

Enclosing Prototypical Variational Autoencoder for Explainable Out-of-Distribution Detection

Conrad Orglmeister¹, Erik Bochinski¹, Volker Eiselein¹, and Elvira Fleig²

¹ Digitale Schiene Deutschland, DB InfraGO AG, Berlin, Germany

² {conrad.orglmeister,erik.bochinski,volker.eiselein}@deutschebahn.com
Communication Systems Group, Technische Universität Berlin, Berlin, Germany
fleig@tu-berlin.de

Abstract. Understanding the decision-making and trusting the reliability of Deep Machine Learning Models is crucial for adopting such methods to safety-relevant applications. We extend self-explainable Prototypical Variational models with autoencoder-based out-of-distribution (OOD) detection: A Variational Autoencoder is applied to learn a meaningful latent space which can be used for distance-based classification, likelihood estimation for OOD detection, and reconstruction. The In-Distribution (ID) region is defined by a Gaussian mixture distribution with learned prototypes representing the center of each mode. Furthermore, a novel restriction loss is introduced that promotes a compact ID region in the latent space without collapsing it into single points. The reconstructive capabilities of the Autoencoder ensure the explainability of the prototypes and the ID region of the classifier, further aiding the discrimination of OOD samples. Extensive evaluations on common OOD detection benchmarks as well as a large-scale dataset from a real-world railway application demonstrate the usefulness of the approach, outperforming previous methods.

Keywords: Out-of-Distribution detection · Explainable AI · Prototypical Variational Autoencoder · Reconstruction · Distance.

1 Introduction

Modern Artificial Neural Networks (ANNs) achieve remarkable results in recognizing patterns. However, due to their complexity and black-box character, their failures are hard to identify [13] which limits their use in safety-critical environments. Additionally, certain common training schemes encourage overconfidence [8]. If Out-of-Distribution (OOD) samples from other distributions than the In-Distribution (ID) training set are encountered in classification tasks, this issue persists. Encountering such samples is often unavoidable in real-world applications, especially when operating in an open world as autonomous transportation systems do. Therefore, OOD detection has arisen as the task of identifying instances not belonging to the training data distribution [25] which often means finding the label distribution but also extends to identifying when the

model might be unable to assess its input reliably. Anomaly detection, Open-Set-Recognition (OSR), and Uncertainty Estimation are closely related to OOD detection and methods can often be applied to the other settings as well [25]. Most importantly, OSR requires explicitly classifying closed-world samples and detecting unknown classes from the open world [25].

Many OOD detection methods rely on post-hoc analysis of output or intermediate features from pre-trained classifiers but models trained solely for discrimination of ID categories may lack relevant features for OOD detection which limits the general usage of such approaches. Integration of OOD detection into the classification framework is thus desirable, rather than applying it afterwards.

In this work, we extend the Prototypical Variational Autoencoder (ProtoVAE) [6] to OOD detection. Instead of the aforementioned post-analysis of application-specific pre-learned features for OOD detection, the feature space is designed to learn to distinguish unknown inputs from the beginning. This is done by estimating the training distribution, learning representations through reconstruction, and designing a distance-based latent space to quantify dissimilarity to ID clusters while also leveraging label information yielding promising results. Additionally, a restriction force is implemented to shape the latent ID region while reconstruction errors are used to identify remaining OOD samples mapped into this region as introduced in [27].

This work proposes the principle of an *enclosing restriction* to decouple the previous trade-off between compression/estimation of the ID region and reconstructive quality to recover the input rather than just reconstruct features, thus alleviating Autoencoder (AE)-based OOD detection by constraining the ID region in the latent space without collapsing it into one point. To enhance the reconstructive power further, Learned Perceptual Image Patch Similarity (LPIPS) – a perceptual metric – is integrated into the framework for the reconstruction loss and OOD score. The generative and reconstructive abilities of the Variational Autoencoder (VAE) framework enable the provision of additional information and explanation about extracted properties of the data distribution and certain samples, rendering the classification and OOD detection transparent. The method is compared to state-of-the-art approaches using the OpenOOD [24] and a custom railway benchmark.

2 Related Work

A ProtoVAE architecture was presented by Gautam *et al.* [6] as a self-explainable model. Distance-based classification makes the decision more transparent and class distributions are divided into clusters. The ability to decode embeddings including prototypes fosters transparency *w.r.t.* learned data distribution. In this work, modifications enable more direct distance-based classification and enforce an *enclosed* ID region making it ideal for OOD detection.

