
Implicit Geometry and Interaction Embeddings Improve Few-Shot Molecular Property Prediction

Christopher Fifty^{★,1}, Joseph M. Paggi^{★,1}, Ehsan Amid²,
Jure Leskovec¹, Ron O. Dror¹

Equal Contribution[★]

Stanford University¹, Google Brain²,

fifty@cs.stanford.com, jpaggi@stanford.edu, rondror@cs.stanford.edu

Abstract

Few-shot learning is a promising approach to molecular property prediction as supervised data is often very limited. However, many important molecular properties depend on complex molecular characteristics — such as the various 3D geometries a molecule may adopt or the types of chemical interactions it can form — that are not explicitly encoded in the feature space and must be approximated from limited data. Learning these characteristics can be difficult, especially for few-shot learning algorithms that are designed for fast adaptation to new tasks. In this work, we develop molecular embeddings that encode complex molecular characteristics to improve the performance of few-shot molecular property prediction. Our approach leverages large amounts of synthetic data, namely the results of molecular docking calculations, and a multi-task learning paradigm to structure the embedding space. The embeddings improve few-shot learning performance on multiple molecular property prediction benchmarks using Multi-Task, MAML, and Prototypical Networks. Our code is available at <https://github.com/cfifty/IGNITE>.

1 Introduction

In stark contrast to common applications in computer vision and natural language processing, little supervised data is available for many important molecular property prediction tasks [1, 24]. Determining the biological effects of a molecule—for example, a drug candidate—often involves a lengthy period of cell cultivation followed by a multi-step assay requiring specialized equipment [16]. Animal studies may even be required to measure certain properties, such as toxicity [18]. The throughput of these assays and studies is limited, and synthesizing molecules to be tested is often expensive [5]. As a result, supervised datasets measuring molecular properties often contain on the order of 10s or 100s of data points, and few-shot learning algorithms are necessary to cope with the scarcity of labeled data [21].

Few-shot learning algorithms are designed for low-data paradigms, but the effectiveness of these approaches for predicting molecular properties, such as the toxicity of a molecule or its binding affinity with a specific protein, is constrained by the inherent complexity of molecular interactions [7, 9]. Molecules are specified as a connected graph of atoms, each with its own feature vector, but this feature space does not directly encode complex molecular characteristics, such as the various 3D geometries a molecule may adopt or the types of chemical interactions it can form in each geometry [27]. These characteristics are essential to the real-world properties of molecules and deterministic of how molecules interact with other molecules, including proteins and nucleic acids [4].

While molecular representations encoding these attributes may be learned from the feature space, doing so injects an additional layer of complexity into the learning process. As the learning dynamics of few-shot learning algorithms are already quite complex — and sometimes even unstable [2] —

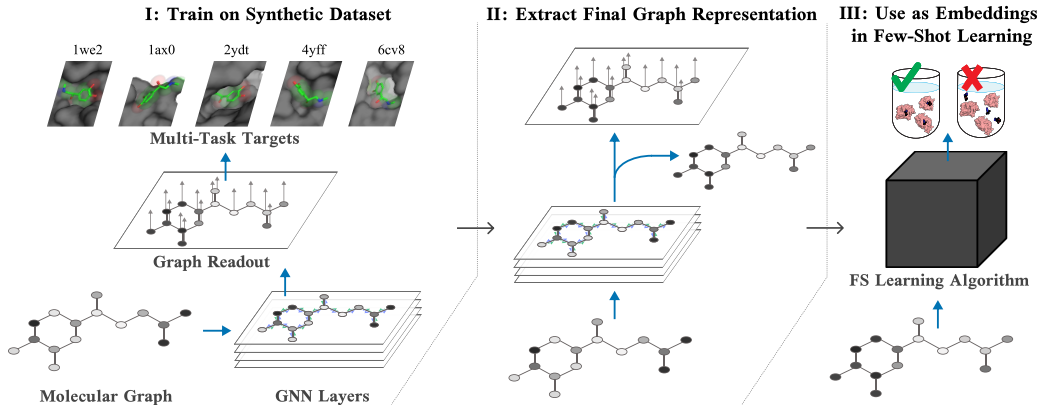


Figure 1: The end-to-end paradigm of developing of **IGNITE** embeddings: (I) train a multi-task GNN to map a small molecule to binding energies (kcal/mol) determined by the interactions between this small molecule and protein targets by physics-based docking. (II) extract the pre-readout graph representation. (III) for a few-shot learning algorithm, replace the default feature space with the learned molecular embeddings. The multi-task targets 1we2, 1ax0, 2ydt, 4yff, and 6cv8, represent Protein Data Bank codes for 5 of the 2,034 targets in our actual dataset.

training a model to develop both the capacity to rapidly adapt to new tasks as well as to approximate complex molecule–target interactions may be difficult, or even impossible, for existing algorithms.

2 IGNITE: Implicit Geometry aNd InTeraction Embedding

We aim to develop molecular embeddings that encode complex molecular characteristics to improve few-shot molecular property prediction. Our approach begins with curating a large-scale synthetic dataset that measures how 1,000s of molecules bind to each one of 1,000s of different proteins, therefore implicitly encoding the various 3D geometries a molecule may adopt as well as the chemical interactions it may form within the modeling objective of this dataset. We then train a MPNN [12] with a multi-task learning scheme—using a single model to predict the interactions among all molecules and targets—to learn a single molecular representation that generalizes to all tasks. Finally, we extract the final-layer learned representation of the multi-task MPNN as a molecular embedding, and use it to initialize the feature space of other few-shot learning algorithms on down-stream molecular property prediction tasks. This approach is visualized in Figure 1, and because of its ability to capture the implicit characteristics of a molecule, we call this method **Implicit Geometry aNd InTeraction Embedding (IGNITE)**.

Synthetic Data Generation. We turn to computational “docking” methods to generate synthetic data that depends on complex molecular characteristics. In particular, we use Glide, a physics-based docking program [11] to curate a large-scale synthetic dataset measuring how 1,000s of molecules bind to each of 1,000s of different proteins. Given the 3D structure of a target protein and the chemical graph of a small molecule, docking methods estimate the binding energy between the two.

Docking matches the process by which real binding energies arise in two key ways. First, docking considers many potential geometries of the protein–small molecule complex, including sampling over the potential internal 3D geometries of the small molecule, and the final prediction is related to the most favorable geometry found. Second, the favorability of each potential geometry is assessed using a physics-inspired scoring function that quantifies the complementarity of molecular geometries and presence of key chemical interactions.

We reason that pre-training on a sufficiently large amount of synthetic data generated by molecular docking would encode molecular attributes deterministic of their interactions into a model’s learned representations. To learn these representations, we turn to multi-task learning.

Multi-Task Learning Paradigm. The synthetic dataset naturally induces a multi-task learning paradigm where each protein represents a separate task. Specifically, for each protein target, we map a small molecule input to its binding energy (kcal/mol) as determined by molecular docking. This process is then repeated across the 1,000s of protein targets in our dataset to compose a multi-task learning system that uses a single deep learning model to predict how a small molecule will bind to each of the proteins in our dataset.

With regard to the model, we use the same model architecture and training procedure outlined in [21] (see Appendix C). The objective function is to minimize the mean squared error between the model’s predictions and the free energy measurements. We randomly split our synthetic dataset into 80% train and 20% validation, and use early stopping with a window size of 10 epochs on the validation dataset. Training spans 100 epochs. We save the best model checkpoint from early stopping to disk.

Usage in Few-Shot Learning. Many few-shot learning algorithms for molecular property prediction operate on a connected graph of atom-level feature vectors. Atom-level feature vectors often include a one-hot encoding of the element type, atomic charge, mass, valency, etc. With **IGNITE**, we simply change the feature space on which few-shot learning algorithms operate from the default featurization—using a one-hot encoding of element type, atomic charge, etc.—to the pre-readout atom representation learned on the synthetic dataset by predicting binding energies.

3 Experiments and Discussion

Our aim is to determine if **IGNITE** improves few-shot learning performance on molecular property prediction. We build our analysis on two conditions reflecting the constraints often imposed on academia and industry where existing models are likely to run efficiently but a capacity — both human and computational — to implement new methods and tune them is limited:

- **Model Agnostic:** **IGNITE** should be applicable and effective across a diverse range of architectures and learning algorithms.
- **Minimal Modification:** replacing the initial feature space with the **IGNITE** embedding space should be effective with minimal modification or hyperparameter tuning.

For our experimental evaluation, we choose FS-Mol [21], a few-shot learning benchmark designed to reflect the size and scope of datasets used in drug discovery, as well as 3 datasets from the MoleculeNet [25] benchmark suite (see Appendix D for two of the MoleculeNet benchmarks).

FS-Mol Description. FS-Mol is a few-shot molecular property prediction benchmark that measures how small molecules interact with protein targets. The dataset spans 5,120 proteins and encompasses 233,786 different small molecules. On average, each protein has 94 small molecule measurements, and the benchmark is formulated to predict if one of these small molecules will inhibit (or not) the protein target. The benchmark uses a unique metric, Δ AUPRC, which reports the relative change in the area under the precision-recall curve (AUPRC) from a random guess baseline.

FS-Mol Experimental Setup. Stanley et al. [21] presents several few-shot learning baselines for this benchmark. We select three of the best: multi-task pre-training with single-task fine-tuning (MT) [6], model-agnostic meta-learning (MAML) [10], and prototypical networks (PROTO) [20]. Our aim is to determine if their performance can be improved by simply modifying the molecular features on which they operate.

This selection of baselines fulfills our “Model Agnostic” criterion as the learning paradigms for MT, MAML, and PROTO are dissimilar. For instance, MT builds molecular representations amenable to all tasks, MAML learns representations that can be quickly adapted, and PROTO leverages a notion of distance—rather than inner product—in its learned representations to predict the label of a molecule. While each baseline is dissimilar, they all use a linear transformation to convert the input feature space to the model’s hidden dimension representation space. We simply replace this first hidden dimension representation space with the **IGNITE** embedding space. This change actually simplifies the few-shot learning algorithm, decreasing computational and runtime complexity, by removing the initial embedding layer.

Adhering to the “Minimum Modification” criterion, we do not perform any additional hyperparameter tuning for few-shot learning approaches using **IGNITE** embeddings. We simply use the hyperparameters found by Stanley et al. [21] for each baseline model. We prefix few-shot learning methods with an **IGNITE**- to indicate this model uses **IGNITE** embeddings.

