

CVAD: A Generic Medical Anomaly Detector Based on Cascade VAE

Xiaoyuan Guo*, Judy Wawira Gichoya†, Saptarshi Purkayastha‡, and Imon Banerjee¶

*Department of Computer Science, Emory University, Georgia, USA

†Department of Radiology and Imaging Sciences, Emory University, Georgia, USA

‡Department of Biomedical Informatics, Emory University, Georgia, USA

§School of Informatics and Computing, Indiana University-Purdue University Indianapolis, IN, USA

¶School of Computing, Informatics, and Decision Systems Engineering, Arizona State University, AZ, USA

Email: {xiaoyuan.guo,judywawira}@emory.edu, saptpurk@iupui.edu, banerjee.imon@mayo.edu

Abstract—Detecting out-of-distribution (OOD) samples in medical imaging plays an important role for downstream medical diagnosis. However, existing OOD detectors are demonstrated on natural images composed of classes with clear inter-class variations and have difficulty generalizing to medical images. The key issue is the granularity of OOD data in the medical domain, where intra-class OOD samples are predominant. We focus on the generalizability of OOD detection for medical images and propose a self-supervised Cascade Variational autoencoder-based Anomaly Detector (CVAD). We use a cascaded variational autoencoder architecture, which combines latent representation at multiple scales, before being fed to a discriminator to distinguish the OOD data from the in-distribution (ID) data. Finally, both the reconstruction error and the OOD probability predicted by the binary discriminator are used to determine the anomalies. We compare the performance with the state-of-the-art deep learning models to demonstrate our model’s efficacy on various open-access medical imaging datasets for both intra- and inter-class OOD. Further extensive results on datasets including common natural datasets show our model’s effectiveness and generalizability.

I. INTRODUCTION

Despite recent advances in deep learning that have contributed to solving various complex real-world problems [1], [2], the safety and reliability of AI technologies remain a big concern in medical applications. Deep learning models for medical tasks are often trained with data from known distributions, and fail to identify out-of-distribution (OOD) inputs and possibly assign high probabilities to the anomalies during inference because of the insensitivity to distribution shifting. Medical anomalies, *a.k.a.*, *OOD data*, *outliers*, can arise due to various reasons such as noise during data acquisition, changes in disease prevalence and incidence (*e.g.*, the evolution of rare cancer types), or inappropriate inputs (*e.g.*, different modalities unseen during training) [3]. To ensure the reliability of deep models’ predictions, it is necessary to identify unknown types of data that are different from the training data distribution. A good anomaly detector should be able to learn good representations of the in-distribution (ID) during training and thus identify the outliers from test datasets. However, the core challenges for medical anomaly detection are – (1) the OOD data is usually unavailable at the time of model training; (2) in theory, there are infinite

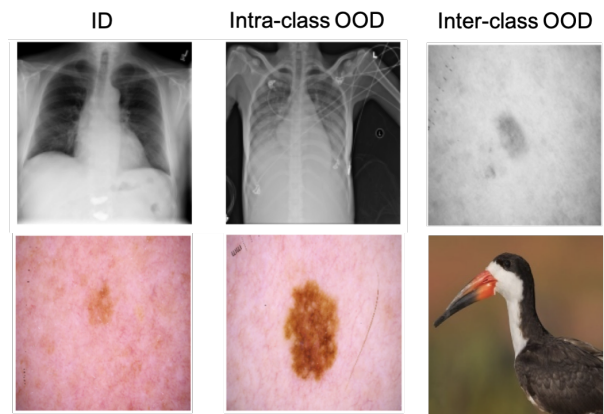


Fig. 1. ID, Intra- and Inter-class OOD examples for medical images. Compared to natural images, medical OOD samples exhibit more subtle intra-class variations (*e.g.*, normal vs pneumonia in the 1st row and benign vs malignant in the 2nd row).

numbers of variations of OOD data; and (3) different types of OOD data can be identified with varying difficulties. In general, the OOD classifications [4] can be refined based on the variation difference by summarizing them as **inter-class** OOD data and **intra-class** OOD data. Inter-class OOD data is in a category different from the ID data¹, *e.g.* a skin cancer image *v.s.* a lung X-ray image; intra-class OOD data belongs to the same category as the ID data but different classes, *e.g.* a normal skin image *v.s.* a skin image with cancer. Therefore, inter-class OOD data often has larger variations from the ID data, whereas the intra-class OOD data is close to ID data, as observed in Figure 1. Thus, identifying intra-class OOD data is more difficult than the inter-class OOD data given subtle differences with ID data.

To cope with the OOD unavailability and uncertainty challenges, we adopt an unsupervised way to design our anomaly detector. For intra-class OOD data, we expect the model can be sensitive to minor variations and thus screen the dissimilar inputs. To acquire such high identification of hard OOD cases, we propose a **Cascade Variational autoencoder based Anomaly Detector (CVAD)** to learn both coarser and finer features

¹By default, we mean a category can contain several classes. For example, a bird category can include owls, woodpeckers, flamingos, etc.

inspired by [5], [6]. With the cascade VAE architecture to model the in-distribution representations, CAVD gains superior reconstructions and learns good-quality features to threshold out the OOD data. To enhance the detection ability of inter-class OOD data, we further train a binary discriminator with the reconstructed data as the fake OOD category. In this paper, our contributions are three-fold:

- We propose a generic medical OOD detector – CVAD. By utilizing a cascade VAE to learn latent variables of in-distribution data, CVAD owns good reconstruction ability of in-distribution inputs and obtains discriminative ability for OOD data based on the reconstruction error.
- We adopt a binary discriminator to further separate the in-distribution data from the OOD data by taking the reconstructed image as fake OOD samples. Thus, our model has better discriminative capability for the inter-class as well as intra-class OOD cases.
- We conduct extensive experiments on multiple public medical image datasets to demonstrate the generalization ability of our proposed model. We evaluate comprehensively against state-of-the-art anomaly detectors in detecting both intra-class and inter-class OOD data, showing improved performance. The implementation technical report including original code and usage instructions has been publicly available in [8].

