective neighborhood aggregators and a sufficient number of layers, FragNets are as powerful as HLG-WL.*

### 3.2 Molecular fragmentation

Apart from the question of how to use a fragmentation, there is the equally important question of how to fragment the graph in the first place. The challenge of fragmentation lies in balancing two goals:



**Figure 1:** Overview of our model and our fragmentation. The molecular graph is fragmented with our rings-paths fragmentation into three cycles, three paths, and a junction node. The figure shows the messages  $m_{F \rightarrow f}^t$ ,  $m_{v \rightarrow f}^t$  to one fragment  $f$ , and the messages  $m_{v \rightarrow v}^t$ ,  $m_{F \rightarrow v}^t$  to one vertex  $v$ .

1. Capture all *important* substructures.
2. Facilitate generalization.

A fragmentation that is too coarse may miss key features, while too fine-grained risks overfitting and hampers finding graph similarities. Obviously, the optimal scheme depends on the application. Existing methods for molecules focus either on a single substructure [54] or use chemical properties that require large vocabularies [10].

Our approach completely fragments molecular graphs using only rings and paths. First, minimal rings are extracted. Next, the remaining edges are connected at nodes of degree two to form paths. Lastly, junction nodes are introduced where three or more fragments meet, reducing cycles in the higher-level graph. Figure 1 illustrates this process.

**Ordinal encoding.** Previous works either use no encoding for the types of fragments [17] or a simple one-hot encoding [4]. However, to facilitate the generalization capabilities of a model, the encodings of similar fragments should also be similar. In our approach, we introduce an ordinal fragment encoding, which accomplishes this by incorporating two embeddings: one for the fragment class (i.e.,  $\text{class}(f) \in \{\text{path}, \text{cycle}, \text{junction}\}$ ) and another that is proportionally scaled based on the fragment size; see Figure 2. More formally:

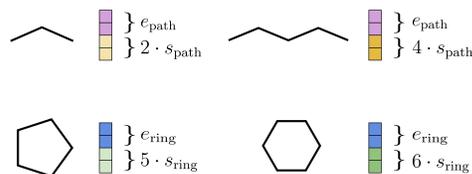
$$h_f^0 = (e_{\text{class}(f)}, |f| \cdot s_{\text{class}(f)}),$$

where  $s$  and  $e$  are different learned embeddings for the classes of cycle, path, and junction fragments. This approach enables the encoding of an infinite vocabulary, accommodating even completely unseen fragments while concurrently supporting effective model generalization.

## 4 Results

While we have theoretically demonstrated that our model attains the highest expressiveness within our Fragment-WL hierarchy, we also empirically evaluate its expressiveness by examining its ability to count substructures. Additionally, we explore its overall predictive effectiveness, and its capacity to generalize. Additional experiments in Appendix F demonstrate FragNet’s improved capability for long-range communication.

**Expressiveness.** To evaluate how well our model can learn to recognize chemically *important* substructures in molecular graphs, we first identify the most common substructures in the ZINC 10k dataset [22] using a chemically-inspired fragmentation scheme, specifically MAGNet [24]. Subsequently, we train our model to predict substructure counts. Our model is able to identify all substructures close to perfection as demonstrated in Table 10 in Appendix F. Notably, our model achieves high accuracy even for intricate substructures not present in our vocabulary. This underscores that our fragmentation, based solely on rings and paths, together with our ordinal encoding and the higher-level message passing, proves sufficient for the model to recognize more complex substructures.



**Figure 2:** Ordinal encoding applied to a 2-path, 4-path, 5-ring, and 6-ring. The encoding comprises two components: one learned embedding  $e$  for every fragment class (i.e., path, cycle, or junction) and another learned embedding  $s$  that is proportionally scaled based on the fragment size.

Type	Model	Peptides-		ZINC	
		Struct (MAE ↓)	Func (AP ↑)	10k (MAE ↓)	Full (MAE ↓)
Transformer	GPS	0.2500 ± 0.0012	0.6535 ± 0.0041	0.070 ± 0.006	-
	GRIT	<b>0.2460</b> ± 0.0012	<b>0.6988</b> ± 0.0082	<b>0.059</b> ± 0.002	<b>0.023</b> ± 0.001
Basic GNNs	GCN	0.3496 ± 0.0013	0.5930 ± 0.0023	0.367 ± 0.011	0.113 ± 0.002
	GIN	0.3547 ± 0.0045	0.5498 ± 0.0079	0.526 ± 0.051	0.088 ± 0.002
Topological	CIN++	0.2523 ± 0.0013	0.6569 ± 0.0117	<b>0.077</b> ± 0.004	0.027 ± 0.007
Fragment-Biased	HIMP	0.2503 ± 0.0008	0.5668 ± 0.0149	0.151 ± 0.006	0.036 ± 0.002
	FragNet (ours)	<b>0.2462</b> ± 0.0021	<b>0.6678</b> ± 0.005	<b>0.0775</b> ± 0.005	<b>0.0237</b> ± 0.00

**Table 1:** Predictive performance for multiple models on Peptides-struct/-func and ZINC. Best Transformer and best GNN are highlighted. A comparison with more models is in Tables 7 and 8.

This is essential for application in, e.g., fragment-based molecule generation, as the task of the encoder is to encode information about such substructures.

**Predictive Performance.** To evaluate the predictive performance on real-world molecular dataset, we use the long-range peptides benchmark [14] and the large-scale molecular benchmark ZINC [46]. The Peptides-struct and Peptides-func datasets are commonly used to benchmark long-range performance of GNNs and transformers. The task in the ZINC dataset is to predict the penalized logP of molecules, a measure of drug-likeness. A summary of all models we compare against, the used hyperparameters of our model, the used datasets and the experimental details for each experiment can be found in Appendix G. All our models adhere to the 500k parameter budget for both datasets. We do not use any additional feature augmentation, such as positional encodings. As shown in Section 4, our model achieves state-of-the-art performance among GNNs on all datasets, additionally surpassing nearly all graph transformers, which typically excel in long-range tasks due to their quadratic complexity. GRIT [36] is the only model that consistently outperforms ours.

**Generalization.** To test the generalization capabilities of our model with the ordinal fragment encoding, we use a test set containing out-of-distribution molecules with completely unseen fragments. For this, we use the ZINC dataset and remove all molecules containing a 7-ring from the training data. After training, we test on all molecules from the test set, thus also containing 7-rings that were not seen during training. The results in Table 11 demonstrate that our model achieves an error 1.8 times lower than GRIT, showcasing the superior generalization capabilities. Our better generalization capabilities can also be seen in the normal ZINC benchmark. In Table 12, we group the ZINC dataset into groups based on the frequency of the rarest fragment. Our model outperforms HIMP everywhere and GRIT for graphs containing rare fragments while GRIT shows better performance for molecules containing frequent fragments. Lastly, we test our model’s capability to transfer the knowledge to a completely different dataset. We train on ZINC and predict the penalized logP on QM9 [48]. FragNet achieves the lowest MAE of 1.12, outperforming GRIT (MAE 1.22) and HIMP (MAE 3.43), suggesting that our model generalizes better to unseen data distributions due to our inductive bias and corresponding fragmentation. In summary, we showcased the generalization capabilities of our model on both a completely unseen dataset and a slightly shifted data distribution. The generalization capabilities also help our model perform better on rare fragments.

## 5 Conclusion and Limitations

In this work, we introduced the Fragment-WL test, a new expressivity measure that provides a hierarchy for fragment-biased GNNs. Using this framework, we developed an expressive model that outperforms all GNN approaches and most transformer architectures on molecular property prediction benchmarks. Our model demonstrates strong generalization capabilities with linear complexity, making it a robust solution for molecular modeling tasks.

However, our method has limitations. The fragmentation and ordinal encoding are tailored to molecules and are less effective on large, densely connected graphs like citation or social networks, where they introduce noise. Additionally, while we outperform GRIT on rare fragments, GRIT performs slightly better on frequent ones (Table 12). Future work could improve fragment-biased models for frequent data and explore their use in other domains.

## Acknowledgements

This project is supported by the Bavarian Ministry of Economic Affairs, Regional Development and Energy with funds from the Hightech Agenda Bayern. Additionally, it is supported by the Helmholtz Association under the joint research school "Munich School for Data Science - MUDS" and by the German Federal Ministry of Education and Research (BMBF) under Grant No. 01IS18036B.

## Impact Statement

Among other contributions, this work presents an approach for predicting the properties of molecules. In the area of machine learning for drug discovery, such methods can sometimes be used for harmful purposes. This also applies to our research, since it might help to discover or create dangerous substances. Despite these concerns, we believe that the benefits of our work outweigh the risks.

## References

- [1] Bemis, G. W. and Murcko, M. A. The properties of known drugs. 1. molecular frameworks. *Journal of Medicinal Chemistry*, 1996.
- [2] Bodnar, C., Frasca, F., Wang, Y. G., Otter, N., Montúfar, G., Liò, P., and Bronstein, M. Weisfeiler and Lehman Go Topological: Message Passing Simplicial Networks, 2021.
- [3] Bodnar, C., Frasca, F., Otter, N., Wang, Y. G., Liò, P., Montúfar, G., and Bronstein, M. Weisfeiler and Lehman Go Cellular: CW Networks, 2022.
- [4] Bouritsas, G., Frasca, F., Zafeiriou, S., and Bronstein, M. M. Improving Graph Neural Network Expressivity via Subgraph Isomorphism Counting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [5] Bresson, X. and Laurent, T. Residual gated graph convnets, 2018.
- [6] Campi, F., Gosch, L., Wollschläger, T., Scholten, Y., and Günnemann, S. Expressivity of Graph Neural Networks Through the Lens of Adversarial Robustness, 2023.
- [7] Chen, Z., Chen, L., Villar, S., and Bruna, J. Can Graph Neural Networks Count Substructures?, 2020.
- [8] Chen, Z., Villar, S., Chen, L., and Bruna, J. On the equivalence between graph isomorphism testing and function approximation with gnns, 2023.
- [9] Degen, J., Wegscheid-Gerlach, C., Zaliani, A., and Rarey, M. On the art of compiling and using 'drug-like' chemical fragment spaces. *ChemMedChem*, 2008.
- [10] Degen, J., Wegscheid-Gerlach, C., Zaliani, A., and Rarey, M. On the art of compiling and using 'drug-like' chemical fragment spaces. *ChemMedChem*, 2008.
- [11] Di Giovanni, F., Giusti, L., Barbero, F., Luise, G., Lio', P., and Bronstein, M. On Over-Squashing in Message Passing Neural Networks: The Impact of Width, Depth, and Topology, 2023.
- [12] Du, Y., Fu, T., Sun, J., and Liu, S. Molgensurvey: A systematic survey in machine learning models for molecule design. *arXiv preprint arXiv:2203.14500*, 2022.
- [13] Dwivedi, V. P. and Bresson, X. A Generalization of Transformer Networks to Graphs, 2021.
- [14] Dwivedi, V. P., Joshi, C. K., Luu, A. T., Laurent, T., Bengio, Y., and Bresson, X. Benchmarking Graph Neural Networks, 2022.
- [15] Falcon, W. and The PyTorch Lightning team. PyTorch Lightning, 2019.
- [16] Fey, M. and Lenssen, J. E. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.

