

# ZERO-SHOT OUTLIER DETECTION VIA SYNTHETICALLY PRE-TRAINED TRANSFORMERS: MODEL SELECTION BYGONE!

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Outlier detection (OD) has a vast literature as it finds numerous applications in environmental monitoring, security, manufacturing, and finance to name a few. Being an inherently *unsupervised* task, model selection is a key bottleneck for OD (both algorithm and hyperparameter selection) without label supervision. There is a long list of techniques to choose from – both classical algorithms and deep neural architectures – and while several studies report their hyperparameter sensitivity, the literature remains quite slim on unsupervised model selection—limiting the effective use of OD in practice. In this paper we present FoMo-OD, for zero/0-shot OD exploring a transformative new direction that *bypasses* the hurdle of model selection altogether (!), thus breaking new ground. The fundamental idea behind FoMo-OD is the Prior-data Fitted Networks, recently introduced by Müller et al. (2022), which trains a Transformer model on a large body of *synthetically* generated data from a prior data distribution. In essence, FoMo-OD is a **pretrained Foundation Model for zero/0-shot OD on tabular data**, which can directly predict the (outlier/inlier) label of any test data at inference time, by merely a *single forward pass*—making obsolete the need for choosing an algorithm/architecture and tuning its associated hyperparameters, besides requiring no training of model parameters when given a new OD dataset. Extensive experiments on **57** public benchmark datasets against **26** baseline methods show that FoMo-OD performs statistically no different from the *2nd* top baseline, while significantly outperforming the majority of the baselines, with an average inference time of **7.7 ms** per test sample.

## 1 INTRODUCTION

Outlier detection (OD) finds applications in many domains such as security, environmental monitoring, finance, and so on. This popularity brings along a large literature that offers a plethora of detection algorithms to choose from given a new OD task. These techniques, however, exhibit several hyperparameters (HPs) that need careful tuning to which they are often quite sensitive (Ma et al., 2023). What makes it notoriously difficult to achieve effective OD performance in practice is *model selection* (both algorithm and HPs) in the *absence of any labels*, as most tasks are unsupervised.<sup>1</sup>

In fact, while deep learning and modern architectures have revolutionized many areas of machine learning (ML), it has not quite been the case for OD—mainly because deep OD models (Pang et al., 2021) exhibit many more HPs (for architecture, regularization, and optimization) that detection performance is sensitive to (Ding et al., 2022), as compared to classical methods with only a few HPs.

Large foundation models have stirred up most recent advances in ML, which are (pre-)trained on massive amounts of data. The most notable progress has been in natural languages and vision, thanks to the admirable quantity and quality of public text and image datasets. In contrast, public (benchmark) datasets for OD is minuscule in comparison (Han et al., 2022; Zhao et al., 2021; Steinbuss and Böhm, 2021). Another obstacle for foundation models for tabular OD has been the non-shared feature spaces of different datasets, unlike the shared pixel or word spaces for images and text.

Recently, the introduction of Prior-data Fitted Networks (PFNs) has marked a milestone as a new approach to ML on tabular data (Müller et al., 2022). PFNs are based on Bayesian non-parametrics

<sup>1</sup>While semi-/supervised settings of OD exist, unsupervised OD is preferable in most domains for the capacity to detect novel/emergent types of anomalies, beyond just the known types.

Table 1: 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-4<sup>s</sup>** baselines with performance avg.’ed over varying HPs (denoted w/ <sup>avg</sup>); 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 (42) w/  $d \leq 100$  and (46) w/  $d \leq 500$  dimensions. FoMo-0D performs as well as (*i.e.*, **statistically no different from**) the **2<sup>nd</sup> best model** ( $k$ NN, 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	$k$ NN	ICL	DTE-C	LOF	CBLOF	Feat.Bag.	SLAD	DDPM	OCSVM	DTE-NP <sup>avg</sup>	$k$ NN <sup>avg</sup>	ICL <sup>avg</sup>	DTE-C <sup>avg</sup>
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
$d \leq 100$	-	0.415	0.700	0.949	<b>0.953</b>	<b>0.970</b>	<b>0.971</b>	<b>0.996</b>	0.876	<b>0.980</b>	<b>0.978</b>	0.752	0.860	<b>0.958</b>	<b>1.000</b>
$d \leq 500$	-	0.220	0.569	0.827	0.894	<b>0.960</b>	<b>0.968</b>	<b>0.994</b>	0.910	<b>0.960</b>	<b>0.979</b>	0.607	0.756	0.846	<b>1.000</b>

and meta-learning on large quantities of *synthetically* simulated data from a data prior. The key idea is to compute a posterior predictive distribution (PPD) for a test point given the training data as input context. To approximate the PPD, a Transformer (Vaswani et al., 2017) is pre-trained to mimic the PPD via simulating numerous training datasets from a (general, complex) data prior. For inference, the fresh training set along with the test samples are passed to the (frozen) pre-trained PFN, which outputs the predictions in a *single forward pass*, requiring no model training or model selection. Variants of PFN are shown to match the performance of tree-based models on small classification datasets (Hollmann et al., 2023) and in time series forecasting with limited data (Dooley et al., 2023).

In this paper, we capitalize on these ideas and introduce FoMo-0D; a prior-data fitted Foundation Model for zero- or 0-shot Outlier Detection (for the “Fear of Missing out”-liers). The implication and “gift” of PFNs for unsupervised OD goes beyond those for supervised learning: it helps *bypass* not only model (parameter) training, but most importantly, the notoriously-hard task of model (hyperparameter) selection altogether. As such, FoMo-0D unlocks zero-shot OD on a new dataset without the need for any algorithm or HP selection. During inference, data is used only as input *context* to FoMo-0D, and *not* for parameter training or HP tuning. Arguably, this is a potential game changer for unsupervised OD, especially for practitioners. Figure 1 illustrates the new FoMo-0D paradigm versus the typical OD setting.

In designing FoMo-0D, we simply use Gaussian mixture models as a simple yet effective tabular data prior, to capture general and diverse inlier data distributions, following current literature (Hollmann et al., 2023; Zhao et al., 2021). We combine these with simulated outlier types common in the real-world; namely local and global subspace outliers (Steinbuss and Böhm, 2021). While the data prior can be extended to comprise more complex data distributions (e.g. through the use of Bayesian Neural Networks (BNNs; (Neal, 2012)) and Structural Causal Models (SCMs; (Pearl, 2009))) as in (Hollmann et al., 2023)), and additional outlier types can be included (e.g. dependency, contextual, etc. outliers), as we show in the experiments, even with the relatively straightforward prior that we employed, FoMo-0D achieves remarkable performance. As shown in Table 1, FoMo-0D pretrained on synthetic datasets with up to 100 dimensions performs *statistically no different* from all 26 state-of-the-art baselines (all  $p$ -values  $> 0.2$ ) on 46 benchmark datasets with dimensionality  $d \leq 500$ , while *significantly outperforming the majority* of the baselines (with  $p > 0.95$ ) (see Appendix Tables 12.1&12.2). Further, FoMo-0D takes a mere average of 7.7 ms to infer a test sample since a new dataset requires a single forward pass for inference and no training overhead.

**Our contributions:** We summarize the main contributions of our work as follows.

- **A Foundation Model for Tabular OD:** We present FoMo-0D, *the first foundation model for zero-shot OD* on tabular datasets. FoMo-0D is a Prior-data Fitted Network (PFN) (Müller et al., 2022) that is pretrained on many synthetically generated datasets drawn from a novel data prior that we introduce to capture various inlier and outlier distributions. The pretrained FoMo-0D can then directly compute the posterior predictive distribution (PPD) of test points in a new dataset.
- **Unsupervised Outlier Model Selection Made Obsolete:** The most outstanding property of FoMo-0D is its *zero-shot inference* on a new dataset via a single forward pass, fully abolishing the need not only for model training on a new dataset, but importantly also the notorious task of algorithm selection and hyperparameter tuning in the absence of labeled data.
- **Scalable Pre-training Design:** To unlock the premise of large-scale pretraining on numerous large datasets, (1) we implement a new mechanism to speed up sample-to-sample attention from

quadratic to *linear time* complexity—enabling *larger datasets*; and (2) we scale up on-the-fly data synthesis through data transformation—enabling *more datasets* in less time.

- **Fast Inference at Detection Time:** Thanks to a pretrained prior-data fitted Transformer, FoMo-OD bypasses both model (parameter) training and selection, both of which can be slow for modern deep OD models with many hyper/parameters. Rather, it takes *fraction of a second* to label a test point through a single forward pass that can be parallelized across test samples. Such speedy inference also unlocks the potential for deploying FoMo-OD in *real time* on data streams.
- **Effectiveness:** We evaluate FoMo-OD on **57** public benchmark datasets (Han et al., 2022) from diverse domains and compare against **26** baselines from classical to modern (Livernoche et al., 2024), where FoMo-OD significantly outperforms the majority of the baselines while performing statistically no different from the top *2nd* baseline, at the fraction of the compute cost.

As FoMo-OD proposes a paradigm shift for OD, abolishing model training and selection altogether, while delivering unreasonable effectiveness on benchmark datasets even with a basic data prior, we expect FoMo-OD will trigger further work in both research and practice. To this end, we make all of our codebase for synthetic data generation, model training, and our pretrained FoMo-OD checkpoints, openly available at <https://anonymous.4open.science/r/PFN40D>.

## 2 PROBLEM AND PRELIMINARIES

### 2.1 SEMI-SUPERVISED OUTLIER DETECTION

Outlier detection (OD) methods can be categorized based on the availability of labeled data. In supervised OD, the task is similar to binary classification with imbalanced classes (as outliers typically make up only a small portion of the overall data). The more difficult unsupervised setting assumes the “contaminated” training data contains both inliers and outliers, but without any labels. A semi-supervised or one-class classification approach lies between these two extremes, where only inlier data is available for training, but unknown outliers may appear during inference. Semi-supervised OD is used in practice where it is easy to gather inlier data, but learning from known, labeled outliers is undesirable because outliers are hard to collect and/or new, unknown outlier types are likely to arise in future test data that renders learning only from the known outliers suboptimal/risky.

Note that semi-supervised OD may be a *misnomer* from the supervised ML perspective, where semi-supervised classification assumes the presence of some labeled instances from **all** classes in the training data. As such, model selection continues to be as difficult for semi-supervised OD as unsupervised OD, where no labeled outliers exist in the input/training data in both settings.

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  $\mathcal{D}_{\text{in}}$ .

### 2.2 BACKGROUND ON PRIOR-DATA FITTED NETWORKS

**Posterior Predictive Distribution (PPD):** In the Bayesian framework for supervised learning, the prior defines a hypotheses space  $\Phi$  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  $\mathbf{x}_{\text{test}}$ , 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  $\Phi$ :

$$p(y_{\text{test}} | \mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}}) = \int_{\Phi} p(y_{\text{test}} | \mathbf{x}_{\text{test}}, \phi) p(\mathcal{D}_{\text{train}} | \phi) p(\phi) d\phi, \quad (1)$$

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

**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 offline on simulated datasets that are generated according to a prior distribution. Specifically, it contains the pre-training and inference stages described as the following.

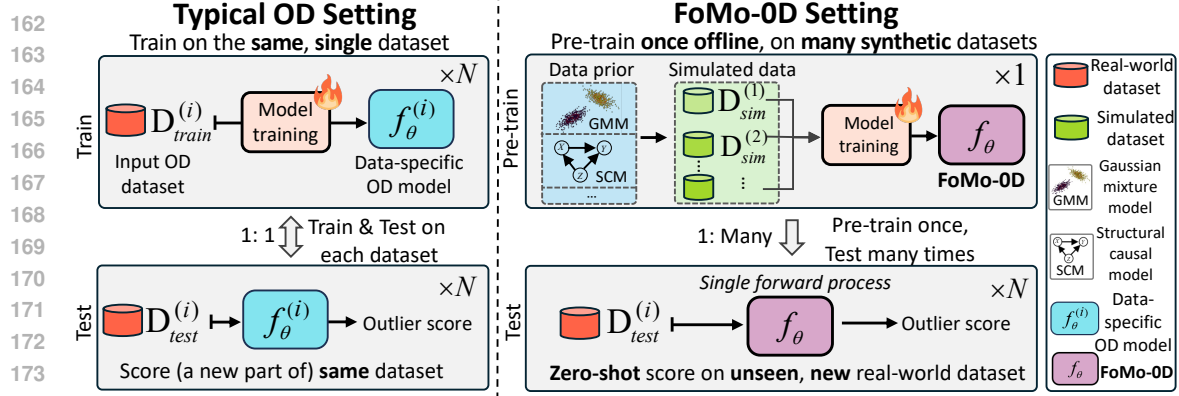


Figure 1: (best in color) Comparison of typical OD vs. the FoMo-OD settings. Given a new un/supervised OD dataset, FoMo-OD not only eliminates the need for model training, but most importantly, also abolishes the onerous task of model selection (algorithm and hyperparameters) w/out labels.

*Pre-training on synthetic data.* At the beginning of the pre-training stage, massive synthetic training datasets are generated, 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 purposes, 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  $\theta$  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 formulated as

$$\mathcal{L} = \mathbb{E}_{(\{\mathbf{x}_{\text{test}}, y_{\text{test}}\}) \cup \mathcal{D}_{\text{train}} \sim p(\mathcal{D})} [-\log q_{\theta}(y_{\text{test}}|\mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})] . \quad (2)$$

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  $q_{\theta}$  is typically implemented by a Transformer-based architecture (Vaswani et al., 2017), 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 output is the conditional class probabilities for  $\mathbf{x}_{\text{test}}$ . As the whole training set  $\mathcal{D}_{\text{train}}$  is passed as input/context to the Transformer, it learns to predict class labels through sample-to-sample attention.

*Inference on real-world data.* In the inference stage, a fresh real-world dataset  $\mathcal{D}_{\text{train}}$  and some test instance  $\mathbf{x}_{\text{test}}$  are 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 process. Importantly, PFNs do not require gradient-based parameter tuning on data observed at inference time, where the training and prediction are delivered through a one-step forward process *in less than a second* (Hollmann et al., 2023).

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

### 3 FoMo-OD: A NEW PFN FOR 0-SHOT OD – MODEL SELECTION BYGONE!

Inspired by the recent PFNs (Müller et al., 2022) and their successful applications in supervised classification (Hollmann et al., 2023) and time series forecasting (Dooley et al., 2023), we propose FoMo-OD, a prior-data fitted Foundation Model for 0-shot Outlier Detection. FoMo-OD is (pre)trained on a large body of synthetically generated OD datasets toward zero-shot inference on a new dataset. Most notable of our zero-shot FoMo-OD is its elimination of the need not only for model training on a new dataset, but especially also for model selection (both algorithm and HPs), which is notoriously hard without any labeled data. By breaking such new ground, and its effectiveness on many benchmark datasets compared to classical and modern baselines, we expect FoMo-OD will become a milestone in future research and practice of OD. The new FoMo-OD paradigm (right) versus the typical OD setting (left) is illustrated in Figure 1.

In the following we describe our OD data prior, training of FoMo-OD on prior-simulated datasets, inference on new datasets, and our specific model architecture and improvements for scalable training.

### 3.1 DESIGNING A DATA PRIOR FOR OUTLIER DETECTION

Arguably, what has triggered the recent breakthroughs in NLP and CV is the massive amounts of datasets available for (pre)training, along with high-capacity model architectures. In comparison to the natural language and image domains, the quantity (and quality) of publicly available tabular OD datasets is minuscule. Even in the presence of large quantities of data, in training their Chronos foundation models for time series forecasting, Ansari et al. (2024) show that using synthetic data in combination with real-world data improves the overall zero-shot performance. For these reasons, we design a new data prior from which we simulate numerous OD datasets for pretraining FoMo-OD.

Ideally the data prior should reflect distributions as general and diverse as seen in real-world datasets, however, “*finding a prior supporting a large enough subset of possible [data generating] functions isn’t trivial*” (Nagler, 2023). Surprisingly, in contrast, our initial attempt has been sufficient to achieve remarkable performance even with a relatively straightforward and simple-to-implement data prior, which we describe next.

**Inlier synthesis:** We simulate inliers by simply drawing from a Gaussian Mixture Model (GMM) with  $m$ -clusters in  $d$ -dimensions, with centers  $\mu_{jk} \in [-5, 5]$ ,  $j \in [m]$ ,  $k \in [d]$  and  $\text{diagonal}^2 \Sigma_j$  with entries in  $[-5, 5]$ . In each step of every epoch during pretraining, we create batch size  $B$  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 90%-ile of the GMM.

**Outlier synthesis:** Following the previous literature on outlier synthesis (Han et al., 2022), we generate *subspace* outliers by first drawing a subset of dimensions  $\mathcal{K}$  at random, where  $|\mathcal{K}| \leq d$ , and then generate  $S$  points from the corresponding “inflated” GMMs, which share the same centers  $\mu_j$ ’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 outside the 90%-ile of the original GMM. We label each sample based on its Mahalanobis distance computed analytically (see Property B.2 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 A.

**Remarks:** We emphasize once again that our model is not trained on **any** real-world data and rather, on purely synthetic data (although future work can combine existing benchmark OD datasets with synthesized data, as was done for Chronos (Ansari et al., 2024)). Notably, our GMM-based data prior can be seen as extremely basic. While it has been our intent to extend our preliminary attempt toward designing a sophisticated data prior for OD, we found to our surprise that even with such an elementary prior, FoMo-OD performs remarkably well against numerous SOTA baselines. Therefore, we present FoMo-OD using this effortless approach for its simplicity to showcase the prowess of PFNs for OD. Future work can employ BNNs and SCMs (Hollmann et al., 2023), and other outlier types (contextual, dependency, etc. (Steinbuss and Böhm, 2021)) toward a more comprehensive data prior.

### 3.2 (PRE)TRAINING AND INFERENCE

**Model (Pre)Training (Once, Offline):** FoMo-OD is a Prior-data Fitted Network (PFN, see Section 2.2) based on the Transformer architecture. In the synthetic prior-data fitting phase, it is trained on datasets drawn from our OD data prior for tabular data that we introduced in Section 3.1. Each dataset is simulated from a different GMM configuration based on randomly drawn parameters, and consists of varying number of training samples and dimensions to capture the diversity in real-world tabular datasets. Detailed steps are outlined in Algo. 1 in Appendix C.2, and described as follows.

Each time, we first draw a hypothesis (i.e. GMM configuration) uniformly at random, that is  $\phi = \{d \in [D], m \in [M], \{\mu_j\}_{j=1}^m \in [-5, 5]^d, \{\Sigma_j\}_{j=1}^m; \text{diag}(\Sigma_j) \in [-5, 5]^d\}$ , and then generate a 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  $\theta$  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. (2)). During

<sup>2</sup>In our early experiments, we found no difference in terms of test performance on synthetic datasets between using diagonal and non-diagonal  $\Sigma$ , however, it is easier to compute the inverse of diagonal  $\Sigma$  for generation.

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 current implementation, we set  $n_L = 500$ , and  $n_U = 5,000$ .

FoMo-OD is trained on 200,000 batches (200 epochs  $\times$  1,000 steps/epoch) of  $B = 8$  generated datasets in each batch. While this pretraining phase can be expensive, it is done *only once, offline*. Moreover, we introduce several scalability improvements to speed up pretraining, as discussed later in Section 3.3. Full details on the training and implementation of FoMo-OD are given in Appendix C.

**Zero-shot Inference (on Unseen Dataset):** During the inference phase, our pretrained-in-advance FoMo-OD can be employed on any unseen real-world dataset. In fact, we apply the same single pretrained network on all benchmark datasets in our experiments in this paper.

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 entire training data is passed as context, FoMo-OD leverages in-context learning (ICL) (Xie et al., 2021; Garg et al., 2022) for inference. The **key** contribution of FoMo-OD goes beyond eliminating gradient-based model training for a new dataset: because no model training is required, one thus neither needs to choose any specific OD model to train, nor grapple with tuning any hyperparameters of the said model—rendering model selection an obsolete concern for the future of OD. Additionally, the speedy, easily parallelizable inference (for *less-than-a-second* per test sample) is then the “icing on the cake”.

For a visual summary, Figure 1 (right) illustrates (top) pretrain & (bottom) test phases of FoMo-OD.

### 3.3 ARCHITECTURE AND SCALABILITY

**Architecture and sample-to-sample attention:** Like existing PFNs in the literature, FoMo-OD is based on the Transformer architecture (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 adaptations of TabPFN (Hollmann et al., 2023), which (1) computes self-attention among all the training samples but only *cross*-attention from test samples to the training samples, (2) enables variable feature dimensionality by zero-padding, and (3) randomly rotates input samples while omitting positional encodings to achieve model invariance to sample permutations in the dataset. We defer the architecture details to the original papers.

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 linear transformations to a key  $\mathbf{k}_i$ , query  $\mathbf{q}_i$  and value  $\mathbf{v}_i$  based on which the  $i$ -th output is computed by weighing all  $\mathbf{v}_j$ ’s by the normalized dot product between  $\mathbf{q}_i$  and all the  $\mathbf{k}_j$ ’s (i.e. sample-to-sample dot product similarity) as

$$\mathbf{z}_i = \sum_{j=1}^n \text{softmax}(\{\langle \mathbf{q}_i, \mathbf{k}_{j'} \rangle\}_{j'=1}^n)_j \cdot \mathbf{v}_j. \quad (3)$$

The sample-to-sample attention is intriguing from the perspective of OD: Many classical OD algorithms (Aggarwal, 2013) are based on nonparametrics; in particular, they make use of the distances to the  $k$  *nearest* neighbors ( $k$ NNs) of a point to compute its outlierness (where  $k$  is a critical hyperparameter (HP)). One can think of FoMo-OD as mimicking non-parametric models but by using parametric attention mechanisms. Interestingly, PFNs are much more robust and flexible than  $k$ NN based OD approaches, for (1) sample-to-sample relations are not pre-specified but rather learned through attention weights, and thus (2) they are not limited to just the nearest neighbors but rather can *learn which* training points are worth attending to, and last but not least (3) as attention is dataset-wide across all points, there is no need for specifying a cut-off HP value like  $k$ , to which most  $k$ NN based OD techniques are sensitive to (Aggarwal and Sathe, 2015; Campos et al., 2016; Goldstein and Uchida, 2016; Ding et al., 2022)—to reiterate, algorithm & HP selection is bygone with FoMo-OD.

While intuitively beneficial for OD, “vanilla” attention among the training samples incurs quadratic complexity. To be able to seize the benefits with scale, we incorporate a scalable architecture to our design, as we describe next. The scale up also unlocks a larger context (i.e. dataset) size for FoMo-OD, enabling its pretraining on larger datasets for potentially better generalization.

**Scaling up attention with “routers”:** The  $\mathcal{O}(n^2)$  quadratic sample complexity at pretraining presents an obstacle for achieving high performance at inference. From dataset size perspective, it limits pretraining to relatively small training datasets. From context size perspective, it limits in-context learning that typically benefits from longer context lengths (Xie et al., 2021).

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

Concretely, the representatives first aggregate information from all samples by serving as query in multi-head self-attention (MSA) and the embedding array of all samples becomes both key and value:

$$\mathbf{M} = \text{MSA}_1(\mathbf{R}, \mathbf{Z}, \mathbf{Z}), \quad (4)$$

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

$$\hat{\mathbf{Z}} = \text{MSA}_2(\mathbf{Z}, \mathbf{M}, \mathbf{M}). \quad (5)$$

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 training samples’ embeddings, computed in the described manner across layers, which finally feed into the prediction head for estimating each test sample’s PPD at the output layer.

**Scaling up (pre)training data synthesis with linear transforms:** Besides the scalability challenge associated with architecture/attention, another computational challenge in pretraining FoMo-0D arises from drawing samples from the data prior. That is, generating samples from a pre-specified data distribution requires considerable time, especially in high dimensions<sup>3</sup>, provided the large number of datasets we sample (concretely, a batch size of 8 datasets over 1,000 steps each for 200 epochs).

To give an idea, sampling a dataset with  $n = 10,000$  points in  $d = 100$  dimensions using 10 CPUs in parallel takes  $\approx 0.4$  seconds (see Appendix Figure 7). Across 200 training epochs with 1,000 steps each, it adds up to more than 177 hours just to generate 1.6 million datasets on-the-fly. Of course, one can trade storage with compute-time by generating all these datasets apriori via massive parallelism. Nevertheless, synthetic data generation demands considerable time (and/or storage).

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

**Second**, and more innovatively, we propose *reuse at dataset level via transformation*: that is, having generated one unique dataset  $\mathbf{X} \in \mathbb{R}^{n \times d}$  from 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  $\mathbf{b} \in \mathbb{R}^d$  (see Appendix B.1).<sup>4</sup> This simple yet efficient transformation creates a new dataset, akin to one being drawn from another 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 materialize these parameters but only transform the dataset. As we show in the following, such transformations preserve the Mahalanobis distances as well as the percentile thresholds for labeling points as inlier/outlier. Details and proofs are given in Appendix B.

<sup>3</sup>This is because the inverse of the  $(d \times d)$  covariance matrix plays a crucial role in the process of generating samples from a GMM, which has  $\mathcal{O}(d^3)$  time complexity. (It is also partly the reason why we use diagonal  $\boldsymbol{\Sigma}_j$ ’s in our data prior.) In addition, Mahalanobis distance for labeling inliers/outliers also requires the inverse.

<sup>4</sup>In practice, we apply the linear transform on the subspace of inflated features only, wherein inliers and outliers are defined, which remains to be a multi-variate GMM.

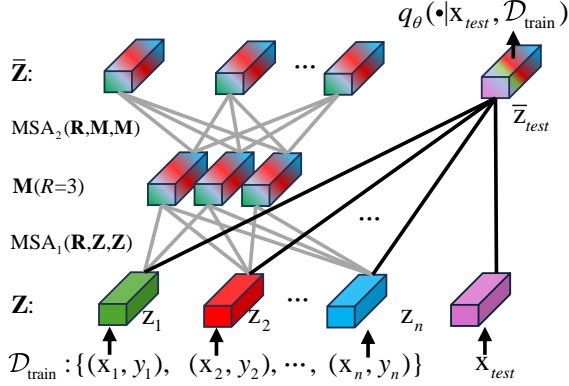


Figure 2: FoMo-0D architecture employs the “router mechanism” for scalable attention.

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

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

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

To reach 8K unique datasets for each epoch, we first generate 500 datasets from different GMMs (with varying configurations), and then employ 15 different linear transformations to each unique dataset by varying  $\mathbf{W}$  and  $\mathbf{b}$ . Drawing each  $(\mathbf{W}, \mathbf{b})$  takes  $\approx 0.02$  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 a transformed dataset offers  $20\times$  speed-up compared to generating one (0.02 vs. 0.4 seconds).

## 4 EXPERIMENTS

### 4.1 SETUP

We present the experiment setup briefly, including important notes on data synthesis, real-world datasets, baselines, metrics and HPs. For additional details, we refer to Appendix D.

**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  $\{\Sigma_j\}_{j=1}^m$  ( $\text{diag}(\Sigma_j) \in [-5, 5]^d$ ). 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.

**Real-world Benchmark Datasets:** While pretraining is purely on synthetic datasets, we evaluate FoMo-OD on 57 real-world datasets from the ADBench benchmark (Han et al., 2022) (see Table 15).

We use 5 train/test splits of each dataset via different seeds and report mean performance and standard deviation. Note that the baselines require model re-training and inference for each  $\mathcal{D}_{\text{train}}/\mathcal{D}_{\text{test}}$  split, while FoMo-OD uses the splits only for inference as  $\mathcal{D}_{\text{train}}$  is merely passed as context.

**Baselines:** We compare FoMo-OD against 26 baselines, from classical/shallow methods to modern/deep models. The baselines are 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. We defer to the original paper for additional details.

**Model Implementation:** We trained our final model for 200,000 steps with a batch size of 8 datasets. That is, our FoMo-OD is trained on 1,600,000 synthetically generated datasets. This training takes about 25 hours on 1 GPU (Nvidia RTX A6000). Each dataset had a fixed size of 10,000 samples, with  $|\mathcal{D}_{\text{train}}| \in [n_L = 500, n_U = 5000]$ , and the rest used as  $\mathcal{D}_{\text{test}}$  with *balanced* number of inliers and outliers. Other implementation details of FoMo-OD, including the training algorithm, model architecture, data synthesis and reuse, and hardware are provided in Appendix C.

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

To compare methods, we compute their rank on each dataset (lower is better), and present **average rank** across datasets. This is an alternative to the average metric (e.g. AUROC), which is not meaningful when task difficulties and hence metric values vary widely. In addition, we perform significance tests to compare two methods statistically, using the one-sided paired Wilcoxon signed rank test (Demšar, 2006) between FoMo-OD and a baseline based on the performances across all datasets and report the  $p$ -values. We consider results to be significant at 0.05 following convention.

