# ON DIFFUSION MODELING FOR ANOMALY DETECTION

**Victor Livernoche** [1 3 *]    **Vineet Jain**[1 3 *]    **Yashar Hezaveh**[2 3]    **Siamak Ravanbakhsh**[1 3]

[1] School of Computer Science, McGill University
[2] Department of Physics, University of Montreal
[3] Mila - Quebec AI Institute
[*] Equal contribution

## ABSTRACT

Known for their impressive performance in generative modeling, diffusion models are attractive candidates for density-based anomaly detection. This paper investigates different variations of diffusion modeling for unsupervised and semi-supervised anomaly detection. In particular, we find that Denoising Diffusion Probability Models (DDPM) are performant on anomaly detection benchmarks yet computationally expensive. By simplifying DDPM in application to anomaly detection, we are naturally led to an alternative approach called Diffusion Time Estimation (DTE).[1] DTE estimates the distribution over diffusion time for a given input and uses the mode or mean of this distribution as the anomaly score. We derive an analytical form for this density and leverage a deep neural network to improve inference efficiency. Through empirical evaluations on the ADBench benchmark, we demonstrate that all diffusion-based anomaly detection methods perform competitively for both semi-supervised and unsupervised settings. Notably, DTE achieves orders of magnitude faster inference time than DDPM, while outperforming it on this benchmark. These results establish diffusion-based anomaly detection as a scalable alternative to traditional methods and recent deep-learning techniques for standard unsupervised and semi-supervised anomaly detection settings.

## 1 INTRODUCTION

Anomaly detection seeks to identify observations that differ from the others to such a large extent that they are likely generated by a different mechanism (Hawkins, 1980). This is a longstanding research problem in machine learning with applications in various fields ranging from medicine (Pachauri & Sharma, 2015; Salem et al., 2013), finance (Ahmed et al., 2016b), security (Ahmed et al., 2016a), manufacturing (Susto et al., 2017), particle physics (Fraser et al., 2022) and geospatial data (Yairi et al., 2006). Despite its significance and potential for impact (e.g., leading to the discovery of new phenomena), to this day traditional anomaly detection methods, such as nearest neighbours, reportedly outperform deep learning techniques on various benchmarks (Han et al., 2022) by a significant margin. This is true for unsupervised, semi-supervised, and supervised anomaly detection tasks. However, the growing number of applications involving high-dimensional data and massive datasets are beginning to challenge the classical, and in particular non-parametric, techniques, and there is a need for scalable, interpretable, and expressive deep learning techniques for anomaly detection.

In recent years, denoising diffusion probabilistic models (DDPMs) (Ho et al., 2020) have received much attention as a powerful class of generative models. While these models have been successfully utilized for anomaly detection in domain-specific image datasets (Wolleb et al., 2022; Zhang et al., 2023a; Wyatt et al., 2022), a comprehensive exploration of their applicability for general-purpose anomaly detection across diverse tabular, image, and natural language datasets is notably absent.

Our starting point is the observation that DDPM exhibits competitive performance compared to previous approaches for unsupervised and semi-supervised anomaly detection. These are some of

---

[1] Code available at https://github.com/vicliv/DTE

the most challenging settings, where either an unlabelled mix of normal and anomalous samples are available for training or, at best, the training data only includes normal samples. However, the expressivity and interpretability of DDPM come with a considerable computational cost. This computational complexity poses challenges for anomaly detection tasks involving large datasets or data streams.

In anomaly detection using DDPM, we deterministically "denoise" the input and measure the distance to its denoised reconstruction; a large distance indicates an anomaly. Since we only use this distance for outlier identification, in order to reduce the complexity of the diffusion-based approach, we propose to directly estimate this distance, which is correlated with diffusion time.

More precisely, we estimate the posterior distribution of diffusion time (or noise variance) for a given input. This estimated distribution serves as a guide for identifying anomalies, as they are anticipated to exhibit higher posterior density at larger time steps compared to normal samples. In particular, we use the mode or mean of this distribution as the anomaly score. We derive an analytical form for this posterior distribution, enabling its non-parametric estimation. We see that the non-parametric approximation produces a ranking for anomalies that is identical to k-Nearest Neighbours (kNN) for anomaly detection. We then propose a parametric model, a deep neural network, allowing us to leverage the generalization capability and efficient inference time of deep learning.
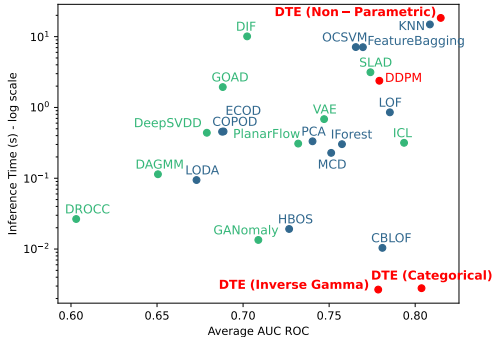


Figure 1: Average inference time vs. average AUC ROC for all 57 ADBench datasets in the semi-supervised setting. Lower right is better (DTE Categorical). Colour scheme: red (diffusion-based), green (deep learning), blue (classical).

We provide an extensive evaluation compared to classical and other deep models for different anomaly detection settings on more than 57 datasets from ADBench (Han et al., 2022). Our empirical results suggest that using a single deep neural network architecture across all datasets and settings makes the diffusion model competitive with classical and other deep models. Figure 1 shows the efficiency and effectiveness of different anomaly detection algorithms across all datasets in ADBench. Notably, our proposed method surpasses the direct application of DDPMs, achieving substantial improvements in inference time.

The contributions of our work are summarized as follows:

- Evaluation of denoising diffusion probabilistic models on various anomaly detection tasks encompassing tabular data and embeddings of images and natural language datasets.
- Development of a simplified approach that models the posterior distribution over diffusion time as a proxy for anomaly detection.
- Derivation of an analytical form of the posterior distribution of diffusion time and development of a non-parametric estimator that leads us to kNN.
- Introduction of a parametric approach utilizing a deep neural network for improved generalization and scalability.
- Implementation of additional baselines and extensive evaluation on 57 datasets from ADBench, showcasing competitive performance compared to classical and existing deep-learning-based anomaly detection algorithms.
- Investigation into the interpretability of diffusion-based methods, including our novel approach, highlighting their strengths and limitations.
- Exploration of optimal representation selections for image datasets with diffusion methods.

## 2 PRELIMINARIES

A classification of anomaly detection methods is based on the availability of labelled data. *Supervised* setting is similar to binary classification with unbalanced classes since the number of anomalies

in the data is generally a small fraction of the total number of samples. This setup is limited to the identification of known anomalies. The more challenging *unsupervised* setting assumes that the data is a mix of normal and anomalies, without access to labels. Methods in this category often make assumptions about the data-generation process. Therefore, embedding techniques and deep generative models are prime candidates. However, a challenge for deep models is the fact that they tend to model the anomalies within the input data more easily, making the task of identifying them harder. A middle ground between supervised and unsupervised is *semi-supervised* or *one-class classification* setting, where one has access to purely normal samples during training, yet anomalies of unknown nature can exist at inference time. Perhaps confusingly, the term semi-supervised is also used when partial labelling of anomalies is available during the training. In this work, we are interested in identifying anomalies with an unknown distribution and therefore do not assume access to any label information for outliers. That is we consider both unsupervised and the one-class classification version of semi-supervised anomaly detection.

## 2.1 DIFFUSION PROBABILISTIC MODELS

A diffusion process is a stochastic process characterized by a probability distribution that evolves over time, governed by the diffusion equation. Diffusion probabilistic models (Sohl-Dickstein et al., 2015; Ho et al., 2020) are latent variable probabilistic models where the state at time steps larger than zero are considered latent variables. Let $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ denote the data and $\mathbf{x}_1, \ldots, \mathbf{x}_T$ denote the corresponding latent variables. The forward diffusion process is generally fixed to add Gaussian noise at each timestep according to a variance schedule $\beta_1, \ldots, \beta_T$. The approximate posterior $q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)$ is given by,

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) := \prod_{t=1}^{T} q(\mathbf{x}_t|\mathbf{x}_{t-1}), \qquad q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \tag{1}$$

Choosing the transitions as Gaussian distributions enables sampling $\mathbf{x}_t$ at any time in closed form. Let $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{s=1}^{t} \alpha_s$, then,

$$q(\mathbf{x}_t|\mathbf{x}_0) := \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1-\bar{\alpha}_t)\mathbf{I}). \tag{2}$$

Diffusion probabilistic models then learn transitions that reverse the forward diffusion process. Starting at $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; 0, \mathbf{I})$, the joint distribution of the reverse process $p_\theta(\mathbf{x}_{0:T})$ is given by,

$$p_\theta(\mathbf{x}_{0:T}) := p(\mathbf{x}_T) \prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t), \qquad p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)) \tag{3}$$

This parameterized Markov chain also called the reverse process, can produce samples matching the data distribution after a finite number of transition steps.

## 3 DIFFUSION TIME ESTIMATION

Denoising diffusion probabilistic models (DDPM), as introduced in (Ho et al., 2020), can be used to generate samples matching the data distribution even in high-dimensional spaces. The reverse diffusion process implicitly learns the score function of the data distribution and can be used for the likelihood-based identification of anomalies. A common approach used in prior works on anomaly detection using diffusion models (Wolleb et al., 2022; Zhang et al., 2023a; Wyatt et al., 2022) is to reconstruct input samples by simulating the reverse diffusion chain and then using the reconstruction distance to identify anomalies. This is particularly useful where anomalies are localized in the image, and the difference between the input and its reconstruction identifies this localized anomaly. While all previous works focus on this scenario in image data, we consider the broader problem of identification of anomalous samples without assumptions on data type or the nature of the anomaly.

Toward this objective, we evaluate the reconstruction-based approach using DDPMs on the AD-Bench benchmark, which comprises 57 datasets, including tabular, image, and natural language data. We observe that the choice of timestep at the start of reverse diffusion is arbitrary, yet it can significantly affect the anomaly detection performance. We found that using 25% of the maximum timestep globally leads to good results; see the Appendix A for an ablation.

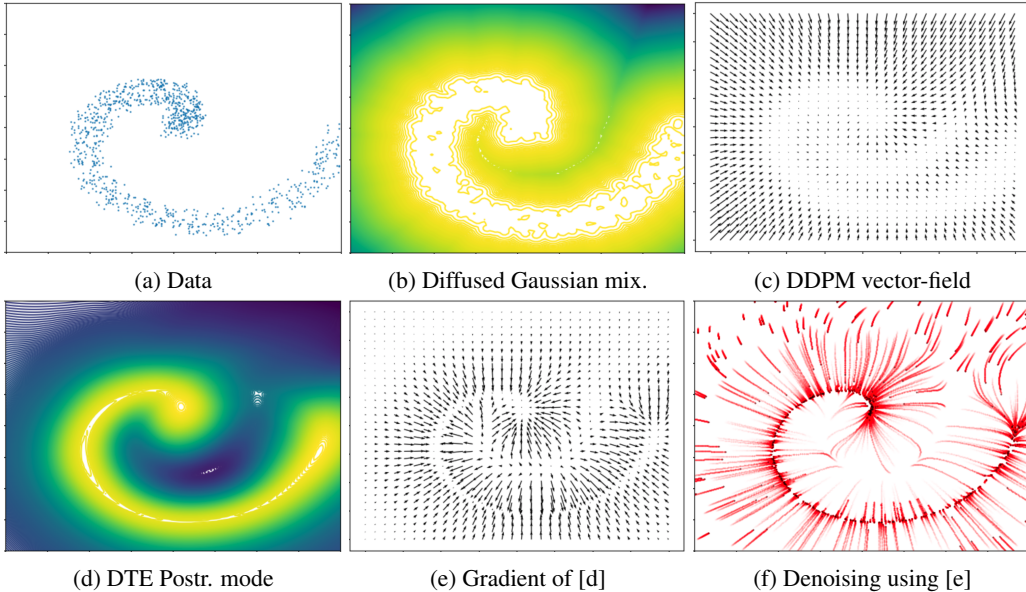|                  |                       |                     |
|------------------|-----------------------|---------------------|
| (a) Data         | (b) Diffused Gaussian mix. | (c) DDPM vector-field |
| (d) DTE Postr. mode | (e) Gradient of [d]   | (f) Denoising using [e] |

Figure 2: DDPM and DTE on a toy dataset shown in (a). (b) shows the Gaussian density function associated with the lowest timestep of DDPM and (c) shows the vector field corresponding to the gradient of this density. (d) plots the mode of the DTE posterior distribution over diffusion time, which we show in subsequent sections is an inverse Gamma distribution. (e) shows the gradient of (d), and (f) shows the flow associated with this gradient, showing that random samples are mapped toward the data manifold.