- [17] Fey, M., Yuen, J.-G., and Weichert, F. Hierarchical Inter-Message Passing for Learning on Molecular Graphs, 2020.
- [18] Frasca, F., Bevilacqua, B., Bronstein, M. M., and Maron, H. Understanding and Extending Subgraph GNNs by Rethinking Their Symmetries, 2022.
- [19] Geisler, S., Li, Y., Mankowitz, D., Cemgil, A. T., Günnemann, S., and Paduraru, C. Transformers meet directed graphs, 2023.
- [20] Geng, Z., Xie, S., Xia, Y., Wu, L., Qin, T., Wang, J., Zhang, Y., Wu, F., and Liu, T.-Y. De novo molecular generation via connection-aware motif mining. *International Conference on Learning Representations*, 2023.
- [21] Giusti, L., Reu, T., Ceccarelli, F., Bodnar, C., and Liò, P. CIN++: Enhancing Topological Message Passing, 2023.
- [22] Gómez-Bombarelli, R., Wei, J. N., Duvenaud, D., Hernández-Lobato, J. M., Sánchez-Lengeling, B., Sheberla, D., Aguilera-Iparraguirre, J., Hirzel, T. D., Adams, R. P., and Aspuru-Guzik, A. Automatic chemical design using a data-driven continuous representation of molecules, 2017.
- [23] He, X., Hooi, B., Laurent, T., Perold, A., Lecun, Y., and Bresson, X. A Generalization of ViT/MLP-Mixer to Graphs. In *International Conference on Machine Learning*, 2023.
- [24] Hetzel, L., Sommer, J., Rieck, B., Theis, F. J., and Günnemann, S. MAGNet: Motif-Agnostic Generation of Molecules from Shapes. *arXiv*, 2023.
- [25] Hu, W., Liu, Y., Chen, X., Chai, W., Chen, H., Wang, H., and Wang, G. Deep learning methods for small molecule drug discovery: A survey. *IEEE Transactions on Artificial Intelligence*, 2023.
- [26] Huang, N. and Villar, S. A short tutorial on the weisfeiler-lehman test and its variants. *International Conference on Acoustics Speech and Signal Processing*, 2021.
- [27] Huang, Y., Peng, X., Ma, J., and Zhang, M. Boosting the cycle counting power of graph neural networks with  $I^2$ -GNNs. *International Conference on Learning Representations*, 2022.
- [28] Inae, E., Liu, G., and Jiang, M. Motif-aware attribute masking for molecular graph pre-training. *arXiv preprint arXiv:2309.04589*, 2023.
- [29] Jin, W., Barzilay, R., and Jaakkola, T. Junction Tree Variational Autoencoder for Molecular Graph Generation, 2019.
- [30] Jin, W., Barzilay, R., and Jaakkola, T. Hierarchical generation of molecular graphs using structural motifs. In *International Conference on Machine Learning*. PMLR, 2020.
- [31] Kiefer, S. and Neuen, D. The power of the weisfeiler-lehman algorithm to decompose graphs. *SIAM Journal on Discrete Mathematics*, 2022.
- [32] Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks, 2017.
- [33] Kong, X., Huang, W., Tan, Z., and Liu, Y. Molecule Generation by Principal Subgraph Mining and Assembling, 2022.
- [34] Kreuzer, D., Beaini, D., Hamilton, W. L., Létourneau, V., and Tossou, P. Rethinking graph transformers with spectral attention, 2021.
- [35] Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2017.
- [36] Ma, L., Lin, C., Lim, D., Romero-Soriano, A., Dokania, P. K., Coates, M., Torr, P., and Lim, S.-N. Graph Inductive Biases in Transformers without Message Passing, 2023.
- [37] Maron, H., Ben-Hamu, H., Serviansky, H., and Lipman, Y. Provably powerful graph networks. *Advances in neural information processing systems*, 2019.

- [38] Maziarz, K., Jackson-Flux, H., Cameron, P., Sirockin, F., Schneider, N., Stiefl, N., Segler, M., and Brockschmidt, M. Learning to Extend Molecular Scaffolds with Structural Motifs, 2022.
- [39] Merlot, C., Domine, D., Cleva, C., and Church, D. J. Chemical substructures in drug discovery. *Drug Discovery Today*, 2003.
- [40] Morris, C., Ritzert, M., Fey, M., Hamilton, W. L., Lenssen, J. E., Rattan, G., and Grohe, M. Weisfeiler and Leman Go Neural: Higher-order Graph Neural Networks, 2021.
- [41] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. Pytorch: An imperative style, high-performance deep learning library, 2019.
- [42] Puny, O., Lim, D., Kiani, B. T., Maron, H., and Lipman, Y. Equivariant polynomials for graph neural networks, 2023.
- [43] Rampásek, L., Galkin, M., Dwivedi, V. P., Luu, A. T., Wolf, G., and Beaini, D. Recipe for a General, Powerful, Scalable Graph Transformer, 2023.
- [44] Singh, R. and Saha, M. Identifying structural motifs in proteins. *Pacific Symposium on Biocomputing*, 2003.
- [45] Sommer, J., Hetzel, L., Lüdke, D., Theis, F. J., and Günnemann, S. The Power of Motifs as Inductive Bias for Learning Molecular Distributions, 2023.
- [46] Sterling, T. and Irwin, J. J. Zinc 15 – ligand discovery for everyone. *Journal of Chemical Information and Modeling*, 2015.
- [47] Thiede, E. H., Zhou, W., and Kondor, R. Autobahn: Automorphism-based graph neural nets. *ArXiv*, 2021.
- [48] Wu, Z., Ramsundar, B., Feinberg, E. N., Gomes, J., Geniesse, C., Pappu, A. S., Leswing, K., and Pande, V. S. Moleculenet: A benchmark for molecular machine learning. *arXiv: Learning*, 2017.
- [49] Xu, K., Hu, W., Leskovec, J., and Jegelka, S. How Powerful are Graph Neural Networks?, 2019.
- [50] Ying, C., Cai, T., Luo, S., Zheng, S., Ke, G., He, D., Shen, Y., and Liu, T.-Y. Do Transformers Really Perform Bad for Graph Representation?, 2021.
- [51] Zang, X., Zhao, X., and Tang, B. Hierarchical molecular graph self-supervised learning for property prediction. *Communications Chemistry*, 2023.
- [52] Zhang, B., Fan, C., Liu, S., Huang, K., Zhao, X., Huang, J., and Liu, Z. The expressive power of graph neural networks: A survey. *arXiv preprint arXiv:2308.08235*, 2023.
- [53] Zhang, Z., Liu, Q., Wang, H., Lu, C., and Lee, C.-K. Motif-based graph self-supervised learning for molecular property prediction. *Advances in Neural Information Processing Systems*, 2021.
- [54] Zhu, J., Wu, K., Wang, B., Xia, Y., Xie, S., Meng, Q., Wu, L., Qin, T., Zhou, W., Li, H., and Liu, T.-Y.  $\mathcal{O}$ -GNN: incorporating ring priors into molecular modeling, 2022.

## A Background

**Notation.** A graph  $G := (\mathcal{V}, \mathcal{E}, \mathbf{X})$  consists of a set of vertices  $\mathcal{V}$ , a set of (undirected) edges  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$  and  $d$  node features  $\mathbf{X} \in \mathbb{R}^{|\mathcal{V}| \times d}$  for every node  $v \in \mathcal{V}$ . The set of nodes that are adjacent to  $v$  is denoted by  $\mathcal{N}(v)$ . Two graphs  $G^1 = (\mathcal{V}^1, \mathcal{E}^1, \mathbf{X}^1)$  and  $G^2 = (\mathcal{V}^2, \mathcal{E}^2, \mathbf{X}^2)$  are *isomorphic* if there exists a bijection  $b : \mathcal{V}^1 \rightarrow \mathcal{V}^2$  that preserves edges and node features, that is,  $\{v, w\} \in \mathcal{E}^1 \Leftrightarrow \{b(v), b(w)\} \in \mathcal{E}^2$  and  $\mathbf{X}_v^1 = \mathbf{X}_{b(v)}^2$ . For a subset of nodes  $\mathcal{U} \subseteq \mathcal{V}$ , we denote the induced subgraph with respect to these nodes by  $G[\mathcal{U}]$ .

**Expressiveness.** We can classify the expressiveness of functions over graphs by their capability to distinguish non-isomorphic graphs. We say that a function  $f$  is (in parts) *more powerful* than a function  $g$  if there exist two non-isomorphic graphs  $G^1, G^2$  such that  $f(G^1) \neq f(G^2)$  whereas  $g(G^1) = g(G^2)$ . The function  $f$  is *strictly more powerful* than  $g$  (we write  $f > g$ ) if  $f$  is more powerful than  $g$  and  $g$  is not (in parts) more powerful than  $f$ .

**Weisfeiler & Leman.** The Weisfeiler & Leman graph isomorphism test is an iterative graph coloring algorithm that bounds the expressive power of MPNNs [31]. In each iteration, it produces a color for each node based on its neighboring nodes’ colors. Starting with a vertex color based only on features  $c_v^0 = \text{HASH}(\mathbf{X}_v)$ , we calculate the update for the color  $c$  of node  $v$  in iteration  $t$ :

$$c_v^{(t)} = \text{HASH} \left( c_v^{(t-1)}, \{ \{ c_w^{(t-1)} \mid w \in \mathcal{N}(v) \} \} \right). \quad (2)$$

The algorithm terminates once the set of unique colors does not increase. Two non-isomorphic graphs can be distinguished if the multiset of colors differs at the end. As this test cannot distinguish all non-isomorphic graphs, it can be extended to strictly more powerful versions,  $k$ -WL, incorporating  $k$ -tuples of nodes to determine the color. For more background information, we refer to Morris et al. [40]. Importantly, 2-WL is equivalent to the previously described WL test [26].

**Fragmentations.** A vocabulary  $\mathcal{Y}$  is a set of graphs (potentially including node features) representing important substructures, e.g., cycles. A fragment of a graph  $G$  is an induced subgraph  $G[f]$  isomorphic to a graph from the vocabulary. We will identify a fragment simply by the subset of nodes  $f \subseteq \mathcal{V}$ . All fragments  $f$  that are isomorphic to the same graph of the vocabulary have the same  $\text{type}(f)$ . A fragmentation scheme  $\mathcal{F}$  is a permutation invariant function that maps a graph  $G$  to a set of fragments  $\mathcal{F}(G) =: F$ , which is called a fragmentation. Note that there might exist subgraphs isomorphic to a graph in  $\mathcal{Y}$  that are *not* in  $\mathcal{F}(G)$ . For example, even if  $\mathcal{Y}$  contains 5-cycles, not all 5-cycles in  $G$  need to be in  $\mathcal{F}(G)$ . If, for all graphs  $G$ ,  $\mathcal{F}(G)$  includes *every* subgraph isomorphic to a graph  $V \in \mathcal{Y}$ , we say that the fragmentation scheme  $\mathcal{F}$  recovers  $V$ .

## B Related work

**Expressiveness of GNNs.** Message-Passing Neural Networks (MPNN)<sup>2</sup> are limited in their expressiveness. Their ability to distinguish between non-isomorphic graphs is confined to the 2-WL algorithm, restricting their discriminative power [49]. Moreover, when it comes to recognizing substructures, MPNNs are unable to accurately count almost all types of substructures [7]. This limitation stems from their reliance on purely local messages, which—despite facilitating excellent linear space and time complexity—renders them blind to higher-level structural information within graphs.

**Higher-order GNNs.** In response to the inability to effectively learn substructures, the introduction of more powerful GNN architectures aims to overcome this limitation and enable comprehensive substructure learning. Morris et al. [40] draw inspiration from the multidimensional  $k$ -WL algorithm and diverge from learning node-specific representations by considering each  $k$ -tuple of nodes instead. Although this improves expressiveness, its complexity increases exponentially. Subgraph GNNs comprise an alternative to improve substructure identification, decomposing a graph into smaller subgraphs for GNN application. The resulting subgraph representations are pooled before a final graph level representation is derived [27; 18]. With some strategies for extracting subgraphs, subgraph GNNs can identify basic substructures such as 4-cycles [27]. Pany et al. [42] extend the WL test for higher-order GNNs to the graph polynomial counting problem as a new expressivity measure, highlighting the importance of more fine-grained tests for GNNs. However, the limitations of higher-order GNNs lie in their inability to effectively learn more intricate substructures, accompanied by an

<sup>2</sup>We use MPNNs and GNNs interchangeably.

increase in time complexity. Recent findings also suggest susceptibility to adversarial attacks and out-of-distribution data, hinting at challenges in robustly learning substructures [6].

**Fragment-Biased GNNs.** Another line of work provides fragment information to GNNs as an explicit inductive bias. These fragment-biased models vary not only in their vocabulary but also in the way fragmentation information is integrated into the model. *Node features:* Bouritsas et al. [4] introduce GSN-v, which uses the number of cycles or cliques as an additional node feature. *Learned fragment representation:* Instead of treating fragmentation information as a fixed feature, other models learn representations for each fragment by aggregating information from the corresponding nodes. Zhu et al. [54] use a vocabulary of only cycles whereas Zang et al. [51] present HiMol, which fragments a molecular graph, based on chemical properties. *Higher-level graph:* A natural extension of the learned fragment representation is a higher-level graph of fragments where neighboring fragments influence each other. Thiede et al. [47] use equivariant computations along the paths of length 3 to 6 and cycles of sizes 5 and 6. Fey et al. [17] build a higher-level junction tree using rings and edges. Yet, none of the existing works compare—theoretically or experimentally—how to encode and use substructure information in the model. Additionally, most works only focus on a single substructure that does not allow to fragment the complete graph.

