

UNSUPERVISED ANOMALY DETECTION WITH LOCAL-SENSITIVE VQVAE AND GLOBAL-SENSITIVE TRANSFORMERS

Mingqing Wang¹, Jiawei Li¹, Zhenyang Li¹, Chengxiao Luo¹, Bin Chen^{2,†}, Shu-Tao Xia¹, Zhi Wang¹

¹Tsinghua Shenzhen International Graduate School, ²Harbin Institute of Technology, Shenzhen
 {wmq20, li-jw15}@mails.tsinghua.edu.cn; chenbin2021@hit.edu.cn

ABSTRACT

Unsupervised anomaly detection (UAD) has been widely implemented in industrial and medical applications, which reduces the cost of manual annotation and improves efficiency in disease diagnosis. Recently, deep auto-encoder with its variants has demonstrated its advantages in many UAD scenarios. Training on the normal data, these models are expected to locate anomalies by producing higher reconstruction error for the abnormal areas than the normal ones. However, this assumption does not always hold because of the uncontrollable generalization capability. To solve this problem, we present LSGS, a method that builds on Vector Quantised-Variational Autoencoder (VQVAE) with a novel aggregated codebook and transformers with global attention. In this work, the VQVAE focus on feature extraction and reconstruction of images, and the transformers fit the manifold and locate anomalies in the latent space. Then, leveraging the generated encoding sequences that conform to a normal distribution, we can reconstruct a more accurate image for locating the anomalies. Experiments on various datasets demonstrate the effectiveness of the proposed method.

Index Terms— Unsupervised Anomaly Detection, VQVAE, Aggregated Codebook, self-supervised training

1. INTRODUCTION

Auto-encoder (AE) with its variants is widely used in unsupervised anomaly detection (UAD) problem [1, 2, 3]. By learning a distribution of normal data, the AE-based UAD approach is expected to reconstruct an abnormal image into a normal image. Comparing the reconstructed image and the input one, it detects and locates anomalies without knowing what the target is, which allows it to be used in a variety of limited scenarios. However, it has been observed that sometimes the auto-encoder “generalizes” so well that it also reconstructs anomalies well, leading to miss detection of anomalies.

Previous studies [4, 5] fit a feature manifold of normal data, and generate normal images close to the input abnormal

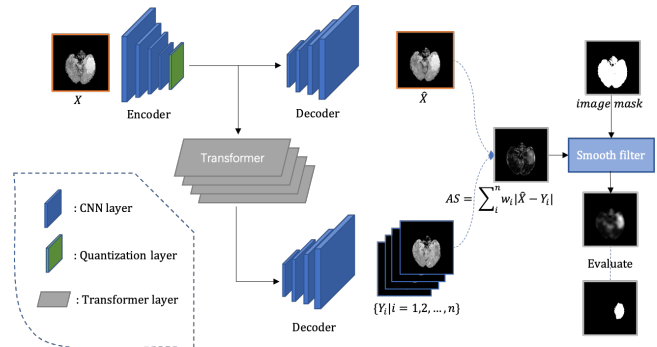


Fig. 1. Anomaly detection with the proposed method.

ones. These approaches mitigate the anomaly reconstruction but also generate unnecessary noise in normal areas, which causes miss detection sometimes. Memory-based approaches [6, 7] propose to augment the auto-encoder with a memory module and restore normal images from learned memory, which work well to repair textures that are not present in the training set. Nevertheless, for the structural anomalies of images, they behave poorly. Some data augmentation-based approaches [8, 9] seem to do better on this problem. [8] introduces prior abnormal patches, to learn a joint representation of an anomalous image and its anomaly-free reconstruction. It requires prior patches to cover all types of anomalies, so it is not an entirely unsupervised method. On the other hand, [9] uses extra natural images for model pre-training to improve its performance. The most relevant work is [10, 11], which combines a VQVAE and an auto-regression model like Pixel-CNN [12] or GPT2 [13]. They abstract the image encodings and repair the encodings one by one. Auto-regression models do well in the sequential generation of ordered sequences but lack integration over global image information.