As anticipated, the expressivity of these models allows them to perform competitively compared to prior work. However, inference for a single data point involves simulating the reverse diffusion chain in its entirety, making this approach computationally expensive. By quantifying the disparity between the reconstructed output and the original input, the objective is to effectively capture the deviations of anomalous samples from the underlying data manifold. We contend that modeling the score function by learning the reverse process is unnecessary if the objective is only the identification of anomalies.

Building upon this idea, we propose a much simpler approach that does not require modeling the reverse diffusion process but instead models the distribution over diffusion time corresponding to noisy input samples. Assuming anomalies are distanced from the data manifold, the density for larger timesteps should have a higher value for anomalies, enabling their probabilistic identification. This can be seen as a direct estimation of reconstruction error.

More concretely, we simulate anomalous samples using a diffusion process and train a neural network to predict the diffusion time corresponding to the noisy samples. Provided that the noisy samples cover the entire feature space, this procedure should also capture potential anomalies. Figure 2 contrasts DDPM and DTE on a toy dataset. The success of our method in using diffusion for anomaly detection is due to the space-filling property of the diffusion process; different regions of the space are sampled at different rates, depending on their proximity to the data manifold. To our knowledge, this is the first setting that uses this property of diffusion beyond its application in learning time-dependent score functions for generative modelling. While in that setting, the estimated score is able to meaningfully approximate the true score over the entire space, we show that we are able to approximate the diffusion time for arbitrary points, including normal or anomalous points.

## 3.1 POSTERIOR DISTRIBUTION OF DIFFUSION TIME

Assuming $\mathbf{x}_s \in \mathbb{R}^d$ is produced through a diffusion process, starting from the data manifold, our goal in this section is to identify the distribution over its diffusion time, as a surrogate for its distance from the manifold. The diffusion process described by Equation (2) specifies a distribution corresponding to each timestep. First, let us assume the dataset consists of a single data point at the

origin. Denote the variance at time $t$ as $\sigma_t^2 = 1 - \bar{\alpha}_t$, and consider the $d$-dimensional zero mean Gaussian distribution at each timestep $\mathcal{N}(\mathbf{0}, \sigma_{\mathbf{t}}^{\mathbf{2}})$. The posterior distribution over $\sigma_t^2$ given $\mathbf{x}_s$ is:

$$p(\sigma_t^2|\mathbf{x}_s) \propto p(\mathbf{x}_s|\sigma_t^2)\,p(\sigma_t^2) = \mathcal{N}(\mathbf{x}_s; \mathbf{0}, \sigma_{\mathbf{t}}^{\mathbf{2}}) \propto \sigma_{\mathbf{t}}^{-\mathbf{d}} \exp\left(-\frac{\|\mathbf{x_s}\|^{\mathbf{2}}}{2\sigma_{\mathbf{t}}^{\mathbf{2}}}\right)$$

This is an *inverse Gamma distribution* $p(\sigma_t^2; a, b) = \frac{b^a}{\Gamma(a)}\left(\frac{1}{\sigma_t^2}\right)^{a+1} \exp\left(-\frac{b}{\sigma_t^2}\right)$ with parameter values $a = d/2 - 1$ and $b = \|\mathbf{x}_s\|^2/2$.
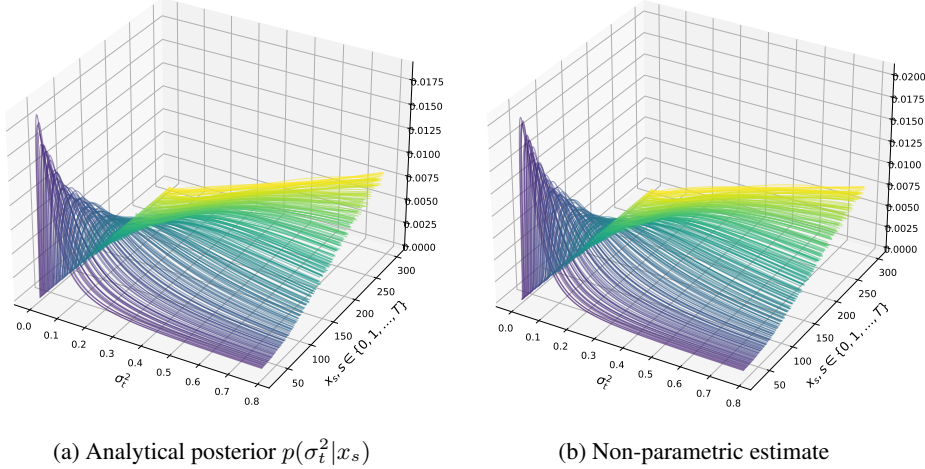


(a) Analytical posterior $p(\sigma_t^2|x_s)$          (b) Non-parametric estimate

Figure 3: Posterior timestep distribution $p(\sigma_t^2|\mathbf{x}_s)$, where $\mathbf{x}_s$ is produced using diffusion with different time steps $s \in \{1, \ldots, T\}$, averaged over the `vertebral` dataset. (a) shows the analytical distribution computed by placing Gaussian distributions of different variances at each point in the dataset, and (b) shows the inverse Gamma distribution with scale parameter value depending on the average distance to the k-nearest neighbours ($k = 32$).

If instead of a single data point at the origin, we have a dataset $\mathcal{D}$, with the corresponding data distribution $p(\mathbf{x})$, we have

$$p(\sigma_t^2|\mathbf{x}_s) \propto p(\mathbf{x}_s|\sigma_t^2)p(\sigma_t^2) = \sum_{\mathbf{x}_0} p(\mathbf{x}_s|\mathbf{x}_0, \sigma_t^2)p(\mathbf{x}_0) = \sum_{\mathbf{x}_0 \in \mathcal{D}} \mathcal{N}(\mathbf{x}_s; \mathbf{x}_0, \sigma_t^2 \mathbf{I}). \tag{4}$$

We refer to Equation (4) as the analytic estimator in subsequent sections since it is the exact posterior distribution. The posterior distribution can be interpreted as adding the likelihoods of Gaussian distributions centered around data points $\mathbf{x}_0 \in \mathcal{D}$ with different (time-dependent) variances. Substituting the Gaussian density function and simplifying, we get

$$p(\sigma_t^2|\mathbf{x}_s) \propto \sum_{\mathbf{x}_0 \in \mathcal{D}} \sigma_t^{-d} \exp\left(-\frac{\|\mathbf{x}_s - \mathbf{x}_0\|^2}{2\sigma_t^2}\right) = \sigma_t^{-d} \exp\left(\log\left(\sum_{\mathbf{x}_0 \in \mathcal{D}} \exp\left(-\frac{\|\mathbf{x}_s - \mathbf{x}_0\|^2}{2\sigma_t^2}\right)\right)\right).$$

We can approximate the log-sum-exp term using `max` function:

$$p(\sigma_t^2|\mathbf{x}_s) \gtrsim \sigma_t^{-d} \exp\left(\max_{\mathbf{x}_0 \in \mathcal{D}} -\frac{\|\mathbf{x}_s - \mathbf{x}_0\|^2}{2\sigma_t^2}\right) = \sigma_t^{-d} \exp\left(-\frac{1}{\sigma_t^2} \min_{\mathbf{x}_0 \in \mathcal{D}} \frac{\|\mathbf{x}_s - \mathbf{x}_0\|^2}{2}\right) \tag{5}$$

The posterior over diffusion time approximately has the form of an inverse Gamma distribution with the shape parameter $a = d/2 - 1$ depending only on the dimensionality of the data and the scale parameter $b = \min_{\mathbf{x}_0 \in \mathcal{D}} \frac{\|\mathbf{x}_s - \mathbf{x}_0\|^2}{2}$ depending on the distance of the input point to the closest point in the dataset. Note that, as $a > 0 \implies d > 2$, this analysis is only valid for three or higher dimensions.

## 3.2 NON-PARAMETRIC MODEL

The posterior over diffusion time given by Equation (5) can potentially be used as a non-parametric approach to anomaly detection. The approximation of log-sum-exp using the maximum value (nearest neighbour) becomes less accurate for larger timesteps, in which a point has a comparable distance to several points in the dataset. We found that instead of setting the scale parameter $b$ based on the distance to the closest point, approximating log-sum-exp using the average distance to k-nearest neighbours of the input point works better in practice. The non-parametric estimator is then:

$$p(\sigma_t^2|\mathbf{x}_s) \propto \sigma_t^{-d} \exp\left(-\frac{1}{\sigma_t^2} \cdot \frac{1}{K} \sum_{\mathbf{x}_0 \in \mathrm{kNN}(\mathbf{x}_s)} \frac{\|\mathbf{x}_s - \mathbf{x}_0\|^2}{2}\right) \tag{6}$$

Figure 3 shows the analytical posterior distribution obtained using Equation (4) and the non-parametric estimator given in Equation (6) for a real dataset.

The upshot is that, given a point $\mathbf{x}_s$, this method approximates the scale parameter of the inverse Gamma distribution using the average distance to its $k$-nearest neighbours. The anomaly score is the mean of this distribution over diffusion time. As seen in Figure 3, points $\mathbf{x}_s$ that are produced using diffusion with larger time-steps also have a higher posterior mean, on average, enabling us to identify them as points that are far from the manifold. Interestingly, this method closely resembles the classical $k$-nearest neighbours (kNN). In fact, the anomaly rankings given by these methods are identical. In our experiments, the difference in score comes from the distance calculation: for DTE non-parametric, we take the mean distance from the k-nearest neighbours as opposed to (a variation of) kNN that takes the distance from the kth-nearest neighbour.

