

# ZERO-SHOT OUTLIER DETECTION VIA PRIOR-DATA FITTED NETWORKS: MODEL SELECTION BYGONE!

Yuchen Shen Haomin Wen Leman Akoglu

Carnegie Mellon University

Pittsburgh, PA 15213, USA

{yuchens3, haominwe, lakoglu}@andrew.cmu.edu

## ABSTRACT

Outlier detection (OD) has a vast literature as it finds numerous applications in environmental monitoring, cybersecurity, finance, and medicine 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 is 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, tuning its associated hyperparameters, and even training any 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 top *2nd* 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 many practical applications in different domains such as security, environmental monitoring, manufacturing, finance, and so on. This popularity brings about a large literature that offers a plethora of detection techniques to choose from given a new OD task. Further, these techniques exhibit several hyperparameters (HPs) that need careful tuning to which they are often quite sensitive Ma et al. (2023). What makes it notoriously difficult for achieving effective OD performance in practical applications is *model selection* (both algorithm and HPs) in the *absence of any labels*, provided most tasks are unsupervised.<sup>1</sup>

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

Most recent advances in ML have been through large foundation models, which are (pre-)trained on massive amounts of data. The most notable progress has been in natural language and image domains, thanks to the admirable quantity and quality of public text and image datasets. On the other hand, 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 point-cloud OD has been

<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>7</sup>** baselines with performance avg.’ed over varying HPs (denoted w/  $^{avg}$ ); followed by  $p$ -values of the pairwise Wilcoxon signed rank test, comparing FoMo-0D to each baseline (from top to bottom) over **All** (57) datasets, those (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	-	<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
$d \leq 100$	-	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>

the non-shared feature spaces of different tabular datasets, unlike e.g. the shared pixel or word spaces for images and text, respectively.

Recently, the introduction of Prior-data Fitted Networks (PFNs) has marked a significant milestone as a new approach to ML on tabular data Müller et al. (2022). PFNs are based on Bayesian non-parametrics 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 observed training data as input context. To approximate the PPD, a Transformer model Vaswani et al. (2017) is pre-trained offline to mimic the PPD via simulating numerous training datasets from a (general and complex) data prior. For inference, the fresh input training set along with the test instances 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 the original PFN has been shown to match the performance of tree-based models on small classification datasets Hollmann et al. (2023) as well as 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 Q-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 is capable of zero-shot OD on a new input 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. This is a game changer (!) for unsupervised OD and in our opinion, a **revolution for the OD literature**. 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 in total to infer a test sample since a new dataset requires a single forward pass for inference and no training time.

Interestingly, PFNs and thus also FoMo-0D exhibit *sample-to-sample attention* by design, as the (training) data of a task is input as context for the test samples. This mimics non-parametric models, by using parametric attention mechanisms. This property is particularly intriguing from the lens of OD: because many OD methods are based on  $k$  nearest neighbors (NNs)<sup>2</sup>, however (*i*) they are

<sup>2</sup>In fact, one of the simplest OD measures is the distance to the  $k$ -NNs (a classical shallow approach), while SOTA deep models such as DTE Livernoche et al. (2024) mimic  $k$ NN.

sensitive to the HP value of  $k$ , (ii) the relations between samples are **not** learnable (unlike attention weights), and (iii) they are limited to nearest neighbors (whereas multi-head attention learns *which* training points are worth attending to, with the capacity to go beyond just the near(est) neighbors).

**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-0D 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 *less than 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-0D in *real time* for streaming OD.
- **Effectiveness:** We evaluate FoMo-0D 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-0D significantly outperforms the majority of the baselines while performing statistically no different from the top 2<sup>nd</sup> baseline, at the fraction of the compute cost.

As FoMo-0D proposes a paradigm shift in the field of OD, abolishing model training and selection altogether, while also delivering unreasonable effectiveness on benchmark datasets even with a straightforward data prior, we expect FoMo-0D will trigger further work in both the research and practical use of OD. To this end, we make all of our codebase, including synthetic data generation, model training, as well as our pretrained FoMo-0D checkpoint, 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.

In this paper, 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 inlier points  $\mathbf{x}_i \in \mathbb{R}^d$ , where  $y_i = 0 \forall i \in [n]$ , and  $\mathcal{D}_{\text{test}}$  depicts the test dataset 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 data  $\mathcal{D}_{\text{in}}$ .

## 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.

*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-0D: 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-0D, a prior-data fitted Foundation Model for 0-shot Outlier Detection. FoMo-0D 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-0D 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 unreasonable effectiveness on many benchmark datasets compared to classical and modern baselines even at its first-attempt stage, we expect FoMo-0D will become a milestone in future research and practice of OD. The new FoMo-0D paradigm (right) versus the typical OD setting (left) is illustrated in Figure 1.

In this section, we present details on our data prior for OD, training of FoMo-0D on prior-simulated datasets and inference on new datasets, and our specific model architecture and improvements for scalable training.

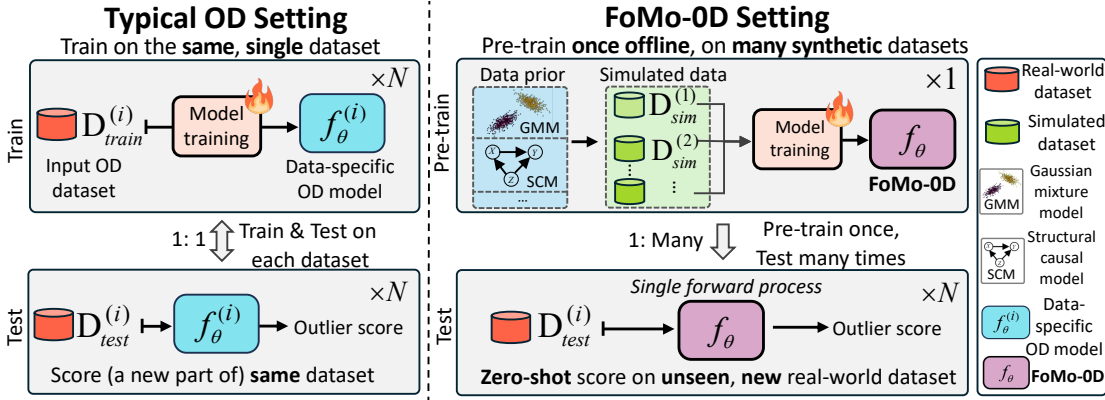


Figure 1: (best in color) Comparison of typical OD vs. the FoMo-OD settings. Given a new un/semi-supervised OD dataset, FoMo-OD not only eliminates the need for model training, but most importantly, also abolishes the onerous task of model selection (algorithm and hyperparameters) w/out labels.

