
FoMo-0D: A Foundation Model for Zero-shot Outlier Detection

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054

Anonymous Authors¹

Abstract

Outlier detection (OD) has a vast literature as it finds numerous real-world applications. Being an unsupervised task, model selection is a key bottleneck for OD without label supervision. Despite a long list of available OD algorithms with tunable hyperparameters, the lack of systematic approaches for unsupervised algorithm and hyperparameter selection limits their effective use in practice. In this paper, we present FoMo-0D, a pre-trained Foundation Model for zero/0-shot OD on tabular data, which bypasses the hurdle of model selection altogether. Having been pre-trained on synthetic data, FoMo-0D can directly predict the (outlier/inlier) label of test samples without parameter fine-tuning—*requiring no labeled data, and no additional training or hyperparameter tuning when given a new task*. Extensive experiments on **57** real-world datasets against **26** baselines show that FoMo-0D is highly competitive; outperforming the majority of the baselines with no statistically significant difference from the *2nd* best method. Further, FoMo-0D is efficient in inference time requiring only **7.7 ms** per sample on average, with at least **7x** speed-up compared to previous methods. To facilitate future research, our implementations for data synthesis and pre-training as well as model checkpoints are openly available at <https://anonymous.4open.science/r/PFN40D>.

1. Introduction

Outlier detection (OD) in tabular data finds numerous applications in critical domains such as security, environmental monitoring, finance, and medicine, to name a few. This popularity brings along a large literature with plethora of detection algorithms to choose from given a new OD task. These algorithms, however, exhibit several hyperparameters

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

(HPs) that need careful tuning (Ma et al., 2023). Since most OD tasks are unsupervised¹, what makes effective detection notoriously difficult is unsupervised model selection (both algorithm and HP selection) in the absence of labels.

While deep learning has revolutionized many areas of machine learning (ML), it is not quite the case for OD—mainly because compared to classical methods, deep OD models (Pang et al., 2021) have many more HPs that detection performance is sensitive to (Ding et al., 2022), rendering model selection even more challenging. While the recent success of large foundation models, pre-trained on massive amounts of data, offers new opportunities for zero-shot OD, thus far the most notable progress has been in NLP and computer vision (Brown et al., 2020; Touvron et al., 2023; Radford et al., 2021). This is thanks to the admirable quantity and quality of public text and image datasets. In comparison, public tabular OD benchmarks remain minuscule (Han et al., 2022; Zhao et al., 2021; Steinbuss & Böhm, 2021).

Recently, Prior-data Fitted Networks (PFNs) has marked a milestone in ML as a new approach to learning on tabular data (Müller et al., 2022). The core idea is to compute a posterior predictive distribution (PPD) for a test point given training data as context. To approximate the PPD, a Transformer (Vaswani et al., 2017) is pre-trained on a large set of synthetic datasets drawn from pre-defined data priors. At inference, the pre-trained PFN is fed with test samples along with some training samples as context for zero-shot prediction, requiring no parameter fine-tuning or model selection on new datasets.

In this paper, we capitalize on these ideas and introduce FoMo-0D: a prior-data fitted Foundation Model for zero/0-shot OD. Once pre-trained on synthetic datasets, FoMo-0D unlocks zero-shot OD on a new dataset where the (unlabeled) input data is fed only as context. As such, FoMo-0D bypasses not only model (parameter) training, but more importantly, the nontrivial task of unsupervised model (algorithm and HP) selection without labeled data. Figure 1 illustrates the new FoMo-0D paradigm versus the typical OD setting. To our knowledge, FoMo-0D is the first pre-trained foundation model for tabular OD.

¹While supervised OD exists, unsupervised setting is preferred in most domains to detect novel, emergent anomalies.

Table 1: p -values of the one-sided Wilcoxon signed rank test, comparing FoMo-0D (with $D = 100$) to **top 10 baselines** with default hyperparameters (HPs), and **top 4^{avg}** baselines³ with **avg.** performance over varying HPs (denoted w/ ^{avg}) over All (57) datasets, those (42) w/ $d \leq 100$ and (46) w/ $d \leq 500$ dimensions. FoMo-0D shows **no statistically significant difference from the 2nd best model** ($k\text{NN}$, w/ $p = 0.106$) over All datasets, while it is comparable to ($p > \alpha$) or significantly better than ($p > 1 - \alpha$) all 26 original + 4^{avg} baselines over datasets w/ $d \leq 100$ (aligned w/ pretraining where $D = 100$) as well as on datasets w/ $d \leq 500$ (generalizing beyond pretraining). Rank, avg.’ed over all 57 datasets by AUROC. (setting: $D = 100$, $R = 500$, train/inference context size=5K, w/ quantile transform) (See Tables 16.1&16.2 for full results.)

	FoMo-0D	DTE-NP	$k\text{NN}$	ICL	DTE-C	LOF	CBLOF	Feat.Bag.	SLAD	DDPM	OCSVM	DTE-NP ^{avg}	$k\text{NN}^{\text{avg}}$	ICL ^{avg}	DTE-C ^{avg}	
$d \leq 100$	-		0.415	0.700	0.949	0.953	0.970	0.971	0.996	0.876	0.980	0.978	0.752	0.860	0.958	1.000
$d \leq 500$	-		0.220	0.569	0.827	0.894	0.960	0.968	0.994	0.910	0.960	0.979	0.607	0.756	0.846	1.000
All	-		<u>0.016</u>	0.106	0.462	0.454	0.585	0.750	0.823	0.759	0.901	0.895	0.112	0.315	0.670	1.000
Rank(avg)	11.886		7.553	9.018	10.851	11.36	12.316	13.342	13.386	12.982	14.061	13.851	9.079	11.105	12.991	22.263

In designing FoMo-0D, we use Gaussian mixture models as a simple yet effective tabular data prior for inlier data distributions (Hollmann et al., 2023; Zhao et al., 2021), which are also employed to simulate outlier types common in the real world; namely, local and global subspace outliers (Steinbuss & Böhm, 2021). While the data prior can be extended to comprise more complex data distributions (Hollmann et al., 2023), and additional outlier types can be included (e.g. dependency, contextual, etc.), as we show in the experiments, even with its relatively straightforward prior, FoMo-0D achieves remarkable performance: As shown in Table 1, FoMo-0D, pre-trained on datasets with $d \leq 100$ dimensions, shows no statistically significant difference from all 26 state-of-the-art baselines (all p -values > 0.4) on 42 benchmark datasets with dimensionality $d \leq 100$ (aligned with pre-training), while our method consistently ranks among the top and outperforms a majority of the baselines with p -value > 0.95 . (See Appendix Tables 16.1&16.2 for full results.) The results remain consistent on (46) benchmarks with $d \leq 500$ dimensions. FoMo-0D is also competitive across all (57) datasets, effectively generalizing beyond its pre-training distributions, with no statistically significant difference from the 2nd best baseline. Further, FoMo-0D takes a mere 7.7 ms to infer a test sample on average with no extra training or tuning overhead on the new dataset.

2. Problem and Preliminaries

Semi-supervised Outlier Detection: We focus on semi-supervised OD. Formally, let $\mathcal{D}_{\text{in}} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ denote the input data containing only inliers $\mathbf{x}_i \in \mathbb{R}^d$, where $y_i = 0 \forall i \in [n]$, and $\mathcal{D}_{\text{test}}$ depicts the test data comprising both inliers and outliers. The task is to assign labels to $\mathbf{x}_i \in \mathcal{D}_{\text{test}}$ given the inlier-only input \mathcal{D}_{in} .

PFNs: (Müller et al., 2023) approximate the posterior predictive distribution (PPD) of a test point using training data as context, formulated as $p(\cdot | \mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})$. A Transformer is pre-trained on synthetic datasets from predefined priors and used at inference for zero-shot prediction with test and

context samples. More details in Appendix C.

Pre-training on synthetic data. Massive synthetic datasets are generated for the pre-training stage, by first sampling a hypothesis (i.e., the generating mechanism) $\phi \sim p(\phi)$, and then sampling a dataset $\mathcal{D} \sim p(\mathcal{D}|\phi)$. For training, each dataset \mathcal{D} can be split as $\mathcal{D}_{\text{test}} \subset \mathcal{D}$ and $\mathcal{D}_{\text{train}} = \mathcal{D} \setminus \mathcal{D}_{\text{test}}$. Thus, the PFN with parameters θ can be optimized by making predictions on data points in $\mathcal{D}_{\text{test}}$. For a test point $(\mathbf{x}_{\text{test}}, y_{\text{test}}) \in \mathcal{D}_{\text{test}}$, the training loss is as follows.

$$\mathcal{L} = \mathbb{E}_{\{(\mathbf{x}_{\text{test}}, y_{\text{test}})\} \cup \mathcal{D}_{\text{train}} \sim p(\mathcal{D})} [-\log q_{\theta}(y_{\text{test}} | \mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})]. \quad (1)$$

Inference on real-world data. At inference, a fresh real-world dataset $\mathcal{D}_{\text{train}}$ and test instance \mathbf{x}_{test} are fed into the pre-trained model, which computes the PPD $q_{\theta}(\cdot | \mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})$ in a single forward pass. Importantly, PFNs do not require gradient-based parameter tuning on new datasets, where prediction is delivered *in less than a second* (Hollmann et al., 2023).

3. FoMo-0D: A Foundation Model for Zero-shot Outlier Detection

3.1. Designing a Data Prior for Outlier Detection

We design a new data prior from which we simulate numerous OD datasets for pre-training FoMo-0D.

Inlier synthesis: We simulate inliers by drawing from a Gaussian Mixture Model (GMM) with m -clusters in d -dimensions, with centers $\boldsymbol{\mu}_{jk} \in [-5, 5]$, $j \in [m]$, $k \in [d]$ and *diagonal*² Σ_j with entries in $(0, 5]$. We create different GMMs with varying $m \leq M$ and $d \leq D$ chosen uniformly at random from $[M]$ and $[D]$, respectively. From each GMM, we draw a set of S inliers, defined as instances within the 90th percentile of the GMM.

Outlier synthesis: Following Han et al. (2022), we gener-

²In early experiments, we found no difference in test performance on synthetic datasets between using diagonal vs. non-diagonal Σ , yet, it is easier to invert diagonal Σ for data synthesis.

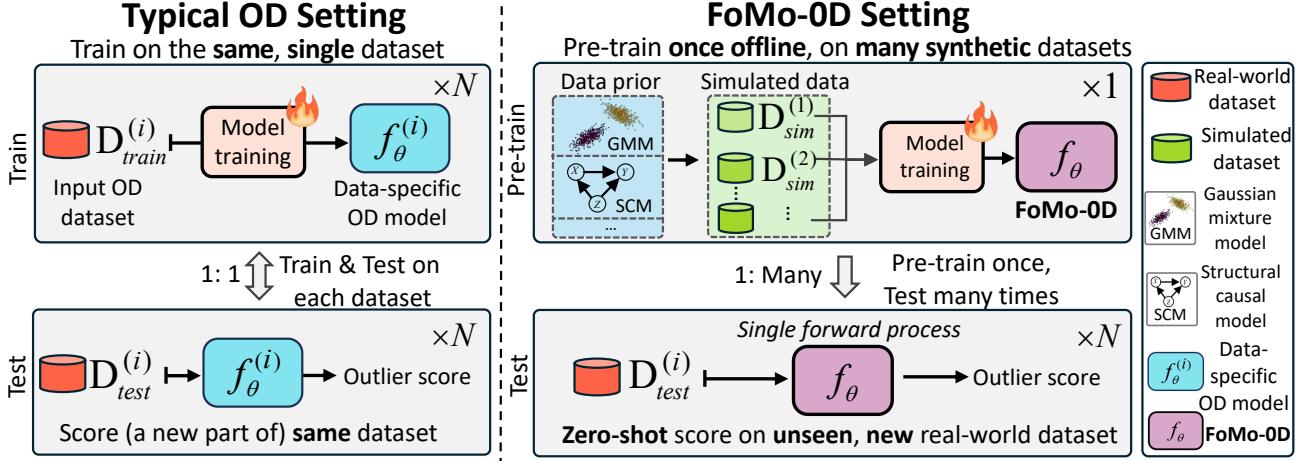


Figure 1: (best in color) Comparison of typical OD vs. the FoMo-OD settings. Given a new unlabeled OD dataset, FoMo-OD not only eliminates the need for model (parameter) training, but most importantly, also abolishes the onerous task of unsupervised model selection (algorithm and hyperparameters).

ate subspace outliers by first drawing a subset of dimensions \mathcal{K} at random, for $|\mathcal{K}| \leq d$, and then generate S points from the “inflated” GMM, which shares the same centers μ_j ’s with the original GMM but with the inflated (diagonal) covariances $5 \times \Sigma_{j,kk}$ ’s for $k \in \mathcal{K}$. Outliers are defined as points sampled outside the 90th percentile of the original GMM, which are labeled based on their Mahalanobis distances (see Property E.6 in the Appendix).

Specifically, we simulate datasets containing $2S = 10,000$ samples (half inlier, half outlier) from the two corresponding GMMs (original and inflated) with up to $M = 5$ clusters and up to $D = 100$ dimensions. Example 2-d synthetic datasets are illustrated in Appendix B.

3.2. (Pre)Training and Inference

Model (Pre)Training (Once, Offline): In the synthetic prior-data fitting phase, we first randomly draw a hypothesis (i.e. GMM configuration) uniformly at random, i.e., $\phi = \{d \in [D], m \in [M], \{\mu_j\}_{j=1}^m \in [-5, 5]^d, \{\Sigma_j\}_{j=1}^m; \text{diag}(\Sigma_j) \in [-5, 5]^d\}$, then generate a synthetic dataset $\mathcal{D} = \{\mathcal{D}_{\text{in}}, \mathcal{D}_{\text{out}}\}$ containing synthetic inlier and outlier samples from the drawn hypothesis and its variance-inflated variant, respectively.

We optimize FoMo-OD’s parameters θ to make predictions on $\mathcal{D}_{\text{test}} = \{\mathcal{D}_{\text{test}}^{\text{in}}, \mathcal{D}_{\text{test}}^{\text{out}}\}$, conditioned on the inlier-only training data $\mathcal{D}_{\text{train}} \subset \mathcal{D}_{\text{in}}$ based on the cross-entropy loss (see Eq. (1)). During training, $\mathcal{D}_{\text{test}}$ contains a *balanced* number of inlier and outlier samples, where $\mathcal{D}_{\text{test}}^{\text{in}} = \mathcal{D}_{\text{in}} \setminus \mathcal{D}_{\text{train}}$, and $\mathcal{D}_{\text{test}}^{\text{out}} \subset \mathcal{D}_{\text{out}}$ contains an equal number of samples as $\mathcal{D}_{\text{test}}^{\text{in}}$. To vary the training data size, we subsample $\mathcal{D}_{\text{train}}$ of randomly drawn size $n \in [n_L, n_U]$, where n_L and n_U denote the lower and upper bounds. In our implementation, we use $n_L = 500$, and $n_U = 5,000$.

Zero-shot Inference (on Unseen/New Dataset): At inference, the pre-trained FoMo-OD can be employed on any unseen real-world dataset. Specifically, for a new semi-supervised OD task with inlier-only training data $\mathcal{D}_{\text{train}}$ and mixed test data $\mathcal{D}_{\text{test}}$, feeding $(\mathcal{D}_{\text{train}}, \mathbf{x}_{\text{test}})$ as input to FoMo-OD (for each $\mathbf{x}_{\text{test}} \in \mathcal{D}_{\text{test}}$ separately) yields the PPD $q_{\theta}(y|\mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})$ in a *single forward pass*. As such, FoMo-OD performs model “training” and prediction simultaneously at test time. In fact, as the training data is passed as context, FoMo-OD leverages in-context learning (ICL) (Xie et al., 2021; Garg et al., 2022) for inference.

3.3. Architecture and Scalability

FoMo-OD is based on Transformer (Vaswani et al., 2017) model, with (1) self-attention among all the training samples and only cross-attention from test samples to the training samples; (2) “router attention mechanism” Zhang & Yan (2023) to reduce the complexity of self-attention to linear (As shown in Figure 2), thus extending the context length. We further propose to scale up (pre)training data synthesis with linear transforms, more details in Appendix D.

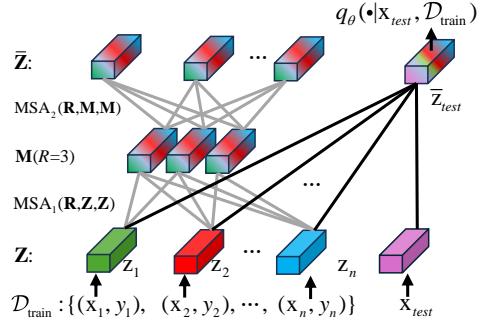


Figure 2: FoMo-OD architecture employs the “router mechanism” for scalable attention.

165 4. Experiments

166 4.1. Setup

168 **Pre-training Dataset Synthesis:** During pre-training, we
 169 generate unique GMM datasets by first drawing a configura-
 170 tion, including dimensionality $d \in [D]$, number of com-
 171 ponents $m \in [M]$, centers $\{\mu_j\}_{j=1}^m$ (each $\mu_j \in [-5, 5]^d$)
 172 and covariances $\{\Sigma_j\}_{j=1}^m$ ($\text{diag}(\Sigma_j) \in [-5, 5]^d$). We set
 173 $M = 5$ and vary $D \in \{20, 100\}$ to study pre-training with
 174 relatively small and high dimensional datasets, respectively.
 175 We synthesize inliers and outliers described in Section 3.1.
 176

177 **Real-world Benchmark Datasets:** While pre-training
 178 is purely on synthetic datasets, we evaluate FoMo-0D on
 179 **57** real-world datasets from ADBench (Han et al., 2022)
 180 (see Table 19 in Appendix N). Following Livernoche et al.
 181 (2024), we use 5 train/test splits of each dataset via different
 182 seeds and report mean performance and standard deviation.

183 **Baselines:** We compare FoMo-0D against **26** baselines,
 184 from classical/shallow methods to modern/deep models,
 185 imported from one of the latest papers that proposed the
 186 SOTA diffusion-based OD model, DTE (Livernoche et al.,
 187 2024). We refer to the original paper for more details.

188 **Model Implementation:** We train FoMo-0D for 200,000
 189 steps with a batch size of 8 datasets (i.e., 1,600,000 synthet-
 190 ically generated datasets), which takes about 25 hours on 1
 191 GPU (Nvidia RTX A6000). Each dataset had a fixed size of
 192 10,000 samples, with $|\mathcal{D}_{\text{train}}| \in [n_L = 500, n_U = 5000]$,
 193 and the rest as $\mathcal{D}_{\text{test}}$ with balanced number of inliers and
 194 outliers. Other details (e.g., the training algorithm, model
 195 architecture, data synthesis) are in Appendix F.

196 **Metrics and Hypothesis Testing:** Detection performance
 197 is w.r.t. 3 widely-used metrics for OD: AUROC; area under
 198 ROC curve, AUPR; area under Precision-Recall curve, and
 199 F1 score; using threshold at the true number of outliers in
 200 the test data (varies by dataset) (Livernoche et al., 2024).

201 To compare different methods on ADBench, we compute
 202 their rank on each dataset (lower is better), and present
 203 average rank across datasets. This is an alternative to the
 204 average metric (e.g. AUROC), which is not meaningful
 205 when tasks vary widely in terms of their difficulties.

206 In addition, we perform significance tests to compare two
 207 methods statistically, using the one-sided paired Wilcoxon
 208 signed rank test (Demšar, 2006) between FoMo-0D and a
 209 baseline based on the performances across all datasets, with
 210 the alternative hypothesis suggesting the “baseline-minus-
 211 FoMo-0D” performance gap is greater than zero. We
 212 consider results significant at $\alpha = 0.05$ following convention.

213 **Hyperparameters (HPs):** Besides comparing FoMo-0D
 214 with the 26 baselines in Livernoche et al. (2024), respec-
 215 tively for AUROC, F1, and AUPR (Livernoche et al., 2024),

216 we also compare to the **top-4**³ best-performing baselines (in
 217 order: DTE-NP, kNN, ICL, and DTE-C) on their *average*
 218 performance across a list of different HP settings. Such
 219 an approach reflects their *expected* performance under HPs
 220 selected at random, in the absence of any other prior knowl-
 221 edge. We annotate the method name with ${}^{\text{avg}}$ for the version
 222 with performance averaged over varying HPs. The list of
 223 HP values for each top baseline is detailed in Appendix G.4.
 224 Overall, we compare FoMo-0D to 30 baselines; 26 from
 225 Livernoche et al. (2024) and ${}^{\text{avg}}$ of the top-4.

226 4.2. Results

227 **Detection performance:** Table 1 presented the compari-
 228 son of FoMo-0D w/ $D = 100$ to all baselines w.r.t. average
 229 rank across datasets as well as pairwise Wilcoxon signed
 230 rank tests based on AUROC, and we present full results on
 231 all datasets and all metrics in Appendix M. Among 30
 232 baselines and 2 variants of FoMo-0D (w/ $D = 100$ and
 233 $D = 20$), FoMo-0D w/ $D = 100$ performs as well as
 234 the 2nd best model (k NN with default HP; $k = 5$) on
 235 all datasets. While DTE-NP outperforms FoMo-0D with
 236 author-recommended $k = 5$, we find that DTE-NP ${}^{\text{avg}}$ is on
 237 par with FoMo-0D.

238 In our tests, $p > \alpha = 0.05$ implies no statistical evidence
 239 for performance difference between two methods. FoMo-0D
 240 w/ $D = 100$ performs statistically no different from **all**
 241 baselines on datasets with $d \leq 100$ (i.e., “at its own game”
 242 when pre-training data dimensions align with real-world
 243 datasets), while it outperforms the majority of baselines
 244 (where $p > 1 - \alpha$). These results continue to hold on
 245 datasets with $d \leq 500$. We present further results, ablation
 246 analyses, generalization analyses in Appendix I, J, K.

247 5. Conclusion

248 This work introduced **FoMo-0D**, the first foundation model
 249 for outlier detection (OD) on tabular data. It capitalizes on
 250 the in-context learning of a Transformer model pre-trained
 251 on a large number of synthetic datasets that can then per-
 252 form **zero-shot** inference on a new dataset, without *any* hyper-
 253 parameter tuning/training. FoMo-0D breaks new ground
 254 by fully abolishing the notoriously-hard model selection
 255 task for unsupervised OD (see Impact Statement). Further,
 256 FoMo-0D offers extremely fast inference thanks to a mere
 257 single forward pass. Against a long list of **26** SOTA base-
 258 lines on **57** public real-world datasets, FoMo-0D performs
 259 on par with the 2nd best baseline, while outperforming the
 260 majority of the baselines. Future work could expand our
 261 data prior and explore similar directions for zero-shot OD
 262 beyond tabular data. For a detailed discussion on limitations
 263 and future directions, we refer to Appendix P.

264 ³ To rank the 26 baselines, we compute the 26×26 p -values of
 265 the pairwise Wilcoxon signed rank test (see Appendix Figure 23),
 266 and order them by their mean p -value against other baselines.

References

- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Adriaensen, S., Rakotoarison, H., Müller, S., and Hutter, F. Efficient bayesian learning curve extrapolation using prior-data fitted networks. *Advances in Neural Information Processing Systems*, 36, 2024.
- Aggarwal, C. C. *Outlier Analysis*. Springer, 2013.
- Aggarwal, C. C. and Sathe, S. Theoretical foundations and algorithms for outlier ensembles. *Acm sigkdd explorations newsletter*, 17(1):24–47, 2015.
- Aggarwal, C. C. and Sathe, S. *Outlier Ensembles - An Introduction*. Springer, 2017.
- Akcay, S., Atapour-Abarghouei, A., and Breckon, T. P. Gandomaly: Semi-supervised anomaly detection via adversarial training. In *ACCV*, pp. 622–637. Springer, 2019.
- Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- Bergman, L. and Hoshen, Y. Classification-based anomaly detection for general data. In *International Conference on Learning Representations*, 2020.
- Breunig, M. M., Kriegel, H.-P., Ng, R. T., and Sander, J. Lof: Identifying density-based local outliers. In *International Conference on Management of Data*, 2000.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901, 2020.
- Campos, G. O., Zimek, A., Sander, J., Campello, R. J., Mincenková, B., Schubert, E., Assent, I., and Houle, M. E. On the evaluation of unsupervised outlier detection: measures, datasets, and an empirical study. *Data mining and knowledge discovery*, 30:891–927, 2016.
- Casella, G. and Berger, R. *Statistical inference*. CRC Press, 2024.
- Chalapathy, R. and Chawla, S. Deep learning for anomaly detection: A survey. *arXiv preprint arXiv:1901.03407*, 2019.
- Demšar, J. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine learning research*, 7:1–30, 2006.
- Ding, X., Zhao, L., and Akoglu, L. Hyperparameter sensitivity in deep outlier detection: Analysis and a scalable hyper-ensemble solution. *Advances in Neural Information Processing Systems*, 35:9603–9616, 2022.
- Ding, X., Zhao, Y., and Akoglu, L. Fast unsupervised deep outlier model selection with hypernetworks. *ACM SIGKDD*, 2024.
- Dolan, E. D. and Moré, J. J. Benchmarking optimization software with performance profiles. *Mathematical programming*, 91:201–213, 2002.
- Dooley, S., Khurana, G. S., Mohapatra, C., Naidu, S. V., and White, C. ForecastPFN: Synthetically-trained zero-shot forecasting. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Du, X., Sun, Y., Zhu, J., and Li, Y. Dream the impossible: Outlier imagination with diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Esmaeilpour, S., Liu, B., Robertson, E., and Shu, L. Zero-shot out-of-distribution detection based on the pre-trained model CLIP. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pp. 6568–6576, 2022.
- Feuer, B., Cohen, N., and Hegde, C. Scaling tabPFN: Sketching and feature selection for tabular prior-data fitted networks. In *NeurIPS 2023 Second Table Representation Learning Workshop*, 2023.
- Feuer, B., Schirrmeister, R. T., Cherepanova, V., Hegde, C., Hutter, F., Goldblum, M., Cohen, N., and White, C. Tunetables: Context optimization for scalable prior-data fitted networks, 2024.
- Garg, S., Tsipras, D., Liang, P. S., and Valiant, G. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 2022.
- Goldstein, M. and Dengel, A. Histogram-based outlier score (hbos): A fast unsupervised anomaly detection algorithm. *KI-2012: poster and demo track*, 1:59–63, 2012.
- Goldstein, M. and Uchida, S. A comparative evaluation of unsupervised anomaly detection algorithms for multivariate data. *PloS one*, 11(4), 2016.

- 275 Goyal, S., Raghunathan, A., Jain, M., Simhadri, H., and
 276 Jain, P. Drock: Deep robust one-class classification. In
 277 *Proceedings of the 37th International Conference on Ma-*
 278 *chine Learning*, ICML’20. JMLR.org, 2020.
- 279
- 280 Gu, Z., Zhu, B., Zhu, G., Chen, Y., Tang, M., and Wang,
 281 J. AnomalyGPT: Detecting industrial anomalies using
 282 large vision-language models. In *Proceedings of the*
 283 *AAAI Conference on Artificial Intelligence*, volume 38,
 284 pp. 1932–1940, 2024.
- 285
- 286 Gut, A. *An Intermediate Course in Probability*. Springer
 287 Texts in Statistics. Springer New York, 2009. ISBN
 288 9781441901620.
- 289
- 290 Han, S., Hu, X., Huang, H., Jiang, M., and Zhao, Y. Ad-
 291 bench: Anomaly detection benchmark. *Advances in Neu-*
 292 *ral Information Processing Systems*, 35, 2022.
- 293
- 294 He, H., Zhang, J., Chen, H., Chen, X., Li, Z., Chen, X.,
 295 Wang, Y., Wang, C., and Xie, L. A diffusion-based frame-
 296 work for multi-class anomaly detection. In *Proceedings*
 297 *of the AAAI Conference on Artificial Intelligence*, vol-
 298 ume 38, pp. 8472–8480, 2024.
- 299
- 300 He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learn-
 301 ing for image recognition. In *Proceedings of the IEEE*
 302 *conference on computer vision and pattern recognition*,
 303 pp. 770–778, 2016.
- 304
- 305 Ho, J., Jain, A., and Abbeel, P. Denoising diffusion proba-
 306 *bilistic models. Advances in Neural Information Process-
 307 ing Systems*, 33:6840–6851, 2020.
- 308
- 309 Hojjati, H., Ho, T. K. K., and Armanfard, N. Self-supervised
 310 anomaly detection: A survey and outlook. *arXiv preprint
 arXiv:2205.05173*, 2022.
- 311
- 312 Hollmann, N., Müller, S., Eggensperger, K., and Hutter, F.
 313 TabPFN: A transformer that solves small tabular classi-
 314 *fication problems in a second. In The Eleventh Interna-*
 315 *tional Conference on Learning Representations*, 2023.
- 316
- 317 Huber-Carol, C., Balakrishnan, N., Nikulin, M., and Mes-
 318 *bab, M. Goodness-of-fit tests and model validity*. Springer
 Science & Business Media, 2012.
- 319
- 320 Jeong, J., Zou, Y., Kim, T., Zhang, D., Ravichandran,
 321 A., and Dabeer, O. Winclip: Zero-/few-shot anomaly
 322 *classification and segmentation. In Proceedings of the*
 323 *IEEE/CVF Conference on Computer Vision and Pattern
 Recognition*, pp. 19606–19616, 2023.
- 324
- 325 Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B.,
 326 Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and
 327 Amodei, D. Scaling laws for neural language models.
 328 *arXiv preprint arXiv:2001.08361*, 2020.
- 329
- Kingma, D. P. Auto-encoding variational bayes. *arXiv
 preprint arXiv:1312.6114*, 2013.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic
 optimization, 2017.
- Li, A., Qiu, C., Kloft, M., Smyth, P., Rudolph, M., and
 Mandt, S. Zero-shot anomaly detection via batch normal-
 ization. In *Thirty-seventh Conference on Neural Infor-*
 329 *mation Processing Systems*, 2023.
- Li, A., Zhao, Y., Qiu, C., Kloft, M., Smyth, P., Rudolph, M.,
 329 and Mandt, S. Anomaly detection of tabular data using
 329 LLMs. *arXiv preprint arXiv:2406.16308*, 2024.
- Liu, F. T., Ting, K. M., and Zhou, Z.-H. Isolation forest.
 In *2008 Eighth IEEE International Conference on Data
 Mining*, pp. 413–422, 2008.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tun-
 329 ing. *Advances in neural information processing systems*,
 329 36, 2024.
- Livernoche, V., Jain, V., Hezaveh, Y., and Ravanbakhsh, S.
 On diffusion modeling for anomaly detection. In *ICLR*,
 329 2024.
- Liznerski, P., Ruff, L., Vandermeulen, R. A., Franks, B. J.,
 Müller, K.-R., and Kloft, M. Exposing outlier exposure:
 What can be learned from few, one, and zero outlier
 images. *arXiv preprint arXiv:2205.11474*, 2022.
- Ma, J., Thomas, V., Yu, G., and Caterini, A. L. In-context
 data distillation with tabPFN. *CoRR*, abs/2402.06971,
 2024.
- Ma, M. Q., Zhao, Y., Zhang, X., and Akoglu, L. The
 need for unsupervised outlier model selection: A review
 and evaluation of internal evaluation strategies. *ACM
 SIGKDD Explorations Newsletter*, 25(1):19–35, 2023.
- Müller, S., Feurer, M., Hollmann, N., and Hutter, F.
 PFNs4BO: in-context learning for bayesian optimization.
 In *International Conference on Machine Learning*, 2023.
- Müller, S., Hollmann, N., Pineda-Arango, S., Grabocka, J.,
 and Hutter, F. Transformers can do bayesian inference.
ICLR, 2022.
- Nagler, T. Statistical foundations of prior-data fitted net-
 works. In *ICML*, volume 202 of *Proceedings of Machine
 Learning Research*. PMLR, 2023.
- Pang, G., Shen, C., Cao, L., and Hengel, A. V. D. Deep
 learning for anomaly detection: A review. *ACM comput-*
 329 *ing surveys (CSUR)*, 54(2):1–38, 2021.