**Hyperparameters (HPs):** Importantly, Livernoche et al. (2024) picked for each baseline the best-performing set of HPs as recommended by the authors in their original paper. As for their own DTE, which behaves similar to  $k$ NN, they use  $k = 5$  and set the *same*  $k$  for the  $k$ NN baseline (Ramswamy et al., 2000) to be consistent. However, it is well known that  $k$ NN is sensitive to the value of  $k$  (Aggarwal and Sathe, 2015), and so are many other OD models to their respective HPs (Campos et al., 2016; Goldstein and Uchida, 2016; Zhao et al., 2021; Ding et al., 2022).



Therefore, we compare to the performance results of these baselines as imported from DTE’s Tables 13, 14 and 15, respectively for AUROC, F1, and AUPR (Livernoche et al., 2024). In addition, we also compare to the **top-4**<sup>5</sup> best performing baselines (in order: DTE-NP,  $k$ NN, ICL, and DTE-C) on their *average* performance across a list of different HP settings (which reflects their *expected* performance under HP values selected at random, in the absence of any other prior knowledge), which is the recommended approach by Goldstein and Uchida (2016) “*to get a fair evaluation when comparing [OD] algorithms*”. We annotate the method name with <sup>avg</sup> for the version with performance averaged over varying HPs. The detailed list of HP values for each top baseline is given in Appendix D.4.

Overall, we compare FoMo-OD to 30 baselines; 26 from Livernoche et al. (2024) and <sup>avg</sup> of the top-4.

## 4.2 RESULTS

**Detection performance:** Table 1 presented the comparison of FoMo-OD w/  $D = 100$  to all baselines w.r.t. average rank across datasets as well as pairwise Wilcoxon signed rank tests based on AUROC (for full results on all datasets and all metrics, see Appendix G). We find that among 30 baselines and 2 variants of FoMo-OD (w/  $D = 100$  and  $D = 20$ ), FoMo-OD w/  $D = 100$  *performs as well as the 2nd best model* ( $k$ NN with default HP;  $k = 5$ ) on all datasets. While DTE-NP outperforms FoMo-OD with author-recommended  $k = 5$ , we find that DTE-NP<sup>avg</sup> is on par with FoMo-OD.

Against all other baselines, we obtain notably large  $p$ -values. Typically,  $p > 0.05$  implies no statistical difference between two methods. On the other hand, the large  $p$ -values we obtain that are often larger than 0.50 suggest that the odds are tilted towards FoMo-OD to outperform.

FoMo-OD w/  $D = 100$  performs statistically no different from **all** baselines on datasets with  $d \leq 100$  (i.e., “at its own game” when pretraining data dimensions align with real-world datasets), while it *outperforms the majority of baselines* ( $p > 0.95$ ). These results also hold on datasets with  $d \leq 500$ .

Table 2 shows similar results for FoMo-OD w/  $D = 20$ , which is pretrained on datasets with considerably fewer dimensions. Even in this limited setting, its performance is remarkable: against 30 baselines, it performs on par with the *3rd* best baseline (ICL, with default HP). The  $p$ -value is even larger (0.437) when compared to ICL<sup>avg</sup>. Moreover, on datasets with  $d \leq 20$  which align with its pretraining data, all  $p$ -values are larger than 0.5, where it outperforms the top *5th* baseline and the majority of others. These are outstanding results for a model pretrained purely on synthetic datasets from a simple data prior in small dimensions, showcasing the prowess of PFNs for OD.

Table 2: Comparison of methods across datasets. (top row) Rank w.r.t. AUROC performance avg.’ed over 57 datasets is presented for FoMo-OD (with  $D = 20$ ), **top-10** baselines with default HPs, and **top-4**<sup>5</sup> baselines with performance **avg.**’ed over varying HPs (denoted w/ <sup>avg</sup>); followed by  $p$ -values of the pairwise Wilcoxon signed rank test, comparing FoMo-OD to each baseline (from top to bottom) over All (57) datasets, those (24) w/  $d \leq 20$  and (38) datasets w/  $d \leq 50$  dimensions, respectively. Even with small  $D = 20$ , FoMo-OD performs as well as (i.e., statistically no different at 0.05 from) *the top 3rd baseline* (ICL, w/  $p = 0.089$ ) across All datasets, while it *outperforms the top 5th (LOF) and onward baselines* over datasets w/  $d \leq 20$  (aligned w/ pretraining where  $D = 20$ ) and  $d \leq 50$  (generalizing beyond pretraining). (setting:  $D = 20$ ,  $P = 50$ ,  $R = 500$ , train/inference context size=5K, no data transformation)

	FoMo-OD	DTE-NP	$k$ NN	ICL	DTE-C	LOF	CBLOF	Feat.Bag.	SLAD	DDPM	OCSVM	DTE-NP <sup>avg</sup>	$k$ NN <sup>avg</sup>	ICL <sup>avg</sup>	DTE-C <sup>avg</sup>
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
All	-	<b>0.001</b>	<b>0.019</b>	0.089	0.159	0.394	0.434	0.703	0.516	0.752	0.679	<b>0.007</b>	0.062	0.437	1.0
$d \leq 20$	-	0.572	0.789	0.968	0.616	<b>0.993</b>	<b>0.989</b>	<b>1.0</b>	<b>0.978</b>	0.906	<b>0.992</b>	0.813	0.924	<b>0.999</b>	<b>1.0</b>
$d \leq 50$	-	0.347	0.794	0.893	0.946	<b>0.997</b>	<b>0.988</b>	<b>1.0</b>	<b>0.963</b>	<b>0.994</b>	<b>0.986</b>	0.574	0.847	<b>0.995</b>	<b>1.0</b>

Figure 3 shows the distribution of ranks across datasets for each of the 32 methods. While paired significant tests are the most conclusive, FoMo-OD achieves relatively small average rank as well as notably low ranks across datasets that is also visually better than the majority of the baselines.

**Running time:** Table 3 presents the total training time and the average inference time per test sample, as measured on our largest benchmark dataset, for FoMo-OD and the top-3 baselines. Given

<sup>5</sup>To rank the baselines, we compute the  $26 \times 26$  pairwise  $p$ -values based on the Wilcoxon signed rank test, as shown in Appendix Figure 16, and rank the baselines w.r.t. their mean  $p$ -value.

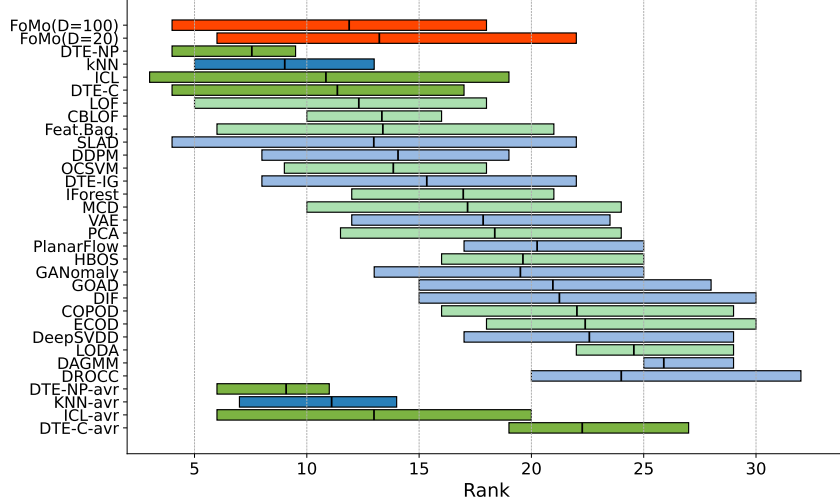


Figure 3: (best in color) Rank (w.r.t. AUROC performance, lower is better) distribution across datasets shown via boxplots for (from top to bottom) FoMo-OD in red, all 26 baselines as ordered by mean  $p$ -value<sup>5</sup> (shallow and deep baselines resp. in green and blue), and top-4 baselines’ <sup>avg</sup> variants.

a new dataset, FoMo-OD bypasses model training (and HP tuning) and directly performs inference, with an average of 7.7 ms per sample (see Appendix Figure 6). In comparison, all baseline methods need to train on each individual dataset preceding inference. This training time can be high for deep learning based models like ICL, and further compounded with training *multiple* models for hyperparameter tuning purposes. Even for non-parametric and/or shallow models like  $k$ NN and DTE-NP (which queries  $k$  nearest neighbors), the training involves various data pre-processing steps such as constructing a tree-like data structure for fast (often approximate)  $k$ NN distance querying.

Table 3: Training and inference time (in milliseconds) comparison between FoMo-OD and the top-3<sup>5</sup> baselines (w/ default HPs, *excluding* the time for model selection/hyperparameter optimization) on our largest dataset (namely, *donors*, see Appendix Table 15).

Method	FoMo-OD	DTE-NP	$k$ NN	ICL
Training time (total)	none	56.83	1433.74	186461.48
Inference time (per sample)	7.7	0.76	0.17	0.01

### 4.3 ABLATION ANALYSES

Due to space limits, we present the detailed ablation analyses in Appendix E. We discuss the effect of  $D$  in E.1, the cost and performance of varying  $R$  in E.2 and E.3, the context size in E.4, the reuse periodicity  $P$  in E.5, the effect of data transformation  $T$  on performance and speed up in E.6 and E.7, data diversity and prolonged training in E.8, and quantile transformation on ADBench in E.9.

## 5 RELATED WORK

Due to space limits, we present the detailed related work in Appendix J.

## 6 CONCLUSION

We introduced FoMo-OD, **the first foundation model for outlier detection (OD)** on tabular data. It capitalizes on the in-context learning ability of a Transformer model pretrained on a large number of synthetic datasets that can then perform zero-shot inference on a new dataset by directly passing it as input context. FoMo-OD breaks new ground by fully abolishing notoriously-hard model selection. Further, FoMo-OD offers extremely fast inference thanks to a mere single forward pass. Against **26** baselines on **57** public datasets from diverse domains, FoMo-OD performs on par with the *2nd* best baseline, while significantly outperforming the majority of baselines. We leave improving the data prior and extending beyond tabular OD, among others, as future directions. For a detailed discussion on limitations and future directions, we refer to Appendix K.

## REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- Steven Adriaensen, Herilalaina Rakotoarison, Samuel Müller, and Frank Hutter. 2024. Efficient bayesian learning curve extrapolation using prior-data fitted networks. *Advances in Neural Information Processing Systems* 36 (2024).
- Charu C. Aggarwal. 2013. *Outlier Analysis*. Springer.
- Charu C. Aggarwal and Saket Sathe. 2015. Theoretical foundations and algorithms for outlier ensembles. *Acm sigkdd explorations newsletter* 17, 1 (2015), 24–47.
- Charu C. Aggarwal and Saket Sathe. 2017. *Outlier Ensembles - An Introduction*. Springer.
- Samet Akcay, Amir Atapour-Abarghouei, and Toby P Breckon. 2019. Ganomaly: Semi-supervised anomaly detection via adversarial training. In *ACCV*. Springer, 622–637.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems* 35 (2022), 23716–23736.
- Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, et al. 2024. Chronos: Learning the language of time series. *arXiv preprint arXiv:2403.07815* (2024).
- Lion Bergman and Yedid Hoshen. 2020. Classification-Based Anomaly Detection for General Data. In *International Conference on Learning Representations*.
- Markus M. Breunig, Hans-Peter Kriegel, Raymond T. Ng, and Jörg Sander. 2000. LOF: Identifying Density-Based Local Outliers. In *International Conference on Management of Data*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, Vol. 33. 1877–1901.
- Guilherme O Campos, Arthur Zimek, Jörg Sander, Ricardo JGB Campello, Barbora Micenková, Erich Schubert, Ira Assent, and Michael E Houle. 2016. On the evaluation of unsupervised outlier detection: measures, datasets, and an empirical study. *Data mining and knowledge discovery* 30 (2016), 891–927.
- George Casella and Roger Berger. 2024. *Statistical inference*. CRC Press.
- Raghavendra Chalapathy and Sanjay Chawla. 2019. Deep learning for anomaly detection: A survey. *arXiv preprint arXiv:1901.03407* (2019).
- Janez Demšar. 2006. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine learning research* 7 (2006), 1–30.
- Xueying Ding, Lingxiao Zhao, and Leman Akoglu. 2022. Hyperparameter sensitivity in deep outlier detection: Analysis and a scalable hyper-ensemble solution. *Advances in Neural Information Processing Systems* 35 (2022), 9603–9616.
- Xueying Ding, Yue Zhao, and Leman Akoglu. 2024. Fast Unsupervised Deep Outlier Model Selection with Hypernetworks. *ACM SIGKDD* (2024).

- Samuel Dooley, Gurnoor Singh Khurana, Chirag Mohapatra, Siddhartha Venkat Naidu, and Colin White. 2023. ForecastPFN: Synthetically-Trained Zero-Shot Forecasting. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Xuefeng Du, Yiyu Sun, Jerry Zhu, and Yixuan Li. 2024. Dream the impossible: Outlier imagination with diffusion models. *Advances in Neural Information Processing Systems* 36 (2024).
- Sepideh Esmaeilpour, Bing Liu, Eric Robertson, and Lei Shu. 2022. Zero-shot out-of-distribution detection based on the pre-trained model CLIP. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 36. 6568–6576.
- Benjamin Feuer, Niv Cohen, and Chinmay Hegde. 2023. Scaling TabPFN: Sketching and Feature Selection for Tabular Prior-Data Fitted Networks. In *NeurIPS 2023 Second Table Representation Learning Workshop*.
- Benjamin Feuer, Robin Tibor Schirrmester, Valeriia Cherepanova, Chinmay Hegde, Frank Hutter, Micah Goldblum, Niv Cohen, and Colin White. 2024. TuneTables: Context Optimization for Scalable Prior-Data Fitted Networks. *arXiv:2402.11137 [cs.LG]*
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. 2022. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems* (2022).
- Markus Goldstein and Andreas Dengel. 2012. Histogram-based outlier score (HBOS): A fast unsupervised anomaly detection algorithm. *KI-2012: poster and demo track 1* (2012), 59–63.
- Markus Goldstein and Seiichi Uchida. 2016. A comparative evaluation of unsupervised anomaly detection algorithms for multivariate data. *PloS one* 11, 4 (2016).
- Sachin Goyal, Aditi Raghunathan, Moksh Jain, Harsha Simhadri, and Prateek Jain. 2020. DROCC: Deep Robust One-Class Classification. In *Proceedings of the 37th International Conference on Machine Learning (ICML’20)*. JMLR.org, Article 348, 11 pages.
- Zhaopeng Gu, Bingke Zhu, Guibo Zhu, Yingying Chen, Ming Tang, and Jinqiao Wang. 2024. AnomalyGPT: Detecting industrial anomalies using large vision-language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 1932–1940.
- A. Gut. 2009. *An Intermediate Course in Probability*. Springer New York.
- Songqiao Han, Xiyang Hu, Hailiang Huang, Minqi Jiang, and Yue Zhao. 2022. Adbench: Anomaly detection benchmark. *Advances in Neural Information Processing Systems* 35 (2022).
- Haoyang He, Jiangning Zhang, Hongxu Chen, Xuhai Chen, Zhishan Li, Xu Chen, Yabiao Wang, Chengjie Wang, and Lei Xie. 2024. A diffusion-based framework for multi-class anomaly detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 8472–8480.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851.
- Hadi Hojjati, Thi Kieu Khanh Ho, and Narges Armanfard. 2022. Self-supervised anomaly detection: A survey and outlook. *arXiv preprint arXiv:2205.05173* (2022).
- Noah Hollmann, Samuel Müller, Katharina Eggenberger, and Frank Hutter. 2023. TabPFN: A Transformer That Solves Small Tabular Classification Problems in a Second. In *The Eleventh International Conference on Learning Representations*.
- Jongheon Jeong, Yang Zou, Taewan Kim, Dongqing Zhang, Avinash Ravichandran, and Onkar Dabeer. 2023. Winclip: Zero-/few-shot anomaly classification and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 19606–19616.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361* (2020).

- Diederik P Kingma. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).
- Diederik P. Kingma and Jimmy Ba. 2017. Adam: A Method for Stochastic Optimization. *arXiv:1412.6980 [cs.LG]*
- Aodong Li, Chen Qiu, Marius Kloft, Padhraic Smyth, Maja Rudolph, and Stephan Mandt. 2023. Zero-Shot Anomaly Detection via Batch Normalization. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Aodong Li, Yunhan Zhao, Chen Qiu, Marius Kloft, Padhraic Smyth, Maja Rudolph, and Stephan Mandt. 2024. Anomaly Detection of Tabular Data Using LLMs. *arXiv preprint arXiv:2406.16308* (2024).
- Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2008. Isolation Forest. In *2008 Eighth IEEE International Conference on Data Mining*. 413–422.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual instruction tuning. *Advances in neural information processing systems* 36 (2024).
- Victor Livernoche, Vineet Jain, Yashar Hezaveh, and Siamak Ravanbakhsh. 2024. On Diffusion Modeling for Anomaly Detection.. In *ICLR*.
- Philipp Liznerski, Lukas Ruff, Robert A Vandermeulen, Billy Joe Franks, Klaus-Robert Müller, and Marius Kloft. 2022. Exposing outlier exposure: What can be learned from few, one, and zero outlier images. *arXiv preprint arXiv:2205.11474* (2022).
- Junwei Ma, Valentin Thomas, Guangwei Yu, and Anthony L. Caterini. 2024. In-Context Data Distillation with TabPFN. *CoRR* abs/2402.06971 (2024).
- Martin Q Ma, Yue Zhao, Xiaorong Zhang, and Leman Akoglu. 2023. The need for unsupervised outlier model selection: A review and evaluation of internal evaluation strategies. *ACM SIGKDD Explorations Newsletter* 25, 1 (2023), 19–35.
- Samuel Müller, Matthias Feurer, Noah Hollmann, and Frank Hutter. 2023. PFNs4BO: in-context learning for Bayesian optimization. In *International Conference on Machine Learning*.
- Samuel Müller, Noah Hollmann, Sebastian Pineda-Arango, Josif Grabocka, and Frank Hutter. 2022. Transformers Can Do Bayesian Inference. *ICLR* (2022).
- Thomas Nagler. 2023. Statistical Foundations of Prior-Data Fitted Networks.. In *ICML (Proceedings of Machine Learning Research, Vol. 202)*. PMLR.
- Radford M Neal. 2012. *Bayesian learning for neural networks*. Vol. 118. Springer Science & Business Media.
- Guansong Pang, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. 2021. Deep learning for anomaly detection: A review. *ACM computing surveys (CSUR)* 54, 2 (2021), 1–38.
- Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. 2021. Deep learning on a data diet: Finding important examples early in training. *Advances in neural information processing systems* 34 (2021), 20596–20607.
- Judea Pearl. 2009. *Causality*. Cambridge University Press.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*.
- Sridhar Ramaswamy, Rajeev Rastogi, and Kyuseok Shim. 2000. Efficient algorithms for mining outliers from large data sets. In *ACM SIGMOD international conference on management of data*.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. In *International conference on machine learning*. Pmlr, 8821–8831.

- Allan Raventós, Mansheej Paul, Feng Chen, and Surya Ganguli. 2024. Pretraining task diversity and the emergence of non-bayesian in-context learning for regression. *Advances in Neural Information Processing Systems* (2024).
- Danilo Rezende and Shakir Mohamed. 2015. Variational Inference with Normalizing Flows. In *Proceedings of the 32nd International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 37)*. PMLR, Lille, France, 1530–1538.
- Lukas Ruff, Jacob R Kauffmann, Robert A Vandermeulen, Grégoire Montavon, Wojciech Samek, Marius Kloft, Thomas G Dietterich, and Klaus-Robert Müller. 2021. A unifying review of deep and shallow anomaly detection. *Proc. IEEE* 109, 5 (2021), 756–795.
- Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. 2018. Deep One-Class Classification. In *International Conference on Machine Learning*. 4393–4402.
- Tom Shenkar and Lior Wolf. 2022. Anomaly Detection for Tabular Data with Internal Contrastive Learning. In *International Conference on Learning Representations*.
- Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. 2022. Beyond neural scaling laws: beating power law scaling via data pruning. *Advances in Neural Information Processing Systems* 35 (2022), 19523–19536.
- Georg Steinbuss and Klemens Böhm. 2021. Benchmarking unsupervised outlier detection with realistic synthetic data. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 15, 4 (2021), 1–20.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*. 5998–6008.
- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2021. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080* (2021).
- Hongzuo Xu, Yijie Wang, Juhui Wei, Songlei Jian, Yizhou Li, and Ning Liu. 2023. Fascinating supervisory signals and where to find them: Deep anomaly detection with scale learning. In *International Conference on Machine Learning*. PMLR, 38655–38673.
- Jaemin Yoo, Tiancheng Zhao, and Leman Akoglu. 2023. Data Augmentation is a Hyperparameter: Cherry-picked Self-Supervision for Unsupervised Anomaly Detection is Creating the Illusion of Success. *Trans. Mach. Learn. Res.* 2023 (2023).
- Sangwoong Yoon, Young-Uk Jin, Yung-Kyun Noh, and Frank Park. 2023. Energy-based models for anomaly detection: A manifold diffusion recovery approach. *Advances in Neural Information Processing Systems* 36 (2023), 49445–49466.
- Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. 2022. Scaling vision transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*.
- Yunhao Zhang and Junchi Yan. 2023. Crossformer: Transformer Utilizing Cross-Dimension Dependency for Multivariate Time Series Forecasting.. In *ICLR*.
- Yue Zhao and Leman Akoglu. 2024. Toward unsupervised outlier model selection. In *International Conference on Automated Machine Learning (AutoML)*.
- Yue Zhao, Ryan Rossi, and Leman Akoglu. 2021. Automatic unsupervised outlier model selection. *Advances in Neural Information Processing Systems* 34 (2021), 4489–4502.
- Yue Zhao, Sean Zhang, and Leman Akoglu. 2022. Toward unsupervised outlier model selection. In *2022 IEEE International Conference on Data Mining (ICDM)*. IEEE, 773–782.
- Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Dae ki Cho, and Haifeng Chen. 2018. Deep Autoencoding Gaussian Mixture Model for Unsupervised Anomaly Detection. In *International Conference on Learning Representations*.

## A ILLUSTRATION OF SYNTHETIC DATA IN 2- $d$

We visualize our synthetic data in Figure 4, with 3 randomly created 2- $d$  GMMs with the number of clusters ( $N = 1, 2, 3$ ). We choose the 80 $th$  percentile as the criterion, such that inliers are samples drawn from the GMM and within the 80 $th$  percentile, and outliers are samples drawn from the inflated GMMs and outside of the 80 $th$  percentile.

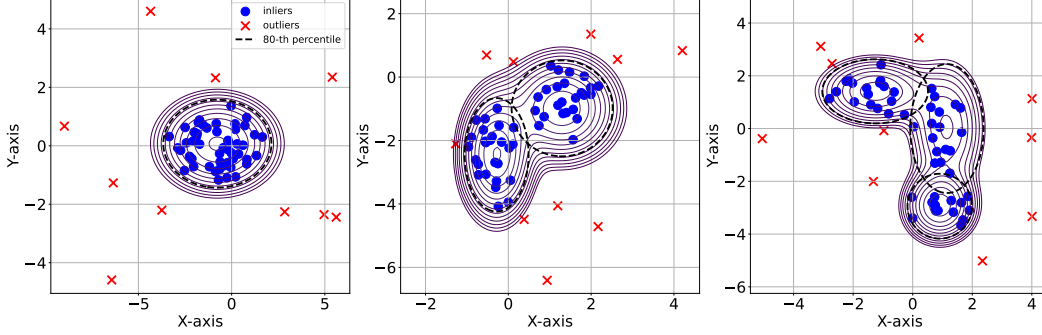


Figure 4: Illustration of synthetic data in 2D with 80 $th$  percentile as the criterion.

## B LINEAR TRANSFORM FOR SCALABLE GMM DATA SYNTHESIS

### B.1 DEFINITIONS

**Definition 1 (Gaussian Mixture Model)** We denote an  $m$ -cluster  $d$ -dimension Gaussian Mixture Model as  $\mathcal{G}_m^d = \{(w_j, \mu_j, \Sigma_j)\}_{j=1}^m$ , which is the weighted sum of  $m$  Gaussian distributions:

$$p(\mathbf{x}) = \sum_{j=1}^m w_j \cdot g(\mathbf{x} | \mu_j, \Sigma_j), \quad (6)$$

where  $w_j \in \mathbb{R}^+$  is the weight for the  $j$ -th Gaussian  $\mathcal{N}(\mu_j, \Sigma_j)$  with  $\sum_{j=1}^m w_j = 1$ , and  $g(\cdot | \mu_j, \Sigma_j)$  is the density of the  $j$ -th component/cluster, with mean/center  $\mu_j \in \mathbb{R}^d$  and covariance  $\Sigma_j \in \mathbb{R}^{d \times d}$  being positive semi-definite, such that  $\mathbf{x}^T \Sigma_j \mathbf{x} \geq 0$ , for all  $\mathbf{x} \in \mathbb{R}^d$ .

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

$$T(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}, \quad (7)$$

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  $T$ .

**Definition 3 (Mahalanobis Distance)** The Mahalanobis distance  $\text{dist}_M$  between a point  $\mathbf{x} \in \mathbb{R}^d$  and a Gaussian distribution  $\mathcal{N}(\mu, \Sigma)$  is defined as:

$$\text{dist}_M(\mathbf{x}) = \sqrt{(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)}. \quad (8)$$

**Definition 4 ( $\chi_d^2$ -distribution)** The Chi-squared distribution  $\chi_d^2$  with  $d$  degrees of freedom is the distribution of the sum of squares of  $d$  independent standard Normal random variables.

### B.2 PROPERTIES

**Property B.1 (Lemma 5.3.2 (Casella and 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  $X_i \sim \chi_1^2$ , then  $\sum_{i=1}^d X_i \sim \chi_d^2$ .

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

*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:

$$\text{dist}_M^2(\mathbf{x}) = \mathbf{z}^T \mathbf{z} = \sum_{i=1}^d z_i^2 \quad (9)$$

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

### B.3 LEMMAS

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

*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 the transformed point  $T(\mathbf{x})$  we have:

$$\text{dist}_M(T(\mathbf{x})) = \sqrt{(T(\mathbf{x}) - (\mathbf{W}\boldsymbol{\mu}_j + \mathbf{b}))^T (\mathbf{W}\boldsymbol{\Sigma}\mathbf{W}^T)^{-1} (T(\mathbf{x}) - (\mathbf{W}\boldsymbol{\mu}_j + \mathbf{b}))} \quad (10)$$

$$= \sqrt{(\mathbf{W}(\mathbf{x} - \boldsymbol{\mu}_j))^T (\mathbf{W}\boldsymbol{\Sigma}\mathbf{W}^T)^{-1} (\mathbf{W}(\mathbf{x} - \boldsymbol{\mu}_j))} \quad (11)$$

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

$$= \sqrt{(\mathbf{x} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_j)} = \text{dist}_M(\mathbf{x}) . \quad (13)$$

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

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

$$\text{Prob}(X \leq \chi_d^2(n)) = \frac{\alpha}{100} . \quad (14)$$

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

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  $\alpha$ -th percentile of  $\mathcal{G}_m^d$ . Since  $\text{dist}_M(\mathbf{x})$  is preserved under  $T$  (see Lemma 1), then we conclude that the linear transform  $T$  with invertible  $\mathbf{W}$  preserves the percentiles of the GMM. ■

## C IMPLEMENTATION DETAILS

### C.1 HARDWARE

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

### C.2 TRAINING AND INFERENCE

We train our models for 200 epochs with the Adam optimizer (Kingma and 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).



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

**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, the number of training epochs  $E$ , the periodicity  $P$ , the number of unique datasets  $q$ , linear transformation  $T$ .