In this work, we present LSGS, a method as shown in Fig. 1 that builds on an improved VQVAE and full-attention transformers. Specifically, we train the VQVAE on the normal images and extract all encodings of images in the training set. Then, the encodings are aggregated into a codebook with a fixed size. The aggregated codebook completely represents the distribution of discrete latent space. Further, we train the

[†] Corresponding Author

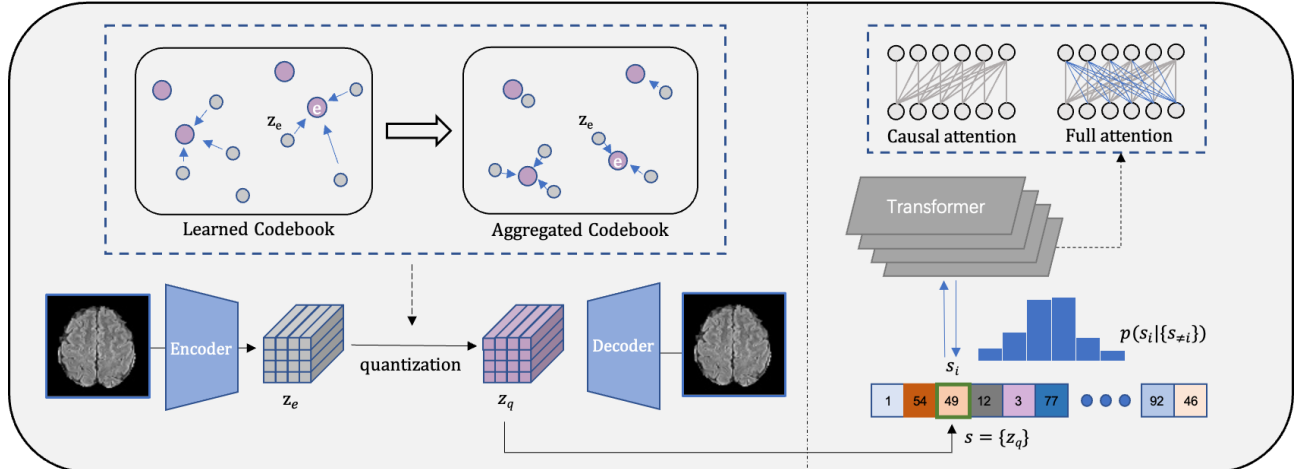


Fig. 2. The proposed method consists of two models: 1) A local-sensitive VQVAE including a CNN encoder, a CNN decoder, and an aggregated codebook. 2) A global-sensitive transformer with full-attention layers trained with a self-supervised strategy.

full-attention transformers on the discrete encoding sequences of images by a self-supervised training strategy with abnormally focal loss, which is sensitive to the global information of images. After completing the training, we utilize the model to repair abnormal encodings in the discrete latent space and reconstruct a normal image for locating anomalies at the pixel level.

The main contributions of our work are summarized as follows: (i) To better represent discrete latent space of images, we exploit the local sensitivity of VQVAE and propose a novel aggregated codebook; (ii) We propose to restore normal image encodings with global-sensitive transformers, and show a novel self-supervised training strategy; (iii) The supposed LSGS that builds on local-sensitive VQVAE and global-sensitive transformers achieve better anomaly-detection performance at the pixel level on both the medical and industrial datasets.

2. METHODOLOGY

2.1. Discrete encoding of images

The goal of the VQVAE is to find a reversible mapping relation, which maps each normal image patch x with a prescribed size of an image to a discrete latent coding z_q in the latent space as shown below:

$$z_q = Q(E(x)) \text{ and } \hat{x} = D(z_q) \quad (1)$$

where \hat{x} is the reconstructed image patch, and $E(\cdot)$, $D(\cdot)$ and $Q(\cdot)$ represent the encoder, the decoder and the learned codebook respectively, as shown in Fig. 2 (left). The elements of the codebook are called the embedding vector, denoted as e .