### 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 first (and final) attempt has been sufficient to achieve astonishing performance even with a relatively straightforward and simple-to-implement data prior, which we describe next. (Informally: we didn’t even try hard!)

**Inlier synthesis:** We designate the data prior for inliers to simply be 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 *diagonal*<sup>3</sup>  $\Sigma_j$  with entries also in  $[-5, 5]$ . For each step of the every epoch of pretraining FoMo-OD, 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  $S$  inlier points, defined as instances within the 90th percentile 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 90th percentile of the original GMM. We decide whether a drawn sample is within/outside the specified percentile based on its Mahalanobis distance computed analytically (see Property A.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 2D synthetic datasets are illustrated in Appendix B.

**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 simple. While it has been our intent to extend our preliminary attempt

<sup>3</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.

toward designing a sophisticated data prior for OD, we found to our big surprise that even with such an elementary prior, FoMo-OD performs remarkable 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. As with previous literature using PFNs, future work can extend our data prior by also employing BNNs and SCMs as discussed in Section 2. Moreover, discrete features (e.g. from multinomials) as well as more complex outlier types such as contextual and dependency outliers Steinbuss and Böhm (2021) can be simulated 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 new OD 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 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 input 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

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

As such, the sample-to-sample attention presents an intriguing intuition from the perspective of OD: Many classical shallow algorithms for OD 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 not need for specifying a cut-off HP like  $k$ , to which most previous OD techniques are extremely 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, therefore, we scale up FoMo-OD’s attention via the “router mechanism” recently developed by Zhang and Yan (2023). As shown in Figure 2, the main idea is to learn a small fixed number ( $R \ll n$ ) of “routers” or representatives, which gather information from all  $n$  samples and then distribute the gathered 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-OD training to **scale linearly** with respect to both dataset 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:

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

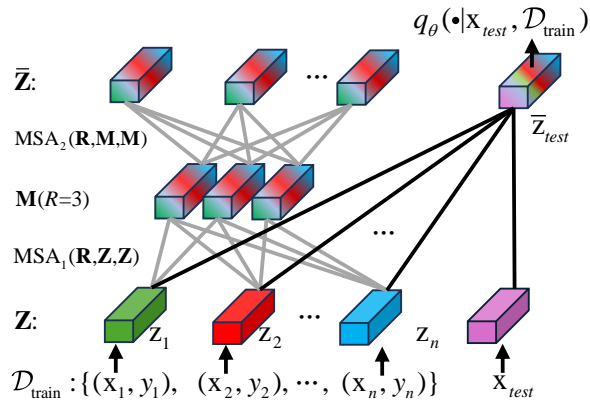


Figure 2: FoMo-OD architecture with the “router mechanism” for scalable attention.

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-OD 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>4</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 12). 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, albeit a longer pretraining time.

**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 A.1).<sup>5</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. All details and proofs are given in Appendix A.

**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 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 unique datasets from different GMMs (with varying configurations) in a typical manner, and then we 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

**Pre-training Dataset Synthesis:** During pretraining, we generate unique GMM datasets by first drawing a configuration, which includes feature dimensionality  $d \in [D]$ , number of components/clusters  $m \in [M]$ , cluster centers  $\{\boldsymbol{\mu}_j\}_{j=1}^m$  (each  $\boldsymbol{\mu}_j \in [-5, 5]^d$ ) and covariances  $\{\boldsymbol{\Sigma}_j\}_{j=1}^m$  ( $\text{diag}(\boldsymbol{\Sigma}_j) \in [-5, 5]^d$ ). We set  $M = 5$  and vary  $D \in \{20, 100\}$  to study pretraining with relatively small and high dimensional datasets, respectively. We 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.  $5 \times \boldsymbol{\Sigma}_{j,kk}$ ’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  $(\mathbf{W}, \mathbf{b})$

<sup>4</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>5</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.



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.

**Real-world Benchmark Datasets:** While pretraining is purely on synthetically generated datasets, we evaluate FoMo-0D on **57** real-world datasets downstream. The datasets are from the ADBench benchmark Han et al. (2022) (see their Table B1), which comprises 47 popular tabular outlier detection datasets, as well as 10 newly-constructed tabular datasets created from images and natural language tasks by using pretrained models to extract embeddings. We defer to the original paper for the details on these benchmark datasets.

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

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

**Baselines:** We compare FoMo-0D against **26** baselines, from classical/shallow methods to modern/deep models. Our baselines include all the baselines imported from one of the latest papers that proposed the SOTA diffusion-based model DTE Livernoche et al. (2024), along with the three variants of DTE, namely, 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 Akcay et al. (2019). We defer to the original paper for additional details on the baselines. Overall, our 26 baselines consist of the most recent, SOTA approaches for OD that span a diverse family (nonparametric, self-supervised, generative, etc.).

**Model Implementation:** We trained our final model for 200,000 steps with a batch size of 8 datasets. That is, our FoMo-0D 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-0D, including the training algorithm, model architecture, data synthesis and reuse, and hardware are provided in Appendix C.