Yang *et al.* [24] categorize OOD detection methods applied post-hoc, requiring training, Outlier Exposure, pre-processing, or data augmentation. Yang *et al.* [25] also distinguish approaches based on outputs of a classifier (classification-

based), modeling the data distribution (density-based/generative), relying on distances in feature space (distance-based), and reconstructing the input measuring a reconstruction error (reconstruction-based). The approach of this work can be considered reconstruction-, distance- and density-based. Maximum Softmax Probability (MSP) as a baseline OOD score was examined by Hendrycks and Gimpel [11]. Hendrycks *et al.* [10] use the maximum logit as a score (post-hoc). Sun *et al.* [20] propose thresholding activations of the penultimate layer thus eliminating overconfidence caused by extreme activations. Wang *et al.* [22] design a virtual logit based on the smallest principle components. Gal and Ghahramani [5] apply *Monte-Carlo dropout* during test time and Lakshminarayanan *et al.* [13] train an ensemble of ANNs. Hendrycks *et al.* [12] propose a training-time augmentation based on fractals (*PixMix*).

Nalisnick *et al.* [15] find that density estimates might assign higher likelihoods to OOD than to ID data. Xiao *et al.* [23] tackle this by retraining a VAE-encoder for a specific test-sample measuring a likelihood discrepancy. Sun *et al.* [19] design a VAE with one Gaussian distribution per class. In contrast to this work, no perceptual metric, distance-based classification, or restriction-scheme for the ID region is used. Moreover, a custom probability is defined for a sample being part of a class distribution. There is a fixed threshold for the latter in contrast to the flexible OOD score fusion used in this work without a fixed threshold for one of the scores alone. *ARPL* [2] generates near-OOD samples for learning adversarial reciprocal points representing individual negative classes.

Reconstructive OOD detection often involves elaborate schemes[3,16,1,27,7] as the reconstruction error alone often cannot separate OOD from ID data [3]. Existing approaches combine reconstruction error with Mahalanobis distance [3], improve ID reconstruction with a deformation transformation[1] or use multiple reconstruction errors [16,7]. In [27], the latent space region of an AE to which ID samples are encoded (*ID region*) is estimated by restricting ID data within the latent space. For OOD samples mapped into this region, the reconstruction error will be higher [27]. In contrast, in this work, an *enclosing restriction* supports the trade-off between reliable estimation of the ID region and reconstruction quality.

Distance-based OOD detection involves Mahalanobis distance[14] and k-Nearest Neighbor (KNN) distance for pre-trained features. Requiring training, Deep SVDD [17] maps ID data into a hypersphere, and SIREN [4] discriminatively shapes representations using prototypes but not reconstruction.

3 Methodology

We introduce the Prototypical Direct-Distance-Classifier VAE (ProtoDistVAE) for explainable OOD detection which extends the ProtoVAE from [6] and further incorporates the principle of AE-based OOD detection from [27]. Following [27], if an AE reconstructs every ID sample sufficiently well and the ID region \mathcal{T}_{ID} can be estimated precisely, a sample can be concluded to be ID by fulfilling two conditions:

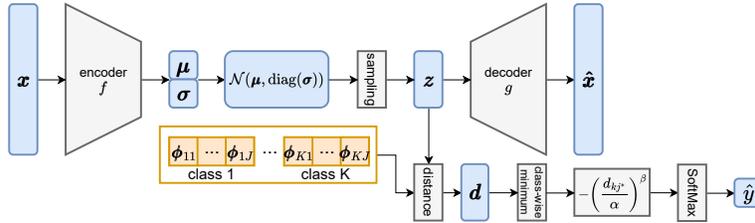


Fig. 1: ProtoDistVAE architecture: The input \mathbf{x} is encoded into a latent Gaussian distribution from which a sample \mathbf{z} is drawn and reconstructed to obtain $\hat{\mathbf{x}}$. Then, in the framework of generalized Gaussians, the SoftMax function returns the predicted probabilities and class estimate \hat{y} for the distances to all prototypes.

1. An ID sample is embedded into \mathcal{T}_{ID} (by definition).
2. An ID sample exhibits a small reconstruction error.

Under the given assumptions, OOD samples should never fulfill both conditions.

Our aim is to model a distribution of data that is representative for a set of prototypes. This means that different classes or parts of classes can be assigned to different sub-distributions during training, thus potentially increasing data diversity and simplifying the OOD detection. A distance metric space is learned where similar samples are in close proximity to each other. Similar to [6], we use an encoder f_ψ , a decoder g_θ and prototypes $\phi_{kj} \in \mathbb{R}^L$ in an end-to-end trainable fashion (see Figure 1). The rows of matrix $\Phi_k \in \mathbb{R}^{J \times L}$ describe the J prototype vectors of class $k \in K$ classes.

Given a training dataset $\mathcal{D} = \{(\mathbf{x}^1, (x^1, y^1)), \dots, (\mathbf{x}^N, (x^N, y^N))\}$ with N labeled samples, the input \mathbf{x}^i itself yields the target variables for reconstruction and a class label y^i . The model is trained as a VAE learning a Gaussian mixture distribution where the encoder embeds the input \mathbf{x}^i to a posterior Gaussian distribution $p(\mathbf{z}|\mathbf{x}^i) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}^i, \text{diag}((\boldsymbol{\sigma}^i)^2))$ in the latent domain. During training, a latent representation \mathbf{z}^i is sampled whereas during inference, the mean value is used for the latent representation the decoder processes into the image space reconstruction $\hat{\mathbf{x}}^i$.