FS-Mol Results. Table 1 indicates that **IGNITE** embeddings significantly improve MT performance on target inhibition prediction, and performance is most improved at smaller support size splits. This result is likely caused from **IGNITE** encoding complex molecular characteristics relevant for protein targets in the test set, characteristics that MT struggles to incorporate into its learned representations from only 16 training examples. While the increase in performance declines as the support size

Table 1: FS-Mol classification performance across the 157 targets in the FS-Mol test split. We report mean Δ AUPRC and standard error from 10 training runs with different random seeds. Performance of MT, and MAML cited from [21]. PROTO is re-evaluated to use only graph-level features.

Method	Support Size				
	16	32	64	128	256
MT	0.112 \pm 0.005	0.144 \pm 0.006	0.177 \pm 0.008	0.223 \pm 0.008	0.237 \pm 0.018
IGNITE -MT	0.157 \pm 0.007	0.191 \pm 0.008	0.221 \pm 0.008	0.259 \pm 0.009	0.281 \pm 0.020
% change	Δ 40.2%	Δ 32.6%	Δ 24.9%	Δ 16.1%	Δ 18.6%
MAML	0.160 \pm 0.008	0.167 \pm 0.008	0.173 \pm 0.008	0.192 \pm 0.009	0.186 \pm 0.019
IGNITE -MAML	0.164 \pm 0.008	0.173 \pm 0.008	0.183 \pm 0.008	0.204 \pm 0.008	0.193 \pm 0.018
% change	Δ 2.5%	Δ 3.6%	Δ 5.8%	Δ 6.3%	Δ 3.8%
PROTO	0.185 \pm 0.008	0.224 \pm 0.009	0.256 \pm 0.009	0.290 \pm 0.009	0.263 \pm 0.018
IGNITE -PROTO	0.195 \pm 0.008	0.232 \pm 0.009	0.263 \pm 0.009	0.297 \pm 0.009	0.283 \pm 0.019
% change	Δ 5.4%	Δ 3.6%	Δ 2.7%	Δ 2.4%	Δ 7.6%

Table 2: Results on the BACE Classification. We report mean AUPRC as well as standard error averaged across 10 training runs with different random seeds.

Method	Support Size				
	16	32	64	128	256
MT	0.547 \pm 0.064	0.600 \pm 0.063	0.652 \pm 0.051	0.696 \pm 0.037	0.738 \pm 0.028
IGNITE -MT	0.574 \pm 0.08	0.674 \pm 0.039	0.719 \pm 0.025	0.739 \pm 0.037	0.784 \pm 0.014
% change	Δ 4.9%	Δ 12.3%	Δ 10.3%	Δ 6.2%	Δ 6.2%
MAML	0.510 \pm 0.026	0.508 \pm 0.035	0.548 \pm 0.065	0.568 \pm 0.075	0.649 \pm 0.016
IGNITE -MAML	0.620 \pm 0.027	0.633 \pm 0.031	0.643 \pm 0.043	0.668 \pm 0.037	0.685 \pm 0.041
% change	Δ 21.6%	Δ 24.6%	Δ 17.3%	Δ 17.6%	Δ 5.5%

increases, it ticks up at 256. This effect is similarly noted by Stanley et al. [21] and is caused by few protein targets in the test set actually containing at least 256 examples to evaluate the 256-size support split.

We also analyze if **IGNITE**—an approach developed with a multi-task learning paradigm—might also improve the performance of other few-shot learning algorithms like Prototypical Networks or MAML. This is a difficult proposition as the inductive biases fundamental to each learning algorithm are highly dissimilar. To our surprise, our findings in Table 1 show consistent improvement of PROTO and MAML across all support splits. It may be possible to achieve even larger performance improvements on PROTO and MAML by learning **IGNITE** embeddings with an approach that mirrors the down-stream few-shot learning algorithm. We leave this exploration to future work.

BACE Classification Description & Setup. BACE is a dataset in the MoleculeNet [25] benchmark that measures if a small molecule inhibits (or not) the human β -secretase-1 enzyme. It contains a single task and measurements for 1,522 small molecules. While BACE is often formulated as a stand-alone molecular property prediction dataset, we adapt it to the few-shot learning setting by meta-training on the larger FS-Mol benchmark, sampling a support from the BACE dataset, and designating all other examples in this dataset as the query set.

BACE Results. Our results on the BACE dataset reflect our findings on FS-Mol: augmenting MT and MAML with **IGNITE** embeddings substantively improves performance across all support sizes (Table 2). However, unlike our findings on FS-Mol, this time MAML benefits the most, improving AUPRC by almost 25% at 32 support size.

4 Visualized Example

We offer qualitative analysis to support the hypothesis that **IGNITE** provides improved performance by implicitly encoding the potential geometries of small molecules. We consider two small molecules, CHEMBL200234 and levisoprenaline, that have dissimilar 2D chemical structures but both bind to the β_1 -adrenergic receptor (B1AR) in a similar binding pose (Figure 2). Moreover, examining

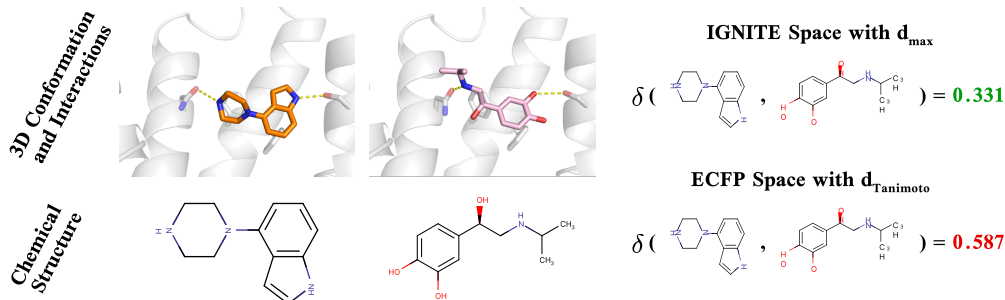


Figure 2: While these two molecules, CHEMBL200234 (left) and levisoprenaline (right), have very different 2D structures, they can adopt similar 3D conformations and form similar interactions with a target protein. Shown are renderings of these molecules in complex with the β_1 -adrenergic receptor, with dashed yellow lines indicating hydrogen bonds. These molecules are relatively close in the embedding space, falling in the 33rd percentile amongst a randomly selected set of molecules, as compared to the 58th percentile using the Tanimoto coefficient in ECFP space. Percentiles are the average rank of each of the active molecules when ordering the molecules by similarity to the other active molecule.

the experimentally determined structure of each one bound to B1AR reveals they adopt similar 3D conformations and form similar chemical interactions (e.g., hydrogen bonds) with the target.

To assess whether the ability of these molecules to adopt similar geometries and thereby bind to the same target protein is captured by the **IGNITE**, we evaluate the distance between these two molecules in the **IGNITE** embedding space, compared to the distance between these two molecules as determined by the Tanimoto distance between their extended connectivity fingerprints. Specifically, we randomly sample 100,000 small molecules to use as decoys and then order all small molecules by their distance from each active compound and then record the rank of the other active compound.

Under Tanimoto distance, the active molecules are ranked at an average percentile of 0.587 whereas, under the **IGNITE** embeddings distance, the active molecules are ranked at an average percentile of 0.331 (lower is closer in embedding space). This example illustrates that molecules even with dissimilar chemical structures may be placed relatively close to one another in the **IGNITE** embedding space. Taken alongside Appendix A, this finding lends support the hypothesis that molecular embeddings learned by pre-training on binding energies implicitly encode complex molecular characteristics.

5 Conclusion

In this work, we develop molecular embeddings that encode complex molecular characteristics—such as the various 3D geometries a molecule may adopt or the types of chemical interactions it can form—to improve the performance of few-shot learning algorithms on molecular property prediction. Our analysis of the **IGNITE** embedding space suggests these characteristics are in fact encoded by the embeddings. Further, our empirical analysis on four molecular property prediction benchmarks, with properties ranging from the inhibition of protein targets to the toxicity of small molecules, suggests **IGNITE** almost universally improves the performance of three popular, but very different, few-shot learning algorithms. Moreover, integrating **IGNITE** into few-shot learning algorithms takes minimum modification and does not require additional hyperparameter tuning.

Instead of using synthetic data, we could have used the same strategy we propose here to learn embeddings from experimental measurements of small molecules binding to an array of proteins. However, synthetic data has two key advantages. First, we can generate as much training data as desired. Second, it overcomes the difficulties in comparing experimental measurements across different assay types. Comparing the measurements from one assay type with another can be difficult, for both medicinal chemists and few-shot learning paradigms.

Acknowledgments and Disclosure of Funding

This work was supported by National Institutes of Health grant R01GM127359 (to R.O.D.).