**Topological GNNs** use higher-level topological structures such as simplicial complexes [2] or CW-Networks [3; 21] in their message-passing schemes. While coming from a different theoretical direction than substructure-biased GNNs, in practice, they use cliques or cycles as learned fragment representations or in a higher-level graph.

**Graph Transformer.** Recently, models such as Graph Transformers [50; 36; 19] and ViT/MLP-Mixers [23] for graphs adapted successful models from other domains to graph data. Their ability to recognize substructures depends on the positional encoding used. Almost all recent models use random walk encodings, which can help to discover simple substructures like cycles.

**Fragmentation Schemes.** Fragmentation methods in the chemical domain aim to divide a molecular graph into subgraphs with distinct structures or properties. There are various strategies to achieve this, such as separating probable reactants [9], categorizing molecules into distinct structural classes [1], or breaking apart acyclic bonds [38; 29]. Unlike these methods, data-driven approaches like those outlined in Kong et al. [33] and Geng et al. [20] focus on deriving subgraphs directly from a dataset without relying on predefined rules for decomposition.

## C Weisfeiler & Leman Go Fragments

Existing fragment-biased MPNNs vary in their underlying fragmentation scheme and how the fragment information is incorporated into the model. This variability makes a direct comparison of the expressiveness of these models difficult. To address this challenge, we propose a new, more fine-grained version of the WL test, called *Fragment-WL* test, that incorporates detailed structural elements. We derive a hierarchy of tests that capture how fragmentation information is incorporated into existing substructure-biased models, while leaving out all variability that does not influence the expressivity.

Our Fragment-WL test also subsumes existing WL variants designed for simplicial complexes and CW-cells [2; 3], providing a more unified framework for assessing the expressiveness of both substructure-biased and topological GNNs. Furthermore, our proposed Fragment-WL test highlights the significance of how fragment information is incorporated into the model, emphasizing that the integration methodology plays a crucial role for determining the model’s expressive power.

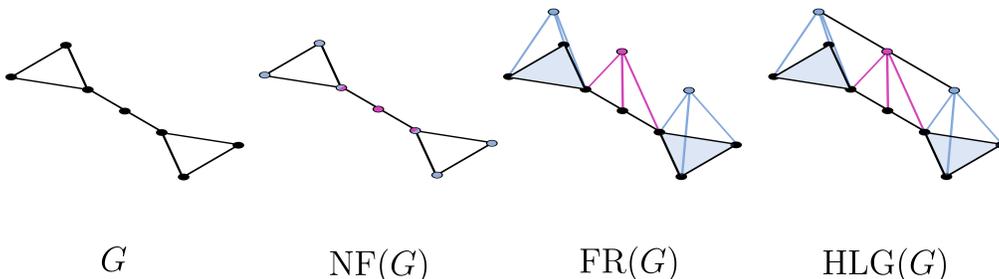
Fragment-WL entails multiple variants with increasing expressiveness in distinguishing isomorphic graphs. In the following, we first provide a general framework and then define the individual Fragment-WL versions that perform the original WL test on different augmented graphs. We start with a definition of WL tests on augmented graphs:

**Definition C.1.** A  $g$ -WL test is a function that performs the WL test on the augmented graph  $g(G)$ , i.e.

$$g\text{-WL}(G) := \text{WL}(g(G))$$

where  $g$  is a function mapping from graphs to graphs, i.e.,  $g : (\mathcal{V}, \mathcal{E}, \mathbf{X}) \mapsto (\mathcal{V}', \mathcal{E}', \mathbf{X}')$ .

There are three ways in which a fragmentation  $F$  is used in existing fragment-biased GNNs: as an additional node feature, as learned fragment representation, and as a higher-level graph. We



**Figure 3:** Example graph  $G$  with corresponding augmented variants.  $\text{NF}(G)$  includes node features,  $\text{FR}(G)$  also includes a representation for each fragment and  $\text{HLG}(G)$  also has connections between neighboring fragment representations.

instantiate  $g$  with the corresponding functions to augment the graph with the respective features. First, we use additional node features. We extend the individual node features with the information of the fragments that the node is contained in. We concatenate this information to the already existing features. Formally, we define this augmentation function in the following way.

**Definition C.2.** We define the node feature function as  $\text{NF}(\mathcal{V}, \mathcal{E}, \mathbf{X}) = (\mathcal{V}, \mathcal{E}, \mathbf{X}^{\text{NF}})$  with:

$$\mathbf{X}_v^{\text{NF}} = X_v \parallel \lambda(\{\{\text{type}(f) \mid v \in f, f \in F\}\}),$$

where  $\lambda$  represents any injective function and  $\parallel$  indicates the concatenation operation. We instantiate  $g$  with  $\text{NF}$  to create the  $\text{NF-WL}$  test.

Another prominent way is using representations for each fragment and messages that are flowing from the lower level nodes to their entailing fragment and backwards. This means that we introduce a new vertex for each fragment and connect it to all its corresponding vertices in the original graph. We depict this graph  $\text{FR}(G)$  in Figure 3. We define this augmentation function in the following.

**Definition C.3.** We define the fragment representation function as  $\text{FR}(\mathcal{V}, \mathcal{E}, \mathbf{X}) = (\mathcal{V}^{\text{FR}}, \mathcal{E}^{\text{FR}}, \mathbf{X}^{\text{FR}})$  with:

$$\begin{aligned} \mathcal{V}^{\text{FR}} &:= \mathcal{V} \cup F, \\ \mathcal{E}^{\text{FR}} &:= \mathcal{E} \cup \{\{f, v\} \mid \forall f \in F, \forall v \in f\}, \\ \mathbf{X}_i^{\text{FR}} &:= \begin{cases} \mathbf{X}_i & i \in \mathcal{V} \\ \text{type}(i) & i \in F \end{cases} \end{aligned}$$

Lastly, we allow messages to be exchanged between neighboring fragments, thus creating a higher-level graph on which information can flow. To this end, we add edges between fragments that have neighboring nodes in  $G$  and thus construct a graph of the higher-level fragments, see  $\text{HLG}(G)$  in Figure 3.

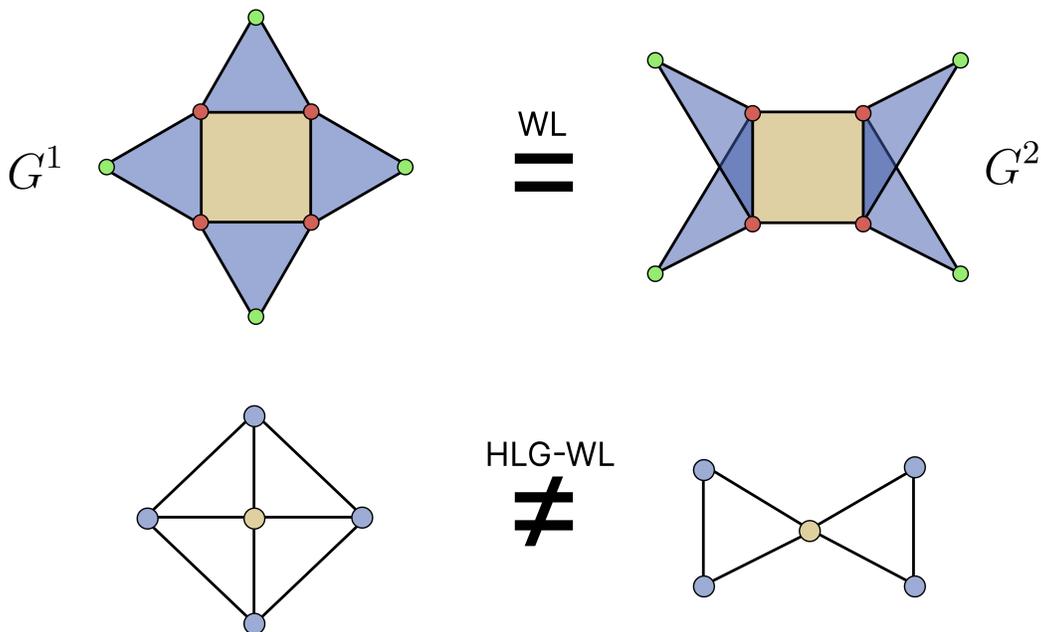
**Definition C.4.** The higher-level graph augmentation functions is  $\text{HLG}(\mathcal{V}, \mathcal{E}, \mathbf{X}) = (\mathcal{V}^{\text{HLG}}, \mathcal{E}^{\text{HLG}}, \mathbf{X}^{\text{HLG}})$  with:

$$\begin{aligned} \mathcal{V}^{\text{HLG}} &:= \mathcal{V}^{\text{FR}}, \quad \mathbf{X}^{\text{HLG}} := \mathbf{X}^{\text{FR}}, \\ \mathcal{E}^{\text{HLG}} &:= \mathcal{E}^{\text{FR}} \cup \{\{f, k\} \mid f, k \in F, f \cap k \neq \emptyset\} \end{aligned}$$

Equipped with the formal definitions, we will now compare the power of performing the WL test on these transformed graphs to the original graph<sup>3</sup>. The power of the Fragment-WL test depends on the fragmentation scheme  $\mathcal{F}$ . With a sufficiently advanced fragmentation scheme, even  $\text{NF-WL}$  can become arbitrarily powerful.

**Theorem C.5.** *There exist fragmentation schemes such that  $\text{NF-WL}$ ,  $\text{FR-WL}$  and  $\text{HLG-WL}$  are all strictly more powerful than  $k$ -WL for any  $k$ .*

<sup>3</sup>All proofs are detailed in Appendix D



**Figure 4:** Graphs  $G^1$  and  $G^2$  with their corresponding higher-level graph of fragments. The edges of the fragment representation to the vertices of  $G^1$  and  $G^2$  are omitted.  $G^1$  and  $G^2$  are indistinguishable by WL, NF-WL and FR-WL but distinguishable by HLG-WL as the higher-level graphs exhibit different connections from the 3-ring nodes to the 4-ring node.

However, in practice, mostly fragmentation schemes with a vocabulary of rings, paths, and cliques are used for fragment-biased GNNs [17; 4; 54]. So, from now on, we will restrict ourselves to such fragmentation in our theoretical analysis. Next, we show that it matters how to incorporate fragment information and that the WL variants become more powerful through higher-level abstraction.

Integrating fragment information from any non-trivial substructure as an additional node feature already increases expressiveness beyond 2-WL.

**Theorem C.6.** *NF-WL is strictly more powerful than 2-WL for fragmentation schemes  $\mathcal{F}$  that recover any substructure with more than two nodes.*

This shows that the classical 2-WL test cannot reveal differences in the expressivity of fragment-biased GNNs since using any substructure as a node feature already increases expressivity beyond 2-WL. Our Fragment-WL test, however, provides a more fine-grained alternative that reveals that higher-level abstraction through a learned fragment representation strictly increases the expressivity compared to node features:

**Theorem C.7.** *FR-WL is strictly more powerful than NF-WL for fragmentation schemes  $\mathcal{F}$  recovering 3-cycles.*

Building a higher-level graph of fragments further increases the expressivity. Figure 4 shows an example of two graphs that are indistinguishable by 2-WL, NF-WL, and FR-WL but distinguishable by HLG-WL. Formally, we express this in the following theorem.

**Theorem C.8.** *HLG-WL is strictly more powerful than FR-WL for fragmentation schemes  $\mathcal{F}$  recovering 3-cycles.*

Hence, we have shown that the expressivity increases strictly monotonically from 2-WL to HLG-WL for fragmentation schemes recovering 3-cycles:

$$2\text{-WL} < \text{NF-WL} < \text{FR-WL} < \text{HLG-WL} \quad (3)$$

Regarding the higher-dimensional  $k$ -WL hierarchy, our Fragment-WL hierarchy cannot be bounded by 3-WL if the fragmentation can recover 5 cycles.