**Output** : A model  $q_\theta$  that will approximate the PPD

```

1 Initialize the neural network  $q_\theta$ ;
2 Initialize the epoch-level collection  $\mathcal{C}_E = []$ ;
3 for  $i \leftarrow 1$  to  $E$  do
4   if  $i \leq P$  then
5     Initialize an empty buffer  $\mathcal{B}_i = []$ ;
6     Initialize the dataset-level collection  $\mathcal{C}_q = []$ ;
7     for  $j \leftarrow 1$  to  $Q$  do
8       if  $j \leq q$  then
9         Step 1: sample  $D_j := \mathcal{D}_{\text{train}} \cup \{(\mathbf{x}_k, y_k)\}_{i=k}^{|\mathcal{D}_{\text{test}}|} \sim p(\mathcal{D})$ ;
10         $\mathcal{C}_q \leftarrow \mathcal{C}_q + [D_j]$ 
11      end
12    else
13       $j \leftarrow j \bmod q$ 
14       $D_j \leftarrow T(\mathcal{C}_q[j])$ 
15    end
16    Step 2: compute stochastic loss approximation  $\bar{\ell}_\theta = \sum_{k=1}^{|\mathcal{D}_{\text{test}}|} (-\log q_\theta(y_k | \mathbf{x}_k, \mathcal{D}_{\text{train}}))$ ;
17    Step 3: update parameters  $\theta$  with stochastic gradient descent on  $\nabla_\theta \bar{\ell}_\theta$ ;
18     $\mathcal{B}_i \leftarrow \mathcal{B}_i + [D_j]$ 
19  end
20   $\mathcal{C}_E \leftarrow \mathcal{C}_E + [\mathcal{B}_i]$ 
21 end
22 else
23    $i \leftarrow i \bmod P$ 
24    $\mathcal{B}_i \leftarrow \mathcal{C}_E[i]$ 
25   for  $j \leftarrow 1$  to  $Q$  do
26      $D_j \leftarrow T(\mathcal{B}_i[j])$ 
27     Perform Step 2 and Step 3
28   end
29 end
30 end

```

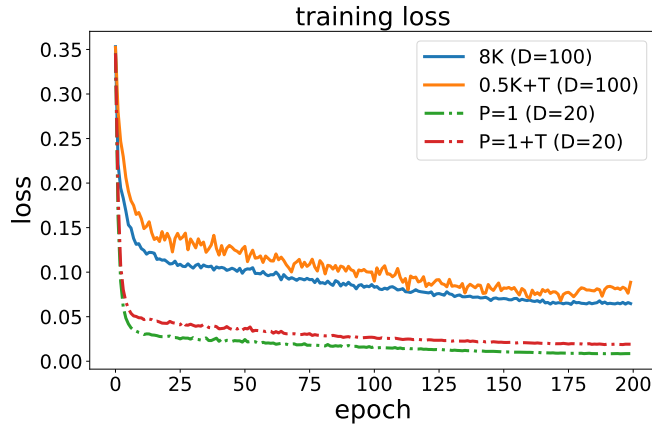


Figure 5: (best in color) Training loss of FoMo-0D ( $D = 100$ ) with 8K unique datasets/epoch (in blue) and using 0.5K unique + 7.5K transformed datasets/epoch (in orange), and FoMo-0D ( $D = 20$ ) with  $P = 1$  (in green) and  $P = 1$  with transformation (in red) over 200 epochs.

**Training loss** In Figure 5, 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 harder and thus result in a higher training loss. The training losses of FoMo-0D with  $D = 100$  are also higher than with  $D = 20$  since the subspace OD tasks are harder in higher dimensions.

**Inference time** Figure 8 (left) showed the inference time of FoMo-0D 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 6 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-0D (with  $R = 500$  routers) uses inference context size of 5K by default, which takes about 7.7 ms per test sample on average.

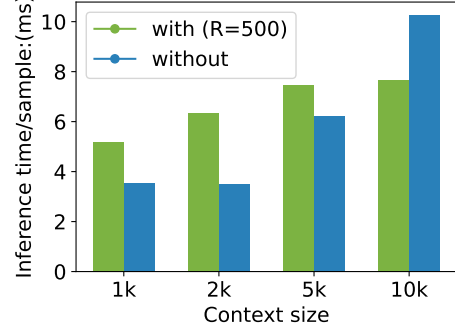


Figure 6: Inference time of FoMo-0D on GPU with vs. w/out router-based attention under varying context size.

## D DETAILED EXPERIMENT SETUP

### D.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  $\{\Sigma_j\}_{j=1}^m$  ( $\text{diag}(\Sigma_j) \in [-5, 5]^d$ ). 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.  $5 \times 100$ ) and transform each one  $15 \times$  using varying parameters  $(\mathbf{W}, \mathbf{b})$  as described in Section 3.3.<sup>6</sup> 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.

### D.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 15). 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.

<sup>6</sup>It 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  $\mathbf{b} \in [-1, 1]^d$ .

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-OD uses the splits only for inference as  $\mathcal{D}_{\text{train}}$  is merely passed as context.

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

### D.3 BASELINES

We compare FoMo-OD 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 and 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 and Hoshen, 2020), ICL (Shenkar and Wolf, 2022), SLAD (Xu et al., 2023), etc.), besides the four additional generative baselines: normalizing planar flows (Rezende and Mohamed, 2015), DDPM (Ho et al., 2020), VAE (Kingma, 2013) and GANomaly (Akçay 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.).

### D.4 HYPERPARAMETERS FOR BASELINES

Table 4 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 4: 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	$k \in \{5, 10, 20, 40, 50\}$

### D.5 RANKING THE 26 BASELINES

Figure 16 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.

### D.6 COMPARISON OF TOP-4 BASELINE VARIANTS WITH VARYING HP CONFIGURATIONS

Figure 17, 18, 19, 20 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  $\text{avg}$  model with the average performance across HPs. (Specifically for ICL,  $\text{learning\_rate}$  ( $\text{lr}$ )  $\in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$ ; and for others,  $\text{\#nearest-neighbors}$   $k \in \{5, 10, 20, 40, 50\}$ ). We find that for ICL,  $\text{lr} = 10^{-3}$  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  $\text{DTE}^{\text{avg}}$  and  $k\text{NN}^{\text{avg}}$  models across HP configurations perform subpar.

## D.7 SAMPLING TIME OF $d$ -DIMENSIONAL GMM

Figure 7 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.

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.

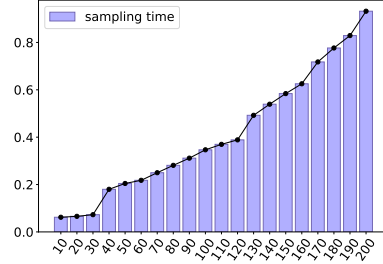


Figure 7: Sampling time (in seconds) of 10,000 points from GMMs with varying number of dimensions.

## E ABLATION ANALYSES

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

### E.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 5. 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 2 ( $D = 20$ ) 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 8.

### E.2 EFFECT OF ROUTERS ON COST

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

Figure 8(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.<sup>7</sup>

Similarly, memory cost with routers is considerably lower when using routers, especially for larger context sizes, as shown in Figure 8(middle).

<sup>7</sup>Note that the inference time is reported on CPUs to show scalability. On GPUs, w/ 5K context size, see Appendix Figure 6, where 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 (4) and (5).

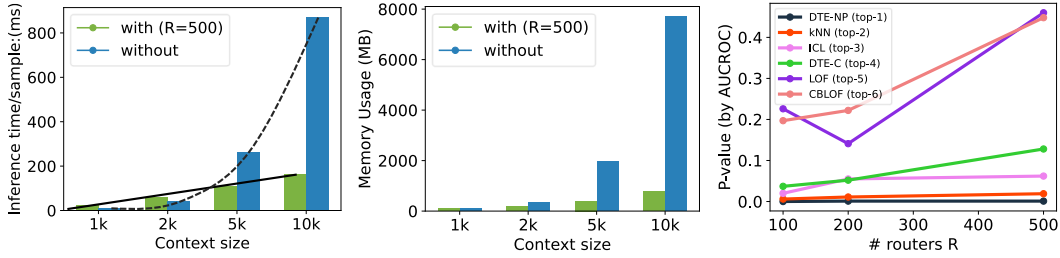


Figure 8: FoMo-0D w/ router mechanism saves time and memory while more #routers perform better, offering a cost-performance trade-off: (left) inference-time (ms) per sample and (middle) memory cost (MB) with & w/out routers by varying context size; (right) performance (based on  $p$ -value against top baselines, higher is better) vs. number of routers. (setting:  $D = 20$ ,  $P = 1$ )

### E.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 8(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.

### E.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 9. Interestingly, when training context size is large enough at 5K, inference with 10K samples generalizes beyond training with no significant difference (at 0.05) from other inference context sizes.

### E.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 3.3. 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  $d \leq 50$ .

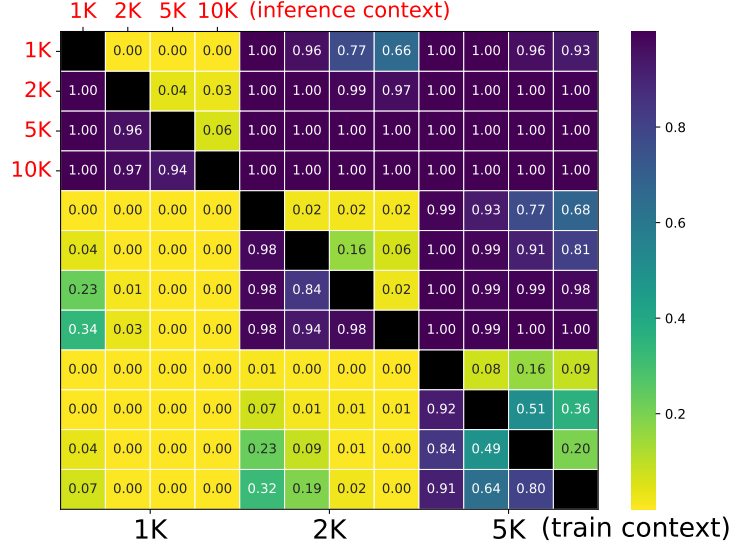


Figure 9:  $p$ -values of the pairwise Wilcoxon signed rank test between models (larger  $p$  implies col-method is better than row-method) w/ different context sizes for **training** (1K/2K/5K, 1st/2nd/3rd four grids, in **black**) and **inference** (1K/2K/5K/10K, every 1st/2nd/3rd/4th grid, in **red**): 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 FoMo-0D with  $D = 20$  remains no different from 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: $8 \times 50 = 400K$ )					$P = 100$ (#unique datasets: $8 \times 100 = 800K$ )				
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	<b>0.998</b>	0.572	0.789	0.968	0.616	<b>0.993</b>	0.439	0.678	0.953	0.550	<b>0.972</b>
$d \leq 50$	0.415	0.750	0.869	<b>0.962</b>	<b>0.999</b>	0.347	0.794	0.893	0.946	<b>0.997</b>	0.293	0.697	0.890	0.924	<b>0.994</b>

## E.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 different 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, where FoMo-0D with  $D = 20$  remains no different from the top 3rd baseline at 0.05 w.r.t. pairwise Wilcoxon signed rank test comparisons. (setting:  $D = 20$ ,  $R = 500$ , context size=5K,  $P = 1$ )

	w/out transformation $T$					w/ transformation $T$				
top-5	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	<b>0.998</b>	0.648	0.708	0.988	0.718	<b>0.955</b>
$d \leq 50$	0.415	0.750	0.869	<b>0.962</b>	<b>0.999</b>	0.264	0.382	0.971	0.900	<b>0.963</b>



## E.7 SPEED UP BY $T$

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

Figure 10 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  $\approx 25$  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 3 and Appendix Figure 6).

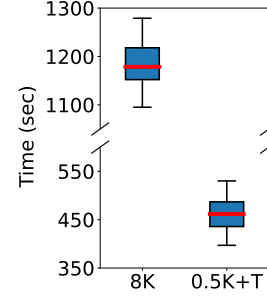


Figure 10: Runtime/epoch dist.n over 100 epochs for FoMo-0D ( $D=100$ ) with (left)  $P=100$ , i.e. 8K unique datasets/epoch vs. (right) 0.5K unique+7.5K transformed datasets/epoch.

## E.8 EFFECT OF DATA DIVERSITY AND PROLONGED TRAINING

### *How does FoMo-0D’s performance change by increasing pre-training 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 12, 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 11 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.

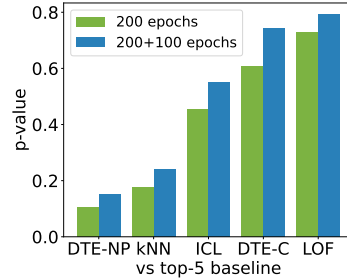


Figure 11:  $p$ -values increase with additional 100 epochs of pretraining, i.e. FoMo-0D w/  $D = 100$  performs better against top-5 baselines on datasets w/  $d \leq 100$ .

## E.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 13 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.

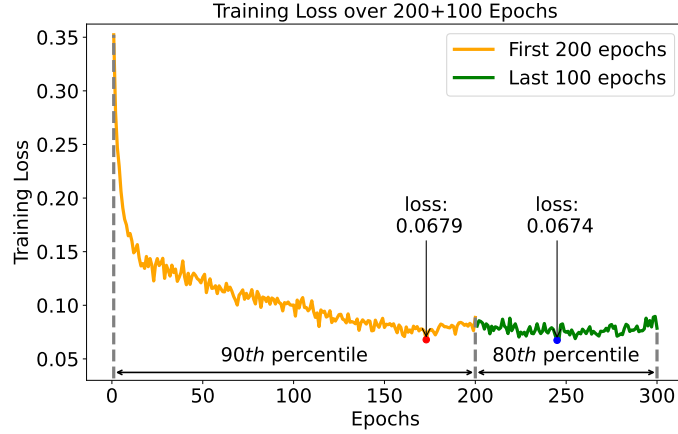


Figure 12: (best in color) Training loss of FoMo-0D ( $D = 100$ ) with 0.5K unique + 7.5K transformed datasets/epoch for 200 epochs (in orange), followed with additional 100 epochs of training (in green). For the first 200 epochs we train with 90th percentile as the inlier/outlier threshold, which we reduce to 80th in the subsequent 100 epochs.

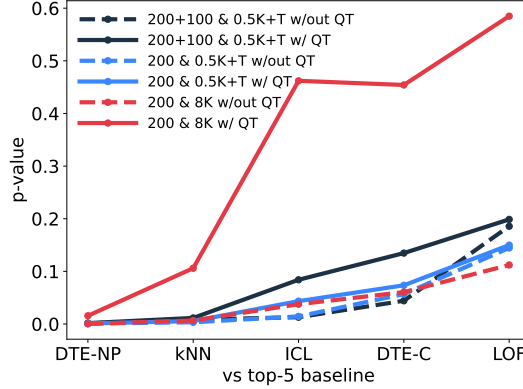


Figure 13:  $p$ -values increase, i.e. FoMo-0D performance improves, against top-5 baselines with quantile transform (QT) preceding inference, for 3 different settings of FoMo-0D w/  $D = 100$ .

Besides the ablation studies, we provide a qualitative case study of sample-to-sample attention in Appendix F, 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.

## F QUALITATIVE ANALYSIS ON SAMPLE-TO-SAMPLE ATTENTION

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 14, 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. 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

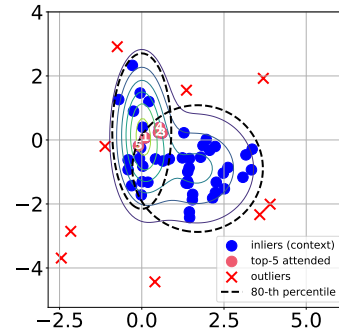


Figure 14: Top-5 attended inliers (all 50 inliers and only part of the outliers are shown for better visualization).



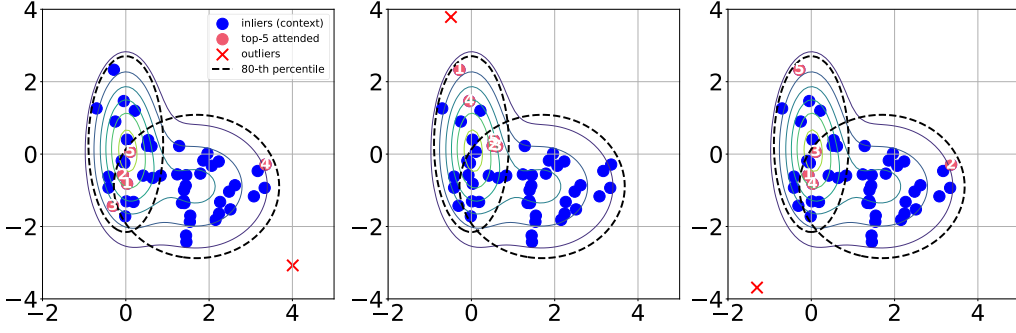


Figure 15: Top-5 attended inliers of 3 outliers at different position of the GMM

Wilcoxon signed rank test between `att` and `rdm` (alternative “less”), `att` and `prob` (alternative “greater”) over all the outliers, with a  $p$ -value of  $4.4 \times 10^{-4}$  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.

We visualize the top-5 attended inliers for 3 outliers at different position of the 2- $d$  GMM in Figure 15. For a specific outlier, there is a similar trend of attending to the center of a Gaussian (as shown in Figure 14), 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), suggesting FoMo-0D is incorporating both boundary and nearest neighbor information dynamically for each outlier.

## G FULL RESULTS

Tables 9.1& 9.2, 10.1 & 10.2, and 11.1 & 11.2 respectively show the AUROC, AUPR and F1 scores of the top-4 baselines, DTE-NP,  $k$ NN, ICL, and DTE-C as well as their corresponding <sup>avg</sup> model with the average performance across HPs, as listed in Table 4.

Tables 12.1&12.2, 13.1&13.2, and 14.1&14.2 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 <sup>avg</sup> variants of top-4 baselines + FoMo-0D w/  $D = 20$ ).

Table 8: 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-4<sup>5</sup>** baselines with performance `avg.`ed over varying HPs (denoted w/ <sup>avg</sup>); 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** ( $k$ NN, 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	$k$ NN	ICL	DTE-C	LOF	CBLOF	Feat.Bag.	SLAD	DDPM	OCSVM	DTE-NP <sup>avg</sup>	$k$ NN <sup>avg</sup>	ICL <sup>avg</sup>	DTE-C <sup>avg</sup>
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	-	<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
$d \leq 20$	-	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

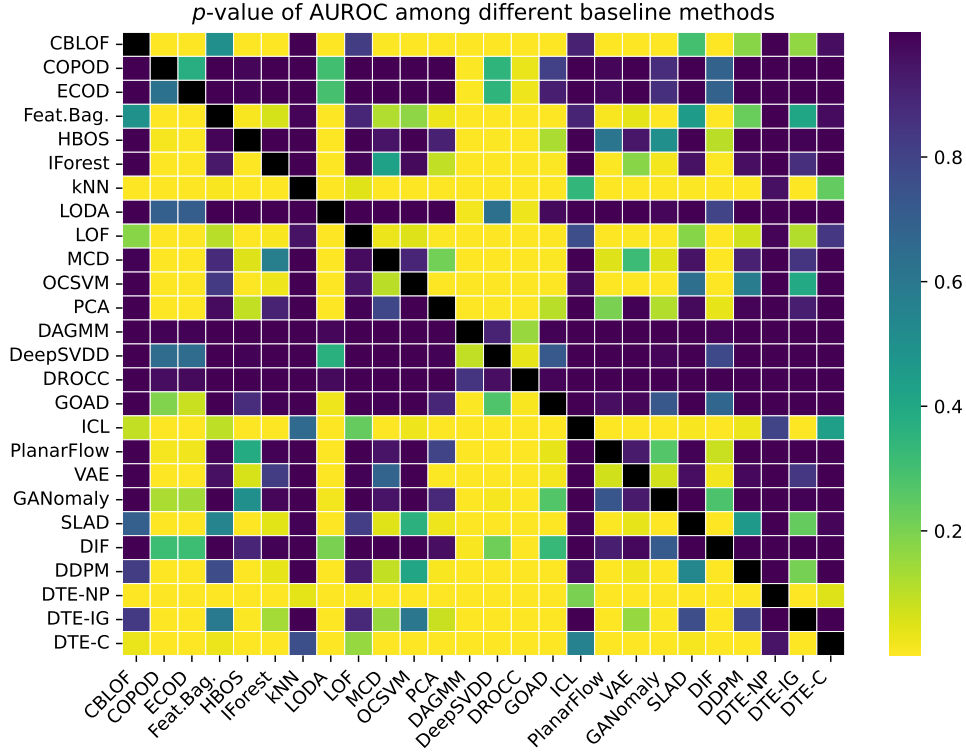


Figure 16: Pairwise *p*-values among baseline methods based on the Wilcoxon signed rank test w.r.t. AUROC performances across datasets.

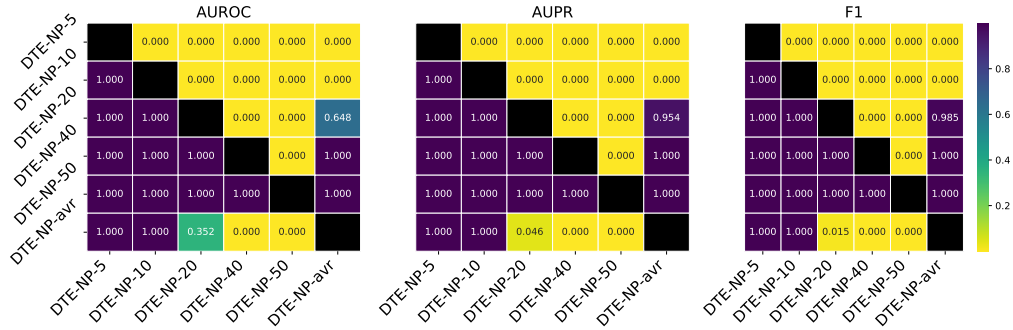


Figure 17: *p*-values w.r.t. AUROC/AUPR/F1 among different HP configurations of DTE-NP (i.e.,  $k \in \{5, 10, 20, 40, 50\}$ ), along with the  $^{avg}$  model with the average performance across HPs.

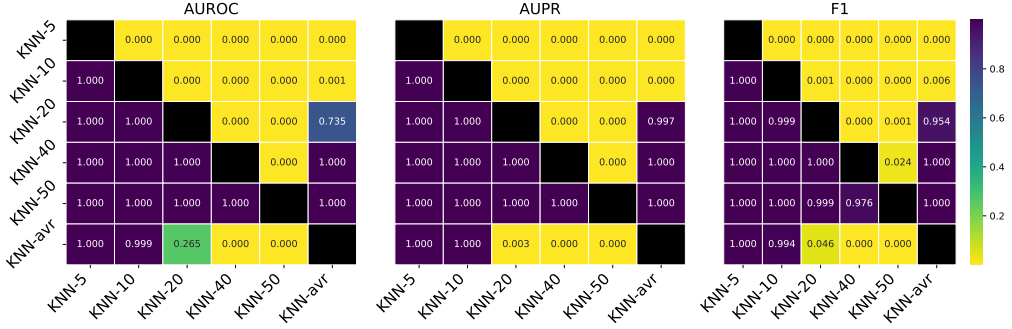


Figure 18:  $p$ -values w.r.t. AUROC/AUPR/F1 among different HP configurations of  $k$ NN (i.e.,  $k \in \{5, 10, 20, 40, 50\}$ ), along with the  $^{avg}$  model with the average performance across HPs.

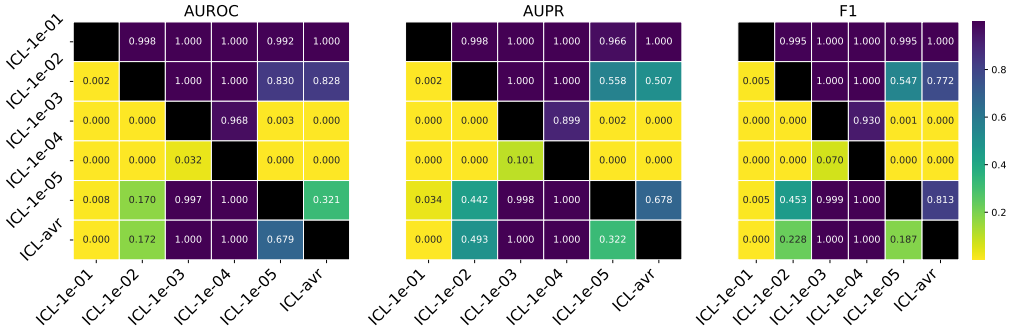


Figure 19:  $p$ -values w.r.t. AUROC/AUPR/F1 among different HP configurations of ICL (i.e.,  $\text{learning\_rate} \in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$ ), along with the  $^{avg}$  model with the average performance across HPs.

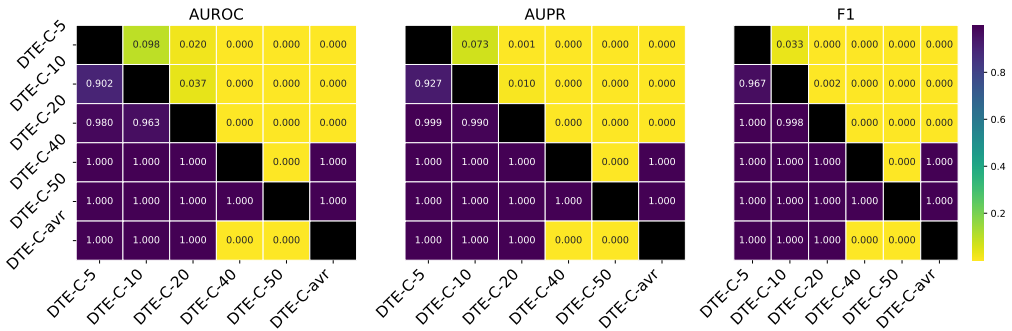


Figure 20:  $p$ -values w.r.t. AUROC/AUPR/F1 among different HP configurations of DTE-C (i.e.,  $k \in \{5, 10, 20, 40, 50\}$ ), along with the  $^{avg}$  model with the average performance across HPs.