**Metrics and Hypothesis Testing:** Detection performance of methods are evaluated with respect to (w.r.t.) three widely-used metrics for OD: AUROC; area under the ROC curve, AUPR; area under the Precision-Recall curve, and F1 score; harmonic mean of Precision and Recall at the threshold of the true number of outliers in the test data (varies by dataset).

To compare methods, we compute their **rank** on each dataset (when ranked by a given metric, lower is better), and present **average rank** across datasets. This is an alternative to the average metric (e.g. AUROC) value, which is not meaningful when task difficulties and hence metric values vary widely. In addition, we perform significance tests to compare two methods statistically. Specifically, we use the one-sided pairwise Wilcoxon signed rank test Demšar (2006) between FoMo-0D 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:** Importantly, Livernoche et al. (2024) picked for each baseline method the best-performing set of hyperparameters (HPs) as recommended by the authors in their original paper. As for their own DTE, which behaves similar to  $k$ NN based approaches, they recommend  $k = 5$  and select the *same*  $k$  for the  $k$ NN baseline Ramaswamy et al. (2000) to be consistent. However, it is

<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  $\mathbf{b} \in [-1, 1]^d$ .

well known that the  $k$ NN baseline is sensitive to the choice of  $k$  across datasets 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>7</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.1.

Overall, we compare FoMo-0D 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-0D 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 F). We find that among 30 baselines and 2 variants of FoMo-0D (w/  $D = 100$  and  $D = 20$ ), FoMo-0D 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-0D with author-recommended  $k = 5$ , we find that DTE-NP<sup>avg</sup> is on par with FoMo-0D.

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-0D to outperform.

FoMo-0D w/  $D = 100$  performs statistically no different from **all** baselines on datasets with  $d \leq 100$  (i.e., “on its own game” when pretraining data dimensions align with real-world datasets), while it *outperforms the majority of the baselines* (with  $p > 0.95$ ). These results continue to hold on datasets with  $d \leq 500$ .

Table 2 shows similar results for FoMo-0D 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-0D (with  $D = 20$ ), **top-10** baselines with default HPs, and **top-4**<sup>7</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/  $d \leq 20$  and (38) datasets w/  $d \leq 50$  dimensions, respectively. Even with small  $D = 20$ , FoMo-0D 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-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)	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>

<sup>7</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 14, and rank the baselines w.r.t. their mean  $p$ -value.

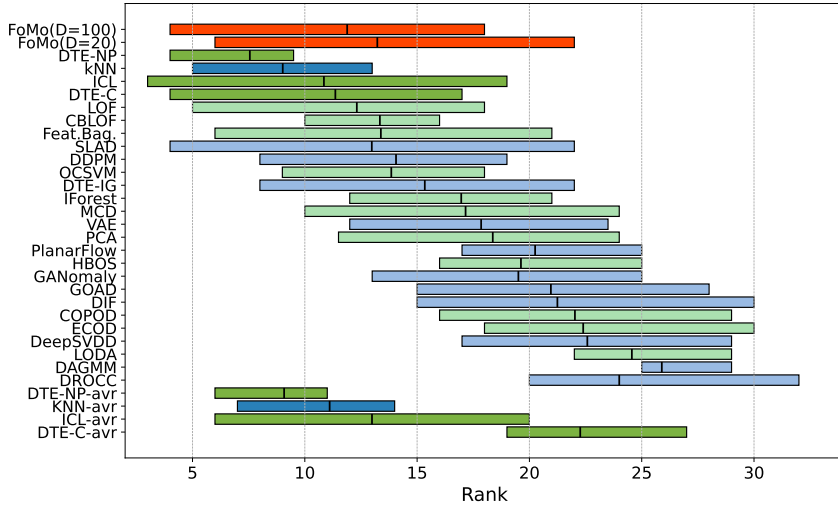


Figure 3: (best in color) Rank (w.r.t. AUROC performance, lower is better) distribution across datasets shown via boxplots for (from top to bottom) FoMo-0D in red, all 26 baselines as ordered by mean  $p$ -value<sup>7</sup> (shallow and deep baselines resp. in green and blue), and top-4 baselines’ avg variants.

Figure 3 shows the distribution of ranks across datasets for each of the 32 methods. While paired significant tests are the most conclusive, FoMo-0D 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-0D and the top-3 baselines. Given a new dataset, FoMo-0D bypasses model training (and HP tuning) and directly performs inference, with an average of 7.7 ms per sample (see Appendix Figure 11). 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-0D and the top-3<sup>7</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-0D	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

In this section, we perform various ablations to study the effect of different design choices in FoMo-0D; namely, **A1.** maximum pretraining data dimensionality  $D$ , **A2.** the number of routers  $R$  on **a.** cost and **b.** performance, **A3.** context size (both for training and inference), **A4.** number of unique datasets used for pretraining (i.e., reuse periodicity  $P$ ), **A5.** data transformation  $T$  during synthesis on **a.** performance and **b.** speed up, **A6.** data diversity and prolonged training, and finally, **A7.** 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.

**A1. Effect of  $D$ :** 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 10. 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.

**A2.a. Effect of routers on cost:** *What is the running time and memory cost of FoMo-0D with & w/out router-based attention?*

Figure 4(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>8</sup>

Similarly, memory cost with routers is considerably lower when using routers, especially for larger context sizes, as shown in Figure 4(middle).