- 330 Paul, M., Ganguli, S., and Dziugaite, G. K. Deep learning on
 331 a data diet: Finding important examples early in training.
 332 *Advances in neural information processing systems*, 34:
 333 20596–20607, 2021.
- 334 Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G.,
 335 Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J.,
 336 et al. Learning transferable visual models from natural
 337 language supervision. In *International conference on machine learning*, 2021.
- 338 Ramaswamy, S., Rastogi, R., and Shim, K. Efficient algo-
 339 rithms for mining outliers from large data sets. In *ACM SIGMOD international conference on management of data*, 2000.
- 340 Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Rad-
 341 ford, A., Chen, M., and Sutskever, I. Zero-shot text-to-
 342 image generation. In *International conference on ma-
 343 chine learning*, pp. 8821–8831. Pmlr, 2021.
- 344 Raventós, A., Paul, M., Chen, F., and Ganguli, S. Pretrain-
 345 ing task diversity and the emergence of non-bayesian
 346 in-context learning for regression. *Advances in Neural
 347 Information Processing Systems*, 2024.
- 348 Rezende, D. and Mohamed, S. Variational inference with
 349 normalizing flows. In *Proceedings of the 32nd Interna-
 350 tional Conference on Machine Learning*, volume 37 of
 351 *Proceedings of Machine Learning Research*, pp. 1530–
 352 1538, Lille, France, 07–09 Jul 2015. PMLR.
- 353 Ruff, L., Vandermeulen, R., Goernitz, N., Deecke, L., Sid-
 354 diqui, S. A., Binder, A., Müller, E., and Kloft, M. Deep
 355 one-class classification. In *International Conference on
 356 Machine Learning*, pp. 4393–4402, 2018.
- 357 Ruff, L., Kauffmann, J. R., Vandermeulen, R. A., Montavon,
 358 G., Samek, W., Kloft, M., Dietterich, T. G., and Müller,
 359 K.-R. A unifying review of deep and shallow anomaly
 360 detection. *Proceedings of the IEEE*, 109(5):756–795,
 361 2021.
- 362 Shenkar, T. and Wolf, L. Anomaly detection for tabular
 363 data with internal contrastive learning. In *International
 364 Conference on Learning Representations*, 2022.
- 365 Sorscher, B., Geirhos, R., Shekhar, S., Ganguli, S., and Mor-
 366 cos, A. Beyond neural scaling laws: beating power law
 367 scaling via data pruning. *Advances in Neural Information
 368 Processing Systems*, 35:19523–19536, 2022.
- 369 Steinbuss, G. and Böhm, K. Benchmarking unsupervised
 370 outlier detection with realistic synthetic data. *ACM Trans-
 371 actions on Knowledge Discovery from Data (TKDD)*, 15
 372 (4):1–20, 2021.
- 373 Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux,
 374 M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro,
 375 E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and
 376 Lample, G. Llama: Open and efficient foundation lan-
 377 guage models, 2023. URL <https://arxiv.org/abs/2302.13971>.
- 378 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,
 379 L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention
 380 is all you need. In *NIPS*, pp. 5998–6008, 2017.
- 381 Xie, S. M., Raghunathan, A., Liang, P., and Ma, T. An
 382 explanation of in-context learning as implicit bayesian
 383 inference. *arXiv preprint arXiv:2111.02080*, 2021.
- 384 Xu, H., Wang, Y., Wei, J., Jian, S., Li, Y., and Liu, N.
 385 Fascinating supervisory signals and where to find them:
 386 Deep anomaly detection with scale learning. In *Inter-
 387 national Conference on Machine Learning*, pp. 38655–
 388 38673. PMLR, 2023.
- 389 Yoo, J., Zhao, T., and Akoglu, L. Data augmentation is
 390 a hyperparameter: Cherry-picked self-supervision for
 391 unsupervised anomaly detection is creating the illusion
 392 of success. *Trans. Mach. Learn. Res.*, 2023, 2023.
- 393 Yoon, S., Jin, Y.-U., Noh, Y.-K., and Park, F. Energy-based
 394 models for anomaly detection: A manifold diffusion re-
 395 covery approach. *Advances in Neural Information Pro-
 396 cessing Systems*, 36:49445–49466, 2023.
- 397 Zhai, X., Kolesnikov, A., Houlsby, N., and Beyer, L. Scaling
 398 vision transformers. In *Proceedings of the IEEE/CVF
 399 conference on computer vision and pattern recognition*,
 400 2022.
- 401 Zhang, J., Yang, J., Wang, P., Wang, H., Lin, Y., Zhang, H.,
 402 Sun, Y., Du, X., Zhou, K., Zhang, W., et al. Openood v1.
 403 5: Enhanced benchmark for out-of-distribution detection.
 404 *arXiv preprint arXiv:2306.09301*, 2023.
- 405 Zhang, Y. and Yan, J. Crossformer: Transformer utilizing
 406 cross-dimension dependency for multivariate time series
 407 forecasting. In *ICLR*, 2023.
- 408 Zhao, Y. and Akoglu, L. Toward unsupervised outlier model
 409 selection. In *International Conference on Automated
 410 Machine Learning (AutoML)*, 2024.
- 411 Zhao, Y., Rossi, R., and Akoglu, L. Automatic unsupervised
 412 outlier model selection. *Advances in Neural Information
 413 Processing Systems*, 34:4489–4502, 2021.
- 414 Zhao, Y., Zhang, S., and Akoglu, L. Toward unsupervised
 415 outlier model selection. In *2022 IEEE International Con-
 416 ference on Data Mining (ICDM)*, pp. 773–782. IEEE,
 417 2022.

385 Zhou, Q., Pang, G., Tian, Y., He, S., and Chen, J. Anoma-
386 lyCLIP: Object-agnostic prompt learning for zero-shot
387 anomaly detection. In *The Twelfth International Confer-*
388 *ence on Learning Representations*, 2024. URL <https://openreview.net/forum?id=buC4E91xZE>.

389
390 Zong, B., Song, Q., Min, M. R., Cheng, W., Lumezanu, C.,
391 ki Cho, D., and Chen, H. Deep autoencoding gaussian
392 mixture model for unsupervised anomaly detection. In
393 *International Conference on Learning Representations*,
394 2018.
395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440 **Table of Contents**441 We detail the contents in the appendix below.
442

- 443
444 • **Appendix A. Related Work** Related literature on outlier detection (OD), unsupervised model selection for OD,
445 prior-data fitted networks, and zero-shot outlier detection.
- 446 • **Appendix B. Illustration of Synthetic Sata in 2-d** illustrates the inlier and outlier data synthesis for pre-training with
447 a 2-dimensional example.
- 448 • **Appendix C. Background on Prior-data Fitted Networks** introduces prior-data fitted networks.
- 449 • **Appendix D. Architecture and Scalability** details the model architecture and methods to scale up pre-training.
- 450 • **Appendix E. Linear Transform for Scalable GMM Data Synthesis** contains the proofs for efficient data synthesis in
451 Appendix D.
- 452 • **Appendix F. Implementation Details** includes the training and inference details of FoMo-0D.
- 453 • **Appendix G. Detailed Experiment Setup** introduces the details of pre-training and inference datasets, baselines, and
454 their hyperparameters.
- 455 • **Appendix H. Qualitative Analysis on Sample-to-Sample Attention** visualizes the attention of FoMo-0D.
- 456 • **Appendix I. Further Results** presents further detection performance and running time results of FoMo-0D.
- 457 • **Appendix J. Ablation Analyses** studies different design choices of FoMo-0D.
- 458 • **Appendix K. Generalization Analyses** studies the generalization ability of FoMo-0D on out-of-distribution synthetic
459 datasets, ADBench, and benchmarks.
- 460 • **Appendix L. Performance Profile Plots** presents a comprehensive comparison of different methods via the cumulative
461 distribution.
- 462 • **Appendix M. Full Results** presents the detailed metric results of FoMo-0D and the baselines, including AUROC,
463 AUPRC, and F1.
- 464 • **Appendix N. Benchmark OD Datasets** shows the details (e.g., number of samples, features) of each dataset in
465 ADBench.
- 466 • **Appendix O. Differences to Prior Work on PFNs for Tabular Data** explains the difference and innovation of
467 FoMo-0D from previous works.
- 468 • **Appendix P. Discussion** provides the summary of our work and discussions on the limitations and future directions of
469 FoMo-0D.
- 470 • **Appendix Q. Reproducibility Statement** details the codebase for FoMo-0D.
- 471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494

495 **A. Related Work**

496 **Outlier Detection (OD):** Thanks to diverse applications in numerous fields, such as security, finance, manufacturing, to
 497 name a few, OD on tabular (or point-cloud) datasets has a vast literature with a long list of techniques. For earlier, shallow
 498 approaches preceding the advances in deep learning, we refer to the books by Aggarwal (2013) and Aggarwal & Sathe
 499 (2017). The modern, deep learning based techniques are surveyed in (Chalapathy & Chawla, 2019; Pang et al., 2021; Ruff
 500 et al., 2021). Most recent deep OD techniques take advantage of newly emerging paradigms, including self-supervised
 501 learning (Hojjati et al., 2022; Yoo et al., 2023) as well as the most recently popularized diffusion-based models (Yoon et al.,
 502 2023; Livernoche et al., 2024; Du et al., 2024; He et al., 2024).

503 **Unsupervised Model Selection for OD:** It is typical of models to exhibit various hyperparameters (HPs) that play a role
 504 in the bias-variance trade-off and hence the generalization performance, and OD models are no exception. Many earlier
 505 work on OD showed the sensitivity of classical (i.e. shallow) OD methods to the choice of their HP(s) (Aggarwal & Sathe,
 506 2015; Campos et al., 2016; Goldstein & Uchida, 2016). Similarly, sensitivity to HPs has also been shown for deep OD
 507 models more recently (Zhao et al., 2021; Ding et al., 2022), as well as for those relying on self-supervised learning/data
 508 augmentation (Yoo et al., 2023).

509 While critical, work on unsupervised outlier model selection (UOMS) is slim as compared to the vast literature on detection
 510 methods. A handful of existing, mostly heuristic strategies has been studied by Ma et al. (2023) reporting discouraging
 511 results; they have shown that existing heuristics are either not significantly different from random selection, or do not
 512 outperform iForest (Liu et al., 2008) with its default HPs (an extremely fast ensemble of randomized trees).

513 More recent UOMS approaches go beyond heuristic measures and instead design scalable hyperensembles (Ding et al.,
 514 2022; 2024), as well as take advantage of meta-learning on historical real-world OD datasets (Zhao et al., 2021; 2022; Zhao
 515 & Akoglu, 2024). These approaches demonstrate the value of learning from many other OD datasets, and transfer these
 516 learnings to a new dataset. While sharing the same spirit on learning from a large collection of (in our case, simulated)
 517 datasets, our FoMo-OD differs from these prior art in a key aspect; FoMo-OD is *not* a model selection technique, but rather, a
 518 foundation model that abolishes model training and selection altogether—unlocking 0-shot inference on a new task.

519 **Prior-data Fitted Networks:** Based on the seminal work by Müller et al. (2022), Prior-data-fitted Networks (PFNs) establish
 520 a new paradigm for machine learning, where a PFN is pretrained on synthetic datasets generated from a data prior, and the
 521 pretrained PFN can then infer the posterior predictive distribution (PPD) for test points in a new dataset in a single forward
 522 pass, through in-context learning (Xie et al., 2021; Garg et al., 2022). It is shown that PFNs provably approximate Bayesian
 523 inference (Müller et al., 2022). Follow-up TabPFN (Hollmann et al., 2023) achieved SOTA classification performance
 524 on small tabular datasets of size up to 1024. Other subsequent works designed LC-PFN (Adriaensen et al., 2024) and
 525 ForecastPFN (Dooley et al., 2023), respectively zero-shot learning curve extrapolation and zero-shot time-series forecasting
 526 models, trained purely on synthetic data. PFN4BO (Müller et al., 2023) employed PFNs for Bayesian optimization, while
 527 Nagler (2023) studied the statistical foundations of PFNs. As training data is passed as context to PFN, others proposed
 528 scaling solutions to enable training on larger pretraining datasets for better generalization (Ma et al., 2024; Feuer et al.,
 529 2023; 2024).

530 Our proposed FoMo-OD differs from these in being the first PFN for OD, using a novel inlier/outlier data prior, employing
 531 linear transform for fast data synthesis, and incorporating the “router” attention mechanism for linear-time scalability w.r.t.
 532 context size. See Appendix O for additional details.

533 **Zero-Shot Outlier Detection:** Foundation models pretrained on massive text and image corpora, such as large language
 534 and/or vision models (L(V)LMs) like OpenAI’s GPT-series (Achiam et al., 2023), DALL-E (Ramesh et al., 2021) and
 535 Flamingo (Alayrac et al., 2022), CLIP (Radford et al., 2021), and LLaVA (Liu et al., 2024) to name a few, have demonstrated
 536 remarkable success on several zero-shot tasks in CV and NLP. Follow-up work extended these models for zero-shot
 537 out-of-distribution detection (Esmaeilpour et al., 2022), zero-shot image OD (Liznerski et al., 2022; Jeong et al., 2023; Zhou
 538 et al., 2024) as well as dialogue-based industrial image anomaly detection (Gu et al., 2024).

539 Foundation models, however, do not exist for tabular data which is widespread across OD applications in the real world,
 540 such as detecting credit card fraud, network intrusion, medical anomalies, and any sensor measurement abnormalities, to
 541 name a few. The recent ACR model by Li et al. (2023) on zero-shot OD does *not* rely on a pretrained foundation model, but
 542 rather is meta-trained on each specific domain using inlier-only datasets from the *same domain*. Concurrent to our work, Li
 543 et al. (2024) apply pretrained LLMs for prompt-based OD on tabular data which they serialize to text. Similar to our work,
 544 they also use *simulated* labeled OD datasets to fine-tune several existing LLMs to improve their performance. Their work,
 545

however, is quite preliminary in several fronts; a key limitation is that they assume independent features and query the LLM one-feature-at-a-time to reach an outlier score. Further, they fine-tune using only 5,000 data batches with up to 100 samples each, subsample 150 points and the first 10 columns of each dataset for evaluation (due to GPU memory constraint), and their testbed includes only two baseline methods. In contrast, FoMo-0D employs and pretrains PFNs at a much larger scale with rigorous evaluation on a much larger testbed.

555

B. Illustration of synthetic data in 2-d

We visualize our synthetic data in Figure 3, with 3 randomly created 2-d GMMs with the number of clusters ($N = 1, 2, 3$). We choose the 80th percentile as the criterion, such that inliers are samples drawn from the GMM and within the 80th percentile, and outliers are samples drawn from the inflated GMMs and outside of the 80th percentile.

561

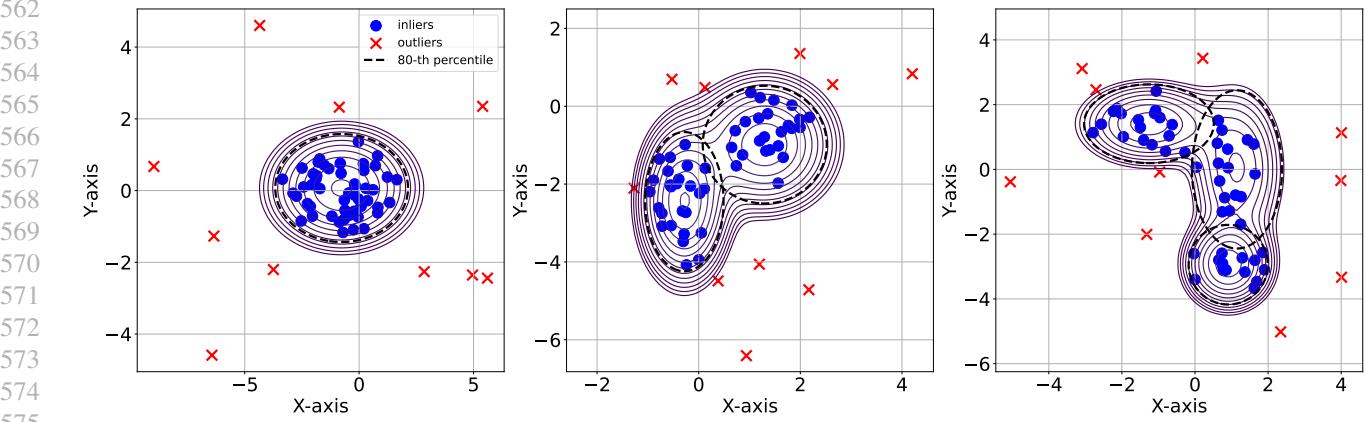


Figure 3: Illustration of synthetic data in 2D with 80th percentile as the criterion.

C. Background on Prior-data Fitted Networks

Posterior Predictive Distribution (PPD): In the Bayesian framework for supervised learning, the prior defines a hypothesis space Φ which expresses our beliefs about the data distribution before seeing any data. Each hypothesis $\phi \in \Phi$ describes a mechanism by which the data is generated. The posterior predictive distribution $p(\cdot | \mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})$ provides a framework for making prediction on new, unseen test data \mathbf{x}_{test} , conditioned on observed training data $\mathcal{D}_{\text{train}} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$. Based on Bayes' Theorem, the PPD can be derived by the integration over the space of hypotheses Φ :

$$p(y_{\text{test}} | \mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}}) = \int_{\Phi} p(y_{\text{test}} | \mathbf{x}_{\text{test}}, \phi) p(\mathcal{D}_{\text{train}} | \phi) p(\phi) d\phi, \quad (2)$$

where $p(\phi)$ denotes the prior probability and $p(\mathcal{D} | \phi)$ is the likelihood of the data \mathcal{D} given ϕ .

PFNs and PPD Approximation: As obtaining the above PPD is generally intractable, Prior-data Fitted Networks (PFNs) are proposed to approximate the PPD (Müller et al., 2022). Unlike traditional machine learning models that are trained directly on observed datasets, PFNs are pre-trained on simulated datasets that are generated according to a prior distribution. Specifically, it contains the pre-training and inference stages described as the following.

Pre-training on synthetic data. Massive synthetic datasets are generated for the pre-training stage, by first sampling a hypothesis (i.e., the generating mechanism) $\phi \sim p(\phi)$, and then sampling a dataset $\mathcal{D} \sim p(\mathcal{D} | \phi)$. For training, each dataset \mathcal{D} can be split as $\mathcal{D}_{\text{test}} \subset \mathcal{D}$ and $\mathcal{D}_{\text{train}} = \mathcal{D} \setminus \mathcal{D}_{\text{test}}$. Thus, the PFN with parameters θ can be optimized by making predictions on data points in D_{test} . For a test point $(\mathbf{x}_{\text{test}}, y_{\text{test}}) \in \mathcal{D}_{\text{test}}$, the training loss is as follows.

$$\mathcal{L} = \mathbb{E}_{\{(\mathbf{x}_{\text{test}}, y_{\text{test}})\} \cup \mathcal{D}_{\text{train}} \sim p(\mathcal{D})} [-\log q_{\theta}(y_{\text{test}} | \mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})]. \quad (3)$$

The above loss can also be interpreted as minimizing the expected KL divergence between $p(\cdot | \mathbf{x}, \mathcal{D})$ and $q_{\theta}(\cdot | \mathbf{x}, \mathcal{D})$ (Müller et al., 2022). In practice, a PFN model q_{θ} is typically implemented by a Transformer-based architecture (Vaswani et al.,

605 which takes $(\mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})$ as input, where $\mathbf{x}_{\text{test}} \in \mathcal{D}_{\text{test}}$ and $\mathcal{D}_{\text{train}}$ contains an arbitrary number of instances. The
606 output is the conditional class probabilities for \mathbf{x}_{test} . As the whole training set $\mathcal{D}_{\text{train}}$ is passed as input/context to the
607 Transformer, it learns to predict class labels through sample-to-sample attention.

608 *Inference on real-world data.* In the inference stage, a fresh real-world dataset $\mathcal{D}_{\text{train}}$ and some test instance \mathbf{x}_{test} are
609 fed into the (frozen) pre-trained model, which computes the PPD $q_{\theta}(\cdot | \mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})$ in a single forward pass. Importantly,
610 PFNs do not require gradient-based parameter tuning on new datasets, where prediction is delivered *in less than a second*
611 ([Hollmann et al., 2023](#)).

612 In summary, PFNs are trained once, and can be used many times for zero-shot inference on new datasets with different
613 characteristics. The main benefit is that **no training or tuning** is required at the inference stage. This type of learning ability
614 is also termed as in-context learning (ICL) ([Xie et al., 2021](#)), which was shown to be effective for various tasks in languages
615 ([Brown et al., 2020](#)). In fact, ICL with PFNs is recently shown to be a promising paradigm for supervised classification on
616 tabular datasets ([Hollmann et al., 2023](#)).

617

D. Architecture and Scalability

618 **Architecture and sample-to-sample attention:** Like existing PFNs, FoMo-0D is based on the Transformer ([Vaswani et al., 2017](#)), encoding each sample’s feature vector as a token, and allowing token representations to attend to each other, hence enabling sample-to-sample attention. We also adopt the three customizations from TabPFN ([Hollmann et al., 2023](#)), which (1) computes self-attention among all the training samples and only cross-attention from test samples to the training samples, (2) enables varying feature dimensionality by zero-padding, and (3) randomly permutes input samples while omitting positional encodings to achieve model invariance in the dataset.

619 Given $\mathcal{D}_{\text{train}} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, each self-attention layer outputs n embeddings $\{\mathbf{z}_i\}_{i=1}^n$; where the i -th token is mapped via
620 linear transformations to a key \mathbf{k}_i , query \mathbf{q}_i and value \mathbf{v}_i , where the i -th output is computed as

$$621 \quad \mathbf{z}_i = \sum_{j=1}^n \text{softmax}(\{\langle \mathbf{q}_i, \mathbf{k}_{j'} \rangle\}_{j'=1}^n)_j \cdot \mathbf{v}_j. \quad (4)$$

622 The sample-to-sample attention is intriguing from the perspective of OD: many classical OD algorithms ([Aggarwal, 2013](#))
623 are based on nonparametrics; in particular, they leverage the distances to the k nearest neighbors (k NNs) of a point to
624 compute its outlierness, where k is a critical hyperparameter. One can think of FoMo-0D as mimicking non-parametric
625 models but by using parametric attention mechanisms. Interestingly, PFNs are much more robust and flexible than k NN
626 based OD approaches, for (1) sample-to-sample relations are not pre-specified but rather learned through attention weights,
627 and thus (2) they are not limited to just the nearest neighbors but rather can *learn which* training points are worth attending
628 to, and (3) as attention is dataset-wide across all points, there is no need for specifying a cut-off HP value like k , to which
629 most k NN based OD techniques are sensitive to ([Aggarwal & Sathe, 2015; Campos et al., 2016; Goldstein & Uchida, 2016;](#)
630 [Ding et al., 2022](#)). We present analyses on sample-to-sample attention in Appendix H.

631 To seize the power of scale, we incorporate a scalable architecture and data synthesis into our design to benefit pre-training
632 and inference, as we describe next. The scale-up unlocks a larger context size for FoMo-0D, enabling pre-training and
633 inference on larger datasets with fast speed.

634 **Scaling up attention with “routers”:** The $\mathcal{O}(n^2)$ quadratic sample complexity at pre-training presents an obstacle for
635 achieving high performance at inference, as it limits pre-training to relatively small training datasets, and degenerates
636 in-context learning that typically benefits from longer context ([Xie et al., 2021](#)).

637 Toward a high-performance model, we scale up FoMo-0D’s attention via the “router mechanism” of [Zhang & Yan \(2023\)](#).
638 As shown in Figure 2, the main idea is to learn a small number ($R \ll n$) of “routers” or representatives, which gather
639 information from all n samples and then distribute the information back to the n output embeddings—reducing complexity
640 from $\mathcal{O}(n^2)$ to $\mathcal{O}(2Rn) = \mathcal{O}(n)$. This design allows FoMo-0D to **scale linearly** with respect to both dimensionality d and
641 dataset size n in pre-training.

642 Concretely, the representatives first aggregate information from all samples by serving as the query in the multi-head
643 self-attention (MSA);

$$644 \quad \mathcal{M} = \text{MSA}_1(\mathbf{R}, \mathbf{Z}, \mathbf{Z}), \quad (5)$$

645 where $\mathbf{R} \in \mathbb{R}^{R \times d}$ depicts the *learnable* vector array of representatives and \mathcal{M} denotes the aggregated messages. Then, the
646 routers distribute the received information among samples by using the sample embeddings as query and the aggregated
647 messages as key and value.

660 messages as both key and value:

$$\hat{\mathbf{Z}} = \text{MSA}_2(\mathbf{Z}, \mathcal{M}, \mathcal{M}) . \quad (6)$$

661 Finally, we obtain $\bar{\mathbf{Z}} = \text{LayerNorm}(\hat{\mathbf{Z}} + \mathbf{Z})$ after layer normalization. Note that the test samples only attend to the
662 training samples' embeddings, computed in the described manner across layers, and are finally fed into the prediction head
663 to estimate the PPD at the output layer.
664

665 **Scaling up (pre)training data synthesis with linear transforms:** Besides the scalability challenge associated with
666 architecture/attention, another computational challenge in pre-training FoMo-0D arises from drawing samples from the data
667 prior, which requires considerable time, especially in high dimensions⁴, provided the large number of datasets we sample
668 (specifically, we utilize a batch size of 8 datasets over 1,000 steps each for 200 epochs).