Table 9.1: 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- <u>avg</u>	KNN-5	KNN-10	KNN-20	KNN-40	KNN-50	KNN- <u>avg</u>
aloi	50.69 $\pm$ 0.00	51.02 $\pm$ 0.00	51.26 $\pm$ 0.00	51.58 $\pm$ 0.00	<b>51.69<math>\pm</math>0.00</b>	51.25 $\pm$ 0.00	51.04 $\pm$ 0.00	51.33 $\pm$ 0.00	51.63 $\pm$ 0.00	51.97 $\pm$ 0.00	<b>52.08<math>\pm</math>0.00</b>	51.61 $\pm$ 0.00
amazon	<b>60.76<math>\pm</math>0.00</b>	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	<b>60.58<math>\pm</math>0.00</b>	60.52 $\pm$ 0.00	60.23 $\pm$ 0.00	60.02 $\pm$ 0.00	59.91 $\pm$ 0.00	60.25 $\pm$ 0.00
amthroid	<b>93.01<math>\pm</math>0.00</b>	92.89 $\pm$ 0.00	92.66 $\pm$ 0.00	92.38 $\pm$ 0.00	92.26 $\pm$ 0.00	92.64 $\pm$ 0.00	<b>92.81<math>\pm</math>0.00</b>	92.60 $\pm$ 0.00	92.34 $\pm$ 0.00	91.98 $\pm$ 0.00	91.79 $\pm$ 0.00	92.30 $\pm$ 0.00
backdoor	<b>94.48<math>\pm</math>0.42</b>	93.72 $\pm$ 0.46	92.67 $\pm$ 0.46	91.20 $\pm$ 0.45	90.68 $\pm$ 0.45	92.58 $\pm$ 0.44	<b>93.71<math>\pm</math>0.46</b>	92.58 $\pm$ 0.46	91.14 $\pm$ 0.47	89.18 $\pm$ 0.47	88.39 $\pm$ 0.51	91.00 $\pm$ 0.47
breast	<b>99.10<math>\pm</math>0.28</b>	98.91 $\pm$ 0.35	98.59 $\pm$ 0.34	98.40 $\pm$ 0.35	98.36 $\pm$ 0.28	98.67 $\pm$ 0.28	<b>99.09<math>\pm</math>0.24</b>	99.11 $\pm$ 0.27	99.16 $\pm$ 0.22	<b>99.21<math>\pm</math>0.18</b>	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	<b>78.93<math>\pm</math>0.00</b>	78.90 $\pm$ 0.00	78.76 $\pm$ 0.00	78.48 $\pm$ 0.00	78.74 $\pm$ 0.00	<b>78.84<math>\pm</math>0.00</b>	78.65 $\pm$ 0.00	<b>94.08<math>\pm</math>0.00</b>	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	<b>93.28<math>\pm</math>0.00</b>	92.47 $\pm$ 0.00	92.00 $\pm$ 0.00	92.44 $\pm$ 0.00	92.99 $\pm$ 0.00	93.85 $\pm$ 0.00	<b>94.08<math>\pm</math>0.00</b>	93.07 $\pm$ 0.00
cardiotocography	60.40 $\pm$ 0.00	61.63 $\pm$ 0.00	63.14 $\pm$ 0.00	65.05 $\pm$ 0.00	<b>65.81<math>\pm</math>0.00</b>	63.21 $\pm$ 0.00	62.11 $\pm$ 0.00	63.39 $\pm$ 0.00	65.12 $\pm$ 0.00	67.85 $\pm$ 0.00	<b>68.91<math>\pm</math>0.00</b>	65.48 $\pm$ 0.00
celeba	70.39 $\pm$ 0.33	72.58 $\pm$ 0.26	74.81 $\pm$ 0.34	76.87 $\pm$ 0.38	<b>77.47<math>\pm</math>0.37</b>	74.42 $\pm$ 0.28	72.23 $\pm$ 0.29	72.36 $\pm$ 0.12	77.50 $\pm$ 0.47	79.14 $\pm$ 0.38	<b>79.68<math>\pm</math>0.37</b>	76.90 $\pm$ 0.35
census	72.18 $\pm$ 0.17	72.34 $\pm$ 0.17	72.28 $\pm$ 0.10	71.93 $\pm$ 0.17	71.80 $\pm$ 0.17	72.11 $\pm$ 0.16	72.23 $\pm$ 0.29	<b>72.36<math>\pm</math>0.12</b>	71.94 $\pm$ 0.16	71.37 $\pm$ 0.21	71.28 $\pm$ 0.16	71.84 $\pm$ 0.15
cover	<b>97.90<math>\pm</math>0.17</b>	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	<b>97.51<math>\pm</math>0.15</b>	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	<b>99.72<math>\pm</math>0.03</b>	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	<b>99.51<math>\pm</math>0.06</b>	99.24 $\pm$ 0.08	98.20 $\pm$ 0.13	97.90 $\pm$ 0.14	<b>61.79<math>\pm</math>0.00</b>	98.74 $\pm$ 0.09
fault	58.34 $\pm$ 0.00	58.37 $\pm$ 0.00	58.70 $\pm$ 0.00	60.00 $\pm$ 0.00	<b>60.43<math>\pm</math>0.00</b>	59.17 $\pm$ 0.00	58.73 $\pm$ 0.00	58.76 $\pm$ 0.00	60.12 $\pm$ 0.00	61.71 $\pm$ 0.00	<b>61.79<math>\pm</math>0.00</b>	60.22 $\pm$ 0.00
fraud	<b>95.70<math>\pm</math>0.90</b>	95.67 $\pm$ 0.93	95.64 $\pm$ 0.93	95.60 $\pm$ 0.92	95.60 $\pm$ 0.92	95.64 $\pm$ 0.92	95.59 $\pm$ 0.97	95.55 $\pm$ 0.98	95.55 $\pm$ 0.89	95.54 $\pm$ 0.92	<b>95.62<math>\pm</math>0.88</b>	95.57 $\pm$ 0.93
glass	<b>96.08<math>\pm</math>0.39</b>	93.04 $\pm$ 1.06	89.82 $\pm$ 1.12	87.89 $\pm$ 1.10	87.31 $\pm$ 1.40	90.83 $\pm$ 0.91	<b>92.13<math>\pm</math>0.94</b>	88.67 $\pm$ 0.98	87.24 $\pm$ 1.18	84.93 $\pm$ 2.92	83.55 $\pm$ 2.61	87.30 $\pm$ 1.59
hepatitis	<b>99.84<math>\pm</math>0.20</b>	99.27 $\pm$ 0.51	96.89 $\pm$ 0.96	99.93 $\pm$ 0.01	99.91 $\pm$ 0.02	99.95 $\pm$ 0.01	<b>100.00<math>\pm</math>0.00</b>	99.99 $\pm$ 0.02	85.30 $\pm$ 2.34	85.46 $\pm$ 1.92	84.88 $\pm$ 2.09	99.96 $\pm$ 0.01
hipp	<b>99.99<math>\pm</math>0.00</b>	99.98 $\pm$ 0.01	99.95 $\pm$ 0.00	99.93 $\pm$ 0.01	99.91 $\pm$ 0.02	99.95 $\pm$ 0.01	<b>100.00<math>\pm</math>0.00</b>	99.99 $\pm$ 0.02	50.23 $\pm$ 0.00	50.23 $\pm$ 0.00	50.23 $\pm$ 0.00	50.18 $\pm$ 0.00
imdb	<b>50.48<math>\pm</math>0.00</b>	50.38 $\pm$ 0.00	50.32 $\pm$ 0.00	50.28 $\pm$ 0.00	50.27 $\pm$ 0.00	50.35 $\pm$ 0.00	<b>50.08<math>\pm</math>0.00</b>	50.04 $\pm$ 0.00	50.29 $\pm$ 0.00	50.29 $\pm$ 0.00	50.23 $\pm$ 0.00	50.18 $\pm$ 0.00
internats	<b>70.96<math>\pm</math>0.00</b>	68.65 $\pm$ 0.00	66.86 $\pm$ 0.00	65.97 $\pm$ 0.00	65.82 $\pm$ 0.00	67.65 $\pm$ 0.00	<b>68.08<math>\pm</math>0.00</b>	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
ionosphere	<b>98.48<math>\pm</math>0.60</b>	98.13 $\pm$ 0.74	97.84 $\pm$ 0.64	96.83 $\pm$ 0.71	96.21 $\pm$ 0.79	97.50 $\pm$ 0.63	<b>97.32<math>\pm</math>0.85</b>	97.62 $\pm$ 0.81	96.33 $\pm$ 0.76	92.80 $\pm$ 1.64	91.33 $\pm$ 1.65	95.12 $\pm$ 0.92
landsat	<b>66.59<math>\pm</math>0.00</b>	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	<b>68.23<math>\pm</math>0.00</b>	64.38 $\pm$ 0.00	64.36 $\pm$ 0.00	62.49 $\pm$ 0.00	61.23 $\pm$ 0.00	64.70 $\pm$ 0.00
letter	<b>56.12<math>\pm</math>0.00</b>	55.86 $\pm$ 0.00	54.78 $\pm$ 0.00	53.72 $\pm$ 0.00	53.40 $\pm$ 0.00	55.43 $\pm$ 0.00	<b>58.43<math>\pm</math>0.00</b>	54.38 $\pm$ 0.00	53.35 $\pm$ 0.00	52.18 $\pm$ 0.00	51.35 $\pm$ 0.00	53.53 $\pm$ 0.00
lymphography	<b>83.81<math>\pm</math>0.25</b>	83.49 $\pm$ 0.12	82.79 $\pm$ 0.32	82.05 $\pm$ 0.11	81.73 $\pm$ 0.10	83.81 $\pm$ 0.00	<b>83.27<math>\pm</math>0.10</b>	82.63 $\pm$ 0.00	81.85 $\pm$ 0.00	80.76 $\pm$ 0.00	80.30 $\pm$ 0.00	81.76 $\pm$ 0.00
magic_gamma	<b>87.65<math>\pm</math>0.00</b>	87.73 $\pm$ 0.00	87.68 $\pm$ 0.00	87.42 $\pm$ 0.00	87.20 $\pm$ 0.00	87.55 $\pm$ 0.00	<b>87.58<math>\pm</math>0.00</b>	87.54 $\pm$ 0.00	87.38 $\pm$ 0.00	86.97 $\pm$ 0.00	86.78 $\pm$ 0.00	87.29 $\pm$ 0.00
mnist	<b>93.93<math>\pm</math>0.00</b>	93.93 $\pm$ 0.00	93.57 $\pm$ 0.00	93.20 $\pm$ 0.00	93.08 $\pm$ 0.00	93.60 $\pm$ 0.00	<b>93.85<math>\pm</math>0.00</b>	93.45 $\pm$ 0.00	93.00 $\pm$ 0.00	92.55 $\pm$ 0.00	92.36 $\pm$ 0.00	93.04 $\pm$ 0.00
mnist_aug	<b>100.00<math>\pm</math>0.00</b>	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	<b>100.00<math>\pm</math>0.00</b>	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
optdigits	<b>95.00<math>\pm</math>0.00</b>	93.97 $\pm$ 0.00	92.52 $\pm$ 0.00	90.87 $\pm$ 0.00	90.28 $\pm$ 0.00	92.53 $\pm$ 0.00	<b>93.72<math>\pm</math>0.00</b>	92.09 $\pm$ 0.00	90.20 $\pm$ 0.00	87.66 $\pm$ 0.00	86.67 $\pm$ 0.00	90.07 $\pm$ 0.00
pageblocks	<b>89.04<math>\pm</math>0.00</b>	89.40 $\pm$ 0.00	<b>89.56<math>\pm</math>0.00</b>	89.37 $\pm$ 0.00	89.23 $\pm$ 0.00	89.32 $\pm$ 0.00	<b>89.65<math>\pm</math>0.00</b>	89.86 $\pm$ 0.00	<b>89.88<math>\pm</math>0.00</b>	89.30 $\pm$ 0.00	89.18 $\pm$ 0.00	89.57 $\pm$ 0.00
pendigits	<b>99.90<math>\pm</math>0.00</b>	99.88 $\pm$ 0.00	99.83 $\pm$ 0.00	99.51 $\pm$ 0.00	99.38 $\pm$ 0.00	99.70 $\pm$ 0.00	<b>99.87<math>\pm</math>0.00</b>	99.79 $\pm$ 0.00	99.21 $\pm$ 0.00	98.58 $\pm$ 0.00	98.39 $\pm$ 0.00	99.17 $\pm$ 0.00
pinna	<b>82.21<math>\pm</math>1.82</b>	82.09 $\pm$ 1.61	77.98 $\pm$ 1.38	77.28 $\pm$ 1.35	77.14 $\pm$ 1.33	78.87 $\pm$ 1.43	<b>77.44<math>\pm</math>2.07</b>	76.14 $\pm$ 1.36	76.39 $\pm$ 1.21	76.58 $\pm$ 1.25	76.38 $\pm$ 1.29	76.58 $\pm$ 1.39
satellite	<b>82.40<math>\pm</math>0.00</b>	82.09 $\pm$ 0.00	81.58 $\pm$ 0.00	80.71 $\pm$ 0.00	80.38 $\pm$ 0.00	81.43 $\pm$ 0.00	<b>82.24<math>\pm</math>0.00</b>	81.56 $\pm$ 0.00	80.56 $\pm$ 0.00	79.22 $\pm$ 0.00	78.76 $\pm$ 0.00	80.47 $\pm$ 0.00
satimage-2	<b>99.68<math>\pm</math>0.00</b>	99.73 $\pm$ 0.00	99.79 $\pm$ 0.00	<b>99.82<math>\pm</math>0.00</b>	99.81 $\pm$ 0.00	99.77 $\pm$ 0.00	<b>99.71<math>\pm</math>0.00</b>	<b>99.80<math>\pm</math>0.00</b>	99.79 $\pm$ 0.00	99.78 $\pm$ 0.00	99.77 $\pm$ 0.00	99.77 $\pm$ 0.00
shuttle	<b>99.94<math>\pm</math>0.00</b>	99.92 $\pm$ 0.00	99.91 $\pm$ 0.00	99.91 $\pm$ 0.00	99.91 $\pm$ 0.00	99.92 $\pm$ 0.00	<b>99.91<math>\pm</math>0.00</b>	99.90 $\pm$ 0.00	99.89 $\pm$ 0.00	99.89 $\pm$ 0.00	99.89 $\pm$ 0.00	99.89 $\pm$ 0.00
skin	<b>99.66<math>\pm</math>0.05</b>	99.52 $\pm$ 0.04	99.28 $\pm$ 0.02	98.77 $\pm$ 0.09	98.54 $\pm$ 0.09	99.15 $\pm$ 0.05	<b>99.51<math>\pm</math>0.06</b>	99.17 $\pm$ 0.05	99.89 $\pm$ 0.00	97.43 $\pm$ 0.05	96.90 $\pm$ 0.09	98.34 $\pm$ 0.04
snmp	92.94 $\pm$ 2.55	92.84 $\pm$ 2.56	<b>93.14<math>\pm</math>2.20</b>	93.14 $\pm$ 2.15	93.14 $\pm$ 2.20	93.04 $\pm$ 2.32	<b>92.90<math>\pm</math>2.48</b>	92.97 $\pm$ 2.52	<b>93.25<math>\pm</math>2.06</b>	93.11 $\pm$ 2.28	93.24 $\pm$ 2.25	93.10 $\pm$ 2.30
speech	<b>41.12<math>\pm</math>0.00</b>	83.49 $\pm$ 0.00	82.97 $\pm$ 0.00	82.50 $\pm$ 0.00	82.36 $\pm$ 0.00	83.07 $\pm$ 0.00	<b>83.36<math>\pm</math>0.00</b>	82.68 $\pm$ 0.00	82.22 $\pm$ 0.00	81.83 $\pm$ 0.00	81.72 $\pm$ 0.00	82.36 $\pm$ 0.00
stamps	<b>97.88<math>\pm</math>0.33</b>	97.04 $\pm$ 0.33	95.60 $\pm$ 0.78	94.59 $\pm$ 1.07	94.29 $\pm$ 1.08	95.88 $\pm$ 0.66	<b>96.19<math>\pm</math>1.24</b>	94.35 $\pm$ 1.35	93.67 $\pm$ 1.09	93.44 $\pm$ 1.21	93.33 $\pm$ 1.25	94.20 $\pm$ 1.17
thyroid	<b>98.58<math>\pm</math>0.00</b>	98.63 $\pm$ 0.00	<b>98.64<math>\pm</math>0.00</b>	98.63 $\pm$ 0.00	98.59 $\pm$ 0.00	98.61 $\pm$ 0.00	<b>98.68<math>\pm</math>0.00</b>	98.67 $\pm$ 0.00	98.69 $\pm$ 0.00	<b>98.70<math>\pm</math>0.00</b>	98.69 $\pm$ 0.00	98.68 $\pm$ 0.00
vertebral	<b>79.59<math>\pm</math>2.23</b>	67.05 $\pm$ 2.09	56.99 $\pm$ 2.07	49.07 $\pm$ 1.60	47.08 $\pm$ 1.43	59.96 $\pm$ 1.83	<b>57.33<math>\pm</math>3.80</b>	49.40 $\pm$ 1.87	45.06 $\pm$ 1.20	41.08 $\pm$ 1.52	40.08 $\pm$ 1.23	46.59 $\pm$ 1.67
vowels	<b>82.11<math>\pm</math>0.00</b>	81.62 $\pm$ 0.00	80.46 $\pm$ 0.00	78.11 $\pm$ 0.00	76.87 $\pm$ 0.00	79.83 $\pm$ 0.00	<b>82.21<math>\pm</math>0.00</b>	80.20 $\pm$ 0.00	77.82 $\pm$ 0.00	72.26 $\pm$ 0.00	69.03 $\pm$ 0.00	76.30 $\pm$ 0.00
waveform	<b>74.42<math>\pm</math>0.00</b>	75.40 $\pm$ 0.00	76.12 $\pm$ 0.00	76.79 $\pm$ 0.00	<b>76.90<math>\pm</math>0.00</b>	75.93 $\pm$ 0.00	<b>82.21<math>\pm</math>0.00</b>	76.04 $\pm$ 0.00	76.78 $\pm$ 0.00	77.47 $\pm$ 0.00	<b>77.60<math>\pm</math>0.00</b>	76.62 $\pm$ 0.00
wbc	<b>99.62<math>\pm</math>0.27</b>	99.36 $\pm$ 0.33	99.17 $\pm$ 0.41	99.12 $\pm$ 0.37	99.13 $\pm$ 0.20	99.27 $\pm$ 0.34	<b>98.88<math>\pm</math>0.50</b>	98.83 $\pm$ 0.25	98.86 $\pm$ 0.38	99.08 $\pm$ 0.34	<b>99.10<math>\pm</math>0.34</b>	98.95 $\pm$ 0.35
wdbc	<b>66.98<math>\pm</math>0.00</b>	63.87 $\pm$ 0.29	60.18 $\pm$ 0.00	59.18 $\pm$ 0.22	55.02 $\pm$ 0.30	60.46 $\pm$ 0.00	<b>63.66<math>\pm</math>0.00</b>	59.09 $\pm$ 0.29	59.08 $\pm$ 0.22	99.08 $\pm$ 0.22	99.08 $\pm$ 0.17	99.11 $\pm$ 0.22
wilt	<b>99.90<math>\pm</math>0.14</b>	99.73 $\pm$ 0.09	99.33 $\pm$ 0.18	98.88 $\pm$ 0.19	98.57 $\pm$ 0.30	99.23 $\pm$ 0.27	<b>99.23<math>\pm</math>0.27</b>	99.09 $\pm$ 0.29	99.08 $\pm$ 0.21	99.08 $\pm$ 0.22	99.08 $\pm$	

Table 9.2: 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  $\text{avg}$  model. We use **bold** and underline respectively to mark the **best** and the worst performance of each model to showcase the variability across different HP settings.

dataset	ICL-0.1	ICL-0.01	ICL-0.001	ICL-0.0001	ICL-1e-05	ICL-avr	DTE-C-5	DTE-C-10	DTE-C-20	DTE-C-40	DTE-C-50	DTE-C-avr
aloi	47.74±0.28	47.12±0.51	46.81±0.47	<b>48.42</b> ±0.24	48.06±0.23	47.63±0.15	50.20±0.21	<b>50.84</b> ±0.10	50.16±0.34	50.26±0.33	50.00±0.00	50.29±0.09
amazon	53.07±0.19	<b>53.44</b> ±0.19	52.73±0.58	53.31±0.21	53.18±0.15	53.17±0.19	<b>56.20</b> ±1.62	55.07±4.20	56.12±2.73	50.00±0.00	50.00±0.00	53.48±1.33
amthroid	84.02±9.46	72.68±3.79	87.29±2.69	<b>88.84</b> ±2.35	88.52±1.40	84.27±2.31	97.47±0.10	97.65±0.11	<b>97.73</b> ±0.15	97.40±0.22	50.00±0.00	88.05±0.05
backdoor	93.03±0.66	92.91±0.70	93.30±0.44	<b>93.93</b> ±0.58	93.32±0.71	93.30±0.48	88.65±1.08	92.06±1.04	<b>92.64</b> ±0.84	50.00±0.00	<b>99.26</b> ±0.07	74.67±0.36
breastui	98.96±0.47	98.97±0.20	<b>99.19</b> ±0.34	99.11±0.18	97.61±0.65	93.44±1.25	93.45±1.34	94.34±1.25	<b>94.61</b> ±0.91	98.85±0.45	50.00±0.00	96.52±0.53
campaign	76.07±0.87	74.61±1.97	78.88±0.42	79.79±0.91	<b>82.30</b> ±0.38	78.33±0.52	<b>79.18</b> ±1.09	78.24±2.09	78.49±1.13	50.00±0.00	50.00±0.00	67.18±0.74
cardio	73.99±7.79	<b>84.98</b> ±4.75	82.30±3.36	78.68±2.18	68.59±0.64	77.71±1.76	<b>88.07</b> ±0.51	87.54±0.69	87.66±0.63	50.00±0.00	50.00±0.00	72.66±0.21
cardiotocography	48.99±2.02	<b>54.18</b> ±2.77	53.24±3.25	50.67±2.55	47.20±1.33	50.85±1.24	60.09±2.19	<b>60.36</b> ±0.51	59.05±1.34	50.00±0.00	50.00±0.00	55.90±0.23
ceiba	79.38±2.17	79.15±2.47	76.13±1.88	78.39±1.78	<b>79.43</b> ±1.44	78.49±1.03	<b>82.95</b> ±1.26	81.59±0.79	80.16±1.26	50.00±0.00	50.00±0.00	68.94±0.44
census	70.41±2.07	67.52±8.79	73.93±0.78	<b>75.03</b> ±0.46	74.37±0.53	72.25±1.63	<b>70.95</b> ±0.91	68.04±0.85	68.05±2.58	50.00±0.00	50.00±0.00	61.41±0.71
cover	95.59±3.09	82.32±9.14	91.92±4.58	<b>94.97</b> ±2.91	94.27±3.40	91.42±1.74	97.57±0.86	<b>97.81</b> ±0.75	96.61±0.63	50.00±0.00	50.00±0.00	78.40±0.14
donors	86.20±10.29	97.76±1.57	98.64±0.55	99.37±0.30	<b>99.51</b> ±0.14	96.30±1.95	<b>98.68</b> ±0.15	97.56±0.53	95.64±0.27	58.94±17.89	50.00±0.00	80.17±3.55
fault	<b>63.37</b> ±3.29	63.04±1.93	61.65±0.61	61.44±0.87	62.83±1.01	62.47±1.14	<b>59.16</b> ±1.69	58.41±1.41	59.60±1.53	50.21±0.42	55.41±0.63	55.41±0.63
fraud	83.24±3.45	93.21±1.25	94.97±0.71	<b>95.84</b> ±0.87	93.90±1.05	94.23±0.87	<b>93.46</b> ±0.91	89.96±2.20	93.50±2.01	66.42±8.60	50.00±0.00	76.21±1.08
glass	84.02±8.38	94.66±3.83	99.25±0.37	<b>99.30</b> ±0.35	99.12±0.55	95.27±1.45	<b>99.54</b> ±0.38	98.90±1.01	98.90±1.08	84.42±1.20	50.00±0.00	76.95±2.58
hepatitis	99.95±0.11	98.76±1.84	99.25±0.37	<b>99.30</b> ±0.35	99.12±0.55	95.27±1.45	<b>99.54</b> ±0.38	98.90±1.01	98.90±1.08	84.42±1.20	50.00±0.00	76.95±2.58
http	99.96±0.07	99.97±0.02	99.98±0.01	<b>100.00</b> ±0.00	99.86±0.29	99.69±0.48	99.54±0.08	99.64±0.28	99.39±0.07	50.00±0.00	50.00±0.00	89.59±0.09
imdb	52.16±0.31	51.88±0.40	52.28±0.15	53.35±0.12	<b>53.06</b> ±0.15	47.68±3.79	99.02±1.49	<b>99.89</b> ±1.90	47.81±2.08	50.00±0.00	50.00±0.00	66.79±0.73
internats	72.61±1.19	73.97±0.44	74.09±1.75	73.74±0.45	68.31±0.74	77.71±1.52	94.67±1.52	<b>95.23</b> ±1.37	85.72±2.48	50.00±0.00	50.00±0.00	83.01±5.81
ionosphere	96.81±2.22	96.13±3.45	98.90±0.41	<b>98.29</b> ±0.31	98.14±0.79	97.78±1.23	94.67±1.52	<b>95.23</b> ±1.37	85.72±2.48	50.00±0.00	50.00±0.00	83.01±5.81
landsat	63.11±2.65	60.92±1.79	63.93±0.76	<b>67.85</b> ±0.68	61.82±1.06	64.78±1.03	<b>51.80</b> ±0.94	50.72±1.62	85.72±2.48	50.00±0.00	50.00±0.00	83.01±5.81
letter	99.95±0.05	99.95±0.05	99.95±0.05	<b>100.00</b> ±0.00	99.95±0.05	99.95±0.05	99.95±0.05	99.95±0.05	99.95±0.05	99.95±0.05	99.95±0.05	99.95±0.05
lymphography	<b>100.00</b> ±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
magic gamma	76.94±3.25	76.94±3.25	77.61±0.70	77.81±0.88	<b>78.36</b> ±1.34	76.09±1.17	86.12±0.31	<b>87.42</b> ±0.35	88.70±0.73	64.88±18.23	50.00±0.00	75.19±3.50
magic gamma manography	79.90±5.13	<b>85.08</b> ±3.65	79.13±1.20	80.75±1.40	77.35±0.77	80.47±0.93	<b>90.21</b> ±0.70	<b>86.02</b> ±1.23	84.85±2.12	85.40±1.03	50.00±0.00	72.55±0.64
manography	75.53±1.43	76.05±3.81	79.13±1.20	80.75±1.40	77.35±0.77	80.47±0.93	<b>90.21</b> ±0.70	<b>86.02</b> ±1.23	84.85±2.12	85.40±1.03	50.00±0.00	72.55±0.64
music	<b>100.00</b> ±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
optdigits	91.22±1.88	94.15±0.97	95.17±0.85	<b>98.47</b> ±0.03	95.92±1.84	97.18±0.37	<b>100.00</b> ±0.00	100.00±0.00	100.00±0.00	50.00±0.00	50.00±0.00	80.00±0.00
pageblocks	89.21±1.10	<b>90.77</b> ±1.73	88.70±1.01	88.58±0.56	88.53±0.27	94.82±0.47	89.96±0.39	<b>90.29</b> ±0.81	89.23±0.26	88.80±0.13	50.00±0.00	81.66±0.20
pendigits	87.01±9.78	96.45±2.40	96.17±2.68	<b>98.28</b> ±0.77	95.10±0.89	94.60±2.36	<b>98.24</b> ±0.47	97.42±0.83	96.89±0.68	50.00±0.00	50.00±0.00	81.66±0.20
pinna	74.66±3.27	73.04±1.77	<b>79.77</b> ±1.92	78.06±0.74	75.09±2.65	76.12±1.04	70.12±1.89	67.22±3.15	68.78±2.78	<b>71.26</b> ±2.38	69.39±5.15	68.95±0.96
satellite	74.62±9.67	83.23±2.89	85.43±0.52	<b>89.24</b> ±0.57	85.35±0.48	83.57±2.48	<b>80.21</b> ±1.18	79.01±0.76	78.54±0.62	56.11±12.22	50.00±0.00	68.77±2.26
satimage-2	94.41±2.46	93.04±12.69	99.99±0.00	99.99±0.00	99.99±0.00	99.99±0.00	<b>99.76</b> ±0.01	99.74±0.01	99.74±0.01	99.28±0.22	50.00±0.00	89.70±0.05
shuttle	99.43±0.50	99.89±0.05	<b>99.90</b> ±0.00	99.90±0.00	99.97±0.01	99.85±0.11	91.69±0.28	<b>92.12</b> ±0.27	91.95±0.22	91.63±0.28	50.00±0.00	83.48±0.19
skin	53.71±32.40	86.42±5.89	<b>86.94</b> ±5.22	89.48±4.52	<b>92.63</b> ±4.06	88.37±5.41	95.06±1.20	95.61±1.24	<b>95.84</b> ±1.23	95.84±1.22	50.00±0.00	86.47±0.96
smtp	81.23±10.57	87.97±4.57	83.16±0.34	<b>83.39</b> ±0.40	78.42±0.49	80.71±0.86	82.98±0.37	83.29±0.49	<b>83.74</b> ±0.38	50.00±0.00	50.00±0.00	70.00±0.04
speech	50.06±3.03	50.96±2.48	<b>53.72</b> ±1.46	51.29±2.33	46.64±1.97	50.53±1.10	38.02±1.99	38.55±1.07	38.72±1.79	<b>50.00</b> ±0.00	50.00±0.00	43.06±0.55
stamps	77.62±9.15	84.33±6.20	<b>97.29</b> ±0.48	96.70±0.85	94.64±3.19	90.11±2.02	<b>93.01</b> ±1.38	91.27±2.53	86.48±2.57	90.17±4.76	50.10±26.50	82.21±5.22
thyroid	94.79±1.60	95.24±1.04	96.20±0.74	<b>96.78</b> ±1.19	93.73±0.80	95.35±0.60	98.75±0.07	98.92±0.02	98.92±0.04	<b>98.94</b> ±0.05	50.00±0.00	89.11±0.02
vertebral	53.80±4.97	58.76±7.63	75.60±6.18	<b>82.96</b> ±2.86	81.92±1.33	70.61±1.11	87.25±1.51	86.66±0.72	<b>87.49</b> ±0.93	50.00±0.00	50.00±0.00	59.97±2.61
vowels	73.59±6.94	79.11±2.66	84.59±3.49	<b>84.81</b> ±3.61	83.99±1.93	81.22±1.57	87.25±1.51	86.66±0.72	<b>87.49</b> ±0.93	50.00±0.00	50.00±0.00	72.28±0.42
waveform	70.94±1.78	<b>73.85</b> ±6.07	70.23±2.82	61.96±1.35	64.69±1.66	68.33±2.15	65.16±1.34	63.97±2.75	<b>66.24</b> ±1.45	50.00±0.00	50.00±0.00	59.07±0.68
wbc	96.66±2.87	99.08±0.62	99.89±0.19	<b>99.90</b> ±0.15	99.29±0.25	99.34±0.42	86.07±4.41	83.96±4.40	<b>86.05</b> ±4.40	50.00±0.00	50.00±0.00	81.10±1.77
wdbc	52.75±14.46	61.09±5.94	83.81±3.04	<b>84.63</b> ±1.21	79.81±0.93	72.42±2.57	98.76±0.46	<b>97.06</b> ±0.69	91.09±1.72	84.90±1.24	50.00±0.00	77.44±0.57
wilt	99.64±0.73	99.73±0.29	<b>99.92</b> ±0.12	99.81±0.39	99.81±0.39	99.79±0.33	99.91±0.11	99.91±0.11	99.91±0.11	99.91±0.11	99.91±0.11	99.91±0.11
wine	91.76±2.53	76.46±9.54	95.40±1.42	<b>95.51</b> ±1.23	92.46±1.77	46.32±1.91	68.59±1.26	68.59±1.26	68.59±1.26	51.85±2.56	50.00±0.00	61.89±0.95
wpbc	45.12±3.69	47.02±1.77	47.25±0.62	<b>47.38</b> ±0.17	46.32±1.91	46.32±1.91	68.59±1.26	68.59±1.26	68.59±1.26	51.85±2.56	50.00±0.00	61.89±0.95
year	54.26±0.42	<b>54.26</b> ±0.51	54.08±0.78	54.07±0.26	52.73±0.21	53.71±0.22	82.00±1.65	83.62±0.43	85.25±0.26	30.03±0.00	50.00±0.00	71.28±0.16
yeap	80.78±0.40	80.78±0.40	<b>85.02</b> ±0.24	84.91±0.11	83.73±0.12	83.46±0.11	90.13±0.19	<b>90.30</b> ±0.13	90.30±0.13	50.00±0.00	50.00±0.00	71.28±0.16
MINST-C	72.40±0.29	90.44±0.09	90.91±0.06	<b>90.95</b> ±0.03	88.57±0.10	89.79±0.10	68.78±0.49	<b>69.56</b> ±0.38	69.56±0.38	50.00±0.00	50.00±0.00	61.25±0.12
FashionMNIST	97.80±0.21	94.12±0.37	95.72±0.09	<b>95.89</b> ±0.27	92.25±0.13	92.25±0.13	63.02±0.30	<b>63.02</b> ±0.30	63.02±0.30	50.00±0.00	50.00±0.00	57.80±0.05
CIFAR10	59.44±0.27	64.12±0.37	65.72±0.09	<b>66.16</b> ±0.08	60.61±0.08	60.61±0.08	69.01±0.03	<b>69.01</b> ±0.03	69.01±0.03	50.00±0.00	50.00±0.00	73.32±0.42
SVHN	93.31±0.36	91.73±0.10	91.73±0.10	<b>94.26</b> ±0.32	92.82±0.44	92.82±0.44	86.78±0.49	<b>86.78</b> ±0.49	86.78±0.49	50.00±0.00	50.00±0.00	58.00±0.60
MTEC-AD	57.95±0.57	59.21±0.72	59.52±0.72	<b>60.25</b> ±0.48	55.58±0.17	55.58±0.17	68.14±1.95	<b>68.14</b> ±1.95	68.14±1.95	50.00±0.00	50.00±0.00	60.06±0.42
20news	56.84±0.12	57.43±0.39	57.59±0.12	<b>57.89</b> ±0.09	56.85±0.16	56.85±0.16	64.75±0.11	<b>64.75</b> ±0.11	64.75±0.11	50.00±0.00	50.00±0.00	60.06±0.42