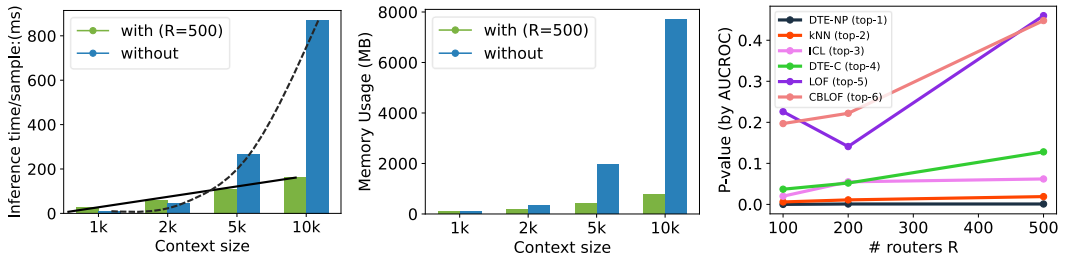


Figure 4: FoMo-0D w/ router mechanism saves time and memory while more #routers perform better, offering a cost-performance trade-off: (left) inference-time (ms) per sample and (middle) memory cost (MB) with & w/out routers by varying context size; (right) performance (based on  $p$ -value against top baselines, higher is better) vs. number of routers. (setting:  $D = 20, P = 1$ )

**A2.b. 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 4(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.

**A3. 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 4 shows the results, where performance is depicted by the average rank of FoMo-0D (the lower, the better).

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

<sup>8</sup>Note that the inference time is reported on CPUs to show scalability. On GPUs, w/ 5K context size, see Appendix Figure 11, 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).

Table 4: 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

statistical tests among the  $3 \times 4 = 12$  models support these observations, as shown in Figure 5. 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.

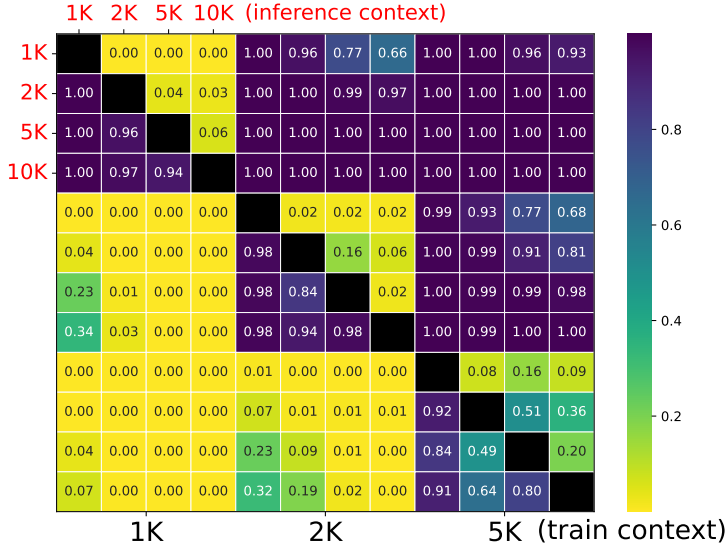


Figure 5:  $p$ -values of the pairwise Wilcoxon signed rank test between models (larger  $p$  implies col-method is better than row-method) w/ different context sizes for **training** (1K/2K/5K, 1st/2nd/3rd four grids, in **black**) and **inference** (1K/2K/5K/10K, every 1st/2nd/3rd/4th grid, in **red**):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

Table 5: 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	<u>0.001</u>	<u>0.019</u>	0.062	0.128	0.460	<u>0.001</u>	<u>0.019</u>	0.089	0.159	0.394	<u>0.001</u>	<u>0.015</u>	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>

#### A4. Effect 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 5, 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  $d \leq 50$ .

**A5.a. Effect of  $T$  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 6 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 6: 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,  $P = 1$ )

top-5	w/out transformation $T$					w/ transformation $T$				
	DTE-NP	kNN	ICL	DTE-C	LOF	DTE-NP	kNN	ICL	DTE-C	LOF
All	<u>0.001</u>	<u>0.019</u>	0.062	0.128	0.460	<u>0.002</u>	<u>0.015</u>	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>

**A5.b. Speed up by  $T$ :** *What is the time saving on data synthesis with linear transformation?*

Figure 6 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 11).

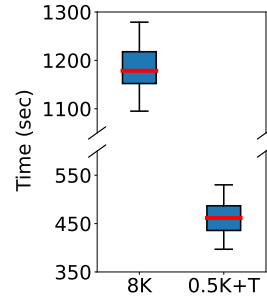


Figure 6: Runtime/epoch dist.n over 100 epochs for FoMo-0D ( $D=100$ ) with (left)  $P=100$ , i.e. 8K unique datasets/epoch vs. (right) 0.5K unique+7.5K transformed datasets/epoch.

**A6. Effect of data diversity and prolonged training:** *How does FoMo-0D’s performance change by increasing pretraining data diversity and number of training epochs?*

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

Toward reducing the loss further, we resume the pretraining for another 100 epochs. Further, to simplify the tasks and thereby increase data diversity, we also decrease the inlier/outlier labeling percentile threshold from 90% to 80% during on-the-fly data generation in the last 100 epochs. We provide the prolonged training loss in Appendix Figure 13.

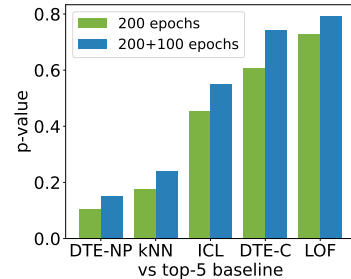


Figure 7:  $p$ -values increase with additional 100 epochs of pretraining, i.e. FoMo-0D w/  $D = 100$  performs better against top-5 baselines on datasets w/  $d \leq 100$ .

Figure 7 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 showcase the benefit of additional training.

**A7. Effect of applying quantile transform on benchmark datasets:** *What is the impact of quantile data transform preceding inference on performance?*

We pretrain FoMo-OD 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-OD for inference, which transforms the features to exhibit a more Gaussian-like probability distribution.

Figure 8 compares the performance of three FoMo-OD 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-OD tends to perform better as suggested by larger  $p$ -values when QT is applied.