669 To give an idea, sampling a dataset with $n = 10,000$ points in $d = 100$ dimensions using 10 CPUs in parallel takes ≈ 0.4
670 seconds (see Appendix Figure 6). Across 200 training epochs with 1,000 steps each, it adds up to more than 177 hours
671 just to generate 1,6 million datasets on-the-fly. Of course, one can trade storage with compute-time by generating all these
672 datasets apriori via massive parallelism. Nevertheless, synthetic data generation demands considerable time (and/or storage).
673

674 To scale up data synthesis, FoMo-0D employs two distinct strategies. **First**, we propose *reuse at epoch level*: that is, one can
675 reuse the same 8K (8×1000) unique datasets at every epoch, or in general, the same $8K \times P$ datasets periodically at every
676 P epochs. A larger P would lead to more diversity in terms of the overall pre-training data used.
677

678 **Second**, we propose *reuse at dataset level via transformation*: that is, having generated one unique dataset $\mathbf{X} \in \mathbb{R}^{n \times d}$ from
679 a GMM, we propose a linear transform $T(\mathbf{x})$ of the form $\mathbf{W}\mathbf{x} + \mathbf{b}$ for randomly drawn parameters $\mathbf{W} \in \mathbb{R}^{d \times d}$ and $\mathbf{b} \in \mathbb{R}^d$
680 (see Appendix E.1).⁵ This simple yet efficient transformation creates a new dataset, akin to one being drawn from another
681 GMM with centers $T(\boldsymbol{\mu}_j) = \mathbf{W}\boldsymbol{\mu}_j + \mathbf{b}$ and covariance $T(\boldsymbol{\Sigma}_j) = \mathbf{W}\boldsymbol{\Sigma}_j\mathbf{W}^T$, $\forall j \in [m]$. Note that we do not actually
682 materialize these parameters but only transform the dataset. As we show in the following, such transformations preserve the
683 Mahalanobis distances as well as the percentile thresholds for labeling points as inlier/outlier. Details and proofs are given
684 in Appendix E.
685

686 **Lemma D.1.** *Linear transform T with invertible \mathbf{W} on \mathcal{G}_m^d preserves Mahalanobis distances.*

687 **Lemma D.2.** *Linear transform T with invertible \mathbf{W} on \mathcal{G}_m^d preserves the percentiles of the GMM.*

688 The implication of these lemmas is that a linear transformation of a dataset from a GMM retains the identity of the inliers
689 and outliers, i.e. no relabeling is required. Moreover, notice that as a byproduct we obtain a transformed dataset as though
690 it is drawn from a GMM with a *non-diagonal* covariance matrix which, besides the time savings, offers a slightly more
691 complex data prior.
692

693 To reach 8K unique datasets for each epoch, we first generate 500 datasets from different GMMs (with varying configurations),
694 then employ 15 different linear transformations to each dataset by varying \mathbf{W} and \mathbf{b} . Drawing each (\mathbf{W}, \mathbf{b}) takes ≈ 0.02
695 seconds, while the matrix-matrix product of \mathbf{X} ($n \times d$) and \mathbf{W} ($d \times d$) takes negligible time (for $d \leq 100$). Thus, obtaining
696 a transformed dataset offers $20 \times$ speed-up compared to generating one (0.02 vs. 0.4 seconds).
697

E. Linear Transform for Scalable GMM Data Synthesis

E.1. Definitions

701 **Definition E.1** (Gaussian Mixture Model). We denote an m -cluster d -dimension Gaussian Mixture Model as $\mathcal{G}_m^d =$
702 $\{(w_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)\}_{j=1}^m$, which is the weighted sum of m Gaussian distributions:
703

$$p(\mathbf{x}) = \sum_{j=1}^m w_j \cdot g(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) , \quad (7)$$

707 where $w_j \in \mathbb{R}^+$ is the weight for the j -th Gaussian $\mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ with $\sum_{j=1}^m w_j = 1$, and $g(\cdot | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ is the density of the
708 j -th component/cluster, with mean/center $\boldsymbol{\mu}_j \in \mathbb{R}^d$ and covariance $\boldsymbol{\Sigma}_j \in \mathbb{R}^{d \times d}$ being positive semi-definite, such that
709

710 ⁴This is because the inverse of the $(d \times d)$ covariance matrix plays a crucial role in the process of drawing samples from GMMs,
711 which has $\mathcal{O}(d^3)$ time complexity. (It is also the reason why diagonal $\boldsymbol{\Sigma}_j$'s are favored in our data prior.) In addition, Mahalanobis
712 distance for labeling inliers/outliers also requires the inverse.

713 ⁵In practice, we apply the linear transform on the subspace of inflated features only, wherein inliers and outliers are defined, which
714 remains to be a multi-variate GMM.

715 $\mathbf{x}^T \boldsymbol{\Sigma}_i \mathbf{x} \geq 0$, for all $\mathbf{x} \in \mathbb{R}^d$.

716 **Definition E.2** (Linear Transform). We denote a linear transformation T in \mathbb{R}^d as:

$$718 \quad T(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}, \quad (8)$$

720 where $\mathbf{x} \in \mathbb{R}^d$, and $\mathbf{W} \in \mathbb{R}^{d \times d}$, $\mathbf{b} \in \mathbb{R}^d$ are the parameters of T .

721 **Definition E.3** (Mahalanobis Distance). The Mahalanobis distance dist_M between a point $\mathbf{x} \in \mathbb{R}^d$ and a Gaussian
722 distribution $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is defined as:

$$724 \quad \text{dist}_M(\mathbf{x}) = \sqrt{(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}. \quad (9)$$

726 **Definition E.4** (χ_d^2 -distribution). The Chi-squared distribution χ_d^2 with d degrees of freedom is the distribution of the sum
727 of squares of d independent standard Normal random variables.

E.2. Properties

731 **Property E.5** (Lemma 5.3.2 (Casella & Berger, 2024)). If $Z \sim \mathcal{N}(0, 1)$, then $Z^2 \sim \chi_1^2$; If X_1, \dots, X_d are independent and
732 $X_i \sim \chi_1^2$, then $\sum_{i=1}^d X_i \sim \chi_d^2$.

733 **Property E.6.** The squared Mahalanobis distance $\text{dist}_M^2(\mathbf{x}) \sim \chi_d^2$, with $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

735 *Proof:* If $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then we have $\mathbf{z} = \boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbf{x} - \boldsymbol{\mu}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ (Gut, 2009), such that:

$$737 \quad \text{dist}_M^2(\mathbf{x}) = \mathbf{z}^T \mathbf{z} = \sum_{i=1}^d z_i^2 \quad (10)$$

740 where z_i are independent standard Normal random variables. We have $\sum_{i=1}^d z_i^2 \sim \chi_d^2$ from Property E.5, which completes
741 the proof.

E.3. Lemmas

745 **Lemma E.7.** Linear transform T with invertible \mathbf{W} on \mathcal{G}_m^d preserves Mahalanobis distances.

747 *Proof:* We denote the transformed GMM as $T(\mathcal{G}_m^d) = \{(w_j, \mathbf{W}\boldsymbol{\mu}_j + \mathbf{b}, \mathbf{W}\boldsymbol{\Sigma}_j\mathbf{W}^T)\}_{j=1}^m$, then with $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$, for
748 the transformed point $T(\mathbf{x})$ we have:

$$749 \quad \text{dist}_M(T(\mathbf{x})) = \sqrt{(T(\mathbf{x}) - (\mathbf{W}\boldsymbol{\mu}_j + \mathbf{b}))^T (\mathbf{W}\boldsymbol{\Sigma}_j\mathbf{W}^T)^{-1} (T(\mathbf{x}) - (\mathbf{W}\boldsymbol{\mu}_j + \mathbf{b}))} \quad (11)$$

$$752 \quad = \sqrt{(\mathbf{W}(\mathbf{x} - \boldsymbol{\mu}_j))^T (\mathbf{W}\boldsymbol{\Sigma}_j\mathbf{W}^T)^{-1} (\mathbf{W}(\mathbf{x} - \boldsymbol{\mu}_j))} \quad (12)$$

$$754 \quad = \sqrt{(\mathbf{x} - \boldsymbol{\mu}_j)^T \mathbf{W}^T (\mathbf{W}^T)^{-1} \boldsymbol{\Sigma}_j^{-1} \mathbf{W} (\mathbf{x} - \boldsymbol{\mu}_j)} \quad (13)$$

$$756 \quad = \sqrt{(\mathbf{x} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j)} = \text{dist}_M(\mathbf{x}). \quad (14)$$

757 \square

759 **Lemma E.8.** Linear transform T with invertible \mathbf{W} on \mathcal{G}_m^d preserves the percentiles of the GMM.

760 *Proof:* Let $\chi_d^2(\alpha)$ denote the α -th percentile of χ_d^2 , such that for $X \sim \chi_d^2$:

$$762 \quad \text{Prob}(X \leq \chi_d^2(\alpha)) = \frac{\alpha}{100}. \quad (15)$$

764 Based on Property E.6, we have $\text{Prob}(\text{dist}_M^2(\mathbf{x}) \leq \chi_d^2(\alpha)) = \frac{\alpha}{100}$.

766 Let $\mathbf{x} \sim \mathcal{G}_m^d$, such that $\text{dist}_M^2(\mathbf{x}) > \chi_d^2(\alpha)$ for all $\mathcal{N}_j(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$, which indicates that \mathbf{x} is outside the α -th percentile of \mathcal{G}_m^d .
767 Since $\text{dist}_M(\mathbf{x})$ is preserved under T (see Lemma E.7), then we conclude that the linear transform T with invertible \mathbf{W}
768 preserves the percentiles of the GMM. \square

770 **Algorithm 1** Prior-fitting of a PFN (Müller et al., 2022) and ours

771 **input** A prior distribution over datasets $p(\mathcal{D})$, from which samples can be drawn and the number of datasets Q to draw for one epoch,
772 the number of training epochs E , the periodicity P , the number of unique datasets q , linear transformation T .

773 **output** A model q_θ that will approximate the PPD

774 1: Initialize the neural network q_θ .

775 2: Initialize the epoch-level collection $\mathcal{C}_E = []$.

776 3: **for** $i \leftarrow 1$ to E **do**

777 4: **if** $i \leq P$ **then**

778 5: Initialize an empty buffer $\mathcal{B}_i = []$.

779 6: Initialize the dataset-level collection $\mathcal{C}_q = []$.

780 7: **for** $j \leftarrow 1$ to Q **do**

781 8: **if** $j \leq q$ **then**

782 9: Step 1: sample $D_j := \mathcal{D}_{\text{train}} \cup \{(\mathbf{x}_k, y_k)\}_{k=1}^{|\mathcal{D}_{\text{test}}|} \sim p(\mathcal{D})$.

783 10: $\mathcal{C}_q \leftarrow \mathcal{C}_q + [D_j]$

784 11: **else**

785 12: $j \leftarrow j \bmod q$

786 13: $D_j \leftarrow T(\mathcal{C}_q[j])$

787 14: **end if**

788 15: Step 2: compute stochastic loss approximation $\bar{\ell}_\theta = \sum_{k=1}^{|\mathcal{D}_{\text{test}}|} (-\log q_\theta(y_k | \mathbf{x}_k, \mathcal{D}_{\text{train}}))$.

789 16: Step 3: update parameters θ with stochastic gradient descent on $\nabla_\theta \bar{\ell}_\theta$.

790 17: $\mathcal{B}_i \leftarrow \mathcal{B}_i + [D_j]$

791 18: **end for**

792 19: $\mathcal{C}_E \leftarrow \mathcal{C}_E + [\mathcal{B}_i]$

793 20: **else**

794 21: $i \leftarrow i \bmod P$

795 22: $\mathcal{B}_i \leftarrow \mathcal{C}_E[i]$

796 23: **for** $j \leftarrow 1$ to Q **do**

797 24: $D_j \leftarrow T(\mathcal{B}_i[j])$

798 25: Perform Step 2 and Step 3

799 26: **end for**

800 27: **end if**

801 28: **end for**

F. Implementation details

F.1. Hardware

We base our experiments on a NVIDIA RTX A6000 GPU with AMD EPYC 7742 64-Core Processors.

F.2. Training and inference

We train our models for 200 epochs with the Adam optimizer (Kingma & Ba, 2017) and a learning_rate = 0.001, and test with the model corresponding to the lowest training loss. The size of our $D = \{20, 100\}$ model is 4.87M and 4.89M parameters, respectively. We show the training process of PFNs and our model in Algorithm 1.

Dealing with varying dimensions and dataset size For an input with d features, we follow Müller et al. (2022) and deal with $d < D$ by rescaling the input with $\frac{D}{d}$ and padding the features to size D with 0, and randomly sample D features out of d if $d > D$. In addition, FoMo-0D uses context size of 5K at inference, where we randomly sample (5K–1) points as $\mathcal{D}_{\text{train}}$ from datasets with $n > 5K$ for each test sample $\mathbf{x} \in \mathcal{D}_{\text{test}}$.

Model architecture We use a 4-layer Transformer with hidden dimension $\text{h_dim} = 256$, a linear layer ($\mathbb{R}^D \rightarrow \mathbb{R}^{\text{h_dim}}$) as the embedding layer and a 2-layer MLP ($\mathbb{R}^{\text{h_dim}} \rightarrow \mathbb{R}^2$) as the classification layer for inlier vs. outlier. For each Transformer layer, we use num_head = 4 for each attention module and $R = 500$ for the router-based attention (Figure 2).

Training loss In Figure 4, we plot the training loss of our $D = 100$ model trained with 8K unique datasets/epoch (denoted as “8K”) versus 0.5K unique + 7.5K transformed datasets/epoch (denoted as “0.5K+T”), together with the $D = 20$ model trained with reuse periodicity $P = 1$ (denoted as “P=1”, reusing the same 8K datasets across epochs) and $P = 1$ with transformation (denoted as “P=1+T”, transforming the 8K datasets across epochs). Notice that the loss with transformation is slightly higher than no transformation (i.e., $D = 100$, “0.5K+T” vs. “8K”, and $D = 20$, “P=1+T” vs. “P=1”) across all 200 epochs, which is reasonable since the transformed datasets have non-diagonal covariances that make the learning task

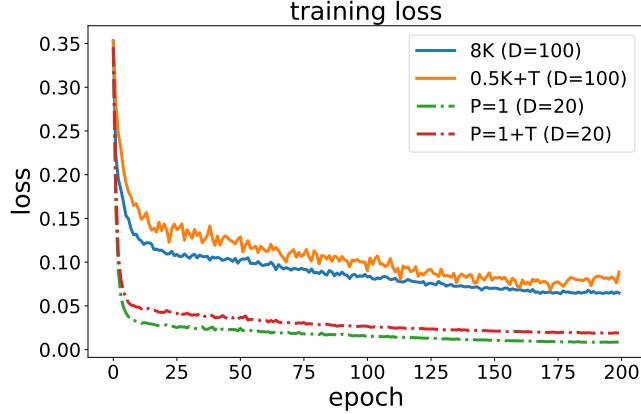


Figure 4: (best in color) Training loss of FoMo-OD ($D = 100$) with 8K unique datasets/epoch (in blue) and using 0.5K unique + 7.5K transformed datasets/epoch (in orange), and FoMo-OD ($D = 20$) with $P = 1$ (in green) and $P = 1$ with transformation (in red) over 200 epochs.

harder and thus result in a higher training loss. The training losses of FoMo-OD with $D = 100$ are also higher than with $D = 20$ since the subspace OD tasks are harder in higher dimensions.

Inference time Figure 10 (left) showed the inference time of FoMo-OD on CPU, comparing typical attention versus the router-based attention (with $R = 500$ routers) under varying context sizes from 1K to 10K. The time is measured on CPU to clearly showcase the scalability trends; *quadratic* without routers and *linear* with routers.

Figure 5 shows the inference time on GPU. Notice that the time is much lower (in milliseconds), thanks to the Transformer architecture taking advantage of GPU parallelism, while the compute time for attention without routers continues to grow faster than that with routers. In implementation, FoMo-OD (with $R = 500$ routers) uses inference context size of 5K by default, which takes about 7.7 ms per test sample on average.

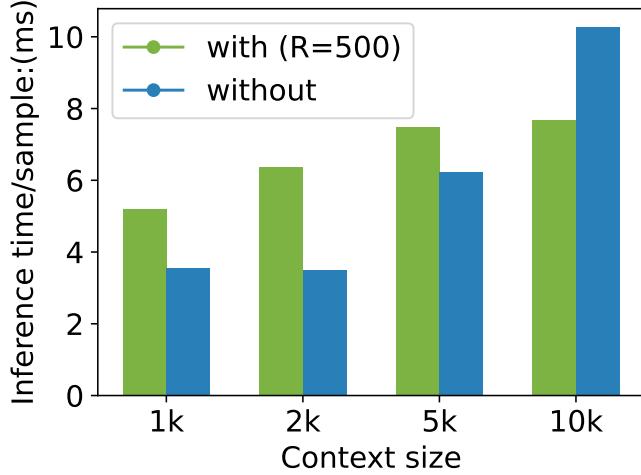


Figure 5: Inference time of FoMo-OD on GPU with vs. w/out router-based attention under varying context size.

G. Detailed Experiment Setup

G.1. Pre-training Dataset Synthesis

During pretraining, we generate unique GMM datasets by first drawing a configuration, including dimensionality $d \in [D]$, number of components $m \in [M]$, centers $\{\mu_j\}_{j=1}^m$ (each $\mu_j \in [-5, 5]^d$) and covariances $\{\Sigma_j\}_{j=1}^m$ ($diag(\Sigma_j) \in [-5, 5]^d$). We set $M = 5$ and vary $D \in \{20, 100\}$ to study pretraining with relatively small and high dimensional datasets, respectively. We synthesize inliers and outliers as described in Section 3.1.

We then sample $S = 5,000$ points that are within the 90th percentile of the GMM. To synthesize outliers, we “inflate” a

subset of dimensions by randomly choosing $|\mathcal{K}| \in [D]$ dimensions and multiplying the corresponding variances by $\times 5$ (following (Han et al., 2022)), i.e. $5 \times \Sigma_{j,kk}$'s for $k \in \mathcal{K}$, and then draw $S = 5,000$ samples from the inflated GMM that are outside the 90th percentile of the original GMM.

To speed up data synthesis via linear transformations, we first draw 500 unique datasets using $m \in [5]$ and $d \in \{1, 2, \dots, 100\}$ (i.e. 5×100) and transform each one $15\times$ using varying parameters (\mathbf{W}, \mathbf{b}) as described in Section D.⁶ This yields 8K unique datasets (500 original and 7,500 transformed) to use at one training epoch (over 1,000 steps with batch size $B = 8$). We repeat this process at each epoch, drawing 500 new datasets and transforming them to reach 8K datasets per epoch.

G.2. Real-world Benchmark Datasets

While pretraining is purely on synthetic datasets, we evaluate FoMo-0D on **57** real-world datasets from the ADBench benchmark (Han et al., 2022) (see Table 19). They consist of 47 popular tabular outlier detection datasets, as well as 10 newly-constructed tabular datasets created from images and natural language tasks by using pretrained models to extract embeddings. We defer to the original paper for the details on these benchmark datasets.

We compare to DTE (Livernoche et al., 2024) and baselines therein as described next, thus, following their semi-supervised OD setup we split each dataset five times into train/test using five different seeds and report the mean performance and its standard deviation. In particular, each random split designates 50% of the inliers as $\mathcal{D}_{\text{train}}$, while $\mathcal{D}_{\text{test}}$ contains the rest of the inliers and all the outlier samples. Note that while the baseline methods require model re-training and inference for each $\mathcal{D}_{\text{train}}/\mathcal{D}_{\text{test}}$ split, FoMo-0D uses the splits only for inference as $\mathcal{D}_{\text{train}}$ is merely passed as context.

Before passing the datasets as input to FoMo-0D, we perform a quantile transform such that the features follow a Normal distribution, to better align with the pretraining data from GMMs.

G.3. Baselines

We compare FoMo-0D against **26** baselines, from classical/shallow methods to modern/deep models. Our baselines include all the baselines imported from one of the latest papers that proposed the SOTA diffusion-based model DTE (Livernoche et al., 2024), and its three variants; DTE-C, DTE-IG, and DTE-NP. Their baselines comprise all those in ADBench (Han et al., 2022); both classical ones (k NN (Ramaswamy et al., 2000), LOF (Breunig et al., 2000), iForest (Liu et al., 2008), HBOS (Goldstein & Dengel, 2012), etc.) and deep models (DeepSVDD (Ruff et al., 2018), DAGMM (Zong et al., 2018), DROCC (Goyal et al., 2020), etc.). They also include more recent approaches based on self-supervised learning (GOAD (Bergman & Hoshen, 2020), ICL (Shenkar & Wolf, 2022), SLAD (Xu et al., 2023), etc.), besides the four additional generative baselines: normalizing planar flows (Rezende & Mohamed, 2015), DDPM (Ho et al., 2020), VAE (Kingma, 2013) and GANomaly (Akcay et al., 2019). We defer to the original paper for additional details. Overall, our 26 baselines consist of the most recent, SOTA approaches for OD that span a diverse family (nonparametric, self-supervised, generative, etc.).

G.4. Hyperparameters for Baselines

Table 2 gives the list of HP values we used to study the HP sensitivity/performance variability of the (from top to bottom) top-4 baselines.

Table 2: Top-4 baselines (from top to bottom) and hyperparameter (HP) configurations.

Baseline	Hyperparameters
DTE-NP	$k \in \{5, 10, 20, 40, 50\}$
k NN	$k \in \{5, 10, 20, 40, 50\}$
ICL	$\text{learning_rate} \in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$
DTE-C	$k \in \{5, 10, 20, 40, 50\}$

⁶It is important to ensure that the eigenvalues of \mathbf{W} (i.e. variances) are not too small such that the dataset does not flatten in any direction. To this end, we draw a random orthonormal basis $\mathbf{U} \in [-1, 1]^{d \times d}$ and a diagonal $\mathbf{\Lambda}$ with eigenvalues $\lambda_{kk} \in ((-1, -0.1] \cup [0.1, 1])^d$, and obtain $\mathbf{W} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$. We also use $\mathbf{b} \in [-1, 1]^d$.

935

G.5. Ranking the 26 baselines

936

Figure 23 presents the visualization of the p -values of the pairwise Wilcoxon signed rank test w.r.t. AUROC among the baseline methods used by Livernoche et al. (2024). We rank these 26 baselines based on their mean p -value (i.e., row-wise average) against the other baselines.

937

938

939

940

941

G.6. Comparison of top-4 baseline variants with varying HP configurations

942

943

944

945

946

947

948

949

Figure 24, 25, 26, 27 give the p -values, respectively comparing the variants of the top-4 baselines (DTE-NP, k NN, ICL, DTE-C) among themselves using different HP configurations, as well as the avg model with the average performance across HPs. (Specifically for ICL, $learning_rate$ (lr) $\in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$; and for others, #nearest-neighbors $k \in \{5, 10, 20, 40, 50\}$). We find that for ICL, $lr = 10^{-3}$ or 10^{-4} are preferable while those that are too small or too large perform poorly. For others, small $k \in \{5, 10\}$ tend to outperform larger $k \in \{40, 50\}$. Note that Livernoche et al. (2024) used $k = 5$ in their paper that proposed DTE (and variants) as well as the k NN baseline for fair comparison, while the DTE avg and k NN avg models across HP configurations perform ~~subpar~~.

950

951

G.7. Sampling time of d -dimensional GMM

952

953

954

955

Figure 6 shows the sampling time of drawing 10,000 points from different GMMs with increasing dimensionality $d = \{10, 20, \dots, 200\}$. We parallelize the sampling process over 10 CPUs, where each CPU draws 1000 samples.

956

957

958

959

960

961

962

We observe that the sampling time grows nonlinearly as the number of dimensions increases, which suggests that it may incur considerable computational overhead to directly draw from the data prior over hundreds of thousands of training steps, motivating the use of our proposed on-the-fly linear transformation T for scalability.

963

964

H. Qualitative Analysis on Sample-to-Sample Attention

965

966

967

968

969

970

971

972

973

974

We sample 50 inliers as context and 100 outliers from a 2-d GMM using the 80th percentile as the labeling threshold, and visualize the top 5 inliers most attended by the 100 outliers based on the average (cross) attention weights over 4 heads from the last layer of FoMo-0D ($D = 100$), which accurately labeled all the 100 outliers. In Figure 7, the most frequently attended inliers are close to either the center of a Gaussian (e.g., 1st, 5th) or the criterion (e.g., 3rd, 4th), suggesting FoMo-0D tends to learn decision boundaries that reflect the prior data generation process.

975

976

977

978

979

980

981

982

983

984

For each outlier, we compute the sum of L2 distances to its top-5 attended inliers (att), the sum of L2 distances to 5 randomly chosen inliers (rdm), and the sum of L2 distances to top-5 inliers with highest likelihood under the GMM ($prob$). We perform Wilcoxon signed rank test between att and rdm (alternative: “less”), att and $prob$ (alternative: “greater”) over all the outliers, with a p -value of 4.4×10^{-4} and 0.99, respectively, suggesting the distances based on attention weights are significantly less than the random distances, and **not** significantly greater than the distances to inliers in high probability region.

985

986

987

988

989

We visualize the top-5 attended inliers for 3 outliers at different position of the 2-d GMM in Figure 8. For a specific outlier, there is a similar trend of attending to the center of a Gaussian (as shown in Figure 7), besides, inliers that reflect the criterion boundary or are close to the outlier are actively attended (e.g., 3rd, 4th in the left, 1st in the middle, 2nd, 5th in the right),

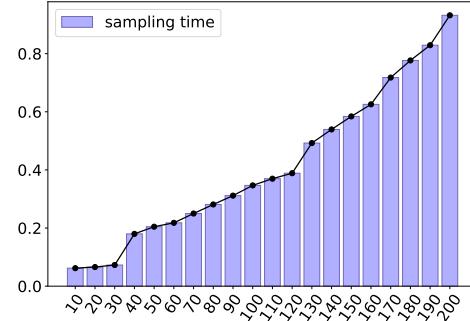


Figure 6: Sampling time (in seconds) of 10,000 points from GMMs with varying number of dimensions.

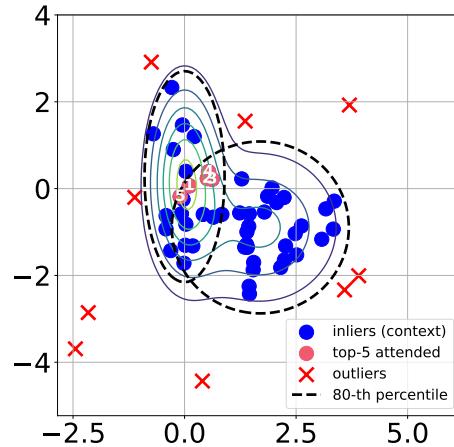
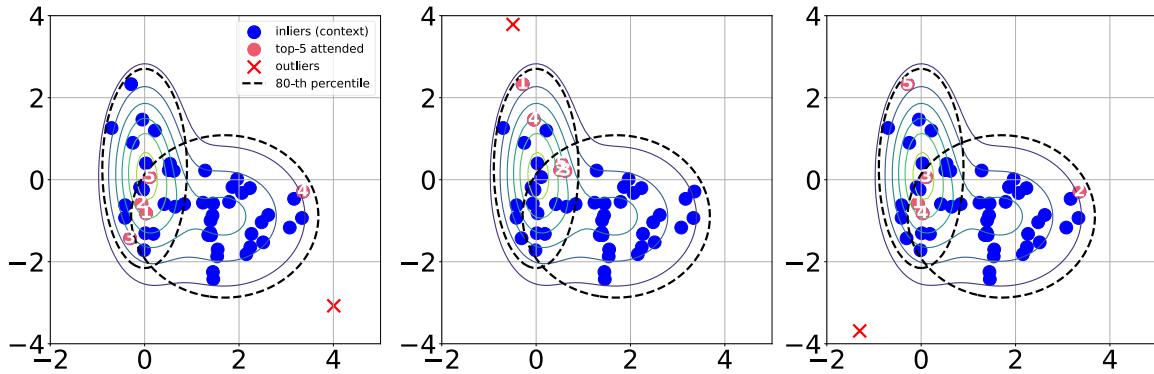


Figure 7: Top-5 attended inliers (all 50 inliers and only part of the outliers are shown for better visualization).

990 suggesting FoMo-0D is incorporating both boundary and nearest neighbor information dynamically for each outlier.
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003



1004 Figure 8: Top-5 attended inliers of 3 outliers at different positions of the GMM
1005
1006

I. Further Results

1007 **Detection performance:** Table 3 shows similar results for FoMo-0D w/ $D = 20$, which is pre-trained on datasets with
1008 considerably fewer dimensions. Even in this limited setting, it performs on par with the 3rd best baseline (ICL, with default
1009 HP) against 30 baselines, with an increased p -value (0.437) when compared to ICL^{avg}. On datasets with $d \leq 20$ which align
1010 with its pre-training data, it outperforms the top 5th baseline and the majority of others. With FoMo-0D pre-trained purely
1011 on synthetic datasets from a simple prior in small dimensions, these results showcase the prowess of PFNs for OD.
1012
1013

1014 Table 3: p -values of the one-sided Wilcoxon signed rank test, comparing FoMo-0D (w/ $D = 20$) to **top 10** baselines with
1015 default HPs, and **top 4^{avg}** baselines³ with **avg.** performance over varying HPs (denoted w/ ^{avg}) over All (57) datasets, those
1016 (24) w/ $d \leq 20$ and (38) datasets w/ $d \leq 50$ dimensions. Although pretrained on datasets w/ small $D = 20$, FoMo-0D
1017 shows **no statistically significant difference from the top 3rd baseline** (ICL, w/ $p = 0.089$) over All datasets, while it
1018 outperforms (w/ $p > 1 - \alpha$) the top 5th (LOF) and onward baselines over datasets w/ $d \leq 20$ (aligned w/ pretraining where
1019 $D = 20$) and on datasets w/ $d \leq 50$ (generalizing beyond pretraining). Rank is avg.'ed over all 57 datasets, where methods
1020 are ranked on each dataset w.r.t. AUROC. (experiment setting: $D = 20$, $P = 50$, $R = 500$, train/inference context size=5K,
1021 no data transformation)
1022

	FoMo-0D	DTE-NP	kNN	ICL	DTE-C	LOF	CBLOF	Feat.Bag.	SLAD	DDPM	OCSVM	DTE-NP ^{avg}	kNN ^{avg}	ICL ^{avg}	DTE-C ^{avg}
$d \leq 20$	-	0.572	0.789	0.968	0.616	0.993	0.989	1.000	0.978	0.906	0.992	0.813	0.924	0.999	1.000
$d \leq 50$	-	0.347	0.794	0.893	0.946	0.997	0.988	1.000	0.963	0.994	0.986	0.574	0.847	0.995	1.000
All	-	<u>0.001</u>	<u>0.019</u>	0.089	0.159	0.394	0.434	0.703	0.516	0.752	0.679	<u>0.007</u>	0.062	0.437	1.000
Rank(avg)	12.59	7.19	8.57	10.34	10.79	11.82	12.81	12.8	12.52	13.50	13.34	8.60	10.63	12.44	21.43

1030 Figure 9 shows the distribution of ranks across datasets for all models. While p-values are the most statistically conclusive,
1031 FoMo-0D achieves a relatively small average rank with notably lower ranks across datasets compared to the majority of the
1032 baselines. Appendix L presents another comparison between detectors through performance profile plots (Dolan & Moré,
1033 2002).