Table 10.1: 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- <u>avg</u>	KNN-5	KNN-10	KNN-20	KNN-40	KNN-50	KNN- <u>avg</u>
aloi	5.95 $\pm$ 0.00	5.99 $\pm$ 0.00	6.02 $\pm$ 0.00	6.06 $\pm$ 0.00	<b>6.07<math>\pm</math>0.00</b>	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	11.69 $\pm$ 0.00	11.69 $\pm$ 0.00	11.70 $\pm$ 0.00	11.70 $\pm$ 0.00	11.60 $\pm$ 0.00	11.59 $\pm$ 0.00	11.65 $\pm$ 0.00
amthroid	<b>67.49<math>\pm</math>0.00</b>	66.73 $\pm$ 0.00	66.04 $\pm$ 0.00	65.39 $\pm$ 0.00	64.87 $\pm$ 0.00	66.11 $\pm$ 0.00	<b>68.07<math>\pm</math>0.00</b>	67.90 $\pm$ 0.00	67.26 $\pm$ 0.00	66.27 $\pm$ 0.00	65.79 $\pm$ 0.00	67.06 $\pm$ 0.00
backdoor	<b>55.90<math>\pm</math>0.99</b>	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	<b>46.70<math>\pm</math>1.22</b>	37.36 $\pm$ 1.35	29.58 $\pm$ 0.41	24.34 $\pm$ 0.41	22.34 $\pm$ 0.53	32.06 $\pm$ 0.76
breast	<b>98.51<math>\pm</math>0.56</b>	98.19 $\pm$ 0.58	97.56 $\pm$ 0.51	97.13 $\pm$ 0.62	97.05 $\pm$ 0.45	97.69 $\pm$ 0.40	98.97 $\pm$ 0.28	99.01 $\pm$ 0.31	99.08 $\pm$ 0.23	99.15 $\pm$ 0.17	<b>99.16<math>\pm</math>0.16</b>	99.08 $\pm$ 0.22
campaign	48.48 $\pm$ 0.00	49.05 $\pm$ 0.00	<b>49.77<math>\pm</math>0.00</b>	49.77 $\pm$ 0.00	49.51 $\pm$ 0.00	49.31 $\pm$ 0.00	49.04 $\pm$ 0.00	49.89 $\pm$ 0.00	<b>50.45<math>\pm</math>0.00</b>	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	<b>79.53<math>\pm</math>0.00</b>	77.22 $\pm$ 0.00	78.33 $\pm$ 0.00	78.33 $\pm$ 0.00	79.14 $\pm$ 0.00	80.67 $\pm$ 0.00	<b>81.15<math>\pm</math>0.00</b>	79.30 $\pm$ 0.00
cardiotocography	56.55 $\pm$ 0.00	57.18 $\pm$ 0.00	58.19 $\pm$ 0.00	59.42 $\pm$ 0.00	<b>59.95<math>\pm</math>0.00</b>	58.26 $\pm$ 0.00	57.43 $\pm$ 0.00	58.37 $\pm$ 0.00	59.44 $\pm$ 0.00	61.41 $\pm$ 0.00	<b>62.19<math>\pm</math>0.00</b>	59.77 $\pm$ 0.00
celeba	10.56 $\pm$ 0.44	11.63 $\pm$ 0.49	12.74 $\pm$ 0.52	13.92 $\pm$ 0.58	<b>14.30<math>\pm</math>0.59</b>	12.63 $\pm$ 0.51	11.99 $\pm$ 0.57	13.26 $\pm$ 0.61	14.50 $\pm$ 0.58	15.70 $\pm$ 0.65	<b>16.10<math>\pm</math>0.68</b>	14.31 $\pm$ 0.60
census	21.14 $\pm$ 0.39	<b>21.38<math>\pm</math>0.54</b>	21.16 $\pm$ 0.43	20.67 $\pm$ 0.41	20.52 $\pm$ 0.42	20.97 $\pm$ 0.43	<b>21.36<math>\pm</math>0.76</b>	21.22 $\pm$ 0.39	20.59 $\pm$ 0.33	20.00 $\pm$ 0.42	19.94 $\pm$ 0.44	20.62 $\pm$ 0.44
cover	<b>63.67<math>\pm</math>3.21</b>	57.85 $\pm$ 3.52	51.55 $\pm$ 3.10	44.58 $\pm$ 2.49	42.11 $\pm$ 2.29	51.95 $\pm$ 2.90	<b>55.15<math>\pm</math>3.45</b>	48.67 $\pm$ 2.84	41.44 $\pm$ 2.32	33.72 $\pm$ 1.51	31.69 $\pm$ 1.35	42.14 $\pm$ 2.20
donors	<b>93.23<math>\pm</math>0.80</b>	91.25 $\pm$ 0.72	88.17 $\pm$ 0.95	83.92 $\pm$ 1.29	82.34 $\pm$ 1.32	87.78 $\pm$ 0.99	<b>89.44<math>\pm</math>0.96</b>	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	61.29 $\pm$ 0.00	61.98 $\pm$ 0.00	<b>62.31<math>\pm</math>0.00</b>	61.84 $\pm$ 0.00	61.98 $\pm$ 0.00	61.16 $\pm$ 0.00	61.92 $\pm$ 0.00	63.67 $\pm$ 0.00	<b>64.06<math>\pm</math>0.00</b>	62.56 $\pm$ 0.00
fraud	40.60 $\pm$ 6.67	<b>43.77<math>\pm</math>5.38</b>	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	<b>44.96<math>\pm</math>3.90</b>	41.19 $\pm$ 3.73	37.33 $\pm$ 4.15	36.42 $\pm$ 4.12	40.45 $\pm$ 4.07
glass	<b>60.15<math>\pm</math>6.89</b>	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	<b>44.08<math>\pm</math>6.38</b>	32.96 $\pm$ 3.74	29.87 $\pm$ 3.75	26.62 $\pm$ 2.65	26.04 $\pm$ 3.61	31.91 $\pm$ 3.73
hepatitis	<b>99.47<math>\pm</math>0.73</b>	97.98 $\pm$ 1.43	91.71 $\pm$ 2.22	81.65 $\pm$ 3.68	78.95 $\pm$ 3.45	89.95 $\pm$ 1.80	<b>91.10<math>\pm</math>4.41</b>	69.28 $\pm$ 3.16	64.12 $\pm$ 5.59	64.29 $\pm$ 5.08	64.33 $\pm$ 6.16	70.62 $\pm$ 4.48
http	<b>98.52<math>\pm</math>0.37</b>	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	<b>100.00<math>\pm</math>0.00</b>	98.01 $\pm$ 3.98	91.24 $\pm$ 1.42	91.44 $\pm$ 1.25	91.28 $\pm$ 1.36	94.39 $\pm$ 1.31
imdb	9.11 $\pm$ 0.00	9.09 $\pm$ 0.00	9.06 $\pm$ 0.00	9.07 $\pm$ 0.00	9.06 $\pm$ 0.00	9.08 $\pm$ 0.00	9.25 $\pm$ 0.00	9.94 $\pm$ 0.00	<b>8.99<math>\pm</math>0.00</b>	8.98 $\pm$ 0.00	8.99 $\pm$ 0.00	8.96 $\pm$ 0.00
internets	49.76 $\pm$ 0.00	48.19 $\pm$ 0.00	47.56 $\pm$ 0.00	47.45 $\pm$ 0.00	47.45 $\pm$ 0.00	49.03 $\pm$ 0.00	<b>49.22<math>\pm</math>0.00</b>	47.29 $\pm$ 0.00	46.93 $\pm$ 0.00	46.93 $\pm$ 0.00	46.94 $\pm$ 0.00	47.47 $\pm$ 0.00
ionosphere	<b>98.72<math>\pm</math>0.48</b>	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	<b>97.86<math>\pm</math>0.60</b>	96.11 $\pm$ 0.52	94.12 $\pm$ 1.45	92.90 $\pm$ 1.65	92.90 $\pm$ 1.65	96.01 $\pm$ 0.78
landsat	<b>56.14<math>\pm</math>0.00</b>	54.25 $\pm$ 0.00	50.75 $\pm$ 0.00	46.43 $\pm$ 0.00	45.17 $\pm$ 0.00	50.55 $\pm$ 0.00	<b>54.85<math>\pm</math>0.00</b>	50.62 $\pm$ 0.00	45.18 $\pm$ 0.00	41.32 $\pm$ 0.00	40.50 $\pm$ 0.00	46.49 $\pm$ 0.00
letter	<b>8.96<math>\pm</math>0.00</b>	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	<b>8.70<math>\pm</math>0.00</b>	8.83 $\pm$ 0.00	8.41 $\pm$ 0.00	8.22 $\pm$ 0.00	8.22 $\pm$ 0.00	8.44 $\pm$ 0.00
lymphography	<b>86.20<math>\pm</math>0.00</b>	85.86 $\pm$ 0.00	85.26 $\pm$ 0.00	84.96 $\pm$ 0.00	84.96 $\pm$ 0.00	85.26 $\pm$ 0.00	<b>85.86<math>\pm</math>0.02</b>	85.25 $\pm$ 0.00	84.51 $\pm$ 0.00	83.70 $\pm$ 0.00	83.25 $\pm$ 0.00	84.77 $\pm$ 0.05
magic_gamma	<b>42.14<math>\pm</math>0.00</b>	41.51 $\pm$ 0.00	40.67 $\pm$ 0.00	40.37 $\pm$ 0.00	40.37 $\pm$ 0.00	41.04 $\pm$ 0.00	<b>41.77<math>\pm</math>0.00</b>	40.55 $\pm$ 0.00	40.24 $\pm$ 0.00	38.97 $\pm$ 0.00	38.10 $\pm$ 0.00	39.83 $\pm$ 0.00
mnist	<b>74.43<math>\pm</math>0.00</b>	73.09 $\pm$ 0.00	71.84 $\pm$ 0.00	70.69 $\pm$ 0.00	70.36 $\pm$ 0.00	72.08 $\pm$ 0.00	<b>72.72<math>\pm</math>0.00</b>	71.40 $\pm$ 0.00	70.09 $\pm$ 0.00	69.02 $\pm$ 0.00	68.60 $\pm$ 0.00	70.36 $\pm$ 0.00
mnist_avg	<b>100.00<math>\pm</math>0.00</b>	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	<b>100.00<math>\pm</math>0.00</b>	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
optdigits	<b>34.44<math>\pm</math>0.00</b>	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	<b>29.11<math>\pm</math>0.00</b>	24.76 $\pm$ 0.00	21.10 $\pm$ 0.00	17.68 $\pm$ 0.00	16.62 $\pm$ 0.00	21.85 $\pm$ 0.00
pageblocks	<b>62.78<math>\pm</math>0.00</b>	62.52 $\pm$ 0.00	62.20 $\pm$ 0.00	61.02 $\pm$ 0.00	60.30 $\pm$ 0.00	61.76 $\pm$ 0.00	<b>67.60<math>\pm</math>0.00</b>	67.74 $\pm$ 0.00	67.87 $\pm$ 0.00	66.41 $\pm$ 0.00	66.13 $\pm$ 0.00	67.15 $\pm$ 0.00
pendigits	<b>97.68<math>\pm</math>0.00</b>	97.31 $\pm$ 0.00	96.28 $\pm$ 0.00	90.01 $\pm$ 0.00	86.69 $\pm$ 0.00	95.59 $\pm$ 0.00	<b>96.99<math>\pm</math>0.00</b>	95.65 $\pm$ 0.00	81.40 $\pm$ 0.00	70.28 $\pm$ 0.00	67.39 $\pm$ 0.00	82.34 $\pm$ 0.00
pinna	<b>80.27<math>\pm</math>1.65</b>	78.05 $\pm$ 2.13	75.87 $\pm$ 2.40	74.73 $\pm$ 2.71	74.49 $\pm$ 2.71	76.68 $\pm$ 2.22	<b>75.66<math>\pm</math>2.91</b>	73.62 $\pm$ 2.98	73.42 $\pm$ 2.98	73.71 $\pm$ 2.79	73.63 $\pm$ 2.86	74.01 $\pm$ 2.75
satellite	<b>85.98<math>\pm</math>0.00</b>	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	<b>86.01<math>\pm</math>0.00</b>	85.31 $\pm$ 0.00	84.02 $\pm$ 0.00	82.19 $\pm$ 0.00	81.56 $\pm$ 0.00	83.82 $\pm$ 0.00
satimage-2	96.10 $\pm$ 0.00	96.64 $\pm$ 0.00	97.02 $\pm$ 0.00	97.39 $\pm$ 0.00	<b>97.42<math>\pm</math>0.00</b>	96.92 $\pm$ 0.00	96.69 $\pm$ 0.00	97.21 $\pm$ 0.00	97.39 $\pm$ 0.00	97.42 $\pm$ 0.00	97.42 $\pm$ 0.00	97.22 $\pm$ 0.00
shuttle	<b>99.16<math>\pm</math>0.00</b>	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	<b>97.86<math>\pm</math>0.00</b>	97.34 $\pm$ 0.00	97.28 $\pm$ 0.00	97.22 $\pm$ 0.00	97.20 $\pm$ 0.00	97.38 $\pm$ 0.00
skin	<b>98.92<math>\pm</math>0.23</b>	98.31 $\pm$ 0.20	96.81 $\pm$ 0.31	94.52 $\pm$ 0.43	93.30 $\pm$ 0.50	96.37 $\pm$ 0.25	<b>98.31<math>\pm</math>0.34</b>	96.30 $\pm$ 0.30	94.09 $\pm$ 0.51	88.39 $\pm$ 0.08	86.43 $\pm$ 0.46	92.71 $\pm$ 0.16
snip	<b>56.70<math>\pm</math>7.16</b>	54.77 $\pm$ 7.80	54.75 $\pm$ 7.81	54.76 $\pm$ 7.81	48.74 $\pm$ 10.23	53.94 $\pm$ 7.83	<b>50.26<math>\pm</math>5.73</b>	50.20 $\pm$ 5.74	50.18 $\pm$ 5.73	50.33 $\pm$ 5.85	<b>50.41<math>\pm</math>5.72</b>	50.27 $\pm$ 5.76
spambase	<b>83.93<math>\pm</math>0.00</b>	83.42 $\pm$ 0.00	83.03 $\pm$ 0.00	82.73 $\pm$ 0.00	82.63 $\pm$ 0.00	83.15 $\pm$ 0.00	<b>83.32<math>\pm</math>0.00</b>	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	<b>3.02<math>\pm</math>0.00</b>	2.89 $\pm$ 0.00	2.70 $\pm$ 0.00	2.76 $\pm$ 0.00	2.70 $\pm$ 0.00	2.82 $\pm$ 0.00	<b>2.80<math>\pm</math>0.00</b>	2.73 $\pm$ 0.00	2.74 $\pm$ 0.00	2.76 $\pm$ 0.00	2.74 $\pm$ 0.00	2.75 $\pm$ 0.00
stamps	<b>82.50<math>\pm</math>3.71</b>	77.11 $\pm$ 4.30	69.99 $\pm$ 5.29	65.85 $\pm$ 6.16	64.64 $\pm$ 6.16	72.02 $\pm$ 4.82	<b>73.26<math>\pm</math>7.70</b>	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	<b>77.22<math>\pm</math>0.00</b>	77.53 $\pm$ 0.00	77.26 $\pm$ 0.00	76.43 $\pm$ 0.00	74.75 $\pm$ 0.00	76.64 $\pm$ 0.00	<b>80.94<math>\pm</math>0.00</b>	81.09 $\pm$ 0.00	81.50 $\pm$ 0.00	81.90 $\pm$ 0.00	<b>81.93<math>\pm</math>0.00</b>	81.47 $\pm$ 0.00
vowels	<b>31.68<math>\pm</math>0.00</b>	30.30 $\pm$ 0.00	29.54 $\pm$ 0.00	27.85 $\pm$ 0.00	27.32 $\pm$ 0.00	28.38 $\pm$ 3.31	<b>25.07<math>\pm</math>3.19</b>	21.55 $\pm$ 2.53	19.65 $\pm$ 2.31	18.09 $\pm$ 1.91	17.76 $\pm$ 2.02	20.42 $\pm$ 2.34
vertebral	<b>26.96<math>\pm</math>0.00</b>	26.71 $\pm$ 0.00	25.68 $\pm$ 0.00	24.67 $\pm$ 0.00	24.49 $\pm$ 0.00	25.70 $\pm$ 0.00	<b>27.00<math>\pm</math>0.00</b>	25.82 $\pm$ 0.00	23.77 $\pm$ 0.00	24.13 $\pm$ 0.00	23.87 $\pm$ 0.00	26.62 $\pm$ 0.00
waveform	<b>96.59<math>\pm</math>2.18</b>	93.20 $\pm$ 4.25	90.32 $\pm$ 6.04	90.14 $\pm$ 5.60	88.96 $\pm$ 6.03	91.84 $\pm$ 4.57	<b>89.48<math>\pm</math>5.58</b>	89.07 $\pm$ 5.41	88.67 $\pm$ 5.41	91.70 $\pm$ 3.59	<b>91.82<math>\pm</math>3.52</b>	90.15 $\pm$ 3.97
wbc	<b>92.08<math>\pm</math>6.52</b>	89.03 $\pm$ 6.77	86.03 $\pm$ 5.81	83.79 $\pm$ 5.59	83.41 $\pm$ 5.19	86.87 $\pm$ 5.84	<b>85.38<math>\pm</math>5.46</b>	83.72 $\pm$ 5.56	82.05 $\pm$ 5.59	82.37 $\pm$ 4.91	82.33 $\pm$ 4.34	83.17 $\pm$ 5.03
wilt	<b>13.43<math>\pm</math>0.00</b>	12.86 $\pm$ 0.00	11.30 $\pm$ 0.00	10.40 $\pm$ 0.00	10.12 $\pm$ 0.00	11.52 $\pm$ 0.00	<b>12.28<math>\pm</math>0.00</b>	11.04 $\pm$ 0.00	10.11 $\pm$ 0.00	9.33 $\pm$ 0.00	9.09 $\pm$ 0.00	10.36 $\pm$ 0.00
wine	<b>75.30<math>\pm</math>0.77</b>	98.32 $\pm$ 0.52	96.09 $\pm$ 1.31</									

Table 10.2: 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  $\text{avg}$  model. We use **bold** and underline respectively to mark the **best** and the worst performance of each model to showcase the variability across different HP settings.

dataset	ICL-0.1	ICL-0.01	ICL-0.001	ICL-0.0001	ICL-1e-05	ICL-avg	DTE-C-5	DTE-C-10	DTE-C-20	DTE-C-40	DTE-C-50	DTE-C-avg
aloi	5.50±0.09	5.39±0.07	5.50±0.08	5.59±0.04	5.46±0.01	5.49±0.02	5.76±0.03	5.82±0.02	5.72±0.05	5.73±0.05	5.91±0.00	5.79±0.01
amazon	10.06±0.10	10.08±0.03	9.91±0.13	10.01±0.05	9.99±0.02	10.00±0.05	11.01±0.43	10.99±1.11	11.05±0.93	11.05±0.93	9.52±0.00	10.42±0.40
amthroid	58.53±13.33	39.69±5.23	53.66±3.08	53.85±5.44	55.94±0.72	52.34±3.72	82.46±0.50	82.25±0.34	83.27±0.63	81.52±1.18	9.52±0.00	68.86±0.20
backdoor	85.81±2.45	88.14±1.38	88.99±1.12	88.96±1.20	86.24±0.75	87.63±1.03	43.90±3.90	61.21±2.58	63.64±1.61	4.83±0.09	4.83±0.09	35.68±0.79
breastw	98.65±0.81	98.46±0.51	98.98±0.57	98.50±0.42	95.32±1.11	97.63±1.03	88.69±2.01	49.49±2.00	94.04±1.49	98.56±0.68	99.22±0.11	94.00±0.69
campaign	47.46±0.41	44.03±2.05	47.94±0.25	49.22±0.91	51.41±0.51	48.01±0.47	49.90±2.01	46.77±1.41	48.40±1.60	20.25±0.00	17.55±0.00	37.11±0.64
cardio	48.81±14.30	65.70±11.26	63.55±5.07	61.75±2.55	29.67±1.62	53.90±3.02	69.32±0.43	53.56±0.73	70.25±0.47	17.78±0.59	36.12±0.00	49.02±0.29
cardiotocography	43.21±5.03	52.56±11.4	50.69±2.12	49.01±1.29	39.79±1.23	47.05±1.34	53.83±0.94	53.56±0.73	49.12±4.17	36.12±0.00	45.75±0.91	45.75±0.91
ceiba	12.89±1.17	13.96±1.49	13.09±1.73	13.50±1.15	12.97±0.55	13.27±0.47	15.16±1.05	13.87±0.59	12.60±1.22	4.31±0.09	10.05±0.49	10.05±0.49
census	20.29±1.27	20.38±2.24	23.32±0.50	23.68±0.70	22.37±0.59	22.01±0.42	18.39±0.83	17.28±0.52	17.44±0.64	11.66±0.20	15.29±0.31	15.29±0.31
cover	22.58±13.63	10.32±4.93	35.90±22.12	47.50±19.76	40.81±22.18	31.42±3.97	76.07±1.84	66.66±3.73	53.74±2.29	1.94±0.07	35.03±0.82	35.03±0.82
donors	39.48±10.00	77.81±8.28	82.59±7.03	91.60±4.15	92.24±2.17	76.75±3.27	63.92±1.40	62.97±0.60	63.62±1.06	18.86±15.32	11.20±0.14	45.30±2.85
fault	65.37±3.40	64.58±1.81	63.83±0.74	63.93±0.94	64.32±0.96	64.41±1.09	68.85±10.35	57.17±4.68	24.97±21.19	0.34±0.03	0.34±0.03	38.79±0.41
fruit	51.45±12.34	49.97±9.27	56.24±9.47	62.41±10.30	72.24±6.51	58.46±7.58	68.85±10.35	57.17±4.68	24.97±21.19	0.34±0.03	0.34±0.03	30.33±4.25
glass	49.70±15.55	69.98±13.15	88.19±6.00	87.25±4.74	87.79±8.75	76.58±1.38	46.09±6.97	36.80±5.13	27.63±1.55	27.63±1.55	7.83±1.03	27.78±2.62
hepatitis	99.85±0.30	97.89±3.11	99.79±0.41	98.94±2.13	99.93±0.14	99.28±1.20	96.64±1.76	95.49±4.81	96.17±4.38	54.16±11.62	27.35±1.71	74.38±1.79
http	94.85±0.30	95.76±3.25	97.86±1.72	99.77±0.25	99.35±0.61	97.56±2.02	56.93±7.39	69.18±22.25	30.00±7.08	48.53±7.22	0.74±0.04	45.08±7.07
imdb	10.09±0.09	10.02±0.07	10.16±0.04	10.18±0.05	10.41±0.04	10.17±0.03	82.71±0.71	82.71±0.71	82.71±0.71	82.71±0.71	82.71±0.71	82.71±0.71
ionosphere	57.32±2.19	60.94±1.44	62.86±1.88	63.83±0.88	47.64±2.20	58.52±1.04	60.45±4.51	96.49±11.77	96.48±0.68	86.50±4.57	36.47±15.94	86.53±4.28
ionosphere	96.24±2.27	97.20±2.57	99.00±0.36	98.91±0.48	97.34±1.39	97.82±1.11	96.22±0.76	36.06±2.35	34.01±2.31	34.32±0.00	34.32±0.00	34.32±0.00
landsat	58.66±11.58	58.41±10.07	55.72±11.20	52.69±0.97	52.69±0.97	53.86±1.06	96.22±0.76	36.06±2.35	34.01±2.31	34.32±0.00	34.32±0.00	34.32±0.00
letter	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
lymphography	73.92±13.64	81.33±2.83	82.64±0.71	82.47±0.98	83.45±0.84	80.74±0.13	89.01±0.45	92.36±2.48	83.45±0.84	66.79±18.08	52.03±0.50	68.82±0.94
magic gamma	33.51±8.47	37.72±7.54	27.06±2.39	26.94±1.92	17.30±1.13	28.51±2.57	35.02±6.35	32.22±4.57	32.22±4.57	35.30±2.06	4.54±0.00	28.91±1.92
magic gamma	49.27±1.39	50.40±3.02	54.46±1.34	63.49±1.45	65.20±1.41	55.56±0.65	60.64±1.13	56.20±2.31	57.17±4.68	16.86±0.00	16.86±0.00	41.15±1.11
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00								