For classification, the Euclidean distances of the latent variable to all prototypes are computed (Equation (1)) and the minimum distance of each class yields the closest prototype. It is important to minimize the distance of an embedding to only one prototype distribution during training. The distances are transformed into logits by the generalized Gaussian distribution for enclosing restriction and are fed into a SoftMax function to obtain a purely distance-based, latent space classification without a learnable classifier.

$$d(\mathbf{z}^i, \phi_{kj}) = d_{kj}^i = \|\mathbf{z}^i - \phi_{kj}\|_2 \quad (1)$$

$$P_\psi(y = k|\mathbf{x}^i) = \frac{\exp(l_k^i)}{\sum_{k'=1}^K \exp(l_{k'}^i)}, \quad l_{k'}^i = -\left(\frac{d_{k'j^*(k')}}{\alpha}\right)^\beta \quad (2)$$

$$j^*(k) = \underset{j}{\operatorname{argmin}}(d_{kj}) \quad (3)$$

The original ProtoVAE architecture uses a linear classifier and distance-based similarity scores [6]. Similarity scores exhibit large gradients for embeddings close to a prototype which potentially leads to embeddings collapsing into the respective prototype position, and thus to degradation of reconstruction quality when different embeddings are not encoded differently. As a remedy, ProtoDistVAE uses an enclosing restriction leading to weaker gradients close to prototypes. Embeddings shall be trapped in a certain *ID region*, but inside, the coding of embeddings shall be unconstrained. For this reason, generalized Gaussian distributions are used in the classification layer where α defines the width of the distribution and $\beta \geq 2$ controls the shape and "enclosedness" of the distribution.

In order to not distort the distance metric space, ProtoDistVAE uses distances more explicitly for classification. The linear classifier which essentially calculates a *sum* of distances is replaced by using only the minimum distances to prototypes per class. These are translated into logits l_k^i , by the framework of generalized Gaussians and probabilities using the SoftMax function (Equation (2)). Cross-entropy is then applied to the modified predicted probabilities. $j^*(k)$ is the nearest prototype within class k while \mathbf{d}^* is the minimum distances vector for every class. Thus, instead of a sum of distances to multiple prototypes, the distance to only one prototype is minimized for a specific embedding.

The overall loss consists of a sum of four terms: The cross-entropy loss $\mathcal{L}'_{\text{cls}}$ shown in Equation (4) provides label information to enable the network to extract useful embeddings for discrimination and minimize the embedding distance to prototypes of the correct class. Each class is modeled by a mixture of J normal distributions centered around the respective class prototypes for VAE-like distribution estimation and Kullback-Leibler divergence (KL divergence) *w.r.t.* the nearest prototype distribution of the correct class is computed to obtain the loss \mathcal{L}'_{KL} (Equation (5)). The reconstruction loss aims to recover the input samples [6] by separating groups of samples near each other for a better reconstruction. We use the LPIPS metric [26] for this task as it gives a more robust similarity between images than traditional metrics as e.g. mean squared error (MSE) by using a calibrated pre-trained network aligned towards human perception [26].

In order to prevent the collapse of prototypes of a class, an orthonormalization loss $\mathcal{L}'_{\text{orth}}$ (Equation (7)) is used to encourage prototypes within a class (after subtracting their mean $\bar{\phi}_k$) to be orthonormal to each other [6]. It is defined as the average of the class-wise Frobenius norms $\|\cdot\|_F$.

$$\mathcal{L}'_{\text{cls}}(\boldsymbol{\psi}, \boldsymbol{\Phi}; \mathbf{x}^i, k) = -\log P_\psi(y = k | \mathbf{x}^i) \quad (4)$$

$$\mathcal{L}'_{\text{KL}}(\boldsymbol{\psi}, \boldsymbol{\Phi}_k; \mathbf{x}^i, k = y^i) = D_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu}^i, \text{diag}((\boldsymbol{\sigma}^i)^2)) \| \mathcal{N}(\boldsymbol{\phi}_{k_{j^*(k)}}, \mathbf{I}_L)) \quad (5)$$

$$\mathcal{L}'_{\text{rec}}(\boldsymbol{\psi}, \boldsymbol{\theta}; \mathbf{x}^i, \hat{\mathbf{x}}^i) = e_{\text{LPIPS}}(\mathbf{x}^i, \hat{\mathbf{x}}^i) \quad (6)$$

$$\mathcal{L}'_{\text{orth}}(\boldsymbol{\Phi}) = \frac{1}{K} \sum_{k=1}^K \|\tilde{\boldsymbol{\Phi}}_k \tilde{\boldsymbol{\Phi}}_k^T - \mathbf{I}_J\|_F^2, \quad \tilde{\boldsymbol{\Phi}}_k = (\boldsymbol{\phi}_{kj} - \bar{\boldsymbol{\phi}}_k)_{j=1..J} \quad (7)$$