There are two classes of errors in the reversible mapping process: (1) reconstruction errors \mathcal{L}_{rec} , i.e. differences between the reconstructed image and the original image are

present in the encoding and decoding process; and (2) quantization errors \mathcal{L}_{VQ} that arise in the quantization process. The optimization objective is described as

$$\arg \min_{E, D, Q} \mathcal{L} = \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{VQ}} \quad (2)$$

In this work, the L_1 loss between x and \hat{x} is employed as the reconstruction loss, as shown in Eq. 3. We employed a straight-through estimator [14, 15] to accomplish gradient back-propagation in the discrete latent space.

$$\mathcal{L}_{\text{rec}} = \|x - \hat{x}\|_1 \quad (3)$$

Following [15], we employ a minimum distance quantizer, which irreversibly maps each latent coding z_e to the embeddings vector e under the closest Euclidean distance. The codebook is updated while training with \mathcal{L}_{VQ} as shown below.

$$\mathcal{L}_{\text{VQ}} = \|\text{sg}[E(x)] - e\|_2^2 + \|\text{sg}[e] - E(x)\|_2^2 \quad (4)$$

where sg is the gradient stop operator.

The model has learned the reversible mapping from images to discrete latent encodings after training on the normal images. However, not all embedding vectors of the codebook participate in this mapping process, which is commonly referred to as a codebook collapse.

Leveraging the local sensitivity of VQVAE, we aggregate the image encodings into a new codebook and further improve its representation capability. Specifically, after a joint training for \mathcal{L}_{rec} and \mathcal{L}_{VQ} , image patches with similar structures are mapped to adjacent latent points in the latent space. And a smaller Euclidean distance between the points represents a higher similarity of the image patches. Therefore, we extract all image encodings z_e in the training set and calculate cluster centers \hat{e} by the k-means algorithm. All embedding vectors are replaced by \hat{e} . Finally, the model is fine-tuned on

the normal images, which ensures the reversible mapping between the new codebook and images. Ablation experiments in Sec. 3.3 demonstrate its effectiveness.

2.2. Learning the distribution of codes

With the local sensitive VQVAE, the image patches are reconstructed from a learned codebook of the normal data. The reconstruction thus tends to be close to a normal sample. The reconstructed errors on local anomalies are strengthened for anomaly detection, which is similar to [6]. Nevertheless, this separate model performs poorly at some points: (1) anomalies on global positions rather than on local textures, and (2) reconstructed anomalies due to the strong generalization of the VQVAE.

Therefore, we employ a full-attention transformer to model the global information of images and rectify abnormal patches based on the prior distribution, which is referred to as a global sensitive transformer as shown in Fig. 2 (right). Specifically, By the learned VQVAE, an image X is represented with a corresponding discrete encoding sequence $s = \{z_q^1, z_q^2, z_q^3, \dots, z_q^n\}$. The prior distribution over the i -th discrete code s_i is a categorical distribution and can be modeled by depending on other codes $\{s_{\neq i}\}$ in the feature map. A full-attention transformer T is trained with all encoding sequences of normal images to fit the prior distribution $\phi(z) = \prod_i p(s_i|\{s_{\neq i}\})$, which is equivalent to maximize the log-likelihood of the data representations:

$$\arg \min_T \mathcal{L} = E_{X \sim \phi(X)} [-\log \phi(z)] \quad (5)$$

A cross-entropy reconstruction loss is commonly used to achieve the optimization objective. Note that the transformer prefers propagating an input encoding to the target output than merging it with other codes in the sequence s . To better combine the global features of the image, we perform a self-supervised training strategy. Specifically, we use random embedding vectors of the codebook to replace part of the encoding sequence. The position of replaced encodings is randomly selected. In our experiments, 10% of the encoding sequence is replaced. The goal of the transformer is to reconstruct the "tampered" encoding sequence into the original one. We train the transformer with an abnormally focal loss:

$$\mathcal{L}_{\text{Transformer}} = (1 - \beta) \sum_{z \in T} H(z) + \beta \sum_{z \notin T} H(z) \quad (6)$$

where T indicates all the "tampered" encodings, $H(*)$ is the cross-entropy loss function and β is a hyperparameter which is 0.01 in our experiments. The first term is the reconstruction loss which fits the prior distribution $\phi(z)$. And the second term is used to speed up the update of model parameters.