Table 11.1: 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- <u>avg</u>	KNN-5	KNN-10	KNN-20	KNN-40	KNN-50	KNN- <u>avg</u>
aloi	5.90 $\pm$ 0.00	5.84 $\pm$ 0.00	5.70 $\pm$ 0.00	5.90 $\pm$ 0.00	5.97 $\pm$ 0.00	5.86 $\pm$ 0.00	5.90 $\pm$ 0.00	5.64 $\pm$ 0.00	5.97 $\pm$ 0.00	6.17 $\pm$ 0.00	6.37 $\pm$ 0.00	6.01 $\pm$ 0.00
amazon	10.80 $\pm$ 0.00	10.80 $\pm$ 0.00	10.20 $\pm$ 0.00	10.20 $\pm$ 0.00	11.40 $\pm$ 0.00	10.60 $\pm$ 0.00	11.40 $\pm$ 0.00	10.20 $\pm$ 0.00	10.60 $\pm$ 0.00	11.20 $\pm$ 0.00	11.20 $\pm$ 0.00	10.92 $\pm$ 0.00
amthroid	62.55 $\pm$ 0.00	61.80 $\pm$ 0.00	60.67 $\pm$ 0.00	58.99 $\pm$ 0.00	58.80 $\pm$ 0.00	60.56 $\pm$ 0.00	61.99 $\pm$ 0.00	60.49 $\pm$ 0.00	58.43 $\pm$ 0.00	58.24 $\pm$ 0.00	56.74 $\pm$ 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	52.53 $\pm$ 1.63	40.37 $\pm$ 2.04	28.71 $\pm$ 1.50	20.21 $\pm$ 0.78	17.52 $\pm$ 0.83	31.87 $\pm$ 1.22
breastw	96.72 $\pm$ 0.64	96.23 $\pm$ 0.39	96.17 $\pm$ 0.47	95.99 $\pm$ 0.39	95.99 $\pm$ 0.39	96.22 $\pm$ 0.43	96.05 $\pm$ 0.44	96.05 $\pm$ 0.44	95.87 $\pm$ 0.28	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	51.57 $\pm$ 0.00	50.93 $\pm$ 0.00	50.37 $\pm$ 0.00	61.93 $\pm$ 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	62.98 $\pm$ 0.00	61.93 $\pm$ 0.00	61.93 $\pm$ 0.00	64.20 $\pm$ 0.00	67.61 $\pm$ 0.00	69.32 $\pm$ 0.00	65.00 $\pm$ 0.00
cardiotocography	44.64 $\pm$ 0.00	45.71 $\pm$ 0.00	47.00 $\pm$ 0.00	48.50 $\pm$ 0.00	49.14 $\pm$ 0.00	47.00 $\pm$ 0.00	46.33 $\pm$ 0.00	46.78 $\pm$ 0.00	47.85 $\pm$ 0.00	50.86 $\pm$ 0.00	51.93 $\pm$ 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	20.27 $\pm$ 0.68	20.48 $\pm$ 0.68	19.11 $\pm$ 0.61
census	22.22 $\pm$ 0.52	21.93 $\pm$ 0.52	21.46 $\pm$ 0.25	21.38 $\pm$ 0.48	21.12 $\pm$ 0.29	21.65 $\pm$ 0.14	22.23 $\pm$ 0.42	21.48 $\pm$ 0.40	21.47 $\pm$ 0.57	21.33 $\pm$ 0.65	21.26 $\pm$ 0.50	21.55 $\pm$ 0.24
cover	69.15 $\pm$ 4.12	66.87 $\pm$ 2.35	63.15 $\pm$ 2.07	55.99 $\pm$ 2.08	53.06 $\pm$ 2.23	61.65 $\pm$ 2.14	65.04 $\pm$ 1.92	61.56 $\pm$ 2.04	52.76 $\pm$ 1.94	42.69 $\pm$ 1.94	39.92 $\pm$ 1.99	52.19 $\pm$ 1.92
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	55.57 $\pm$ 0.00	56.91 $\pm$ 0.00	57.36 $\pm$ 0.00	56.32 $\pm$ 0.00	55.57 $\pm$ 0.00	55.87 $\pm$ 0.00	57.06 $\pm$ 0.00	57.50 $\pm$ 0.00	58.25 $\pm$ 0.00	56.85 $\pm$ 0.00
fraud	48.18 $\pm$ 4.56	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 $\pm$ 3.53	49.29 $\pm$ 3.57	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	49.16 $\pm$ 2.75	27.98 $\pm$ 4.21	18.37 $\pm$ 2.40	17.81 $\pm$ 2.98	29.42 $\pm$ 2.61	29.87 $\pm$ 2.62	22.55 $\pm$ 6.87	18.58 $\pm$ 4.19	17.23 $\pm$ 3.73	17.23 $\pm$ 3.73	21.09 $\pm$ 5.16
hepatitis	98.94 $\pm$ 1.41	94.16 $\pm$ 2.12	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 $\pm$ 4.51	62.31 $\pm$ 6.09	60.39 $\pm$ 6.21	60.04 $\pm$ 7.30	66.15 $\pm$ 4.71
hnp	98.50 $\pm$ 0.38	95.10 $\pm$ 2.50	88.26 $\pm$ 1.38	82.85 $\pm$ 3.80	78.96 $\pm$ 6.40	88.73 $\pm$ 2.54	100.00 $\pm$ 0.00	98.57 $\pm$ 2.86	92.67 $\pm$ 0.91	92.67 $\pm$ 0.91	92.67 $\pm$ 0.91	95.32 $\pm$ 0.75
imdb	5.40 $\pm$ 0.00	5.40 $\pm$ 0.00	5.40 $\pm$ 0.00	5.20 $\pm$ 0.00	5.20 $\pm$ 0.00	5.28 $\pm$ 0.00	5.40 $\pm$ 0.00	5.40 $\pm$ 0.00	5.40 $\pm$ 0.00	5.00 $\pm$ 0.00	5.00 $\pm$ 0.00	5.24 $\pm$ 0.00
internets	48.37 $\pm$ 0.00	48.37 $\pm$ 0.00	46.47 $\pm$ 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	45.11 $\pm$ 0.00	46.68 $\pm$ 0.00
ionosphere	51.63 $\pm$ 0.00	51.63 $\pm$ 0.00	91.63 $\pm$ 1.09	90.41 $\pm$ 1.35	89.19 $\pm$ 1.72	91.12 $\pm$ 1.24	91.40 $\pm$ 1.86	91.29 $\pm$ 1.87	85.06 $\pm$ 1.91	85.06 $\pm$ 1.91	82.75 $\pm$ 3.12	87.86 $\pm$ 1.86
landsat	52.29 $\pm$ 1.17	51.24 $\pm$ 0.00	49.06 $\pm$ 0.00	45.99 $\pm$ 0.00	45.39 $\pm$ 0.00	48.79 $\pm$ 0.00	54.25 $\pm$ 0.00	49.21 $\pm$ 0.00	45.76 $\pm$ 0.00	42.69 $\pm$ 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	76.29 $\pm$ 0.00	76.29 $\pm$ 0.00	75.40 $\pm$ 0.00	74.84 $\pm$ 0.00	71.95 $\pm$ 0.34	75.41 $\pm$ 0.53	76.17 $\pm$ 0.36	75.13 $\pm$ 0.36	74.60 $\pm$ 0.36	73.49 $\pm$ 0.36	73.00 $\pm$ 0.36	74.48 $\pm$ 0.00
magic_gamma	41.92 $\pm$ 0.00	41.15 $\pm$ 0.00	44.23 $\pm$ 0.00	44.23 $\pm$ 0.00	44.73 $\pm$ 0.00	43.15 $\pm$ 0.00	40.38 $\pm$ 0.00	43.46 $\pm$ 0.00	43.88 $\pm$ 0.00	43.88 $\pm$ 0.00	43.46 $\pm$ 0.00	43.04 $\pm$ 0.00
mask	72.71 $\pm$ 0.00	71.57 $\pm$ 0.00	70.43 $\pm$ 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	69.21 $\pm$ 0.00	69.71 $\pm$ 0.00	70.09 $\pm$ 0.00
mnist	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
optdigits	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	21.33 $\pm$ 0.00	12.00 $\pm$ 0.00	6.67 $\pm$ 0.00	2.67 $\pm$ 0.00	2.00 $\pm$ 0.00	8.93 $\pm$ 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	56.08 $\pm$ 0.00	56.27 $\pm$ 0.00	58.16 $\pm$ 0.00
pendigits	94.23 $\pm$ 0.00	92.31 $\pm$ 0.00	91.03 $\pm$ 0.00	80.13 $\pm$ 0.00	78.21 $\pm$ 0.00	87.18 $\pm$ 0.00	90.38 $\pm$ 0.00	90.38 $\pm$ 0.00	73.72 $\pm$ 0.00	64.74 $\pm$ 0.00	62.82 $\pm$ 0.00	76.41 $\pm$ 0.00
pinna	74.73 $\pm$ 2.13	71.96 $\pm$ 2.46	71.48 $\pm$ 1.96	71.44 $\pm$ 2.36	71.93 $\pm$ 1.99	72.50 $\pm$ 1.97	71.03 $\pm$ 2.53	70.18 $\pm$ 2.12	71.58 $\pm$ 2.16	71.67 $\pm$ 2.10	72.03 $\pm$ 2.02	71.30 $\pm$ 2.00
satellite	72.20 $\pm$ 0.00	71.76 $\pm$ 0.00	70.68 $\pm$ 0.00	69.79 $\pm$ 0.00	69.20 $\pm$ 0.00	70.73 $\pm$ 0.00	71.81 $\pm$ 0.00	70.73 $\pm$ 0.00	69.84 $\pm$ 0.00	67.93 $\pm$ 0.00	67.29 $\pm$ 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	90.14 $\pm$ 0.00	91.55 $\pm$ 0.00	92.96 $\pm$ 0.00	92.96 $\pm$ 0.00	92.96 $\pm$ 0.00	92.11 $\pm$ 0.00
shuttle	98.35 $\pm$ 0.00	98.23 $\pm$ 0.00	98.15 $\pm$ 0.00	98.12 $\pm$ 0.00	98.12 $\pm$ 0.00	98.19 $\pm$ 0.00	98.23 $\pm$ 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	96.83 $\pm$ 0.38	96.14 $\pm$ 0.17	94.73 $\pm$ 0.23	94.30 $\pm$ 0.27	95.82 $\pm$ 0.14	96.71 $\pm$ 0.41	95.85 $\pm$ 0.17	94.56 $\pm$ 0.29	91.73 $\pm$ 0.16	90.60 $\pm$ 0.28	93.89 $\pm$ 0.13
smtp	68.05 $\pm$ 5.12	68.05 $\pm$ 5.12	69.59 $\pm$ 3.95	69.59 $\pm$ 3.95	63.34 $\pm$ 8.31	67.73 $\pm$ 4.99	69.59 $\pm$ 3.95	69.59 $\pm$ 3.95	69.59 $\pm$ 3.95	69.59 $\pm$ 3.95	69.59 $\pm$ 3.95	69.59 $\pm$ 3.95
stamps	80.88 $\pm$ 0.00	80.29 $\pm$ 0.00	80.23 $\pm$ 0.00	79.81 $\pm$ 0.00	79.45 $\pm$ 0.00	80.13 $\pm$ 0.00	80.52 $\pm$ 0.00	79.93 $\pm$ 0.00	79.15 $\pm$ 0.00	78.98 $\pm$ 0.00	78.80 $\pm$ 0.00	79.48 $\pm$ 0.00
speech	3.28 $\pm$ 0.00	3.28 $\pm$ 0.00	1.64 $\pm$ 0.00	1.64 $\pm$ 0.00	3.28 $\pm$ 0.00	2.62 $\pm$ 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 $\pm$ 2.66	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	60.52 $\pm$ 7.08	61.78 $\pm$ 6.58	66.49 $\pm$ 7.15
thyroid	75.27 $\pm$ 0.00	75.27 $\pm$ 0.00	74.19 $\pm$ 0.00	74.19 $\pm$ 0.00	74.19 $\pm$ 0.00	74.62 $\pm$ 0.00	75.27 $\pm$ 0.00	73.12 $\pm$ 0.00	72.04 $\pm$ 0.00	73.12 $\pm$ 0.00	74.19 $\pm$ 0.00	73.55 $\pm$ 0.00
vowels	28.00 $\pm$ 0.00	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	17.18 $\pm$ 2.01	12.19 $\pm$ 3.52	9.07 $\pm$ 3.30	8.51 $\pm$ 2.65	13.99 $\pm$ 2.95
vertebral	49.03 $\pm$ 4.93	28.00 $\pm$ 0.00	28.00 $\pm$ 0.00	30.00 $\pm$ 0.00	30.00 $\pm$ 0.00	28.80 $\pm$ 0.00	26.00 $\pm$ 0.00	26.00 $\pm$ 0.00	30.00 $\pm$ 0.00	26.00 $\pm$ 0.00	26.00 $\pm$ 0.00	26.80 $\pm$ 0.00
waveform	26.00 $\pm$ 0.00	26.00 $\pm$ 0.00	27.00 $\pm$ 0.00	28.00 $\pm$ 0.00	28.00 $\pm$ 0.00	27.00 $\pm$ 0.00	27.00 $\pm$ 0.00	27.00 $\pm$ 0.00	27.00 $\pm$ 0.00	27.00 $\pm$ 0.00	27.00 $\pm$ 0.00	27.00 $\pm$ 0.00
waveform	89.35 $\pm$ 3.08	88.05 $\pm$ 3.35	88.05 $\pm$ 3.35	89.16 $\pm$ 2.38	89.16 $\pm$ 2.38	88.75 $\pm$ 2.48	83.63 $\pm$ 4.53	84.82 $\pm$ 3.86	83.72 $\pm$ 5.17	89.16 $\pm$ 2.38	89.16 $\pm$ 2.38	86.10 $\pm$ 2.43
wbc	85.05 $\pm$ 6.11	83.36 $\pm$ 7.16	81.00 $\pm$ 4.08	81.00 $\pm$ 4.08	81.00 $\pm$ 4.08	82.28 $\pm$ 4.91	80.92 $\pm$ 4.47	78.72 $\pm$ 5.61	79.49 $\pm$ 5.02	77.95 $\pm$ 6.52	76.84 $\pm$ 4.81	78.79 $\pm$ 5.12
wilt	3.50 $\pm$ 0.00	2.33 $\pm$ 0.00	1.36 $\pm$ 0.00	0.78 $\pm$ 0.00	0.78 $\pm$ 0.00	1.79 $\pm$ 0.00	2.33 $\pm$ 0.00	1.36 $\pm$ 0.00	0.78 $\pm$ 0.00	0.78 $\pm$ 0.00	0.78 $\pm$ 0.00	1.25 $\pm$ 0.00
wilt	3.50 $\pm$ 0.00	2.33 $\pm$ 0.00	1.36 $\pm$ 0.00	0.78 $\pm$ 0.00	0.78 $\pm$ 0.00	1.79 $\pm$ 0.0						



Table 11.2: 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  $\text{avg}$  model. We use **bold** and underline respectively to mark the **best** and the **worst** performance of each model to showcase the variability across different HP settings.

dataset	ICL-0.1	ICL-0.01	ICL-0.0001	ICL-1e-05	ICL-avr	DTE-C-5	DTE-C-10	DTE-C-20	DTE-C-40	DTE-C-50	DTE-C-avr
amazon	4.51±0.69	4.34±0.42	5.28±0.47	4.68±0.30	4.94±0.07	4.75±0.27	11.96±1.68	4.28±0.10	4.51±0.17	0.00±0.00	3.56±0.03
aloi	10.44±0.46	9.76±0.34	9.92±0.84	10.08±0.35	9.59±0.32	11.48±0.97	11.96±1.68	11.60±0.27	0.00±0.00	0.00±0.00	7.01±0.72
amthroid	54.87±13.24	42.25±5.45	53.45±4.13	54.72±5.45	57.53±2.97	77.23±0.25	77.53±0.85	75.43±0.93	0.00±0.00	0.00±0.00	61.63±0.20
backdoor	87.17±0.98	87.32±0.99	87.11±0.99	86.85±0.95	85.37±0.11	86.76±1.00	83.03±2.14	84.50±1.00	0.00±0.00	0.00±0.00	42.75±1.54
breastw	95.98±0.34	96.07±0.94	96.80±0.66	96.11±0.75	96.48±0.28	51.98±0.59	90.10±1.35	92.46±1.78	95.31±0.70	96.11±0.44	92.50±0.79
campaign	48.12±0.36	46.81±1.72	50.68±0.66	51.37±0.85	53.40±0.51	50.07±0.49	51.98±0.59	52.33±1.00	0.00±0.00	0.00±0.00	31.35±0.44
cardio	49.09±11.28	61.93±5.57	58.86±1.59	57.95±2.30	40.57±4.28	53.68±2.83	58.30±0.58	57.84±0.43	58.07±0.23	0.34±0.68	34.91±0.09
cardiotocography	36.14±1.28	41.97±1.73	39.18±4.49	35.36±2.14	32.66±1.59	36.88±1.27	39.91±1.05	39.48±1.54	37.73±1.73	62.02±0.00	48.23±0.39
celeba	15.42±2.29	17.97±2.55	17.20±1.92	17.46±1.17	16.17±0.65	16.84±0.88	19.18±2.74	17.12±1.45	14.31±2.28	0.00±0.00	10.12±1.04
census	22.72±1.73	24.06±2.05	25.80±1.34	27.15±1.12	24.06±1.31	24.76±0.50	26.51±2.37	68.92±4.22	46.54±3.93	0.00±0.00	10.31±0.52
cover	26.77±15.24	15.53±8.17	42.68±17.00	53.70±16.88	44.34±20.24	36.61±2.59	76.51±2.37	75.05±7.18	63.08±1.70	11.80±23.60	38.39±0.62
donors	43.71±11.52	81.85±3.56	83.59±4.47	89.28±2.66	92.77±1.39	78.24±2.61	59.13±1.36	55.16±0.81	55.33±1.66	97.03±0.00	47.58±4.61
fault	60.33±3.36	59.52±2.09	58.57±1.11	58.37±0.79	58.87±1.51	62.30±4.91	75.61±7.76	54.25±5.90	22.55±24.73	0.00±0.00	72.00±0.47
glass	57.54±10.13	48.97±8.02	58.90±6.77	66.88±4.88	79.18±3.51	62.30±4.91	75.61±7.76	54.25±5.90	22.55±24.73	0.00±0.00	30.48±4.66
hepatitis	43.53±20.21	57.05±16.03	84.05±6.11	87.24±5.04	82.50±5.41	70.87±7.76	75.61±7.76	54.25±5.90	22.55±24.73	0.00±0.00	24.15±2.73
hyp	99.64±0.71	94.69±7.07	97.69±1.51	99.36±0.19	99.14±0.33	97.23±2.45	96.40±3.01	94.63±4.90	92.51±3.12	51.68±9.78	18.97±10.60
imdb	10.52±0.65	10.44±0.54	9.84±0.37	9.76±0.15	10.40±0.33	10.19±0.23	8.40±1.23	7.32±1.04	6.95±2.24	0.00±0.00	4.47±0.43
internats	55.92±2.66	57.45±0.61	57.77±1.10	58.26±0.50	59.02±1.45	55.68±0.94	67.99±1.98	68.87±0.95	64.95±2.42	41.85±0.00	44.70±0.87
ionosphere	92.64±4.66	91.41±4.67	93.86±1.63	94.48±0.71	94.49±1.56	93.38±2.21	89.67±1.44	89.41±1.53	78.12±2.07	45.44±16.82	79.23±4.14
landsat	49.50±1.24	47.92±1.88	54.51±0.52	54.25±0.74	47.97±1.09	53.45±0.79	55.47±3.76	31.07±3.98	50.29±8.34	54.46±0.00	41.47±2.39
letter	6.80±4.21	4.00±1.55	3.60±1.36	3.20±0.75	11.06±2.06	5.84±1.11	7.40±1.50	3.20±1.11	0.00±0.00	0.00±0.00	17.21±0.37
lymphography	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00
magic_gamma	36.68±2.76	69.99±0.87	69.99±0.87	70.21±0.80	71.64±0.90	69.34±1.69	70.87±1.48	80.94±0.33	80.19±0.73	89.72±0.81	96.10±0.00
mnist	45.37±1.84	46.20±1.18	27.69±1.80	29.08±2.63	18.15±1.23	29.94±2.59	32.69±5.73	34.62±2.71	35.08±2.68	39.69±2.19	85.36±1.69
mnist	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	89.36±1.65	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	60.00±0.00
optdigits	29.87±5.68	41.20±3.49	44.67±5.56	71.73±0.53	46.00±1.81	46.69±2.20	62.24±0.98	62.71±1.24	61.25±0.32	60.47±0.69	49.33±0.37
pageblocks	62.16±2.84	64.08±2.96	62.31±2.65	63.88±1.08	61.03±1.05	62.93±0.29	63.97±6.36	54.36±7.40	43.46±7.30	0.00±0.00	32.36±1.43
pendigits	46.03±17.81	60.51±10.58	66.03±5.18	62.03±1.57	51.03±1.78	66.14±2.35	63.89±4.28	64.71±3.39	67.99±2.34	65.96±4.44	65.74±2.10
pinna	68.77±4.96	68.29±2.03	74.97±1.27	71.40±1.67	68.83±2.08	73.62±2.69	72.06±0.60	72.97±0.77	72.63±0.32	91.51±9.02	81.04±1.96
satellite	65.94±10.44	72.95±2.96	76.27±0.81	78.70±0.90	74.22±0.57	69.69±7.51	78.31±3.03	60.85±5.87	46.20±5.07	26.20±32.11	42.31±7.14
satimage-2	38.03±8.59	72.68±32.26	91.83±0.56	89.58±1.13	96.34±1.28	98.31±0.24	98.00±0.01	97.98±0.00	97.73±0.10	92.83±2.84	77.31±0.58
skin	97.17±1.11	98.27±0.10	98.83±0.14	98.91±0.12	98.74±1.31	56.43±9.55	82.15±0.39	82.23±0.37	81.66±0.57	80.48±0.35	76.25±0.32
smtp	39.28±18.81	72.09±8.03	67.99±9.92	54.29±12.30	49.74±10.83	69.59±3.95	69.59±3.95	69.59±3.95	49.07±5.05	27.32±14.61	43.11±3.67
spambase	76.97±2.12	76.88±3.76	80.38±0.48	80.23±0.62	74.93±0.44	77.87±0.46	80.19±0.18	80.27±0.41	80.56±0.29	87.61±0.00	83.25±0.07
speech	2.95±1.23	3.61±1.23	3.28±1.04	2.95±2.41	4.59±1.61	3.48±0.92	3.93±0.80	2.95±1.61	3.28±1.80	0.00±0.00	2.03±0.25
stamps	34.84±12.18	42.45±13.97	83.03±3.34	76.24±6.92	73.47±7.50	62.00±6.23	63.62±0.79	58.37±5.83	57.63±9.24	56.93±11.83	14.60±29.20
thyroid	68.39±4.17	61.29±1.80	61.08±4.63	63.87±3.57	33.76±5.75	57.68±1.74	73.76±2.21	76.13±1.72	75.27±0.00	78.28±0.80	50.23±10.89
vertebal	23.17±6.96	29.15±9.05	52.53±11.60	64.81±3.19	54.89±6.76	44.91±2.82	43.88±5.22	36.26±11.73	32.59±10.94	24.68±8.71	3.48±6.96
vowels	22.40±17.59	18.80±6.52	30.80±7.44	30.80±12.43	24.00±7.04	25.36±2.94	40.00±5.51	40.80±3.25	45.20±0.98	0.00±0.00	25.20±1.36
waveform	47.80±3.06	38.20±17.57	35.40±7.47	11.60±1.62	28.40±4.22	32.28±5.79	11.20±1.72	11.80±2.14	12.20±1.60	0.00±0.00	7.04±0.23
wbc	78.86±7.89	84.03±7.23	96.89±4.35	95.71±4.90	90.32±5.25	89.16±3.80	40.59±7.11	34.81±6.65	40.96±12.93	67.88±5.79	36.84±4.92
wdbc	64.32±18.21	80.42±5.68	84.10±6.72	87.13±4.26	78.33±5.25	78.86±2.34	78.62±5.69	70.31±4.12	77.59±7.50	10.43±20.87	47.39±3.91
wilt	61.57±7.00	11.46±4.93	37.04±8.60	39.61±4.50	37.59±2.50	26.35±2.35	19.53±2.50	19.53±2.50	5.37±2.97	16.96±5.74	12.36±2.38
wine	97.95±4.09	96.72±3.06	98.28±2.23	98.18±3.64	97.67±3.23	97.76±2.03	98.86±1.44	95.05±5.74	97.27±3.64	41.73±22.73	66.38±4.38
wpc	81										
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87	49.23±1.21	49.55±2.12	49.55±2.12	48.68±1.91	96.22±0.00	59.93±5.76
yeast	46.31±22.91	48.99±1.78	49.35±0.54	49.31±0.62	48.41±0.87</						

Table 12.1: 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  $\text{avg}$  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  $\text{avg}$  variants of top-4 baselines + FoMo-0D w/  $D = 20$ ).

Dialect	FAKED(=20)	PRE-SP	ASN	CTL	DTLC	LOF	POWERSHIFT	SLAD	DDPH	OCSTM	DTFLG	BLIND	MCD
antennas	54.74(1.4188)	54.76(1.2341)	54.81(1.0414)	54.81(1.0812)	54.71(1.0228)	57.84(0.9026)	54.60(1.1650)	50.76(1.2018)	55.09(1.3125)	50.97(1.1316)	50.97(1.1316)	50.97(1.1316)	50.97(1.1316)
antennas	83.17(0.0224)	82.82(=0.001)	82.51(0.051)	82.14(1.0228)	80.57(1.1510)	75.84(0.901)	58.17(1.0810)	59.94(1.0071)	55.09(1.3125)	58.45(=0.001)	57.63(1.3260)	58.45(1.0512)	58.45(1.0512)
backdoor	78.90(1.7117)	78.97(1.88626)	79.75(1.0585)	80.11(1.1710)	79.55(1.2010)	88.33(0.415)	90.14(1.0810)	90.14(1.0810)	90.14(1.0810)	88.45(1.051)	87.64(1.4461)	90.28(1.1520)	90.28(1.1520)
backdoor	93.31(1.771)	93.31(1.771)	93.31(1.771)	93.31(1.771)	93.31(1.771)	93.31(1.771)	93.31(1.771)	93.31(1.771)	93.31(1.771)	93.31(1.771)	93.31(1.771)	93.31(1.771)	93.31(1.771)
cardio	65.54(0.25)	65.54(0.25)	65.54(0.25)	65.54(0.25)	65.54(0.25)	65.54(0.25)	65.54(0.25)	65.54(0.25)	65.54(0.25)	65.54(0.25)	65.54(0.25)	65.54(0.25)	65.54(0.25)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	84.71(0.6613)	84.71(0.6613)	84.71(0.6613)	84.71(0.6613)	84.71(0.6613)	84.71(0.6613)	84.71(0.6613)	84.71(0.6613)	84.71(0.6613)	84.71(0.6613)	84.71(0.6613)	84.71(0.6613)	84.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)	94.71(0.6613)
catapult	94.71(0												

Table 12.2: 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.