1034
1035
1036 **Running time:** Table 4 presents the total training time and the average inference time per test sample as measured on the
1037 largest dataset for FoMo-0D and the top-3 baselines. Given a new dataset, FoMo-0D bypasses model training (and HP tuning)
1038 and directly performs inference, with an average of 7.7 ms per sample (see Appendix Figure 5). In comparison, all baseline
1039 methods need to train on each individual dataset preceding inference. This training time can be high for deep learning based
1040 models like ICL, and further compounded with training multiple models for hyperparameter tuning purposes. Even for
1041 non-parametric and/or shallow models like kNN and DTE-NP (which queries k nearest neighbors), the training involves
1042 various data pre-processing steps such as constructing a tree-like data structure for fast (often approximate) kNN distance
1043 querying.
1044

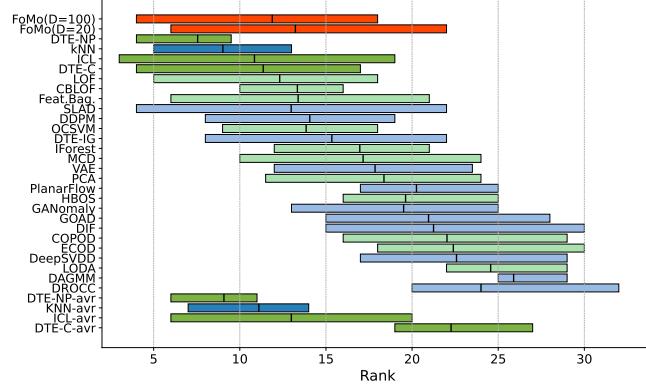


Figure 9: (best in color) Rank (w.r.t. AUROC, lower is better) distribution across all 57 real-world datasets shown via boxplots for (from top to bottom) FoMo-0D in red, all 26 baselines ordered by mean p -value³ (shallow and deep baselines in green and blue), and top 4 baselines' avg variants.

Table 4: Train-time and Inference-time (in milliseconds) of FoMo-0D and the top-3³ baselines (w/ default HPs, excluding the time for model selection/hyperparameter optimization) on our largest dataset donors (see Appendix Table 19). FoMo-0D skips any model training or fine-tuning and takes a mere forward pass for inference out-of-the-box.

Method	FoMo-0D	DTE-NP	kNN	ICL
Train-t (total)	none/0-shot	56.83	1433.74	186461.48
Infer-t (per sample)	7.7	0.76	0.17	0.01

J. Ablation Analyses

In this section, we perform various ablations to study the effect of different design choices in FoMo-0D; namely, **J.1** maximum pretraining data dimensionality D , the number of routers R on **J.2** cost and **J.3** performance, **J.4** context size (both for training and inference), **J.5** number of unique datasets used for pretraining (i.e., reuse periodicity P), data transformation T during synthesis on **J.6** performance and **J.7** speed up, **J.8** data diversity and prolonged training, and finally, **J.9** quantile transforming the benchmark datasets preceding inference.

Unless stated otherwise, most ablation results are performed using FoMo-0D with $D = 20$, as it is faster to pretrain under these many varying settings.

J.1. Effect of pretraining dimensionality D

How does FoMo-0D's generalization performance change by increasing dimensionality of the pretraining data?

We start by comparing FoMo-0D pretrained on datasets with up to $D = 20$ versus $D = 100$ dimensions. Note that learning on higher dimensional datasets is harder, as evident from the relatively larger pretraining loss as shown in Appendix Figure 4. While the statement is accurate in general, it is also partly because subspace outliers “hide” better in higher dimensions.

Comparing Table 1 ($D = 100$) with Table 3 ($D = 20$) w.r.t. p -values over All datasets, we find that FoMo-0D at larger scale does better, where all p -values are larger for $D = 100$ than $D = 20$. We find that FoMo-0D with $D = 20$ performs well on datasets with $d \leq 20$ (i.e., “on its own game”), however beyond its pretraining setting, e.g. on datasets with $d \leq 50$, $D = 100$ is superior to $D = 20$ as shown in Appendix Table 12.

J.2. Effect of routers on cost

What is the running time and memory cost of FoMo-0D with & w/out router-based attention?

Figure 10(left) shows the average inference time per test sample, comparing FoMo-0D using a router-based attention mechanism with $R = 500$ routers (in green) versus FoMo-0D using typical attention without any routers (in blue). As inference context size increases, running time for traditional attention grows quadratically while router mechanism scales linearly.⁷

⁷Note that the inference time is reported on CPUs to show scalability. On GPUs, w/ 5K context size, see Appendix Figure 5, where

Similarly, memory cost with routers is considerably lower when using routers, especially for larger context sizes, as shown in Figure 10(middle).

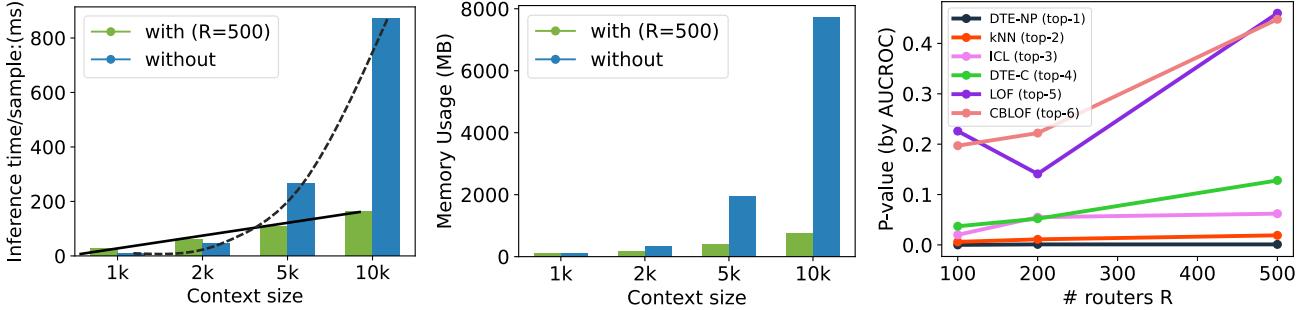


Figure 10: FoMo-0D w/ router mechanism saves time and memory while more #routers perform better, offering a cost-performance trade-off: (left) inference-time (ms) per sample and (middle) memory cost (MB) with & w/out routers by varying context size; (right) performance (based on p -value against top baselines, higher is better) vs. number of routers. (setting: $D = 20$, $P = 1$)

J.3. Effect of routers on performance

What is the impact of the number R of routers (or representatives) on performance?

Router-based mechanism allows to trade-off running time with expressiveness of the attention and hence performance. Figure 10(right) shows the p -values of the Wilcoxon signed rank test as the number of routers R is increased from 100 to 200 and 500, comparing FoMo-0D to each of the top-6 baselines. We notice that FoMo-0D performance tends to increase monotonically with more routers.

J.4. Effect of context size

What is the impact of context size, both during model pretraining as well as during inference?

To study how performance changes by context size, we train FoMo-0D with varying context size in {1K,2K,5K} and employ each pretrained model for inference with varying context size in {1K,2K,5K,10K}. Table 5 shows the results, where performance is depicted by the average rank of FoMo-0D (the lower, the better).

Table 5: Average rank (based on comparison to 30 baselines w.r.t. AUROC) of FoMo-0D across datasets under different context sizes for training and inference. Smaller ranks imply better performance. (setting: $D = 20$, $R = 500$, $P = 1$)

	Infer:1K	Infer:2K	Infer:5K	Infer:10K
Train:1K	13.816	14.623	15.193	15.439
Train:2K	13.079	13.219	13.439	13.561
Train:5K	13.088	13.211	13.307	13.430

We find that training with a larger context improves performance at any inference context size. On the other hand, perhaps counter-intuitively, FoMo-0D with smaller inference context size does better. We conjecture that is because the #routers-to-context size ratio increases with a larger context size at inference, limiting the expressive power of the “bottleneck” attention mechanism. The pairwise statistical tests among the $3 \times 4 = 12$ models support these observations, as shown in Figure 11. Interestingly, when the training context size is large enough at 5K, inference with 10K samples generalizes beyond training with no statistical evidence for performance difference (at 0.05) from other inference context sizes.

typical attention takes advantage of parallelism (6.5ms), while router-based attention is slightly slower (7.7 ms w/ 500 routers) due to its two sequential self-attentions; see Eq.s (5) and (6).

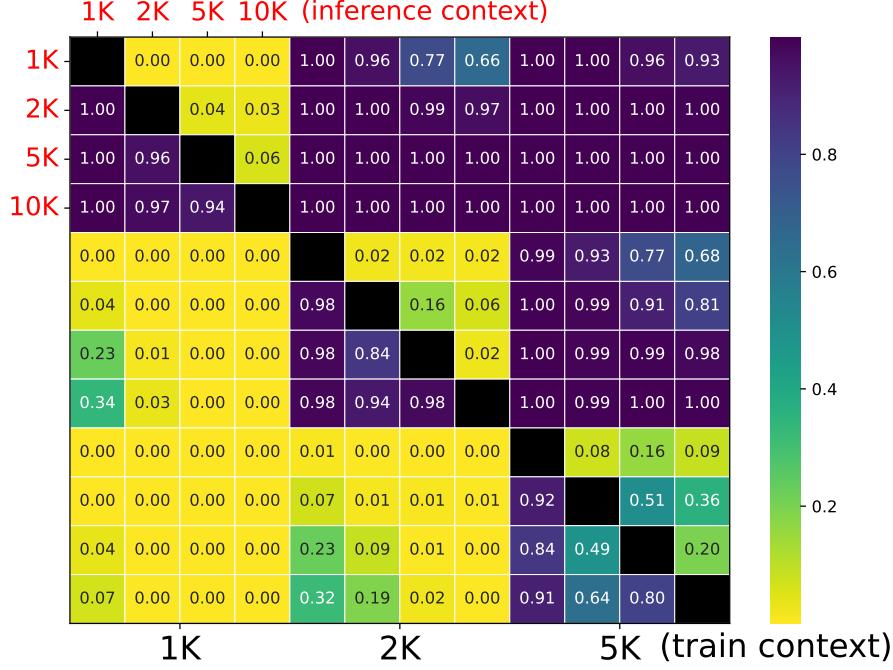


Figure 11: p -values of the pairwise Wilcoxon signed rank test between models (larger p implies col-method is better than row-method) w/ different context sizes for **training** (1K/2K/5K, 1st/2nd/3rd four grids, in **black**) and **inference** (1K/2K/5K/10K, every 1st/2nd/3rd/4th grid, in **red**): Larger training context improves overall performance, while smaller inference context is preferable.

Table 6: Ablation results on dataset reuse across epochs with varying $P \in \{1, 50, 100\}$ show stable p -values against the top-5 baselines, where there is no statistical evidence to suggest performance difference between FoMo-0D with $D = 20$ and the top 3rd baseline at 0.05 w.r.t. pairwise Wilcoxon signed rank test comparisons, while it continues to significantly outperform the top 5th baseline (LOF) when $d \leq 50$. (setting: $D=20$, $R=500$, context size=5K, w/out transformation T)

	$P = 1$ (#unique datasets: 8K)					$P = 50$ (#unique datasets: $8 \times 50 = 400$ K)					$P = 100$ (#unique datasets: $8 \times 100 = 800$ K)				
top-5	DTE-NP	kNN	ICL	DTE-C	LOF	DTE-NP	kNN	ICL	DTE-C	LOF	DTE-NP	kNN	ICL	DTE-C	LOF
All	0.001	0.019	0.062	0.128	0.460	0.001	0.019	0.089	0.159	0.394	0.001	0.015	0.072	0.121	0.290
$d \leq 20$	0.583	0.755	0.943	0.736	0.998	0.572	0.789	0.968	0.616	0.993	0.439	0.678	0.953	0.550	0.972
$d \leq 50$	0.415	0.750	0.869	0.962	0.999	0.347	0.794	0.893	0.946	0.997	0.293	0.697	0.890	0.924	0.994

J.5. Effect of number of unique datasets

How do FoMo-0D performances compare when pretrained on unique vs. reused datasets, via varying periodicity P ?

Next we study the effect of dataset *reuse at epoch level* (w/out transformation) on performance as presented in Section D. We vary reuse periodicity P in $\{1, 50, 100\}$, and accordingly, increase the number of unique datasets used for pretraining across epochs. As shown in Table 6, FoMo-0D (w/ $D = 20$) performs similarly with varying dataset reuse. In fact, it is competitive even with $P = 1$, remaining no different from the 3rd best baseline (ICL) across All (57) datasets, while significantly outperforming the top 5th (LOF) across (24) datasets with $d \leq 20$ as well as (38) with $d \leq 50$.

J.6. Effect of transformation T for synthesis

How do FoMo-0D performances compare when pretrained on datasets with vs. w/out linear transformation?

Setting $P = 1$, we next study the impact of linear transformation T . Table 7 presents the results, where we compare reuse of the *same* 8K unique datasets across epochs (w/out T), versus *transforming* these datasets with T at every epoch with different parameters (w/ T). FoMo-0D performance remains stable; no statistical evidence for performance difference from the top 3rd model on All datasets, while significantly outperforming the top 5th across those with $d \leq 20$ and $d \leq 50$. This suggests that T can be employed without sacrificing performance to save time during pretraining.

Table 7: Ablation results on performance w/ & w/out linear transformation T show stable p -values against the top-5 baselines, with no statistical evidence for performance difference between FoMo-0D with $D = 20$ and the top 3rd baseline at 0.05 w.r.t. pairwise Wilcoxon signed rank test comparisons. (setting: $D = 20$, $R = 500$, context size=5K, $P = 1$)

top-5	w/out transformation T					w/ transformation T				
	DTE-NP	kNN	ICL	DTE-C	LOF	DTE-NP	kNN	ICL	DTE-C	LOF
All	0.001	0.019	0.062	0.128	0.460	0.002	0.015	0.226	0.210	0.280
$d \leq 20$	0.583	0.755	0.943	0.736	0.998	0.648	0.708	0.988	0.718	0.955
$d \leq 50$	0.415	0.750	0.869	0.962	0.999	0.264	0.382	0.971	0.900	0.963

J.7. Speed up by T

What is the time saving on data synthesis with linear transformation?

Figure 12 shows the distribution of pretraining running-time per epoch with and w/out data transformation. Specifically, we compare (left) generating 8K unique datasets/epoch on-the-fly and (right) first generating 500 unique datasets on-the-fly and then transforming each one 15 times using T with different parameters to reach 8K datasets at each epoch.

Notice that pretraining with T takes about 450 sec./epoch on average, while without T it requires 1200 sec./epoch to generate 8K unique datasets and gradient descent across 1000 steps. Different from other ablation results, which are based on the $D = 20$ model, here we report the running times for our $D = 100$ model. Overall, our final FoMo-0D took **≈25 hours** for pre-training (450 sec. $\times 200$ epochs). Importantly, this is a one-time cost that amortizes across many downstream tasks with as low as **7.7 ms inference time** per test sample (see Table 4 and Appendix Figure 5).

J.8. Effect of data diversity and prolonged training

How does FoMo-0D’s performance change by increasing pretraining data diversity and number of training epochs?

Originally we have trained FoMo-0D w/ $D = 100$ using 0.5K unique + 7.5K transformed datasets over 200 epochs. As mentioned earlier, learning in higher dimensions tends to incur a larger loss in general but also specifically here, as subspace outliers are harder to detect in high dimensions.

Toward reducing the loss further, we resume the pretraining for another 100 epochs. Further, to simplify the tasks and thereby increase data diversity, we also decrease the inlier/outlier labeling percentile threshold from 90% to 80% during on-the-fly data generation in the last 100 epochs. In Figure 13, we present the training loss of FoMo-0D ($D = 100$) trained with 0.5K unique + 7.5K transformed datasets/epoch over 200 epochs (90th percentile as labeling threshold) and then 100 additional epochs (80th percentile as the threshold) to show how data diversity and amount affect model performance.

Figure 14 compares FoMo-0D’s performance (w/ $D = 100$) to top-5 baselines w.r.t. p -values of the paired Wilcoxon signed rank test on datasets with $d \leq 100$, after the first 200 epochs versus after 300 epochs. The increase in all the p -values showcases the benefit of additional training.

J.9. Effect of applying quantile transform on benchmark datasets

What is the impact of quantile data transform preceding inference on performance?

We pretrain FoMo-0D on synthetic datasets from a simple data prior based on GMMs. The real-world benchmark datasets, on the other hand, may exhibit features with distributions different from Gaussians. To close the gap, we apply a quantile transform (denoted QT) on the benchmark datasets prior to feeding them to FoMo-0D for inference, which transforms the features to exhibit a more Gaussian-like probability distribution.

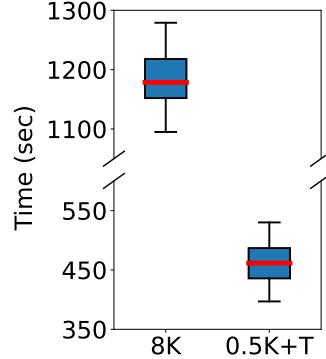


Figure 12: Runtime/epoch dist.n over 100 epochs for FoMo-0D ($D=100$) with (left) $P=100$, i.e. 8K unique datasets/epoch vs. (right) 0.5K unique+7.5K transformed datasets/epoch.

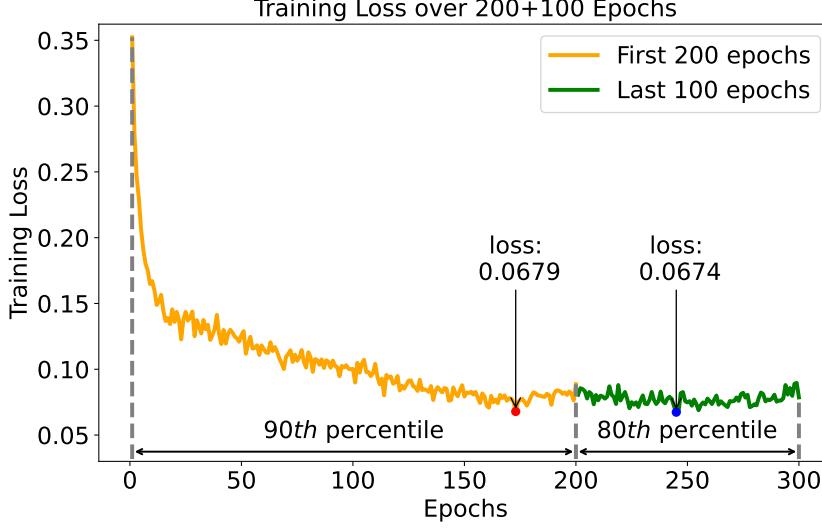


Figure 13: (best in color) Training loss of FoMo-0D ($D = 100$) with 0.5K unique + 7.5K transformed datasets/epoch for 200 epochs (in orange), followed with additional 100 epochs of training (in green). For the first 200 epochs we train with 90th percentile as the inlier/outlier threshold, which we reduce to 80th in the subsequent 100 epochs.

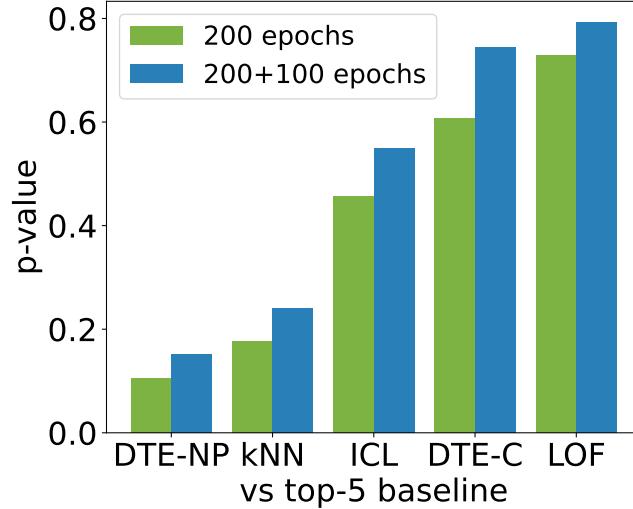


Figure 14: p -values increase with additional 100 epochs of pretraining, i.e. FoMo-0D w/ $D = 100$ performs better against top-5 baselines on datasets w/ $d \leq 100$.

Figure 15 compares the performance of three FoMo-0D w/ $D = 100$ variants with and w/out QT against the top-5 baselines w.r.t. the p -values of the paired Wilcoxon signed rank test. FoMo-0D tends to perform better as suggested by larger p -values when QT is applied.

Besides the ablation studies, we provide a qualitative case study of sample-to-sample attention in Appendix H, showing that an outlier attends to the points in context that are within a short distance significantly more than random points, suggesting that PFNs tend to mimic non-parametrics.

K. Generalization Analyses

K.1. Generalization to Out of Distribution Synthetic Datasets

We conduct analyses to understand FoMo's ability to generalize on out-of-distribution synthetic GMM datasets. Besides the in-distribution setting for pre-training (i.e., $\mu \in [-5, 5]$, $\Sigma \in (0, 5]$, $m \leq 5$, $d \leq 100$), we consider the following

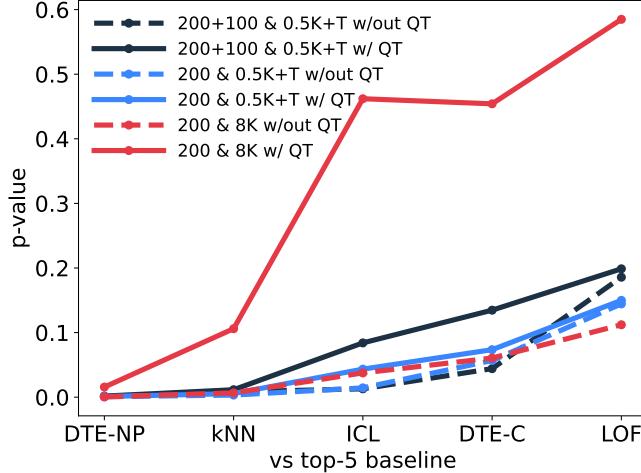


Figure 15: p -values increase, i.e. FoMo-0D performance improves, against top-5 baselines with quantile transform (QT) preceding inference, for 3 different settings of FoMo-0D w/ $D = 100$.

out-of-distribution settings: **(a)** mean and covariance significantly out of range, with $\mu \in [-50, -5] \cup [5, 50]$, $\Sigma \in [5, 50]$, denoted as “ $|\mu|, |\Sigma| \in [5, 50]$ ”; **(b)** number of clusters significantly out of range, denoted as “ $m \in [5, 50]$ ”; **(c)** number of dimensions significantly out of range, denoted as “ $d \in [100, 500]$ ”; **(d)** binary outliers with values either 0 or 1 in one dimension from the sub-dimensions, denoted as “binary”; **(e)** “all”, which combines all the variants above. For each setting, we generate 1000 datasets with random seeds from 0 to 999, where on each dataset, we simulate 1000 test points with an outlier rate of 5% and evaluate FoMo-0D with a context length of 5000. We present the results with averaged performance over 1000 datasets for each setting in Table 8.

Table 8: Average metric score \pm standard dev. over 1000 seeds for different out-of-distribution (OOD) synthetic GMMs. FoMo-0D remains robust against OOD test datasets as in (a)–(d), maintaining similar performance to in-distribution performance (top). Performance is affected more when datasets are OOD w.r.t. multiple factors combined as in (e).

Dataset	AUROC	AUCPR	F1
ID: in-distribution	98.55 ± 2.73	91.17 ± 13.07	86.74 ± 15.43
(a) OOD w.r.t. $ \mu , \Sigma \in [5, 50]$	94.79 ± 7.53	80.62 ± 21.19	76.32 ± 19.85
(b) OOD w.r.t. $m \in [5, 50]$	97.69 ± 3.59	86.72 ± 15.57	81.20 ± 16.23
(c) OOD w.r.t. $d \in [100, 500]$	96.22 ± 9.01	86.37 ± 23.27	83.23 ± 22.08
(d) OOD w.r.t. binary variable	100.00 ± 0.00	100.00 ± 0.06	99.97 ± 0.34
(e) OOD w.r.t. all combined	85.44 ± 16.96	64.17 ± 35.07	63.99 ± 33.53

We can observe different extents of performance degradation when applying out-of-distribution variations. Compared to other single variations, FoMo-0D seems to suffer more from inflating the mean and covariances, as due to the significant deviation in the parameters of the GMMs, inliers generated under such a setting are seemingly “outliers” w/o any reference points. Surprisingly, although FoMo-0D is only trained on continuous data, it can almost perfectly classify binary outliers hidden in one of the sub-dimensions, suggesting FoMo-0D could potentially generalize to discrete data at test time.

However, with all variations added, FoMo-0D becomes less capable compared to one single out-of-distribution variation, although there might exist some signals (e.g., binary labels) in favor of its decision-making process, for which training a powerful model with more comprehensive priors could possibly alleviate the issue.

K.2. Generalization to Out of GMM-Distribution Real-World Datasets

To understand how FoMo-0D generalizes from the pre-training GMM data priors to complex real-world datasets when performing zero-shot OD, we conduct the goodness of fit test (Huber-Carol et al., 2012) for datasets in ADBench. We fit GMMs to each real-world dataset D_{real} with up to 5 components (as with our pre-training datasets), then sample D_{syn} from

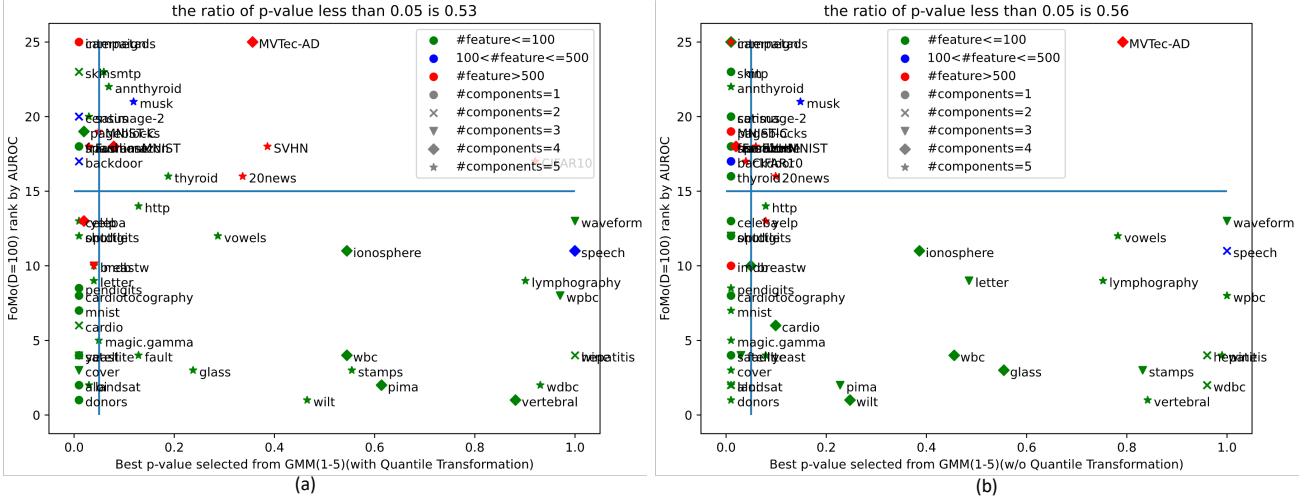


Figure 16: Goodness of fit test results for 57 datasets in ADBench. The x-axis is the p -value of a dataset, where a small p -value indicates GMMs are not a good fit for the dataset, and the y-axis is the rank of FoMo-0D ($D=100$) on that dataset (the smaller the better). We plot p -value = 0.05 (vertical) and rank = 15 (horizontal) as references. The left figure (a) is with quantile transformation while figure (b) is without quantile transformation. We use different colors to represent datasets at different dimensions, and use different markers to represent different numbers of clusters.

the best-parameter-fitted GMM and perform a two-sample test⁸ on D_{real} and D_{syn} , with the null hypothesis that they come from the same distribution. A smaller p -value (≤ 0.05) of such a test provides evidence toward rejecting the null, which suggests GMM is not a good fit for the dataset (i.e., the pre-training data distribution is different from the test data).

We present the results in Figure 16, depicting the p -value (of the goodness-of-GMM-fit test) vs. FoMo-0D's performance rank (the lower the better) among 30 baselines. We report the result both with quantile transformation in Figure 16(a) and without quantile transformation in 16(b). Since the two figures are highly similar, our next analysis will primarily focus on the figure with quantile transformation to align with our model's implementation. We plot the vertical and horizontal lines as p -value = 0.05 and rank = 15. For p -value ≥ 0.05 and rank < 15 , we observe that performance is good on datasets with relatively large p -value where we cannot reject the null (i.e. GMM is a relatively good fit). This is where arguably FoMo-0D recalls its data prior distribution and generalizes to datasets similar to those seen during pretraining. We also see, for p -value < 0.05 and rank ≥ 15 , datasets with relatively poor performance where we can reject the null (i.e. GMM is not a good fit). These can be attributed to falling short in generalization to OOD datasets.

On the other hand, we observe many datasets concentrate on p -value < 0.05 and rank < 15 , where p -value is small (GMM not a good fit) yet the performance is competitive — those are the datasets on which FoMo-0D is likely to have achieved out-of-distribution generalization. It remains an open (theoretical) question to understand what (algorithm, if any) FoMo-0D might have learned that generalizes to out-of-distribution datasets. It is also an open (empirical) quest to explore whether a more complex data prior, beyond GMMs, could further push the performance up and by how much.