## 3.3 PARAMETRIC MODEL

The non-parametric estimator of diffusion time becomes compute and memory-intensive when dealing with large datasets due to the need to find the k-nearest neighbours for each input sample in the entire dataset. To tackle the scalability problem, we employ deep neural networks to estimate the posterior distribution, which also enhances generalization capabilities. The full training procedure for both parametric models is available in Appendix D.2.
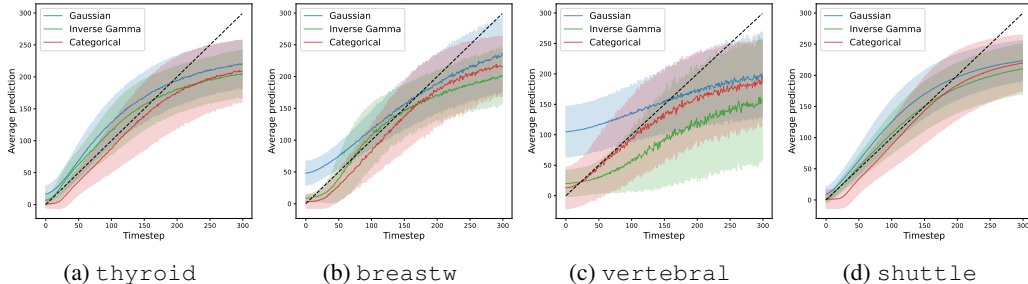


|     |     |     |     |
| :-: | :-: | :-: | :-: |
| (a) `thyroid` | (b) `breastw` | (c) `vertebral` | (d) `shuttle` |

Figure 4: Predicted diffusion time against ground truth diffusion time for Gaussian model ($\ell_2$-regression), Inverse Gamma model, and categorical model (with seven bins) on the test set for various datasets. The maximum length of the diffusion Markov chain is $T = 300$. The shaded region indicates the standard deviation in predictions across the dataset.

**Inverse Gamma model** In Section 3.1 we saw that the posterior distribution over time-dependent variance has the form of an inverse Gamma distribution. We train a deep neural network parameterized by $\theta$, which we denote by $f_\theta$, to predict the scale parameter $b$ of the inverse Gamma distribution, given the noisy sample $\mathbf{x}_t$. Since the shape parameter $a$ depends only on the dimensionality of the data, it is a known fixed parameter. We minimize the negative log-likelihood given by:

$$\mathcal{L}(\theta) := -\mathbb{E}_{t,\mathbf{x}_0}\left[a \log f_\theta(\mathbf{x}_t) - (a+1)\log \sigma_t^2 - f_\theta(\mathbf{x}_t)/\sigma_t^2\right] \tag{7}$$

The expectation is over data samples $\mathbf{x}_0 \sim p(\mathbf{x})$ and timesteps $t \sim \mathcal{U}[1, T]$. The mode of the distribution is used as anomaly score.

Figure 4 shows the predicted timestep for the inverse Gamma model applied to different datasets, with the length of Markov chain $T = 300$. Compared to standard $\ell_2$ regression which assumes that the output variable is Gaussian distributed, the inverse Gamma model has a much lower bias for diffusion time prediction for smaller timesteps, which empirically validates our analysis. However, this model suffers from high bias and high variance for larger timesteps. The high bias can be attributed to the approximation error of log-sum-exp using k-nearest neighbours, which becomes inaccurate for larger timesteps. The high variance is a consequence of the shape of the inverse Gamma distribution, which becomes flat for large values of the scale parameter (see Figure 3).

**Categorical model** The inverse Gamma model while analytically accurate, can restrict the expressivity of the neural network. In order to provide more flexibility in learning the diffusion time distribution, we can model it as a categorical distribution over $T$ classes, where $T$ is the length of the Markov chain associated with the diffusion process. This approach does not assume any parametric distribution over diffusion time and requires the model to accurately predict the full distribution. Let $y_t \in \{0, 1\}^T$ denote the one-hot vector with one at coordinate $t$, and $f_\theta$ denote the deep neural network that predicts the class probabilities, $f_\theta : \mathcal{X} \to [0, 1]^T$. We minimize the cross-entropy loss function, which is equivalent to maximizing the log-likelihood of the categorical distribution:

$$\mathcal{L}(\theta) := \mathbb{E}_{t,x_0} \left[ -\sum_{k=0}^{K} y_t^{(k)} \log \left( f_\theta(\mathbf{x}_t)^{(k)} \right) \right] \tag{8}$$

In practice, we simplify the learning task by combining timesteps into bins and training a model to predict the correct bin. If $B$ denotes the number of bins, then the corresponding bin for a timestep $t$ would be $\lfloor \frac{t \cdot B}{T} \rfloor$. Figure 4 shows the predicted timestep for the categorical model on different datasets. Compared to the inverse Gamma model, it suffers from significantly less bias across the entire range of timesteps. The score calculation is described in Appendix D.3 with the training algorithm in Appendix D.2.

## 4 EXPERIMENTS

**Setting** We perform experiments on the ADBench benchmark (Han et al., 2022), which comprises a set of popular tabular anomaly detection datasets as well as newly created tabular datasets made from images and natural language tasks, all described in Appendix D.1. The implementation details are provided in Appendix D, with the training algorithm, model architecture, hyperparameters, and comparison of the run-time. Some ablation studies are in Appendix A. We implement and compare the results of the various approaches proposed in Section 3: the non-parametric, the parametric inverse Gamma, and the parametric categorical DTE.

**Baselines** We compare against all the unsupervised learning methods included in ADBench. These include classical methods, namely CBLOF (He et al., 2003), COPOD (Li et al., 2020), ECOD (Li et al., 2022), FeatureBagging (Lazarevic & Kumar, 2005), HBOS (Goldstein & Dengel, 2012), IForest (Liu et al., 2008), kNN (Ramaswamy et al., 2000), LODA (Pevný, 2016), LOF (Breunig et al., 2000), MCD (Fauconnier & Haesbroeck, 2009), OCSVM (Schölkopf et al., 1999), and PCA (Shyu et al., 2003). The deep learning-based methods include DeepSVDD (Ruff et al., 2018), and DAGMM (Zong et al., 2018). Outside of ADBench, we also compare against some more recently proposed deep learning-based approaches such as DROCC (Goyal et al., 2020), GOAD (Bergman & Hoshen, 2020), ICL (Shenkar & Wolf, 2022), SLAD (Xu et al., 2023b) and DIF (Xu et al., 2023a); see Section 5 for a brief overview. For each method, we picked the best-performing set of hyperparameters given in their original paper. We also have four additional generative baselines: normalizing flows with planar flows (Rezende & Mohamed, 2015) to identify anomalies based on the log-likelihood, DDPM, VAE (Kingma & Welling, 2013) and GAN (Goodfellow et al., 2014) to reconstruct the input and compare it with the original input to identify anomalies.

**Results** Figure 5 shows the overall performance of these different methods on 57 tasks in AD-Bench, each limited to 50,000 data points. The results for each individual dataset are provided in Appendix F. We report the mean AUC ROC and its standard deviation over five different seeds for each method. For the unsupervised setting, we used bootstrapping over the whole dataset for training, while inference is made on the full dataset. For the semi-supervised setting, we used 50% of
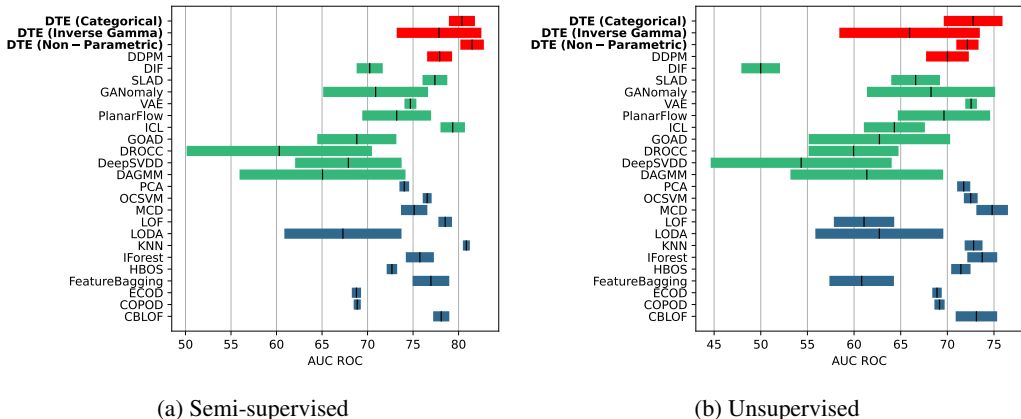
(a) Semi-supervised

(b) Unsupervised

Figure 5: AUC ROC means and standard deviations on the 57 datasets from ADBench over five different seeds for a) the semi-supervised setting using normal samples only for training and b) the unsupervised setting with bootstrapped training instances. Colour scheme: red (diffusion-based), green (deep learning methods), blue (classical methods). DTE outperforms all baselines for the semi-supervised setting apart from kNN. It is also competitive in the unsupervised setting.

the normal samples for training, while the test set contains the rest of the normal samples and all anomalous samples. The proposed method is among the few competitive in both semi-supervised and unsupervised settings. In particular, our method outperforms all previous deep learning-based approaches in both settings significantly and also outperforms the DDPM model. Unsurprisingly, deep learning methods have a higher variance than non-parametric methods. Using bagging can be a way to help reduce the variance at the cost of more training and inference time.

Figure 1 compares our method's performance and inference time with the other baselines. In some applications, such as medical and network monitoring, fast inference time is crucial as the algorithm must detect anomalies in real time. Our method uses a forward pass through a simple neural network for predictions, which gives it the shortest inference time over all the methods considered here. Training time, inference time and compute amounts are available in Appendix D.4.

**Choice of representation** ADBench's image datasets use vector representation derived from pre-trained ImageNet embeddings. We investigated the impact of representation quality for semi-supervised anomaly detection across several datasets: VisA (Zou et al., 2022), CIFAR-10, and MNIST. We observe that different methods, including DDPM, kNN and DTE, perform better when applied to image embeddings rather than raw images. In particular, embeddings produced through self-supervision are generally of higher quality when compared to those produced for classification, and the embeddings that are specialized or fined-tuned to the target dataset produce the best results. The results are reported in Appendix E.

We also observe that kNN remains a top-performing algorithm for anomaly detection, where its only disadvantage remains its scalability. As explained in Section 3.2, the non-parametric method gives the same anomaly ranking as kNN. DTE can thus be approximately interpreted as a parametric $k$-nearest neighbours algorithm which can be beneficial for large datasets that require smaller inference time. To understand the anomalies, both DDPM and DTE are able to identify a "denoised" data point; DDPM depends on an initial time step hyper-parameter, whereas DTE does not, by using deterministic ODE flow. However, DDPM outperforms in denoising, being explicitly trained for it. Further interpretability discussion, illustrated with a toy example, is in Appendix B.

## 5  RELATED WORK

We refer the reader to the following surveys for a comprehensive review (Pang et al., 2021; Chandola et al., 2009; Ruff et al., 2021; Hodge & Austin, 2004). Although recently, the spotlight has shifted towards deep learning methodologies, classical techniques such as kNN (Ramaswamy et al., 2000)

persistently exhibit strong performance. We compared our method with some of these techniques in Section 4. Clustering and nearest neighbour algorithms use the distance to score instances, making them easily interpretable. Clustering algorithms, such as CBLOF (He et al., 2003) and k-means (MacQueen, 1967), assume that anomalies are either not part of cluster, are part of smaller clusters than normal instances, or lie further away from the cluster centroid. In contrast, nearest neighbour algorithms use the distance between points or relative density with respect to their neighbourhood.

As anomalies can be more difficult to detect in high-dimensional spaces and complex data distributions (Pang et al., 2021), the development of deep anomaly detection algorithms has been increasing over the past few years (Ruff et al., 2021). In particular, several works combine autoencoders with other classical techniques (Zhou & Paffenroth, 2017; Kim et al., 2020; An & Cho, 2015; Erfani et al., 2016; Sakurada & Yairi, 2014; Xia et al., 2015). Other notable methods include DeepSVDD (Ruff et al., 2018), DAGMM (Zong et al., 2018); Lunar (Goodge et al., 2022), DROCC (Goyal et al., 2020), GOAD (Bergman & Hoshen, 2020), SO-GAAL and MO-GAAL (Liu et al., 2019), SLAD (Xu et al., 2023b) and DIF (Xu et al., 2023a). Deep kNN methods (Pang et al., 2018; Sun et al., 2022) learn representations to apply kNN. ICL (Shenkar & Wolf, 2022), which uses contrastive representation learning reported competitive results for ODDS datasets, for the semi-supervised setting.