References

- [1] Han Altae-Tran, Bharath Ramsundar, Aneesh S Pappu, and Vijay Pande. Low data drug discovery with one-shot learning. *ACS central science*, 3(4):283–293, 2017.
- [2] Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your maml. *arXiv preprint arXiv:1810.09502*, 2018.
- [3] Thomas Bachlechner, Bodhisattwa Prasad Majumder, Henry Mao, Gary Cottrell, and Julian McAuley. Rezero is all you need: Fast convergence at large depth. In *Uncertainty in Artificial Intelligence*, pages 1352–1361. PMLR, 2021.
- [4] Caterina Bissantz, Bernd Kuhn, and Martin Stahl. A medicinal chemist’s guide to molecular interactions. *Journal of medicinal chemistry*, 53(14):5061–5084, 2010.
- [5] David C Blakemore, Luis Castro, Ian Churcher, David C Rees, Andrew W Thomas, David M Wilson, and Anthony Wood. Organic synthesis provides opportunities to transform drug discovery. *Nature chemistry*, 10(4):383–394, 2018.
- [6] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.
- [7] Kangway V Chuang, Laura M Gunsalus, and Michael J Keiser. Learning molecular representations for medicinal chemistry: miniperspective. *Journal of Medicinal Chemistry*, 63(16):8705–8722, 2020.
- [8] Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal neighbourhood aggregation for graph nets. *Advances in Neural Information Processing Systems*, 33:13260–13271, 2020.
- [9] Evan N Feinberg, Elizabeth Joshi, Vijay S Pande, and Alan C Cheng. Improvement in admet prediction with multitask deep featurization. *Journal of medicinal chemistry*, 63(16):8835–8848, 2020.
- [10] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017.
- [11] Richard A Friesner, Jay L Banks, Robert B Murphy, Thomas A Halgren, Jasna J Klicic, Daniel T Mainz, Matthew P Repasky, Eric H Knoll, Mee Shelley, Jason K Perry, et al. Glide: a new approach for rapid, accurate docking and scoring. 1. method and assessment of docking accuracy. *Journal of medicinal chemistry*, 47(7):1739–1749, 2004.
- [12] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International conference on machine learning*, pages 1263–1272. PMLR, 2017.
- [13] Matthew P. Jacobson, Richard A. Friesner, Zhexin Xiang, and Barry Honig. On the role of the crystal environment in determining protein side-chain conformations. *Journal of Molecular Biology*, 320(3):597–608, 2002. ISSN 0022-2836. doi: 10.1016/S0022-2836(02)00470-9. URL <http://www.sciencedirect.com/science/article/pii/S0022283602004709>.
- [14] Maurice G Kendall. *A Course in the Geometry of n Dimensions*. Courier Corporation, 2004.
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [16] Olga Maria Lage, María C Ramos, Rita Calisto, Eduarda Almeida, Vitor Vasconcelos, and Francisca Vicente. Current screening methodologies in drug discovery for selected human diseases. *Marine drugs*, 16(8):279, 2018.
- [17] Stephen Merity. Single headed attention rnn: Stop thinking with your head. *arXiv preprint arXiv:1911.11423*, 2019.

- [18] S Parasuraman. Toxicological screening. *Journal of pharmacology & pharmacotherapeutics*, 2(2):74, 2011.
- [19] David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *Journal of chemical information and modeling*, 50(5):742–754, 2010.
- [20] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017.
- [21] Megan Stanley, John F Bronskill, Krzysztof Maziarsz, Hubert Misztela, Jessica Lanini, Marwin Segler, Nadine Schneider, and Marc Brockschmidt. Fs-mol: A few-shot learning dataset of molecules. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- [22] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 32–42, 2021.
- [23] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [24] Michael J Waring, John Arrowsmith, Andrew R Leach, Paul D Leeson, Sam Mandrell, Robert M Owen, Garry Pairaudeau, William D Pennie, Stephen D Pickett, Jibo Wang, et al. An analysis of the attrition of drug candidates from four major pharmaceutical companies. *Nature reviews Drug discovery*, 14(7):475–486, 2015.
- [25] Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine learning. *Chemical science*, 9(2):513–530, 2018.
- [26] Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tieyan Liu. On layer normalization in the transformer architecture. In *International Conference on Machine Learning*, pages 10524–10533. PMLR, 2020.
- [27] Kevin Yang, Kyle Swanson, Wengong Jin, Connor Coley, Philipp Eiden, Hua Gao, Angel Guzman-Perez, Timothy Hopper, Brian Kelley, Miriam Mathea, et al. Analyzing learned molecular representations for property prediction. *Journal of chemical information and modeling*, 59(8):3370–3388, 2019.

Appendix

A Molecular Space Analysis

While our MPNN model may fit the docking data, it is unclear if relevant biophysical knowledge is actually encoded into the learned **IGNITE** embeddings. Answering this question is the focus of our analysis.

Our analysis commences with an exploration of various molecular spaces, and how the molecular space structure implicitly defines a notion of similarity among molecules within that space. We present 3 ways to measure distance between molecules: docking-based, embedding-based, and fingerprint-based. Later, we use the notion of distance among molecules—distance being a proxy for similarity—to quantify if the molecular characteristics implicitly captured by molecular docking are encoded into the **IGNITE** embedding space. We employ the Kendall Tau measure [14] to objectively quantify this effect.

Docking-Based Distance. Physics-based docking software implicitly encodes a notion of similarity among small molecules. For instance, one may define “similar” small molecules as those that manifest similar binding energies across many different protein targets; such molecules typically adopt similar 3D conformations and form similar chemical interactions with the same target. We quantify the implicitly defined similarity between a pair of molecules m_1, m_2 from docking as the relative difference in docking energies across all targets:

$$d(m_1, m_2) = \sum_{t \in \text{targets}} \frac{|t(m_1) - t(m_2)|}{\max_m \{t(m)\} - \min_m \{t(m)\}}. \quad (1)$$

We compute the relative difference in binding energy as opposed to the absolute difference as different targets have different ranges of binding energies, and we wish to avoid lending more weight to targets with larger ranges than to targets with smaller ranges.

Embedding-Based Distance. Next we define the distance between molecules in the **IGNITE** embedding space. Our MPNN embedding model uses a combined graph readout composed of an element-wise max, learned weighted sums, and learned weighted means. Applying any of these pooling operations to the per-atom **IGNITE** embeddings will output a component of the vector used to predict the binding energies associated with this target. Accordingly, we choose the simplest operation — element-wise max pooling — to motivate the distance between two arbitrary molecules in our embedding space.

Let m_1, m_2 be two molecules parameterized by a per-atom, learned feature matrix and an adjacency matrix in the embedding space: $m_1^{(e)} = (\mathcal{V}_1^{(e)}, \mathcal{E}_1^{(e)})$. We define ρ to be the max-pooling operation across the embedding dimension for all nodes in the graph.

$$\text{For atom embeddings } \mathcal{V}^{(e)} = \begin{bmatrix} v_{1,1} & \dots & v_{1,d} \\ \dots & & \\ v_{n,1} & \dots & v_{n,d} \end{bmatrix} \in \mathbb{R}^{n \times d},$$
$$\rho(\mathcal{V}^{(e)}) \in \mathbb{R}^d, \text{ where } [\rho(\mathcal{V}^{(e)})]_j = \max_i v_{ij}.$$

We can now define the distance between two molecules in the embedding space as the l_2 norm of their element-wise maximum difference:

$$d_\rho(m_1, m_2) = \|\rho(\mathcal{V}_1) - \rho(\mathcal{V}_2)\|_2 \quad (2)$$

One may also apply Equation 2 to the initial feature space used as input to the model:

$$m_1^{(f)} = (\mathcal{V}_1^{(f)}, \mathcal{E}_1^{(f)})$$

A model’s feature space is simply the initial numerical representation of a molecule. Machine learning models often define an atom-level feature space as a numerical vector encoding an atom’s element type atomic charge, atomic mass, valency and/or other atomic-level quantities.

Fingerprint-Based Distance. A third metric of the distance between two molecules is the Tanimoto distance d_τ between their extended-connectivity fingerprints (ECFPs). An ECFP is a binary string

indicating whether or not each of many 2D chemical substructures is present in a given molecule. This Tanimoto distance—which is one of the most commonly used notions of molecular similarity in both academia and industry [19]—is defined as the number of chemical sub-structures shared between a pair of molecules divided by the total number of chemical sub-structures present in either molecule. A high value indicates that a pair of molecules share the same components but does not necessarily indicate that those components are connected to one another in the same way or positioned similarly in 3D space.

Kendall Tau Measure. While one may define various metrics to determine the distance among molecules within a molecular space, comparing distances across molecular spaces is more involved. One may abstract the notion of distance between two pairs of molecules to an ordering of all molecules ranked by distance to a single anchor molecule. By ranking the similarity of all molecules to an anchor molecule in one molecular space, we may quantify the extent to which other molecular spaces produce a (dis)similar ordering with respect to this same anchor molecule.

We employ the Kendall Tau rank distance [14], also known as the bubble-sort distance, to quantify how similar one ranking of molecular similarity is to the docking-based ordering induced by Equation 1. The Kendall Tau rank distance is given by:

$$K_d(\tau_1, \tau_2) = |(i, j) : i < j, \{[\tau_1(i) \wedge \tau_2(i) > \tau_2(j)] \\ \vee [\tau_1(i) > \tau_1(j) \wedge \tau_2(i) < \tau_2(j)]\}|$$

where $\tau_1(i)$ and $\tau_2(i)$ are the rankings of molecule i given by the rankings τ_1 and τ_2 , respectively. Intuitively, Kendall Tau counts the number of swaps made during a bubble sort to convert the ordering in the **IGNITE** embedding space, feature space, or ECFP space to the ordering in the docking space. The smaller the Kendall Tau distance between two molecular spaces with respect to the same anchor molecule, the more “similar” they are in the sense of molecular space structure with relation to that anchor molecule.

Findings. A naïve Kendall Tau analysis would compare the product space of space of (32,547 \times 32,547), with each unique molecule in our synthetic dataset serving as an anchor molecule for all other molecules. As this comparison exceeds our computational resources, we randomly sample 100 anchor molecules instead to compare the relative orderings induced by various molecule spaces with the ordering induced by the docking-based distance in Equation 1. This analysis approximates how “similar” the various molecular spaces are to the docking-based space, the space that directly encodes the various 3D geometries a molecule may adopt as well as the chemical interactions it can form.

Our findings are summarized in Table 3 and indicate ranking molecules by distance in the **IGNITE** embedding space is most similar to the rankings induced by distance in the docking space. Orderings among molecules in ECFP space requires, on average, 4.5 million more swaps than orderings in the **IGNITE** embedding space. Similarly, the feature space requires 77 million more swaps to reorder its structure to match the docking-space than does the **IGNITE** embedding space. The “Random” entry in Table 3 represents a randomly structured molecular space and serves as a reference value for comparison.

Table 3: Kendall Tau Analysis. Scores are averaged across 100 randomly selected small molecules. The number of swaps indicates the average number of swaps made by KT to transform an ordering to the docking-based ordering. Normalized Kendall Tau is the number of swaps divided by the number of molecules in an ordering.

Molecular Space	Number of Swaps (↓)	Normalized Kendall Tau (↓)
Docking	0	0
IGNITE	2.96×10^7	0.386
Feature	3.73×10^7	0.488
ECFP	3.28×10^7	0.429
Random	3.83×10^7	0.50

B GLIDE Protocol

To produce synthetic data, we used the physics-based docking method GLIDE [11] (Schrödinger Release 2022-2: Maestro, Schrödinger, LLC, New York, NY, 2022). Running GLIDE can be

separated into three steps: preparation of small molecules, preparation of protein structures, and the docking of the prepared small molecules into the prepared protein structures.