**Theorem C.9.** *HLG-WL is in parts more powerful than 3-WL for fragmentation schemes  $\mathcal{F}$  recovering 5-cycles.*

Our developed Fragment-WL hierarchy models the different ways in which fragment information is used in most fragment-biased and topological GNNs. This new measure of expressiveness allows the comparison and ordering of these existing methods; see Table 2 for an overview.

In summary, our Fragment-WL test provides a new alternative measure of expressivity compared to the original WL test. Our developed hierarchy reveals that it matters how to incorporate fragmentation information, i.e., higher-level abstraction increases expressivity. Additionally, it allows for a comparison of the expressiveness of most existing fragment-biased and topological GNNs.

## D Proofs

This chapter presents the proofs for the theorems introduced in Appendix C, and the expressiveness analysis of existing fragment-biased and topological GNNs. We will first introduce general concepts that will later help to bound the expressiveness of different models and our Fragment WL test.

### D.1 Color refinement and expressiveness: Useful definitions and lemmas

To prove the power of different graph coloring algorithms, it will be useful to first introduce the definition of color refinement. The intuition is that a "finer" coloring contains more information than a "coarser" coloring.

**Definition D.1.** Let  $c, d$  be colorings of a graph  $G$ . The coloring  $c$  refines  $d$  (we write  $c \sqsubseteq d$ ) if there exists a function  $h$  such that  $h(c_v) = d_v$  for all  $v \in \mathcal{V}$ .

We will sometimes write  $c_v \sqsubseteq d_v$  if  $c \sqsubseteq d$  and the set of vertices  $\mathcal{V}$  is clear from the context.

**Example D.1.** Let  $c^{(t)}$  be the coloring of iteration  $t$  of the WL test. Then one can easily show that  $c_v^{(t)} \sqsubseteq c_v^{(l)}$  for  $l \leq t$  as  $c_v^{(t)}$  contains the information of all previous colorings  $c_v^{(l)}$ .

Note that an alternative definition of color refinement [2] is:  $c \sqsubseteq d$  iff  $c_v = c_w$  implies  $d_v = d_w$  for all  $v, w \in \mathcal{V}$ . It is easy to see that the two definitions are equivalent. Next, we extend our definition to arbitrary functions and not just colorings.

**Definition D.2.** Let  $a, b$  be two functions over the set of vertices  $\mathcal{V}$ . Then  $a \sqsubseteq b$  if there exists a function  $h$  such that  $h(a(v)) = b(v)$  for all  $v \in \mathcal{V}$ .

Intuitively,  $a \sqsubseteq b$  means that we can compute  $b(v)$  from the result of  $a(v)$ . So  $a(v)$  contains more or the same information as  $b(v)$ . Again, we will sometimes simply write  $a_v \sqsubseteq b_v$  with  $a_v := a(v)$ ,  $b_v := b(v)$  if  $a \sqsubseteq b$  and if the set of vertices  $\mathcal{V}$  is clear from the context.

**Example D.2.** Let  $c^{(t)}$  be the coloring of iteration  $t$  of the WL test. Then, because of the injectiveness of the hash function HASH in Equation (2):

$$c_v^{(t)} \sqsubseteq \{\{c_w^{(t-1)} \mid w \in \mathcal{N}(v)\}\}.$$

Note that we use the simplified notation here. The right and left-hand side are actually the functions that map from  $v \in \mathcal{V}$  to these terms. Intuitively, this shows that one can compute the previous color of all neighbors from the color of a node.

It is easy to see that the refinement relation is transitive, i.e.,  $a \sqsubseteq b$  and  $b \sqsubseteq c$  imply  $a \sqsubseteq c$ .

We will now formally define the expressive power of an algorithm with respect to the ability to distinguish non-isomorphic subgraphs.

**Definition D.3.** A function  $f$  is (in parts) *more powerful* than a function  $g$  if there exist two non-isomorphic graphs  $G^1, G^2$  such that  $f$  can distinguish them

$$f(G^1) \neq f(G^2)$$

whereas  $g$  cannot distinguish them

$$g(G^1) = g(G^2).$$

Note that this relation is *not* anti-symmetric, i.e.  $f$  can be (in parts) more powerful than  $g$ , and conversely,  $g$  can also be (in parts) more powerful than  $f$ . Hence, we introduce the following stronger anti-symmetric relation:

**Definition D.4.** A function  $f$  is *strictly more powerful* than  $g$  (denoted as  $f > g$ ) if

1.  $f$  is more powerful than  $g$
2. and  $g$  is not (in parts) more powerful than  $f$ .

Additionally, we write we write  $g \leq f$  if a function  $g$  is not more powerful than  $f$ .

Next, we will prove a connection between color refinement and expressiveness: a function that always produces a finer coloring cannot be less powerful than a function with a coarser coloring.

**Lemma D.5.** *Let  $f, g$  be functions with  $f \sqsubseteq g$  for all graphs. Then,  $g$  is not more powerful than  $f$ , i.e.,  $g \leq f$ .*

*Proof.* Assume for the sake of contradiction that there exist non-isomorphic graphs  $G^1, G^2$  that can be distinguished by  $g$  but not by  $f$ . Let  $d$  be the coloring obtained with  $g$ , and  $c$  be the coloring obtained with  $f$ . The multiset of colors  $d$  has to differ for  $G^1$  and  $G^2$ , i.e. there exists a color  $\alpha$  with

$$\begin{aligned} D_\alpha^1 &:= \{v \mid d_v = \alpha, v \in \mathcal{V}^1\} \\ D_\alpha^2 &:= \{v \mid d_v = \alpha, v \in \mathcal{V}^2\} \end{aligned}$$

such that

$$|D_\alpha^1| \neq |D_\alpha^2|. \quad (4)$$

Since  $c$  refines  $d$  no node  $v$  in  $D_\alpha^1$  and  $D_\alpha^2$  can share a color  $c_v$  with another node not in  $D_\alpha^1$  and  $D_\alpha^2$ . Hence, the set of colors  $c$  of nodes in  $D_\alpha^1$  and  $D_\alpha^2$  is disjoint from the set of colors  $c$  for the other nodes in the graph. But because of 4 there has to exist a color  $\beta$  with

$$\begin{aligned} C_\beta^1 &:= \{v \mid c_v = \beta, v \in \mathcal{V}^1\} \subseteq D_\alpha^1 \\ C_\beta^2 &:= \{v \mid c_v = \beta, v \in \mathcal{V}^2\} \subseteq D_\alpha^2 \end{aligned}$$

and

$$|C_\beta^1| \neq |C_\beta^2|.$$

This contradicts the initial assumption.  $\square$

Note that all our augmentation functions introduced in Appendix C only add information to the graph, or more formally:

**Definition D.6.** A function  $g$  from graphs to graphs is called *additive* if the set of nodes and edges does not decrease, i.e.,  $g(\mathcal{V}, \mathcal{E}, \mathbf{X}) = (\mathcal{V}', \mathcal{E}', \mathbf{X}')$  with  $\mathcal{V} \subseteq \mathcal{V}'$  and  $\mathcal{E} \subseteq \mathcal{E}'$ .

But by adding too much information, one could completely destroy the initial structure of the graph (e.g., make every graph a complete graph). Hence, we need the additional condition that it should be possible to recover the original graph from the augmented graph.

**Definition D.7.** An additive function  $g$  from graphs to graphs (i.e.  $g(G) = G'$ ) is called *reversible* if there exists a function  $h$  for vertices  $v \in \mathcal{V}'$  such that

$$h(\mathbf{X}'_v) = \begin{cases} \mathbf{X}_v & v \in \mathcal{V} \\ \perp & v \notin \mathcal{V} \end{cases} \quad (5)$$

and a function  $\varepsilon$  for edges  $e = \{u, v\} \in \mathcal{E}'$  such that

$$\varepsilon(\mathbf{X}'_u, \mathbf{X}'_v) = \begin{cases} 1 & e \in \mathcal{E} \\ 0 & e \notin \mathcal{E}. \end{cases} \quad (6)$$

By only adding such reversible information the WL test cannot become less powerful:

**Lemma D.8.** *Let  $g$  be a reversible function. Then  $g$ -WL is not less powerful than WL.*

*Proof.* Let  $c^{(t)}$  be the coloring obtained by the  $t$ -th iteration of the WL-test, and  $d^{(t)}$  the coloring of  $g$ -WL. We will show by induction that there exists a function  $h$  such that  $h(d_v^{(t)}) = c_v^{(t)}$  for  $v \in \mathcal{V}$ , i.e.  $d_v^{(t)} \sqsubseteq c_v^{(t)}$ . For  $t = 0$ , this follows immediately from Equation (5). For the induction step, note

$$d_v^{(t)} \sqsubseteq \{\{d_u^{(t-1)} \mid u \in \mathcal{N}_{G'}(v)\}\}.$$

Now note that the function  $\varepsilon$  (Equation (6)) makes it possible to reconstruct the neighborhood in  $G$  based on the features  $\mathbf{X}'$  and neighborhood in  $G'$ . Since  $d^{(t)}$  only refines the features  $\mathbf{X}'$ , we can also reconstruct the neighborhood in  $G$  based on  $d^{(t)}$  and, hence,

$$\{\{d_u^{(t-1)} \mid u \in \mathcal{N}_{G'}(v)\}\} \sqsubseteq \{\{d_u^{(t-1)} \mid u \in \mathcal{N}_G(v)\}\}$$

and by induction hypothesis

$$\{\{d_u^{(t-1)} \mid u \in \mathcal{N}_G(v)\}\} \sqsubseteq \{\{c_u^{(t-1)} \mid u \in \mathcal{N}_G(v)\}\}.$$

So, taken together, we have

$$d_v^{(t)} \sqsubseteq \{\{c_u^{(t-1)} \mid u \in \mathcal{N}_G(v)\}\}. \quad (7)$$

Additionally, note that

$$d_v^{(t)} \sqsubseteq d_v^{(t-1)} \stackrel{\text{IH}}{\sqsubseteq} c_v^{(t-1)}. \quad (8)$$

By combining Equation (7) and Equation (8), we get

$$d_v^{(t)} \sqsubseteq \left( c_v^{(t-1)}, \{\{c_u^{(t-1)} \mid u \in \mathcal{N}_G(v)\}\} \right) \sqsubseteq c_v^{(t)}.$$

Hence,  $d$  refines  $c$ , and by Lemma D.5  $g$ -WL is not less powerful than WL.  $\square$

Now that we have found a lower bound of the expressiveness, we will also give an upper bound of the expressiveness. The idea is that a graph augmentation function does not increase the power of a coloring algorithm  $f$  (e.g., the WL test) if all the information added by  $g$  can also be computed from the coloring obtained with  $f$ . Or, to put it differently, a graph augmentation function  $g$  does not increase the expressiveness if it is possible to compute the set of colors on  $g(G)$  from the set of colors on  $G$ .

**Lemma D.9.** *Let  $g$  be a function that augments a graph, i.e., a function from graphs to graphs. Let  $f$  be a coloring function. If there exists a function  $h$  such that*

$$f(g(G)) = h(f(G))$$

*then  $f \circ g$  is not more powerful than  $f$ .*

*Proof.* Let  $G^1, G^2$  be two non-isomorphic graphs that are distinguishable by  $f \circ g$ . Then

$$\begin{aligned} f(g(G^1)) &\neq f(g(G^2)) \\ h(f(G^1)) &\neq h(f(G^2)) \end{aligned}$$

It follows that  $f(G^1) \neq f(G^2)$ . Hence,  $f \circ g$  is not more powerful than  $f$ .  $\square$

Note that the condition in the lemma is similar to the definition of color refinement. However, we cannot use color refinement directly because the function  $g$  could add or delete nodes, making a direct comparison of nodes between  $G^1$  and  $G^2$  impossible. Consequently, we have to use a more global view rather than the more localized approach of color refinement.

## D.2 Graph augmentation functions

We will now analyze the change in expressiveness with some graph augmentation functions that model message-passing schemes that are frequently used in practice.

### D.2.1 Fragment augmentations

We will first give the proofs for the augmentation functions from Appendix C that incorporate fragment information.