**Diffusion-based Techniques** While diffusion models have been previously used for anomaly detection in image and video (Yan et al., 2023; Flaborea et al., 2023; Tur et al., 2023) data for a one-class setting (semi-supervised), their application in the context of tabular data and the unsupervised setting was unexplored. Wolleb et al. (2022) proposed an encoding method using a diffusion process followed by a denoising procedure guided by a classifier. Zhang et al. (2023a) synthesizes anomaly samples to train the denoising network for anomaly repair. AnoDDPM employs a specific diffusion noise to train a denoising network for normal image reconstruction (Wyatt et al., 2022). Similarly, Graham et al. (2023) utilized a DDPM to reconstructs an image for multiple different timesteps combined together to make anomaly scores. Liu et al. (2023) introduced a diffusion method that reconstruct an image by in-painting the input masked by a checkerboard pattern. Lastly, Zhang et al. (2023b) used a latent diffusion model trained with simulated anomalous samples on images.

# 6 CONCLUSION

This paper investigates the applicability of diffusion modelling for unsupervised and semi-supervised anomaly detection. We observe that specific design choices in DDPMs, although somewhat arbitrary, significantly influence their performance. Despite the expressivity and interpretability of DDPMs, they come with notable computational overhead compared to existing parametric techniques. For anomaly detection, DDPM essentially estimates the distance between the input and its "denoised reconstruction"; we observe that one could directly produce this estimate, or equivalently estimate the diffusion time. We first observe that the distribution of diffusion time given a noisy input, follows an inverse Gamma distribution. This forms the basis for our non-parametric approach that accurately predicts the diffusion time and turns out to create the same anomaly score ranking as kNN. A subsequent parametric strategy leverages a deep neural network, harnessing its generalization and rapid inference capabilities for large datasets. We evaluate the effectiveness of DTE on ADBench, a benchmark comprising popular anomaly detection datasets. Our results demonstrate competitive performance compared to prior work while improving the inference time by several orders of magnitude. Furthermore, we find that using pre-trained embeddings for images considerably improves the performance of diffusion-based methods, showing the potential advantage of using latent space diffusion.

# 7 LIMITATIONS AND FUTURE WORK

While our approach, DTE, achieves excellent performance with low inference time, it is important to acknowledge that in terms of interpretability, DTE falls behind DDPM as we explain in Appendix B and Section 4. This may pose challenges for practitioners seeking to understand the underlying mechanisms and behaviours of the data. Evaluating DTE in handling larger and more complex real-world datasets remains an avenue for future exploration. While here, we only address point anomalies, applications of diffusion modelling for group and contextual anomalies remain a high-impact unexplored area that we plan to investigate in the future.

## 8 REPRODUCIBILITY STATEMENT

We have made efforts to ensure that our method is reproducible. Appendix D.1 provides a description of all datasets included in ADBench, along with the preprocessing steps. Appendix D.2 presents a formal algorithm for parametric DTE and Appendix D.3 provides a detailed description of the network architecture and hyperparameters. We provide full results for both the unsupervised and semi-supervised settings with additional metrics, for all individual datasets and baselines in Appendix F as a reference for researchers to reproduce our experimental results. We are releasing the code as part of the supplemental material with detailed explanations to run the experiments.

## ACKNOWLEDGEMENTS

## REFERENCES

Mohiuddin Ahmed, Abdun Naser Mahmood, and Jiankun Hu. A survey of network anomaly detection techniques. *Journal of Network and Computer Applications*, 60:19–31, 2016a.

Mohiuddin Ahmed, Abdun Naser Mahmood, and Md Rafiqul Islam. A survey of anomaly detection techniques in financial domain. *Future Generation Computer Systems*, 55:278–288, 2016b.

Jinwon An and Sungzoon Cho. Variational autoencoder based anomaly detection using reconstruction probability. 2015.

Lion Bergman and Yedid Hoshen. Classification-based anomaly detection for general data. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=H1lK_lBtvS.

Markus M. Breunig, Hans-Peter Kriegel, Raymond T. Ng, and Jörg Sander. Lof: Identifying density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*, SIGMOD '00, pp. 93–104, New York, NY, USA, 2000. Association for Computing Machinery. ISBN 1581132174. doi: 10.1145/342009.335388. URL https://doi.org/10.1145/342009.335388.

Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM Comput. Surv.*, 41(3), jul 2009. ISSN 0360-0300. doi: 10.1145/1541880.1541882. URL https://doi.org/10.1145/1541880.1541882.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*, 2019. URL https://api.semanticscholar.org/CorpusID:52967399.

Sarah M. Erfani, Sutharshan Rajasegarar, Shanika Karunasekera, and Christopher Leckie. High-dimensional and large-scale anomaly detection using a linear one-class svm with deep learning. *Pattern Recognition*, 58:121–134, 2016. ISSN 0031-3203. doi: https://doi.org/10.1016/j.patcog.2016.03.028. URL https://www.sciencedirect.com/science/article/pii/S0031320316300267.

C. Fauconnier and Gentiane Haesbroeck. Outliers detection with the minimum covariance determinant estimator in practice. *Statistical Methodology*, 6:363–379, 07 2009. doi: 10.1016/j.stamet.2008.12.005.

Alessandro Flaborea, Luca Collorone, Guido Maria D'Amely di Melendugno, Stefano D'Arrigo, Bardh Prenkaj, and Fabio Galasso. Multimodal motion conditioned diffusion model for skeleton-based video anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 10318–10329, October 2023.

Katherine Fraser, Samuel Homiller, Rashmish K. Mishra, Bryan Ostdiek, and Matthew D. Schwartz. Challenges for unsupervised anomaly detection in particle physics. *Journal of High Energy Physics*, 2022(3), mar 2022. doi: 10.1007/jhep03(2022)066. URL `https://doi.org/10.1007%2Fjhep03%282022%29066`.

Markus Goldstein and Andreas Dengel. Histogram-based outlier score (hbos): A fast unsupervised anomaly detection algorithm. 09 2012.

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014.

Adam Goodge, Bryan Hooi, See Kiong Ng, and Wee Siong Ng. Lunar: Unifying local outlier detection methods via graph neural networks. 2022.

Yury Gorishniy, Ivan Rubachev, Valentin Khrulkov, and Artem Babenko. Revisiting deep learning models for tabular data. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021. URL `https://openreview.net/forum?id=i_Q1yrOegLY`.

Sachin Goyal, Aditi Raghunathan, Moksh Jain, Harsha Vardhan Simhadri, and Prateek Jain. DROCC: Deep robust one-class classification. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 3711–3721. PMLR, 13–18 Jul 2020. URL `https://proceedings.mlr.press/v119/goyal20c.html`.

Mark S. Graham, Walter H.L. Pinaya, Petru-Daniel Tudosiu, Parashkev Nachev, Sebastien Ourselin, and Jorge Cardoso. Denoising diffusion models for out-of-distribution detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pp. 2947–2956, June 2023.

Songqiao Han, Xiyang Hu, Hailiang Huang, Minqi Jiang, and Yue Zhao. ADBench: Anomaly detection benchmark. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. URL `https://openreview.net/forum?id=foA_SFQ9zo0`.

Douglas M Hawkins. *Identification of outliers*, volume 11. Springer, 1980.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. URL `https://arxiv.org/abs/1512.03385`.

Zengyou He, Xiaofei Xu, and Shengchun Deng. Discovering cluster-based local outliers. *Pattern Recogn. Lett.*, 24(9–10):1641–1650, jun 2003. ISSN 0167-8655. doi: 10.1016/S0167-8655(03)00003-5. URL `https://doi.org/10.1016/S0167-8655(03)00003-5`.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.

Victoria Hodge and Jim Austin. A survey of outlier detection methodologies. *Artificial Intelligence Review*, 22:85–126, 10 2004. doi: 10.1023/B:AIRE.0000045502.10941.a9.

Ki Hyun Kim, Sangwoo Shim, Yongsub Lim, Jongseob Jeon, Jeongwoo Choi, Byungchan Kim, and Andre S. Yoon. Rapp: Novelty detection with reconstruction along projection pathway. In *ICLR*. OpenReview.net, 2020. URL `http://dblp.uni-trier.de/db/conf/iclr/iclr2020.html#KimSLJCKY20`.

Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2013. URL `https://api.semanticscholar.org/CorpusID:216078090`.

Akim Kotelnikov, Dmitry Baranchuk, Ivan Rubachev, and Artem Babenko. Tabddpm: Modelling tabular data with diffusion models, 2022. URL `https://arxiv.org/abs/2209.15421`.

Aleksandar Lazarevic and Vipin Kumar. Feature bagging for outlier detection. In *Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery in Data Mining*, KDD '05, pp. 157–166, New York, NY, USA, 2005. Association for Computing Machinery. ISBN 159593135X. doi: 10.1145/1081870.1081891. URL `https://doi.org/10.1145/1081870.1081891`.

Zheng Li, Yue Zhao, Nicola Botta, Cezar Ionescu, and Xiyang Hu. COPOD: Copula-based outlier detection. In *2020 IEEE International Conference on Data Mining (ICDM)*. IEEE, nov 2020. doi: 10.1109/icdm50108.2020.00135. URL `https://doi.org/10.1109%2Ficdm50108.2020.00135`.

Zheng Li, Yue Zhao, Xiyang Hu, Nicola Botta, Cezar Ionescu, and George Chen. ECOD: Unsupervised outlier detection using empirical cumulative distribution functions. *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–1, 2022. doi: 10.1109/tkde.2022.3159580. URL `https://doi.org/10.1109%2Ftkde.2022.3159580`.

Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In *2008 Eighth IEEE International Conference on Data Mining*, pp. 413–422, 2008. doi: 10.1109/ICDM.2008.17.

Yezheng Liu, Zhe Li, Chong Zhou, Yuanchun Jiang, Jianshan Sun, Meng Wang, and Xiangnan He. Generative adversarial active learning for unsupervised outlier detection, 2019.

Zhenzhen Liu, Jinjie Zhou, Yufan Wang, and Kilian Q. Weinberger. Unsupervised out-of-distribution detection with diffusion inpainting. In *International Conference on Machine Learning*, 2023. URL `https://api.semanticscholar.org/CorpusID:257050245`.

J. B. MacQueen. Some methods for classification and analysis of multivariate observations. In L. M. Le Cam and J. Neyman (eds.), *Proc. of the fifth Berkeley Symposium on Mathematical Statistics and Probability*, volume 1, pp. 281–297. University of California Press, 1967.

Girik Pachauri and Sandeep Sharma. Anomaly detection in medical wireless sensor networks using machine learning algorithms. *Procedia Computer Science*, 70:325–333, 2015. ISSN 1877-0509. doi: https://doi.org/10.1016/j.procs.2015.10.026. URL `https://www.sciencedirect.com/science/article/pii/S1877050915031907`. Proceedings of the 4th International Conference on Eco-friendly Computing and Communication Systems.