To prepare a set of small molecules to be docked, we began with a selection of 32,547 small molecules from ChEMBL. Importantly, we used ChEMBL simply to provide examples of small molecules; we took the SMILES strings and did not use the annotated activities in any way. We then used the Schrodinger ligprep tool to enumerate relevant tautomeric states (e.g. protonation states) and, if not specified in the SMILES strings, up to 32 stereoisomers.

To prepare a set of protein structures, we began with a set of 1,601 structures from the PDBind database. Schrodinger prepwizard was used to add hydrogens, choose tautomeric states, and perform a constrained minimization (heavy atoms within 0.3 Å of starting coordinates) using default parameters [13]. The small molecule annotated in PDBind was used to center the docking site, but was otherwise not used in any way.

Docking was run using GLIDE in SP mode using default parameters. The docking score for a given small molecule–protein pair was defined as the most favorable score encountered for any tautomeric state of the small molecule, taking into account penalties for choosing an unfavorable tautomeric state. Due to this, it might be that **IGNITE** encodes an implicit representation of a molecule’s potential tautomeric states in addition to potential 3D geometries and chemical interactions. Docking calculations were attempted for all small molecule–protein pairs but in some cases docking failed to produce any reasonable poses or the job failed for other reasons so these pairs weren’t considered.

C Experimental Design

We offer additional details related to the training of **IGNITE**. While this description can be valuable to some readers, we direct those interested in reproducing our findings or building on our framework to access the code directly. Our code is released in the supplementary material download on OpenReview; however, due to the synthetic data occupying 77 GB on disk, we are unable to upload (or link to) the dataset without breaking anonymity.

C.1 Model Architecture

We use a 10-layer MPNN [12] augmented with PNA [8] and using a graph readout composed of an element-wise maximum, learned multi-headed weighted sum, and learned multi-headed weighted mean. Each layer in the MPNN uses a Pre-Norm [26] transformer-like [23] residual structure with ReZero weighing [3] using a vector formulation as described in [22]. Moreover, a BOOM-layer [17] is used for a second residual connection after message-passing. The post-readout molecular representation is then passed through a shallow MLP that is shared among all tasks before a task-specific linear projection is used to predict the binding energy.

Our optimization protocol uses a linear warmup schedule with 100 steps and Adam [15] optimizer.

C.2 Additional Training Details

Our architecture employs a learning rate of 5e-5 for the shared parameters and a learning rate of 1e-4 for the task-specific parameters. We use a linear warm-up scheduler for 100 steps (starting at 0 and ending at the specified learning rate) for both shared and task-specific learning rates. The training process also leverages Adam [15] with default Pytorch parameters and uses a batch size of 256 molecules. With regards to the task-specific parameters in the **IGNITE** multi-task

training paradigm, the output from the graph readout is passed through a shared MLP of hidden dimension 512 and then uses a task-specific projection layer of dimension $[512 \times 1]$ for each target. Table 4 offers a comparison of the total number of small molecules as well as the number of tasks in each benchmark that we use for evaluation. For our experiments, we used a single 80 GB Nvidia A100 GPU, with the model taking approximately 10 GB of memory on the GPU.

Table 4: Benchmark comparison.

Dataset	# tasks	# compounds
FS-Mol Test	157	27520
BACE	1	1522
Tox21	4	7831
HIV	1	41127

Table 5: Tox21 classification performance across 4 targets that measure the toxicity of a small molecule with a certain biological pathway or nuclear receptor. We report mean AUPRC as well as standard error averaged across 10 training runs with different random seeds.

Method	Support Size				
	16	32	64	128	256
MT	0.120 \pm 0.006	0.140 \pm 0.013	0.155 \pm 0.015	0.166 \pm 0.019	0.198 \pm 0.025
IGNITE -MT	0.150 \pm 0.017	0.148 \pm 0.017	0.173 \pm 0.019	0.194 \pm 0.031	0.242 \pm 0.034
% change	Δ25.1%	Δ5.9%	Δ11.7%	Δ16.6%	Δ22.8%
MAML	0.131 \pm 0.009	0.130 \pm 0.009	0.131 \pm 0.009	0.131 \pm 0.010	0.134 \pm 0.008
IGNITE -MAML	0.124 \pm 0.009	0.134 \pm 0.010	0.136 \pm 0.011	0.142 \pm 0.012	0.163 \pm 0.020
% change	Δ-5.2%	Δ3.1%	Δ3.8%	Δ8.2%	Δ21.7%

D Additional Experimental Results

To offer additional insight into the FS-Mol empirical findings in section 3, we present a per-target breakdown of the first 50 targets in the test set of FS-Mol across support size splits of 16, 32, 64, 128, and 256. Table 7 depicts these results for the Multi-Task Learning baseline, Table 8 depicts these results for the Prototypical Networks baseline, and Table 9 depicts these results from the MAML baseline. Highlighted values indicate the highest result for this target at the given support size.

MoleculeNet Description & Setup. MoleculeNet [25] is a benchmark suite spanning many disparate molecular properties. While each dataset can stand alone, we adapt it to the few-shot learning setting by meta-training on the larger FS-Mol benchmark, sampling a support set from the MoleculeNet dataset, and designating all other examples in this dataset as the query set. As FS-Mol is a binary prediction benchmark, we limit ourselves to binary prediction tasks from the MoleculeNet benchmark and select BACE, HIV, and TOX21 for our experimental analysis.

BACE measures if a small molecule inhibits (or not) the human β -secretase-1 enzyme. It contains a single task and measurements for 1,522 small molecules. Tox21 measures the toxicity of small molecules for several different biological targets such as nuclear receptors and stress pathways. We filter out the targets with high class imbalance, and the resulting dataset contains 7,831 measures of toxicity across 4 targets. HIV measures if a small molecule will inhibit HIV replication. It is composed of a single task and spans 41,127 small molecules. This dataset is extremely unbalanced, and as both MT and MAML require at least two points from each class to be present—one for training and the other for early stopping when fine-tuning on the support set—we are unable to evaluate HIV at smaller support sizes.

As models are pre-trained on FS-Mol, we use the same fine-tuning hyperparameters and approach as in our earlier evaluation; no additional hyperparameter tuning is performed for **IGNITE**. We report AUPRC, or the area under the precision-recall curve, to align our evaluation with the MoleculeNet benchmark [25]. Moreover, due to implementation challenges, we drop our evaluation of PROTO to focus on MT and MAML.

MoleculeNet Results. Our results on MoleculeNet datasets reflect our findings on FS-Mol: **IGNITE** appears to almost universally improve few-shot learning performance. On BACE, we find augmenting MT and MAML with **IGNITE** embeddings substantively improves performance across all support sizes (Table 2). However, unlike our findings on FS-Mol, this time MAML benefits the most, improving AUPRC by almost 25% at 32 support size. On Tox 21, the performance of MT and MAML is improved across all support splits except at the 16 support size for MAML (Table 5). On HIV, **IGNITE** substantively increases AUPRC for both MT and MAML (Table 6). This dataset is especially

Table 6: Results on the HIV dataset from the MoleculeNet benchmark suite. We report mean AUPRC as well as standard error averaged across 10 training runs with different random seeds.

Method	Support Size	
	128	256
MT	0.062 \pm 0.035	0.078 \pm 0.033
IGNITE -MT	0.066 \pm 0.029	0.087 \pm 0.035
% change	Δ6.5%	Δ11.5%
MAML	0.053 \pm 0.015	0.060 \pm 0.023
IGNITE -MAML	0.070 \pm 0.021	0.068 \pm 0.016
% change	Δ32.1%	Δ13.3%

challenging given its extreme label imbalance. This closely mirrors real-life drug discovery where exceedingly few molecules manifest a desired property. It is exciting that **IGNITE** appears to improve the performance of few-shot learning methods on these types of data distributions.

Finally, on the HIV MoleculeNet benchmark, we find **IGNITE** substantively increases AUPRC for both MT and MAML. This dataset is especially challenging given its extreme unbalance. We are unable to evaluate this dataset at smaller support splits due to the constraint that at least two points from each class are in the support set for each training run, but this scenario actually closely mirrors real-life drug discovery where very few molecules have the desired property. As a result, It is particularly exciting that **IGNITE** improves the performance of few-shot learning methods on this paradigm.

Table 7: Multi-Task Results measuring Δ AUPRC on the first 50 tasks in the test set of FS-Mol.