II. RELATED WORK

Although there have been extensive research on outlier detection [1], [9], effective medical image OOD detectors are still lacking due to complicated data types (e.g., various modalities and protocols, difference in acquisition devices) and user-defined application situations (e.g., disease types). Without OOD data available during training, unsupervised anomaly detection becomes the mainstream research direction, which CVAD also belongs to. Recent unsupervised anomaly detection approaches can be roughly classified as two main categories - generative and objective.

A. Generative methods

Deep generative models appear to be promising in detecting OOD data since they can learn latent features of training data and generate synthetic data with similar features to known classes [10]. Thus, the compressed latent features can be used to distinguish OOD data from ID data. Two major families of deep generative models are Variational Autoencoders (VAEs) [11] and Generative Adversarial Networks (GANs) [12].

VAEs: Traditional AutoEncoders [13] (AEs) can reconstruct input images well and be used to detect anomalies [43], but risk learning the identity of deep image features. Comparatively, VAEs generate contents by regularizing the latent feature distribution representations. With this trait, VAE [11] and its modifications have been used widely in generating realistic synthetic images [7], [14], [15]. Although VAEs are theoretically elegant and easy to train with nice manifold representations, they usually produce blurry images that lack

detailed information [6], [14]. To improve the image reconstruction quality, pchVAE [7] adds a conditional hierarchical VAE branch to learn lower-level image components. The improved reconstructions of VAEs are adopted for detecting OOD samples based on the reconstruction quality [16]. Other approaches seek to enhance the reliable uncertainty estimation of VAE for better performance [1], [17], [18], [19], [20]. Reference [18] applies an improved noise contrastive prior (INCP) to acquiring reliable uncertainty estimate for standard VAEs; whereas Bayesian VAE [1] detects OOD by estimating a full posterior distribution over the decoder parameters using stochastic gradient Markov chain Monte Carlo. Nonetheless, most of the VAE-based OOD detections are only evaluated on natural image datasets (MNIST [21], FashionMNIST [22], CIFAR10 [23], SVHN [24], *etc.*), which are with small image size (e.g., 32×32) and clear intra- and inter-class variations.

GANs: Compared with VAEs, GANs usually generate much sharper images but face challenges in training stability and sampling diversity, especially when synthesizing high-resolution images [14]. Still, GANs remain popular in outlier detection, such as, AnoGAN [25], f-AnoGAN [44], ADGAN [26], GANomaly [27] to detect OOD samples using GAN architectures. Besides standard architectures, there are hybrid models that detect anomalies by combining a VE/VAE with a GAN [5], [6], [28]. In order to acquire competitive OOD discriminative ability, OCGAN [28] integrates four components: a denoising auto-encoder, two discriminators and a classifier with complicate training process. Generally, such hybrid networks are not competitive for image datasets with clear class variations, as reported in [5]. Their experiments are often done with small-sized images and may fail when experimenting on large-sized medical images.

B. Objective methods

Objective anomaly detectors learn identifying OOD data via specific optimization functions and auxiliary transformations. Such OOD detection approaches include classifier-based and transformation-based methods [29], [30], [31].

Classifier-based method: ODIN [31] uses temperature scaling and adds small perturbations to input data for separating the softmax score distributions between ID and OOD images. Similar separation via a multi-class classifier is also followed by [32]. However, the prerequisite of balanced multiple classes is not always applicable in medical applications. Comparatively, the one-vs-rest setup [33] is much more common and useful in medical OOD detection, which treats one-class as in-distribution data and evaluates performance on the left OOD data. Following the setting, the anomaly detection reduces to a one-class classification (OCC) problem [34]. Representative one-class classifiers are DeepSVDD [29], OCSVM [35].

Transformation-based method: Most of the anomaly detectors are unsupervised given the assumption the anomalies are unavailable during training. Hence, good detection performance largely depends on the learning of high-quality in-distribution features. Self-augmentation with transformations on training data not only enriches the training diversity but