**Proof of Theorem C.5.**

**Theorem C.5.** *There exist fragmentation schemes such that NF-WL, FR-WL and HLG-WL are all strictly more powerful than  $k$ -WL for any  $k$ .*

*Proof.* Consider a fragmentation scheme that decomposes a graph into every possible subgraph (so the vocabulary is the set of all possible graphs). It is known that for every  $k$  there exist two non-isomorphic graphs  $G^1$  and  $G^2$  that are indistinguishable by  $k$ -WL. But as  $G^1$  and  $G^2$  are themselves part of the vocabulary and the fragmentation, they are trivially distinguishable by NF-WL, FR-WL, and HLG-WL.  $\square$

**Proof of Theorem C.6.**

**Theorem C.6.** *NF-WL is strictly more powerful than 2-WL for fragmentation schemes  $\mathcal{F}$  that recover any substructure with more than two nodes.*

*Proof.* Chen et al. [7] showed that the WL test cannot count the number of (induced) subgraphs with at least three nodes, i.e. for any substructure  $S$  with more than three nodes, there exist non-isomorphic graphs  $G^1$  and  $G^2$  such that 2-WL cannot distinguish them but they have a different count of  $S$ . So, let  $X$  be the substructure that  $\mathcal{F}$  recovers. Then, WL cannot count the number of occurrences of  $X$ , whereas NF-WL can trivially count the number of  $X$  in a graph. Hence, NF-WL is more powerful than WL. Since NF is a reversible function, NF-WL is by Lemma D.8 also not less powerful than WL, making NF-WL strictly more powerful than WL.  $\square$

**Proof of Theorem C.7.** Before coming to the proof of Theorem C.7, we will prove the following useful lemma:

**Lemma D.10.** *Two graphs  $G^1 = (\mathcal{V}^1, \mathcal{E}^1, \mathbf{X}^1)$ ,  $G^2 = (\mathcal{V}^2, \mathcal{E}^2, \mathbf{X}^2)$  are undistinguishable by WL if the set of node features is the same*

$$\mathbf{X}^1 = \mathbf{X}^2 := \mathbf{X}$$

and all nodes with the same node feature have the same neighborhood

$$\begin{aligned} \forall i, j \in \mathcal{V}^1 \cup \mathcal{V}^2 : \\ \mathbf{X}_i = \mathbf{X}_j \Rightarrow \{\mathbf{X}_n \mid n \in \mathcal{N}(i)\} = \{\mathbf{X}_m \mid m \in \mathcal{N}(j)\}. \end{aligned} \quad (9)$$

*Proof.* We will show by induction over  $t$  that the color of all nodes with the same node features is the same:

$$\begin{aligned} \forall i, j \in \mathcal{V}^1 \cup \mathcal{V}^2 : \\ \mathbf{X}_i = \mathbf{X}_j \Rightarrow c_i^t = c_j^t. \end{aligned}$$

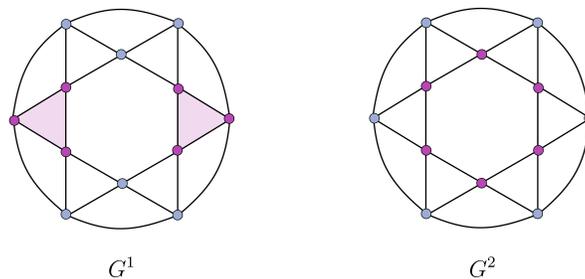
For  $t = 0$ , this follows immediately from  $c_i^0 = \text{HASH}(\mathbf{X}_i)$ .

For  $t > 0$ , we have for nodes  $i, j$  with  $\mathbf{X}_i = \mathbf{X}_j$

$$\begin{aligned} c_i^t &= \text{HASH}(c_i^{t-1}, \{c_n^{t-1} \mid n \in \mathcal{N}(i)\}) \\ &= \text{HASH}(c_j^{t-1}, \{c_n^{t-1} \mid n \in \mathcal{N}(i)\}) && \text{(by IH)} \\ &= \text{HASH}(c_j^{t-1}, \{c_m^{t-1} \mid m \in \mathcal{N}(j)\}) && \text{(by 9 and IH)} \\ &= c_j^t \end{aligned}$$

As both graphs have the same node features  $\mathbf{X}^1 = \mathbf{X}^2$  the set of colors is also the same in each iteration  $t$  of the WL-test. Hence, the graphs are indistinguishable by the WL-test  $\square$

**Theorem C.7.** *FR-WL is strictly more powerful than NF-WL for fragmentation schemes  $\mathcal{F}$  recovering 3-cycles.*



**Figure 5:** Graph  $G^1$  and graph  $G^2$  that are indistinguishable by NF-WL but distinguishable by FR-WL. Node features are represented by the color of the nodes.

*Proof.* With Lemma D.10 we can now prove Theorem C.7: We first show that FR-WL is more powerful than NF-WL. Consider the two graphs  $G^1, G^2$  depicted in Figure 5 with two different node features colored violet and blue. Note that when not considering the node features the graphs are identical and each node is isomorphic to every other node. As every fragmentation scheme has to be permutation invariant, every node is assigned the same additional node feature in NF (which, thus, holds no additional information to distinguish the graphs). Now, observe that the two graphs fulfill the conditions of Lemma D.10. Therefore, the graphs are indistinguishable by NF-WL. In contrast, FR-WL can distinguish the two graphs as  $G^1$  contains 3-cycles with three violet nodes whereas  $G^2$  does not contain such 3-cycles. Hence, the fragment representation for these 3-cycles differs and distinguishes the two graphs.

Now, it remains to show that NF-WL is not more powerful than FR-WL. Let  $c_v^t, d_v^t$  be the colorings of the  $t$ -th iteration of FR-WL and NF-WL, respectively. We will show by induction over  $t$  that  $c_v^{t+1} \sqsubseteq d_v^t$ . To simplify the notation, let  $\mathcal{N}_\uparrow(v) := \{f \mid v \in f, f \in F\}$  be the set of fragments that  $v$  is part of. For  $t = 0$ , note that after the first iteration of the FR-WL test each vertex  $v \in G$  receives information about the fragment it is part of

$$\begin{aligned}
 c_v^1 &= \text{HASH}(c_v^0, \{\{c_n^0 \mid n \in \mathcal{N}_G(v) \uplus \mathcal{N}_\uparrow(v)\}\}) \\
 &= \text{HASH}(c_v^0, \{\{c_n^0 \mid n \in \mathcal{N}(v)\}\} \uplus \{\{c_f^0 \mid f \in \mathcal{N}_\uparrow(v)\}\}) \\
 &\sqsubseteq (c_v^0, \{\{c_f^0 \mid f \in \mathcal{N}_\uparrow(v)\}\}) \\
 &= (\mathbf{X}_v, \{\{\text{type}(f) \mid f \in \mathcal{N}_\uparrow(v)\}\}) \\
 &\sqsubseteq d_v^0
 \end{aligned}$$

For the induction step  $(t-1) \rightarrow t$ , we have

$$\begin{aligned}
 c_v^t &= \text{HASH}(c_v^{t-1}, \{\{c_n^{t-1} \mid n \in \mathcal{N}_G(v) \uplus \mathcal{N}_\uparrow(v)\}\}) \\
 &\sqsubseteq (c_v^{t-1}, \{\{c_n^{t-1} \mid n \in \mathcal{N}_G(v)\}\}) \\
 &\sqsubseteq (d_v^{t-2}, \{\{d_n^{t-2} \mid n \in \mathcal{N}_G(v)\}\}) \\
 &\sqsubseteq d_v^{t-1}
 \end{aligned}$$

This concludes the induction step. Hence, we have  $c \sqsubseteq d$ , and by Lemma D.8 NF-WL is not more powerful than FR-WL.  $\square$

### Proof of Theorem C.8.

**Theorem C.8.** *HLG-WL is strictly more powerful than FR-WL for fragmentation schemes  $\mathcal{F}$  recovering 3-cycles.*

*Proof.* Consider the two graphs  $G^1, G^2$  depicted in Figure 4 with two different node features colored green and red. Note that the two graphs fulfill the conditions of Lemma D.10. Furthermore, even the graphs  $\text{FR}(G^1)$  and  $\text{FR}(G^2)$  fulfill these conditions. Hence, the graphs  $G^1, G^2$  cannot be distinguished by FR-WL.

In the higher level graph of  $G^1$  each 3-cycle representation is connected to two other 3-cycle representations. Contrarily, in the higher level graph of  $G^2$  each 3-cycle representation is connected

to only one 3-cycle representations. Hence, the coloring will differ and HLG-WL distinguishes the two graphs.

Now, it remains to show that FR-WL is not more powerful than HLG-WL. This follows immediately from Lemma D.8 and the fact that the change in the graph from FR to HLG is reversible.  $\square$

**Proof of Theorem C.9.**

**Theorem C.9.** *HLG-WL is in parts more powerful than 3-WL for fragmentation schemes  $\mathcal{F}$  recovering 5-cycles.*

*Proof.* The proof follows a similar proof by Bodnar et al. [3]. Consider the Rook’s 4x4 and Shrikhande graph (both in the family of strongly regular graphs  $\text{SR}(16, 6, 2, 2)$ ). The Shrikhande graph possesses 5-cycles while the Rook’s graph does not. Hence, HLG-WL can trivially distinguish those two graphs with a fragmentation recovering 5-cycles. However, it is known that 3-WL cannot distinguish those two graphs [3].  $\square$

### D.3 Additional graph augmentation function

We will now consider two additional graph augmentation functions that are often used in practice: a learned representation for each edge and a learned representation for the complete graph that is connected to all other nodes. While these augmentations might be beneficial in practice, we will show that they do not increase expressiveness.

**Edge representation.** We will first formally define the edge representation augmentation (ER). For every edge, we introduce a new node that is connected to its two endpoints.

**Definition D.11.** The edge representation graph augmentation function is  $\text{ER}(\mathcal{V}, \mathcal{E}, \mathbf{X}) = (\mathcal{V}^{\text{ER}}, \mathcal{E}^{\text{ER}}, \mathbf{X}^{\text{ER}})$  with

$$\begin{aligned}\mathcal{V}^{\text{ER}} &:= \mathcal{V} \cup \mathcal{E}, \\ \mathcal{E}^{\text{ER}} &:= \mathcal{E} \cup \{\{e, v\} \mid e = \{u, v\} \in \mathcal{E}\} \\ \mathbf{X}_i^{\text{ER}} &:= \begin{cases} \mathbf{X}_i & i \in \mathcal{V} \\ \alpha & i \in \mathcal{E} \end{cases}\end{aligned}$$

where  $\alpha$  is some new label.

Now, we can show that this augmentation does not increase expressiveness compared to the WL test.

**Lemma D.12.** *ER-WL is as powerful as WL.*

*Proof.* We will first show that ER-WL is not more powerful than WL: Let  $c^{(i)}, d^{(i)}$  be the colorings of the WL test and of the ER-WL test, respectively. Then we will show that

We will use Lemma D.9 and give a function  $h$  that maps a coloring  $c^{(i)}$  of the WL test without edge representation to a coloring  $d^{(i)}$  of ER-WL: Note that with the color  $c_v^{(i)}$  of a node  $v$  one can compute the colors  $c_u^{(i-1)}$  of all neighboring nodes  $u$ . Hence, we can determine from  $c^{(i)}$  the following multiset

$$\{\{(c_v^{(i)}, c_u^{(i-1)}) \mid e = (u, v) \in \mathcal{E}\}\}$$

which allows us to compute the corresponding edge representations  $d_e^{(i)}$ .

We will now show that ER-WL is not less powerful than WL. This follows directly from Lemma D.8 as ER is a reversible function.  $\square$

**Graph representation.** We will now formally define the learned graph representation (GR), sometimes called virtual node.

**Definition D.13.** The graph representation augmentation function is  $\text{GR}(\mathcal{V}, \mathcal{E}, \mathbf{X}) = (\mathcal{V}^{\text{GR}}, \mathcal{E}^{\text{GR}}, \mathbf{X}^{\text{GR}})$  with

$$\begin{aligned}\mathcal{V}^{\text{GR}} &:= \mathcal{V} \cup \{g\}, \\ \mathcal{E}^{\text{GR}} &:= \mathcal{E} \cup \{\{v, g\} \mid v \in \mathcal{V}\} \\ \mathbf{X}_i^{\text{GR}} &:= \begin{cases} \mathbf{X}_i & i \in \mathcal{V} \\ \alpha & i = g \end{cases}\end{aligned}$$

**Table 2:** Overview of the vocabulary and expressiveness of existing topological and fragment-biased models. The bounds for GSN-v,  $\mathcal{O}$ -GNN, and HIMP are tight, i.e. when using a sufficient number of layers and injective neighborhood aggregators, the models are as powerful as the corresponding Fragment-WL test.