Guansong Pang, Longbing Cao, Ling Chen, and Huan Liu. Learning representations of ultrahigh-dimensional data for random distance-based outlier detection. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '18, pp. 2041–2050, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450355520. doi: 10.1145/3219819.3220042. URL `https://doi.org/10.1145/3219819.3220042`.

Guansong Pang, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. Deep learning for anomaly detection. *ACM Computing Surveys*, 54(2):1–38, mar 2021. doi: 10.1145/3439950. URL `https://doi.org/10.1145%2F3439950`.

Tomáš Pevný. Loda: Lightweight on-line detector of anomalies. *Mach. Learn.*, 102(2):275–304, feb 2016. ISSN 0885-6125. doi: 10.1007/s10994-015-5521-0. URL `https://doi.org/10.1007/s10994-015-5521-0`.

Sridhar Ramaswamy, Rajeev Rastogi, and Kyuseok Shim. Efficient algorithms for mining outliers from large data sets. *SIGMOD Rec.*, 29(2):427–438, may 2000. ISSN 0163-5808. doi: 10.1145/335191.335437. URL `https://doi.org/10.1145/335191.335437`.

Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 1530–1538, Lille, France, 07–09 Jul 2015. PMLR. URL `https://proceedings.mlr.press/v37/rezende15.html`.

Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. Deep one-class classification. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 4393–4402. PMLR, 10–15 Jul 2018. URL `https://proceedings.mlr.press/v80/ruff18a.html`.

Lukas Ruff, Jacob R. Kauffmann, Robert A. Vandermeulen, Gregoire Montavon, Wojciech Samek, Marius Kloft, Thomas G. Dietterich, and Klaus-Robert Muller. A unifying review of deep and shallow anomaly detection. *Proceedings of the IEEE*, 109(5):756–795, may 2021. doi: 10.1109/jproc.2021.3052449. URL `https://doi.org/10.1109%2Fjproc.2021.3052449`.

Mayu Sakurada and Takehisa Yairi. Anomaly detection using autoencoders with nonlinear dimensionality reduction. In *Proceedings of the MLSDA 2014 2nd Workshop on Machine Learning for Sensory Data Analysis*, MLSDA'14, pp. 4–11, New York, NY, USA, 2014. Association for Computing Machinery. ISBN 9781450331593. doi: 10.1145/2689746.2689747. URL `https://doi.org/10.1145/2689746.2689747`.

Osman Salem, Alexey Guerassimov, Ahmed Mehaoua, Anthony Marcus, and Borko Furht. Sensor fault and patient anomaly detection and classification in medical wireless sensor networks. In *2013 IEEE International Conference on Communications (ICC)*, pp. 4373–4378, 2013. doi: 10.1109/ICC.2013.6655254.

Bernhard Schölkopf, Robert Williamson, Alex Smola, John Shawe-Taylor, and John Platt. Support vector method for novelty detection. volume 12, pp. 582–588, 01 1999.

Tom Shenkar and Lior Wolf. Anomaly detection for tabular data with internal contrastive learning. In *International Conference on Learning Representations*, 2022. URL `https://openreview.net/forum?id=_hszZbt46bT`.

Mei-Ling Shyu, Shu-Ching Chen, Kanoksri Sarinnapakorn, and Liwu Chang. A novel anomaly detection scheme based on principal component classifier. 01 2003.

Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pp. 2256–2265. PMLR, 2015.

Yiyou Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest neighbors. *ICML*, 2022.

Gian Antonio Susto, Matteo Terzi, and Alessandro Beghi. Anomaly detection approaches for semiconductor manufacturing. *Procedia Manufacturing*, 11:2018–2024, 2017.

Anil Osman Tur, Nicola Dall'Asen, Cigdem Beyan, and Elisa Ricci. Exploring diffusion models for unsupervised video anomaly detection. *2023 IEEE International Conference on Image Processing (ICIP)*, pp. 2540–2544, 2023. URL `https://api.semanticscholar.org/CorpusID:258079336`.

Julia Wolleb, Florentin Bieder, Robin Sandkühler, and Philippe C Cattin. Diffusion models for medical anomaly detection. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part VIII*, pp. 35–45. Springer, 2022.

Julian Wyatt, Adam Leach, Sebastian M. Schmon, and Chris G. Willcocks. Anoddpm: Anomaly detection with denoising diffusion probabilistic models using simplex noise. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 649–655, 2022. doi: 10.1109/CVPRW56347.2022.00080.

Yan Xia, Xudong Cao, Fang Wen, Gang Hua, and Jian Sun. Learning discriminative reconstructions for unsupervised outlier removal. *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 1511–1519, 2015.

Hongzuo Xu, Guansong Pang, Yijie Wang, and Yongjun Wang. Deep isolation forest for anomaly detection. *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–14, 2023a. doi: 10.1109/TKDE.2023.3270293.

Hongzuo Xu, Yijie Wang, Juhui Wei, Songlei Jian, Yizhou Li, and Ning Liu. Fascinating supervisory signals and where to find them: Deep anomaly detection with scale learning. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org, 2023b.

T. Yairi, Y. Kawahara, R. Fujimaki, Y. Sato, and K. Machida. Telemetry-mining: a machine learning approach to anomaly detection and fault diagnosis for space systems. In *2nd IEEE International Conference on Space Mission Challenges for Information Technology (SMC-IT'06)*, pp. 8 pp.–476, 2006. doi: 10.1109/SMC-IT.2006.79.

Cheng Yan, Shiyu Zhang, Yang Liu, Guansong Pang, and Wenjun Wang. Feature prediction diffusion model for video anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 5527–5537, October 2023.

Hui Zhang, Zheng Wang, Zuxuan Wu, and Yu-Gang Jiang. Diffusionad: Denoising diffusion for anomaly detection, 2023a.

Xinyi Zhang, Naiqi Li, Jiawei Li, Tao Dai, Yong Jiang, and Shu-Tao Xia. Unsupervised surface anomaly detection with diffusion probabilistic model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 6782–6791, October 2023b.

Chong Zhou and Randy C. Paffenroth. Anomaly detection with robust deep autoencoders. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '17, pp. 665–674, New York, NY, USA, 2017. Association for Computing Machinery. ISBN 9781450348874. doi: 10.1145/3097983.3098052. URL https://doi.org/10.1145/3097983.3098052.

Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Dae ki Cho, and Haifeng Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In *International Conference on Learning Representations*, 2018.

Yang Zou, Jongheon Jeong, Latha Pemula, Dongqing Zhang, and Onkar Dabeer. Spot-the-difference self-supervised pre-training for anomaly detection and segmentation. In *European Conference on Computer Vision*, pp. 392–408. Springer, 2022.

# A    ABLATION STUDIES

We perform several ablation studies to understand DDPM and the proposed DTE method.
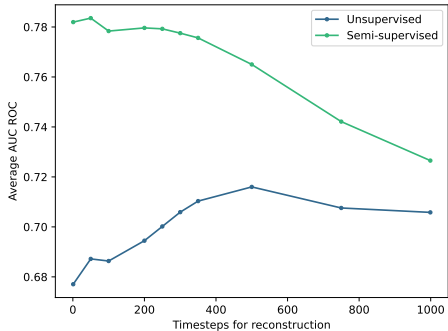


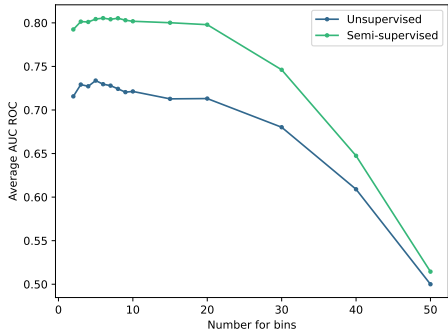Figure 6: Average AUC ROC over the 57 ADBench datasets for different reconstruction timesteps of the DDPM model.

Figure 7: Average AUC ROC over the 57 AD-Bench datasets for different number of bins of the DTE categorical model.

**Reconstruction timestep in DDPM**    When using DDPMs for anomaly detection based on the reconstruction distance, the denoising model requires an input timestep to create the reconstruction. We found that this somewhat arbitrary hyperparameter choice can significantly affect performance as shown in Figure 6.

For the unsupervised setting, we found that a value close to 50% of the maximum timestep results in the highest AUC ROC score on average. For the semi-supervised, the AUC ROC decreases as we increase the reconstruction timestep. Since the model is trained only on normal samples, the anomalies are sufficiently distanced from the learned data manifold for minor changes to result in a large reconstruction error while a larger timestep decreases the precision on normal samples.

**Number of bins in categorical DTE**    As discussed in Section 3.3, we implement categorical DTE by combining multiple timesteps into bins. This turns out to be an important hyperparameter as it affects the final performance significantly. Figure 7 shows that a low number of bins leads to better performance. This can be attributed to the fact that we calculate the mean of the predicted timestep distribution rather than the mode to calculate anomaly scores and that adding more bins increases the complexity of the learning task.

**Maximum timestep in DTE**    We study the effect of changing the maximum timestep in the noising diffusion process. As seen in Figure 8, the maximum timestep affects performance until roughly $T = 250$, since for very low values of $T$, the noisy samples might not resemble standard Gaussian noise and might not cover all potential anomalies in the dataset. We also note categorical DTE is more robust to the value of $T$ than the inverse Gamma DTE.

Figure 9 shows the value of standard deviation versus timestep as we increase the maximum timestep $T$. We observe that for values of $T \geq 250$, the final timestep corresponds to a standard deviation close enough to 1.0 that the data resembles samples drawn from a standard Gaussian distribution.
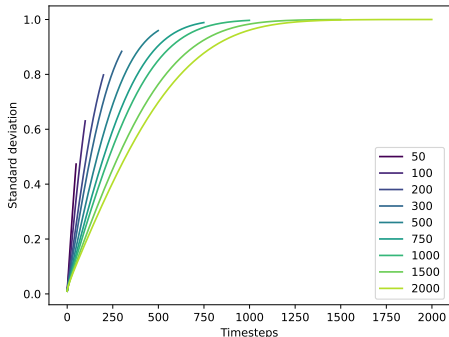


Figure 9: Standard deviations versus timestep for different values of the maximum timestep $T$.

15

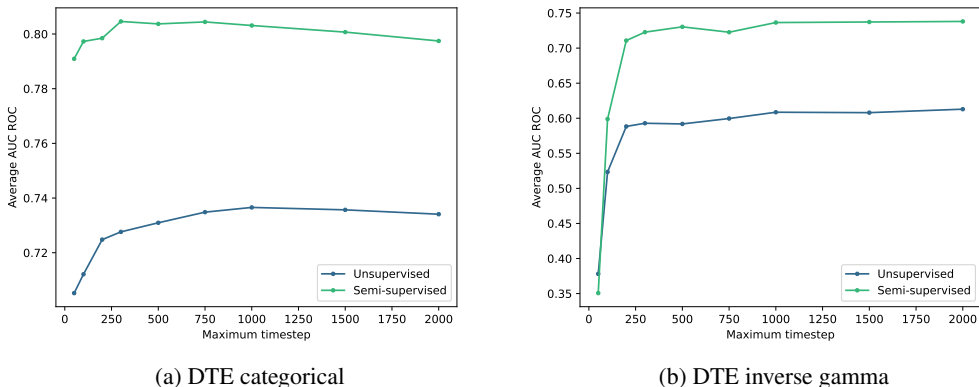(a) DTE categorical          (b) DTE inverse gamma

Figure 8: Average AUC ROC over the 57 ADBench datasets for different maximum timestep T for the categorical and inverse gamma DTE models on both semi-supervised and unsupervised settings.