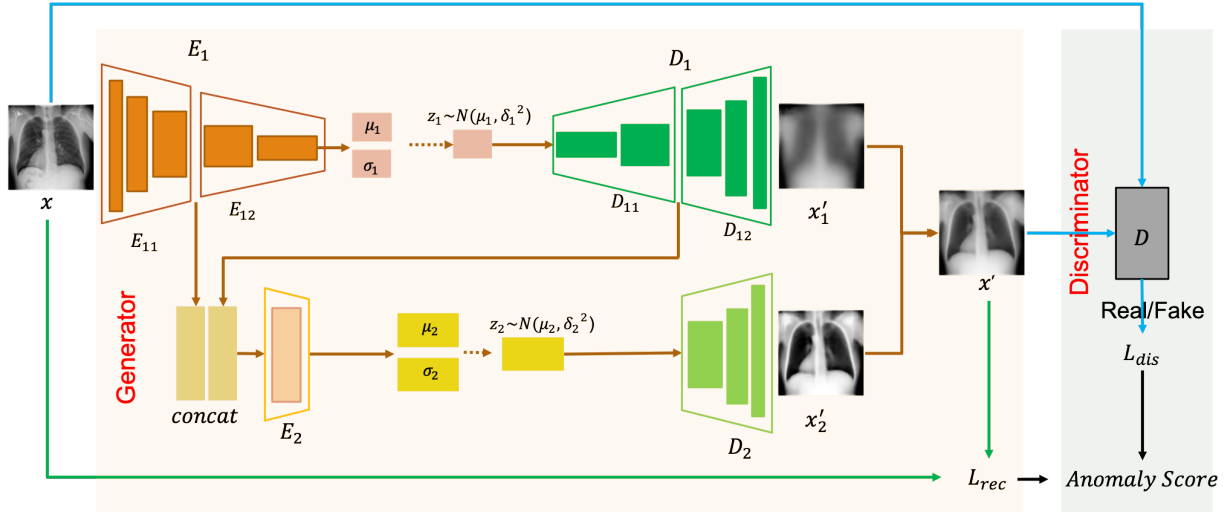


Fig. 2. Proposed CVAD architecture - a cascade VAE as the generator and a separate binary classifier (D) as the discriminator. The main VAE pipeline is composed by the encoder E_1 shown as the orange part and the decoder D_1 in the dark green part; the branch VAE has the pink part as the encoder E_2 and the light green for its decoder D_2 . Given an input image x , the main VAE learns to reconstruct x'_1 via latent representations μ_1 and σ_1 ; the branch VAE takes the outputs of the results of the main VAE encoder intermediate part E_{11} and the intermediate decoder D_{11} as inputs and feeds the concatenated features to E_2 to formulate the branch latent variables μ_2 and σ_2 , which gives a low-level reconstruction x'_2 via the corresponding decoder D_2 . By adding the two reconstructions - x'_1 and x'_2 together with a sigmoid function, a final reconstruction x' is generated and later treated as fake OOD data as compared to the original input x . The binary discriminator D will learn to distinguish them.

also introduces discriminative knowledge. For example, [30] proposes contrasting shifted instances for anomaly detection. Nevertheless, the augmentations are with limited transformations and consume more time to train as more generated data are fed as fake OOD data. Our model CVAD has no additional augmentations but still captures high-quality representations of in-distribution data. Besides, there are many other approaches contributing to OOD detection, such as GradCon [36], generalized ODIN [37] and FSSD [38]. Please refer to the papers for more details.

III. METHODS

Anomaly detection includes both intra- and inter-class OOD identification, of which medical intra-class OOD data is much more challenging because of the minute dissimilarity compared to ID data. With no prior knowledge available and no sophisticated pre-processing, we utilize a variational autoencoder to learn the “normality” of in-distribution inputs via image reconstruction and enhance the discriminative ability for both two OOD classes via a binary discriminator. Both the reconstruction and discrimination contribute to accurate intra- and inter-class OOD detection.

A. CVAD architecture

Figure 2 shows the design of CVAD. Inspired by the GAN’s architecture, we adopt a cascade VAE architecture as the “generator” for modeling ID representations and a separate classifier as the “discriminator” to strengthen OOD discrimination.

A standard VAE module consists of two neural networks: an encoder and a decoder [11], with the encoder $q_\phi(z|x)$ (parameterized by ϕ) mapping the visible variables x to the latent variables z and the decoder $p_\theta(x|z)$ (parameterized by

θ) sampling the visible variables x given the latent variables z [14]. Given a dataset $D = \{x_i\}_{i=1}^N$ with N input vectors drawn from some underlying data distribution $p^*(x)$, ϕ and θ are then learned by maximizing the variational lower bound (ELBO) $L(\phi, \theta)$, which is a lower bound to the marginal log-likelihood $\log p(x|\theta)$ [1]. However, a vanilla VAE exhibits limited potential in distinguishing unseen distributions due to the blurry reconstructions for large-size images. Thus, we learn from pchVAE [7] and tailored it as “generator” to acquire high-quality reconstruction and better latent representations.

Generator: Different from the standard VAE, our “generator” has two encoders E_1, E_2 and two decoders D_1, D_2 . To learn the high-level features, a deep and standard VAE architecture constructed by E_1 and D_1 formulates the deep latent variables z_1 by sampling parameters μ_1 and σ_1 of size K . Meanwhile, the low-level features are learnt by the branch VAE. Instead of using the original input, branch VAE utilizes the concatenation of two intermediate features from E_{11} and D_{11} . Given original input variables x , the input of branch VAE can be represented as $f(x)$. The encoder of branch VAE E_2 is simpler than E_1 whereas the decoder D_2 owns the same architecture as D_{12} . This branch VAE formulates latent Gaussian distributions with parameters μ_2, σ_2 of size $4K$. After sampling, two sets of latent variables, i.e., z_1, z_2 are acquired and decoded to image contexts x'_1 and finer details x'_2 respectively. x' is the combination of x'_1 and x'_2 .

Discriminator: Since the “generator” itself has no awareness of distinguishing outliers, we add a binary discriminator D to distinguish the reconstructed image x' from the original input image x . As x' shares very similar features with x after the first-stage training of the image generator, the discriminator is much more sensitive to minor differences from the in-distribution data, enhancing the accuracy of identifying both

intra-class OOD data and inter-class OOD data.

B. Network training

Instead of training CVAD in an adversarial way, we train the generator and the discriminator in two stages. The reason is that training with adversarial losses often leads to much sharper reconstructions but ignores the low-level information of ID data, incurring high reconstruction errors and potential dangerous decisions for medical applications. Therefore, CAVD is designed to first train the image generator and then the binary discriminator to detect OOD data. This non-adversarial training enables CVAD to inherit the merit of VAEs [11] and avoid the instability of GANs [12].

To optimize the generator, we minimize two objectives for the primary VAE part in Eqn. 1 and the branch VAE part in Eqn. 2, KL refers to Kullback-Leibler divergence.

$$L(x; \phi_1, \theta_1) = -E_{z_1 \sim q_{\phi_1}(z_1|x)}[\log p_{\theta_1}(x|z_1)] + D_{KL}(q_{\phi_1}(z_1|x) || p_{\theta_1}(z_1)) \quad (1)$$

$$L(x; \phi_2, \theta_2) = -E_{z_2 \sim q_{\phi_2}(z_2|f(x))}[\log p_{\theta_2}(x|z_2)] + D_{KL}(q_{\phi_2}(z_2|f(x)) || p_{\theta_2}(z_2)) \quad (2)$$

Therefore, the ‘‘generator’’ loss can be formulated as Eqn. 3. α_1 and α_2 to balance the weights of the two individual terms.

$$L_G = \alpha_1 L(x; \phi_1, \theta_1) + \alpha_2 L(x; \phi_2, \theta_2) \quad (3)$$

The binary discriminator is trained to distinguish true/fake images using binary cross entropy.