K.3. Generalization to Out of Distribution (OOD) Detection Tasks

We further evaluate FoMo-0D on more complex datasets (e.g., ImageNet-level). Specifically, we employ OpenOOD (Zhang et al., 2023), an out-of-distribution (OOD) detection benchmark, where models are trained on labeled in-distribution datasets, with K known classes, and then evaluated on out-of-distribution datasets, aiming to detect $K + 1, K + 2, \dots$ novel classes. Although OOD detection is inherently different from OD, we can construct an OD dataset from OOD datasets, treating all K class samples as inliers and the $K + 1, K + 2, \dots$ OOD samples as outliers. For the in-distribution datasets, we choose ImageNet1K, which contains 1000 categories of images, and ImageNet200, a subset of ImageNet1K containing 200 categories. We further choose SSB-hard, NINCO, iNaturalist, Textures, and OpenImage-O as the out-of-distribution datasets, which gives us a total number of $2 \times 5 = 10$ datasets that are ImageNet-level complex.

⁸We use e-test from <https://www.rdocumentation.org/packages/energy/versions/1.7-11/topics/eqdist.eetest>

Following Han et al. (2022), we create 10 new OD datasets from OpenOOD containing 10,000 samples with 5% outliers, and use the embedding from the last average pooling layer of ResNet18 (He et al., 2016) as the feature (512) for each sample. Comparing FoMo-0D with the top-4 (on our original testbed) baselines in the order of: DTE-NP, kNN, ICL, DTE-C, we follow Livernoche et al. (2024) and report mean (standard dev.) over 5 runs (seed=0/1/2/3/4) on each dataset. We present the results with in-distribution datasets being ImageNet200 and ImageNet1K in Table 9 and 10, respectively.

Table 9: Average AUROC score \pm standard dev. over five seeds for in-distribution dataset being **ImageNet200**. We use blue and green respectively to mark the top-1 and the top-2 method.

dataset	DTE-NP	kNN	ICL	DTE-C	FoMo-0D
ssb-hard	58.03 \pm 0.00	58.14 \pm 0.00	60.52 \pm 0.25	60.74 \pm 1.88	58.34 \pm 1.55
ninco	53.28 \pm 0.00	54.14 \pm 0.00	59.56 \pm 0.63	58.83 \pm 1.54	55.16 \pm 2.19
inaturalist	29.38 \pm 0.00	29.51 \pm 0.00	35.96 \pm 1.10	41.77 \pm 2.84	38.85 \pm 3.29
textures	59.28 \pm 0.00	59.91 \pm 0.00	66.40 \pm 0.69	70.33 \pm 3.18	59.89 \pm 2.07
openimageo	52.82 \pm 0.00	53.79 \pm 0.00	55.20 \pm 0.69	59.09 \pm 1.50	54.77 \pm 1.19
average	50.56	51.10	55.53	58.15	53.40

Table 10: Average AUROC score \pm standard dev. over five seeds for in-distribution dataset being **ImageNet1K**. We use blue and green respectively to mark the top-1 and the top-2 method.

dataset	DTE-NP	kNN	ICL	DTE-C	FoMo-0D
ssb-hard	55.63 \pm 0.00	55.94 \pm 0.00	58.79 \pm 1.20	59.17 \pm 1.82	56.73 \pm 2.65
ninco	48.23 \pm 0.00	49.10 \pm 0.00	55.25 \pm 0.87	57.60 \pm 3.93	52.70 \pm 2.70
inaturalist	30.24 \pm 0.00	30.28 \pm 0.00	35.03 \pm 1.42	41.96 \pm 3.13	38.94 \pm 4.59
textures	54.38 \pm 0.00	55.43 \pm 0.00	61.30 \pm 0.95	63.10 \pm 3.72	55.18 \pm 2.92
openimageo	54.31 \pm 0.00	54.91 \pm 0.00	54.02 \pm 0.43	58.71 \pm 2.08	56.95 \pm 3.89
average	48.56	49.13	52.88	56.11	52.10

We further report the p -value of the Wilcoxon sign test between the baselines and FoMo-0D on the 10 datasets from OpenOOD, as well as on the expanded benchmark combining those 10 with our original ADBench (10+57) in Table 11. In terms of metric values, FoMo-0D performs 2nd or 3rd best across OOD datasets. p -values show that it significantly outperforms DTE-NP and kNN ($p > 0.95$, such that the p -value < 0.05 for rejecting the null hypothesis and accepting the alternative hypothesis that the “baseline-minus-FoMo-0D” gap is smaller than zero) and is no different from ICL (2nd best after DTE-C). These results demonstrate that FoMo-0D generalizes beyond OD datasets and maintains strong zero-shot OD performance on complex, ImageNet-level OOD benchmarks.

Table 11: p -value of the Wilcoxon sign test (alternative: “greater”) between baselines and FoMo-0D on OpenOOD and combined benchmark on AUROC. A small p -value (≤ 0.05) means that there is statistical evidence for the alternative hypothesis such that baselines achieve higher metric performance than FoMo-0D.

method	DTE-NP	kNN	ICL	DTE-C
OpenOOD	1	0.9951	0.1875	0.0009
OpenOOD+ADBench	0.1271	0.3308	0.3153	0.1265

Interestingly, we observe that ICL and DTE-C outperform DTE-NP and kNN on the OpenOOD datasets, whereas on ADBench, DTE-NP and kNN are the top-2 methods outperforming ICL and DTE-C. We hypothesize this is because it is harder for non-parametric methods like DTE-NP and kNN to estimate meaningful decision boundaries in high dimensions (e.g., 512). In contrast, the performance of FoMo-0D is consistently competitive, where the p -values on the combined testbed (OpenOOD+ADBench) show that FoMo-0D is as competitive as all the top baselines across 67 diverse datasets, while maintaining zero-shot detection ability.

L. Performance Profile Plots

To enable a comprehensive comparison of different methods, we adopt τ performance profile plots as described in (Dolan & Moré, 2002). These plots display the cumulative distribution of the τ metric—which quantifies suboptimality relative to the best-performing method. By computing sorted τ values along with their cumulative probabilities, we then use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

Figure 17, Figure 18, and Figure 19 illustrate performance profile plots of FoMo-0D and other baselines across all datasets. The results show that FoMo-0D ($D=100$) ranks at **top-5 (w.r.t. AUROC)**, **top-3 (w.r.t. AUPR)** and **top-1 (w.r.t. F1)**, respectively, outperforming many baselines.

Moreover, the performance of FoMo-0D ($D=100$) is even better (i.e., ranked within top-2) when tested on datasets with dimensions less than 100. As shown in Figure 20, Figure 21, and Figure 22, the area under the curve of FoMo-0D ($D=100$) ranks at **top-1 (w.r.t. AUROC)**, **top-2 (w.r.t. AUPR)** and **top-2 (w.r.t. F1)**, respectively, under this setting.

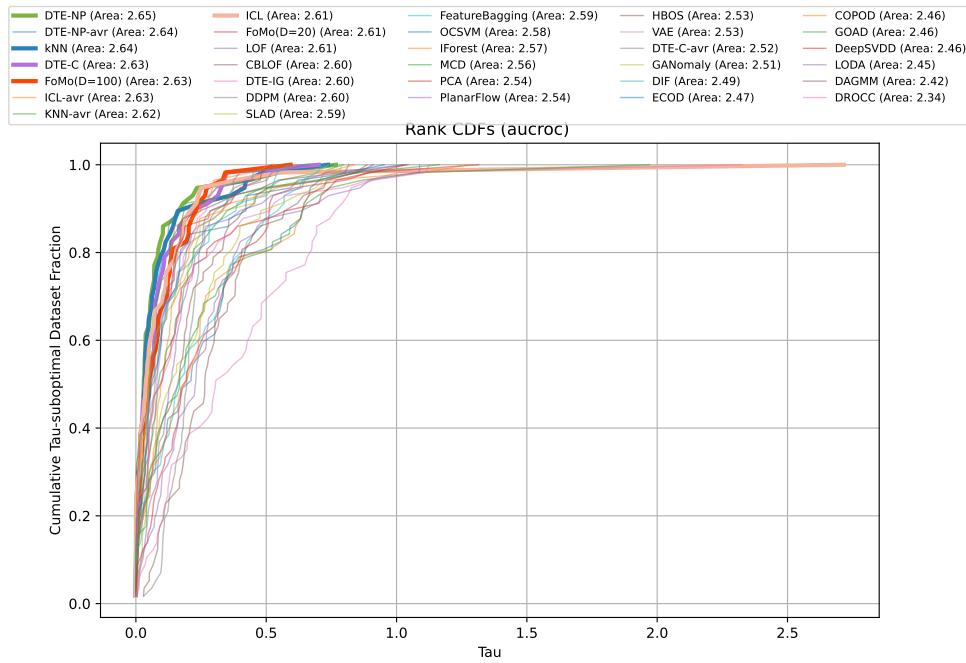


Figure 17: FoMo-0D ranks in **top-5** based on the performance profile plots of all detectors w.r.t. **AUROC** across **all datasets**. In the plot, x-axis represents the τ values—performance ratios that compare each method’s metric to the best performance achieved, while y-axis displays the cumulative fraction of test datasets for which a method’s performance is within the τ value. We use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

M. Full Results

Tables 13.1& 13.2, 14.1 & 14.2, and 15.1 & 15.2 respectively show the AUROC, AUPR, and F1 scores of the top-4 baselines, DTE-NP, kNN, ICL, and DTE-C as well as their corresponding avg model with the average performance across HPs, as listed in Table 2.

Tables 16.1&16.2, 17.1&17.2, and 18.1&18.2 respectively show the AUROC, AUPR, and F1 scores of all methods across all benchmark datasets. In all these tables, the last four rows show the avg_rank of methods across datasets, and p -values of the Wilcoxon signed rank test comparing FoMo-0D w/ $D = 100$ with other baselines. The preceding four rows are the same for FoMo-0D w/ $D = 20$, when ranking 31 models (26 baselines + 4 avg variants of top-4 baselines + FoMo-0D w/ $D = 20$).

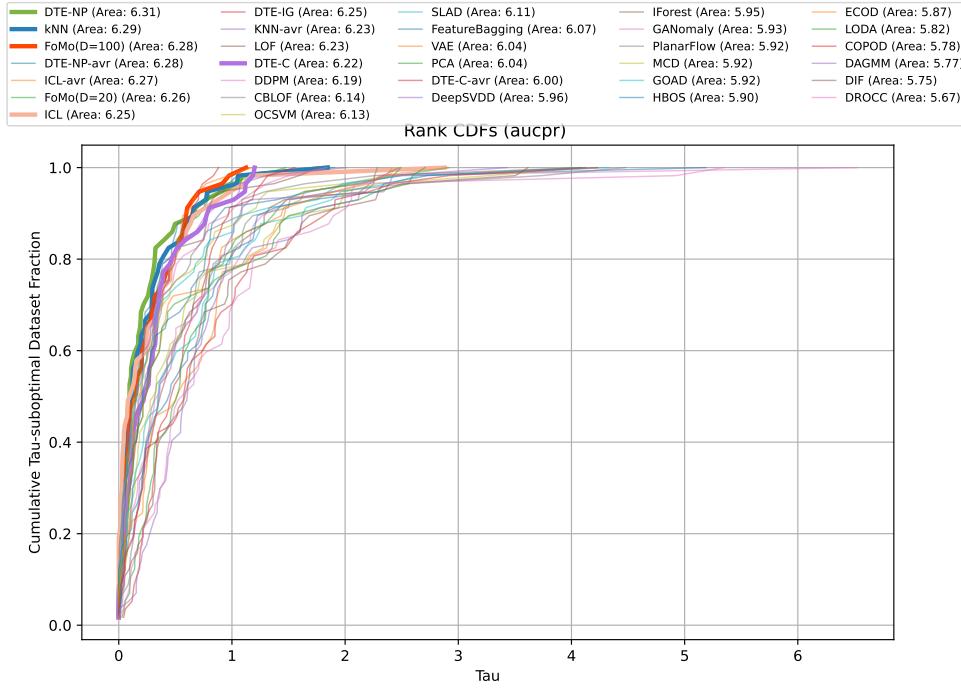


Figure 18: FoMo-0D ranks at **top-3** based on the performance profile plots of all detectors w.r.t. **AUPR** across **all datasets**. In the plot, x-axis represents the τ values—performance ratios that compare each method’s metric to the best performance achieved, while y-axis displays the cumulative fraction of test datasets for which a method’s performance is within the τ value. We use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

Table 12: Comparison of methods across datasets. (top row) Rank w.r.t. AUROC performance avg.’ed over 57 datasets is presented for FoMo-0D (with $D = 100$), **top-10 baselines** with default HPs, and **top-4³** baselines with performance avg.’ed over varying HPs (denoted w/ ^{avg}); followed by p -values of the pairwise Wilcoxon signed rank test, comparing FoMo-0D to each baseline (from top to bottom) over All (57) datasets, those (24) w/ $d \leq 20$, (38) w/ $d \leq 50$, (42) w/ $d \leq 100$ and (46) datasets w/ $d \leq 500$ dimensions. FoMo-0D performs as well as (**i.e., statistically no different from**) the **2nd best model** (kNN , w/ $p = 0.106$) across All datasets, while it is **comparable to** ($p > 0.05$) or **better than** ($p > 0.95$) **all baselines** over datasets w/ $d \leq 100$ (aligned w/ pretraining where $D = 100$) and $d \leq 500$ (generalizing beyond pretraining).

	FoMo-0D	DTE-NP	kNN	ICL	DTE-C	LOF	CBLOF	Feat.Bag.	SLAD	DDPM	OCSVM	DTE-NP ^{avg}	kNN^{avg}	ICL ^{avg}	DTE-C ^{avg}
Rank(avg)	11.886	7.553	9.018	10.851	11.36	12.316	13.342	13.386	12.982	14.061	13.851	9.079	11.105	12.991	22.263
All	-	0.016	0.106	0.462	0.454	0.585	0.750	0.823	0.759	0.901	0.895	0.112	0.315	0.670	1.000
$d \leq 20$	-	0.428	0.665	0.987	0.727	0.911	0.940	0.987	0.868	0.758	0.968	0.781	0.868	0.990	1.000
$d \leq 50$	-	0.734	0.923	0.992	0.973	0.989	0.987	0.999	0.948	0.985	0.986	0.948	0.967	0.989	1.000
$d \leq 100$	-	0.415	0.700	0.949	0.953	0.970	0.971	0.996	0.876	0.980	0.978	0.752	0.860	0.958	1.000
$d \leq 200$	-	0.315	0.605	0.923	0.919	0.944	0.977	0.990	0.904	0.970	0.983	0.663	0.789	0.937	1.000
$d \leq 500$	-	0.220	0.569	0.827	0.894	0.960	0.968	0.994	0.910	0.960	0.979	0.607	0.756	0.846	1.000

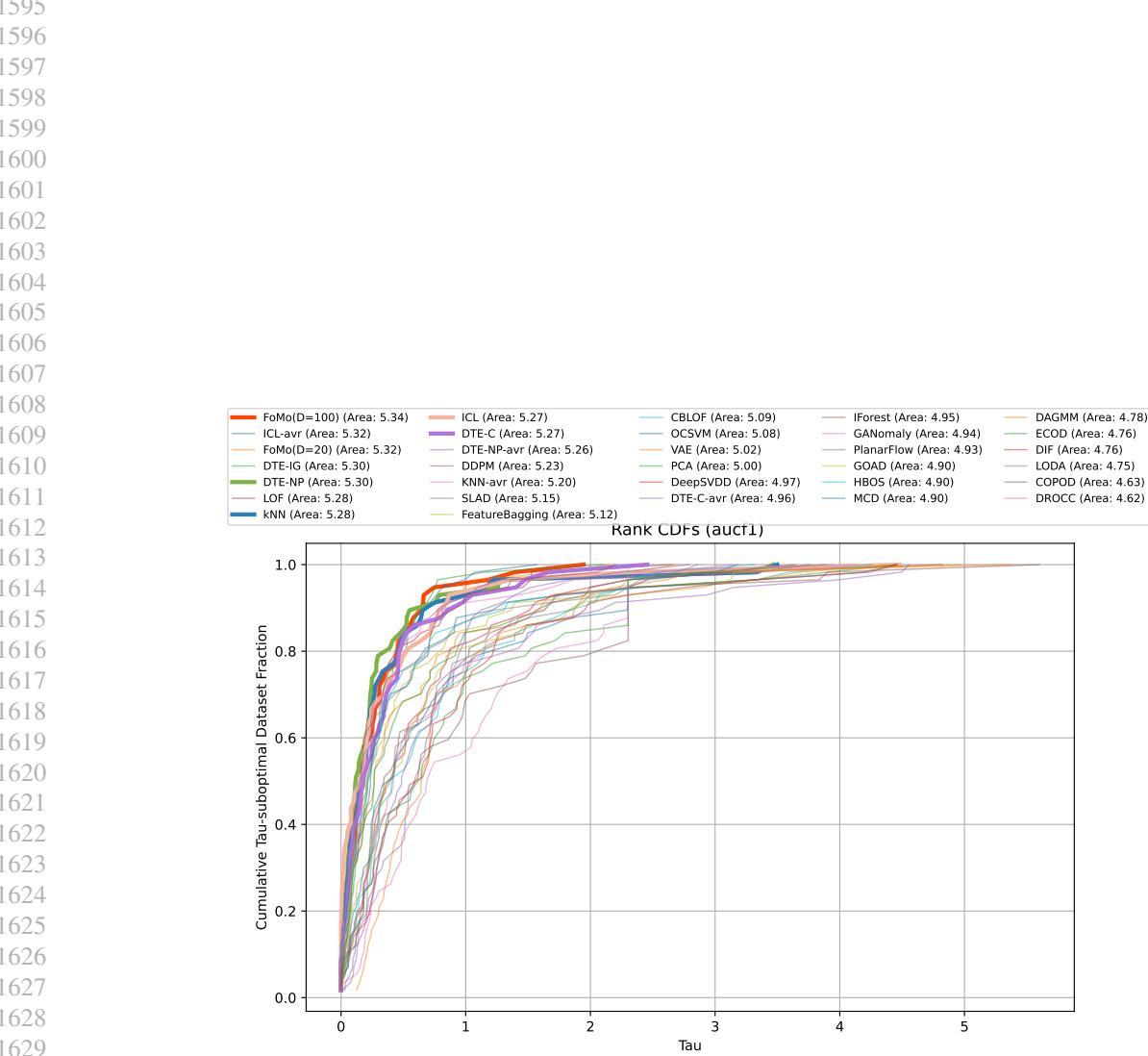


Figure 19: FoMo-0D ranks at **top-1** based on the performance profile plots of all detectors w.r.t. F1 across all datasets. In the plot, x-axis represents the τ values—performance ratios that compare each method’s metric to the best performance achieved, while y-axis displays the cumulative fraction of test datasets for which a method’s performance is within the τ value. We use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

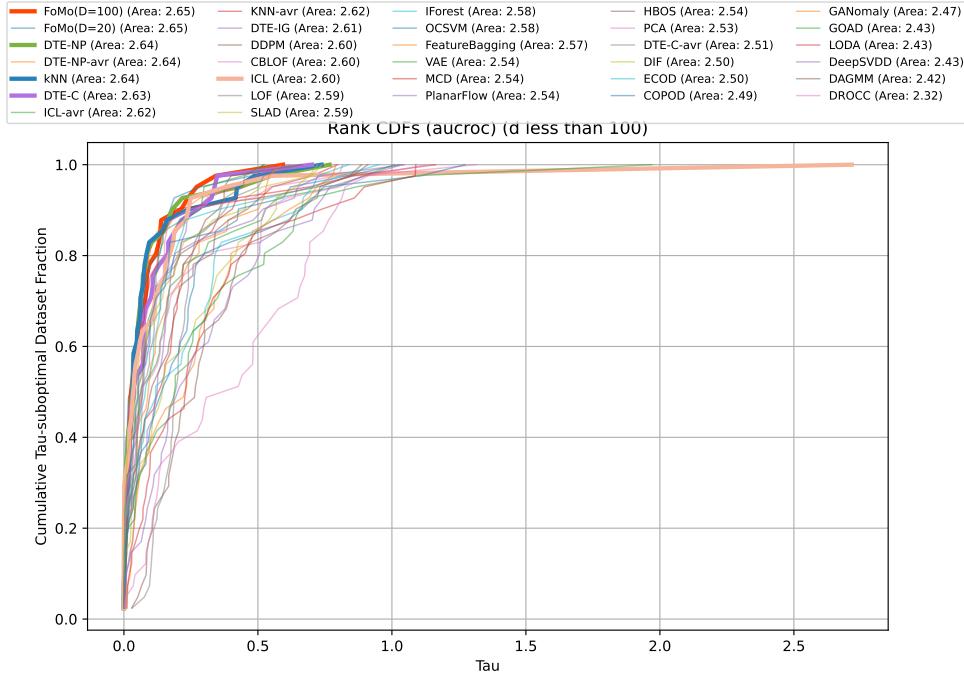


Figure 20: FoMo-0D ($D=100$) ranks at **top-1** based on the performance profile plots of all detectors w.r.t. AUROC in datasets with dimensions less than $d \leq 100$. In the plot, x-axis represents the τ values—performance ratios that compare each method’s metric to the best performance achieved, while y-axis displays the cumulative fraction of test datasets for which a method’s performance is within the τ value. We use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

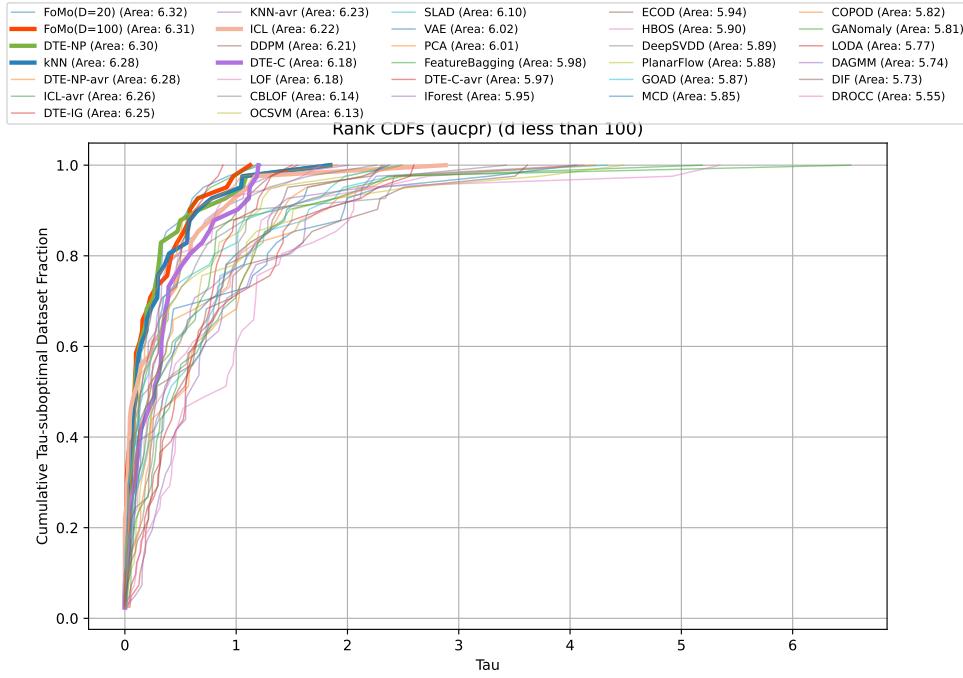


Figure 21: FoMo-0D ($D=100$) ranks at **top-2** based on the performance profile plots of all detectors w.r.t. **AUPR** in datasets with dimensions less than $d \leq 100$. In the plot, x-axis represents the τ values—performance ratios that compare each method’s metric to the best performance achieved, while y-axis displays the cumulative fraction of test datasets for which a method’s performance is within the τ value. We use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

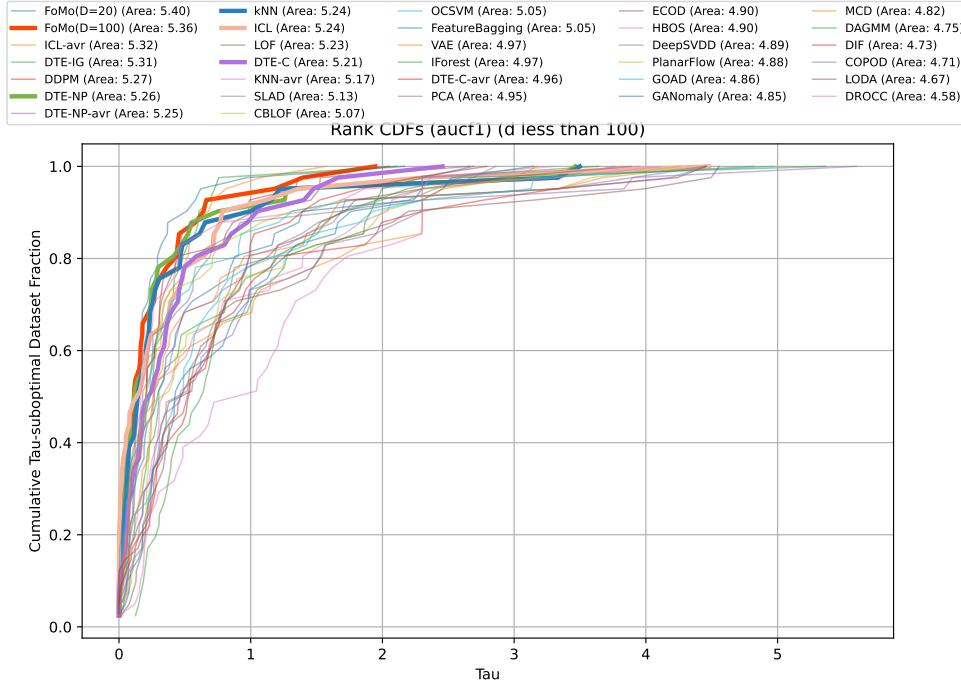


Figure 22: FoMo-OD (D=100) ranks at **top-2** based on the performance profile plots of all detectors w.r.t. **F1** in *datasets with dimensions less than $d \leq 100$* . In the plot, x-axis represents the τ values—performance ratios that compare each method’s metric to the best performance achieved, while y-axis displays the cumulative fraction of test datasets for which a method’s performance is within the τ value. We use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

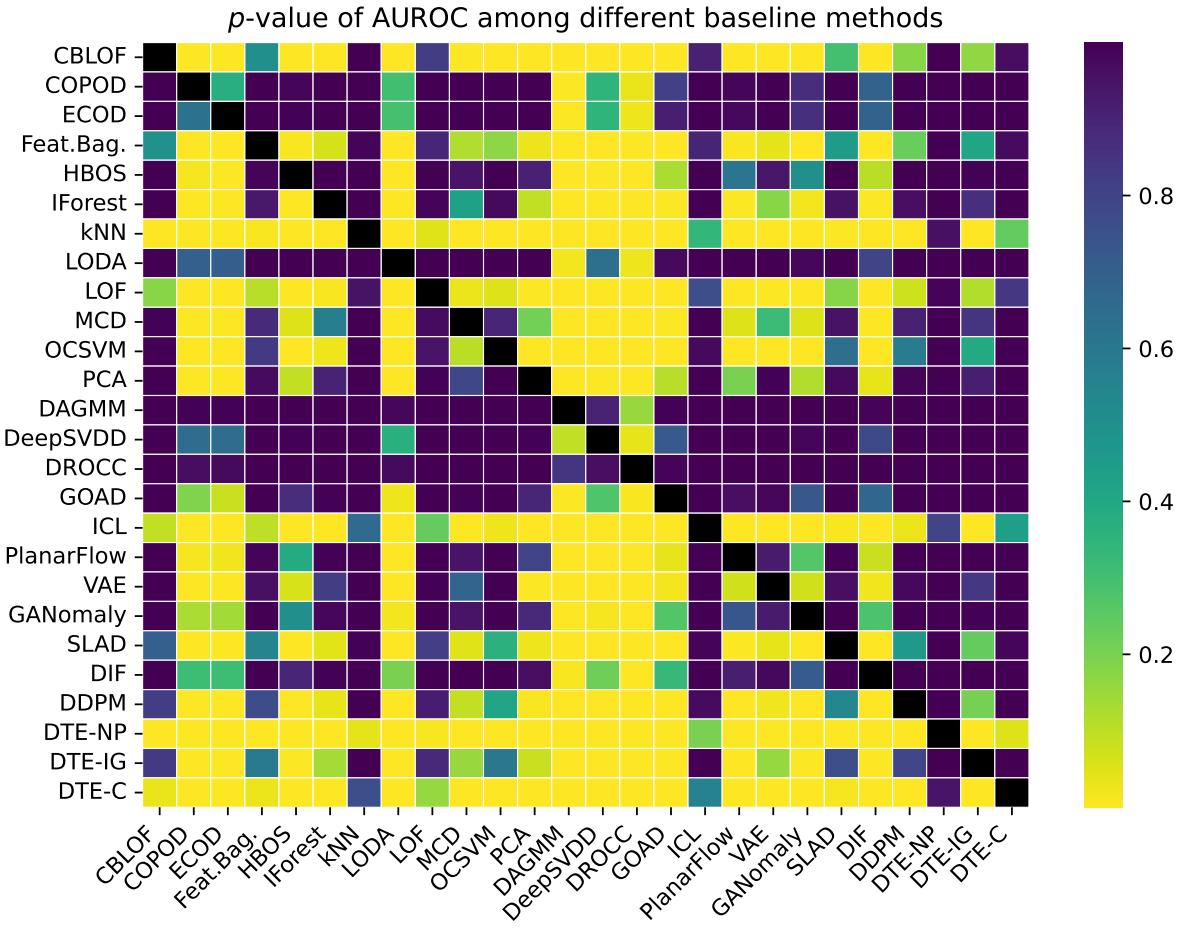


Figure 23: Pairwise p-values among baseline methods based on the Wilcoxon signed rank test w.r.t. AUROC performances across datasets.