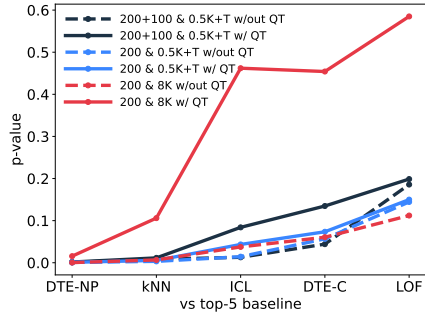


Figure 8:  $p$ -values increase, i.e. FoMo-OD performance improves, against top-5 baselines with quantile transform (QT) preceding inference, for 3 different settings of FoMo-OD w/  $D = 100$ .

5 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-0D 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 H 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 Esmaeilpour 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-0D employs and pretrains PFNs at a much larger scale with rigorous evaluation on a much larger testbed.

## 6 CONCLUSION AND DISCUSSION

**Summary:** We introduced FoMo-0D, **the first foundation model for outlier detection** (OD) on tabular data. FoMo-0D 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-0D 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-0D 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-0D 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-0D takes a mere **7.7 ms** per test sample for inference only.

**Discussion, Limitations and Future Directions:** FoMo-0D 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-0D 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.

**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 lightning-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.

**Reproducibility:** 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>.

## ACKNOWLEDGEMENT

We thank Frank Hutter for planting the seeds of our ideas during a discussion at the SIGKDD 2023 conference, as well as for their spearheading of the literature on PFNs.

## 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, Siddartha 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 Esmailpour, 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 Schirrmeyer, 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 Liversoche, 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 LINEAR TRANSFORM FOR SCALABLE GMM DATA SYNTHESIS

### A.1 DEFINITIONS

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

$$p(\mathbf{x}) = \sum_{j=1}^m w_j \cdot g(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j), \quad (6)$$

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

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

**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.

### A.2 PROPERTIES

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

*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 A.1, which completes the proof.

### A.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:

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

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

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

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

■

**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  $X \sim \chi_d^2$ :

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

Based on Property A.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 Theorem 1), then we conclude that the linear transform  $T$  with invertible  $\mathbf{W}$  preserves the percentiles of the GMM. ■

## B ILLUSTRATION OF SYNTHETIC DATA IN 2D

We visualize our synthetic data in Figure 9, with 3 randomly created 2d-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.

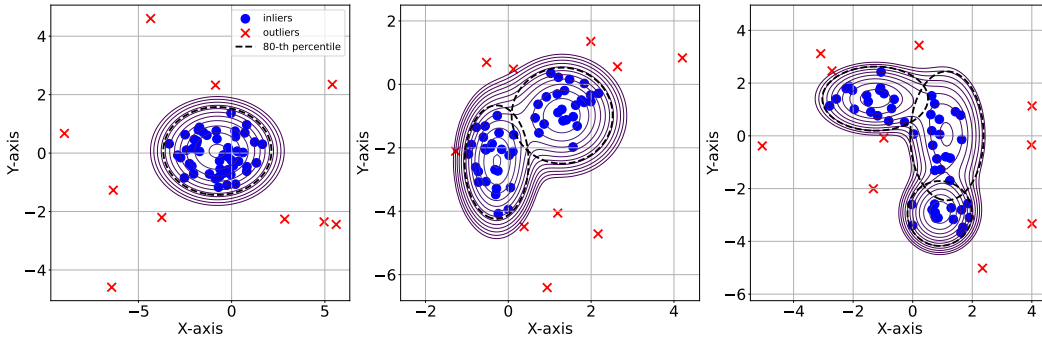


Figure 9: Illustration of synthetic data in 2D with 80-th percentile as the criterion.

## 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 `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

---

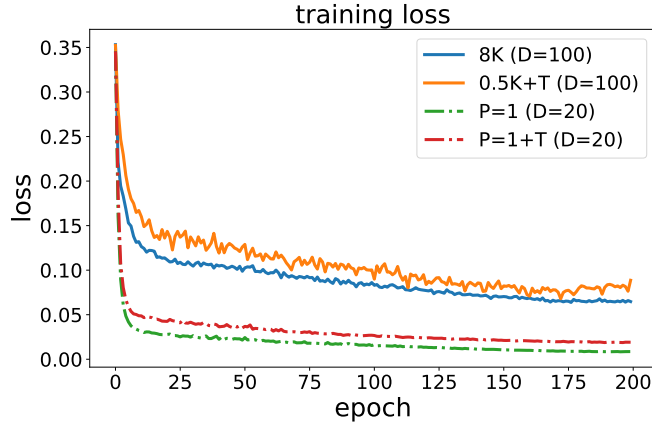


Figure 10: (best in color) Training loss of FoMo-0D ( $D = 100$ ) with 8K unique datasets/epoch (in blue) and using 0.5K unique + 7.5K transformed datasets/epoch (in orange), and FoMo-0D ( $D = 20$ ) with  $P = 1$  (in green) and  $P = 1$  with transformation (in red) over 200 epochs.

**Training loss** In Figure 10, 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 4 (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 11 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.

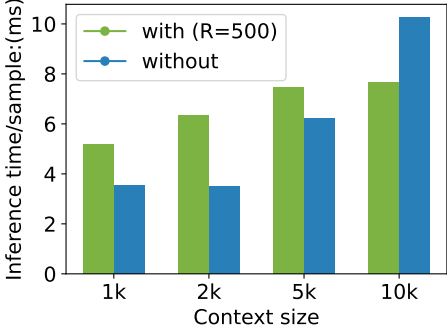


Figure 11: Inference time of FoMo-0D on GPU with vs. w/out router-based attention under varying context size.

## D DETAILS ON EXPERIMENT SETUP

### D.1 HYPERPARAMETERS FOR BASELINES

Table 7 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 7: 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.2 RANKING THE 26 BASELINES

Figure 14 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.3 COMPARISON OF TOP-4 BASELINE VARIANTS WITH VARYING HP CONFIGURATIONS

Figure 15, 16, 17, 18 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.4 SAMPLING TIME OF $d$ -DIMENSIONAL GMM

Figure 12 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 non-linearly 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 transformations.