TASK-ID	16 (IGNITE-MT)	32 (IGNITE-MT)	64 (IGNITE-MT)	128 (IGNITE-MT)	256 (IGNITE)	16 (MT)	32 (MT)	64 (MT)	128 (MT)	256 (MT)
1066005	0.57 \pm 0.05	0.62 \pm 0.04	0.63 \pm 0.05	0.62 \pm 0.05	nan	0.54 \pm 0.03	0.54 \pm 0.03	0.55 \pm 0.03	0.56 \pm 0.06	nan
1066254	0.71 \pm 0.07	0.73 \pm 0.06	0.81 \pm 0.06	0.83 \pm 0.10	nan	0.65 \pm 0.05	0.64 \pm 0.09	0.77 \pm 0.07	0.86 \pm 0.11	nan
1119333	0.69 \pm 0.08	0.73 \pm 0.05	0.78 \pm 0.03	0.82 \pm 0.03	0.85 \pm 0.02	0.69 \pm 0.07	0.72 \pm 0.05	0.76 \pm 0.02	0.79 \pm 0.04	0.84 \pm 0.04
1243967	0.62 \pm 0.08	0.72 \pm 0.09	0.74 \pm 0.05	0.80 \pm 0.05	nan	0.63 \pm 0.07	0.66 \pm 0.05	0.71 \pm 0.04	0.76 \pm 0.05	nan
1243970	0.66 \pm 0.06	0.70 \pm 0.04	0.75 \pm 0.05	0.77 \pm 0.04	nan	0.63 \pm 0.06	0.65 \pm 0.04	0.70 \pm 0.02	0.75 \pm 0.05	nan
1613777	0.52 \pm 0.03	0.53 \pm 0.04	0.55 \pm 0.02	0.57 \pm 0.03	0.59 \pm 0.02	0.52 \pm 0.04	0.53 \pm 0.02	0.54 \pm 0.02	0.57 \pm 0.02	0.59 \pm 0.02
1613800	0.42 \pm 0.02	0.43 \pm 0.03	0.43 \pm 0.02	0.45 \pm 0.02	0.47 \pm 0.02	0.42 \pm 0.03	0.42 \pm 0.01	0.43 \pm 0.03	0.45 \pm 0.02	0.47 \pm 0.01
1613898	0.53 \pm 0.03	0.54 \pm 0.06	0.56 \pm 0.04	0.60 \pm 0.07	nan	0.55 \pm 0.04	0.54 \pm 0.06	0.58 \pm 0.04	0.61 \pm 0.08	nan
1613907	0.62 \pm 0.06	0.62 \pm 0.08	0.63 \pm 0.08	0.66 \pm 0.14	nan	0.60 \pm 0.05	0.62 \pm 0.07	0.67 \pm 0.06	0.70 \pm 0.15	nan
1613926	0.67 \pm 0.07	0.69 \pm 0.06	0.74 \pm 0.06	0.85 \pm 0.14	nan	0.61 \pm 0.04	0.60 \pm 0.06	0.69 \pm 0.06	0.73 \pm 0.19	nan
1613949	0.47 \pm 0.08	0.46 \pm 0.06	0.54 \pm 0.04	0.62 \pm 0.12	nan	0.45 \pm 0.06	0.45 \pm 0.06	0.52 \pm 0.06	0.56 \pm 0.12	nan
1614027	0.54 \pm 0.03	0.56 \pm 0.04	0.60 \pm 0.03	0.66 \pm 0.02	0.69 \pm 0.03	0.53 \pm 0.03	0.56 \pm 0.04	0.59 \pm 0.03	0.63 \pm 0.02	0.66 \pm 0.02
1614153	0.36 \pm 0.01	0.38 \pm 0.03	0.38 \pm 0.03	0.40 \pm 0.02	0.40 \pm 0.02	0.36 \pm 0.02	0.37 \pm 0.03	0.37 \pm 0.02	0.38 \pm 0.03	0.39 \pm 0.01
1614292	0.37 \pm 0.03	0.37 \pm 0.02	0.38 \pm 0.02	0.38 \pm 0.02	0.39 \pm 0.01	0.36 \pm 0.02	0.37 \pm 0.01	0.38 \pm 0.02	0.38 \pm 0.01	0.38 \pm 0.02
1614423	0.66 \pm 0.12	0.68 \pm 0.10	0.72 \pm 0.05	0.77 \pm 0.04	0.82 \pm 0.02	0.52 \pm 0.05	0.54 \pm 0.05	0.59 \pm 0.05	0.67 \pm 0.04	0.74 \pm 0.04
1614433	0.45 \pm 0.04	0.47 \pm 0.06	0.48 \pm 0.04	0.50 \pm 0.02	0.55 \pm 0.04	0.45 \pm 0.04	0.47 \pm 0.05	0.47 \pm 0.04	0.49 \pm 0.03	0.53 \pm 0.04
1614466	0.47 \pm 0.04	0.48 \pm 0.05	0.49 \pm 0.03	0.49 \pm 0.05	0.50 \pm 0.03	0.46 \pm 0.05	0.47 \pm 0.05	0.47 \pm 0.04	0.48 \pm 0.04	0.48 \pm 0.02
1614503	0.47 \pm 0.08	0.48 \pm 0.09	0.52 \pm 0.05	0.73 \pm 0.21	nan	0.44 \pm 0.05	0.44 \pm 0.06	0.48 \pm 0.08	0.66 \pm 0.22	nan
1614508	0.74 \pm 0.10	0.85 \pm 0.02	0.87 \pm 0.02	0.90 \pm 0.04	nan	0.76 \pm 0.07	0.85 \pm 0.04	0.87 \pm 0.03	0.91 \pm 0.05	nan
1614522	0.58 \pm 0.04	0.61 \pm 0.04	0.62 \pm 0.03	0.66 \pm 0.01	0.66 \pm 0.03	0.54 \pm 0.04	0.55 \pm 0.02	0.56 \pm 0.02	0.58 \pm 0.03	0.60 \pm 0.02
1737951	0.64 \pm 0.10	0.67 \pm 0.06	0.79 \pm 0.05	0.85 \pm 0.08	nan	0.55 \pm 0.06	0.60 \pm 0.08	0.63 \pm 0.05	0.72 \pm 0.07	nan
1738079	0.50 \pm 0.02	0.50 \pm 0.04	0.49 \pm 0.03	0.48 \pm 0.03	nan	0.52 \pm 0.04	0.52 \pm 0.04	0.53 \pm 0.03	0.56 \pm 0.05	nan
1738362	0.50 \pm 0.10	0.46 \pm 0.08	0.59 \pm 0.09	0.75 \pm 0.20	nan	0.52 \pm 0.07	0.56 \pm 0.11	0.62 \pm 0.05	0.84 \pm 0.14	nan
1738395	0.51 \pm 0.04	0.49 \pm 0.05	0.50 \pm 0.04	0.54 \pm 0.04	nan	0.49 \pm 0.04	0.48 \pm 0.03	0.50 \pm 0.05	0.51 \pm 0.05	nan
1738485	0.56 \pm 0.02	0.55 \pm 0.03	0.58 \pm 0.06	0.59 \pm 0.04	0.65 \pm 0.06	0.54 \pm 0.04	0.56 \pm 0.03	0.58 \pm 0.04	0.60 \pm 0.06	0.68 \pm 0.06
1738502	0.47 \pm 0.06	0.50 \pm 0.06	0.53 \pm 0.05	0.53 \pm 0.03	0.58 \pm 0.03	0.37 \pm 0.03	0.39 \pm 0.02	0.42 \pm 0.03	0.45 \pm 0.02	0.49 \pm 0.04