In summary, ProtoDistVAE introduces LPIPS [26] as reconstruction loss and replaces the linear classifier layer as well as similarity scores by direct minimum distances and the framework of generalized Gaussians to implement an enclosing

restriction loss. The complete loss function is:

$$\mathcal{L} = w_{\text{cls}}\mathcal{L}_{\text{cls}} + w_{\text{KL}}\mathcal{L}_{\text{KL}} + w_{\text{rec}}\mathcal{L}_{\text{rec}} + w_{\text{orth}}\mathcal{L}_{\text{orth}} \quad (8)$$

For OOD detection, a distance-based OOD score and the LPIPS reconstruction error are merged. During experimentation, we found that the minimum distance to the next prototype can be improved by using the MSP score $\lambda_{\text{MSP}} = \max_k P_{\psi}(y = k|\mathbf{x}^i)$ in the ProtoDistVAE context which is the probability that an embedding belongs to the most likely generalized Gaussian under condition that it is ID. As ProtoDistVAE relies on distances for classification, MSP is also distance-based. Also the $\lambda_{\text{DistRatio}} = \sum_j d_{\hat{k}j} / (\sum_k \sum_j d_{kj})$ is applied where \hat{k} indicates the predicted class. We assume these scores perform better than the minimum distance because the class distribution in the latent space might be skewed and OOD samples are embedded between different class regions.

For fusion of scores, one distance score and one reconstruction error are normalized *w.r.t.* to their validation set distributions to make them comparable using a lower and upper percentile of the score distribution to obtain the normalized score $\tilde{\lambda}(\mathbf{x}) = (\lambda(\mathbf{x}) - \lambda_{\text{lower}}) / (\lambda_{\text{upper}} - \lambda_{\text{lower}})$. Both score types are combined into one score using L_2 or L_∞ norm: $\lambda_{L_p}(\mathbf{x}) = \|(\tilde{\lambda}_1(\mathbf{x}), \tilde{\lambda}_2(\mathbf{x}))^T\|_p$ where p denominates the degree. The L_∞ norm tends to reflect a hard decision (e.g. at least one score is above its threshold) and the L_2 norm a flexible decision (one score is too high or both together are rather high and therefore indicates an OOD sample). This type of fusion means that no probabilities need to be modeled explicitly and thus avoids any need for modeling assumptions.

4 Experimental Results

For numerical evaluation, we compare our approach to the state-of-the-art based on the OpenOOD benchmark [24] and a non-public dataset from the railway domain (DBS dataset). A general advantage of the proposed method is that it allows human insights into the training distribution and decision-making of the network by reconstructing samples, prototypes, and distances in the latent space which supports its usage in safety-critical domains.

General Experimental Setup The OpenOOD benchmark provides implementations of state-of-the-art approaches for comparison and defines sub-benchmarks according to the ID datasets MNIST, CIFAR10, CIFAR100, and ImageNet. Another dataset is then used as OOD data. Datasets are labeled as near OOD or far OOD according to their ID similarity, e.g. if they have similar color distributions. Open Set Recognition (OSR) is also provided by partitioning a dataset into ID and OOD classes. M-6 benchmark is based on MNIST, C-6 on CIFAR-10, C-50 on CIFAR-100, and T-20 on TinyImageNet with the numeral representing the number of ID classes.

The DBS dataset was collected from video recordings of a camera mounted on a commuter train in a typical operation. Object proposals were automatically

Table 1: OOD detection performance (AUROC in %) on OpenOOD benchmark and CIFAR-100 ID accuracy (%) for different approaches: Best performances marked in bold. Results from other methods taken from [24].

Method	M-6	C-6	C-50	T-20	MNIST		CIFAR-10		CIFAR-100		ImageNet		Acc
	osr	osr	osr	osr	near	far	near	far	near	far	near	far	CIFAR-100
ARPL [2]	-	-	-	-	93.9	99.0	87.2	88.0	74.9	74.0	-	-	71.7
Mahalanobis [14]	89.8	42.9	55.1	57.6	98.0	98.1	66.5	88.8	51.4	70.1	68.3	94.0	75.8
KNN [21]	97.5	86.9	83.4	74.1	96.5	96.7	90.5	92.8	79.9	82.2	80.8	98.0	77.1
ReAct [20]	82.9	85.9	80.5	74.6	90.3	97.4	87.6	89.0	79.5	80.5	79.3	95.2	75.8
MaxLogit [10]	98.0	84.8	82.7	75.5	92.5	99.1	86.1	88.8	81.0	78.6	73.6	92.3	77.1
Dropout [5]	96.2	84.5	81.1	73.6	91.5	97.1	87.3	90.4	80.1	79.4	-	-	77.1
DeepEnsemble [13]	97.2	87.8	83.1	76.0	96.1	99.4	90.6	93.2	82.7	80.7	-	-	80.5
PixMix [12]	93.9	90.9	78.0	73.5	93.7	99.5	93.1	95.7	79.6	85.5	-	-	77.1
ProtoDistVAE	98.4	76.6	69.0	62.4	99.9	100.0	80.0	90.6	65.3	74.2	69.6	80.1	48.8