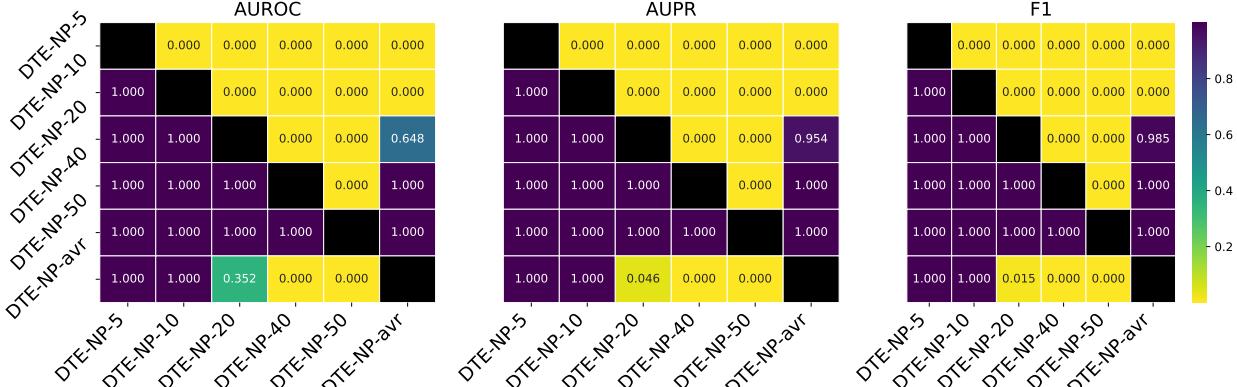


Figure 24: p-values w.r.t. AUROC/AUPR/F1 among different HP configurations of DTE-NP (i.e., $k \in \{5, 10, 20, 40, 50\}$), along with the avg model with the average performance across HPs.

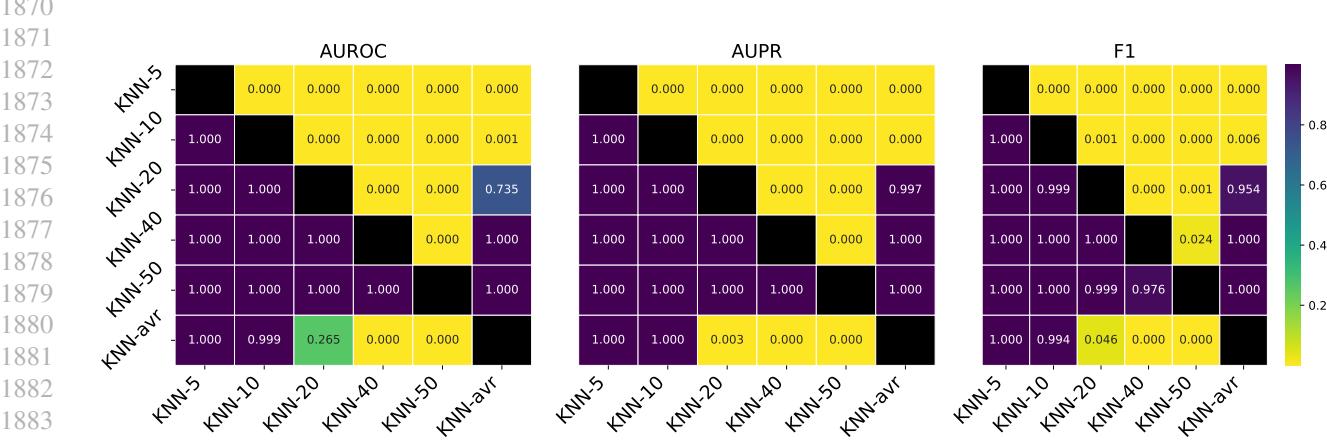


Figure 25: p -values w.r.t. AUROC/AUPR/F1 among different HP configurations of **kNN** (i.e., $k \in \{5, 10, 20, 40, 50\}$), along with the ^{avg} model with the average performance across HPs.

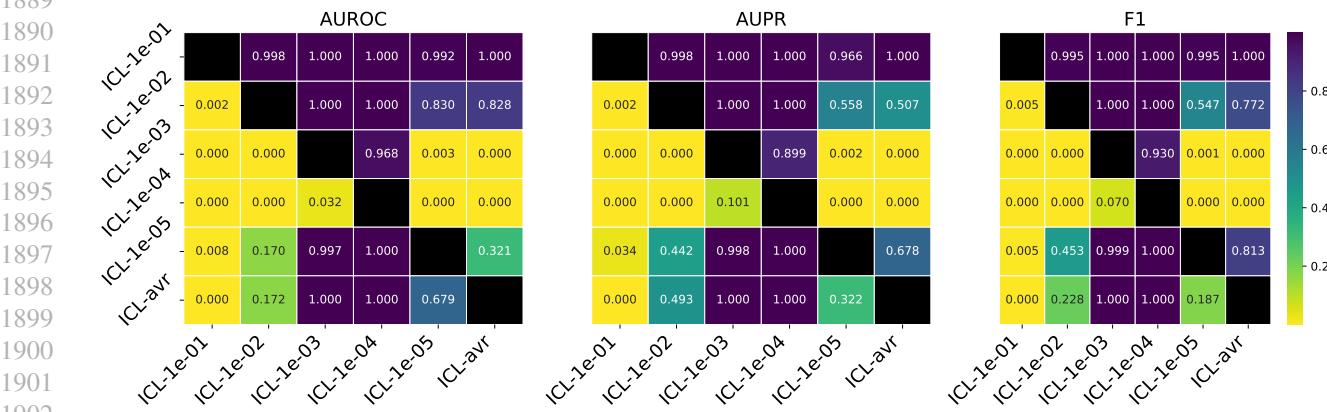


Figure 26: p -values w.r.t. AUROC/AUPR/F1 among different HP configurations of **ICL** (i.e., $\text{learning_rate} \in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$), along with the ^{avg} model with the average performance across HPs.

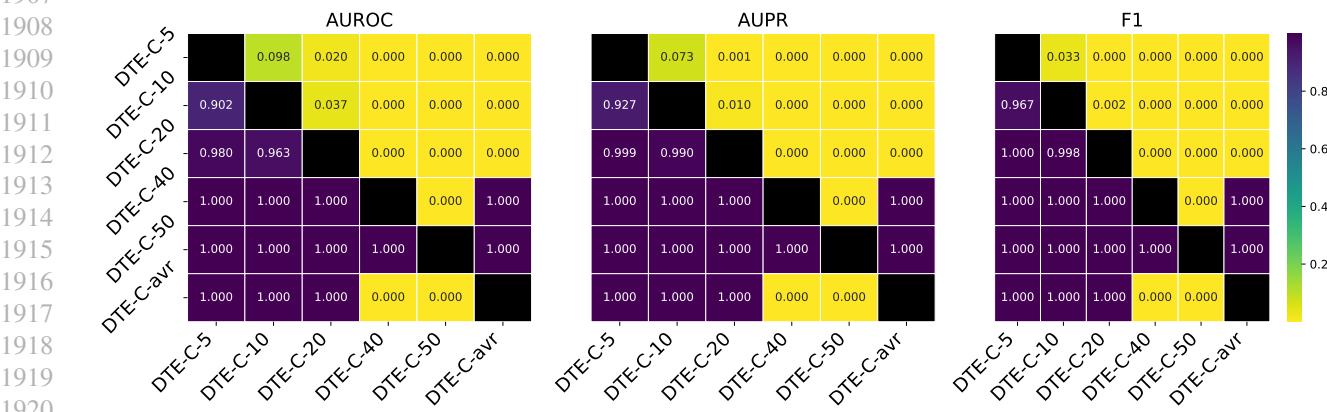


Figure 27: p -values w.r.t. AUROC/AUPR/F1 among different HP configurations of **DTE-C** (i.e., $k \in \{5, 10, 20, 40, 50\}$), along with the ^{avg} model with the average performance across HPs.

Table 13.1: Average AUROC \pm standard dev. over five seeds for the semi-supervised setting of DTE-NP, k NN with varying hyperparameter (HP) values; $k \in \{5, 10, 20, 40, 50\}$. Also reported is the avg model. We use **bold** and underline respectively to mark the **best** and the **worst** performance of each model to showcase the variability of performance across different HP settings.

dataset	DTE-NP-5	DTE-NP-10	DTE-NP-20	DTE-NP-40	DTE-NP-50	DTE-NP-avr	<u>k</u> NN-5	<u>k</u> NN-10	<u>k</u> NN-20	<u>k</u> NN-40	<u>k</u> NN-50	KNN-avr	
abt	50.69 \pm 0.00	51.02 \pm 0.00	51.26 \pm 0.00	51.58 \pm 0.00	51.69 \pm 0.00	51.25 \pm 0.00	51.33 \pm 0.00	51.63 \pm 0.00	51.97 \pm 0.00	52.08 \pm 0.00	51.61 \pm 0.00	51.61 \pm 0.00	
amazon	60.76 \pm 0.00	60.69 \pm 0.00	60.53 \pm 0.00	60.17 \pm 0.00	60.22 \pm 0.00	60.47 \pm 0.00	60.58 \pm 0.00	60.52 \pm 0.00	60.92 \pm 0.00	60.92 \pm 0.00	60.25 \pm 0.00	60.25 \pm 0.00	
android	93.01 \pm 0.00	92.89 \pm 0.00	92.38 \pm 0.00	92.26 \pm 0.00	92.64 \pm 0.00	92.68 \pm 0.45	92.55 \pm 0.44	93.71 \pm 0.46	92.58 \pm 0.46	91.14 \pm 0.46	89.18 \pm 0.47	88.39 \pm 0.51	91.00 \pm 0.47
backdoor	94.48 \pm 0.42	93.72 \pm 0.46	92.67 \pm 0.46	91.20 \pm 0.45	90.68 \pm 0.45	92.55 \pm 0.44	93.71 \pm 0.46	92.58 \pm 0.46	91.14 \pm 0.46	90.91 \pm 0.46	90.21 \pm 0.18	90.21 \pm 0.17	90.16 \pm 0.21
breast	99.10 \pm 0.28	98.91 \pm 0.35	98.59 \pm 0.34	98.66 \pm 0.34	98.36 \pm 0.28	98.67 \pm 0.28	98.71 \pm 0.28	98.99 \pm 0.24	98.11 \pm 0.27	99.16 \pm 0.22	99.21 \pm 0.18	99.21 \pm 0.17	99.16 \pm 0.21
campaign	78.34 \pm 0.00	78.71 \pm 0.00	78.91 \pm 0.00	78.93 \pm 0.00	78.90 \pm 0.00	78.76 \pm 0.00	78.48 \pm 0.00	78.74 \pm 0.00	78.63 \pm 0.00	78.63 \pm 0.00	78.66 \pm 0.00	78.66 \pm 0.00	78.66 \pm 0.00
cardio	91.53 \pm 0.00	92.03 \pm 0.00	92.46 \pm 0.00	93.06 \pm 0.00	93.28 \pm 0.00	92.47 \pm 0.00	92.09 \pm 0.00	92.44 \pm 0.00	92.99 \pm 0.00	93.88 \pm 0.00	94.08 \pm 0.00	93.07 \pm 0.00	93.07 \pm 0.00
cardiography	60.40 \pm 0.00	61.65 \pm 0.00	72.58 \pm 0.26	74.81 \pm 0.34	76.87 \pm 0.38	77.47 \pm 0.37	74.42 \pm 0.38	72.91 \pm 0.29	75.24 \pm 0.40	77.30 \pm 0.47	79.14 \pm 0.38	79.68 \pm 0.37	76.90 \pm 0.35
celeba	70.39 \pm 0.33	71.18 \pm 0.16	72.34 \pm 0.14	72.28 \pm 0.16	71.89 \pm 0.17	71.37 \pm 0.29	71.36 \pm 0.12	71.36 \pm 0.12	71.28 \pm 0.16	71.84 \pm 0.15	71.28 \pm 0.16	71.28 \pm 0.16	71.28 \pm 0.16
census	cover	97.72 \pm 0.17	97.72 \pm 0.14	97.40 \pm 0.18	96.99 \pm 0.23	96.84 \pm 0.24	97.37 \pm 0.19	97.51 \pm 0.15	97.19 \pm 0.15	96.75 \pm 0.22	96.21 \pm 0.28	96.00 \pm 0.31	96.73 \pm 0.22
donors	99.72 \pm 0.03	99.61 \pm 0.03	99.43 \pm 0.06	99.14 \pm 0.09	99.02 \pm 0.10	99.38 \pm 0.06	99.51 \pm 0.06	99.24 \pm 0.08	98.85 \pm 0.10	98.20 \pm 0.13	97.90 \pm 0.14	98.74 \pm 0.09	98.74 \pm 0.09
fault	58.34 \pm 0.00	58.57 \pm 0.00	60.02 \pm 0.00	60.00 \pm 0.00	59.17 \pm 0.00	59.64 \pm 0.02	59.53 \pm 0.09	59.53 \pm 0.09	59.54 \pm 0.09	59.54 \pm 0.09	59.57 \pm 0.03	60.22 \pm 0.00	60.22 \pm 0.00
fraud	95.08 \pm 0.90	95.67 \pm 0.93	95.64 \pm 0.93	95.60 \pm 0.92	95.60 \pm 0.92	95.81 \pm 0.92	95.39 \pm 0.97	95.39 \pm 0.97	95.86 \pm 0.97	95.86 \pm 0.97	95.54 \pm 0.92	95.54 \pm 0.92	95.57 \pm 0.93
glass	96.08 \pm 0.39	93.04 \pm 1.06	89.82 \pm 1.12	87.89 \pm 1.10	87.31 \pm 1.40	89.83 \pm 0.91	92.13 \pm 0.94	88.67 \pm 0.98	87.21 \pm 0.98	88.94 \pm 0.92	88.94 \pm 0.92	88.94 \pm 0.92	88.94 \pm 0.92
hepatitis	99.84 \pm 0.20	99.27 \pm 0.51	96.89 \pm 0.96	93.15 \pm 1.69	91.97 \pm 1.76	96.22 \pm 0.88	96.77 \pm 1.47	86.88 \pm 2.21	85.50 \pm 2.34	85.46 \pm 1.92	84.88 \pm 2.21	84.88 \pm 2.09	87.90 \pm 1.75
http	99.99 \pm 0.00	99.98 \pm 0.01	99.95 \pm 0.00	99.93 \pm 0.01	99.93 \pm 0.02	99.95 \pm 0.01	99.95 \pm 0.00	99.99 \pm 0.02	99.95 \pm 0.01	99.95 \pm 0.01	99.95 \pm 0.01	99.95 \pm 0.01	99.96 \pm 0.01
imdb	50.48 \pm 0.00	50.38 \pm 0.00	50.72 \pm 0.00	50.28 \pm 0.00	50.22 \pm 0.00	50.28 \pm 0.00	50.35 \pm 0.00	50.08 \pm 0.00	50.23 \pm 0.00	50.23 \pm 0.00	50.18 \pm 0.00	50.18 \pm 0.00	50.18 \pm 0.00
intemrads	70.96 \pm 0.00	68.65 \pm 0.00	66.68 \pm 0.00	65.97 \pm 0.00	65.82 \pm 0.00	67.65 \pm 0.00	68.08 \pm 0.00	65.48 \pm 0.00	65.02 \pm 0.00	65.04 \pm 0.00	65.04 \pm 0.00	65.73 \pm 0.00	65.73 \pm 0.00
ionosphere	98.48 \pm 0.60	98.13 \pm 0.74	97.84 \pm 0.64	96.83 \pm 0.79	97.50 \pm 0.79	97.50 \pm 0.79	97.62 \pm 0.81	96.33 \pm 0.76	97.32 \pm 0.85	97.80 \pm 0.76	97.80 \pm 0.76	95.12 \pm 0.92	95.12 \pm 0.92
landsat	68.99 \pm 0.00	68.02 \pm 0.00	66.46 \pm 0.00	64.73 \pm 0.00	64.16 \pm 0.00	66.47 \pm 0.00	68.25 \pm 0.00	66.48 \pm 0.00	64.36 \pm 0.00	62.49 \pm 0.00	61.92 \pm 0.00	60.70 \pm 0.00	60.70 \pm 0.00
letter	36.12 \pm 0.00	35.66 \pm 0.00	34.78 \pm 0.00	33.48 \pm 0.00	33.72 \pm 0.00	34.74 \pm 0.00	35.43 \pm 0.00	34.54 \pm 0.00	33.71 \pm 0.00	32.11 \pm 0.00	31.69 \pm 0.00	33.39 \pm 0.00	33.39 \pm 0.00
lymphography	83.98 \pm 0.25	99.79 \pm 0.32	99.79 \pm 0.32	99.76 \pm 0.31	99.76 \pm 0.31	99.76 \pm 0.31	99.80 \pm 0.30	99.87 \pm 0.10	99.85 \pm 0.05	99.88 \pm 0.08	99.88 \pm 0.08	99.88 \pm 0.08	99.88 \pm 0.08
magic gamma	83.91 \pm 0.00	83.49 \pm 0.00	82.87 \pm 0.00	82.05 \pm 0.00	81.73 \pm 0.00	82.81 \pm 0.00	83.27 \pm 0.00	82.64 \pm 0.00	81.85 \pm 0.00	80.76 \pm 0.00	80.76 \pm 0.00	81.76 \pm 0.00	81.76 \pm 0.00
mammography	87.65 \pm 0.00	87.75 \pm 0.00	87.68 \pm 0.00	87.42 \pm 0.00	87.29 \pm 0.00	87.55 \pm 0.00	87.58 \pm 0.00	87.75 \pm 0.00	87.38 \pm 0.00	86.97 \pm 0.00	86.78 \pm 0.00	87.29 \pm 0.00	87.29 \pm 0.00
mnist	93.93 \pm 0.00	93.57 \pm 0.00	93.20 \pm 0.00	93.08 \pm 0.00	93.00 \pm 0.00	93.60 \pm 0.00	93.85 \pm 0.00	93.45 \pm 0.00	92.55 \pm 0.00	92.55 \pm 0.00	92.36 \pm 0.00	93.04 \pm 0.00	93.04 \pm 0.00
musik	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00
opdigits	95.00 \pm 0.00	93.97 \pm 0.00	92.22 \pm 0.00	90.57 \pm 0.00	90.28 \pm 0.00	92.53 \pm 0.00	93.72 \pm 0.00	92.99 \pm 0.00	92.99 \pm 0.00	92.99 \pm 0.00	92.66 \pm 0.00	90.67 \pm 0.00	90.67 \pm 0.00
pageblocks	89.04 \pm 0.00	89.40 \pm 0.00	89.56 \pm 0.00	89.37 \pm 0.00	89.23 \pm 0.00	89.32 \pm 0.00	89.65 \pm 0.00	89.86 \pm 0.00	89.88 \pm 0.00	89.30 \pm 0.00	89.18 \pm 0.00	89.34 \pm 0.00	89.34 \pm 0.00
pendigits	99.90 \pm 0.00	99.88 \pm 0.00	99.83 \pm 0.00	99.51 \pm 0.00	99.58 \pm 0.00	99.70 \pm 0.00	99.70\pm0.00	99.70 \pm 0.00	99.70 \pm 0.00				
pima	82.21 \pm 1.82	79.74 \pm 1.61	77.98 \pm 1.38	77.29 \pm 1.35	77.14 \pm 1.33	78.87 \pm 1.43	77.44\pm2.07	76.14 \pm 1.36	76.39 \pm 1.21	76.55 \pm 1.25	76.58 \pm 1.29	76.58 \pm 1.29	76.58 \pm 1.29
satellite	82.40 \pm 0.00	82.09 \pm 0.00	81.55 \pm 0.00	80.71 \pm 0.00	80.38 \pm 0.00	80.38 \pm 0.00	80.43 \pm 0.00	80.24 \pm 0.00	81.56 \pm 0.00	80.56 \pm 0.00	80.76 \pm 0.00	80.47 \pm 0.00	80.47 \pm 0.00
satimage-2	99.68 \pm 0.00	99.73 \pm 0.00	99.92 \pm 0.00	99.91 \pm 0.00	99.91 \pm 0.00	99.92 \pm 0.00	99.92 \pm 0.00	99.90 \pm 0.00	99.90 \pm 0.00	99.90 \pm 0.00	99.77 \pm 0.00	99.77 \pm 0.00	99.77 \pm 0.00
skim	99.66 \pm 0.05	99.52 \pm 0.04	98.72 \pm 0.04	98.57 \pm 0.04	98.54 \pm 0.05	99.15 \pm 0.05	99.51 \pm 0.06	97.43 \pm 0.05	97.43 \pm 0.05	96.90 \pm 0.09	99.89 \pm 0.00	99.89 \pm 0.00	99.89 \pm 0.00
smnlp	41.12 \pm 0.00	39.09 \pm 0.00	37.59 \pm 0.00	37.37 \pm 0.00	37.01 \pm 0.00	38.42 \pm 0.00	37.56 \pm 0.00	36.74 \pm 0.00	36.15 \pm 0.00	36.25 \pm 0.00	36.31 \pm 0.00	36.31 \pm 0.00	36.31 \pm 0.00
stamps	97.88 \pm 0.33	98.63 \pm 0.00	98.64 \pm 0.00	98.59 \pm 0.00	98.61 \pm 0.00	98.59 \pm 0.00	98.61 \pm 0.00	98.61 \pm 0.00	98.61 \pm 0.00	98.61 \pm 0.00	98.68 \pm 0.00	98.68 \pm 0.00	98.68 \pm 0.00
thyroid	98.38 \pm 0.00	98.63 \pm 0.14	99.33 \pm 0.18	98.88 \pm 0.19	98.57 \pm 0.30	99.56 \pm 0.00	99.28 \pm 0.12	98.67 \pm 0.00	98.67 \pm 0.00	98.67 \pm 0.00	98.69 \pm 0.00	98.69 \pm 0.00	98.69 \pm 0.00
vertebral	79.59 \pm 2.23	81.05 \pm 2.09	80.59 \pm 2.07	80.07 \pm 1.60	79.88 \pm 1.43	82.21 \pm 1.43	82.73 \pm 3.80	84.40 \pm 1.87	84.40 \pm 1.87				
vowels	82.11 \pm 0.												

A Foundation Model for Zero-shot Outlier Detection

Table 13.2: Average AUROC \pm standard dev. over five seeds for the semi-supervised setting of ICL and DTE-C baselines with varying hyperparameter (HP) values. For ICL, the learning rate $\in \{0.1, 0.02, 0.001, 0.0001, 1e-05\}$, for DTE-C, $k \in \{5, 10, 20, 40, 50\}$. Also reported is the avg model. We use **bold** and underline respectively to mark the **best** and the **worst** performance of each model to showcase the variability of performance across different HP settings.

Table 14.1: Average AUPR \pm standard dev. over five seeds for the semi-supervised setting of DTE-NP, k NN baselines with varying hyperparameter (HP) values; $k \in \{5, 10, 20, 40, 50\}$. Also reported is the avg model. We use **bold** and underline respectively to mark the **best** and the worst performance of each model to showcase the variability of performance across different HP settings.