1738573	0.51 \pm 0.02	0.52 \pm 0.03	0.54 \pm 0.03	0.55 \pm 0.02	0.57 \pm 0.01	0.52 \pm 0.02	0.53 \pm 0.01	0.53 \pm 0.01	0.54 \pm 0.01	0.56 \pm 0.02
1738579	0.58 \pm 0.05	0.59 \pm 0.05	0.59 \pm 0.04	0.67 \pm 0.04	nan	0.52 \pm 0.03	0.54 \pm 0.04	0.56 \pm 0.03	0.62 \pm 0.04	nan
1738633	0.59 \pm 0.06	0.68 \pm 0.04	0.71 \pm 0.06	0.70 \pm 0.08	nan	0.57 \pm 0.05	0.60 \pm 0.08	0.60 \pm 0.05	0.71 \pm 0.10	nan
1794324	0.52 \pm 0.04	0.54 \pm 0.03	0.55 \pm 0.03	0.57 \pm 0.01	0.60 \pm 0.01	0.53 \pm 0.03	0.54 \pm 0.03	0.56 \pm 0.03	0.58 \pm 0.02	0.60 \pm 0.02
1794504	0.72 \pm 0.06	0.75 \pm 0.05	0.81 \pm 0.04	0.85 \pm 0.24	nan	0.67 \pm 0.04	0.67 \pm 0.06	0.73 \pm 0.05	0.90 \pm 0.21	nan
1794519	0.69 \pm 0.08	0.73 \pm 0.05	0.79 \pm 0.02	0.81 \pm 0.05	nan	0.70 \pm 0.06	0.72 \pm 0.05	0.80 \pm 0.03	0.77 \pm 0.07	nan
1794557	0.52 \pm 0.04	0.53 \pm 0.03	0.56 \pm 0.04	0.58 \pm 0.04	0.60 \pm 0.02	0.51 \pm 0.03	0.53 \pm 0.02	0.53 \pm 0.03	0.53 \pm 0.01	0.54 \pm 0.03
1963701	0.54 \pm 0.03	0.54 \pm 0.02	0.58 \pm 0.04	0.61 \pm 0.04	0.65 \pm 0.03	0.53 \pm 0.03	0.54 \pm 0.04	0.56 \pm 0.03	0.58 \pm 0.03	0.63 \pm 0.03
1963705	0.67 \pm 0.09	0.75 \pm 0.05	0.78 \pm 0.03	0.82 \pm 0.03	0.84 \pm 0.02	0.65 \pm 0.07	0.70 \pm 0.06	0.73 \pm 0.04	0.78 \pm 0.03	0.82 \pm 0.02
1963715	0.61 \pm 0.05	0.59 \pm 0.08	0.68 \pm 0.04	0.72 \pm 0.04	0.75 \pm 0.02	0.50 \pm 0.06	0.50 \pm 0.06	0.56 \pm 0.04	0.60 \pm 0.03	0.65 \pm 0.03
1963721	0.45 \pm 0.08	0.48 \pm 0.08	0.55 \pm 0.06	0.57 \pm 0.07	0.71 \pm 0.09	0.38 \pm 0.05	0.42 \pm 0.04	0.44 \pm 0.03	0.50 \pm 0.05	0.61 \pm 0.08
1963723	0.70 \pm 0.10	0.75 \pm 0.09	0.80 \pm 0.03	0.84 \pm 0.02	0.86 \pm 0.02	0.61 \pm 0.05	0.69 \pm 0.05	0.74 \pm 0.04	0.78 \pm 0.03	0.81 \pm 0.02
1963731	0.74 \pm 0.06	0.79 \pm 0.03	0.82 \pm 0.02	0.87 \pm 0.02	0.89 \pm 0.02	0.63 \pm 0.08	0.70 \pm 0.04	0.75 \pm 0.02	0.79 \pm 0.03	0.82 \pm 0.02
1963741	0.60 \pm 0.05	0.65 \pm 0.05	0.70 \pm 0.04	0.75 \pm 0.03	0.78 \pm 0.02	0.51 \pm 0.08	0.57 \pm 0.06	0.59 \pm 0.05	0.64 \pm 0.04	0.71 \pm 0.03
1963756	0.63 \pm 0.09	0.75 \pm 0.04	0.77 \pm 0.03	0.80 \pm 0.02	0.83 \pm 0.02	0.37 \pm 0.07	0.67 \pm 0.05	0.69 \pm 0.06	0.75 \pm 0.02	0.79 \pm 0.02
1963773	0.59 \pm 0.06	0.61 \pm 0.07	0.65 \pm 0.04	0.69 \pm 0.03	0.72 \pm 0.02	0.54 \pm 0.06	0.57 \pm 0.06	0.61 \pm 0.05	0.65 \pm 0.03	0.70 \pm 0.03
1963799	0.58 \pm 0.07	0.64 \pm 0.04	0.68 \pm 0.04	0.71 \pm 0.03	0.74 \pm 0.04	0.48 \pm 0.06	0.55 \pm 0.05	0.58 \pm 0.07	0.63 \pm 0.05	0.68 \pm 0.07
1963810	0.78 \pm 0.06	0.81 \pm 0.05	0.85 \pm 0.02	0.86 \pm 0.02	0.88 \pm 0.01	0.67 \pm 0.06	0.73 \pm 0.05	0.76 \pm 0.06	0.80 \pm 0.04	0.84 \pm 0.02
1963818	0.71 \pm 0.05	0.74 \pm 0.04	0.78 \pm 0.03	0.80 \pm 0.03	0.82 \pm 0.02	0.66 \pm 0.07	0.68 \pm 0.06	0.72 \pm 0.03	0.76 \pm 0.03	0.79 \pm 0.02
1963819	0.61 \pm 0.07	0.67 \pm 0.06	0.73 \pm 0.06	0.77 \pm 0.05	0.82 \pm 0.01	0.48 \pm 0.05	0.56 \pm 0.07	0.60 \pm 0.08	0.68 \pm 0.06	0.77 \pm 0.02
1963824	0.59 \pm 0.05	0.66 \pm 0.05	0.70 \pm 0.03	0.74 \pm 0.03	0.77 \pm 0.05	0.54 \pm 0.04	0.59 \pm 0.05	0.64 \pm 0.04	0.64 \pm 0.04	0.69 \pm 0.03
1963825	0.62 \pm 0.06	0.69 \pm 0.04	0.70 \pm 0.05	0.75 \pm 0.05	0.81 \pm 0.04	0.57 \pm 0.04	0.61 \pm 0.04	0.64 \pm 0.04	0.65 \pm 0.04	0.69 \pm 0.03
1963827	0.74 \pm 0.09	0.78 \pm 0.06	0.84 \pm 0.03	0.86 \pm 0.01	0.88 \pm 0.02	0.71 \pm 0.06	0.72 \pm 0.08	0.79 \pm 0.05	0.82 \pm 0.02	0.84 \pm 0.02
1963831	0.61 \pm 0.10	0.69 \pm 0.07	0.76 \pm 0.02	0.78 \pm 0.04	0.81 \pm 0.05	0.55 \pm 0.11	0.61 \pm 0.07	0.69 \pm 0.06	0.74 \pm 0.05	0.77 \pm 0.04
1963838	0.56 \pm 0.04	0.57 \pm 0.05	0.59 \pm 0.06	0.67 \pm 0.05	nan	0.54 \pm 0.05	0.55 \pm 0.04	0.58 \pm 0.04	0.62 \pm 0.06	nan