collected and classified into trains and persons. The annotations were manually checked and OOD samples (i.e. false positive detections) were placed in a separate category. In our evaluation, we used 8351 samples of people, 8340 samples of trains, and 5001 non-objects labeled as OOD, all rescaled to size 64×64 . Person and train samples were divided equally into training (60%), validation (10%), and test (30%) splits (OOD samples used only for testing). We use $J=1$ prototype per class in all experiments as a higher number did not improve the performance.

The generalized Gaussian parameters α and β were both set to 2 for all experiments. The encoder was chosen as ResNet-50 [9] for ImageNet and as ResNet-18 for all benchmarks with 64×64 sized images (including the DBS dataset) and 32×32 sized images. A convolutional encoder with five layers was used for all 28×28 sized images, for the decoder a five-layered network using subpixel-convolutions [18] is used. For ImageNet the decoder consists of seven layers and for all other benchmarks, it consists of six layers. The latent dimensionality L is chosen as $1/3$, $1/24$ or $1/96$ of the input dimensionality. After training, ID validation data were used for normalization of the OOD scores, which are used afterwards for score fusion during testing. For evaluation, ID classification performance is measured in accuracy and OOD detection performance in Area Under the Receiver Operating Characteristic (AUROC). AUROC is a threshold-independent metric and measures how well a score separates ID and OOD.

4.1 OOD Detection Performance

Table 1 shows the OOD detection performance in terms of AUROC compared to state-of-the-art methods. ProtoDistVAE was trained using only LPIPS reconstruction loss with weight $w_{\text{rec}} = 1$. Cross-entropy and KL divergence loss were used similarly with a weight of $w_{\text{cls}} = w_{\text{KL}} = 1$. Distance ratio $\lambda_{\text{DistRatio}}$ and LPIPS λ_{LPIPS} were used as scores to be fused by L_{∞} norm. The latent space dimensionality L was chosen as $1/24$ of the input dimensionality.

Compared to the other methods, ProtoDistVAE performs best on the MNIST-based benchmarks. This is likely due to its low diversity, making it easier to learn a latent distribution. For CIFAR10, ProtoDistVAE performs on par with other

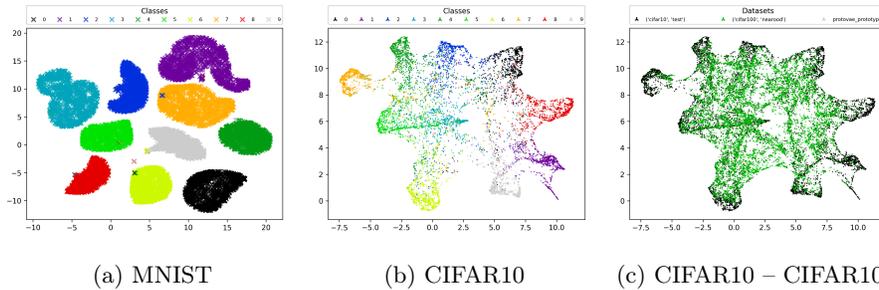


Fig. 2: UMAP visualization of the latent space embeddings of trained ProtoDistVAEs. (a) On MNIST, color-coded classes are clearly separated. (b) On CIFAR10, clusters blend into each other. (c) ID (CIFAR10) versus OOD (CIFAR100): embedding of OOD samples appears mainly between class prototypes.

methods. However, the performance for highly diverse datasets with a large number of classes decreases as ID estimation and classification are performed in the same latent space and may impair each other. Similarly, higher resolutions lead to difficulties for ProtoDistVAE in detecting OOD samples, likely due to the increased complexity of reconstruction.

Figure 2 shows further insights through an Uniform Manifold Approximation and Projection (UMAP) visualization of the latent space and illustrates how our method allows understanding its decision-making. The method works best in cases of clearly separable datasets and performs worse if data cannot be attributed well to clusters extracted. However, it should be mentioned that CIFAR10 vs. CIFAR100 is generally a hard OOD benchmark. ID samples in the space between prototypes might be interesting for further analysis since these exhibit a higher uncertainty and could be exploited by active learning or for identifying a sample with very different attributes within a class.