Anomaly score: An anomaly score S is defined in Eqn. 4 based on errors during inference and includes two parts: the reconstruction error L_G output by the ‘‘generator’’ and the probability of being the anomaly class S_D output by the discriminator. Instead of simply adding the two parts together, we first scale the ‘‘generator’’ reconstruction errors into $[0,1]$ for the whole dataset and get the average score value to avoid assigning imbalanced weights between the two parts:

$$S = 0.5 * \left(\frac{L_G - L_{Gmin}}{L_{Gmax} - L_{Gmin}} + S_D \right) \quad (4)$$

C. Network Details

As illustrated in Figure 2, our generator has a standard VAE part which consists of E_{11} , E_{12} , D_{11} and D_{12} and a branch VAE composed by a shallow encoder E_2 and a decoder D_2 . The primary VAE is a symmetric network with five 4×4 convolutions with stride 2 and padding 1 followed by five transposed convolutions. Respectively, E_{11} stands for the first three convolution layers; E_{12} refers the last two convolution layers; D_{11} is for the first three transposed convolution layers and D_{12} means the last two transposed convolution layers. The input of the branch VAE is the intermediate features of E_{11} and the middle decoded features of D_{11} . E_2 here is a convolution layer which has a same 4×4 kernel with stride 2 and padding 1. D_2 shares the same decoder architecture as the standard VAE, namely, $D_2 = D_{11} + D_{12}$. All convolutions and transposed-convolutions are followed by batch normalization and leaky ReLU (with slope 0.2) operations. We used a base

TABLE I
THE SELECTION DETAILS OF ID AND OOD DATA

Dataset	Details
RSNA	<i>In-class</i> : normal (8,851)
	<i>Intra-class</i> : pneumonia (9,555), abnormal (11,821)
	<i>InterClass1</i> : BIRD (37,715)
	<i>InterClass2</i> : SIIM (33,125)
IVC-Filter	<i>InterClass3</i> : IVC-Filter (1,258)
	<i>In-class</i> : type 11 (196)
	<i>Intra-class</i> : type 0-10, 12,13 (1,062)
	<i>InterClass1</i> : BIRD (37,715)
SIIM	<i>InterClass2</i> : SIIM (33,125)
	<i>InterClass3</i> : RSNA (30,227)
	<i>In-class</i> : benign (32,541)
	<i>Intra-class</i> : malignant (584)
	<i>InterClass1</i> : BIRD (37,715)
	<i>InterClass2</i> : IVC-Filter (1,258)
	<i>InterClass3</i> : RSNA (30,227)

channel size of 16 and increased number of channels by a factor of 2 with every encoder layer and decreased the number of channels to half for each decoder layer. The latent dimension K of z_1 is set as 512 and z_2 is with $4K$, i.e., 2048 dimensions. The binary discriminator is composed of five convolution layers with the same settings as above and a final fully connected layer to make a binary prediction. After a sigmoid function, the final ID/OOD class probability is obtained.

IV. EXPERIMENTS

We conducted extensive experiments, verifying the generalizability and effectiveness of our approach on multiple open-access medical image datasets for intra- and inter-class OOD detection. In total, we used four independent datasets, including three medical image datasets – RSNA Pneumonia dataset [40], inferior vena cava filters (IVC-Filter in short) on radiographs [41] and SIIM-ISIC Melanoma dataset [42] (identify melanoma in lesion images) and one natural image datasets – Bird Species². Among the medical datasets, RSNA and SIIM datasets have binary classes – normal and abnormal, whereas IVC-Filter dataset has 14 distinct types (classes). Table I lists the class information and number of images for each dataset and the corresponding usage in the **Detail** column. Bird dataset, which contains 270 bird species with 38,518 training images, was only used as inter-class OOD for detection validation. To unify the OOD detection pipeline and facilitate evaluation, we resized both the medical images and the validation inter-class OOD images to a unified $256 \times 256 \times channel$ size, where IVC-Filter and RSNA datasets are in gray scale with *channel* as 1 and the SIIM images are in RGB format and have *channel* 3. To train the anomaly detectors, we split the ID data into training and valuation parts in the ratio of 80% v.s. 20%. All the OOD data will only be used during evaluation phase.

We implemented our model using Pytorch 1.5.0, Python 3.6. α_1, α_2 were equal to 1. We ran the models on 4 NVIDIA

²<https://www.kaggle.com/gpiosenka/100-bird-species>

TABLE II

INTRA-CLASS OOD DETECTION RESULTS (FPR, TPR AND AUC VALUES) OF VARIOUS ANOMALY DETECTORS TRAINED ON RSNA, IVC-FILTER AND SIIM DATASETS. BEST RESULTS ARE HIGHLIGHTED. STANDARD DEVIATIONS ARE CALCULATED VIA 10 ROUNDS OF BOOTSTRAPPING ESTIMATIONS.