Dataset	VAE	PCA	PlanarFlow	HBOS	GANomaly	GOAD	DIF	COPOD	ECOD	DeepSVDD	LODA	DAGMM	DROC	DTE-NP <sup>++</sup>	KNN <sup>++</sup>	ICL <sup>++</sup>	DTE-C <sup>++</sup>
air	54.04 $\pm$ 0.0033	54.04 $\pm$ 0.0033	48.52 $\pm$ 0.34(29)	52.23 $\pm$ 0.0071	55.94 $\pm$ 1.65(5)	48.01 $\pm$ 0.92(30)	51.11 $\pm$ 0.00(13)	50.99 $\pm$ 2.05(15)	50.99 $\pm$ 2.05(15)	50.99 $\pm$ 2.05(15)	49.24 $\pm$ 2.78(25)	50.84 $\pm$ 2.97(17)	50.00 $\pm$ 0.00(22)	51.25 $\pm$ 0.00(11)	51.61 $\pm$ 0.00(9)	47.63 $\pm$ 0.15(31)	50.29 $\pm$ 0.09(21)
amazon	85.44 $\pm$ 0.0020	85.44 $\pm$ 0.0020	87.52 $\pm$ 0.14(3)	86.02 $\pm$ 0.0031	87.52 $\pm$ 0.14(3)	81.01 $\pm$ 0.16(25)	83.34 $\pm$ 0.16(25)	76.77 $\pm$ 0.00(30)	78.45 $\pm$ 0.00(26)	78.45 $\pm$ 0.00(26)	77.35 $\pm$ 0.15(27)	72.35 $\pm$ 0.15(29)	72.35 $\pm$ 0.15(29)	92.34 $\pm$ 0.00(7)	92.34 $\pm$ 0.00(7)	84.27 $\pm$ 0.31(2)	88.05 $\pm$ 0.09(21)
backdoor	64.57 $\pm$ 0.0524	64.57 $\pm$ 0.0524	67.03 $\pm$ 0.14(3)	66.02 $\pm$ 0.0031	67.03 $\pm$ 0.14(3)	52.94 $\pm$ 0.48(28)	83.73 $\pm$ 0.00(16)	50.02 $\pm$ 0.00(30)	50.02 $\pm$ 0.00(30)	50.02 $\pm$ 0.00(30)	74.62 $\pm$ 0.63(32)	74.62 $\pm$ 0.63(32)	74.62 $\pm$ 0.63(32)	92.34 $\pm$ 0.00(7)	92.34 $\pm$ 0.00(7)	93.33 $\pm$ 0.00(16)	74.67 $\pm$ 0.36(20)
breastw	69.22 $\pm$ 0.18(7)	69.22 $\pm$ 0.18(7)	69.22 $\pm$ 0.18(7)	69.22 $\pm$ 0.18(7)	69.22 $\pm$ 0.18(7)	98.86 $\pm$ 0.33(15)	56.74 $\pm$ 1.19(31)	99.46 $\pm$ 0.09(3)	99.46 $\pm$ 0.09(3)	99.46 $\pm$ 0.09(3)	98.33 $\pm$ 0.35(21)	89.33 $\pm$ 0.19(27)	98.33 $\pm$ 0.35(21)	98.67 $\pm$ 0.02(18)	99.16 $\pm$ 0.02(19)	98.75 $\pm$ 0.18(16)	96.52 $\pm$ 0.53(24)
campagin	77.77 $\pm$ 0.00(11)	77.77 $\pm$ 0.00(11)	77.77 $\pm$ 0.00(11)	77.77 $\pm$ 0.00(11)	77.77 $\pm$ 0.00(11)	47.89 $\pm$ 0.22(32)	57.87 $\pm$ 0.00(30)	78.86 $\pm$ 0.00(3)	78.86 $\pm$ 0.00(3)	78.86 $\pm$ 0.00(3)	58.88 $\pm$ 0.48(29)	61.47 $\pm$ 0.27(28)	50.00 $\pm$ 0.00(31)	78.76 $\pm$ 0.00(3)	78.66 $\pm$ 0.00(4)	78.33 $\pm$ 0.52(7)	67.18 $\pm$ 0.74(24)
cardio	96.55 $\pm$ 0.00(1)	96.55 $\pm$ 0.00(1)	89.04 $\pm$ 0.31(8)	80.7 $\pm$ 0.00(24)	86.06 $\pm$ 0.31(8)	96.01 $\pm$ 0.22(3)	68.53 $\pm$ 0.00(30)	93.16 $\pm$ 0.00(9)	94.95 $\pm$ 0.00(5)	94.95 $\pm$ 0.00(5)	77.92 $\pm$ 0.88(26)	71.79 $\pm$ 0.27(28)	50.00 $\pm$ 0.00(31)	92.47 $\pm$ 0.00(11)	93.07 $\pm$ 0.00(10)	77.71 $\pm$ 0.76(27)	72.66 $\pm$ 0.21(29)
cardiography	78.89 $\pm$ 0.00(3)	78.89 $\pm$ 0.00(3)	62.77 $\pm$ 0.62(18)	61.24 $\pm$ 0.00(21)	67.78 $\pm$ 0.62(18)	76.06 $\pm$ 1.4(4)	41.79 $\pm$ 0.00(32)	66.35 $\pm$ 0.00(12)	79.3 $\pm$ 0.00(1)	79.3 $\pm$ 0.00(1)	72.79 $\pm$ 0.76(7)	67.11 $\pm$ 0.90(11)	45.98 $\pm$ 0.16(31)	63.21 $\pm$ 0.00(17)	65.48 $\pm$ 0.00(13)	50.85 $\pm$ 1.24(28)	55.9 $\pm$ 0.23(24)
cardio	78.89 $\pm$ 0.00(3)	78.89 $\pm$ 0.00(3)	62.77 $\pm$ 0.62(18)	61.24 $\pm$ 0.00(21)	67.78 $\pm$ 0.62(18)	76.06 $\pm$ 1.4(4)	41.79 $\pm$ 0.00(32)	66.35 $\pm$ 0.00(12)	79.3 $\pm$ 0.00(1)	79.3 $\pm$ 0.00(1)	72.79 $\pm$ 0.76(7)	67.11 $\pm$ 0.90(11)	45.98 $\pm$ 0.16(31)	63.21 $\pm$ 0.00(17)	65.48 $\pm$ 0.00(13)	50.85 $\pm$ 1.24(28)	55.9 $\pm$ 0.23(24)
cover	70.82 $\pm$ 0.02(9)	70.82 $\pm$ 0.02(9)	68.07 $\pm$ 0.56(14)	62.5 $\pm$ 0.00(16)	68.07 $\pm$ 0.56(14)	35.24 $\pm$ 0.19(32)	61.85 $\pm$ 0.00(18)	50.00 $\pm$ 0.00(30)	50.00 $\pm$ 0.00(30)	50.00 $\pm$ 0.00(30)	51.12 $\pm$ 0.19(29)	52.84 $\pm$ 1.09(28)	53.36 $\pm$ 0.62(5)	72.11 $\pm$ 0.16(6)	71.84 $\pm$ 0.15(6)	72.25 $\pm$ 0.16(8)	61.41 $\pm$ 0.71(9)
donors	94.85 $\pm$ 0.15(16)	94.85 $\pm$ 0.15(16)	76.33 $\pm$ 0.11(24)	71.11 $\pm$ 0.82(27)	76.33 $\pm$ 0.11(24)	13.83 $\pm$ 0.13(32)	97.09 $\pm$ 0.24(28)	88.2 $\pm$ 0.27(21)	88.2 $\pm$ 0.27(21)	88.2 $\pm$ 0.27(21)	94.93 $\pm$ 0.07(14)	75.98 $\pm$ 0.16(9(31)	74.18 $\pm$ 0.20(9(28)	97.37 $\pm$ 0.00(19)	96.73 $\pm$ 0.00(20)	91.42 $\pm$ 0.17(41)	78.4 $\pm$ 0.14(23)
fault	88.66 $\pm$ 0.25(19)	88.66 $\pm$ 0.25(19)	81.64 $\pm$ 0.32(15)	81.19 $\pm$ 0.62(27)	75.34 $\pm$ 1.11(9(27)	33.57 $\pm$ 0.16(32)	90.04 $\pm$ 1.79(16)	81.5 $\pm$ 0.21(24)	81.5 $\pm$ 0.21(24)	81.5 $\pm$ 0.21(24)	63.52 $\pm$ 0.27(30)	62.15 $\pm$ 0.16(9(31)	74.18 $\pm$ 0.20(9(28)	97.37 $\pm$ 0.00(19)	96.3 $\pm$ 0.00(17)	96.3 $\pm$ 0.00(17)	80.17 $\pm$ 0.55(26)
fraud	95.87 $\pm$ 0.00(20)	95.87 $\pm$ 0.00(20)	57.51 $\pm$ 0.25(15)	53.06 $\pm$ 0.00(26)	59.54 $\pm$ 0.25(15)	69.75 $\pm$ 0.31(31)	82.58 $\pm$ 0.27(29)	94.3 $\pm$ 0.16(28)	94.3 $\pm$ 0.16(28)	94.3 $\pm$ 0.16(28)	50.27 $\pm$ 0.18(29)	52.83 $\pm$ 0.17(27)	55.73 $\pm$ 0.33(23)	59.17 $\pm$ 0.00(12)	55.73 $\pm$ 0.33(23)	95.57 $\pm$ 0.03(4)	94.23 $\pm$ 0.87(17)
glass	95.87 $\pm$ 0.00(20)	95.87 $\pm$ 0.00(20)	90.72 $\pm$ 0.26(25)	95.02 $\pm$ 0.00(7)	92.52 $\pm$ 0.85(21)	69.75 $\pm$ 0.31(31)	82.58 $\pm$ 0.27(29)	94.3 $\pm$ 0.16(28)	94.3 $\pm$ 0.16(28)	94.3 $\pm$ 0.16(28)	50.27 $\pm$ 0.18(29)	52.83 $\pm$ 0.17(27)	55.73 $\pm$ 0.33(23)	59.17 $\pm$ 0.00(12)	55.73 $\pm$ 0.33(23)	95.57 $\pm$ 0.03(4)	94.23 $\pm$ 0.87(17)
higgs	72.81 $\pm$ 1.66(25)	72.81 $\pm$ 1.66(25)	73.44 $\pm$ 2.22(24)	82.59 $\pm$ 0.33(18)	77.77 $\pm$ 0.72(20)	90.03 $\pm$ 1.24(4(32)	96.53 $\pm$ 1.39(5)	78.01 $\pm$ 0.94(23)	78.01 $\pm$ 0.94(23)	78.01 $\pm$ 0.94(23)	67.34 $\pm$ 0.39(28)	65.33 $\pm$ 0.53(30)	64.89 $\pm$ 0.54(31)	90.83 $\pm$ 0.09(19)	87.3 $\pm$ 1.59(14)	95.27 $\pm$ 0.14(5)	76.95 $\pm$ 2.38(22)
insects	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)
inb	47.91 $\pm$ 0.00(28)	47.91 $\pm$ 0.00(28)	40.23 $\pm$ 0.25(22)	49.94 $\pm$ 0.00(16)	51.58 $\pm$ 0.00(16)	48.46 $\pm$ 0.65(26)	51.4 $\pm$ 0.49(4)	51.05 $\pm$ 0.00(8)	46.08 $\pm$ 0.00(32)	46.08 $\pm$ 0.00(32)	47.23 $\pm$ 0.22(41)	48.6 $\pm$ 0.30(25)	51.35 $\pm$ 0.20(5)	50.35 $\pm$ 0.00(13)	50.35 $\pm$ 0.00(13)	53.38 $\pm$ 0.12(1)	49.27 $\pm$ 1.12(21)
internals	65.12 $\pm$ 0.00(25)	65.12 $\pm$ 0.00(25)	70.87 $\pm$ 0.84(9)	70.87 $\pm$ 0.84(9)	70.87 $\pm$ 0.84(9)	66.01 $\pm$ 0.00(16)	72.96 $\pm$ 0.32(13)	66.01 $\pm$ 0.00(16)	66.01 $\pm$ 0.00(16)	66.01 $\pm$ 0.00(16)	78.73 $\pm$ 0.34(24)	49.47 $\pm$ 0.58(28)	53.47 $\pm$ 0.58(27)	67.65 $\pm$ 0.00(13)	67.65 $\pm$ 0.00(13)	72.53 $\pm$ 0.05(4)	66.79 $\pm$ 0.73(14)
ionosphere	89.78 $\pm$ 1.24(24)	89.78 $\pm$ 1.24(24)	90.85 $\pm$ 1.21(9)	70.68 $\pm$ 2.85(31)	91.89 $\pm$ 0.20(20)	91.54 $\pm$ 0.55(22)	93.59 $\pm$ 0.06(17)	78.32 $\pm$ 0.13(28)	71.77 $\pm$ 1.43(30)	71.77 $\pm$ 1.43(30)	85.56 $\pm$ 0.33(29)	56.27 $\pm$ 0.34(61)	61.14 $\pm$ 0.28(52)	97.51 $\pm$ 0.63(5)	95.12 $\pm$ 0.92(16)	97.78 $\pm$ 1.23(3)	83.01 $\pm$ 0.81(27)
landsat	54.24 $\pm$ 0.24(50)	54.24 $\pm$ 0.24(50)	50.85 $\pm$ 1.21(9)	70.68 $\pm$ 2.85(31)	91.89 $\pm$ 0.20(20)	91.54 $\pm$ 0.55(22)	93.59 $\pm$ 0.06(17)	78.32 $\pm$ 0.13(28)	71.77 $\pm$ 1.43(30)	71.77 $\pm$ 1.43(30)	85.56 $\pm$ 0.33(29)	56.27 $\pm$ 0.34(61)	61.14 $\pm$ 0.28(52)	97.51 $\pm$ 0.63(5)	95.12 $\pm$ 0.92(16)	97.78 $\pm$ 1.23(3)	83.01 $\pm$ 0.81(27)
letter	30.3 $\pm$ 0.00(31)	30.3 $\pm$ 0.00(31)	38.73 $\pm$ 0.53(13)	35.91 $\pm$ 0.00(1)	34.02 $\pm$ 0.74(23)	31.08 $\pm$ 0.55(26)	74.07 $\pm$ 0.00(1)	36.53 $\pm$ 0.00(1)	42.01 $\pm$ 0.00(32)	42.01 $\pm$ 0.00(32)	85.66 $\pm$ 0.33(29)	56.27 $\pm$ 0.34(61)	61.14 $\pm$ 0.28(52)	97.51 $\pm$ 0.63(5)	95.12 $\pm$ 0.92(16)	97.78 $\pm$ 1.23(3)	83.01 $\pm$ 0.81(27)
lymphography	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)
magnetic gamma	70.84 $\pm$ 0.00(25)	70.84 $\pm$ 0.00(25)	74.12 $\pm$ 0.75(20)	74.12 $\pm$ 0.75(20)	74.12 $\pm$ 0.75(20)	91.81 $\pm$ 1.63(2)	69.46 $\pm$ 0.60(28)	63.86 $\pm$ 0.00(28)	63.86 $\pm$ 0.00(28)	63.86 $\pm$ 0.00(28)	67.04 $\pm$ 0.18(31)	67.04 $\pm$ 0.18(31)	67.04 $\pm$ 0.18(31)	99.84 $\pm$ 0.00(1)	99.84 $\pm$ 0.00(1)	75.19 $\pm$ 0.51(7)	75.19 $\pm$ 0.51(7)
mnist	90.21 $\pm$ 0.00(15)	90.21 $\pm$ 0.00(15)	81.9 $\pm$ 0.27(20)	62.34 $\pm$ 0.00(29)	77.96 $\pm$ 0.44(23)	90.07 $\pm$ 0.35(13)	50.21 $\pm$ 0.00(30)	50.00 $\pm$ 0.00(30)	50.00 $\pm$ 0.00(30)	50.00 $\pm$ 0.00(30)	63.71 $\pm$ 0.03(27)	72.19 $\pm$ 1.62(5)	83.13 $\pm$ 1.64(19)	93.4 $\pm$ 0.00(3)	93.4 $\pm$ 0.00(3)	91.04 $\pm$ 0.04(1)	72.55 $\pm$ 0.64(24)
nuk	100.0 $\pm$ 0.00(5)	100.0 $\pm$ 0.00(5)	76.65 $\pm$ 1.82(31)	89.92 $\pm$ 0.00(1)	100.0 $\pm$ 0.00(5)	100.0 $\pm$ 0.00(5)	97.76 $\pm$ 0.00(24)	99.71 $\pm$ 0.00(17)	99.71 $\pm$ 0.00(17)	99.71 $\pm$ 0.00(17)	99.71 $\pm$ 0.00(17)	99.71 $\pm$ 0.00(17)	99.71 $\pm$ 0.00(17)	100.0 $\pm$ 0.00(5)	100.0 $\pm$ 0.00(5)	97.18 $\pm$ 0.37(25)	80.0 $\pm$ 0.00(30)
optdigits	86.16 $\pm$ 0.00(20)	86.16 $\pm$ 0.00(20)	84.85 $\pm$ 1.19(23)	65.62 $\pm$ 0.00(32)	72.83 $\pm$ 0.94(31)	88.05 $\pm$ 1.19(13)	87.29 $\pm$ 0.00(16)	80.85 $\pm$ 0.00(26)	80.85 $\pm$ 0.00(26)	80.85 $\pm$ 0.00(26)	83.62 $\pm$ 0.26(24)	86.29 $\pm$ 0.14(42)	92.33 $\pm$ 0.18(1)	89.32 $\pm$ 0.00(8)	89.32 $\pm$ 0.00(8)	90.77 $\pm$ 0.00(10)	67.73 $\pm$ 1.71(20)
pageshops	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)	94.37 $\pm$ 0.00(18)
pendigits	73.18 $\pm$ 1.93(12)	73.18 $\pm$ 1.93(12)	72.28 $\pm$ 1.19(14)	72.28 $\pm$ 1.19(14)	72.28 $\pm$ 1.19(14)	60.45 $\pm$ 0.10(28)	89.7 $\pm$ 0.20(25)	90.74 $\pm$ 0.00(23)	92.95 $\pm$ 0.00(14)	92.95 $\pm$ 0.00(14)	92.13 $\pm$ 1.03(22)	56.29 $\pm$ 0.14(42)	92.33 $\pm$ 0.18(1)	89.32 $\pm$ 0.00(8)	89.32 $\pm$ 0.00(8)	90.77 $\pm$ 0.00(10)	67.73 $\pm$ 1.71(20)
pima	74.44 $\pm$ 0.02(19)	74.44 $\pm$ 0.02(19)	72.33 $\pm$ 1.73(25)	72.33 $\pm$ 1.73(25)	72.33 $\pm$ 1.73(25)	68.76 $\pm$ 0.92(28)	66.91 $\pm$ 0.00(30)	68.34 $\pm$ 0.00(29)	62.22 $\pm$ 0.00(32)	62.22 $\pm$ 0.00(32)	69.73 $\pm$ 1.15(26)	72.79 $\pm$ 0.20(23)	73.38 $\pm$ 0.44(21)	81.43 $\pm$ 0.00(9)	80.47 $\pm$ 0.00(10)	83.57 $\pm$ 0.48(5)	68.77 $\pm$ 0.26(27)
satellite	98.88 $\pm$ 0.00(4)	98.88 $\pm$ 0.00(4)	96.67 $\pm$ 0.57(27)	97.95 $\pm$ 0.00(22)	97.57 $\pm$ 1.19(24)	98.99 $\pm$ 0.06(17)	86.94 $\pm$ 0.00(31)	97.92 $\pm$ 0.00(23)	97.92 $\pm$ 0.00(23)	97.92 $\pm$ 0.00(23)	98.67 $\pm$ 0.57(27)	91.82 $\pm$ 0.44(30)	99.22 $\pm$ 0.02(15)	99.77 $\pm$ 0.00(3)	99.77 $\pm$ 0.00(3)	96.99 $\pm$ 0.25(26)	83.3 $\pm$ 0.47(32)
shuttle	99.83<																



Table 13.2: 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.

Dataset	VAE	PCA	PlanFlow	HBOs	GANomaly	GQAD	DIF	COPOD	ECOD	DapSVDD	LODA	DAGMM	DRDCC	DTE-NP <sub>2</sub>	RNN <sup>sg</sup>	ICL <sup>sg</sup>	DTE-C <sup>sg</sup>
aloi	6.54 $\pm$ 0.05	6.54 $\pm$ 0.05	5.84 $\pm$ 0.332	6.42 $\pm$ 0.08	8.09 $\pm$ 1.2701	5.71 $\pm$ 0.2428	5.8 $\pm$ 0.024	5.72 $\pm$ 0.027	6.06 $\pm$ 0.041	6.23 $\pm$ 0.35(105)	5.93 $\pm$ 0.5(21)	6.07 $\pm$ 0.55(31)	5.91 $\pm$ 0.022	6.02 $\pm$ 0.0(6.5)	6.09 $\pm$ 0.0(12)	5.49 $\pm$ 0.02(31)	5.79 $\pm$ 0.0(125)
amazon	10.74 $\pm$ 0.075	10.72 $\pm$ 0.0075	9.84 $\pm$ 0.59(40)	11.1 $\pm$ 0.0095	8.9 $\pm$ 1.2701	5.71 $\pm$ 0.2428	5.8 $\pm$ 0.024	5.72 $\pm$ 0.027	6.06 $\pm$ 0.041	6.23 $\pm$ 0.35(105)	5.93 $\pm$ 0.5(21)	6.07 $\pm$ 0.55(31)	5.91 $\pm$ 0.022	6.02 $\pm$ 0.0(6.5)	6.09 $\pm$ 0.0(12)	5.49 $\pm$ 0.02(31)	5.79 $\pm$ 0.0(125)
aminoacid	56.73 $\pm$ 0.018	56.57 $\pm$ 0.018	56.15 $\pm$ 0.83(8)	56.73 $\pm$ 0.018	56.73 $\pm$ 0.018	56.73 $\pm$ 0.018	56.73 $\pm$ 0.018	56.73 $\pm$ 0.018	56.73 $\pm$ 0.018	56.73 $\pm$ 0.018	56.73 $\pm$ 0.018	56.73 $\pm$ 0.018	56.73 $\pm$ 0.018	56.73 $\pm$ 0.018	56.73 $\pm$ 0.018	56.73 $\pm$ 0.018	56.73 $\pm$ 0.018
bird	99.17 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015	99.14 $\pm$ 0.015
campan	48.84 $\pm$ 0.011	48.84 $\pm$ 0.011	42.77 $\pm$ 0.2920	40.69 $\pm$ 0.024	67.11 $\pm$ 0.2522	23.09 $\pm$ 0.7183	24.99 $\pm$ 0.030	31.05 $\pm$ 0.001	40.51 $\pm$ 0.005	36.98 $\pm$ 1.2(25)	29.75 $\pm$ 8.87(9)	32.56 $\pm$ 4.69(28)	20.25 $\pm$ 0.81(30)	78.33 $\pm$ 0.0(4)	79.33 $\pm$ 0.0(6)	48.01 $\pm$ 0.0(4)	37.11 $\pm$ 0.0(28)
cardio	86.25 $\pm$ 0.011	86.17 $\pm$ 0.002	68.92 $\pm$ 1.77(21)	88.87 $\pm$ 0.024	67.11 $\pm$ 0.2522	23.09 $\pm$ 0.7183	24.99 $\pm$ 0.030	31.05 $\pm$ 0.001	40.51 $\pm$ 0.005	36.98 $\pm$ 1.2(25)	29.75 $\pm$ 8.87(9)	32.56 $\pm$ 4.69(28)	20.25 $\pm$ 0.81(30)	78.33 $\pm$ 0.0(4)	79.33 $\pm$ 0.0(6)	48.01 $\pm$ 0.0(4)	37.11 $\pm$ 0.0(28)
celeba	20.95 $\pm$ 1.1(11.5)	20.95 $\pm$ 1.0621	59.27 $\pm$ 0.03(12)	50.73 $\pm$ 0.024	54.94 $\pm$ 3.83(20)	67.52 $\pm$ 0.82(4)	33.51 $\pm$ 0.0024	16.48 $\pm$ 0.08(29)	68.98 $\pm$ 0.0(4)	45.78 $\pm$ 5.1(28)	60.56 $\pm$ 7.01(9)	9.04 $\pm$ 2.73(21)	7.65 $\pm$ 0.1(22)	58.36 $\pm$ 0.0(4)	59.77 $\pm$ 0.0(10)	47.05 $\pm$ 1.34(27)	45.75 $\pm$ 0.91(29)
census	19.82 $\pm$ 0.58(11)	20.03 $\pm$ 0.59(10)	14.68 $\pm$ 1.56(20)	14.01 $\pm$ 0.32(25)	17.49 $\pm$ 2.1(14)	14.66 $\pm$ 1.25(21)	11.73 $\pm$ 0.29(30.5)	11.73 $\pm$ 0.29(30.5)	11.73 $\pm$ 0.29(30.5)	15.35 $\pm$ 1.07(17)	13.42 $\pm$ 0.92(27)	13.17 $\pm$ 0.92(28)	14.25 $\pm$ 1.1(22)	20.92 $\pm$ 0.43(6)	20.62 $\pm$ 0.44(7)	14.31 $\pm$ 0.6(10)	10.05 $\pm$ 0.49(19)
cover	16.05 $\pm$ 0.88(21)	16.17 $\pm$ 0.86(20)	1.98 $\pm$ 0.63(1)	5.42 $\pm$ 0.61(27)	25.03 $\pm$ 3.24(16)	1.09 $\pm$ 0.18(32)	12.26 $\pm$ 0.08(23)	12.26 $\pm$ 0.08(23)	12.26 $\pm$ 0.08(23)	2.96 $\pm$ 1.53(29)	2.96 $\pm$ 1.53(29)	2.96 $\pm$ 1.53(29)	3.13 $\pm$ 0.85(15)	51.96 $\pm$ 2.9(10)	42.14 $\pm$ 2.2(11)	31.42 $\pm$ 0.97(14)	35.03 $\pm$ 0.82(12)
donors	36.01 $\pm$ 0.72(23)	35.22 $\pm$ 1.2(24)	49.31 $\pm$ 14.82(13)	36.33 $\pm$ 1.85(22)	23.91 $\pm$ 1.25(40)	9.01 $\pm$ 0.93(32)	37.26 $\pm$ 0.14(21)	33.5 $\pm$ 0.08(25)	41.27 $\pm$ 0.97(19)	42.75 $\pm$ 2.7(48)	42.75 $\pm$ 2.7(48)	42.75 $\pm$ 2.7(48)	30.27 $\pm$ 1.77(27)	87.78 $\pm$ 0.99(5)	79.92 $\pm$ 1.13(7)	76.75 $\pm$ 3.27(8)	45.33 $\pm$ 2.85(16)
donors	28.74 $\pm$ 5.54(26)	26.93 $\pm$ 1.52(28)	62.81 $\pm$ 9.37(3)	32.25 $\pm$ 5.32(22)	60.24 $\pm$ 1.15(51)	29.44 $\pm$ 2.66(25)	2.08 $\pm$ 0.83(31)	38.43 $\pm$ 3.99(18)	33.24 $\pm$ 4.68(21)	48.33 $\pm$ 1.71(32)	36.99 $\pm$ 1.51(31)	15.57 $\pm$ 0.13(30)	0.33 $\pm$ 0.03(32)	41.22 $\pm$ 5.02(15)	40.45 $\pm$ 0.07(6)	58.46 $\pm$ 7.58(7)	30.33 $\pm$ 4.25(23)
glass	18.51 $\pm$ 3.73(30)	20.96 $\pm$ 5.90(26)	30.93 $\pm$ 6.50(21)	27.61 $\pm$ 6.50(20)	26.04 $\pm$ 1.35(51)	67.02 $\pm$ 1.3(16)	20.09 $\pm$ 0.44(28)	23.02 $\pm$ 0.94(23)	25.02 $\pm$ 0.94(23)	35.35 $\pm$ 1.04(7)	15.55 $\pm$ 2.58(32)	18.62 $\pm$ 1.24(30)	23.14 $\pm$ 1.38(24)	31.91 $\pm$ 3.73(15)	31.91 $\pm$ 3.73(15)	41.38 $\pm$ 4.61(10)	27.78 $\pm$ 2.62(19)
hepatitis	64.78 $\pm$ 3.1(21)	64.85 $\pm$ 5.90(20)	80.63 $\pm$ 3.08(12)	63.49 $\pm$ 5.99(22)	73.25 $\pm$ 1.75(127)	65.77 $\pm$ 5.43(19)	89.16 $\pm$ 0.51(13)	56.08 $\pm$ 0.45(25)	45.84 $\pm$ 0.47(29)	98.73 $\pm$ 1.05(7)	50.15 $\pm$ 2.5(32)	54.37 $\pm$ 0.05(27)	34.91 $\pm$ 1.31(52)	89.28 $\pm$ 1.8(11)	70.62 $\pm$ 4.48(18)	99.28 $\pm$ 1.2(6)	74.38 $\pm$ 1.79(16)
hup	90.42 $\pm$ 1.67(12)	91.69 $\pm$ 1.53(10)	52.33 $\pm$ 4.21(22)	38.95 $\pm$ 1.2(626)	94.19 $\pm$ 0.16(29)	88.38 $\pm$ 0.0(18)	50.34 $\pm$ 1.1(24)	46.31 $\pm$ 1.1(24)	25.18 $\pm$ 0.82(28)	36.09 $\pm$ 7.9(31)	7.46 $\pm$ 9.79(31)	9.22 $\pm$ 0.31(13)	9.89 $\pm$ 0.02(2)	9.08 $\pm$ 0.0(1.55)	94.39 $\pm$ 1.31(8)	97.56 $\pm$ 2.2(4)	45.08 $\pm$ 7.07(25)
imb	8.71 $\pm$ 0.02(5)	8.71 $\pm$ 0.02(5)	9.53 $\pm$ 0.67(10)	9.01 $\pm$ 0.0(9)	94.19 $\pm$ 0.16(29)	88.38 $\pm$ 0.0(18)	50.34 $\pm$ 1.1(24)	46.31 $\pm$ 1.1(24)	25.18 $\pm$ 0.82(28)	36.09 $\pm$ 7.9(31)	7.46 $\pm$ 9.79(31)	9.22 $\pm$ 0.31(13)	9.89 $\pm$ 0.02(2)	9.08 $\pm$ 0.0(1.55)	94.39 $\pm$ 1.31(8)	97.56 $\pm$ 2.2(4)	45.08 $\pm$ 7.07(25)
interreads	46.72 $\pm$ 0.02(25)	46.97 $\pm$ 0.02(25)	47.57 $\pm$ 1.08(17)	47.57 $\pm$ 1.08(17)	52.89 $\pm$ 1.61(6)	47.43 $\pm$ 0.86(20)	30.56 $\pm$ 0.37(31)	67.74 $\pm$ 0.02(3)	67.74 $\pm$ 0.02(3)	9.68 $\pm$ 1.47(9)	8.73 $\pm$ 0.39(28)	9.22 $\pm$ 0.31(13)	9.89 $\pm$ 0.02(2)	9.08 $\pm$ 0.0(1.55)	94.39 $\pm$ 1.31(8)	97.56 $\pm$ 2.2(4)	45.08 $\pm$ 7.07(25)
labeled	40.79 $\pm$ 7.81(15)	40.79 $\pm$ 7.81(15)	34.19 $\pm$ 0.94(27)	60.12 $\pm$ 0.0(3)	37.14 $\pm$ 8.32(6)	31.21 $\pm$ 4.03(8)	37.37 $\pm$ 0.0(9)	33.82 $\pm$ 0.08(29)	31.09 $\pm$ 0.0(3)	49.43 $\pm$ 4.4(11)	35.74 $\pm$ 4.4(11)	40.28 $\pm$ 2.4(16)	37.55 $\pm$ 1.93(18)	50.81 $\pm$ 0.0(4)	46.49 $\pm$ 0.0(13)	55.86 $\pm$ 0.6(5)	34.98 $\pm$ 0.81(25)
letter	90.59 $\pm$ 0.0(32)	90.59 $\pm$ 0.0(32)	9.19 $\pm$ 0.81(13)	8.73 $\pm$ 0.0(3)	8.45 $\pm$ 0.09(23)	8.73 $\pm$ 0.0(3)	8.73 $\pm$ 0.0(3)	8.73 $\pm$ 0.0(3)	8.73 $\pm$ 0.0(3)	8.73 $\pm$ 0.0(3)	8.73 $\pm$ 0.0(3)	8.73 $\pm$ 0.0(3)	8.73 $\pm$ 0.0(3)	8.73 $\pm$ 0.0(3)	8.73 $\pm$ 0.0(3)	8.73 $\pm$ 0.0(3)	8.73 $\pm$ 0.0(3)
lymphography	98.99 $\pm$ 1.0(12)	98.99 $\pm$ 1.0(12)	96.24 $\pm$ 4.22(19)	96.55 $\pm$ 1.1(18)	90.53 $\pm$ 0.9(23)	98.73 $\pm$ 0.82(28)	98.1 $\pm$ 2.64(16)	93.93 $\pm$ 2.85(31)	94.38 $\pm$ 1.43(21)	96.82 $\pm$ 0.61(7)	24.13 $\pm$ 1.29(32)	73.47 $\pm$ 1.67(8)	30.87 $\pm$ 5.63(31)	96.86 $\pm$ 0.43(20)	86.74 $\pm$ 0.0(21)	98.57 $\pm$ 0.0(25)	68.82 $\pm$ 5.96(30)
magnetic_gamma	75.57 $\pm$ 0.02(5)	75.57 $\pm$ 0.02(5)	78.51 $\pm$ 2.91(18)	77.15 $\pm$ 0.02(6)	65.83 $\pm$ 2.36(40)	76.13 $\pm$ 2.42(23)	65.76 $\pm$ 0.03(31)	72.22 $\pm$ 0.00(2)	67.92 $\pm$ 0.00(2)	69.54 $\pm$ 0.73(28)	75.78 $\pm$ 1.08(24)	64.54 $\pm$ 0.62(32)	83.19 $\pm$ 0.66(12)	85.26 $\pm$ 0.0(10)	84.53 $\pm$ 0.0(12)	80.74 $\pm$ 1.03(14)	77.29 $\pm$ 3.52(20)
mnist	41.76 $\pm$ 0.02(6)	41.65 $\pm$ 0.00(7)	18.52 $\pm$ 9.52(29)	21.32 $\pm$ 0.02(6)	37.08 $\pm$ 2.89(15)	27.82 $\pm$ 3.84(22)	11.77 $\pm$ 0.01(31)	54.63 $\pm$ 0.00(2)	5.2 $\pm$ 0.00(1)	27.54 $\pm$ 1.44(28)	43.21 $\pm$ 1.04(2)	27.24 $\pm$ 2.32(24)	27.24 $\pm$ 2.32(24)	41.08 $\pm$ 0.0(10)	39.83 $\pm$ 0.0(12)	28.51 $\pm$ 2.57(21)	28.91 $\pm$ 1.92(20)
mnist	64.99 $\pm$ 0.01(25)	64.99 $\pm$ 0.01(25)	52.22 $\pm$ 3.32(1)	22.21 $\pm$ 0.00(29)	47.97 $\pm$ 4.78(23)	65.09 $\pm$ 0.57(11)	20.19 $\pm$ 0.00(31)	16.86 $\pm$ 0.00(31)	16.86 $\pm$ 0.00(31)	46.09 $\pm$ 0.35(25)	34.07 $\pm$ 0.67(27)	46.06 $\pm$ 0.77(24)	59.72 $\pm$ 1.28(25)	72.08 $\pm$ 0.0(3)	70.36 $\pm$ 0.0(5)	56.56 $\pm$ 0.65(17)	41.15 $\pm$ 1.11(26)
mnist	10.00 $\pm$ 0.08(5)	10.00 $\pm$ 0.08(5)	32.68 $\pm$ 3.48(31)	10.00 $\pm$ 0.08(5)	10.00 $\pm$ 0.08(5)	10.00 $\pm$ 0.08(5)	10.00 $\pm$ 0.08(5)	10.00 $\pm$ 0.08(5)	10.00 $\pm$ 0.08(5)	99.91 $\pm$ 0.17(17)	90.81 $\pm$ 0.84(25)	70.61 $\pm$ 0.94(27)	15.62 $\pm$ 1.19(61)	10.00 $\pm$ 0.08(5)	10.00 $\pm$ 0.08(5)	89.48 $\pm$ 1.94(24)	62.46 $\pm$ 0.02(29)
mnist	99.91 $\pm$ 0.02(1)	99.91 $\pm$ 0.02(1)	38.26 $\pm$ 5.02(21)	22.8 $\pm$ 0.02(3)	46.98 $\pm$ 1.72(29)	63.54 $\pm$ 1.71(3)	59.14 $\pm$ 0.02(3)	41.51 $\pm$ 0.03(3)	59.54 $\pm$ 0.03(3)	57.06 $\pm$ 0.43(23)	48.57 $\pm$ 3.88(28)	69.26 $\pm$ 1.23(41)	33.46 $\pm$ 1.43(28)	61.75 $\pm$ 0.0(1)	67.15 $\pm$ 0.0(1)	64.91 $\pm$ 1.31(6)	54.53 $\pm$ 1.0(26)
mnist	39.14 $\pm$ 0.0(9)	38.63 $\pm$ 0.0(21)	34.47 $\pm$ 8.88(29)	42.33 $\pm$ 0.0(2)	14.65 $\pm$ 1.5(2027)	33.35 $\pm$ 1.72(2)	22.36 $\pm$ 0.02(2)	30.86 $\pm$ 0.0(2)	30.86 $\pm$ 0.0(2)	93.44 $\pm$ 7.85(2)	37.23 $\pm$ 9.42(1)	11.71 $\pm$ 9.83(1)	14.57 $\pm$ 4.32(8)	82.34 $\pm$ 0.0(6)	82.34 $\pm$ 0.0(6)	58.53 $\pm$ 0.70(13)	29.79 $\pm$ 0.96(25)
mnist	71.49 $\pm$ 3.69(13)	71.18 $\pm$ 3.39(15.5)	71.23 $\pm$ 2.99(14)	75.88 $\pm$ 2.56(6)	61.06 $\pm$ 4.66(27)	65.15 $\pm$ 8.82(4)	56.78 $\pm$ 3.69(30)	69.07 $\pm$ 2.47(19)	64.77 $\pm$ 2.3(25)	59.75 $\pm$ 7.5(28)	59.36 $\pm$ 7.6(29)	52.98 $\pm$ 3.33(31)	53.42 $\pm$ 1.35(32)	74.01 $\pm$ 2.75(9)	74.01 $\pm$ 2.75(9)	75.31 $\pm$ 2.15(8)	67.61 $\pm$ 2.81(23)
mnist	81.04 $\pm$ 0.1(19)	77.79 $\pm$ 0.00(25)	77.86 $\pm$ 2.47(24)	86.49 $\pm$ 0.05(5)	81.83 $\pm$ 0.75(16)	78.96 $\pm$ 0.46(23)	63.25 $\pm$ 0.03(2)	73.33 $\pm$ 0.00(31)	69.57 $\pm$ 0.0(31)	81.14 $\pm$ 1.97(18)	79.77 $\pm$ 0.93(22)	75.98 $\pm$ 3.34(28)	77.46 $\pm$ 3.34(28)	83.82 $\pm$ 0.0(4)	83.82 $\pm$ 0.0(4)	85.84 $\pm$ 2.16(8)	71.75 $\pm$ 2.86(30)
mnist	92.94 $\pm$ 0.28(13)	91.92 $\pm$ 0.01(5)	62.47 $\pm$ 5.18(29)	87.68 $\pm$ 0.02(0)	80.25 $\pm$ 1.59(23)	95.89 $\pm$ 0.1(8)	90.87 $\pm$ 0.0(21)	85.27 $\pm$ 0.00(31)	79.66 $\pm$ 0.0(24)	76.28 $\pm$ 1.2(26)	93.74 $\pm$ 0.69(12)	47.48 $\pm$ 0.34(30)	79.32 $\pm$ 1.38(25)	96.92 $\pm$ 0.03(5)	97.22 $\pm$ 0.00(2)	71.56 $\pm$ 7.98(27)	43.96 $\pm$ 5.34(31)
mnist	96.27 $\pm$ 0.01(9)	96.27 $\pm$ 0.01(9)	51.66 $\pm$ 1.93(30)	97.49 $\pm$ 0.01(6)	51.66 $\pm$ 1.93(30)	97.49 $\pm$ 0.01(6)	90.87 $\pm$ 0.0(21)	90.87 $\pm$ 0.0(21)	90.87 $\pm$ 0.0(21)	90.87 $\pm$ 0.0(21)	90.87 $\pm$ 0.0(21)	90.87 $\pm$ 0.0(21)	90.87 $\pm$ 0.0(21)	90.87 $\pm$ 0.0(21)	90.87 $\pm$ 0.0(21)	90.87 $\pm$ 0.0(21)	90.87 $\pm$ 0.0(21)
mnist	40.14 $\pm$ 0.32(7)	36.39 $\pm$ 0.33(28)	74.74 $\pm$ 7.46(6)	33.57 $\pm$ 0.59(30)	31.88 $\pm$ 2.33(40)	42.18 $\pm$ 1.84(26)	63.01 $\pm$ 1.72(15)	25.69 $\pm$ 0.18(32)	25.69 $\pm$ 0.18(32)	99.91 $\pm$ 0.17(17)	90.81 $\$						