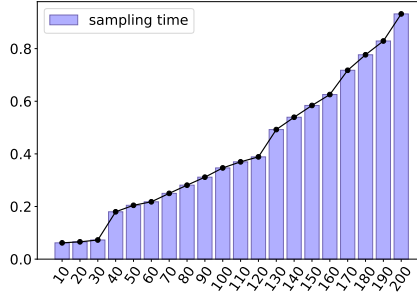


Figure 12: Sampling time (in seconds) of 10,000 points from GMMs with varying #dimensions.

### E ADDITIONAL DETAILS ON ABLATION STUDY

#### E.1 LOSS FOR PROLONGED TRAINING

In Figure 13, we present the training loss of FoMo-0D ( $D = 100$ ) trained with 0.5K unique + 7.5K transformed datasets/epoch over 200 epochs (90% percentile as labeling threshold) and then 100 additional epochs (80% percentile as the threshold) to show how data diversity and amount affect model performance.

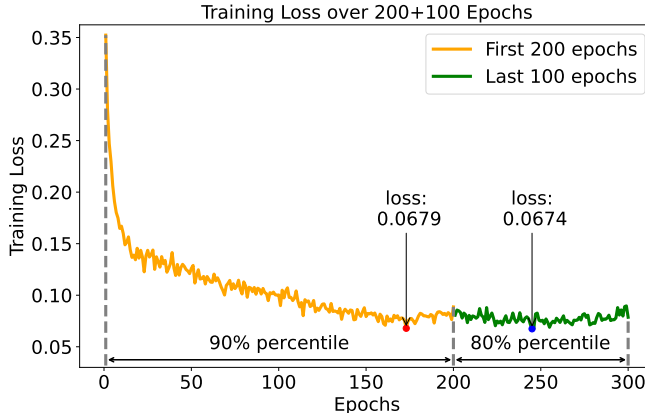


Figure 13: (best in color) Training loss of FoMo-0D ( $D = 100$ ) with 0.5K unique + 7.5K transformed datasets/epoch for 200 epochs (in orange), followed with additional 100 epochs of training (in green). For the first 200 epochs we train with 90% percentile as the inlier/outlier threshold, which we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

### F 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  $^{avg}$  model with the average performance across HPs, as listed in Table 7.

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  $^{avg}$  variants of top-4 baselines + FoMo-0D w/  $D = 20$ ).

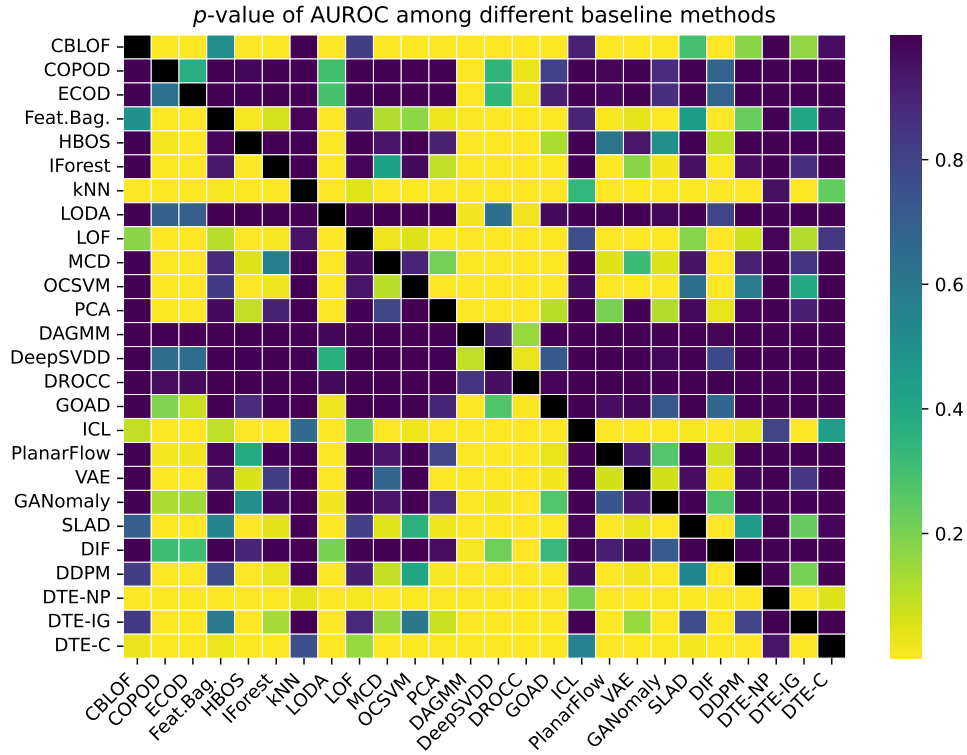


Figure 14: Pairwise  $p$ -values among baseline methods based on the Wilcoxon signed rank test w.r.t. AUROC performances across datasets.

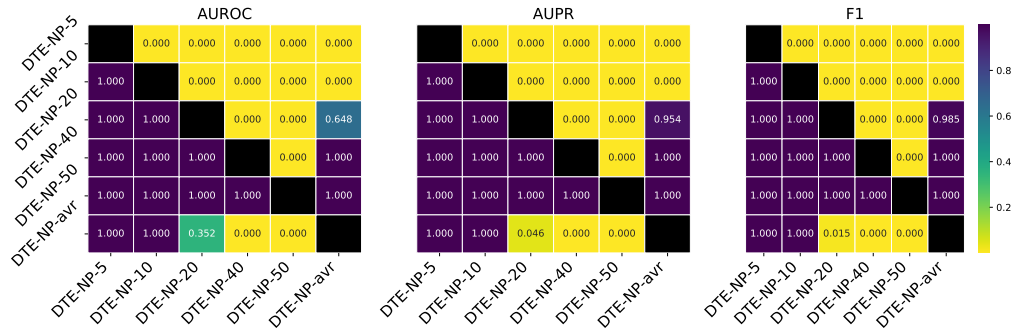


Figure 15:  $p$ -values w.r.t. AUROC/AUPR/F1 among different HP configurations of **DTE-NP** (i.e.,  $k \in \{5, 10, 20, 40, 50\}$ ), along with the  $^{avg}$  model with the average performance across HPs.