Methods	RSNA			IVC-Filter			SIIM		
	\downarrow FPR	\uparrow TPR	\uparrow AUC	\downarrow FPR	\uparrow TPR	\uparrow AUC	\downarrow FPR	\uparrow TPR	\uparrow AUC
AE [39]	0.318 \pm 0.014	0.461 \pm 0.009	0.566 \pm 0.004	0.198 \pm 0.104	0.350 \pm 0.075	0.436 \pm 0.040	0.420 \pm 0.024	0.714 \pm 0.030	0.673 \pm 0.006
VAE [16]	0.473 \pm 0.001	0.462 \pm 0.001	0.487 \pm 0.001	0.489 \pm 0.097	0.707 \pm 0.076	0.542 \pm 0.080	0.442 \pm 0.008	0.740 \pm 0.006	0.676 \pm 0.023
pchVAE [7]	0.501 \pm 0.018	0.731 \pm 0.030	0.600 \pm 0.007	0.590 \pm 0.072	0.620 \pm 0.013	0.472 \pm 0.038	0.378 \pm 0.045	0.558 \pm 0.040	0.616 \pm 0.012
DeepSVDD [29]	0.508 \pm 0.021	0.413 \pm 0.023	0.421 \pm 0.009	0.503 \pm 0.106	0.672 \pm 0.042	0.500 \pm 0.075	0.276 \pm 0.036	0.683 \pm 0.050	0.740 \pm 0.010
GANomaly [27]	0.524 \pm 0.005	0.678 \pm 0.015	0.576 \pm 0.005	0.446 \pm 0.172	0.627 \pm 0.227	0.518 \pm 0.103	0.553 \pm 0.103	0.495 \pm 0.108	0.418 \pm 0.016
f-AnoGAN [44]	0.365 \pm 0.033	0.541 \pm 0.029	0.614 \pm 0.005	0.419 \pm 0.077	0.611 \pm 0.054	0.544 \pm 0.042	0.381 \pm 0.000	0.624 \pm 0.033	0.721 \pm 0.015
CVAD (ours)	0.327 \pm 0.016	0.646 \pm 0.017	0.696 \pm 0.005	0.541 \pm 0.094	0.706 \pm 0.091	0.582 \pm 0.031	0.376 \pm 0.020	0.766 \pm 0.021	0.749 \pm 0.010

Quadro RTX 6000 GPUs with 24 GB memory each. In our model training, we used Adam optimizer with a learning rate of 0.001, and each network was trained for 100-350 epochs.

We evaluated our anomaly detection model performance in terms of standard statistical metrics - (i) area under the receiver operating characteristic (AUROC, AUC in short); (ii) True Positive rate (TPR); (iii) False positive rate (FPR). To classify ID and OOD classes, a threshold should be defined for the anomaly scores. Notably, the AUC value is threshold-invariant, while the TPR and FPR are determined by the selection of the anomaly threshold. We adopted the Geometric Mean (G-Mean) method to determine an optimal threshold for the ROC curve by tuning the decision thresholds and reported the resulting FPR and TPR values. To be fair and thorough, we ran all the experiments on both intra-class OOD and inter-class OOD to further analyze the performance of anomaly detectors on the specific type of OOD detection.

V. RESULTS

A. Quantitative Results

We set the vanilla AE and VAE architectures as baselines and compared our CVAD model with several representative models with varying architectures – pchVAE [7], a classifier-based approach DeepSVDD [29], and two GAN-based methods, i.e., GANomaly [27] and f-AnoGAN [44]. Table II shows the models’ performance for the intra-class OOD detection and Table III primarily presents the inter-class OOD performance.

1) *Results for Intra-class OOD Detection:* Intra-class OOD images are the most challenging outliers to identify since they often share similarity to the ID data but belong to a different class with unique characteristics. Still, CVAD exhibits its superiority in detecting intra-class OOD for medical images. On the RSNA dataset, CVAD achieves the best AUC score 0.696 (+0.275 from DeepSVDD’s AUC score 0.421, +0.120 from GANomaly’s AUC score 0.576, +0.082 from f-AnoGAN’s AUC score 0.614); for IVC-Filter, CVAD obtains the highest AUC values 0.582; for SIIM dataset, although DeepSVDD and f-AnoGAN show competitive performance, CVAD acquires the optimal AUC score 0.749. Overall, CVAD performs stably and effectively for intra-class OOD detection.

2) *Results for Inter-class OOD Detection:* To fairly evaluate all the models, we tested them on multiple inter-class OOD data types and presented the corresponding AUC scores in Table. III. As the OOD image datasets may have different

TABLE III

AUC SCORES PREDICTED BY OOD DETECTORS FOR INTER-CLASS IDENTIFICATION ON RSNA, IVC-FILTER AND SIIM DATASETS. BOLD INDICATES THE BEST PERFORMANCE.

Dataset	Methods	AUROC score		
		InterClass1	InterClass2	InterClass3
RSNA	AE [39]	0.677 \pm 0.006	0.608 \pm 0.005	0.616 \pm 0.004
	VAE [16]	0.752 \pm 0.004	0.604 \pm 0.007	0.613 \pm 0.006
	pchVAE [7]	0.790 \pm 0.006	0.776 \pm 0.005	0.632 \pm 0.007
	DeepSVDD [29]	0.838 \pm 0.005	0.834 \pm 0.004	0.604 \pm 0.006
	GANomaly [27]	0.733 \pm 0.005	0.816 \pm 0.004	0.597 \pm 0.007
	f-AnoGAN [44]	0.842 \pm 0.001	0.693 \pm 0.001	0.682 \pm 0.002
	CVAD (ours)	0.863 \pm 0.003	0.803 \pm 0.004	0.703 \pm 0.005
IVC-Filter	AE [39]	0.372 \pm 0.051	0.342 \pm 0.041	0.237 \pm 0.051
	VAE [16]	0.666 \pm 0.026	0.400 \pm 0.039	0.706 \pm 0.027
	pchVAE [7]	0.885 \pm 0.022	0.732 \pm 0.033	0.905 \pm 0.026
	DeepSVDD [29]	0.861 \pm 0.051	0.724 \pm 0.060	0.883 \pm 0.102
	GANomaly [27]	0.803 \pm 0.018	0.827 \pm 0.190	0.922 \pm 0.072
	f-AnoGAN [44]	0.911 \pm 0.020	0.625 \pm 0.043	0.864 \pm 0.042
	CVAD (ours)	0.984 \pm 0.002	0.911 \pm 0.017	0.985 \pm 0.001
SIIM	AE [39]	0.572 \pm 0.004	0.013 \pm 0.000	0.752 \pm 0.005
	VAE [16]	0.712 \pm 0.006	0.021 \pm 0.002	0.759 \pm 0.003
	pchVAE [7]	0.943 \pm 0.002	0.992 \pm 0.000	0.684 \pm 0.004
	DeepSVDD [29]	0.980 \pm 0.001	0.992 \pm 0.000	0.804 \pm 0.002
	GANomaly [27]	0.688 \pm 0.005	0.989 \pm 0.000	0.442 \pm 0.006
	f-AnoGAN [44]	0.951 \pm 0.001	0.924 \pm 0.002	0.606 \pm 0.003
	CVAD (ours)	0.983 \pm 0.001	0.978 \pm 0.001	0.869 \pm 0.003