Table 14.2: 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.

Dataset	VAE	PCA	PlanetFlow	HBOs	GANomaly	GOAD	DIF	COPOD	RECOD	DeepSVDD	LODA	DAGMM	DROCC	DTE-NP <sup>ss</sup>	KNN <sup>ss</sup>	ICL <sup>ss</sup>	DTE-C <sup>ss</sup>
doi	7.63 $\pm$ 0.0055	7.43 $\pm$ 0.007	3.83 $\pm$ 0.7229	7.43 $\pm$ 0.007	9.35 $\pm$ 0.2271	5.73 $\pm$ 0.4818	5.73 $\pm$ 0.4818	4.58 $\pm$ 0.0024	4.44 $\pm$ 0.0025	5.17 $\pm$ 0.9220	6.61 $\pm$ 0.5811	5.98 $\pm$ 1.6714	6.01 $\pm$ 0.0032	5.98 $\pm$ 0.0016	5.98 $\pm$ 0.0016	4.59 $\pm$ 0.0723	3.56 $\pm$ 0.0330
amazon	50.19 $\pm$ 0.020	50.00 $\pm$ 0.021	8.01 $\pm$ 0.767	35.96 $\pm$ 0.009	34.27 $\pm$ 0.4800	57.57 $\pm$ 0.6213	57.57 $\pm$ 0.6213	31.65 $\pm$ 0.031	38.99 $\pm$ 0.028	23.33 $\pm$ 0.1232	46.78 $\pm$ 0.9426	46.78 $\pm$ 0.9426	46.78 $\pm$ 0.9426	59.18 $\pm$ 0.088	59.18 $\pm$ 0.088	59.18 $\pm$ 0.088	50.03 $\pm$ 0.023
audiology	8.51 $\pm$ 0.1022	8.51 $\pm$ 0.1022	36.73 $\pm$ 2.221	36.73 $\pm$ 2.221	29.93 $\pm$ 0.6125	47.99 $\pm$ 0.4727	47.99 $\pm$ 0.4727	96.41 $\pm$ 0.031	96.41 $\pm$ 0.031	82.96 $\pm$ 1.3851	82.96 $\pm$ 1.3851	82.96 $\pm$ 1.3851	82.96 $\pm$ 1.3851	42.86 $\pm$ 1.3212	42.86 $\pm$ 1.3212	42.86 $\pm$ 1.3212	96.76 $\pm$ 1.002
breast	48.84 $\pm$ 0.014	48.84 $\pm$ 0.014	49.09 $\pm$ 1.622	49.09 $\pm$ 1.622	49.09 $\pm$ 1.622	96.63 $\pm$ 0.1245	96.63 $\pm$ 0.1245	96.63 $\pm$ 0.1245	96.63 $\pm$ 0.1245	91.85 $\pm$ 0.7124	91.85 $\pm$ 0.7124	91.85 $\pm$ 0.7124	91.85 $\pm$ 0.7124	96.22 $\pm$ 0.5532	96.22 $\pm$ 0.5532	96.22 $\pm$ 0.5532	92.51 $\pm$ 0.7923
cardio	76.14 $\pm$ 0.0015	76.14 $\pm$ 0.0015	59.77 $\pm$ 1.9420	56.25 $\pm$ 0.024	58.75 $\pm$ 0.3722	74.89 $\pm$ 0.0513	74.89 $\pm$ 0.0513	49.27 $\pm$ 0.0031	48.88 $\pm$ 0.015	37.81 $\pm$ 0.1929	30.74 $\pm$ 0.4729	34.13 $\pm$ 0.4327	40.00 $\pm$ 0.032	50.93 $\pm$ 0.005	51.1 $\pm$ 0.002	50.93 $\pm$ 0.005	31.83 $\pm$ 0.4428
cardiography	61.59 $\pm$ 0.0025	61.59 $\pm$ 0.0025	49.43 $\pm$ 0.5111	41.42 $\pm$ 0.022	46.35 $\pm$ 0.6025	99.96 $\pm$ 1.084	99.96 $\pm$ 1.084	48.28 $\pm$ 0.0155	48.28 $\pm$ 0.0155	37.08 $\pm$ 0.5626	55.11 $\pm$ 0.718	52.32 $\pm$ 0.9999	33.95 $\pm$ 0.1239	62.95 $\pm$ 0.0145	65.0 $\pm$ 0.011	62.95 $\pm$ 0.0145	34.91 $\pm$ 0.0931
celeba	27.04 $\pm$ 0.413	27.04 $\pm$ 0.413	17.91 $\pm$ 0.4912	22.68 $\pm$ 0.039	11.07 $\pm$ 0.9923	3.98 $\pm$ 0.9830	3.98 $\pm$ 0.9830	22.78 $\pm$ 0.018	22.78 $\pm$ 0.018	8.43 $\pm$ 0.527	13.26 $\pm$ 0.3921	14.19 $\pm$ 0.5519	8.62 $\pm$ 0.8426	17.87 $\pm$ 0.5713	19.11 $\pm$ 0.6105	17.87 $\pm$ 0.5713	48.76 $\pm$ 0.3901
census	20.76 $\pm$ 0.2710	20.82 $\pm$ 0.3391	13.86 $\pm$ 2.422	10.77 $\pm$ 0.624	18.27 $\pm$ 3.9141	4.96 $\pm$ 2.1829	4.96 $\pm$ 2.1829	0.04 $\pm$ 0.0315	0.04 $\pm$ 0.0315	19.28 $\pm$ 1.4413	14.15 $\pm$ 0.8121	14.46 $\pm$ 2.6919	15.55 $\pm$ 1.4408	21.62 $\pm$ 0.146	21.55 $\pm$ 0.247	21.62 $\pm$ 0.146	10.31 $\pm$ 0.5226
cover	16.21 $\pm$ 1.6822	16.24 $\pm$ 1.5312	2.59 $\pm$ 2.4830	10.75 $\pm$ 1.2626	25.69 $\pm$ 3.3716	0.00 $\pm$ 0.029	0.00 $\pm$ 0.029	18.82 $\pm$ 0.8120	18.82 $\pm$ 0.8120	3.43 $\pm$ 0.4429	24.09 $\pm$ 1.2619	12.16 $\pm$ 1.5624	41.87 $\pm$ 1.5624	61.65 $\pm$ 2.1406	52.19 $\pm$ 1.9211	61.65 $\pm$ 2.1406	38.91 $\pm$ 0.6214
donors	53.24 $\pm$ 0.0825	53.24 $\pm$ 0.0825	57.65 $\pm$ 1.4859	53.24 $\pm$ 0.0825	53.24 $\pm$ 0.0825	95.96 $\pm$ 0.1835	95.96 $\pm$ 0.1835	95.96 $\pm$ 0.1835	95.96 $\pm$ 0.1835	51.17 $\pm$ 0.8940	51.17 $\pm$ 0.8940	51.17 $\pm$ 0.8940	51.17 $\pm$ 0.8940	56.23 $\pm$ 0.6106	56.23 $\pm$ 0.6106	56.23 $\pm$ 0.6106	59.31 $\pm$ 0.1106
fraud	34.46 $\pm$ 3.2225	33.76 $\pm$ 1.7027	66.59 $\pm$ 1.809	41.51 $\pm$ 4.7823	61.17 $\pm$ 1.0771	37.28 $\pm$ 2.8524	37.28 $\pm$ 2.8524	4.99 $\pm$ 3.7831	4.99 $\pm$ 3.7831	58.11 $\pm$ 14.779	45.06 $\pm$ 11.5919	50.28 $\pm$ 22.3230	0.04 $\pm$ 0.0032	47.78 $\pm$ 3.5314	45.61 $\pm$ 3.8817	47.78 $\pm$ 3.5314	30.84 $\pm$ 6.6628
glass	18.03 $\pm$ 3.1124	15.77 $\pm$ 8.6627	19.56 $\pm$ 11.222	27.68 $\pm$ 10.3914	20.15 $\pm$ 11.0031	60.41 $\pm$ 4.5555	60.41 $\pm$ 4.5555	15.77 $\pm$ 8.6627	15.77 $\pm$ 8.6627	45.39 $\pm$ 21.8777	14.66 $\pm$ 4.7131	13.72 $\pm$ 16.1321	15.48 $\pm$ 13.1429	29.42 $\pm$ 2.6111	21.09 $\pm$ 5.1619	29.42 $\pm$ 2.6111	24.15 $\pm$ 12.7106
higgs	60.51 $\pm$ 6.921	60.56 $\pm$ 7.5920	79.75 $\pm$ 3.3713	58.08 $\pm$ 1.2422	65.51 $\pm$ 10.6619	57.86 $\pm$ 7.7723	57.86 $\pm$ 7.7723	37.63 $\pm$ 0.8831	37.63 $\pm$ 0.8831	93.82 $\pm$ 1.2977	46.76 $\pm$ 9.928	47.51 $\pm$ 7.5427	29.29 $\pm$ 17.2832	84.54 $\pm$ 2.2310	66.15 $\pm$ 4.7108	84.54 $\pm$ 2.2310	70.87 $\pm$ 3.415
hup	91.94 $\pm$ 1.2511	92.71 $\pm$ 1.4101	14.43 $\pm$ 12.3125	36.4 $\pm$ 3.6827	18.99 $\pm$ 4.1324	56.39 $\pm$ 1.7108	56.39 $\pm$ 1.7108	2.05 $\pm$ 1.2929	2.05 $\pm$ 1.2929	25.02 $\pm$ 22.6021	1.05 $\pm$ 0.9630	48.95 $\pm$ 3.3719	0.04 $\pm$ 0.0315	88.73 $\pm$ 2.5414	95.32 $\pm$ 0.518	88.73 $\pm$ 2.5414	77.24 $\pm$ 12.8822
indb	4.64 $\pm$ 0.0025	4.64 $\pm$ 0.0025	5.53 $\pm$ 1.4583	4.64 $\pm$ 0.0025	4.64 $\pm$ 0.0025	5.53 $\pm$ 1.4583	5.53 $\pm$ 1.4583	6.04 $\pm$ 0.0031	6.04 $\pm$ 0.0031	9.99 $\pm$ 0.1707	7.04 $\pm$ 1.0631	7.04 $\pm$ 1.0631	7.04 $\pm$ 1.0631	5.24 $\pm$ 0.0028	5.24 $\pm$ 0.0028	5.24 $\pm$ 0.0028	4.47 $\pm$ 0.0025
isip	4.54 $\pm$ 0.0025	4.54 $\pm$ 0.0025	5.53 $\pm$ 1.4583	4.54 $\pm$ 0.0025	4.54 $\pm$ 0.0025	5.53 $\pm$ 1.4583	5.53 $\pm$ 1.4583	6.04 $\pm$ 0.0031	6.04 $\pm$ 0.0031	9.99 $\pm$ 0.1707	7.04 $\pm$ 1.0631	7.04 $\pm$ 1.0631	7.04 $\pm$ 1.0631	5.24 $\pm$ 0.0028	5.24 $\pm$ 0.0028	5.24 $\pm$ 0.0028	4.47 $\pm$ 0.0025
isoprene	79.82 $\pm$ 6.240	78.99 $\pm$ 2.8726	90.83 $\pm$ 1.6106	69.49 $\pm$ 1.8629	86.17 $\pm$ 2.520	83.42 $\pm$ 5.9823	83.42 $\pm$ 5.9823	69.49 $\pm$ 1.8629	69.49 $\pm$ 1.8629	91.07 $\pm$ 1.2533	77.34 $\pm$ 4.5427	69.33 $\pm$ 4.2330	69.49 $\pm$ 1.8629	91.21 $\pm$ 1.2409	87.86 $\pm$ 1.8617	91.21 $\pm$ 1.2409	79.23 $\pm$ 4.1425
landsat	38.73 $\pm$ 4.5920	38.73 $\pm$ 4.5920	31.82 $\pm$ 3.9135	61.0 $\pm$ 0.08	1.21 $\pm$ 0.8422	12.10 $\pm$ 0.4522	12.10 $\pm$ 0.4522	40.00 $\pm$ 0.15	40.00 $\pm$ 0.15	5.04 $\pm$ 1.5801	2.40 $\pm$ 0.5422	8.64 $\pm$ 3.965	13.6 $\pm$ 4.8622	1.04 $\pm$ 0.0725	1.04 $\pm$ 0.0725	1.04 $\pm$ 0.0725	5.84 $\pm$ 1.1109
lymphography	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025
magnetic	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025	9.02 $\pm$ 0.0025
gamma	65.53 $\pm$ 0.024	65.53 $\pm$ 0.024	91.07 $\pm$ 9.3216	88.49 $\pm$ 3.8820	83.31 $\pm$ 1.4325	91.13 $\pm$ 0.4522	91.13 $\pm$ 0.4522	88.49 $\pm$ 3.8820	88.49 $\pm$ 3.8820	5.04 $\pm$ 1.5801	2.40 $\pm$ 0.5422	8.64 $\pm$ 3.965	13.6 $\pm$ 4.8622	1.04 $\pm$ 0.0725	1.04 $\pm$ 0.0725	1.04 $\pm$ 0.0725	5.84 $\pm$ 1.1109
mnist	45.0 $\pm$ 0.06	44.62 $\pm$ 0.007	22.01 $\pm$ 0.8828	16.92 $\pm$ 0.30	36.85 $\pm$ 2.7818	35.62 $\pm$ 5.719	35.62 $\pm$ 5.719	62.86 $\pm$ 0.026	59.72 $\pm$ 0.030	59.88 $\pm$ 0.299	67.82 $\pm$ 1.4827	67.82 $\pm$ 1.4827	67.82 $\pm$ 1.4827	75.61 $\pm$ 0.011	74.48 $\pm$ 0.012	75.61 $\pm$ 0.011	85.36 $\pm$ 1.691
mnist	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	35.05 $\pm$ 2.7311	24.14 $\pm$ 0.029	41.83 $\pm$ 0.025	63.91 $\pm$ 1.3111	63.91 $\pm$ 1.3111	0.04 $\pm$ 0.0315	0.04 $\pm$ 0.0315	99.18 $\pm$ 1.3117	90.72 $\pm$ 5.4121	70.22 $\pm$ 1.4432	12.16 $\pm$ 1.7732	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	60.04 $\pm$ 0.028
mnist	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	35.05 $\pm$ 2.7311	24.14 $\pm$ 0.029	41.83 $\pm$ 0.025	63.91 $\pm$ 1.3111	63.91 $\pm$ 1.3111	0.04 $\pm$ 0.0315	0.04 $\pm$ 0.0315	99.18 $\pm$ 1.3117	90.72 $\pm$ 5.4121	70.22 $\pm$ 1.4432	12.16 $\pm$ 1.7732	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	60.04 $\pm$ 0.028
mnist	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	35.05 $\pm$ 2.7311	24.14 $\pm$ 0.029	41.83 $\pm$ 0.025	63.91 $\pm$ 1.3111	63.91 $\pm$ 1.3111	0.04 $\pm$ 0.0315	0.04 $\pm$ 0.0315	99.18 $\pm$ 1.3117	90.72 $\pm$ 5.4121	70.22 $\pm$ 1.4432	12.16 $\pm$ 1.7732	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	60.04 $\pm$ 0.028
mnist	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	35.05 $\pm$ 2.7311	24.14 $\pm$ 0.029	41.83 $\pm$ 0.025	63.91 $\pm$ 1.3111	63.91 $\pm$ 1.3111	0.04 $\pm$ 0.0315	0.04 $\pm$ 0.0315	99.18 $\pm$ 1.3117	90.72 $\pm$ 5.4121	70.22 $\pm$ 1.4432	12.16 $\pm$ 1.7732	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	60.04 $\pm$ 0.028
mnist	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	35.05 $\pm$ 2.7311	24.14 $\pm$ 0.029	41.83 $\pm$ 0.025	63.91 $\pm$ 1.3111	63.91 $\pm$ 1.3111	0.04 $\pm$ 0.0315	0.04 $\pm$ 0.0315	99.18 $\pm$ 1.3117	90.72 $\pm$ 5.4121	70.22 $\pm$ 1.4432	12.16 $\pm$ 1.7732	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	60.04 $\pm$ 0.028
mnist	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	35.05 $\pm$ 2.7311	24.14 $\pm$ 0.029	41.83 $\pm$ 0.025	63.91 $\pm$ 1.3111	63.91 $\pm$ 1.3111	0.04 $\pm$ 0.0315	0.04 $\pm$ 0.0315	99.18 $\pm$ 1.3117	90.72 $\pm$ 5.4121	70.22 $\pm$ 1.4432	12.16 $\pm$ 1.7732	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	60.04 $\pm$ 0.028
mnist	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	35.05 $\pm$ 2.7311	24.14 $\pm$ 0.029	41.83 $\pm$ 0.025	63.91 $\pm$ 1.3111	63.91 $\pm$ 1.3111	0.04 $\pm$ 0.0315	0.04 $\pm$ 0.0315	99.18 $\pm$ 1.3117	90.72 $\pm$ 5.4121	70.22 $\pm$ 1.4432	12.16 $\pm$ 1.7732	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	60.04 $\pm$ 0.028
mnist	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	35.05 $\pm$ 2.7311	24.14 $\pm$ 0.029	41.83 $\pm$ 0.025	63.91 $\pm$ 1.3111	63.91 $\pm$ 1.3111	0.04 $\pm$ 0.0315	0.04 $\pm$ 0.0315	99.18 $\pm$ 1.3117	90.72 $\pm$ 5.4121	70.22 $\pm$ 1.4432	12.16 $\pm$ 1.7732	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	60.04 $\pm$ 0.028
mnist	100.0 $\pm$ 0.0085	100.0 $\pm$ 0.0085	35.05 $\pm$ 2.7311	24.14 $\pm$ 0.029	41.83 $\pm$ 0.025	63.91 $\pm$ 1.											

## H BENCHMARK OD DATASETS

Table 15: Description of all datasets in ADBench Livernoche et al. (2024).

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
smtpt	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



## I DIFFERENCES TO PRIOR WORK ON PFNS FOR TABULAR DATA

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

Here we highlight the differences of our proposed FoMo-OD from these existing PFNs.

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

## J RELATED WORK

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

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

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

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

**Prior-data Fitted Networks:** Based on the seminal work by Müller et al. (2022), Prior-data-fitted Networks (PFNs) establish a new paradigm for machine learning, where a PFN is pretrained on

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

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

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

Foundation models, however, do not exist for tabular data which is widespread across OD applications in the real world, such as detecting credit card fraud, network intrusion, medical anomalies, and any sensor measurement abnormalities, to 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 rather is meta-trained on each specific domain using inlier-only datasets from the *same domain*. Concurrent to our work, Li et al. (2024) apply pretrained LLMs for prompt-based OD on tabular data which they serialize to text. Similar to our work, they also use *simulated* labeled OD datasets to fine-tune several existing LLMs to improve their performance. Their work, 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-OD employs and pretrains PFNs at a much larger scale with rigorous evaluation on a much larger testbed.

## K DISCUSSION

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

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

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

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

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

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

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

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

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

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

## L BROADER IMPACT STATEMENT

FoMo-OD is zero-shot, abolishing not only parameter training but also model selection given a new dataset. This is a radical paradigm shift for OD literature, which historically focused on designing new models and recently also effective ways for unsupervised model selection. Obviating the need for either, we expect FoMo-OD to route attention of the community from new OD model design and selection to designing better data priors and gathering datasets for PFN pretraining, along with better and more scalable architectures for PFN.

From the applied perspective, a zero-shot OD model like FoMo-OD is a game-changer for practitioners! Given the plethora of OD algorithms to choose from, which often come with a list of hyperparameters to set, and not having the tools for effective and efficient model selection, the practitioners are burdened with a “choice paralysis”. With FoMo-OD, practitioners can not only bypass such dilemmas on one dataset, but thanks to the “train once, use many times” nature of pretrained models, they can do so for any dataset such as those arriving over time. In fact, provided its lightening-fast inference

via a single forward pass, FoMo-OD is amenable to deploy in real time on streaming datasets, such that each (test) sample over a stream can be inferred with the preceding samples passed as context.

## M REPRODUCIBILITY STATEMENT

We expect that the disruptive nature of FoMo-OD will trigger future innovations in the OD literature, as well as a widespread adoption by practitioners thanks to its key desirable properties. To foster future research and accessibility in practice, we make all resources (our codebase used for prior data synthesis, data transformation, and pretraining as well as our pretrained model checkpoints) publicly available at <https://anonymous.4open.science/r/PFN40D>. Further, full implementation details are provided in Appendix C.