dataset	DTE-NP-5	DTE-NP-10	DTE-NP-20	DTE-NP-40	DTE-NP-50	DTE-NP-avr	<u>k</u> NN-5	<u>k</u> NN-10	<u>k</u> NN-20	<u>k</u> NN-40	<u>k</u> NN-50	KNN-avr
abt	5.95 \pm 0.00	5.99 \pm 0.00	6.02 \pm 0.00	6.06 \pm 0.00	6.07 \pm 0.00	6.02 \pm 0.00	6.02 \pm 0.00	6.07 \pm 0.00	6.09 \pm 0.00	6.13 \pm 0.00	6.15 \pm 0.00	6.09 \pm 0.00
amazon	11.68 \pm 0.00	11.68 \pm 0.00	11.68 \pm 0.00	11.61 \pm 0.00	11.62 \pm 0.00	11.65 \pm 0.00	11.65 \pm 0.00	11.70 \pm 0.00	11.66 \pm 0.00	11.66 \pm 0.00	11.59 \pm 0.00	11.65 \pm 0.00
android	67.49 \pm 0.00	66.73 \pm 0.00	66.04 \pm 0.00	65.39 \pm 0.00	66.11 \pm 0.00	66.75\pm0.00	67.20 \pm 0.00	67.30 \pm 0.00	67.39 \pm 0.00	67.39 \pm 0.00	67.39 \pm 0.00	67.06 \pm 0.00
backdoor	55.51 \pm 0.56	98.19 \pm 0.58	97.69 \pm 0.51	47.16 \pm 1.45	38.31 \pm 1.02	31.44 \pm 0.47	29.58 \pm 0.37	40.48 \pm 0.81	46.70 \pm 1.22	29.38 \pm 1.35	24.34 \pm 0.41	22.34 \pm 0.53
breast	98.51 \pm 0.00	98.10 \pm 0.00	98.05 \pm 0.00	49.05 \pm 0.00	49.77 \pm 0.00	49.77 \pm 0.00	49.55 \pm 0.45	97.13 \pm 0.62	97.05 \pm 0.40	97.69 \pm 0.28	99.15 \pm 0.31	99.08 \pm 0.22
campaign	48.48 \pm 0.00	49.05 \pm 0.00	49.77 \pm 0.00	49.77 \pm 0.00	49.51 \pm 0.00	49.31 \pm 0.00	49.89 \pm 0.00	49.89 \pm 0.00	50.45 \pm 0.00	49.47 \pm 0.00	49.33 \pm 0.00	49.64 \pm 0.00
cardio	76.90 \pm 0.00	77.73 \pm 0.00	78.30 \pm 0.00	79.19 \pm 0.00	79.53 \pm 0.00	78.33 \pm 0.00	77.22 \pm 0.00	78.33 \pm 0.00	79.14 \pm 0.00	80.67 \pm 0.00	81.15 \pm 0.00	79.30 \pm 0.00
cardiography	56.55 \pm 0.00	57.18 \pm 0.00	59.42 \pm 0.00	59.59 \pm 0.00	59.59 \pm 0.00	58.26 \pm 0.00	57.23 \pm 0.00	58.37 \pm 0.00	59.44 \pm 0.00	61.41 \pm 0.00	62.19 \pm 0.00	59.77 \pm 0.00
celeba	10.36 \pm 0.44	11.63 \pm 0.49	12.74 \pm 0.52	13.92 \pm 0.59	14.30 \pm 0.59	12.63 \pm 0.51	11.99 \pm 0.57	13.26 \pm 0.61	14.30 \pm 0.58	15.70 \pm 0.65	16.10 \pm 0.68	14.31 \pm 0.60
census	21.14 \pm 0.39	21.38 \pm 0.54	21.16 \pm 0.43	20.67 \pm 0.41	20.52 \pm 0.42	20.97 \pm 0.43	21.61 \pm 0.39	21.61 \pm 0.39	20.00 \pm 0.42	19.94 \pm 0.44	20.62 \pm 0.44	19.94 \pm 0.44
cover	63.67 \pm 3.21	57.58 \pm 3.52	51.55 \pm 3.10	44.88 \pm 2.49	42.11 \pm 2.29	51.95 \pm 2.90	55.15 \pm 3.45	48.67 \pm 2.84	41.44 \pm 2.04	33.72 \pm 1.51	31.69 \pm 1.35	42.14 \pm 2.20
donors	93.33 \pm 0.80	91.25 \pm 0.72	88.17 \pm 0.95	83.92 \pm 1.29	83.34 \pm 1.32	87.78 \pm 0.99	89.44 \pm 0.96	85.33 \pm 1.15	80.15 \pm 1.33	73.68 \pm 1.32	71.00 \pm 1.30	79.92 \pm 1.13
fault	62.03 \pm 0.00	61.58 \pm 0.00	62.13 \pm 0.00	61.84 \pm 0.00	61.98 \pm 0.00	61.92 \pm 0.00	63.67 \pm 0.00	63.67 \pm 0.00	64.06 \pm 0.00	62.56 \pm 0.00	64.06 \pm 0.00	62.56 \pm 0.00
fraud	40.60 \pm 0.67	43.77 \pm 5.38	43.03 \pm 4.92	39.91 \pm 4.75	38.80 \pm 4.94	41.22 \pm 5.02	42.35 \pm 5.61	44.96 \pm 3.90	41.19 \pm 3.73	37.33 \pm 4.15	36.42 \pm 4.12	40.45 \pm 4.07
hepatitis	60.15 \pm 6.89	47.75 \pm 5.62	37.27 \pm 4.92	31.23 \pm 3.12	30.48 \pm 2.85	41.38 \pm 4.46	44.07 \pm 6.38	32.96 \pm 3.74	26.62 \pm 2.65	26.04 \pm 3.61	31.91 \pm 5.73	26.04 \pm 3.61
http	98.52 \pm 0.37	95.38 \pm 2.26	88.66 \pm 1.01	84.43 \pm 2.65	80.40 \pm 4.55	89.48 \pm 1.90	91.10 \pm 4.41	69.28 \pm 4.16	64.12 \pm 5.59	64.29 \pm 5.08	64.33 \pm 6.16	70.62 \pm 4.48
imdb	9.11 \pm 0.00	9.06 \pm 0.00	9.06 \pm 0.00	9.06 \pm 0.00	9.06 \pm 0.00	9.06 \pm 0.00	9.06 \pm 0.00	9.06 \pm 0.00	9.06 \pm 0.00	9.06 \pm 0.00	9.06 \pm 0.00	9.06 \pm 0.00
intertials	52.20 \pm 0.00	49.76 \pm 0.00	48.19 \pm 0.00	47.56 \pm 0.00	47.45 \pm 0.00	49.03 \pm 0.00	49.22 \pm 0.00	47.29 \pm 0.00	46.93 \pm 0.00	46.93 \pm 0.00	47.47 \pm 0.00	46.93 \pm 0.00
ionosphere	98.72 \pm 0.48	98.46 \pm 0.54	98.27 \pm 0.42	97.44 \pm 0.50	96.93 \pm 0.61	97.96 \pm 0.46	97.86 \pm 0.60	98.11 \pm 0.52	97.04 \pm 0.60	94.12 \pm 1.45	92.59 \pm 1.65	96.01 \pm 0.78
landsat	56.14 \pm 0.00	54.25 \pm 0.00	50.75 \pm 0.00	46.43 \pm 0.00	45.17 \pm 0.00	50.45 \pm 0.00	50.62 \pm 0.00	45.18 \pm 0.00	41.32 \pm 0.00	40.50 \pm 0.00	40.50 \pm 0.00	46.49 \pm 0.00
letter	8.86 \pm 0.00	8.78 \pm 0.00	8.67 \pm 0.00	8.54 \pm 0.00	8.50 \pm 0.00	8.67 \pm 0.00	8.70 \pm 0.00	8.58 \pm 0.00	8.41 \pm 0.00	8.27 \pm 0.00	8.22 \pm 0.00	8.44 \pm 0.00
lymphography	97.27 \pm 5.45	96.07 \pm 6.79	96.07 \pm 6.59	95.68 \pm 6.60	95.68 \pm 6.43	98.61 \pm 1.02	98.45 \pm 0.52	98.45 \pm 0.52	98.45 \pm 0.52	98.70 \pm 0.83	98.57 \pm 0.65	98.57 \pm 0.65
magic gamma	86.40 \pm 0.00	85.86 \pm 0.00	85.28 \pm 0.00	84.56 \pm 0.00	84.29 \pm 0.00	85.26 \pm 0.00	85.86 \pm 0.00	85.25 \pm 0.00	84.51 \pm 0.00	83.61 \pm 0.00	83.61 \pm 0.00	84.50 \pm 0.00
mammography	42.14 \pm 0.00	41.51 \pm 0.00	40.67 \pm 0.00	40.37 \pm 0.00	40.50 \pm 0.00	41.04 \pm 0.00	41.27 \pm 0.00	40.53 \pm 0.00	41.24 \pm 0.00	38.97 \pm 0.00	38.10 \pm 0.00	39.83 \pm 0.00
mnist	74.43 \pm 0.00	73.09 \pm 0.00	71.84 \pm 0.00	70.69 \pm 0.00	70.36 \pm 0.00	70.08 \pm 0.00	72.08 \pm 0.00	71.40 \pm 0.00	70.09 \pm 0.00	69.03 \pm 0.00	68.60 \pm 0.00	70.36 \pm 0.00
musik	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00
opdigits	34.44 \pm 0.00	30.53 \pm 0.00	26.28 \pm 0.00	22.67 \pm 0.00	21.61 \pm 0.00	27.11 \pm 0.00	29.11 \pm 0.00	26.76 \pm 0.00	24.76 \pm 0.00	21.10 \pm 0.00	17.68 \pm 0.00	16.62 \pm 0.00
pageblocks	62.78 \pm 0.00	62.52 \pm 0.00	62.02 \pm 0.00	61.02 \pm 0.00	60.30 \pm 0.00	67.60 \pm 0.00	67.60 \pm 0.00	67.74 \pm 0.00	67.74 \pm 0.00	67.87 \pm 0.00	67.87 \pm 0.00	67.88 \pm 0.00
pendigits	97.68 \pm 0.00	97.31 \pm 0.00	96.28 \pm 0.00	96.01 \pm 0.00	95.99 \pm 0.00	93.59 \pm 0.00	90.01 \pm 0.00	87.70 \pm 0.00	87.70 \pm 0.00	87.24 \pm 0.00	87.39 \pm 0.00	87.34 \pm 0.00
satellite	85.98 \pm 0.00	85.74 \pm 0.00	85.17 \pm 0.00	84.15 \pm 0.00	83.72 \pm 0.00	84.95 \pm 0.00	86.01 \pm 0.00	85.31 \pm 0.00	84.02 \pm 0.00	82.15 \pm 0.00	81.56 \pm 0.00	83.82 \pm 0.00
satimage-2	99.16 \pm 0.00	96.64 \pm 0.00	97.02 \pm 0.00	97.39 \pm 0.00	97.42 \pm 0.00	96.92 \pm 0.00	96.92 \pm 0.00	97.21 \pm 0.00	97.33 \pm 0.00	97.39 \pm 0.00	97.42 \pm 0.00	97.22 \pm 0.00
shuttle	98.92 \pm 0.23	98.76 \pm 0.00	98.72 \pm 0.00	98.78 \pm 0.00	98.77 \pm 0.00	98.84 \pm 0.00	97.96 \pm 0.00	97.36 \pm 0.00	97.28 \pm 0.00	97.22 \pm 0.00	97.20 \pm 0.00	97.38 \pm 0.00
skin	56.70 \pm 7.16	54.77 \pm 7.80	54.75 \pm 6.59	54.76 \pm 7.81	54.76 \pm 7.81	54.76 \pm 7.83	50.26 \pm 5.73	50.26 \pm 5.74	50.18 \pm 5.75	50.33 \pm 5.85	50.41 \pm 5.76	50.27 \pm 5.76
spambase	83.93 \pm 0.00	83.42 \pm 0.00	83.03 \pm 0.00	82.73 \pm 0.00	82.63 \pm 0.00	83.31 \pm 0.00	83.32 \pm 0.00	82.70 \pm 0.00	82.41 \pm 0.00	82.17 \pm 0.00	82.11 \pm 0.00	82.54 \pm 0.00
speech	3.92 \pm 0.00	2.99 \pm 0.00	2.70 \pm 0.00	2.07 \pm 0.00	2.07 \pm 0.00	2.80 \pm 0.00	2.70 \pm 0.00	2.70 \pm 0.00	2.70 \pm 0.00	2.70 \pm 0.00	2.70 \pm 0.00	2.75 \pm 0.00
stamps	82.30 \pm 3.70	77.11 \pm 4.30	69.99 \pm 5.29	65.83 \pm 6.16	64.64 \pm 6.16	72.02 \pm 4.82	73.26 \pm 7.70	65.58 \pm 7.71	63.12 \pm 6.69	62.09 \pm 7.20	61.57 \pm 7.18	65.12 \pm 7.10
thyroid	43.77 \pm 0.00	77.53 \pm 0.00	24.99 \pm 2.97	21.10 \pm 2.37	20.29 \pm 2.40	28.38 \pm 3.31	57.07 \pm 3.19	21.55 \pm 2.53	19.65 \pm 2.31	18.08 \pm 1.91	18.08 \pm 1.91	18.08 \pm 1.91
vertebral	75.30 \pm 5.50	61.19 \pm 1.49	51.53 \pm 2.17	46.62 \pm 2.23	45.63 \pm 2.17	56.05 \pm 1.44	47.07 \pm 2.57	43.16 \pm 2.41	42.33 \pm 2.53	42.28 \pm 2.60	42.11 \pm 2.34	42.24 \pm 2.34
yeast	48.37 \pm 0.00	47.91 \pm 0.00	47.52 \pm 0.00	47.26 \pm 0.00	47.20 \pm 0.00	47.65 \pm 0.00	48.26 \pm 0.00	47.48 \pm 0.00	47.24 \pm 0.00	46.74 \pm 0.00	46.48 \pm 0.00	47.24 \pm 0.00
yelp	16.05 \pm 0.00	15.78 \pm 0.00	15.40 \pm 0.00	15.01 \pm 0.00	14.89 \pm 0.00	15.43 \pm 0.00	16.03 \pm 0.00	15.63 \pm 0.00	15.17 \pm 0.00	14.77 \pm 0.00	14.77	

Table 14.2: Average AUPR \pm standard dev. over five seeds for the semi-supervised setting of ICL and DTE-C baselines with varying hyperparameter (HP) values; For ICL, the learning rate $\in \{0.1, 0.02, 0.001, 0.0001, 1e-05\}$, for DTE-C, $k \in \{5, 10, 20, 40, 50\}$. Also reported is the avg model. We use **bold** and underline respectively to mark the **best** and the **worst** performance of each model to showcase the variability of performance across different HP settings.

2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199

Table 15.1: Average F1 score \pm standard dev. over five seeds for the semi-supervised setting of DTE-NP, k NN baselines with varying hyperparameter (HP) values; $k \in \{5, 10, 20, 40, 50\}$. Also reported is the avg model. We use **bold** and underline respectively to mark the **best** and the worst performance of each model to showcase the variability of performance across different HP settings.

dataset	DTE-NP-5	DTE-NP-10	DTE-NP-20	DTE-NP-40	DTE-NP-50	DTE-NP-avr	<u>k</u> NN-5	<u>k</u> NN-10	<u>k</u> NN-20	<u>k</u> NN-40	<u>k</u> NN-50	KNN-avr
abt	5.90 \pm 0.00	5.84 \pm 0.00	5.70 \pm 0.00	5.90 \pm 0.00	<u>5.97</u> \pm 0.00	5.86 \pm 0.00	5.90 \pm 0.00	<u>5.64</u> \pm 0.00	5.97 \pm 0.00	6.17 \pm 0.00	<u>6.37</u> \pm 0.00	6.0 \pm 0.00
amazon	10.80 \pm 0.00	10.80 \pm 0.00	10.20 \pm 0.00	10.20 \pm 0.00	<u>11.00</u> \pm 0.00	10.60 \pm 0.00	<u>11.40</u> \pm 0.00	10.20 \pm 0.00	11.20 \pm 0.00	11.20 \pm 0.00	10.92 \pm 0.00	10.92 \pm 0.00
android	<u>62.55</u> \pm 0.00	61.80 \pm 0.00	58.99 \pm 0.00	60.56 \pm 0.00	58.80 \pm 0.00	60.49 \pm 0.00	58.43 \pm 0.00	58.43 \pm 0.00	56.14 \pm 0.00	56.14 \pm 0.00	59.18 \pm 0.00	59.18 \pm 0.00
backdoor	64.15 \pm 1.04	52.30 \pm 1.87	40.62 \pm 1.46	30.25 \pm 1.34	26.96 \pm 1.20	42.86 \pm 1.32	53.53 \pm 1.63	40.37 \pm 2.04	28.71 \pm 1.50	20.21 \pm 0.78	17.52 \pm 1.83	31.87 \pm 1.22
breast	96.72 \pm 0.64	96.73 \pm 0.39	96.73 \pm 0.47	96.73 \pm 0.39	96.73 \pm 0.39	96.73 \pm 0.43	96.73 \pm 0.44	96.05 \pm 0.33	95.99 \pm 0.39	95.99 \pm 0.39	95.93 \pm 0.32	95.97 \pm 0.31
campaign	49.94 \pm 0.00	50.62 \pm 0.00	51.14 \pm 0.00	51.38 \pm 0.00	51.70 \pm 0.00	50.93 \pm 0.00	50.37 \pm 0.00	50.86 \pm 0.00	51.27 \pm 0.00	51.70 \pm 0.00	51.29 \pm 0.00	51.10 \pm 0.00
cardio	63.64 \pm 0.00	61.36 \pm 0.00	61.93 \pm 0.00	63.64 \pm 0.00	64.20 \pm 0.00	62.95 \pm 0.00	61.93 \pm 0.00	64.78 \pm 0.00	67.61 \pm 0.00	69.32 \pm 0.00	65.00 \pm 0.00	65.00 \pm 0.00
cardiography	44.64 \pm 0.00	45.71 \pm 0.00	47.18 \pm 0.00	48.75 \pm 0.00	49.14 \pm 0.00	46.85 \pm 0.00	47.00 \pm 0.00	47.85 \pm 0.00	50.86 \pm 0.00	48.76 \pm 0.00	49.13 \pm 0.00	48.76 \pm 0.00
celeba	15.83 \pm 0.69	17.05 \pm 0.43	18.17 \pm 0.61	19.02 \pm 0.69	19.30 \pm 0.60	17.87 \pm 0.57	17.08 \pm 0.58	18.41 \pm 0.65	19.30 \pm 0.81	21.27 \pm 0.68	20.48 \pm 0.68	19.11 \pm 0.61
census	22.22 \pm 0.54	21.93 \pm 0.52	21.62 \pm 0.25	21.38 \pm 0.48	21.12 \pm 0.29	21.46 \pm 0.40	21.48 \pm 0.40	21.47 \pm 0.65	21.33 \pm 0.65	21.26 \pm 0.50	21.55 \pm 0.24	21.55 \pm 0.24
cover	69.15 \pm 2.12	66.87 \pm 2.07	63.15 \pm 2.35	63.15 \pm 2.29	55.99 \pm 2.08	53.06 \pm 2.23	61.65 \pm 2.14	65.04 \pm 1.92	65.04 \pm 2.04	52.67 \pm 1.83	42.68 \pm 1.94	39.92 \pm 1.99
donors	97.27 \pm 0.36	96.20 \pm 0.45	94.49 \pm 0.55	91.70 \pm 0.90	90.57 \pm 0.86	94.05 \pm 0.60	94.98 \pm 0.62	92.36 \pm 0.59	88.71 \pm 0.99	80.36 \pm 1.90	76.62 \pm 1.50	86.60 \pm 0.98
fault	56.02 \pm 0.00	55.72 \pm 0.00	56.19 \pm 0.00	57.36 \pm 0.00	56.32 \pm 0.00	55.52 \pm 0.00	57.36 \pm 0.00	55.87 \pm 0.00	57.36 \pm 0.00	58.25 \pm 0.00	56.85 \pm 0.00	56.85 \pm 0.00
fraud	48.18 \pm 4.00	49.60 \pm 3.28	49.00 \pm 3.95	46.66 \pm 3.52	45.46 \pm 3.60	47.78 \pm 3.53	47.78 \pm 3.09	49.39 \pm 3.24	46.64 \pm 4.18	42.58 \pm 3.16	41.64 \pm 3.36	45.61 \pm 3.48
glass	47.81 \pm 5.78	35.14 \pm 2.75	27.98 \pm 4.21	18.37 \pm 2.40	17.81 \pm 2.98	29.42 \pm 3.61	27.19 \pm 3.87	27.19 \pm 3.62	27.19 \pm 3.73	17.23 \pm 3.73	21.09 \pm 5.16	21.09 \pm 5.16
hepatitis	98.94 \pm 1.41	94.16 \pm 2.13	81.68 \pm 4.18	75.95 \pm 4.78	71.99 \pm 4.14	84.54 \pm 2.23	81.98 \pm 4.50	66.04 \pm 4.53	62.31 \pm 6.09	60.39 \pm 6.21	60.04 \pm 7.30	66.15 \pm 7.30
http	98.50 \pm 0.38	95.10 \pm 2.50	88.26 \pm 1.38	82.85 \pm 3.80	82.85 \pm 2.54	88.73 \pm 2.54	90.00 \pm 0.00	98.57 \pm 2.86	92.67 \pm 0.91	92.67 \pm 0.91	93.32 \pm 0.75	93.32 \pm 0.75
imdb	<u>5.20</u> \pm 0.00	5.40 \pm 0.00	5.16 \pm 0.00	5.20 \pm 0.00	5.28 \pm 0.00	5.40 \pm 0.00	5.24 \pm 0.00	5.40 \pm 0.00	5.00 \pm 0.00	5.00 \pm 0.00	5.24 \pm 0.00	5.24 \pm 0.00
intemrads	<u>51.16</u> \pm 0.00	51.63 \pm 0.00	48.87 \pm 0.00	46.47 \pm 0.00	46.20 \pm 0.00	49.57 \pm 0.00	51.90 \pm 0.00	46.20 \pm 0.00	45.11 \pm 0.00	45.11 \pm 0.00	46.68 \pm 0.00	46.68 \pm 0.00
ionosphere	92.33 \pm 1.17	92.05 \pm 1.63	91.63 \pm 1.09	90.41 \pm 1.35	89.19 \pm 1.72	91.12 \pm 1.24	91.81 \pm 1.87	90.23 \pm 1.86	87.12 \pm 1.86	87.12 \pm 1.86	87.86 \pm 1.86	87.86 \pm 1.86
landsat	52.29 \pm 0.00	51.24 \pm 0.00	49.03 \pm 0.00	45.99 \pm 0.00	45.39 \pm 0.00	48.79 \pm 0.00	51.46 \pm 0.00	49.29 \pm 0.00	45.76 \pm 0.00	42.68 \pm 0.00	41.34 \pm 0.00	46.11 \pm 0.00
letter	1.00 \pm 0.00	1.00 \pm 0.00	1.00 \pm 0.00	1.00 \pm 0.00	1.00 \pm 0.00	1.00 \pm 0.00	1.00 \pm 0.00	1.00 \pm 0.00	1.00 \pm 0.00	1.00 \pm 0.00	1.00 \pm 0.00	1.00 \pm 0.00
lymphography	97.89 \pm 4.21	93.61 \pm 7.97	94.76 \pm 6.53	92.71 \pm 6.69	91.66 \pm 7.34	91.57 \pm 5.93	91.57 \pm 5.96	90.69 \pm 3.65	90.69 \pm 3.65	92.96 \pm 4.97	92.96 \pm 4.97	91.77 \pm 3.69
magic gamma	76.79 \pm 0.00	76.20 \pm 0.00	75.49 \pm 0.00	75.49 \pm 0.00	75.49 \pm 0.00	75.61 \pm 0.00	75.61 \pm 0.00	76.17 \pm 0.00	75.13 \pm 0.00	74.60 \pm 0.00	73.49 \pm 0.00	73.49 \pm 0.00
mammography	41.92 \pm 0.00	41.92 \pm 0.00	41.23 \pm 0.00	44.23 \pm 0.00	44.23 \pm 0.00	43.15 \pm 0.00	43.15 \pm 0.00	43.46 \pm 0.00	43.83 \pm 0.00	43.83 \pm 0.00	43.00 \pm 0.00	43.00 \pm 0.00
mnist	72.71 \pm 0.00	72.29 \pm 0.00	71.57 \pm 0.00	70.43 \pm 0.00	69.86 \pm 0.00	71.37 \pm 0.00	71.86 \pm 0.00	71.29 \pm 0.00	69.86 \pm 0.00	69.71 \pm 0.00	69.71 \pm 0.00	69.49 \pm 0.00
musik	30.00 \pm 0.00	24.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00
opdigits	30.00 \pm 0.00	24.00 \pm 0.00	12.00 \pm 0.00	7.33 \pm 0.00	6.67 \pm 0.00	16.00 \pm 0.00	59.02 \pm 0.00	59.02 \pm 0.00	12.00 \pm 0.00	12.00 \pm 0.00	12.00 \pm 0.00	12.00 \pm 0.00
pageblocks	59.41 \pm 0.00	59.22 \pm 0.00	59.61 \pm 0.00	58.43 \pm 0.00	58.43 \pm 0.00	59.02 \pm 0.00	59.02 \pm 0.00	59.22 \pm 0.00	60.20 \pm 0.00	60.20 \pm 0.00	56.08 \pm 0.00	58.16 \pm 0.00
pendigits	94.23 \pm 0.00	92.31 \pm 0.00	71.94 \pm 2.46	71.48 \pm 2.36	71.44 \pm 2.36	71.93 \pm 1.99	72.50 \pm 1.97	70.18 \pm 2.53	70.18 \pm 2.12	71.58 \pm 2.16	71.67 \pm 2.10	71.73 \pm 2.02
satellite	72.20 \pm 0.00	71.76 \pm 0.00	60.68 \pm 0.00	69.79 \pm 0.00	69.36 \pm 0.00	70.73 \pm 0.00	71.81 \pm 0.00	70.73 \pm 0.00	69.84 \pm 0.00	69.73 \pm 0.00	69.52 \pm 0.00	69.52 \pm 0.00
satimage-2	90.14 \pm 0.00	90.14 \pm 0.00	90.14 \pm 0.00	92.96 \pm 0.00	92.96 \pm 0.00	91.27 \pm 0.00	91.27 \pm 0.00	91.55 \pm 0.00	92.06 \pm 0.00	92.06 \pm 0.00	92.11 \pm 0.00	92.11 \pm 0.00
shuttle	98.35 \pm 0.00	98.23 \pm 0.00	98.12 \pm 0.00	98.12 \pm 0.00	98.19 \pm 0.00	98.19 \pm 0.00	98.26 \pm 0.00	98.06 \pm 0.00	98.06 \pm 0.00	98.09 \pm 0.00	98.09 \pm 0.00	98.11 \pm 0.00
skin	97.12 \pm 0.46	96.93 \pm 0.38	94.14 \pm 0.23	94.30 \pm 0.27	95.82 \pm 0.14	94.73 \pm 1.41	95.83 \pm 0.17	94.73 \pm 1.27	96.76 \pm 0.29	96.76 \pm 0.29	93.89 \pm 0.13	93.89 \pm 0.13
spambase	80.88 \pm 0.00	80.29 \pm 0.00	79.81 \pm 0.00	79.45 \pm 0.00	80.13 \pm 0.00	80.52 \pm 0.00	79.13 \pm 0.00	79.45 \pm 0.00	79.13 \pm 0.00	78.98 \pm 0.00	78.80 \pm 0.00	79.48 \pm 0.00
speech	3.28 \pm 0.00	1.64 \pm 0.00	3.28 \pm 0.00	1.64 \pm 0.00	3.28 \pm 0.00	3.28 \pm 0.00	3.28 \pm 0.00	3.28 \pm 0.00	3.28 \pm 0.00	3.28 \pm 0.00	3.28 \pm 0.00	3.28 \pm 0.00
stamps	85.93 \pm 1.66	80.63 \pm 2.92	74.04 \pm 4.45	68.12 \pm 7.14	66.98 \pm 6.79	75.14 \pm 4.45	78.57 \pm 8.41	68.89 \pm 8.23	62.67 \pm 6.19	59.52 \pm 7.08	67.55 \pm 0.00	66.49 \pm 5.15
thyroid	75.27 \pm 0.00	75.27 \pm 0.00	74.19 \pm 0.00	74.19 \pm 0.00	74.62 \pm 0.00	74.62 \pm 0.00	74.12 \pm 0.00	74.12 \pm 0.00	74.12 \pm 0.00	74.12 \pm 0.00	74.19 \pm 0.00	74.19 \pm 0.00
vertebral	49.03 \pm 4.93	32.73 \pm 5.11	24.06 \pm 3.90	15.07 \pm 1.64	12.51 \pm 2.44	26.68 \pm 3.24	23.02 \pm 5.17					

2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254

Table 15.2: Average F1 score \pm standard dev. over five seeds for the semi-supervised setting of ICL and DTE-C baselines with varying hyperparameter (HP) values; For ICL, the learning rate $\in \{0.1, 0.02, 0.001, 0.0001, 1e-05\}$, for DTE-C, $k \in \{5, 10, 20, 40, 50\}$. Also reported is the avg model. We use **bold** and underline respectively to mark the **best** and the worst performance of each model to showcase the variability of performance across different HP settings.