image channels and image sizes from the ID training images, we adjusted the image channels and resized the images to ensure consistent input data format for evaluation³. CVAD obtains the highest AUC values on RSNA and SIIM datasets (except for inter-class2), and performs the best for IVC-Filter dataset across three inter-class OOD detection evaluations. Generally, the inter-class OOD detection of CVAD is satisfied with stable performance.

3) *Effectiveness of CVAD’s Components:* We here demonstrate the importance of each component of CVAD. Table IV shows the performance difference under the intra-class and three inter-class OOD data situations. CVAD_G represents the “generator”, CVAD_D stands for only using the predictions of the discriminator. CVAD balances the two components’ prediction. As can be observed, CVAD_G and CVAD_D show certain variations for different cases. For example, CVAD_D generally works better than CVAD_G for RSNA dataset but behaves worse than CVAD_G in SIIM scenario. Nevertheless, each component owns its unique OOD discriminative ability,

³For example, to evaluate trained models on RSNA, we converted the BIRD and SIIM images to grayscale mode and resized them to the same in-distribution image size.

TABLE IV
AUC SCORES PREDICTED BY THE “GENERATOR” CVAD_G, THE DISCRIMINATOR CVAD_D AND CVAD FOR INTER-CLASS IDENTIFICATION ON RSNA, IVC-FILTER AND SIIM DATASETS RESPECTIVELY.

Dataset	Methods	AUROC score			
		IntraClass	InterClass1	InterClass2	InterClass3
RSNA	CVAD_G (ours)	0.602±0.006	0.854±0.003	0.517±0.004	0.601±0.005
	CVAD_D (ours)	0.672±0.005	0.793±0.003	0.809±0.003	0.679±0.005
	CVAD (ours)	0.696±0.005	0.863±0.003	0.803±0.004	0.703±0.005
IVC-Filter	CVAD_G (ours)	0.568±0.031	0.981±0.003	0.787±0.023	0.983±0.002
	CVAD_D (ours)	0.543±0.041	0.661±0.018	0.925±0.011	0.834±0.013
	CVAD (ours)	0.582±0.031	0.984±0.002	0.911±0.017	0.985±0.001
SIIM	CVAD_G (ours)	0.746±0.010	0.995±0.000	0.995±0.000	0.827±0.004
	CVAD_D (ours)	0.724±0.008	0.874±0.002	0.055±0.001	0.862±0.005
	CVAD (ours)	0.749±0.010	0.983±0.001	0.978±0.001	0.869±0.003

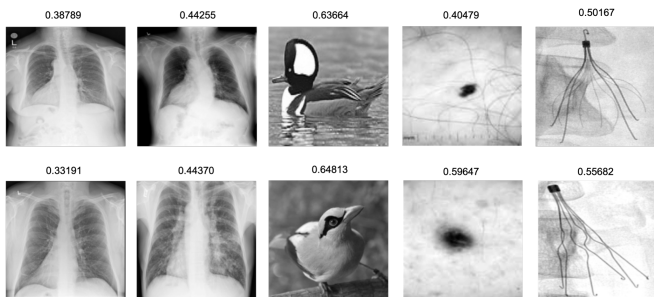


Fig. 3. Anomaly scores output by CVAD for different types of input data (experiments for RSNA dataset). Columns from left to right, ID, intra-class OOD, inter-class OOD1, inter-class OOD2, inter-class OOD3.

and combining their advantages entitles CVAD the capability of capturing both intra-class and inter-class dissimilarities. For which sake, CVAD has better generalization and can perform well and stably under different situations.

B. Qualitative Results

1) *Anomaly Detection*: Figure 3 shows two experimental results for RSNA dataset. Each row stands for one case and each column represents a specific type of input data. From left to right, they are in-distribution data, intra-class OOD data, inter-class OOD1 data, inter-class OOD2 data and inter-class OOD3 data, respectively. The corresponding anomaly score predicted by CVAD is on top of each example. Higher anomaly scores mean more likely the inputs are OOD. As can be seen in Figure 3, the two intra-class OOD samples (2nd column) are alike as the in-distribution data but the inter-class OOD examples show very different appearance from in-distribution data. Correspondingly, the anomaly scores of intra-class OOD are close to the scores of ID samples and difficult to separate whereas the intra-class OOD cases with clear variations are assigned higher anomaly scores and are easy to identify. This phenomenon further demonstrates the challenges of identifying intra-class OOD data.