Table 2a shows some results on the DBS dataset. Here, an increased weight on LPIPS ($w_{\text{rec}} = 100$) was used to improve the OOD detection performance without harming classification accuracy. The accuracy is on par with other methods, likely due to only two classes being available. For OOD detection, PixMix and ProtoDistVAE perform best, while VIM and KNN also show good results. Combining λ_{LPIPS} with λ_{MSP} further improves the results with a gain of 0.9%.

ProtoDistVAE performs well on the DBS dataset due to its composition. The data samples are often quite similar as trains and persons are captured from the same angles and there are little variations e.g. in perspective, weather, lighting, and color. In comparison, ImageNet shows more inconsistent data with more diverse appearances across the same class. ProtoDistVAE benefits from a reduced intra-class variance and “complete” data distribution which allows it to model the data more easily. Hypothetically, it is easier for the network to recognize systematics in the data. PixMix augmentation seems to benefit from a complete distribution and even further increases the diversity of the data. However, the data distribution is not represented in the model and classification is not transparent. Other methods perform worse: Ensembling shows a lower-

Table 2: Experimental results of OOD detection in AUROC (%) and ID accuracy (%): (a) DBS dataset results of state-of-the-art methods (parameterized as in [24]) compared to ProtoDistVAE with LPIPS score combined by L_∞ fusion with DistRatio and MSP, respectively. (b) ProtoVAE vs. ProtoDistVAE. (c) Influence of reconstruction loss when using LPIPS as OOD score.

(a) DBS dataset			(c) OpenOOD benchmark (partial): reconstruction loss											
Method	AUC	Acc	L	Loss	mnist		cifar10		cifar100		ImageNet		cifar10	cifar100
					near	far	near	far	near	far	near	far	ID	ID
					AUC	AUC	AUC	AUC	AUC	AUC	AUC	AUC	Acc	Acc
ARPL [2]	80.2	99.7												
Mahalanobis [14]	75.4	93.6												
KNN [21]	84.9	99.6												
MSP [11]	83.6	99.6												
MaxLogit [10]	83.6	99.6												
ReAct [20]	81.8	99.6												
VIM [22]	85.4	99.6												
Dropout [5]	79.8	99.7												
DeepEnsemble [13]	83.4	99.7												
Pixmix [12]	89.3	99.6												
LPIPS+DistRatio	87.9	99.5												
LPIPS+MSP	88.8	99.5												

(b) OpenOOD benchmark: ProtoVAE vs. ProtoDistVAE using MSP score															
arch	L	mnist6	cifar6	cifar50	tin20	mnist		cifar10		cifar100		ImageNet		cifar10	cifar100
		osr	osr	osr	osr	near	far	near	far	near	far	near	far	ID	ID
		AUC	Acc	Acc											
ProtoVAE	1/3	94.1	66.1	65.3	62.4	92.1	98.6	71.2	75.9	61.5	57.3	55.6	30.7	76.1	39.8
ProtoVAE	1/24	95.7	68.6	66.2	61.4	92.3	99.1	71.7	75.1	62.9	56.9	49.9	49.9	80.1	43.0
ProtoVAE	1/96	96.5	65.4	64.8	62.5	94.9	98.4	71.5	76.8	62.5	61.5	50.1	50.1	79.4	41.0
ProtoDistVAE	1/3	97.9	76.5	63.7	60.5	96.7	99.0	75.6	73.9	63.9	51.1	60.7	70.7	81.3	41.0
ProtoDistVAE	1/24	96.4	76.0	67.4	59.4	91.4	93.4	76.3	76.3	64.5	62.2	67.3	78.6	79.4	46.7
ProtoDistVAE	1/96	96.1	77.3	66.8	63.1	94.5	96.5	75.4	76.9	65.1	56.4	50.1	50.1	78.7	46.5

than-usual performance as it depends on variations in the prediction of individual networks and these variations are weaker due to low data diversity in this dataset. Methods depending purely on classification-based schemes might suffer from overconfidence due to easier classification across only two classes and low data diversity. ProtoDistVAE, however, does not overfit for classification and aims to learn a representation of the data. In addition, the reconstruction error helps it to identify overconfidently classified samples mapped into its ID-region.

4.2 Ablation Study: ProtoVAE vs. ProtoDistVAE

Comparing the proposed ProtoDistVAE architecture to the base ProtoVAE, the reconstruction loss was set to a constant level. This does not change reconstruction error-based OOD detection according to the observed data. Table 2b shows detection results for ProtoVAE and ProtoDistVAE using the distance-based MSP score based on the predicted probabilities. Note that an improved distance-based score potentially increases performance even further when fused with a reconstruction error score. ProtoDistVAE outperforms ProtoVAE in almost all benchmarks for OOD detection and for different values of the latent dimension L which can be explained by the direct use of distances for classification and the enclosing restriction used during training. The latter actively shapes the ID-region by trapping the ID embeddings in the proximity of the



Fig. 3: Comparison of MSE and LPIPS loss: CIFAR10 (ID) and FashionMNIST (OOD). From top to bottom: Input, reconstruction (MSE), and reconstruction (LPIPS). ($L = 32$)

class-specific prototypes. Furthermore, the results display the importance of the latent dimensionality L for both networks. Different values for L are optimal for different levels of complexity reflected in different datasets. Too low values reduce the information coded in the representation while too high values inhibit a clear assignment of samples to class prototypes.