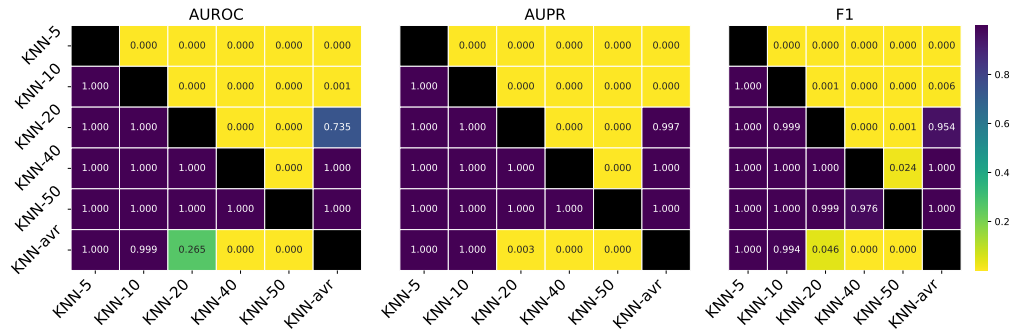


Figure 16:  $p$ -values w.r.t. AUROC/AUPR/F1 among different HP configurations of  $k$ NN (i.e.,  $k \in \{5, 10, 20, 40, 50\}$ ), along with the  $^{avg}$  model with the average performance across HPs.

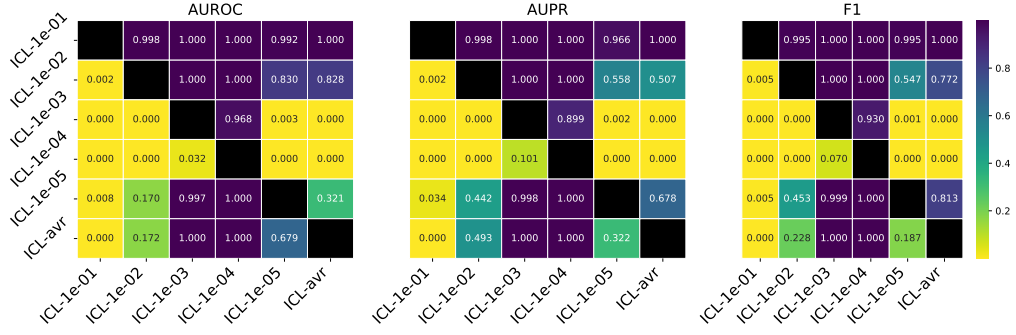


Figure 17:  $p$ -values w.r.t. AUROC/AUPR/F1 among different HP configurations of **ICL** (i.e.,  $\text{learning\_rate} \in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$ ), along with the  $\text{avg}$  model with the average performance across HPs.

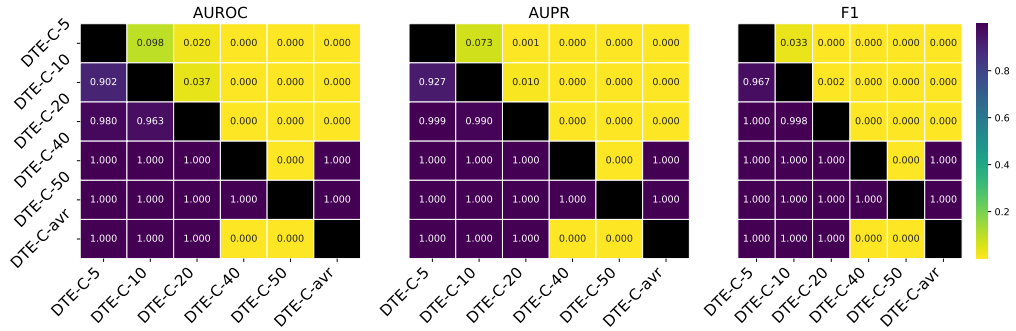


Figure 18:  $p$ -values w.r.t. AUROC/AUPR/F1 among different HP configurations of **DTE-C** (i.e.,  $k \in \{5, 10, 20, 40, 50\}$ ), along with the  $\text{avg}$  model with the average performance across HPs.

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** baselines with performance **avg.**’ed over varying HPs (denoted w/  $\text{avg}$ ); followed by  $p$ -values of the pairwise Wilcoxon signed rank test, comparing FoMo-0D to each baseline (from top to bottom) over **All** (57) datasets, those (24) w/  $d \leq 20$ , (38) w/  $d \leq 50$ , (42) w/  $d \leq 100$  and (46) datasets w/  $d \leq 500$  dimensions. FoMo-0D performs as well as (i.e., **statistically no different from**) the **2nd best model** ( $k$ 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 $\text{avg}$	$k$ NN $\text{avg}$	ICL $\text{avg}$	DTE-C $\text{avg}$
Rank(avg)	11.886	7.553	9.018	10.851	11.36	12.316	13.342	13.386	12.982	14.061	13.851	9.079	11.105	12.991	22.263
All	-	<b>0.016</b>	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















Table 12.1: Average AUROC ± standard dev. over five seeds for the semi-supervised setting on ADBench. Rank of each model among 32 models (26 baselines + 4 avg variants of top-4 baselines + 2 FoMo-0D variants w/ D = 100 and D = 20) per dataset is provided (in parentheses) (the lower, the better). We use blue and green respectively to mark the top-1 and the top-2 method. Last four rows show avg\_rank of methods across datasets, and p-values of the Wilcoxon signed rank test comparing FoMo-0D (D = 100) with other baselines. The previous four rows are the same for FoMo-0D (D = 20), when ranking 31 models (26 baselines + 4 avg variants of top-4 baselines + FoMo-0D w/ D = 20).