2) *Visualization of Reconstruction Effects*: CVAD gains good latent in-distribution features via its “generator”, which learns both low-level and high-level representations. To demonstrate the effectiveness, we took RSNA dataset as a representative and showcased the reconstruction details in Figure 4, with the first column for branch VAE reconstruction

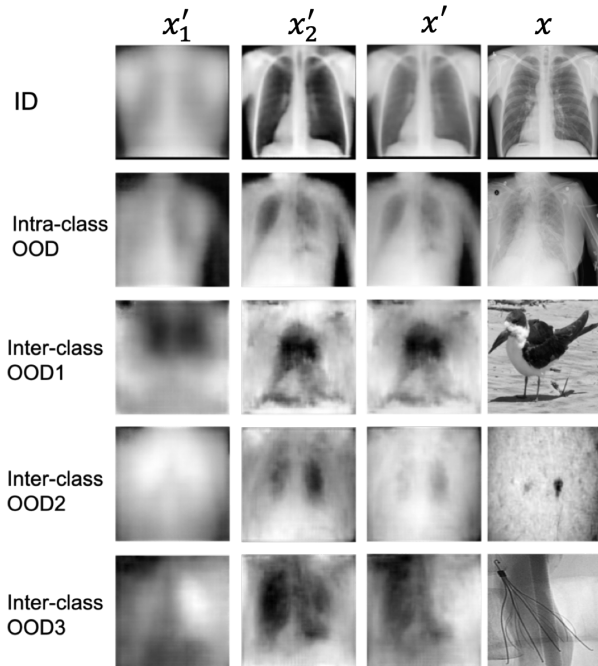


Fig. 4. Reconstruction details visualization of CVAD’s “generator” trained on RSNA dataset for different data types.

x'_2 , the second column for standard VAE part reconstruction x'_1 , the third column for ultimate reconstruction x' and the last column for the original input image x (following the same notations indicated in Figure 2). To further reveal the effects of “generator” on different OOD samples, we also presented example images for ID (i.e., normal class, 1st row), intra-class OOD (i.e., pneumonia or with opacity, 2nd row), inter-class OOD1 (i.e., gray-scale bird images, 3rd row), inter-class OOD2 (i.e., skin cancer images from SIIM dataset, 4th row) and inter-class OOD3 (i.e., images from IVC-Filter dataset, 5th row) in Figure 4. Compared with the intra-class medical OOD data, reconstructions on inter-class OOD inputs are more messy and dissimilar to the original OOD data, which leads to larger reconstruction errors and thus easier to distinguish. This observation reveals the varying difficulties of detecting different types of OOD data – intra-class OOD is much more challenging than inter-class OOD.

VI. CONCLUSION

We propose an effective medical anomaly detector CVAD that can reconstruct coarse and fine image components by learning multi-scale latent representations. The high quality of generated images enhances the discriminative ability of the binary discriminator in identifying unknown OOD data. We demonstrate the OOD detection efficacy for both intra-class and inter-class OOD data on various medical and natural image datasets. Our model has no prior assumptions on the input images and application scenarios for OOD, thus can be applied to detect OOD samples in a generic way for multiple scenarios. A detailed technical report about the code implementation and parameter usages of CVAD has been publicly available for easy reproduction.