Model	Bounded by	Vocabulary
GSN-v <sup>4</sup>	$\leq$ NF-WL	Cliques or Rings
$\mathcal{O}$ -GNN	$\leq$ FR-WL	Rings
HIMP	$\leq$ HLG-WL	Rings
MPSN	$\leq$ HLG-WL	Simplicial complexes (in practice cliques)
CIN	$\leq$ HLG-WL	CW complexes (in practice rings & edges)
CIN++	$\leq$ HLG-WL	CW complexes (in practice rings & edges)

where  $\alpha$  is some new label.

Similar to the edge representation, the graph representation does not increase expressiveness:

**Lemma D.14.** *GR-WL is as powerful as WL.*

*Proof.* We will first show that GR-WL is not more powerful than WL. We will use Lemma D.9 and give a function  $h$  that maps a coloring  $c^{(i)}$  of the WL test without graph representation to a coloring  $d^{(i)}$  of GR-WL: We will show this by induction over  $i$ . For  $i = 0$ , this follows immediately from the definition of GR. For the induction step, assume that there exists such a function from  $c^{(t-1)}$  to  $d^{(t-1)}$ . Note that the graph representation  $d_g^{(t)}$  is computed as:

$$\begin{aligned} d_g^{(t)} &= \text{HASH}(d_g^{(t-1)}, \{\{d_v^{(t-1)} \mid v \in \mathcal{N}_{G'}(g)\}\}) \\ &= \text{HASH}(d_g^{(t-1)}, \{\{d_v^{(t-1)} \mid v \in \mathcal{V}\}\}). \end{aligned}$$

which can be derived from  $c^{(t-1)}$  by induction hypothesis. With this we can trivially compute  $d_v^{(t)}$  from  $c^{(t)}$  for all other nodes  $v$ , too.

Since GR is a reversible function, by Lemma D.8 GR-WL is not less powerful than WL. Hence, GR-WL is as powerful as WL.  $\square$

## D.4 Expressiveness of existing models

We will use our Fragment-WL tests to compare the expressiveness of existing fragment-biased GNN models. Table 2 gives an overview of the vocabulary of existing fragment-biased and topological GNNs. Additionally, it shows the expressiveness in our Fragment-WL hierarchy.

### D.4.1 GSN-v

GSN-v [4] incorporate fragment information as an additional node feature. The additional node features consist of the counts of fragment types a node is part of. Their framework also differentiates between different (non-symmetric) positions inside the fragment (e.g., first node in path vs. second node in path) that correspond to different orbits. While their framework can use any fragmentation scheme, in all real-world experiments, they only use rings or cliques. Note that for rings and cliques, no different orbits exist, i.e., each node in the substructure has the same orbit. Hence, this information becomes irrelevant.

<sup>4</sup>We evaluate the GSN-v that is used in practice, i.e. with a vocabulary of cliques and rings. Note that the theory of the authors allows for potentially more expressive instantiations.

**Theorem D.15.** *GSN-v using rings and/or cliques as vocabulary is at most as powerful as NF-WL. Additionally, when using injective neighborhood aggregators and a sufficient number of layers, GNS are as powerful as NF-WL with a fragmentation scheme based on rings and cliques.*

*Proof.* GSN-v appends the node features by the counts of substructures and the respective orbits each node is part of. After that, a standard GNN is applied to the graph.

Note that in a ring or a clique, each node has exactly the same orbit. So, for a vocabulary based on cliques and rings the appended information degenerates to solely the substructure counts. Further, note that this substructure count function is an injective function  $\lambda$  as defined in Definition C.2. Hence, when using injective neighborhood aggregators and an MLP update function with a sufficiently large number of layers such that it can approximate the HASH function, GSN-v exactly models the NF-WL test. Hence, GSN-v is exactly as powerful as NF-WL.  $\square$

#### D.4.2 $\mathcal{O}$ -GNNs

Besides representations for nodes,  $\mathcal{O}$ -GNNs [54] use explicit representation for rings, edges, and the whole graph.

**Theorem D.16.**  *$\mathcal{O}$ -GNNs [54] are at most as powerful as FR-WL. Additionally, when using injective neighborhood aggregators and a sufficient number of layers,  $\mathcal{O}$ -GNNs are as powerful as FR-WL with a fragmentation scheme based solely on rings.*

*Proof.*  $\mathcal{O}$ -GNNs (when using injective neighborhood aggregators instead of their original sum aggregators and an MLP with a sufficiently large number of layers such that it can approximate the HASH function) models performing the WL test on an FR, ER, and GR augmented graph. As shown in Lemmas D.12 and D.14, the edge representation and the graph representation do not influence the expressivity. Hence,  $\mathcal{O}$ -GNNs are exactly as powerful as FR-WL with a fragmentation scheme based solely on rings.  $\square$

#### D.4.3 HIMP

HIMP [17] builds a higher-level junction tree based on rings and edges for message passing on the original graph, the higher-level junction tree, and between those two.

**Theorem D.17.** *HIMP is at most as powerful as HLG-WL. Additionally, when using injective neighborhood aggregators and a sufficient number of layers, HIMP is as powerful as HLG-WL with a fragmentation scheme based on rings and edges.*

*Proof.* HIMP (when using injective neighborhood aggregators instead of their original sum aggregators and an MLP with a sufficiently large number of layers such that it can approximate the HASH function) exactly models performing the WL test on an HLG augmented graph. Hence, HIMP is exactly as powerful as FR-WL with a fragmentation scheme based on rings and edges.  $\square$

#### D.4.4 FragNet

Next, we consider the expressiveness of our FragNet model:

**Theorem E.1.** *FragNets are at most as powerful as HLG-WL. Additionally, when using injective neighborhood aggregators and a sufficient number of layers, FragNets are as powerful as HLG-WL.*

*Proof.* For the proof, we rely on Lemma D.12, the finding that an explicit edge representation does not augment expressiveness. Notice that our model, when using injective neighborhood aggregators and an MLP with a sufficiently large number of layers such that it can approximate the HASH function, exactly models performing the WL test on an HLG and ER augmented graph. As shown in Lemma D.12, ER does not change the expressiveness. Hence, our model is exactly as powerful as HLG-WL.  $\square$

### D.4.5 Topological GNNs

We will now consider topological GNNs. We will start by comparing HLG-WL with CWL, a variant of the WL test operating on CW complexes [3]. In the CWL framework, every graph is (permutation invariantly) mapped to a set of cells  $\mathcal{X}$ , a CW complex (using a skeleton-preserving lifting map). Let  $\mathcal{X}_i$  denote the set of cells with dimension  $i$ . Then,  $\mathcal{X}_0$  corresponds to all vertices  $\mathcal{V}$  and  $\mathcal{X}_1$  to all edges  $\mathcal{E}$ . For higher dimensions, the results depend on the particular cellular lifting map. For instance,  $\mathcal{X}_2$  could correspond to all cycles.

**Theorem D.18.** *HLG-WL is not less powerful than CWL, with a fragmentation scheme  $\mathcal{F}$  that corresponds to the cellular lifting map used by CWL.*

*Proof.* Let  $F$  be the fragmentation that corresponds to  $\mathcal{X}$  without the vertices:  $F = \mathcal{X} \setminus \mathcal{V}$ . Let  $c^{(t)}, b^{(t)}$  be the coloring of iteration  $t$  of HLG-WL and CWL, respectively. We will show that  $b^{(t)} \sqsubseteq c^{(2t)}$  which implies that  $b \sqsubseteq c$  after termination.

We will show this by induction over  $t$ . For  $t = 0$ , this follows immediately from the fact that the node features in HLG-WL are finer than the features of cells in CWL.

Now we will show  $c^{(t)} \sqsubseteq b^{(2t)}$  assuming  $c^{(t-1)} \sqsubseteq b^{(2t-2)}$ :

The idea of the induction step is that the hash update function HASH receives more information in HLG-WL compared to CWL. Let us first consider vertices, i.e.,  $\mathcal{X}_0$ . The update function in CWL for  $v \in \mathcal{V} = \mathcal{X}_0$  is:

$$b_v^{(t)} = \text{HASH}(b_v^{(t-1)}, \{\{(b_w^{(t-1)}, b_e^{(t-1)}) \mid e = \{v, w\} \in \mathcal{E}\}\})$$

Now note that the update function in HLG-WL for  $v \in \mathcal{V}$  is:

$$\begin{aligned} c_v^{(2t)} &= \text{HASH}(c_v^{(2t-1)}, \{\{c_w^{(2t-1)} \mid w \in \mathcal{N}_{\text{HLG}(G)}(v)\}\}) \\ &\sqsubseteq \text{HASH}(c_v^{(2t-1)}, \{\{c_e^{(2t-1)} \mid e = \{v, w\} \in \mathcal{E}\}\}) \\ &\sqsubseteq \text{HASH}(c_v^{(2t-1)}, \{\{(c_e^{(2t-2)}, c_w^{(2t-2)}) \mid e = \{v, w\} \in \mathcal{E}\}\}) \\ &\sqsubseteq \text{HASH}(b_v^{(t-1)}, \{\{(b_e^{(t-1)}, b_w^{(t-1)}) \mid e = \{v, w\} \in \mathcal{E}\}\}) \\ &= b_v^{(t)} \end{aligned}$$

The first step follows from  $e \in \mathcal{N}_{\text{HLG}(G)}(v)$  as the edges are part of the fragmentation  $F$ . The second step follows from  $c_e^{(2t-1)} \sqsubseteq (c_e^{(2t-2)}, c_w^{(2t-2)}, c_v^{(2t-2)})$ . The third step uses the induction hypothesis.

Now, we will consider a cell  $x \in \mathcal{X}_k \subseteq F$ . The update function in CWL is

$$b_x^{(t)} = \text{HASH}(b_x^{(t-1)}, \{\{(b_u^{(t-1)}, b_o^{(t-1)}) \mid x \prec u, o \prec u\}, \{\{b_l^{(t-1)} \mid l \prec x\}\}\})$$

where  $x \prec y$  means that  $x$  with dimension  $k$  is part of the cell  $y$  of dimension  $k + 1$ . For example,  $e \prec r$  if  $e$  is an edge in a ring  $r \in \mathcal{X}_{k+1}$ . For details, we refer to Bodnar et al. [2].

The update function in HLG-WL for a fragment  $x \in F \cap \mathcal{X}_k$  in  $G' := \text{HLG}(G)$  is:

$$\begin{aligned} c_x^{(2t)} &= \text{HASH}(c_x^{(2t-1)}, \{\{c_w^{(2t-1)} \mid w \in \mathcal{N}_{G'}(x)\}\}) \\ &\sqsubseteq \text{HASH}(c_x^{(2t-1)}, \{\{c_u^{(2t-1)} \mid u \in \mathcal{N}_{G'}(x) \cap \mathcal{X}_{k+1}\}, \{\{c_l^{(2t-1)} \mid l \in \mathcal{N}_{G'}(x) \cap \mathcal{X}_{k-1}\}\}\}) \\ &\sqsubseteq \text{HASH}(c_x^{(2t-1)}, \{\{c_u^{(2t-1)} \mid x \prec u\}, \{\{c_l^{(2t-1)} \mid l \prec x\}\}\}) \\ &\sqsubseteq \text{HASH}(c_x^{(2t-1)}, \{\{(c_u^{(2t-2)}, c_o^{(2t-2)}) \mid x \prec u, o \prec u\}, \{\{c_l^{(2t-1)} \mid l \prec x\}\}\}) \\ &\sqsubseteq \text{HASH}(b_x^{(t-1)}, \{\{(b_u^{(t-1)}, b_o^{(t-1)}) \mid x \prec u, o \prec u\}, \{\{b_l^{(t-1)} \mid l \prec x\}\}\}) \\ &= b_x^{(t)} \end{aligned}$$

The steps are very similar to the vertex case above. This concludes the proof that  $c$  refines  $b$ . By Lemma D.5 this implies that HLG-WL is not less powerful than CWL using a fragmentation that corresponds to the cellular complex.  $\square$

As CWL bounds the expressiveness of CIN [3] and CIN++ [21], we get the following corollary:

**Corollary D.19.** *CIN and CIN++ are at most as powerful as HLG-WL with a fragmentation scheme that corresponds to the cellular lifting map.*

Additionally, CWL subsumes the WL version, SWL, introduced by Bodnar et al. [2] for simplicial complexes. The cellular complex just corresponds to all cliques of the graph.