4.3 Reconstruction

Table 2c shows OOD detection performance using the LPIPS score based on ProtoDistVAE trained with either MSE or LPIPS loss. In contrast to using the MSE score which showed a generally lower performance (results not shown), the LPIPS score can achieve good detection results, even when training with MSE reconstruction loss. However, using LPIPS as reconstruction loss outperforms MSE loss. A special case is the ImageNet benchmark which is different due to image size and data diversity. The reconstruction performance for MSE and LPIPS loss on the CIFAR10 benchmark is depicted in Figure 3. ProtoDistVAE trained with MSE shows significant blur, regardless of ID or OOD samples. Training with LPIPS helps to preserve more semantic information and leads to differences when reconstructing OOD samples.

Figure 4 displays reconstructions of the DBS dataset. ProtoDistVAE appears to have learned the data distribution and can reconstruct ID better than OOD in most cases. It successfully distinguishes the class distributions of persons and trains and can show the features associated with a certain sample. For example, images of train stations and regular structures are often associated with trains, whereas background images are often reconstructed into person-like images. The learned prototypes of ProtoDistVAE can also be reconstructed. As Figure 5 shows, prototypes can be better extracted from datasets with low-variance datasets like MNIST and the DBS dataset while for datasets with higher diversity like CIFAR10, prototypes are harder to extract and images are less expressive. Human observers can thus assess which properties the network extracted from the data and evaluate features associated across classes.

5 Conclusion

In this work, ProtoDistVAE was applied to OOD detection. Its classification based on data density estimation, reconstruction, and prototype distances makes

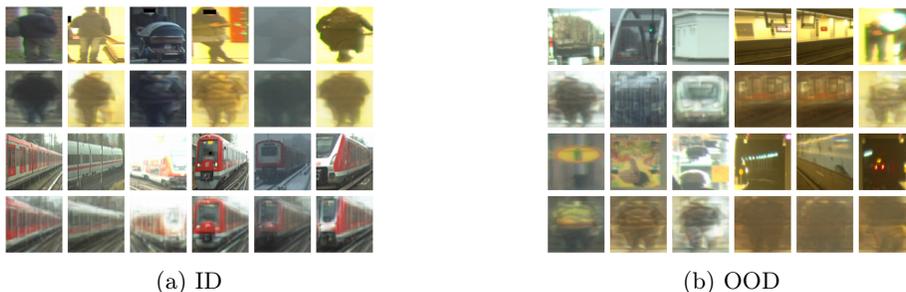


Fig. 4: DBS samples and reconstructions: ID and OOD ($L=4096$, $w_{LPIPS}=100$)

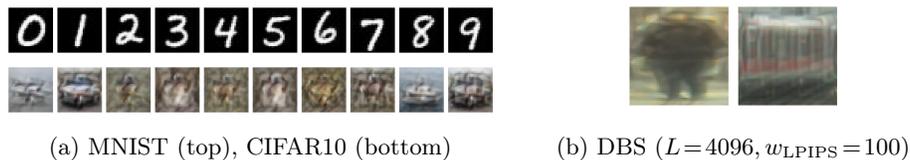


Fig. 5: Prototype reconstructions of different benchmarks

the approach robust and enables human insights into the classification and OOD detection process. The reconstruction of embeddings and prototypes enhances the model transparency w.r.t. properties learned from the data additionally.

The proposed method is extensively evaluated against state-of-the-art approaches using OpenOOD [24] and real-world railway data. Combining the estimation of the ID region with a perceptual reconstruction metric as well as integrating the proposed principle of enclosing restriction prove to be a good basis while the simple score fusion results in a gradual OOD score without imposing too many modeling assumptions. As empirically shown for the real-world railway dataset, the model learns In-Distribution training data and can reliably exploit it for OOD detection. The reconstruction error is an additional cue to identify system failure. Future work could involve combining an enclosing restriction with other reconstruction-based approaches and extending the presented framework to the task of object detection. Adversarial training may be used to increase the overall robustness.