REFERENCES

- [1] E. Daxberger, and J.M. Hernández-Lobato, "Bayesian variational autoencoders for unsupervised out-of-distribution detection," *arXiv:1912.05651*, 2019.
- [2] Jiali Duan and C-C Jay Kuo, "Bridging Gap between Image Pixels and Semantics via Supervision: A Survey," *arXiv preprint arXiv:2107.13757*, 2021.
- [3] T. Fernando, H. Gammulle, S. Denman, S. Sridharan, and C. Fookes, "Deep Learning for Medical Anomaly Detection—A Survey," *ACM Comput. Surv. (CSUR)*, vol. 54, no. 7, pp. 1–37, 2021.
- [4] T. Cao, C.W. Huang, D. Y. T. Hui, and J. P. Cohen, "A benchmark of medical out of distribution detection," *arXiv:2007.04250*, 2020.
- [5] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther. "Autoencoding beyond pixels using a learned similarity metric," in *Proc. 33rd Int. Conf. Mach. Learn. (ICML)*, Jun. 2016, pp. 1558–1566.
- [6] J. Bao, D. Chen, F. Wen, H. Li, and G. Hua. "CVAE-GAN: fine-grained image generation through asymmetric training," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2745–2754.
- [7] D. Zimmerer, J. Petersen, and K. Maier-Hein, "High-and Low-level image component decomposition using VAEs for improved reconstruction and anomaly detection," *arXiv:1911.12161*, 2019.
- [8] X. Guo, J.W. Gichoya, S. Purkayastha, and I. Banerjee, "CVAD/An unsupervised image anomaly detector," *Software Impacts*, 2021, pp. 100195.
- [9] D. Abati, A. Porrello, S. Calderara, and R. Cucchiara, "Latent space autoregression for novelty detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 481–490.
- [10] D. Li, D. Chen, J. Goh, and S. K. Ng, "Anomaly detection with generative adversarial networks for multivariate time series," *arXiv:1809.04758*, 2018.
- [11] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv:1312.6114*, 2013.
- [12] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville and Y. Bengio. "Generative adversarial networks," *arXiv:1406.2661*, 2014.
- [13] M. A. Kramer, "Nonlinear principal component analysis using auto-associative neural networks," *AICHE journal*, vol. 37, no. 2, pp. 233–243, 1991.
- [14] H. Huang, Z. Li, R. He, Z. Sun, and T. Tan, "Introvae: Intropective variational autoencoders for photographic image synthesis," *arXiv:1807.06358*, 2018.
- [15] C. K. Sønderby, T. Raiko, L. Maaløe, S.K. Sønderby, and O. Winther, "Ladder variational autoencoders," *arXiv:1602.02282*, 2016.
- [16] J. An, and S. Cho, "Variational autoencoder based anomaly detection using reconstruction probability," *Special Lecture on IE*, vol. 2, no. 1, pp. 1–18, 2015.
- [17] Z. Xiao, Q. Yan, and Y. Amit, "Likelihood regret: An out-of-distribution detection score for variational auto-encoder," *arXiv:2003.02977*, 2020.
- [18] X. Ran, M. Xu, L. Mei, Q. Xu, and Q. Liu, "Detecting out-of-distribution samples via variational auto-encoder with reliable uncertainty estimation," *arXiv:2007.08128*, 2020.
- [19] A. A. Pol, V. Berger, C. Germain, G. Cerminara, and M. Pierini, "Anomaly detection with conditional variational autoencoders," in *Proc. 18th Int. Conf. Mach. Learn. & App. (ICMLA)*, Dec. 2019, pp. 1651–1657.
- [20] G. Somepalli, and Y. Wu, Y. Balaji, B. Vinzamuri, and S. Feizi, "Unsupervised Anomaly Detection with Adversarial Mirrored AutoEncoders," *arXiv:2003.10713*, 2020.
- [21] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition" *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [22] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv:1708.07747*, 2017.
- [23] A. Krizhevsky, and G. Hinton, "Learning Multiple Layers of Features from Tiny Images," 2009.
- [24] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu and A.Y. Ng. "Reading digits in natural images with unsupervised feature learning," in *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- [25] T. Schlegl, P. Seeböck, S.M. Waldstein, U. Schmidt-Erfurth and G. Langs. "Unsupervised anomaly detection with generative adversarial networks to guide marker discovery," in *Proc. Int. Conf. Inf. Pro. Med. Imag. (IPMI)*, Jun. 2017, pp. 146–157.
- [26] L. Deecke, R. Vandermeulen, L. Ruff, S. Mandt, and M. Kloft, "Image anomaly detection with generative adversarial networks," *Joint Eur. Conf. Mach. Learn. & Principle & Practice Knowledge Discovery Database (ECML-PKDD)*, Sep. 2018, pp. 3–17.
- [27] S. Akcay, A. Atapour-Abarghouei, and T. P. Breckon, "Ganomaly: Semi-supervised anomaly detection via adversarial training," in *Proc. 14th As. Conf. Comput. Vis. (ACCV)*, Dec. 2018, pp. 622–637.
- [28] P. Perera, R. Nallapati, and B. Xiang, "Ocgan: One-class novelty detection using gans with constrained latent representations," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 2898–2906.
- [29] L. Ruff, R. Vandermeulen, N. Goernitz, L. Deecke, S.A. Siddiqui, A. Binder, E. Müller and M. Kloft. "Deep one-class classification," in *Proc. 35th Int. Conf. Mach. Learn. (ICML)*, Jul. 2018, pp. 4393–4402.
- [30] J. Tack, S. Mo, J. Jeong, and J. Shin. "Csi: Novelty detection via contrastive learning on distributionally shifted instances," *arXiv:2007.08176*, 2020.
- [31] S. Liang, Y. Li, and R. Srikant. "Enhancing the reliability of out-of-distribution image detection in neural networks," *arXiv:1706.02690*, 2017.
- [32] Q. Yu, and K. Aizawa. "Unsupervised out-of-distribution detection by maximum classifier discrepancy," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019, pp. 9518–9526.
- [33] P. Liznerski, L. Ruff, R.A. Vandermeulen, B.J. Franks, M. Kloft and K.R. Müller. "Explainable Deep One-Class Classification," *arXiv:2007.01760*, 2020.
- [34] S. S. Khan, and M. G. Madden, "One-class classification: taxonomy of study and review of techniques," *The Knowledge Engineering Review*, 2014, vol. 29, no. 3, pp. 345–374.
- [35] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, "Estimating the support of a high-dimensional distribution," *Neural computation*, 2001, vol. 13, no. 7, pp. 1443–1471.
- [36] G. Kwon, M. Prabhushankar, D. Temel, and G. AlRegib, "Backpropagated Gradient Representations for Anomaly Detection," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, Aug. 2020, pp. 206–226.
- [37] Y. C. Hsu, Y. Shen, H. Jin, and Z. Kira, "Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2020, pp. 10951–10960.
- [38] H. Huang, Z. Li, L. Wang, S. Chen, B. Dong and X. Zhou. "Feature space singularity for out-of-distribution detection," *arXiv:2011.14654*, 2020.
- [39] M. Sakurada, and T. Yairi, "Anomaly detection using autoencoders with nonlinear dimensionality reduction," in *Proc. MLSDA 2nd Workshop on Mach. Learn. Sensory Data Analysis*, 2014, pp. 4–11.
- [40] X. Wang, Y. Peng, L. Lu, Z. Lu, M. Bagheri, and R. M. Summers, "Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 2097–2106.
- [41] J.C. Ni, K. Shpanskaya, M. Han, E.H. Lee, B.H. Do, W.T. Kuo, K.W. Yeom and D.S. Wang. "Deep learning for automated classification of inferior vena cava filter types on radiographs," *Journal of Vascular and Interventional Radiology*, 2020, vol. 30, no. 1, pp. 66–73.
- [42] V. Rotemberg, N. Kurtansky, B. Betz-Stablein, L. Caffery, E. Chousakos, N. Codella, M. Combalia, S. Dusza, P. Guitera, D. Gutman and A. Halpern. "A patient-centric dataset of images and metadata for identifying melanomas using clinical context," *Scientific data*, 2021, vol. 8, no. 1, pp. 1–8.
- [43] C. Baur, S. Denner, B. Wiestler, N. Navab and S. Albarqouni. "Autoencoders for unsupervised anomaly segmentation in brain MR images: a comparative study," *Medical Image Analysis*, Elsevier, 2021, pp. 101952.
- [44] T. Schlegl, P. Seeböck, S.M. Waldstein, G. Langs and U. Schmidt-Erfurth. "f-AnoGAN: Fast unsupervised anomaly detection with generative adversarial networks," in *Medical image analysis*, Elsevier 2019, vol. 54, pp. 30–44.