**Corollary D.20.** *HLG-WL, with a fragmentation scheme recovering cliques, is not less powerful than SWL.*

As MPSNs [2] are bounded by SWL, we have the following result for MPSNs:

**Corollary D.21.** *MPSNs are at most as powerful as HLG-WL with a fragmentation scheme based on cliques.*

## E Model

In the following we explicitly describe the update and output mechanisms of our model.

The learned representation  $h_v^t$  at step or layer  $t$  of a node receives a message from neighboring nodes  $m_{\mathcal{V} \rightarrow v}^t$  and fragments  $m_{F \rightarrow v}^t$ . Similarly, the learned representation  $h_f^t$  of a fragment receives a message from neighboring fragments  $m_{F \rightarrow f}^t$  and nodes  $m_{\mathcal{V} \rightarrow f}^t$ .

$$\begin{aligned} m_{\mathcal{V} \rightarrow v}^t &= \text{AGG}(\{\{\text{MLP}(h_u^{t-1}, h_e^{t-1}) \mid e = \{u, v\} \in \mathcal{E}\}\}) & m_{F \rightarrow v}^t &= \text{AGG}(\{\{h_f^{t-1} \mid v \in f, f \in F\}\}) \\ m_{F \rightarrow f}^t &= \text{AGG}(\{\{h_g^{t-1} \mid g \in N_F(f)\}\}) & m_{\mathcal{V} \rightarrow f}^t &= \text{AGG}(\{\{h_v^{t-1} \mid v \in f, v \in \mathcal{V}\}\}) \end{aligned}$$

where  $N_F(f)$  denotes the neighbors of fragment  $f$  in the higher-level graph of fragments. The hidden representations are updated by combining the incoming messages with the previous hidden representation:

$$h_v^t = \text{MLP}(h_v^{t-1}, m_{\mathcal{V} \rightarrow v}^t, m_{F \rightarrow v}^t) \quad h_f^t = \text{MLP}(h_f^{t-1}, m_{\mathcal{V} \rightarrow f}^t, m_{F \rightarrow f}^t) \quad h_e^t = \text{MLP}(h_e^{t-1}, h_u^{t-1}, h_v^{t-1})$$

The final graph-level readout after  $T$  layers is computed by aggregating the multiset of node representation, edge representations, and fragment representations:

$$\text{OUT}(\{\{h_v^T \mid v \in \mathcal{V}\}\}, \{\{h_e^T \mid e \in \mathcal{E}\}\}, \{\{h_f^T \mid f \in F\}\})$$

Note that the complexity of our FragNet model is linear in the number of nodes and fragments (assuming that each node is only part of a constant number of fragments).

Additionally, our FragNet model achieves the highest expressiveness in our Fragment-WL hierarchy and also compared to other fragment-biased GNNs.

**Theorem E.1.** *FragNets are at most as powerful as HLG-WL. Additionally, when using injective neighborhood aggregators and a sufficient number of layers, FragNets are as powerful as HLG-WL.*

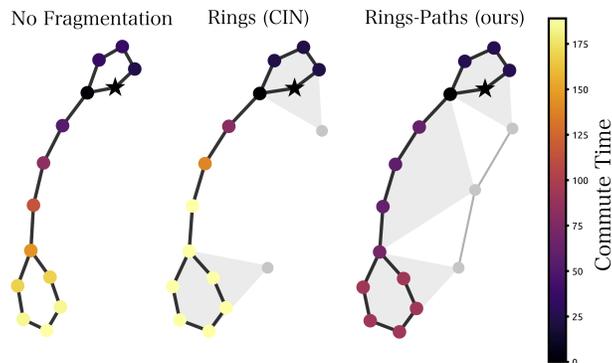
## F Further experiments

### F.1 Long-Range tests

We provide more experiments to measure the long-range capabilities of FragNet.

#### Commute time

In addition to the recovery rates in ??, we also consider commute times. The commute time between the nodes  $a$  and  $b$  is the expected time for a random walker from  $a$  to reach  $b$  and return again to  $a$ . Di Giovanni et al. [11] have proposed the commute time as a measure for over-squashing. To compute and compare commute times across different fragmentations, we connected all nodes in each fragmentation that could exchange a message within one layer. Figure 6 shows the commute from the star node to every other node for the same graph as in ?. The close alignment between commute time and recovery rate supports the theoretical findings by Di Giovanni et al. [11] and further emphasizes the potentially enhanced long-range capabilities of our model. Additionally, we



**Figure 6:** Commute time from the star node to all other nodes. The first graph has no fragmentation, the second one a rings fragmentation (like in CIN/CIN++), the third a rings and paths fragmentation (like our model).

also compute commute times on a molecule from the ZINC dataset that contains more fragments (see Figure 7).

Quantitatively, we compute average commute times on a random sample of molecules from the peptides dataset for a model without any fragmentation and for FragNet (RingsPaths fragmentation with a higher-level graph). We observe that the addition of a higher-level graph reduces commute times by 16%.

**Table 3:** Average commute times between all nodes on a random sample of 50 molecular graphs from the peptides dataset with and without the higher-level graph.

Normal Molecular Graph	5056
Molecular Graph + HLG	4253

## F.2 Distribution of Fragments

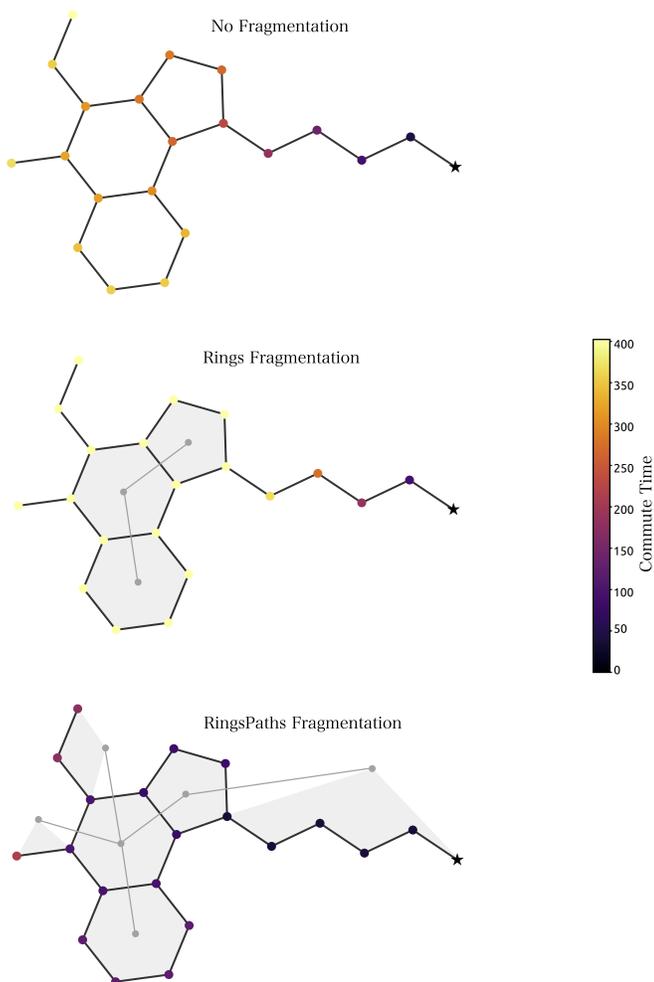
Figure 8 illustrates the distribution of fragment sizes, i.e., path lengths and ring lengths, extracted by our RingsPaths fragmentation method across the ZINC-10k, ZINC-full, and peptides datasets. It is worth noting that the peptides dataset features some exceptionally large rings.

## F.3 Ablation Studies

In the following, we test the design choices of our model and fragmentation. First, we test FragNet without the different fragment information or ordinal encoding on ZINC and Peptides. We show the results in Table 4. We observe that a reduction in expressiveness generally leads to a reduction in performance. An exception is the use of fragment representations (FR-WL, i.e. FragNet - Higher-level graph) which shows a higher error on ZINC 10k and Peptides Struct. This is similar to the pattern that we observe in the message reconstruction toy experiment we show in ?? where the additional fragment representations increase the importance of the current substructure and do not contain a message from other parts of the molecule.

Furthermore, we compare different fragmentation schemes in combination with FragNet. In Table 5, we observe that our RingsPath fragmentation scheme performs the best across the different datasets.

In Figure 9, we look at how large the vocabulary size has to be per fraction of fragmented atoms. That is, for an increasing vocabulary, we observe how many atoms belong to a fragment. The steeper the increase the better. We can observe that on ZINC-10k BBB, BRICS and Rings are not able to assign a fragment to each atom no matter how large the vocabulary size. Magnet achieves full fragmentation but slower compared to our RingsPaths which is the most vocabulary efficient. We further show the necessary vocabulary size on ZINC-Full and Peptides for RingsPaths in Table 6.



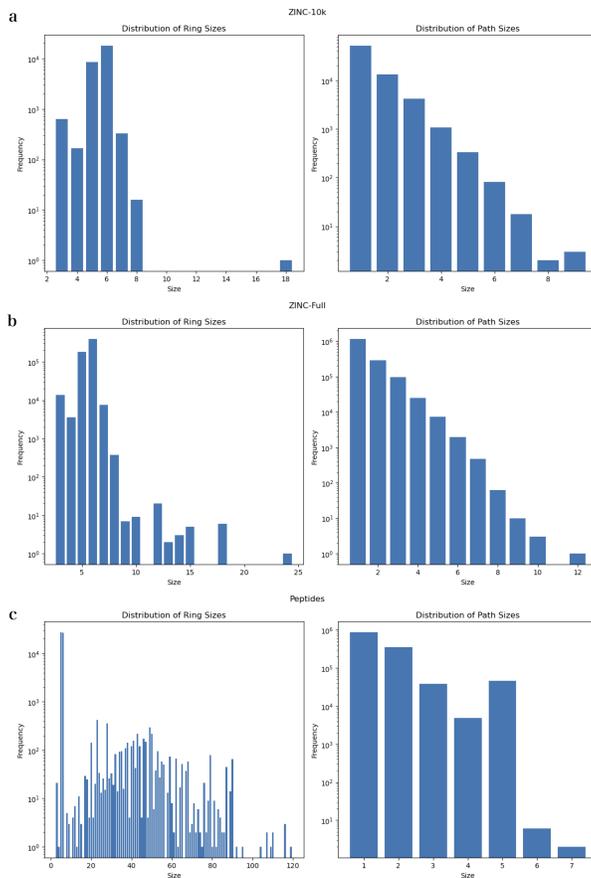
**Figure 7:** Commute time from the star node to all other nodes. The first graph has no fragmentation, the second one a rings fragmentation (like in CIN/CIN++), the third a rings and paths fragmentation (like our model).

## G Experimental details

In the following, we will describe details for all our experiments. Unless otherwise stated, for our FragNet, we use a 2-layer fully connected neural network with ReLU activations and batch norm as the MLP update function. For the aggregation method AGG, we use a sum aggregation for messages within the original or higher-level graph and a mean aggregation for messages between the original graph and the higher-level graph. For training, we use the AdamW [35] optimizer and gradient clipping with a value of 1. The model has been implemented in PyTorch [41] using the PyTorchGeometric [16] and the PyTorch Lightning [15] library. It is in parts adapted from HIMP [17]. All results for other methods are taken from Rampášek et al. [43] and Giusti et al. [21].

### G.1 ZINC and peptides

We compare our model against standard GNNs, like **GCN** [32], **GIN** [49] or **GatedGCN** [5], higher-order GNNs such as **RingGNN** [8] or **3WLGNN** [37], topological GNNs, especially **CIN** and **CIN++** [2; 3] and other fragment-biased GNNs such as **HIMP** [17] and **GSN** [4]. While our main focus lies on the evaluation against other message-passing GNNs, especially the group of fragment-biased GNNs, we compare against Transformer architectures for completeness. We test **Graphormer** [50],



**Figure 8:** Distribution of sizes of path and ring fragments for the a) ZINC-10k, b) ZINC-full, and c) Peptides dataset.