Table 8: PROTO Results measuring Δ AUPRC on the first 50 tasks in the test set of FS-Mol.

TASK-ID	16 (IGNITE-PROTO)	32 (IGNITE-PROTO)	64 (IGNITE-PROTO)	128 (IGNITE-PROTO)	256 (IGNITE-PROTO)	16 (PROTO)	32 (PROTO)	64 (PROTO)	128 (PROTO)	256 (PROTO)
1006005	0.61 ± 0.05	0.63 ± 0.04	0.66 ± 0.04	0.69 ± 0.06	nan	0.56 ± 0.06	0.61 ± 0.05	0.65 ± 0.03	0.66 ± 0.07	nan
1066254	0.72 ± 0.07	0.77 ± 0.06	0.84 ± 0.04	0.87 ± 0.10	nan	0.71 ± 0.08	0.78 ± 0.06	0.82 ± 0.06	0.85 ± 0.08	nan
1119333	0.73 ± 0.05	0.74 ± 0.06	0.78 ± 0.03	0.82 ± 0.02	0.84 ± 0.04	0.71 ± 0.06	0.72 ± 0.04	0.78 ± 0.03	0.81 ± 0.02	0.83 ± 0.04
1243967	0.78 ± 0.03	0.81 ± 0.04	0.83 ± 0.04	0.84 ± 0.03	nan	0.73 ± 0.05	0.76 ± 0.04	0.80 ± 0.03	0.81 ± 0.03	nan
1243970	0.75 ± 0.05	0.79 ± 0.04	0.81 ± 0.03	0.85 ± 0.03	nan	0.71 ± 0.07	0.74 ± 0.04	0.81 ± 0.03	0.82 ± 0.05	nan
1613777	0.52 ± 0.02	0.54 ± 0.03	0.56 ± 0.03	0.59 ± 0.02	0.61 ± 0.02	0.52 ± 0.04	0.54 ± 0.04	0.57 ± 0.03	0.60 ± 0.04	0.62 ± 0.03
1613800	0.42 ± 0.02	0.44 ± 0.02	0.46 ± 0.02	0.48 ± 0.02	0.50 ± 0.01	0.41 ± 0.03	0.43 ± 0.03	0.45 ± 0.02	0.46 ± 0.02	0.48 ± 0.02
1613898	0.52 ± 0.02	0.55 ± 0.04	0.54 ± 0.04	0.61 ± 0.07	nan	0.55 ± 0.04	0.55 ± 0.05	0.55 ± 0.03	0.57 ± 0.08	nan
1613907	0.56 ± 0.04	0.58 ± 0.06	0.63 ± 0.04	0.64 ± 0.14	nan	0.55 ± 0.06	0.62 ± 0.06	0.66 ± 0.07	0.68 ± 0.10	nan
1613926	0.70 ± 0.04	0.74 ± 0.06	0.77 ± 0.04	0.88 ± 0.12	nan	0.65 ± 0.08	0.70 ± 0.06	0.76 ± 0.04	0.86 ± 0.15	nan
1613949	0.57 ± 0.08	0.63 ± 0.05	0.67 ± 0.06	0.67 ± 0.10	nan	0.53 ± 0.05	0.57 ± 0.04	0.58 ± 0.05	0.63 ± 0.10	nan
1614027	0.53 ± 0.02	0.57 ± 0.02	0.60 ± 0.03	0.64 ± 0.02	0.67 ± 0.02	0.55 ± 0.02	0.57 ± 0.04	0.61 ± 0.04	0.64 ± 0.02	0.67 ± 0.02
1614153	0.35 ± 0.02	0.37 ± 0.01	0.36 ± 0.02	0.39 ± 0.01	0.41 ± 0.01	0.34 ± 0.03	0.37 ± 0.02	0.37 ± 0.02	0.39 ± 0.01	0.41 ± 0.01
1614292	0.36 ± 0.02	0.38 ± 0.03	0.39 ± 0.02	0.40 ± 0.02	0.41 ± 0.02	0.36 ± 0.02	0.39 ± 0.02	0.41 ± 0.02	0.42 ± 0.03	0.44 ± 0.02
1614423	0.68 ± 0.06	0.72 ± 0.04	0.76 ± 0.02	0.78 ± 0.02	0.79 ± 0.04	0.54 ± 0.07	0.58 ± 0.05	0.67 ± 0.04	0.73 ± 0.02	0.77 ± 0.02
1614433	0.49 ± 0.04	0.50 ± 0.04	0.51 ± 0.03	0.53 ± 0.02	0.54 ± 0.02	0.46 ± 0.04	0.48 ± 0.02	0.48 ± 0.03	0.52 ± 0.03	0.54 ± 0.04
1614466	0.47 ± 0.03	0.48 ± 0.04	0.49 ± 0.01	0.51 ± 0.03	0.53 ± 0.03	0.48 ± 0.05	0.53 ± 0.05	0.53 ± 0.04	0.57 ± 0.02	0.58 ± 0.01
1614503	0.49 ± 0.03	0.51 ± 0.07	0.56 ± 0.06	0.70 ± 0.22	nan	0.49 ± 0.07	0.55 ± 0.05	0.59 ± 0.02	0.76 ± 0.22	nan
1614508	0.84 ± 0.02	0.85 ± 0.03	0.86 ± 0.03	0.89 ± 0.04	nan	0.86 ± 0.02	0.86 ± 0.02	0.87 ± 0.02	0.88 ± 0.04	nan
1614522	0.58 ± 0.04	0.61 ± 0.03	0.60 ± 0.03	0.64 ± 0.02	0.67 ± 0.03	0.58 ± 0.03	0.60 ± 0.03	0.61 ± 0.03	0.64 ± 0.02	0.67 ± 0.02
1737951	0.65 ± 0.08	0.71 ± 0.06	0.76 ± 0.04	0.80 ± 0.06	nan	0.63 ± 0.06	0.67 ± 0.05	0.71 ± 0.04	0.75 ± 0.07	nan
1738079	0.49 ± 0.03	0.49 ± 0.03	0.50 ± 0.04	0.49 ± 0.04	nan	0.49 ± 0.02	0.49 ± 0.03	0.50 ± 0.03	0.54 ± 0.06	nan
1738362	0.46 ± 0.06	0.55 ± 0.08	0.59 ± 0.06	0.64 ± 0.22	nan	0.49 ± 0.10	0.56 ± 0.09	0.65 ± 0.05	0.82 ± 0.17	nan
1738395	0.53 ± 0.05	0.51 ± 0.04	0.52 ± 0.03	0.57 ± 0.04	nan	0.53 ± 0.04	0.53 ± 0.04	0.54 ± 0.05	0.58 ± 0.06	nan
1738485	0.54 ± 0.04	0.57 ± 0.04	0.58 ± 0.05	0.63 ± 0.04	0.69 ± 0.05	0.58 ± 0.03	0.59 ± 0.03	0.61 ± 0.04	0.62 ± 0.03	0.64 ± 0.07
1738502	0.41 ± 0.07	0.44 ± 0.05	0.47 ± 0.03	0.49 ± 0.03	0.53 ± 0.03	0.42 ± 0.07	0.45 ± 0.05	0.48 ± 0.03	0.49 ± 0.02	0.52 ± 0.02
1738573	0.52 ± 0.03	0.53 ± 0.03	0.54 ± 0.02	0.56 ± 0.02	0.57 ± 0.02	0.53 ± 0.02	0.52 ± 0.02	0.53 ± 0.02	0.55 ± 0.01	0.55 ± 0.01
1738579	0.52 ± 0.04	0.55 ± 0.05	0.58 ± 0.02	0.65 ± 0.06	nan	0.55 ± 0.06	0.59 ± 0.03	0.59 ± 0.03	0.64 ± 0.04	nan
1738633	0.56 ± 0.04	0.58 ± 0.04	0.59 ± 0.05	0.63 ± 0.08	nan	0.56 ± 0.05	0.59 ± 0.05	0.59 ± 0.04	0.62 ± 0.13	nan
1794324	0.50 ± 0.03	0.51 ± 0.03	0.53 ± 0.03	0.57 ± 0.02	0.60 ± 0.02	0.53 ± 0.03	0.54 ± 0.02	0.56 ± 0.03	0.59 ± 0.02	0.62 ± 0.02
1794504	0.70 ± 0.06	0.72 ± 0.06	0.80 ± 0.16	0.80 ± 0.16	nan	0.64 ± 0.09	0.72 ± 0.09	0.95 ± 0.16	nan	nan
1794519	0.65 ± 0.08	0.72 ± 0.06	0.78 ± 0.04	0.82 ± 0.05	nan	0.73 ± 0.06	0.78 ± 0.03	0.81 ± 0.02	0.79 ± 0.04	nan
1794557	0.54 ± 0.02	0.54 ± 0.02	0.55 ± 0.02	0.57 ± 0.03	0.60 ± 0.02	0.52 ± 0.02	0.53 ± 0.02	0.53 ± 0.03	0.54 ± 0.03	0.56 ± 0.02
1963701	0.58 ± 0.05	0.59 ± 0.06	0.62 ± 0.03	0.67 ± 0.02	0.69 ± 0.02	0.56 ± 0.04	0.58 ± 0.05	0.60 ± 0.04	0.62 ± 0.03	0.66 ± 0.02
1963705	0.74 ± 0.06	0.77 ± 0.04	0.81 ± 0.02	0.83 ± 0.01	0.84 ± 0.01	0.71 ± 0.05	0.74 ± 0.03	0.78 ± 0.03	0.80 ± 0.02	0.81 ± 0.02
1963715	0.63 ± 0.03	0.68 ± 0.05	0.72 ± 0.02	0.76 ± 0.02	0.78 ± 0.02	0.57 ± 0.02	0.61 ± 0.03	0.65 ± 0.03	0.69 ± 0.02	0.72 ± 0.02
1963721	0.47 ± 0.07	0.52 ± 0.04	0.54 ± 0.04	0.58 ± 0.06	0.63 ± 0.08	0.45 ± 0.06	0.47 ± 0.04	0.51 ± 0.04	0.53 ± 0.03	0.61 ± 0.10
1963723	0.77 ± 0.05	0.80 ± 0.04	0.84 ± 0.01	0.86 ± 0.01	0.87 ± 0.01	0.72 ± 0.08	0.79 ± 0.03	0.81 ± 0.01	0.83 ± 0.01	0.85 ± 0.01
1963731	0.79 ± 0.06	0.84 ± 0.02	0.87 ± 0.02	0.89 ± 0.01	0.90 ± 0.01	0.80 ± 0.05	0.83 ± 0.02	0.86 ± 0.01	0.87 ± 0.01	0.88 ± 0.01
1963741	0.64 ± 0.06	0.70 ± 0.04	0.74 ± 0.02	0.76 ± 0.02	0.78 ± 0.02	0.59 ± 0.03	0.63 ± 0.04	0.68 ± 0.04	0.71 ± 0.02	0.73 ± 0.02
1963756	0.74 ± 0.06	0.77 ± 0.03	0.79 ± 0.02	0.82 ± 0.01	0.84 ± 0.01	0.63 ± 0.06	0.68 ± 0.04	0.73 ± 0.02	0.76 ± 0.02	0.77 ± 0.03
1963773	0.62 ± 0.06	0.64 ± 0.06	0.69 ± 0.03	0.74 ± 0.03	0.75 ± 0.03	0.65 ± 0.06	0.69 ± 0.04	0.73 ± 0.02	0.76 ± 0.03	0.78 ± 0.03
1963799	0.65 ± 0.07	0.67 ± 0.05	0.69 ± 0.02	0.74 ± 0.02	0.76 ± 0.03	0.56 ± 0.07	0.63 ± 0.04	0.67 ± 0.03	0.69 ± 0.02	0.72 ± 0.05
1963810	0.80 ± 0.04	0.83 ± 0.03	0.86 ± 0.01	0.86 ± 0.01	0.88 ± 0.01	0.76 ± 0.06	0.79 ± 0.04	0.82 ± 0.02	0.85 ± 0.01	0.87 ± 0.01
1963818	0.71 ± 0.06	0.75 ± 0.05	0.78 ± 0.03	0.82 ± 0.02	0.84 ± 0.02	0.71 ± 0.03	0.73 ± 0.04	0.78 ± 0.02	0.81 ± 0.01	0.82 ± 0.01
1963819	0.65 ± 0.08	0.68 ± 0.06	0.74 ± 0.03	0.78 ± 0.03	0.81 ± 0.02	0.62 ± 0.12	0.70 ± 0.05	0.73 ± 0.05	0.79 ± 0.02	0.81 ± 0.02
1963824	0.60 ± 0.06	0.63 ± 0.04	0.67 ± 0.04	0.72 ± 0.04	0.75 ± 0.03	0.52 ± 0.04	0.54 ± 0.03	0.57 ± 0.03	0.61 ± 0.02	0.63 ± 0.03
1963825	0.66 ± 0.06	0.70 ± 0.05	0.72 ± 0.03	0.74 ± 0.02	0.78 ± 0.04	0.62 ± 0.07	0.68 ± 0.03	0.71 ± 0.04	0.73 ± 0.02	0.77 ± 0.03
1963827	0.78 ± 0.06	0.81 ± 0.03	0.86 ± 0.02	0.88 ± 0.00	0.88 ± 0.01	0.77 ± 0.04	0.80 ± 0.04	0.84 ± 0.02	0.86 ± 0.01	0.87 ± 0.01
1963831	0.62 ± 0.10	0.69 ± 0.07	0.74 ± 0.03	0.77 ± 0.03	0.78 ± 0.05	0.61 ± 0.08	0.69 ± 0.05	0.76 ± 0.03	0.78 ± 0.02	0.79 ± 0.06
1963838	0.52 ± 0.03	0.55 ± 0.02	0.58 ± 0.03	0.60 ± 0.03	nan	0.55 ± 0.03	0.57 ± 0.05	0.61 ± 0.04	0.64 ± 0.03	nan