## B  INTERPRETABILITY

In certain applications, the mere identification of anomalies in the dataset is insufficient; it is imperative to understand the underlying reasons for flagging specific data points as anomalies. Both DDPM and DTE can provide interpretability by identifying a corresponding "denoised" or normal data point. In DDPM this is achieved using the deterministic ODE flow, which is (rather arbitrarily) initialized at some large time step. We found the initial time step to be an important hyper-parameter, which impacts both anomaly detection and interpretability for DDPM. In practice, $T' = .25 \times T$ performs well as the initial time-step. DTE has the benefit of avoiding such hyper-parameters, where one could use the gradient flow associated with the mode of the posterior to denoise a given input; see Figure 2 (d, e, and f).

Figure 10 shows another example, this time using the categorical likelihood on the MNIST dataset. We artificially introduce a gray patch as an anomaly (Figure 10 (a)) and perform the gradient descent procedure reducing the mean



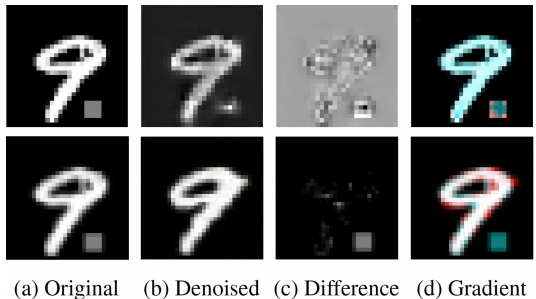(a) Original   (b) Denoised   (c) Difference   (d) Gradient

Figure 10: Interpretability in DTE (first row) and DDPM (second row) on MNIST. Visual interpretation of a gray patch anomaly on an MNIST image using the categorical diffusion model with a simple convolution network on the first row and a DDPM on the second for comparison. a) original anomalous image, b) the denoised version using gradient descent c) difference between the original and the denoised image, d) visualization of the gradient on top of the original image.

of the posterior density. We observe that this procedure indeed partially eliminates the patch (Figure 10 (b)). We also note that since it is explicitly trained to remove the noise from a noisy input, DDPM performs better in removing the patch.

As detailed in Section 3.2, the non-parametric DTE yields the same anomaly score as kNN. Thus, the parametric DTE can be viewed as an approximate parametric kNN algorithm. This perspective enhances DTE's interpretability: the neural network's score represents the estimated distance of a point to the manifold. Although we can't pinpoint which training set instance most closely matches an input, interpreting the score as a distance to a certain neighbourhood offers a straightforward insight into the method's functioning.

16

## C NON-PARAMETRIC ESTIMATION OF TIMESTEP DISTRIBUTION

In Figure 3, we visualize the analytical posterior distribution along with the non-parametric estimate. The difference between these distributions is shown in Figure 11. While the two distributions are quite similar, their shape is very peaked for low values of the diffusion timestep. The slight misalignment between the peaks of the analytical and the non-parametric estimate gives rise to the spiky shape seen in the difference. For higher values of the diffusion timestep, the difference is very close to zero, demonstrating that the non-parametric estimate based on k-nearest neighbours is a very close approximation to the true posterior distribution of timestep.



(a) Analytical posterior $p(\sigma_t^2|x_s)$    (b) Non-parametric estimate    (c) Difference

Figure 11: Posterior timestep distribution $p(\sigma_t^2|\mathbf{x}_s)$, where $\mathbf{x}_s$ is produced using diffusion with different time steps $s \in \{1, \ldots, T\}$, averaged over the `vertebral` dataset. (a) shows the analytical distribution computed by placing Gaussian distributions of different variances at each point in the dataset, (b) shows the inverse Gamma distribution with scale parameter value depending on the average distance to the k-nearest neighbours ($k = 32$), and (c) shows the difference between (a) and (b).

## D IMPLEMENTATION DETAILS

### D.1 DATASETS AND PREPROCESSING

**Datasets description** We show the results from our methods and baselines over multiple datasets from ADBench (Han et al., 2022) described in Table 1. There are 47 tabular datasets ranging from multiple different applications. There are also five datasets composed of extracted representations of images after the last average polling layer from a Resnet-18 (He et al., 2015) model pre-trained on ImageNet. Similarly, there are five datasets composed of extracted embedding of NLP tasks from BERT (Devlin et al., 2019). We also show results on VisA (Zou et al., 2022), which is a dataset composed of images of 12 different objects where the anomalies are various flaws on the objects.

**Training and test data configuration** For ADBench, the semi-supervised setting, we use half of the normal data in the training set, and the other half is in the test set with all the anomalies. For the unsupervised setting, we sample the whole dataset with replacement for the training data, while the test data is the whole dataset. This bootstrapping method allows us to test the variance over the training dataset for each method.

**Preprocessing** We standardize the input samples based on the mean and standard deviation calculated over the training data, to ensure consistency across the input values and mitigate the impact of potential outliers or scale variations. For VisA, 90% of the normal instances are making the training data, while the anomalies and the remaining 10% are in the test set. For CIFAR-10 and MNIST, One class is set as the anomaly while the others are part of the training data. 80% of the normal instances are in the training data while the remaining 20% and the anomalies are in the test data. For ADBench, CIFAR-10, MNIST-C, SVHN, and FashionMNIST are made up of one class for the normal sample, while the anomalies are the rest of the classes downsampled to make up 5% of the total data.

**On the importance of standardization for diffusion models** Throughout the course of our investigations, we discovered the critical importance of standardization. This is due to the fact that the incorporated Gaussian noise operates under the assumption that each feature is centered at zero with unit standard deviation. Consequently, implementing standard scaling facilitates the comprehensive

Table 1: Description of all datasets in ADBench

| 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 |

coverage of the anomaly detection space by the noise. This proved to be an essential component of the proposed anomaly detection method.

## D.2 ALGORITHM

---

**Algorithm 1** Training Process for parametric DTE

---

    **Parameters:** $T$ : maximum timestep, $\lambda$ : learning rate
    **Input:** Training data $\mathcal{D}$
 1: $\theta \leftarrow \theta_0$                                              ▷ Initialize weights of the model
 2: $\beta_0, \beta_1, ..., \beta_{T-1} \leftarrow \text{linear}(0, 0.01)$          ▷ Define the $\beta$ schedule for forward diffusion
 3: **for all** $t < T$ **do**
 4:      $\bar{\alpha}_t \leftarrow \prod_{s=1}^{t}(1 - \beta_s)$                            ▷ Compute the $\bar{\alpha}$
 5:      $\sigma_t \leftarrow \sqrt{1 - \bar{\alpha}_t}$                    ▷ Set standard deviation for each timestep
 6: **end for**
 7: **for** num_epochs **do**
 8:      **for all** $\mathbf{x}_0$ in $\mathcal{D}$ **do**
 9:          $t \sim \mathcal{U}(0, T-1)$                       ▷ Sample timestep $t$ uniformly
10:          $\epsilon \sim \mathcal{N}(0, 1)$                        ▷ Sample standard Gaussian noise
11:          $\mathbf{x}_t \leftarrow \mathbf{x}_0 + \sigma_t \epsilon$              ▷ Compute noisy sample of $x$ at timestep $t$
12:          $\mathcal{L} \leftarrow \text{loss}(f_\theta(\mathbf{x}_t))$    ▷ Equation (8) for inverse Gamma or Equation (9) for categorical
13:          $\theta \leftarrow \theta - \lambda \nabla_\theta \mathcal{L}$                     ▷ Update model parameters
14:      **end for**
15: **end for**

---

## D.3 MODEL ARCHITECTURE AND HYPERPARAMETERS

We first found the hyperparameters using different training splits for the semi-supervised setting on the shuttle and thyroid datasets (network architecture, maximum timestep, batch size, number of epochs). We then tuned some of them over all the datasets using different training seeds than the ones used for the final results (number of bins and learning rate). This is the case for the diffusion methods and the normalizing flow method. For the other baselines, we picked the set of hyperparameters from the original papers that provided the best results over the whole benchmark.

**DTE** For the non-parametric DTE, the score is calculated based on the approximate posterior distribution in Equation (6) with $k = 5$ for the semi-supervised setting and $k = 32$ for the unsupervised setting. The anomaly score is the mean of the posterior to avoid having an anomaly score that is restricted by the maximum variance using the mode. The be consistent, we selected the same $k$ for the kNN baseline.

For the DTE parametric approach, we employ a multi-layer perceptron (MLP) neural network. We use a common architecture and set of hyperparameters across all datasets. When training on images, we used a ResNet-50 architecture.

For the categorical model, we found that using the mean over each output probability bin provided the best results. That is, the anomaly score for each individual $\mathbf{x}$ is computed as follow:

$$score = f_\theta(\mathbf{x}) \begin{bmatrix} 0 \\ 1 \\ 2 \\ \vdots \\ B-1 \end{bmatrix} \tag{9}$$

where $B$ is the number of bins and $f_\theta(\mathbf{x}_t)$ is the output probability vector of the network using a softmax, which is an $N \times B$ matrix, where the sum across each row equates to one and $N$ is the batch size. The score for each instance will be a value between 0 and $B - 1$. The higher the score is, the more anomalous an instance is.

Employing the mode as a measurement metric proved suboptimal given the disproportionate representation of the first bin, a pattern that remained consistent even among anomalous instances.

Consequently, it was observed that while the probabilities could be diffusely distributed across the remaining bins, the mode predominantly remained in the first bin. In contrast, utilizing the mean allowed us to effectively account for this distribution characteristic, enabling an inclusive weighting scheme across all bins. Additionally, the mean offered a continuous scoring system as opposed to the integer values provided by the mode, thereby affording a more nuanced understanding of the anomalous data.

Table 2: Hyperparameters for parametric DTE model

| Hyperparameter | Value |
| --- | --- |
| Hidden layer sizes | [256, 512, 256] |
| Activation function | ReLU |
| Optimizer | Adam |
| Learning rate | 0.0001 |
| Dropout | 0.5 |
| Batch size | 64 |
| Number of epochs | 400 |
| Maximum timestep | 300 |
| Number of bins | 7 |

**DDPM** For the DDPM model, we used a modified ResNet for tabular data (Gorishniy et al., 2021) with added time embedding before each block, inspired by the work done for TabDDPM (Kotelnikov et al., 2022). Recognizing that learning noise at each timestep presents a considerably complex task, the necessity for a more sophisticated model than a simple MLP became evident to optimize the efficacy of our method. Furthermore, the lack of research on diffusion models for tabular data has constrained our ability to apply a model of comparable strength to the U-net model typically used for images, to our benchmark datasets. This presents an interesting direction for further research, with the potential to significantly enhance the performance of machine learning models on tabular datasets. In contrast to prior work (Wyatt et al., 2022; Wolleb et al., 2022), we do not add noise to the data point before reconstructing it as we found that it leads to overall slightly better results. This is a minor change, one intuition for the boost of performance could be that adding noise can modify the images toward anomalous data, thus increasing the amount of false positives.

Table 3: Hyperparameters for DDPM model

| Hyperparameter | Value |
| --- | --- |
| Number of blocks | 3 |
| Main layer size | 128 |
| Hidden layer size | 256 |
| Time embedding dimensions | 256 |
| Optimizer | Adam |
| Learning rate | 0.0001 |
| Dropout layer 1 | 0.4 |
| Dropout layer 2 | 0.1 |
| Batch size | 64 |
| Number of epochs | 400 |
| Maximum timestep | 1000 |
| Reconstruction timestep | 250 |

**Normalizing Flows Baseline** We compare our diffusion methods with a normalizing flows baseline that uses planar flows (Rezende & Mohamed, 2015). Normalizing flows allow to compute the exact likelihoods of data point, which allow to easily assign anomaly scores. Once trained, the model can estimate the density of any data point in the input space. This is done by passing the data point through the inverse of the learned transformation and then computing the density of the transformed point under the simple target distribution. The density of the original point under the complex data distribution can be computed from this using the change-of-variables formula.