**Table 4:** Ablation of different expressivity choices for FragNet. Additionally, we ablate the ordinal encoding.

Model	ZINC	Peptides	
	10k (MAE ↓)	Struct (MAE ↓)	Func (AP ↑)
FragNet (=HLG-WL)	<b>0.0775</b> ± 0.004	<b>0.246</b> ± 0.002	<b>0.668</b> ± 0.003
– Higher-level graph (=FR-WL)	0.0872 ± 0.004	0.256 ± 0.003	0.661 ± 0.005
– Fragment representation (=NF-WL)	0.0994 ± 0.007	0.247 ± 0.003	0.654 ± 0.005
– All fragment information (=WL)	0.1609 ± 0.003	0.249 ± 0.001	0.652 ± 0.005
FragNet – ordinal encoding	0.0945 ± 0.006	0.249 ± 0.001	0.666 ± 0.004

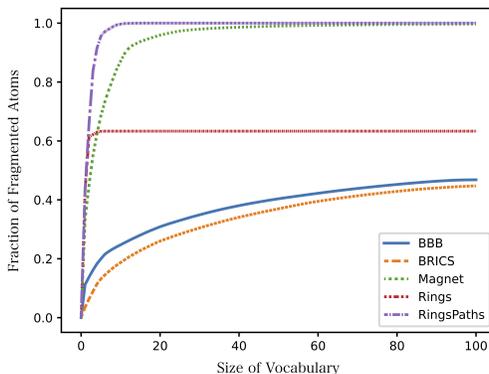
**GPS** [43], **GT** [13], **SAN** [34] and **GRIT** [36]. The detailed results can be found in Table 7 and Table 8. The hyperparameters of our model for ZINC (10k and full) and peptides (struct and func) can be found in 9. Note that we adhere to the 500K parameter budget. Each experiment is repeated over three different seeds except for the ZINC-full experiment, where we only have a single run because of computational and time limitations.

## G.2 Expressiveness

We use the MagNet [24] fragmentation to fragment all graphs in the ZINC-subset dataset. We sort the fragments by number of occurrences in the training set. For each of the 28 most common substructure we train our model to predict the counts of these substructures. As model parameters we

**Table 5:** Performance of FragNet with different fragmentation schemes.

Fragmentation Scheme	ZINC	Peptides	
	10k (MAE ↓)	Struct (MAE ↓)	Func (AP ↑)
BBB	0.127	0.252 ± 0.002	0.637 ± 0.003
BRICS	0.127	0.247 ± 0.008	0.658 ± 0.011
MagNet	0.098	-	-
Rings	0.078	0.249 ± 0.001	0.659 ± 0.007
RingsPaths (ours)	<b>0.077</b>	<b>0.246</b> ± 0.002	<b>0.668</b> ± 0.005



**Figure 9:** Fraction of atoms in ZINC-10k dataset that are part of a fragment as a function of vocabulary size. A fraction of 1 indicates that all molecules in the dataset can be completely fragmented. We compare the chemically inspired fragmentation schemes BBB, BRICS, and MagNet with a fragmentation based just on rings and our RingsPaths fragmentation. The substructures in the vocabulary are sorted by the frequency in which they appear in the molecules.

use three layers of message passing with a hidden dimension of 120. For the final readout function we use a sum aggregation and a two layer MLP. We train our model using the MAE loss for 200 epochs with a learning rate of 0.001 and a reduce-on-plateau learning rate scheduling. We report the accuracy (percentage of graphs where rounded prediction equals the ground-truth count) on the test set. Table 10 shows the complete table of all substructures.

### G.3 Long-range Interaction: Recovery rate

In our synthetic long-range experiment, we consider a graph consisting of two rings connected by a path  $??$ . One node in the graph is the designated source node (marked by a star). The feature of the source node is initialized with one-hot-encoding of one of 10 different classes. All other node features are initialized with a constant encoding. For every node  $t$  in the graph, we train a separate model to predict the class of the source node  $s$ , i.e. the target node  $t$  has to reconstruct a message from the source node. The number of layers of the models is  $\max(d(s, t), 3)$ , ensuring that the target can receive messages from the source. We train the model with the cross entropy loss between the prediction at the target node and the true class of the source. We compare the results of models that have no fragmentation, a ring fragmentation and a ring-path fragmentation. We use our model without batchnorm and a hidden dimension of 64. We train the model for a maximum of 200 epochs with a starting learn rate of 0.001 and average the results over at least five seeds.

### G.4 Generalization: Leave 7-rings out

Table 11 shows the training and test errors of the experiment where we remove all molecules containing 7-rings from the training set.

**Table 6:** Vocabulary sizes for RingsPaths on different datasets.

	ZINC-10k	ZINC-Full	Peptides
<b>Vocabulary Size</b>	18	28	100

**Table 7:** Predictive performance for multiple models on Peptides-struct and func. Best Transformer and best GNN are highlighted.

Type	Model	Peptides-	
		Struct (MAE ↓)	Func (AP ↑)
Transformer	GPS	0.2500 ± 0.0012	0.6535 ± 0.0041
	SAN+LapPE	0.2683 ± 0.0043	0.6384 ± 0.0121
	SAN+RWSE	0.2545 ± 0.0012	0.6439 ± 0.0075
	GRIT	<b>0.2460</b> ± 0.0012	<b>0.6988</b> ± 0.0082
Basic GNNs	GCN	0.3496 ± 0.0013	0.5930 ± 0.0023
	GIN	0.3547 ± 0.0045	0.5498 ± 0.0079
	GatedGCN	0.3420 ± 0.0013	0.5864 ± 0.0035
	GatedGCN+RWSE	0.3357 ± 0.0006	0.6069 ± 0.0035
Topological	CIN++	0.2523 ± 0.0013	0.6569 ± 0.0117
Fragment-Biased	HIMP	0.2503 ± 0.0008	0.5668 ± 0.0149
	FragNet (ours)	<b>0.2462</b> ± 0.0021	<b>0.6678</b> ± 0.005

### G.5 Generalization: Rarity

For the experiment in Table 12, we report the MAE of the ZINC-full validation set grouped by the frequency of the rarest fragment in the molecule. The frequency of a fragment is defined as the fraction of molecules that contain the fragment. As the fragmentation scheme, we use the simple Rings fragmentation.

### G.6 Generalization: QM9

To perform our generalization experiment on QM9, we transform the edge and node features of the molecular graphs in QM9 so that they have the same node features and edge features as the graphs in the ZINC dataset. Additionally, we do not use any molecular graphs that contain atom types that do not appear in the ZINC dataset. We calculate penalized logP as ground truth. Then, we trained our model and GRIT on ZINC-full and tested them on the transformed QM9 dataset.

## H Downstream tasks using substructures

Many other molecular tasks beyond property prediction can benefit from substructure information, highlighting the broader potential applications of our model.

**Motifs for Drug Discovery** Motifs and specific substructures are important inductive biases in molecular generation, optimization, and scaffolding tasks [25; 45; 12]. Employing a set of fragments can simplify the generation process and increase the chemical validity of the generated molecules. Given a fragmentation procedure, the fragments are aggregated into a vocabulary of motifs through complete enumeration of the dataset [29; 30; 20], top-k selection [33; 38] or consolidation into Murcko scaffolds [24]. Encoders for molecule generation often integrate motifs as node features or via additional higher-level encoder networks, as the decoder is explicitly tasked with reconstructing the set of motifs from a given embedding.

**Table 8:** Predictive performance for multiple models on ZINC 10k and ZINC full. Best Transformer and best GNN are highlighted.

Type	Model	ZINC	
		10k (MAE ↓)	Full (MAE ↓)
Transformer	Graphormer	0.122 ± 0.006	0.052 ± 0.005
	GPS	0.070 ± 0.006	-
	GT	0.226 ± 0.014	-
	SAN	0.139 ± 0.006	-
	Graphormer-URPE	0.086 ± 0.007	0.028 ± 0.002
	Graphormer-GD	0.081 ± 0.009	0.025 ± 0.004
	GRIT	<b>0.059</b> ± 0.002	<b>0.023</b> ± 0.001
Basic GNNs	GCN	0.367 ± 0.011	0.113 ± 0.002
	GIN	0.526 ± 0.051	0.088 ± 0.002
	GAT	0.384 ± 0.007	0.111 ± 0.002
	GraphSAGE	0.398 ± 0.002	0.126 ± 0.003
Higher-order	RingGNN	0.353 ± 0.019	-
	3WLGNN	0.303 ± 0.068	-
Topological	CIN-Small	0.094 ± 0.004	0.044 ± 0.003
	CIN++	<b>0.077</b> ± 0.004	0.027 ± 0.007
Fragment-Biased	HIMP	0.151 ± 0.006	0.036 ± 0.002
	GSN	0.115 ± 0.012	-
	Autobahn	0.106 ± 0.004	0.029 ± 0.001
	FragNet (ours)	<b>0.0775</b> ± 0.005	<b>0.0237</b> ± 0.00

**Pretraining** In the context of using GNNs for drug discovery, incorporating motifs as part of a pretraining phase has been shown to improve representation learning capabilities. Zang et al. [51] integrates higher-level structures as nodes in a graph and leverages the graph’s hierarchy for self-supervised pretraining. Similarly, Zhang et al. [53] propose a GNN that operates on a two-tiered graph and predicts the sequence of motifs during network pretraining. To improve the encoding of higher-level structures, Inae et al. [28] suggest a motif-aware pretraining technique, which masks entire motifs during the pretraining phase.

## I Future Work

Future work could address several limitations of our current method. First, our fragmentation approach and ordinal encoding, while effective for molecular graphs, are not well-suited to large, densely connected graphs such as citation or social networks, where the generation of numerous meaningless fragments can introduce significant noise. Second, while our generalization experiments demonstrate superior performance compared to GRIT, the results in Table 12 indicate that GRIT outperforms our method on molecules with frequent fragments. Future efforts could focus on refining fragment-biased models to enhance their performance in these scenarios.

	peptides-struct	peptides-func	ZINC-10k	ZINC-full
num_layers	3	2	5	3
hidden_channels	110	128	64	120
num_layers_out	3	3	3	2
frag-reduction	sum	sum	max	max
out-reduction	mean	mean	mean	mean
dropout	0.05	0.15	0	0
lr	0.001	0.001	0.001	0.001
weight decay	0	0	0.001	0
ema decay	0.99	0.99	0.99	0.99
scheduler	Cosine	ReduceonPlateau	Cosine	ReduceonPlateau
patience	-	30	-	15
factor	-	0.5	-	0.9
batchsize	32	128	32	128
max epochs	300	400	2000	1000
num parameters	440K	440K	221K	494K

**Table 9:** Hyperparameter configuration of our model for the ZINC and peptides benchmarks**Table 10:** Fragment counts for the 42 most common MagNet fragments in ZINC and accuracy scores of our model in predicting the counts.

Fragment														
Count	12862	7548	6198	5629	3904	2204	1799	1772	1348	1330	1071	741	573	375
Accuracy	1.0	0.997	0.999	0.986	0.99	0.969	0.999	1.0	0.963	0.997	0.997	0.933	0.999	0.954
Fragment														
Count	208	204	176	156	113	113	90	80	77	66	54	45	37	32
Accuracy	0.999	0.988	0.983	0.982	0.996	0.995	0.993	0.991	0.998	0.995	0.996	0.996	0.996	0.998
Fragment														
Count	32	31	28	25	23	19	19	18	18	17	15	15	15	13
Accuracy	0.998	0.998	0.996	0.997	1.0	0.998	0.997	0.998	1.0	1.0	0.998	0.998	1.0	0.999

**Table 11:** We remove all molecules that contain 7-rings from the training set and test on all molecules, i.e., also the ones with 7-rings.

Model	ZINC 10k	
	training (MAE ↓)	test (MAE ↓)
GRIT	0.02	0.61
FragNet (ours)	0.08	<b>0.34</b>

**Table 12:** Comparison of the MAE of GRIT and our model on ZINC-full. Graphs are grouped by the frequency of their rarest fragment.

Frequency of rarest Fragment	< 0.1%	< 1%	< 10%	≥ 10%
HIMP (MAE)	14.4	0.48	0.15	0.030
GRIT (MAE)	9.5	0.26	0.026	0.018
FragNet (ours) (MAE)	5.3	0.15	0.045	0.021