References

1. Bercea, C., Rueckert, D., Schnabel, J.: What do we learn? debunking the myth of unsupervised outlier detection. arXiv preprint arXiv:2206.03698 (2022)
2. Chen, G., Peng, P., Wang, X., Tian, Y.: Adversarial reciprocal points learning for open set recognition. TPAMI pp. 8065–8081 (2022)
3. Denouden, T., Salay, R., Czarnecki, K., Abdelzad, V., Phan, B., Vernekar, S.: Improving reconstruction autoencoder out-of-distribution detection with mahalanobis distance. arXiv preprint arXiv:1812.02765 (2018)
4. Du, X., Gozum, G., Ming, Y., Li, Y.: Siren: Shaping representations for detecting out-of-distribution objects. In: NeurIPS. vol. 35, pp. 20434–20449 (2022)
5. Gal, Y., Ghahramani, Z.: Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In: ICML. pp. 1050–1059 (2016)

6. Gautam, S., Boubekki, A., Hansen, S., Salahuddin, S., Jenssen, R., Höhne, M., Kampffmeyer, M.: Protovae: A trustworthy self-explainable prototypical variational model. In: *NeurIPS*. pp. 17940–17952 (2022)
7. Graham, M.S., Pinaya, W.H.L., Tudosi, P.D., Nachev, P., Ourselin, S., Cardoso, M.J.: Denoising diffusion models for out-of-distribution detection. In: *CVPR Workshops*. pp. 2948–2957 (2023)
8. Guo, C., Pleiss, G., Sun, Y., Weinberger, K.Q.: On calibration of modern neural networks. In: *ICML*. pp. 1321–1330 (2017)
9. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *CVPR*. pp. 770–778 (2016)
10. Hendrycks, D., Basart, S., Mazeika, M., Zou, A., Kwon, J., Mostajabi, M., Steinhardt, J., Song, D.: Scaling out-of-distribution detection for real-world settings. In: *ICML*. pp. 8759–8773 (2022)
11. Hendrycks, D., Gimpel, K.: A baseline for detecting misclassified and out-of-distribution examples in neural networks. In: *ICLR* (2017)
12. Hendrycks, D., Zou, A., Mazeika, M., Tang, L., Li, B., Song, D., Steinhardt, J.: Pixmix: Dreamlike pictures comprehensively improve safety measures. In: *CVPR*. pp. 16762–16771 (2022)
13. Lakshminarayanan, B., Pritzel, A., Blundell, C.: Simple and scalable predictive uncertainty estimation using deep ensembles. *NIPS* p. 6405–6416 (2017)
14. Lee, K., Lee, K., Lee, H., Shin, J.: A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *NeurIPS* p. 7167–7177 (2018)
15. Nalisnick, E., Matsukawa, A., Teh, Y.W., Gorur, D., Lakshminarayanan, B.: Do deep generative models know what they don’t know? In: *ICLR* (2019)
16. Oza, P., Patel, V.M.: C2ae: Class conditioned auto-encoder for open-set recognition. In: *CVPR*. pp. 2307–2316 (2019)
17. Ruff, L., Vandermeulen, R., Goernitz, N., Deecke, L., Siddiqui, S.A., Binder, A., Müller, E., Kloft, M.: Deep one-class classification. In: *ICML*. pp. 4393–4402 (2018)
18. Shi, W., Caballero, J., Huszár, F., Totz, J., Aitken, A.P., Bishop, R., Rueckert, D., Wang, Z.: Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In: *CVPR*. pp. 1874–1883 (2016)
19. Sun, X., Yang, Z., Zhang, C., Ling, K.V., Peng, G.: Conditional gaussian distribution learning for open set recognition. In: *CVPR*. pp. 13477–13486 (2020)
20. Sun, Y., Guo, C., Li, Y.: React: Out-of-distribution detection with rectified activations. *NeurIPS* pp. 144–157 (2021)
21. Sun, Y., Ming, Y., Zhu, X., Li, Y.: Out-of-distribution detection with deep nearest neighbors. In: *ICML*. pp. 20827–20840 (2022)
22. Wang, H., Li, Z., Feng, L., Zhang, W.: Vim: Out-of-distribution with virtual-logit matching. In: *CVPR*. pp. 4921–4930 (2022)
23. Xiao, Z., Yan, Q., Amit, Y.: Likelihood regret: An out-of-distribution detection score for variational auto-encoder. *NeurIPS* pp. 20685–20696 (2020)
24. Yang, J., Wang, P., Zou, D., Zhou, Z., Ding, K., Peng, W., Wang, H., Chen, G., Li, B., Sun, Y., et al.: Openood: Benchmarking generalized out-of-distribution detection. *NeurIPS* pp. 32598–32611 (2022)
25. Yang, J., Zhou, K., Li, Y., Liu, Z.: Generalized out-of-distribution detection: A survey. *arXiv preprint arXiv:2110.11334* (2021)
26. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: *CVPR*. pp. 586–595 (2018)
27. Zhou, Y.: Rethinking reconstruction autoencoder-based out-of-distribution detection. In: *CVPR*. pp. 7369–7377 (2022)