Table 4: Hyperparameters for PlanarFlow model

| Hyperparameter | Value |
|---|---|
| Number of transformations | 10 |
| Optimizer | Adam |
| Learning rate | 0.002 |
| Batch size | 64 |
| Number of epochs | 200 |

## D.4 COMPUTE

The total amount of compute required to reproduce our experiments with five seeds, including all of the baselines and the proposed DTE model amounts to 473 GPU-hours for the unsupervised setting and 225 GPU-hours for the semi-supervised setting on an RTX8000 GPU with 48 gigabytes of memory for running the ADBench datasets.

Figure 12 shows the training and inference times averaged over all datasets in ADBench over five seeds for all methods discussed in Section 4. As expected, deep learning-based methods have significantly higher training times compared to classical methods but comparable inference times. In particular, the inference time for the parametric DTEs is orders of magnitude lower than all other methods. The non-parametric variant of DTE has no training phase, so we show the inference time in both plots.



(a) Training time



(b) Inference time

Figure 12: Mean training and inference time on the 57 datasets from ADBench over five different seeds for the semi-supervised setting using normal samples only for training. Colour scheme: red (diffusion-based), green (deep learning methods), blue (classical methods).

# E  CHOICE OF REPRESENTATION FOR IMAGES

In this section, we compare the effect of choice of representation on the performance of diffusion-based anomaly detection techniques. Three choices considered are 1) pixel space representation, 2) self-supervised embedding, and 3) embedding produced by a classifier. Results for three image datasets are reported in Table 5. The datasets and preprocessing are described in Appendix D.1 and the full results are in Appendix F. As expected, using pre-trained embeddings leads to better results than pixel space for all methods considered. Tables 9 to 12 report other experiments that lead to a similar conclusion.

In particular, using self-supervised embedding for CIFAR-10, significantly improved the anomaly detection performance as the pre-training was done on CIFAR-10 itself. Note that all the other pre-training were supervised classification using ResNet-34 on ImageNet and not directly on the datasets. Overall, pre-training improves the results for all methods and all datasets with the exception of kNN and the non-parametric DTE (DTE-NP) on MNIST. This result can be attributed to the simplicity of the MNIST dataset when adapted to anomaly detection tasks. As a reminder, DTE-NP is equivalent to kNN, but corresponds to the variation that uses the mean distance of the k-nearest neighbours instead of the distance to the kth-nearest neighbour.

Zou et al. (2022) highlighted the advantages of tailoring specialized self-supervised learning techniques to specific datasets, exemplified by their method for VisA. As our methods are not explicitly designed for these datasets, our results for all diffusion-based methods reported here lag behind those of methods specialized to this dataset. In particular, VisA dataset contains images that are quite similar with the exception of highly localized anomalies.

Table 5: Average AUC ROC and standard deviations for the different subsets of each dataset, average across 5 runs, semi-supervised setting using different pre-training algorithms.

|  | DTE-NP | DTE-C | DDPM | kNN |
|---|---|---|---|---|
| VisA, supervised ImageNet pre-training | 83.63(10.50) | 81.07(11.01) | 80.47(12.47) | 83.26(10.64) |
| VisA, VicReg ImageNet pre-training | 83.36(12.44) | 81.89(12.26) | 83.14(13.76) | 83.68(13.54 |
| VisA, no pre-training | 75.96(10.54) | 64.53(19.61) | 57.85(21.74) | 75.40(9.85) |
| CIFAR10, supervised ImageNet pre-training | 53.91(7.16) | 52.57(5.53) | 52.96(7.05) | 54.42(7.56) |
| CIFAR10, VicReg pre-training | 80.92(10.81) | 63.36(11.92) | 54.22(10.26) | 79.01(11.53) |
| CIFAR10, no pre-training | 51.53(14.81) | 50.25(3.34) | 50.50(7.67) | 51.64(14.90) |
| MNIST, supervised ImageNet pre-training | 78.07(12.48) | 64.34(11.93) | 60.62(10.26) | 76.86(11.54) |
| MNIST, no pre-training | 81.94(16.46) | 49.02(16.51) | 51.29(18.85) | 84.14(15.60) |

# F  FULL RESULTS

We provide the full table of results corresponding to the AUC ROC box-plots in Section 4. We report additional metrics including F1 score and area under the precision-recall curve (AUC PR) along with the corresponding box-plots. All results are shown averaged across five seeds along with standard deviations in brackets for all 57 datasets in ADBench. In the subsequent tables, DTE-NP refers to the non-parametric DTE estimator, DTE-IG refers to the parametric inverse Gamma model, and DTE-C refers to the parametric categorical model. Tables 9 to 12 show the results for three methods when using pre-trained embeddings on CIFAR-10 and SVHN compared to trained directly on the images, as it is set up in ADBench. The difference with Table 5 is that instead of having one class as an anomaly, here we have one class as normal while the rest of the classes are downsampled to produce the anomalies.

### F.1 SEMI-SUPERVISED SETTING
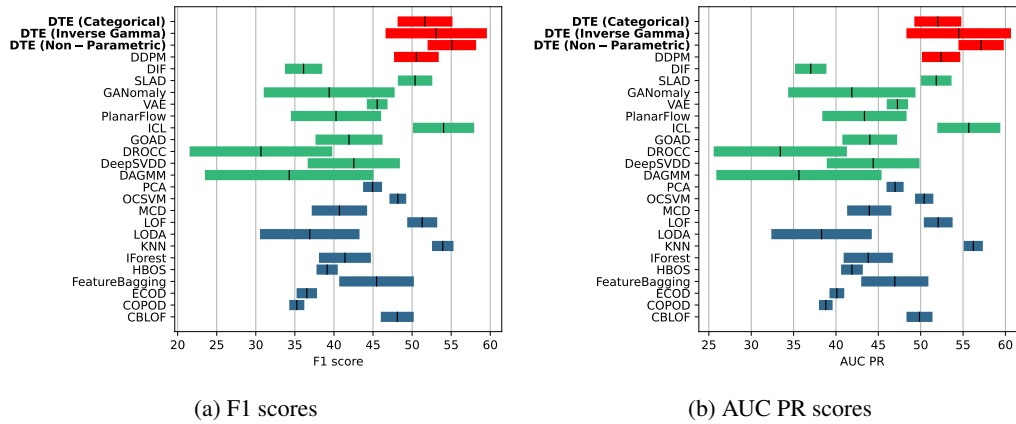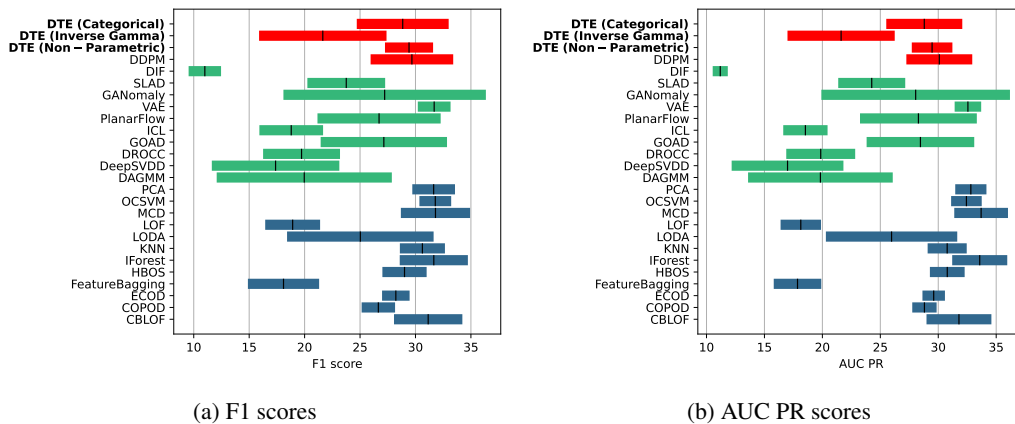


(a) F1 scores

(b) AUC PR scores

Figure 13: F1 score and AUC PR means and standard deviations on the 57 datasets from ADBench over five different seeds for the semi-supervised setting using normal samples only for training. Colour scheme: red (diffusion-based), green (deep learning methods), blue (classical methods).

Table 6: Average AUC ROC and standard deviations for 5 runs of the VisA dataset, semi-supervised setting using embeddings of supervised ResNet-34 pre-trained on ImageNet with the same training split.

|  | **DTE-NP** | **DTE-C** | DDPM | kNN |
|---|---|---|---|---|
| candle | 90.88(0.0) | 89.26(3.37) | 87.37(0.26) | 90.76(0.0) |
| capsules | 62.77(0.0) | 56.04(3.51) | 66.65(0.6) | 62.67(0.0) |
| cashew | 93.7(0.0) | 87.17(2.95) | 89.69(0.25) | 93.24(0.0) |
| chewinggum | 93.88(0.0) | 94.69(1.53) | 92.55(0.26) | 93.52(0.0) |
| fryum | 87.38(0.0) | 81.25(3.91) | 85.84(0.17) | 87.68(0.0) |
| macaroni1 | 70.27(0.0) | 71.9(4.94) | 64.33(0.69) | 69.34(0.0) |
| macaroni2 | 67.65(0.0) | 66.77(2.36) | 51.6(0.64) | 66.35(0.0) |
| pcb1 | 90.14(0.0) | 85.49(1.72) | 94.38(0.37) | 90.67(0.0) |
| pcb2 | 88.88(0.0) | 81.32(2.76) | 83.24(0.4) | 87.77(0.0) |
| pcb3 | 80.83(0.0) | 81.72(1.98) | 79.61(0.2) | 81.25(0.0) |
| pcb4 | 93.77(0.0) | 91.05(2.46) | 88.7(0.56) | 93.07(0.0) |
| pipe fryum | 83.24(0.0) | 86.19(2.38) | 81.72(0.36) | 82.86(0.0) |
| mean | 83.63(10.50) | 81.07(11.01) | 80.47(12.47) | 83.26(10.64) |

Table 7: Average AUC ROC and standard deviations for 5 runs of the VisA dataset, semi-supervised setting using embeddings of VicReg pre-trained on ImageNet with the same training split.

|  | DTE-NP | DTE-C | DDPM | kNN |
|---|---|---|---|---|
| candle | 82.97(0.0) | 84.72(0.8) | 85.27(0.07) | 85.44(0.0) |
| capsules | 65.63(0.0) | 68.24(1.08) | 69.01(0.6) | 65.92(0.0) |
| cashew | 90.7(0.0) | 82.75(5.28) | 90.26(0.43) | 90.74(0.0) |
| chewinggum | 97.98(0.0) | 97.78(0.12) | 97.78(0.08) | 98.0(0.0) |
| fryum | 88.88(0.0) | 79.62(2.69) | 89.13(0.23) | 88.82(0.0) |
| macaroni1 | 70.03(0.0) | 68.17(1.19) | 64.09(0.17) | 69.52(0.0) |
| macaroni2 | 52.06(0.0) | 55.81(2.06) | 51.93(0.27) | 52.16(0.0) |
| pcb1 | 93.02(0.0) | 91.07(0.5) | 93.19(0.08) | 93.22(0.0) |
| pcb2 | 85.7(0.0) | 83.77(0.96) | 83.68(0.17) | 85.69(0.0) |
| pcb3 | 83.26(0.0) | 81.57(0.75) | 82.03(0.13) | 83.05(0.0) |
| pcb4 | 98.6(0.0) | 98.21(0.35) | 98.35(0.03) | 98.66(0.0) |
| pipe fryum | 91.44(0.0) | 90.98(1.29) | 93.05(0.05) | 92.96(0.0) |
| mean | 83.36(12.44) | 81.89(12.26) | 83.14(13.76) | 83.68(13.54) |