Table 9: MAML results measuring Δ AUPRC on the first 50 tasks in the test set of FS-Mol.

TASK-ID	16 (IGNITE-MAML)	32 (IGNITE-MAML)	64 (IGNITE-MAML)	128 (IGNITE-MAML)	256 (IGNITE-MAML)	16 (MAML)	32 (MAML)	64 (MAML)	128 (MAML)	256 (MAML)
1006005	0.47 ± 0.03	0.49 ± 0.05	0.49 ± 0.04	0.54 ± 0.06	nan	0.49 ± 0.02	0.51 ± 0.04	0.52 ± 0.03	0.56 ± 0.05	nan
1066254	0.58 ± 0.07	0.61 ± 0.05	0.65 ± 0.07	0.65 ± 0.13	nan	0.56 ± 0.07	0.58 ± 0.04	0.59 ± 0.07	0.63 ± 0.15	nan
1119333	0.73 ± 0.03	0.74 ± 0.03	0.77 ± 0.04	0.76 ± 0.05	0.77 ± 0.06	0.70 ± 0.03	0.72 ± 0.04	0.74 ± 0.05	0.72 ± 0.02	0.74 ± 0.05
1243967	0.72 ± 0.03	0.76 ± 0.03	0.75 ± 0.04	0.77 ± 0.04	nan	0.72 ± 0.05	0.75 ± 0.01	0.74 ± 0.04	0.76 ± 0.04	nan
1243970	0.70 ± 0.04	0.71 ± 0.01	0.70 ± 0.02	0.71 ± 0.03	nan	0.71 ± 0.03	0.72 ± 0.03	0.72 ± 0.03	0.68 ± 0.03	nan
1613777	0.53 ± 0.01	0.53 ± 0.01	0.53 ± 0.01	0.54 ± 0.02	0.54 ± 0.02	0.50 ± 0.01	0.51 ± 0.01	0.51 ± 0.01	0.53 ± 0.02	0.53 ± 0.02
1613800	0.49 ± 0.01	0.40 ± 0.01	0.41 ± 0.01	0.42 ± 0.01	0.42 ± 0.01	0.40 ± 0.01	0.40 ± 0.01	0.40 ± 0.01	0.41 ± 0.01	0.41 ± 0.01
1613898	0.56 ± 0.04	0.55 ± 0.04	0.58 ± 0.04	0.57 ± 0.08	nan	0.57 ± 0.02	0.56 ± 0.03	0.56 ± 0.03	0.56 ± 0.04	nan
1613907	0.54 ± 0.03	0.55 ± 0.04	0.54 ± 0.06	0.55 ± 0.07	nan	0.53 ± 0.03	0.51 ± 0.02	0.53 ± 0.05	0.55 ± 0.13	nan
1613926	0.67 ± 0.04	0.65 ± 0.06	0.71 ± 0.06	0.78 ± 0.20	nan	0.57 ± 0.09	0.59 ± 0.10	0.65 ± 0.12	0.77 ± 0.12	nan
1613949	0.46 ± 0.04	0.47 ± 0.04	0.47 ± 0.05	0.45 ± 0.10	nan	0.48 ± 0.03	0.48 ± 0.04	0.48 ± 0.06	0.43 ± 0.09	nan
1614027	0.51 ± 0.02	0.52 ± 0.02	0.54 ± 0.02	0.57 ± 0.03	0.60 ± 0.02	0.52 ± 0.01	0.52 ± 0.01	0.53 ± 0.03	0.54 ± 0.03	0.58 ± 0.04
1614153	0.33 ± 0.01	0.34 ± 0.01	0.33 ± 0.01	0.34 ± 0.01	0.34 ± 0.01	0.33 ± 0.01	0.34 ± 0.01	0.34 ± 0.01	0.35 ± 0.01	0.36 ± 0.01
1614292	0.35 ± 0.02	0.36 ± 0.02	0.37 ± 0.02	0.37 ± 0.02	0.39 ± 0.01	0.35 ± 0.01	0.35 ± 0.01	0.35 ± 0.01	0.35 ± 0.01	0.35 ± 0.01
1614423	0.44 ± 0.07	0.52 ± 0.10	0.64 ± 0.03	0.69 ± 0.05	0.71 ± 0.04	0.45 ± 0.05	0.47 ± 0.07	0.51 ± 0.07	0.60 ± 0.07	0.68 ± 0.02
1614433	0.46 ± 0.02	0.46 ± 0.02	0.47 ± 0.02	0.49 ± 0.03	0.52 ± 0.04	0.44 ± 0.01	0.46 ± 0.03	0.47 ± 0.03	0.48 ± 0.04	0.49 ± 0.04
1614466	0.46 ± 0.01	0.46 ± 0.02	0.47 ± 0.02	0.47 ± 0.02	0.51 ± 0.03	0.47 ± 0.01	0.48 ± 0.02	0.49 ± 0.02	0.49 ± 0.02	0.51 ± 0.03
1614503	0.35 ± 0.04	0.36 ± 0.06	0.37 ± 0.04	0.36 ± 0.24	nan	0.37 ± 0.02	0.36 ± 0.03	0.38 ± 0.04	0.57 ± 0.20	nan
1614508	0.70 ± 0.12	0.77 ± 0.13	0.80 ± 0.05	0.82 ± 0.09	nan	0.66 ± 0.08	0.79 ± 0.07	0.77 ± 0.09	0.83 ± 0.09	nan
1614522	0.52 ± 0.05	0.53 ± 0.03	0.54 ± 0.06	0.61 ± 0.02	0.62 ± 0.02	0.50 ± 0.03	0.52 ± 0.02	0.52 ± 0.04	0.56 ± 0.03	0.58 ± 0.03
1737951	0.49 ± 0.09	0.46 ± 0.07	0.61 ± 0.03	0.67 ± 0.06	nan	0.57 ± 0.04	0.57 ± 0.03	0.58 ± 0.04	0.58 ± 0.07	nan
1738079	0.48 ± 0.01	0.47 ± 0.01	0.47 ± 0.02	0.49 ± 0.04	nan	0.52 ± 0.01	0.51 ± 0.01	0.51 ± 0.03	0.52 ± 0.03	nan
1738362	0.31 ± 0.03	0.34 ± 0.05	0.34 ± 0.07	0.64 ± 0.23	nan	0.37 ± 0.01	0.38 ± 0.03	0.40 ± 0.04	0.46 ± 0.19	nan
1738395	0.49 ± 0.03	0.50 ± 0.03	0.49 ± 0.03	0.50 ± 0.05	nan	0.50 ± 0.03	0.48 ± 0.03	0.48 ± 0.03	0.50 ± 0.05	nan
1738485	0.49 ± 0.01	0.51 ± 0.02	0.52 ± 0.04	0.54 ± 0.04	0.55 ± 0.06	0.49 ± 0.01	0.49 ± 0.01	0.48 ± 0.01	0.50 ± 0.06	0.54 ± 0.06
1738502	0.39 ± 0.02	0.40 ± 0.02	0.40 ± 0.02	0.41 ± 0.02	0.41 ± 0.02	0.41 ± 0.02	0.41 ± 0.02	0.41 ± 0.02	0.41 ± 0.02	0.41 ± 0.02
1738573	0.51 ± 0.02	0.50 ± 0.01	0.51 ± 0.02	0.52 ± 0.02	0.53 ± 0.02	0.52 ± 0.01	0.51 ± 0.01	0.52 ± 0.01	0.52 ± 0.01	0.52 ± 0.01
1738579	0.54 ± 0.03	0.56 ± 0.02	0.56 ± 0.03	0.60 ± 0.05	nan	0.53 ± 0.03	0.55 ± 0.03	0.55 ± 0.04	0.58 ± 0.04	nan
1738633	0.55 ± 0.04	0.56 ± 0.04	0.60 ± 0.04	0.57 ± 0.09	nan	0.59 ± 0.02	0.60 ± 0.02	0.60 ± 0.04	0.63 ± 0.12	nan
1794324	0.52 ± 0.02	0.52 ± 0.01	0.52 ± 0.01	0.53 ± 0.02	0.54 ± 0.02	0.51 ± 0.00	0.51 ± 0.01	0.52 ± 0.01	0.52 ± 0.01	0.53 ± 0.01
1794504	0.68 ± 0.01	0.70 ± 0.04	0.71 ± 0.04	0.85 ± 0.24	nan	0.65 ± 0.05	0.68 ± 0.05	0.71 ± 0.04	0.85 ± 0.24	nan
1794522	0.64 ± 0.07	0.66 ± 0.08	0.69 ± 0.08	0.75 ± 0.13	nan	0.64 ± 0.07	0.65 ± 0.07	0.65 ± 0.07	0.67 ± 0.11	nan
1794557	0.53 ± 0.02	0.54 ± 0.01	0.54 ± 0.02	0.53 ± 0.02	0.53 ± 0.01	0.55 ± 0.01	0.55 ± 0.01	0.55 ± 0.01	0.55 ± 0.01	0.55 ± 0.01
1963701	0.60 ± 0.01	0.59 ± 0.02	0.59 ± 0.02	0.60 ± 0.01	0.60 ± 0.02	0.60 ± 0.01	0.59 ± 0.01	0.59 ± 0.01	0.59 ± 0.01	0.59 ± 0.02
1963705	0.71 ± 0.03	0.70 ± 0.02	0.72 ± 0.01	0.72 ± 0.01	0.71 ± 0.02	0.77 ± 0.01	0.77 ± 0.02	0.77 ± 0.01	0.77 ± 0.03	0.76 ± 0.04
1963715	0.63 ± 0.00	0.62 ± 0.03	0.63 ± 0.01	0.64 ± 0.01	0.64 ± 0.02	0.63 ± 0.01	0.61 ± 0.04	0.64 ± 0.01	0.64 ± 0.01	0.64 ± 0.02
1963721	0.50 ± 0.02	0.50 ± 0.02	0.50 ± 0.02	0.51 ± 0.04	0.55 ± 0.07	0.52 ± 0.02	0.51 ± 0.01	0.52 ± 0.01	0.51 ± 0.03	0.54 ± 0.01
1963725	0.72 ± 0.03	0.72 ± 0.03	0.73 ± 0.02	0.74 ± 0.02	0.74 ± 0.02	0.72 ± 0.03	0.74 ± 0.02	0.74 ± 0.02	0.74 ± 0.02	0.74 ± 0.02
1963731	0.70 ± 0.02	0.70 ± 0.01	0.70 ± 0.02	0.76 ± 0.05	0.76 ± 0.05	0.71 ± 0.03	0.74 ± 0.02	0.74 ± 0.02	0.76 ± 0.02	0.79 ± 0.03
1963741	0.65 ± 0.02	0.66 ± 0.00	0.66 ± 0.01	0.67 ± 0.01	0.67 ± 0.01	0.65 ± 0.01	0.65 ± 0.01	0.66 ± 0.01	0.65 ± 0.01	0.65 ± 0.01
1963756	0.76 ± 0.05	0.75 ± 0.04	0.76 ± 0.04	0.77 ± 0.02	0.77 ± 0.01	0.71 ± 0.02	0.72 ± 0.03	0.70 ± 0.04	0.72 ± 0.02	0.72 ± 0.03
1963773	0.60 ± 0.01	0.62 ± 0.01	0.60 ± 0.03	0.60 ± 0.02	0.60 ± 0.02	0.62 ± 0.02	0.63 ± 0.02	0.63 ± 0.02	0.64 ± 0.03	0.63 ± 0.02
1963781	0.63 ± 0.01	0.59 ± 0.04	0.61 ± 0.02	0.61 ± 0.02	0.60 ± 0.02	0.64 ± 0.02	0.66 ± 0.01	0.66 ± 0.02	0.67 ± 0.03	0.67 ± 0.02
1963810	0.78 ± 0.01	0.77 ± 0.02	0.78 ± 0.01	0.78 ± 0.01	0.78 ± 0.01	0.77 ± 0.01	0.78 ± 0.03	0.76 ± 0.02	0.77 ± 0.02	0.77 ± 0.02
1963818	0.71 ± 0.02	0.72 ± 0.01	0.72 ± 0.03	0.73 ± 0.02	0.73 ± 0.03	0.70 ± 0.04	0.72 ± 0.01	0.72 ± 0.03	0.72 ± 0.03	0.73 ± 0.03
1963819	0.59 ± 0.04	0.61 ± 0.03	0.59 ± 0.01	0.61 ± 0.03	0.61 ± 0.03	0.60 ± 0.02	0.61 ± 0.03	0.62 ± 0.03	0.63 ± 0.04	0.63 ± 0.03
1963824	0.62 ± 0.04	0.63 ± 0.02	0.65 ± 0.04	0.65 ± 0.04	0.69 ± 0.05	0.62 ± 0.03	0.64 ± 0.04	0.64 ± 0.04	0.63 ± 0.03	0.63 ± 0.04
1963825	0.73 ± 0.01	0.72 ± 0.02	0.73 ± 0.01	0.74 ± 0.02	0.75 ± 0.04	0.77 ± 0.02	0.77 ± 0.02	0.79 ± 0.02	0.79 ± 0.03	0.80 ± 0.04
1963826	0.78 ± 0.01	0.80 ± 0.03	0.80 ± 0.03	0.80 ± 0.03	0.80 ± 0.02	0.78 ± 0.04	0.79 ± 0.04	0.79 ± 0.04	0.77 ± 0.04	0.76 ± 0.02
1963831	0.69 ± 0.01	0.69 ± 0.02	0.69 ± 0.02	0.69 ± 0.02	0.68 ± 0.06	0.68 ± 0.02	0.68 ± 0.01	0.70 ± 0.02	0.71 ± 0.03	0.70 ± 0.02
1963838	0.49 ± 0.01	0.49 ± 0.01	0.51 ± 0.04	0.55 ± 0.03	nan	0.54 ± 0.02	0.53 ± 0.01	0.55 ± 0.02	0.52 ± 0.04	nan