Table with columns: Dataset, FoMo(D=100), FoMo(D=20), DTE-NP, RNN, ICL, DTE-C, LOF, CBLOF, FeatureBugging, SLAD, DDPM, OC-SVM, DTE-IG, HForest, MCD. Rows include datasets like aloi, amazon, amthropy, backdoor, breastw, campaign, cardio, cardiology, celeba, census, cover, donors, fruit, fraud, glass, hepatitis, htp, imdb, internet, isosphere, landat, letter, lymphography, magnetic, namemag, mnist, musk, oodigit, paghebooks, pendigits, pima, satellite, skin, slant2, spam, speech, stamps, thyroid, vertebral, vowels, waveform, vbc, vulbe, wilt, wine, wpc, yeast, yelp, FashionMINST, CIFAR10, SVHN, MVTe-AD, 20news, agnews, and Benchmark. Each cell contains a value and a rank in parentheses.

Table 12.2: Average AUROC ± standard dev. over five seeds for the semi-supervised setting on ADBench. Rank of each model per dataset is provided (in parentheses) (the lower, the better). We use blue and green respectively to mark the top-1 and the top-2 method.

Table with columns: Dataset, PCA, PlainFlow, HBOS, GANomaly, GOAD, DJF, COPOD, ECOD, DeepSVDD, LODA, DAGMM, DRDCC, DTE-NP%, KNN%, ICL%, DTE-C%.

Table 13.1: Average AUPR  $\pm$  standard dev. over five seeds for the semi-supervised setting on ADBench. Rank of each model per dataset is provided (in parentheses) (the lower, the better). We use blue and green respectively to mark the top-1 and the top-2 method.

Table with 15 columns: Dataset, FoMo(D=100), FoMo(D=20), DTE-NP, KNN, ICL, DTE-C, LOF, CBLof, FeatureBagging, SLAD, DDPM, OC5VM, DTE-IG, IPovert, MCD. Rows include datasets like aloi, amazon, amblyoid, backdoor, breast, campaign, card, cardiology, cat5, census, cover, donors, fruit, fruit, glass, hepatitis, http, imdb, internets, inosphere, landsat, letter, lymphography, magic\_gamma, mammography, mist, must, outdigs, pageblinks, pendigits, pima, satellite, satellite, satellite-2, silhouette, skin, snp, spanish, speech, stamps, thyroid, vowels, waveform, wave, webc, wine, yeast, yeast, yield, FashionMNIST, CIFAR10, SVHN, MVTe-AD, 20news, agnews, Rank(erg), All, d <= 20, Rank(erg), All, d <= 100, Rank(erg), All, d <= 500.

Table 13.2: Average AUPR ± standard dev. over five seeds for the semi-supervised setting on ADBench. Rank of each model per dataset is provided (in parentheses) (the lower, the better). We use blue and green respectively to mark the top-1 and the top-2 method.

Table with 15 columns: Dataset, VAE, PCA, PlainFlow, HBGS, GNNomaly, GQAD, DIF, COPOD, ECOD, DeqSVDD, LODA, DAGMM, DRCC2, DTE-NP, KNN, ICL, DTE-C. Rows list various datasets like aloi, amazon, amynoid, backdoor, bias, campaign, cardio, celeba, etc.



Table 14.2: Average F1 score ± standard dev. over five seeds for the semi-supervised setting on ADBench. Rank of each model per dataset is provided (in parentheses) (the lower, the better). We use blue and green respectively to mark the top-1 and the top-2 method.

Table with columns: Dataset, VAE, PCA, PhantFlow, HRBS, GANomaly, GOAD, DIF, COPOD, RECORD, DeepSVDD, LODA, DNGMM, DROCC, DTE-NPvs, KNNvs, ICLvs, DTE-Cvs. Rows list various datasets like aoi, aneurysm, ant, etc., and their corresponding model performance metrics.



## G 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
smtp	95156	3	30	0.03	Web
SpamBase	4207	57	1679	39.91	Document
speech	3686	400	61	1.65	Linguistics
Stamps	340	9	31	9.12	Document
thyroid	3772	6	93	2.47	Healthcare
vertebral	240	6	30	12.50	Biology
vowels	1456	12	50	3.43	Linguistics
Waveform	3443	21	100	2.90	Physics
WBC	223	9	10	4.48	Healthcare
WDBC	367	30	10	2.72	Healthcare
Wilt	4819	5	257	5.33	Botany
wine	129	13	10	7.75	Chemistry
WPBC	198	33	47	23.74	Healthcare
yeast	1484	8	507	34.16	Biology
CIFAR10	5263	512	263	5.00	Image
FashionMNIST	6315	512	315	5.00	Image
MNIST-C	10000	512	500	5.00	Image
MVTec-AD	5354	512	1258	23.50	Image
SVHN	5208	512	260	5.00	Image
Agnews	10000	768	500	5.00	NLP
Amazon	10000	768	500	5.00	NLP
Imdb	10000	768	500	5.00	NLP
Yelp	10000	768	500	5.00	NLP
20newsgroups	11905	768	591	4.96	NLP

## H 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-0D, 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-0D 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-0D 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.

## I REPRODUCIBILITY STATEMENT

We open-source all of our codebase used for prior data synthesis, data transformation, pretraining, as well as our final model checkpoint at <https://anonymous.4open.science/r/PFN40D>. Further, full implementation details are provided in Appendix C.