Table 8: Average AUC ROC and standard deviations for 5 runs of the VisA dataset, semi-supervised setting using the images directly with the same training split.

|  | DTE-NP | DTE-C | DDPM | kNN |
|---|---|---|---|---|
| candle | 77.48(0.0) | 83.03(4.42) | 51.96(6.2) | 77.38(0.0) |
| capsules | 63.75(0.0) | 72.61(7.91) | 33.19(0.37) | 68.02(0.0) |
| cashew | 90.14(0.0) | 79.5(27.78) | 96.26(0.56) | 93.32(0.0) |
| chewinggum | 66.92(0.0) | 56.99(4.66) | 68.82(1.02) | 65.66(0.0) |
| fryum | 74.32(0.0) | 77.28(10.99) | 25.24(1.22) | 74.5(0.0) |
| macaroni1 | 68.67(0.0) | 54.52(20.66) | 74.7(1.14) | 70.11(0.0) |
| macaroni2 | 74.04(0.0) | 54.48(8.11) | 37.04(0.48) | 77.02(0.0) |
| pcb1 | 83.59(0.0) | 51.53(16.31) | 72.02(0.97) | 80.53(0.0) |
| pcb2 | 87.4(0.0) | 74.04(15.82) | 77.56(0.55) | 78.87(0.0) |
| pcb3 | 71.75(0.0) | 40.86(8.75) | 68.11(2.12) | 66.03(0.0) |
| pcb4 | 94.5(0.0) | 73.61(11.82) | 28.46(2.18) | 92.94(0.0) |
| pipe fryum | 58.96(0.0) | 56.04(19.1) | 60.95(5.78) | 60.42(0.0) |
| mean | 75.96(10.54) | 64.53(19.61) | 57.85(21.74) | 75.40(9.85) |

Table 9: Mean AUC ROC and standard deviation over 5 seeds for different methods trained on the images directly versus trained on embeddings generated by a pre-trained ResNet-18 on ImageNet for the unsupervised setting on the CIFAR-10 dataset.

|  | DDPM | DTE-C | kNN |
|---|---|---|---|
| Images | 54.72(4.55) | 48.95(8.33) | 57.45(1.55) |
| Embeddings | 66.34(0.14) | 62.87(1.57) | 66.17(0.33) |

Table 10: Mean AUC ROC and standard deviation over 5 seeds for different methods trained on the images directly versus trained on embeddings generated by a pre-trained ResNet-18 on ImageNet for the unsupervised setting on the SVHN dataset.

|  | DDPM | DTE-C | kNN |
|---|---|---|---|
| Images | 54.97(2.41) | 49.07(3.06) | 56.29(1.22) |
| Embeddings | 61.48(0.24) | 59.96(1.24) | 61.17(0.28) |

Table 11: Mean AUC ROC and standard deviation over 5 seeds for different methods trained on the images directly versus trained on embeddings generated by a pre-trained ResNet-18 on ImageNet for the semi-supervised setting on the CIFAR-10 dataset.

|  | DDPM | DTE-C | kNN |
|---|---|---|---|
| Images | 55.96(4.69) | 52.66(6.28) | 59.10(1.80) |
| Embeddings | 67.91(0.13) | 68.53(1.59) | 67.53(0.0) |

Table 12: Mean AUC ROC and standard deviation over 5 seeds for different methods trained on the images directly versus trained on embeddings generated by a pre-trained ResNet-18 on ImageNet for the semi-supervised setting on the SVHN dataset.

|  | DDPM | DTE-C | kNN |
|---|---|---|---|
| Images | 57.28(2.75) | 48.78(4.23) | 55.92(1.17) |
| Embeddings | 61.37(0.08) | 62.91(1.1) | 61.69(0.0) |

## F.2 UNSUPERVISED SETTING



(a) F1 scores

(b) AUC PR scores

Figure 14: F1 score and AUC PR means and standard deviations on the 57 datasets from ADBench over five different seeds for the unsupervised setting with bootstrapped training instances. Colour scheme: red (diffusion-based), green (deep learning methods), blue (classical methods).

Table 13: Average AUC ROC and standard deviations over five seeds for the semi-supervised setting on ADBench.

Table 14: Average F1 score and standard deviations over five seeds for the semi-supervised setting on ADBench.

| | CBLOF | COPOD | ECOD | FeatureBagging | HBOS | IForest | KNN | LODA | LOF | MCD | OCSVM | PCA | DAGMM | DeepSVDD | DROCC | GOAD | ICL | PlanarFlow | VAE | GANomaly | SLAD | DIF | DDPM | DTE-NP | DTE-IG | DTE-C |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 53.69(0.15) | 49.5 alot | 6.74(0.08) | 4.58(0.0) | 4.44(0.0) | 8.93(0.57) | 7.43(0.0) | 4.2(0.26) | 5.9(0.0) | 6.6(1.58) | 8.1(6.0.0) | 3.41(0.14) | 7.29(0.0) | 7.63(0.0) | 5.98(1.67) | 5.17(0.92) | 0.0(0.0) | 5.73(1.45) | 4.91(0.56) | 3.83(0.72) | 7.63(0.0) | 9.35(2.27) | 5.32(0.11) | 3.91(0.0) | 6.76(0.19) | 5.82(0.07) |

*(Full per-dataset F1 score and standard deviation values for datasets aloi, amazon, annthyroid, backdoor, breastw, campaign, cardio, cardiotocography, celeba, census, cover, donors, fault, fraud, glass, hepatitis, http, imdb, internetads, ionosphere, landsat, letter, lymphography, magic.gamma, mammography, mnist, musk, optdigits, pageblocks, pendigits, pima, satellite, satimage-2, shuttle, skin, smtp, spambase, speech, stamps, thyroid, vertebral, vowels, waveform, wbc, wdbc, wilt, wine, wpbc, yeast, yelp, MNIST-C, FashionMNIST, CIFAR10, SVHN, MVTec-AD, 20news, agnews across all listed methods.)*

Table 15: Average AUC PR and standard deviations over five seeds for the semi-supervised setting on ADBench.

Table 16: Average AUC ROC and standard deviations over five seeds for the unsupervised setting on ADBench.

Table 17: Average F1 score and standard deviations over five seeds for the unsupervised setting on ADBench.

Table 18: Average AUC PR and standard deviations over five seeds for the unsupervised setting on ADBench.

| | CBLOF | COPOD | ECOD | FeatureBagging | HBOS | IForest | KNN | LODA | LOF | MCD | OCSVM | PCA | DAGMM | DeepSVDD | DROCC | GOAD | ICL | PlanarFlow | VAE | GANomaly | SLAD | DIF | DDPM | DTE-NP | DTE-IG | DTE-C |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| aloi | | | | | | | | | | | | | | | | | | | | | | | | | | |
| amazon | | | | | | | | | | | | | | | | | | | | | | | | | | |
| annthyroid | | | | | | | | | | | | | | | | | | | | | | | | | | |
| backdoor | | | | | | | | | | | | | | | | | | | | | | | | | | |
| breastw | | | | | | | | | | | | | | | | | | | | | | | | | | |
| campaign | | | | | | | | | | | | | | | | | | | | | | | | | | |
| cardio | | | | | | | | | | | | | | | | | | | | | | | | | | |
| cardiotocography | | | | | | | | | | | | | | | | | | | | | | | | | | |
| celeba | | | | | | | | | | | | | | | | | | | | | | | | | | |
| census | | | | | | | | | | | | | | | | | | | | | | | | | | |
| cover | | | | | | | | | | | | | | | | | | | | | | | | | | |
| donors | | | | | | | | | | | | | | | | | | | | | | | | | | |
| fault | | | | | | | | | | | | | | | | | | | | | | | | | | |
| fraud | | | | | | | | | | | | | | | | | | | | | | | | | | |
| glass | | | | | | | | | | | | | | | | | | | | | | | | | | |
| hepatitis | | | | | | | | | | | | | | | | | | | | | | | | | | |
| http | | | | | | | | | | | | | | | | | | | | | | | | | | |
| imdb | | | | | | | | | | | | | | | | | | | | | | | | | | |
| internetads | | | | | | | | | | | | | | | | | | | | | | | | | | |
| ionosphere | | | | | | | | | | | | | | | | | | | | | | | | | | |
| landsat | | | | | | | | | | | | | | | | | | | | | | | | | | |
| letter | | | | | | | | | | | | | | | | | | | | | | | | | | |
| lymphography | | | | | | | | | | | | | | | | | | | | | | | | | | |
| magic.gamma | | | | | | | | | | | | | | | | | | | | | | | | | | |
| mammography | | | | | | | | | | | | | | | | | | | | | | | | | | |
| mnist | | | | | | | | | | | | | | | | | | | | | | | | | | |
| musk | | | | | | | | | | | | | | | | | | | | | | | | | | |
| optdigits | | | | | | | | | | | | | | | | | | | | | | | | | | |
| pageblocks | | | | | | | | | | | | | | | | | | | | | | | | | | |
| pendigits | | | | | | | | | | | | | | | | | | | | | | | | | | |
| pima | | | | | | | | | | | | | | | | | | | | | | | | | | |
| satellite | | | | | | | | | | | | | | | | | | | | | | | | | | |
| satimage-2 | | | | | | | | | | | | | | | | | | | | | | | | | | |
| shuttle | | | | | | | | | | | | | | | | | | | | | | | | | | |
| skin | | | | | | | | | | | | | | | | | | | | | | | | | | |
| smtp | | | | | | | | | | | | | | | | | | | | | | | | | | |
| spambase | | | | | | | | | | | | | | | | | | | | | | | | | | |
| speech | | | | | | | | | | | | | | | | | | | | | | | | | | |
| stamps | | | | | | | | | | | | | | | | | | | | | | | | | | |
| thyroid | | | | | | | | | | | | | | | | | | | | | | | | | | |
| vertebral | | | | | | | | | | | | | | | | | | | | | | | | | | |
| vowels | | | | | | | | | | | | | | | | | | | | | | | | | | |
| waveform | | | | | | | | | | | | | | | | | | | | | | | | | | |
| wbc | | | | | | | | | | | | | | | | | | | | | | | | | | |
| wine | | | | | | | | | | | | | | | | | | | | | | | | | | |
| wpbc | | | | | | | | | | | | | | | | | | | | | | | | | | |
| yeast | | | | | | | | | | | | | | | | | | | | | | | | | | |
| yelp | | | | | | | | | | | | | | | | | | | | | | | | | | |
| MNIST-C | | | | | | | | | | | | | | | | | | | | | | | | | | |
| CIFAR10 | | | | | | | | | | | | | | | | | | | | | | | | | | |
| FashionMNIST | | | | | | | | | | | | | | | | | | | | | | | | | | |
| SVHN | | | | | | | | | | | | | | | | | | | | | | | | | | |
| MVTec-AD | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 20news | | | | | | | | | | | | | | | | | | | | | | | | | | |
| agnews | | | | | | | | | | | | | | | | | | | | | | | | | | |