dataset	ICL-0.1	ICL-0.01	ICL-0.001	ICL-le-0.5	ICL-le-0.05	ICL-le-0.001	DTE-C-5	DTE-C-10	DTE-C-20	DTE-C-40	DTE-C-50	DTE-C-avr
abt	4.51 \pm 0.69	4.34 \pm 0.42	5.28 \pm 0.47	4.68 \pm 0.30	4.16 \pm 0.38	4.59 \pm 0.07	4.75 \pm 0.27	4.27 \pm 0.19	4.28 \pm 0.10	4.51 \pm 0.17	4.00 \pm 0.00	3.56 \pm 0.03
amazon	10.44 \pm 0.46	9.76 \pm 0.34	9.92 \pm 0.84	10.08 \pm 0.35	<u>9.52</u> \pm 0.43	9.94 \pm 0.32	11.48 \pm 0.97	11.96 \pm 1.68	11.04 \pm 2.07	0.00 \pm 0.00	0.00 \pm 0.00	7.01 \pm 0.72
android	54.87 \pm 1.38	42.25 \pm 2.97	57.53 \pm 3.45	52.56 \pm 3.28	77.23 \pm 2.97	53.57 \pm 2.97	86.76 \pm 1.00	46.19 \pm 18.39	83.03 \pm 2.14	84.50 \pm 0.60	0.00 \pm 0.00	61.63 \pm 0.20
backdoor	87.17 \pm 0.99	<u>87.32</u> \pm 0.99	87.11 \pm 1.09	86.85 \pm 0.95	83.57 \pm 1.01	86.76 \pm 1.00	90.48 \pm 0.78	90.10 \pm 1.59	90.10 \pm 1.35	95.31 \pm 0.70	96.11 \pm 0.44	92.50 \pm 0.79
breast	95.98 \pm 0.34	96.07 \pm 0.94	97.44 \pm 0.55	96.80 \pm 0.75	97.44 \pm 0.55	96.11 \pm 0.78	97.44 \pm 0.55	97.44 \pm 0.55	97.44 \pm 0.55	97.44 \pm 0.55	97.44 \pm 0.55	97.44 \pm 0.55
campaign	48.12 \pm 0.36	46.81 \pm 1.72	50.68 \pm 0.66	51.37 \pm 0.85	53.40 \pm 0.51	50.07 \pm 0.49	51.98 \pm 1.07	52.45 \pm 1.07	52.33 \pm 1.00	0.00 \pm 0.00	0.00 \pm 0.00	31.35 \pm 0.44
cardio	49.09 \pm 1.18	61.93 \pm 5.57	58.86 \pm 1.59	57.95 \pm 2.30	<u>40.57</u> \pm 4.28	53.68 \pm 2.83	58.30 \pm 0.58	57.84 \pm 0.43	58.07 \pm 0.23	0.34 \pm 0.68	0.00 \pm 0.00	34.91 \pm 0.99
cardiography	36.14 \pm 1.28	41.07 \pm 1.28	36.14 \pm 1.28	36.36 \pm 1.34	<u>26.66</u> \pm 1.59	36.88 \pm 1.27	59.91 \pm 1.05	14.31 \pm 7.73	14.31 \pm 7.73	48.23 \pm 0.39	0.00 \pm 0.00	10.12 \pm 1.04
celeba	15.42 \pm 2.29	17.97 \pm 2.55	17.20 \pm 1.92	17.46 \pm 1.17	16.17 \pm 0.65	16.84 \pm 0.88	19.18 \pm 2.74	17.12 \pm 1.45	14.31 \pm 2.28	0.00 \pm 0.00	0.00 \pm 0.00	10.31 \pm 0.52
census	22.72 \pm 1.73	24.06 \pm 2.05	24.06 \pm 1.51	24.06 \pm 1.79	24.06 \pm 1.51	17.85 \pm 1.42	24.06 \pm 1.79	24.06 \pm 1.79	24.06 \pm 1.79	24.06 \pm 1.79	24.06 \pm 1.79	24.06 \pm 1.79
cover	26.77 \pm 1.54	15.53 \pm 8.17	42.68 \pm 17.00	53.70 \pm 20.24	36.61 \pm 2.37	36.61 \pm 2.37	68.92 \pm 4.22	46.54 \pm 3.93	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	38.39 \pm 0.62
donors	43.71 \pm 1.52	81.85 \pm 3.56	83.59 \pm 2.36	89.28 \pm 2.66	89.28\pm2.66	78.24 \pm 6.61	87.99 \pm 1.87	78.24 \pm 6.71	78.24 \pm 6.71	11.80 \pm 23.60	0.00 \pm 0.00	47.58 \pm 1.61
fault	60.35 \pm 3.56	59.52 \pm 2.09	58.57 \pm 1.11	58.57 \pm 1.51	59.13 \pm 1.36	55.57 \pm 1.51	55.16\pm0.81	55.16 \pm 0.81	55.16 \pm 0.81	96.91 \pm 0.24	72.00 \pm 0.47	72.00 \pm 0.47
fraud	57.54 \pm 0.13	48.97 \pm 8.02	58.90 \pm 6.77	66.88 \pm 4.88	79.18 \pm 3.21	62.30 \pm 9.91	75.61 \pm 7.76	54.25 \pm 5.90	22.55 \pm 24.73	0.00 \pm 0.00	0.00 \pm 0.00	30.48 \pm 4.66
hepatitis	43.53 \pm 20.21	57.05 \pm 16.03	84.05 \pm 6.11	87.24\pm5.04	82.50 \pm 5.41	80.87 \pm 3.12	85.43 \pm 5.07	54.48 \pm 4.93	19.45 \pm 6.69	24.15 \pm 2.73	0.00 \pm 0.00	24.15 \pm 2.73
http	99.64 \pm 0.71	94.69 \pm 7.81	99.64 \pm 0.71	99.64 \pm 0.71	99.64 \pm 0.71	99.64 \pm 0.71	98.65 \pm 2.11	96.40 \pm 3.01	94.63 \pm 4.90	92.51 \pm 3.12	51.68 \pm 9.78	18.79 \pm 16.60
imdb	10.52 \pm 0.65	10.44 \pm 0.54	10.44 \pm 0.54	10.44 \pm 0.54	10.44 \pm 0.54	9.76 \pm 0.54	9.76 \pm 0.54	9.76 \pm 0.54	9.76 \pm 0.54	9.76 \pm 0.54	9.76 \pm 0.54	9.76 \pm 0.54
intemrads	55.92 \pm 2.66	57.45 \pm 0.61	57.77 \pm 1.10	58.26 \pm 0.50	49.02 \pm 1.45	58.87 \pm 1.56	93.38 \pm 2.79	89.67 \pm 1.44	89.41 \pm 1.53	78.12 \pm 2.13	78.12 \pm 2.13	79.23 \pm 4.14
ionosphere	92.64 \pm 4.66	91.41 \pm 4.66	91.41 \pm 4.66	94.48 \pm 0.71	94.49 \pm 1.56	93.38 \pm 2.27	93.76 \pm 0.79	93.76 \pm 0.79	93.76 \pm 0.79	49.44 \pm 16.82	49.44 \pm 16.82	49.44 \pm 16.82
landsat	49.50 \pm 1.24	47.94 \pm 1.88	54.51 \pm 1.02	54.25 \pm 0.74	47.97 \pm 1.09	50.83 \pm 0.68	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	41.47 \pm 3.39
letter	6.80 \pm 4.21	4.00 \pm 1.55	3.60 \pm 1.36	3.20 \pm 0.75	11.60 \pm 2.06	5.84 \pm 1.11	2.40 \pm 1.50	3.00 \pm 1.17	3.20 \pm 1.17	1.72 \pm 0.37	0.00 \pm 0.00	1.72 \pm 0.37
lymphography	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	77.11 \pm 6.79	79.84 \pm 1.53	79.84 \pm 1.53	100.00 \pm 0.00	0.00 \pm 0.00	58.34 \pm 6.75
magic gamma	64.38 \pm 0.50	69.99 \pm 2.98	69.99 \pm 2.87	69.99 \pm 2.87	70.12 \pm 0.50	71.64 \pm 0.43	99.34 \pm 0.93	99.34 \pm 0.93	99.34 \pm 0.93	100.00 \pm 0.00	0.00 \pm 0.00	85.36 \pm 1.69
mammography	36.62 \pm 8.76	38.18 \pm 8.41	27.69 \pm 8.87	29.08 \pm 8.63	52.79 \pm 2.09	29.08 \pm 2.09	73.94 \pm 2.39	67.69 \pm 2.39	67.69 \pm 2.39	35.48 \pm 2.36	35.48 \pm 2.36	35.48 \pm 2.36
mnist	45.62 \pm 1.84	46.20 \pm 3.18	50.54 \pm 1.26	50.20 \pm 1.30	62.66 \pm 2.09	46.80 \pm 8.26	89.36 \pm 1.65	100.00 \pm 0.00	100.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	38.47 \pm 1.48
musik	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	100.00 \pm 0.00	97.73 \pm 0.00	97.73 \pm 0.00	97.73 \pm 0.00	60.00 \pm 0.00	0.00 \pm 0.00	77.31 \pm 0.38
opdigits	29.87 \pm 5.68	41.20 \pm 3.49	44.07 \pm 5.56	71.73 \pm 0.53	46.69 \pm 5.01	46.69 \pm 5.01	14.59 \pm 1.81	62.24 \pm 9.08	82.23 \pm 3.37	61.14 \pm 0.52	49.33 \pm 0.57	61.14 \pm 0.52
pageblocks	62.16 \pm 2.84	64.08 \pm 2.95	62.63 \pm 2.65	63.88 \pm 1.08	62.20 \pm 1.05	62.93 \pm 0.92	62.56 \pm 5.32	61.03 \pm 7.50	61.03 \pm 7.50	49.37 \pm 0.95	49.37 \pm 0.95	49.37 \pm 0.95
pendigits	46.03 \pm 1.78	60.51 \pm 10.58	60.51 \pm 3.15	60.51 \pm 3.57	66.03 \pm 5.52	51.03 \pm 3.57	63.56 \pm 5.32	63.56 \pm 5.32	63.56 \pm 5.32	32.36 \pm 1.43	32.36 \pm 1.43	32.36 \pm 1.43
pima	68.77 \pm 4.16	76.88 \pm 2.12	76.88 \pm 2.07	71.40 \pm 0.77	71.40 \pm 0.77	70.45 \pm 1.78	66.14 \pm 2.35	63.82 \pm 2.19	63.82 \pm 2.19	28.42 \pm 1.40	28.42 \pm 1.40	28.42 \pm 1.40
satellite	65.94 \pm 0.44	72.95 \pm 2.96	72.95 \pm 2.96	72.76 \pm 0.57	73.47 \pm 2.07	73.47 \pm 2.07	73.62 \pm 2.67	73.62 \pm 2.67	73.62 \pm 2.67	73.27 \pm 0.32	73.27 \pm 0.32	73.27 \pm 0.32
satimage-2	2.05 \pm 1.23	3.01 \pm 1.23	3.28 \pm 1.04	4.20 \pm 1.23	4.59 \pm 1.61	3.48 \pm 0.92	5.93 \pm 0.80	2.95 \pm 1.61	2.95 \pm 1.61	0.00 \pm 0.00	0.00 \pm 0.00	8.01 \pm 0.96
shuttle	38.03 \pm 8.59	71.68 \pm 32.26	91.83 \pm 3.54	89.58 \pm 5.56	89.58 \pm 5.56	89.58 \pm 5.56	93.34 \pm 12.38	56.82 \pm 6.23	56.82 \pm 6.23	56.82 \pm 6.23	56.82 \pm 6.23	56.82 \pm 6.23
spams	2.97 \pm 1.11	3.82 \pm 0.50	3.82 \pm 0.50	3.82 \pm 0.50	3.82 \pm 0.50	3.82 \pm 0.50	3.82 \pm 0.50	3.82 \pm 0.50	3.82 \pm 0.50	0.00 \pm 0.00	0.00 \pm 0.00	2.03 \pm 0.25
thyroid	68.39 \pm 4.17	61.29 \pm 1.80	61.08 \pm 4.53	61.87 \pm 5.57	57.68 \pm 6.77	61.76 \pm 7.57	73.76 \pm 2.21	76.13 \pm 3.03	76.13 \pm 3.03	60.85 \pm 5.87	60.85 \pm 5.87	60.85 \pm 5.87
vertebral	23.17 \pm 6.96	29.15 \pm 9.05	68.13 \pm 11.60	68.13 \pm 3.19	54.89 \pm 6.76	44.91 \pm 1.28	43.88 \pm 5.22	36.26 \pm 11.73	36.26 \pm 11.73	78.28 \pm 8.71	78.28 \pm 8.71	78.28 \pm 8.71
vowels	22.40 \pm 7.59	18.38 \pm 6.52	18.38 \pm 7.44	30.80 \pm 4.04	24.00 \pm 7.04	24.00 \pm 7.04	25.36 \pm 4.22	30.00 \pm 5.51	40.80 \pm 3.25	82.18 \pm 9.6	82.18 \pm 9.6	82.18 \pm 9.6
yelp	47.80 \pm 3.06	38.20 \pm 7.57	35.40 \pm 4.47	11.60 \pm 1.62	38.40 \pm 4.47	32.28 \pm 5.79	11.20 \pm 1.72	11.20 \pm 1.72	11.20 \pm 1.72	12.20 \pm 1.60	0.00 \pm 0.00	7.20 \pm 0.32
wbc	28.86 \pm 3.89	84.03 \pm 7.36	76.89 \pm 4.35	95.71 \pm 4.90	90.32 \pm 5.34	89.61 \pm 5.34						

Table 16.1: Average AUROC \pm standard dev. over five seeds for the semi-supervised setting on ADBench. Rank of each model among 32 models (26 baselines + 4 avg variants of top-4 baselines + 2 FoMo-0D variants w/ $D = 100$ and $D = 20$) per dataset is provided (in parentheses) (the lower, the better). We use blue and green respectively to mark the top-1 and the top-2 method. Last four rows show avg_rank of methods across datasets, and p -values of the Wilcoxon signed rank test comparing FoMo-0D ($D = 100$) with other baselines. The previous four rows are the same for FoMo-0D ($D = 20$), when ranking 31 models (26 baselines + 4 avg variants of top-4 baselines + FoMo-0D w/ $D = 20$).

Table 16.2: Average AUROC \pm standard dev. over five seeds for the semi-supervised setting on ADBench. Rank of each model per dataset is provided (in parentheses) (the lower, the better). We use blue and green respectively to mark the top-1 and the top-2 method.

Table 17.1: Average AUPR \pm standard dev. over five seeds for the semi-supervised setting on ADBench. Rank of each model per dataset is provided (in parentheses) (the lower, the better). We use blue and green respectively to mark the top-1 and the top-2 method.

A Foundation Model for Zero-shot Outlier Detection

Table 18.1: Average F1 score \pm standard dev. over five seeds for the semi-supervised setting on ADBench. Rank of each model per dataset is provided (in parentheses) (the lower, the better). We use blue and green respectively to mark the top-1 and the top-2 method.

2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584

A Foundation Model for Zero-shot Outlier Detection

Table 18.2: Average F1 score \pm standard dev. over five seeds for the semi-supervised setting on AD-Bench. Rank of each model per dataset is provided (in parentheses) (the lower, the better). We use blue and green respectively to mark the top-1 and the top-2 method.

Dataset	VAE	PCA	PlanarFlow	HBOS	GAnomaly		GOAD		DIF		COPOD		ECOD		DeepSVD		LDA		DAGMM		DROCC		ICL-PW		KNN ^{avg}		DT-EPPs		DT-ICL			
					3.91	±1.75	3.45	±1.88	0.40	±0.25	5.84	±0.24	11.40	±0.07	0.48	±0.12	9.67	±0.20	60.0	±0.32	11.76	±1.93	0.40	±0.23	11.40	±0.13	3.56	±0.38	3.56	±0.38		
alei	7.63±0.0(5.5)	7.61±0.0(5.5)	7.43±0.0(7.29)	7.43±0.0(7.7)	9.32±0.0(2.29)	9.32±0.0(2.65)	11.56±0.1(0.85)	11.56±0.1(0.85)	11.44±0.0(0.75)	11.44±0.0(0.75)	5.84	±0.24	10.0	±0.02	11.81	±0.17	10.64±1.18(1.7)	11.76±1.19(2.0)	5.98	±1.80(1)	11.81	±0.17	4.59	±0.07	4.59	±0.07	3.25	±0.32	3.94	±0.32	3.94	±0.32
amazon	1.10±0.0(2.5)	1.10±0.0(2.5)	9.48±0.6(22.6)	9.48±0.6(22.6)	8.96±1.4(7.40)	8.96±1.4(7.40)	58.99±0.0(3.9)	58.99±0.0(3.9)	58.99±0.0(3.9)	58.99±0.0(3.9)	31.65±0.1(0.31)	31.65±0.1(0.31)	31.65±0.1(0.31)	31.65±0.1(0.31)	31.65±0.1(0.31)	31.65±0.1(0.31)	94.88±1.3(2.6)	94.88±1.3(2.6)	45.66±1.6(4.47)	45.66±1.6(4.47)	85.44±1.1(14.3)	85.44±1.1(14.3)	42.86±1.2(12.2)	42.86±1.2(12.2)	3.87±1.2(2.05)	3.87±1.2(2.05)	25.51±0.5(13)	25.51±0.5(13)	25.51±0.5(13)	25.51±0.5(13)		
anomaly	50.19±0.0(2.0)	50.19±0.0(2.0)	6.37±0.2(17.14)	6.37±0.2(17.14)	7.19±0.4(41.06)	7.19±0.4(41.06)	4.79±0.2(3.72)	4.79±0.2(3.72)	9.65±0.6(3.42)	9.65±0.6(3.42)	50.96±1.4(5.31)	50.96±1.4(5.31)	27.11±0.0(3.0)	27.11±0.0(3.0)	48.33±0.1(0.15)	48.33±0.1(0.15)	91.85±1.7(7.24)	91.85±1.7(7.24)	95.67±1.0(4.78)	95.67±1.0(4.78)	60.56±2.5(5.32)	60.56±2.5(5.32)	95.42±0.3(11.0)	95.42±0.3(11.0)	50.93±0.3(10.2)	50.93±0.3(10.2)	50.93±0.3(10.2)	50.93±0.3(10.2)				
backdoor	8.5±1.2(27.21)	8.3±1.2(27.21)	9.57±0.4(16.21)	9.57±0.4(16.21)	7.23±0.2(20.17)	7.23±0.2(20.17)	0.00±0.0(0.01)	0.00±0.0(0.01)	0.00±0.0(0.01)	0.00±0.0(0.01)	90.07±1.4(3.72)	90.07±1.4(3.72)	30.71±0.0(12)	30.71±0.0(12)	48.33±0.1(0.15)	48.33±0.1(0.15)	31.73±1.2(2.6)	31.73±1.2(2.6)	95.64±1.4(5.28)	95.64±1.4(5.28)	42.86±1.2(12.2)	42.86±1.2(12.2)	42.75±1.0(5.13)	42.75±1.0(5.13)	42.75±1.0(5.13)	42.75±1.0(5.13)						
breastw	96.12±0.4(7.10)	95.78±0.4(5.15)	90.09±1.6(69.22)	90.09±1.6(69.22)	42.11±2.2(88.21)	42.11±2.2(88.21)	47.70±0.1(0.12)	47.70±0.1(0.12)	49.27±0.1(0.12)	49.27±0.1(0.12)	30.27±0.0(3.0)	30.27±0.0(3.0)	91.73±1.4(3.72)	91.73±1.4(3.72)	30.67±0.0(12)	30.67±0.0(12)	48.33±0.1(0.15)	48.33±0.1(0.15)	31.73±1.2(2.6)	31.73±1.2(2.6)	95.64±1.4(5.28)	95.64±1.4(5.28)	42.75±1.0(5.13)	42.75±1.0(5.13)	42.75±1.0(5.13)	42.75±1.0(5.13)						
cannabis	4.88±0.0(0.14)	4.88±0.0(0.14)	7.61±4.0(9.0)	7.61±4.0(9.0)	4.79±0.4(5.22)	4.79±0.4(5.22)	6.07±0.0(0.15)	6.07±0.0(0.15)	6.07±0.0(0.15)	6.07±0.0(0.15)	30.27±0.0(3.0)	30.27±0.0(3.0)	91.73±1.4(3.72)	91.73±1.4(3.72)	30.67±0.0(12)	30.67±0.0(12)	48.33±0.1(0.15)	48.33±0.1(0.15)	31.73±1.2(2.6)	31.73±1.2(2.6)	95.64±1.4(5.28)	95.64±1.4(5.28)	42.75±1.0(5.13)	42.75±1.0(5.13)	42.75±1.0(5.13)	42.75±1.0(5.13)						
cardio	7.63±4.4(0.15)	7.63±4.4(0.15)	56.25±0.0(2.49)	56.25±0.0(2.49)	46.33±0.9(5.21)	46.33±0.9(5.21)	81.14±0.0(0.25)	81.14±0.0(0.25)	81.14±0.0(0.25)	81.14±0.0(0.25)	37.08±0.5(6.26)	37.08±0.5(6.26)	5.17±0.1(0.01)	5.17±0.1(0.01)	14.19±0.6(1.19)	14.19±0.6(1.19)	44.66±1.4(4.48)	44.66±1.4(4.48)	33.43±1.4(2.26)	33.43±1.4(2.26)	12.65±2.1(1.40)	12.65±2.1(1.40)	52.09±1.4(2.0)	52.09±1.4(2.0)	52.09±1.4(2.0)	52.09±1.4(2.0)						
chaingraph	61.59±0.0(0.25)	61.59±0.0(0.25)	49.44±5.4(3.01)	49.44±5.4(3.01)	41.42±0.6(1.22)	41.42±0.6(1.22)	7.07±0.1(0.01)	7.07±0.1(0.01)	7.07±0.1(0.01)	7.07±0.1(0.01)	9.07±0.5(6.26)	9.07±0.5(6.26)	7.22±0.1(0.18)	7.22±0.1(0.18)	50.96±1.4(5.32)	50.96±1.4(5.32)	37.80±0.5(6.26)	37.80±0.5(6.26)	12.65±2.1(1.40)	12.65±2.1(1.40)	52.09±1.4(2.0)	52.09±1.4(2.0)	52.09±1.4(2.0)	52.09±1.4(2.0)								
celiba	27.07±0.4(4.13)	27.07±0.4(4.13)	22.88±0.3(0.59)	22.88±0.3(0.59)	3.98±2.9(3.08)	3.98±2.9(3.08)	1.07±0.5(0.23)	1.07±0.5(0.23)	1.07±0.5(0.23)	1.07±0.5(0.23)	10.79±1.4(2.24)	10.79±1.4(2.24)	22.78±0.1(0.18)	22.78±0.1(0.18)	4.28±0.4(4.47)	4.28±0.4(4.47)	10.75±1.4(2.24)	10.75±1.4(2.24)	12.65±2.1(1.40)	12.65±2.1(1.40)	52.09±1.4(2.0)	52.09±1.4(2.0)	52.09±1.4(2.0)	52.09±1.4(2.0)								
census	20.76±0.2(7.10)	20.76±0.2(7.10)	10.75±2.6(4.22)	10.75±2.6(4.22)	10.75±2.6(4.22)	10.75±2.6(4.22)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.19±0.5(8.51)	1.19±0.5(8.51)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)					
cover	36.77±0.5(3.22)	36.77±0.5(3.22)	16.21±1.8(2.21)	16.21±1.8(2.21)	12.03±0.1(0.18)	12.03±0.1(0.18)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.19±0.5(8.51)	1.19±0.5(8.51)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)					
covertype	37.75±0.5(3.22)	37.75±0.5(3.22)	16.21±1.8(2.21)	16.21±1.8(2.21)	12.03±0.1(0.18)	12.03±0.1(0.18)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.19±0.5(8.51)	1.19±0.5(8.51)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)					
covertype	37.75±0.5(3.22)	37.75±0.5(3.22)	16.21±1.8(2.21)	16.21±1.8(2.21)	12.03±0.1(0.18)	12.03±0.1(0.18)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.19±0.5(8.51)	1.19±0.5(8.51)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)					
covertype	37.75±0.5(3.22)	37.75±0.5(3.22)	16.21±1.8(2.21)	16.21±1.8(2.21)	12.03±0.1(0.18)	12.03±0.1(0.18)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.19±0.5(8.51)	1.19±0.5(8.51)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)					
covertype	37.75±0.5(3.22)	37.75±0.5(3.22)	16.21±1.8(2.21)	16.21±1.8(2.21)	12.03±0.1(0.18)	12.03±0.1(0.18)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.19±0.5(8.51)	1.19±0.5(8.51)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)					
covertype	37.75±0.5(3.22)	37.75±0.5(3.22)	16.21±1.8(2.21)	16.21±1.8(2.21)	12.03±0.1(0.18)	12.03±0.1(0.18)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.19±0.5(8.51)	1.19±0.5(8.51)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)					
covertype	37.75±0.5(3.22)	37.75±0.5(3.22)	16.21±1.8(2.21)	16.21±1.8(2.21)	12.03±0.1(0.18)	12.03±0.1(0.18)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.19±0.5(8.51)	1.19±0.5(8.51)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)					
covertype	37.75±0.5(3.22)	37.75±0.5(3.22)	16.21±1.8(2.21)	16.21±1.8(2.21)	12.03±0.1(0.18)	12.03±0.1(0.18)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.19±0.5(8.51)	1.19±0.5(8.51)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)					
covertype	37.75±0.5(3.22)	37.75±0.5(3.22)	16.21±1.8(2.21)	16.21±1.8(2.21)	12.03±0.1(0.18)	12.03±0.1(0.18)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.19±0.5(8.51)	1.19±0.5(8.51)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)					
covertype	37.75±0.5(3.22)	37.75±0.5(3.22)	16.21±1.8(2.21)	16.21±1.8(2.21)	12.03±0.1(0.18)	12.03±0.1(0.18)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.19±0.5(8.51)	1.19±0.5(8.51)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32)	1.00±0.0(0.32																

N. Benchmark OD Datasets

Table 19: Description of all datasets in ADBench (Livernoche et al., 2024). Datasets in blue are image and text datasets that are vectorized through pretrained encoders. We refer to the original paper for details.

Dataset Name	# Samples	# Features	# Anomaly	% Anomaly	Category
ALOI	49534	27	1508	3.04	Image
annthyroid	7200	6	534	7.42	Healthcare
backdoor	95329	196	2329	2.44	Network
breastw	683	9	239	34.99	Healthcare
campaign	41188	62	4640	11.27	Finance
cardio	1831	21	176	9.61	Healthcare
Cardiotocography	2114	21	466	22.04	Healthcare
celeba	202599	39	4547	2.24	Image
census	299285	500	18568	6.20	Sociology
cover	286048	10	2747	0.96	Botany
donors	619326	10	36710	5.93	Sociology
fault	1941	27	673	34.67	Physical
fraud	284807	29	492	0.17	Finance
glass	214	7	9	4.21	Forensic
Hepatitis	80	19	13	16.25	Healthcare
http	567498	3	2211	0.39	Web
InternetAds	1966	1555	368	18.72	Image
Ionosphere	351	32	126	35.90	Oryctognosy
landsat	6435	36	1333	20.71	Astronautics
letter	1600	32	100	6.25	Image
Lymphography	148	18	6	4.05	Healthcare
magic.gamma	19020	10	6688	35.16	Physical
mammography	11183	6	260	2.32	Healthcare
mnist	7603	100	700	9.21	Image
musk	3062	166	97	3.17	Chemistry
optdigits	5216	64	150	2.88	Image
PageBlocks	5393	10	510	9.46	Document
pendigits	6870	16	156	2.27	Image
Pima	768	8	268	34.90	Healthcare
satellite	6435	36	2036	31.64	Astronautics
satimage-2	5803	36	71	1.22	Astronautics
shuttle	49097	9	3511	7.15	Astronautics
skin	245057	3	50859	20.75	Image
smtp	95156	3	30	0.03	Web
SpamBase	4207	57	1679	39.91	Document
speech	3686	400	61	1.65	Linguistics
Stamps	340	9	31	9.12	Document
thyroid	3772	6	93	2.47	Healthcare
vertebral	240	6	30	12.50	Biology
vowels	1456	12	50	3.43	Linguistics
Waveform	3443	21	100	2.90	Physics
WBC	223	9	10	4.48	Healthcare
WDBC	367	30	10	2.72	Healthcare
Wilt	4819	5	257	5.33	Botany
wine	129	13	10	7.75	Chemistry
WPBC	198	33	47	23.74	Healthcare
yeast	1484	8	507	34.16	Biology
CIFAR10	5263	512	263	5.00	Image
FashionMNIST	6315	512	315	5.00	Image
MNIST-C	10000	512	500	5.00	Image
MVTec-AD	5354	512	1258	23.50	Image
SVHN	5208	512	260	5.00	Image
Agnews	10000	768	500	5.00	NLP
Amazon	10000	768	500	5.00	NLP
Imdb	10000	768	500	5.00	NLP
Yelp	10000	768	500	5.00	NLP
20newsgroups	11905	768	591	4.96	NLP

2640 O. Differences to Prior Work on PFNs for Tabular Data

2641 There exist applications of PFNs (originally developed by Müller et al. (2022)) that pre-date our proposed FoMo-0D,
 2642 namely, TabPFN (Hollmann et al., 2023) for supervised classification, LC-PFN (Adriaensen et al., 2024) for learning curve
 2643 extrapolation, PFN4BO (Müller et al., 2023) for Bayesian optimization, and ForecastPFN (Dooley et al., 2023) for time
 2644 series forecasting.

2645 Here we highlight the differences of our proposed FoMo-0D from these existing PFNs.

- 2646 1. **First PFN4OD:** We employ prior-data fitted networks (PFNs) for outlier detection (OD) for the first time.
- 2647 2. **First large-scale pretrained OD model:** FoMo-0D is the first model for zero-shot OD that is pretrained at large
 2650 scale on a large collection of (synthetic) datasets, due to the minuscule nature of existing real-world OD benchmark
 2651 datasets.
- 2652 3. **New data prior:** Thanks to PFN’s reliance on synthetically generated datasets, we establish a new data prior for
 2653 OD, specifically for outlier synthesis.
- 2654 4. **Data transformation for scale:** While drawing samples from a data prior may be relatively fast, pretraining a large
 2655 foundation model requires many such draws for every step of each epoch. To speed up data synthesis on-the-fly, we
 2656 are the first to leverage a linear transformation.
- 2657 5. **Router-based attention for scale:** PFNs ingest the entire training dataset as context for in-context learning at
 2659 inference time. To accommodate larger datasets at both training (for better generalization) and inference (for
 2660 large-scale real-world datasets), we leveraged a “bottleneck” architecture for scalable self-attention, and in turn,
 2661 larger context size.

2663 P. Discussion

2664 **Summary:** We introduced FoMo-0D, **the first foundation model for outlier detection** (OD) on tabular data. FoMo-0D is
 2665 a prior-data fitted network (PFN), pretrained on a large number of *synthetic* datasets generated from a new data prior for OD,
 2666 which can infer the posterior predictive distribution for test points in a new dataset in a **zero-shot** fashion where the training
 2667 data is input as context, capitalizing on *in-context learning*.

2668 Zero-shot OD implies **no additional OD model training or model selection**, given a new OD task. That is a revolution
 2669 for OD (!), for which algorithm and hyperparameter selection are notoriously-hard *without any labeled data*, and also
 2670 computationally taxing especially for today’s modern deep OD models with numerous parameters *and* a long list of
 2671 hyperparameters. What is more, FoMo-0D provides **extremely fast inference** thanks to a mere *single forward pass*, making
 2672 it amenable for OD on data streams.

2673 Building on the PFN paradigm (Müller et al., 2022), FoMo-0D breaks new ground not only conceptually by abolishing the
 2674 burden of model training and selection, but also empirically: Against **26** different (both classical and modern) baselines on
 2675 **57** public benchmark datasets from diverse domains, FoMo-0D performs on par with the top *2nd* baseline, while significantly
 2676 outperforming the majority of the baselines. Without the need to train any, let alone multiple models for HP tuning, FoMo-0D
 2677 takes a mere **7.7 ms** per test sample for inference only.

2678 **Limitations and Future Directions:** FoMo-0D employs a simple straightforward data prior based on GMMs. While
 2679 it is remarkable to see how far one can go with synthetic data from such a simple prior, future work can design more
 2680 comprehensive data priors, inclusive of discrete features as well as other possible outlier types. We have also pretrained
 2681 FoMo-0D solely on synthetic datasets, while future work can augment both synthetic and real-world datasets for pretraining.

2682 Besides the lack of massive real-world datasets for tabular OD, a motivation for a data prior to pretrain purely on synthetic
 2683 datasets comes from neural scaling laws (Kaplan et al., 2020; Zhai et al., 2022). Interestingly, the scaling laws for large
 2684 Transformer models have shown that their generalization error tends to drop as a power law with the amount of training
 2685 data (also, with number of parameters and amount of compute), but the power law exponent is very small—suggesting that
 2686 acquiring more colossal real-world datasets would be a slow, if not expensive approach to advancing ML/AI. Others have
 2687 proposed ways to subset-select smaller, non-redundant “foundation datasets” (Sorscher et al., 2022; Paul et al., 2021), and
 2688 emphasized the importance of task/dataset diversity in pretraining (Raventós et al., 2024). Arguably, synthetic data from a
 2689 complex and diverse data prior is a potential gateway to obtaining non-redundant and diverse datasets for pretraining large
 2690 foundation models like FoMo-0D. On the other hand, designing such a data prior requires a level of domain/prior knowledge.

2695 Another improvement could be scaling up to even larger context (i.e. dataset) size and dimensionality. While FoMo-OD
2696 generalizes beyond pretrained context sizes and dimensionality, it is limited to and performs particularly well on downstream
2697 datasets of similar nature as our experiments showed. A promising direction for size generalization is using PFNs as
2698 extremely fast ensemble components at inference; since “*PFNs are quick enough to be used as ensemble members. The size*
2699 *constraints could therefore be overcome by boosting and bagging techniques*” (Nagler, 2023).

2700 Further, our work focused on semi-supervised OD with clean/inlier-only training data. Future work can study the unsupervised
2701 OD setting and pretraining with mixed/“contaminated” data in this transductive setting, where the unlabeled test
2702 data is the same as training data. In addition, we performed offline evaluation of FoMo-OD on static datasets, while its fast
2703 inference lends itself to streaming OD, which future work can explore. Technically, both extensions (unsupervised OD and
2704 streaming OD) are straightforward from the implementation perspective.

2705 Our current work is limited to OD for tabular (or point-cloud) data. Our ideas can be extended to other data modalities, such
2706 as image, graph, and text outliers, to comprise other domains with critical OD applications such as video surveillance, fraud
2707 detection and LLM hallucination detection. To that end, the design of novel inlier/outlier priors would be an open direction.
2708 A promising approach here could be the use of pretrained generative models to draw synthesized image/text/etc. datasets for
2709 pretraining the PFN, in place of manually-designed data priors.

2710 Finally, our quest here has been mainly experimental. Theoretically understanding why these models work as well as they
2711 do and investigating their failure cases are important yet open questions.

2712 As the first foundation model for OD, FoMo-OD inspires many promising directions for future research that could lead to
2713 fruition for additional practical applications.

2714

2715 **Q. Reproducibility Statement**

2716 We expect that the disruptive nature of FoMo-OD will trigger future innovations in the OD literature, as well as a widespread
2717 adoption by practitioners thanks to its key desirable properties. To foster future research and accessibility in practice,
2718 we make all resources (our codebase used for prior data synthesis, data transformation, and pretraining as well as our
2719 pretrained model checkpoints) publicly available at <https://anonymous.4open.science/r/PFN4OD>. Further,
2720 full implementation details are provided in Appendix F.

2721

2722

2723

2724

2725

2726

2727

2728

2729

2730

2731

2732

2733

2734

2735

2736

2737

2738

2739

2740

2741

2742

2743

2744

2745

2746

2747

2748

2749