Note that compared to our work, [10, 11] employ an auto-regression model to solve the mentioned problems. However, it fails to integrate global information of images, which

leads to poor detection results sometimes. Sec. 3.3 demonstrates the comparison between the auto-regression structure (i.e. transformers with casual attention) and our method.

2.3. Anomaly detection with LSGS

As shown in Fig. 1, the discrete encoding sequence of an abnormal image is resampled according to the prior distribution fit by the transformer. Next, multiple normal images $\{Y_i | i = 1, 2, \dots, n\}$ are reconstructed from the the generated sequence. Afterward, the consolidated pixel-wise anomaly score (AS) is estimated as shown below:

$$AS = \sum_i^n w_i |\hat{X} - Y_i| \quad (7)$$

where \hat{X} is reconstructed from the original sequence which reduces perturbation of VQVAE reconstruction error, and $w_i = \text{softmax}(k/\|\hat{X} - Y_i\|_1)$ reduces the weight of restorations which have lost consistency. Finally, the anomaly score is fused with an image mask extracted from the original image and smoothed with a 3x3 MinPooling filter followed by a 7x7 AveragePooling filter as [10] done.

3. EXPERIMENTS

3.1. Experiments Setting

Dataset. To demonstrate the high effectiveness of the proposed method on images of different distributions, we evaluate the anomaly score on BraTS2018 and MVTEC-AD.

The BraTS2018 dataset, derived from the BRATS challenge, is a 3D MRI dataset. In the experiments, we use the flair attenuated inversion recovery data consisting of 163 samples. Each sample with a size of 240x240x155 is sliced into 155 images with a size of 128x128. The anomaly-free slices are used for training, and the remaining slices with anomalies are used for evaluation.

The MVTEC-AD dataset is commonly used for industrial anomaly detection. It consists of 15 different categories. Each category contains an anomaly-free training set, and a test set consisting of both normal and abnormal samples. We train the proposed method on training data resized to 128x128 of all categories and evaluate it on individual category data.

Implementation Details. In this work, the spatial downsampling rate of the VQVAE encoder is 8, which means that the image is encoded into a discrete feature map with a size of 16x16. The codebook size n is set to 1024 while training, and the channel dimension of encoded features is 512. The transformer consists of 12 multi-head attention layers, where the channel dimension is 768.

Evaluation Metrics. For quantitative evaluations, we measure the anomaly score from three metrics: Average Precision Score (AP), Area Under the Receiver Operating Characteristic Curve (AUROC), and Dice similarity coefficient (Dice).

Table 1. Quantitative comparison of pixel-wise anomaly detection in AP and Dice on the BraTS2018 dataset.

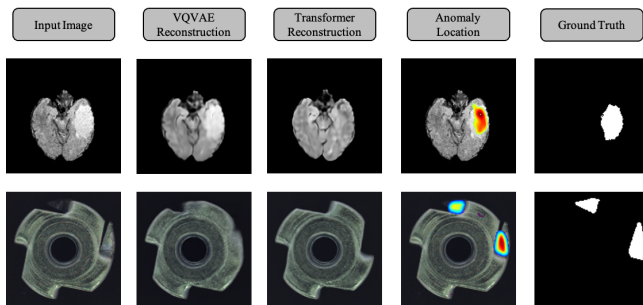
Metric	AE	VAE	GMVAE	fAnoGAN [16]	IS-cycle [2]	PraNet [17]	CaraNet [18]	LSGS (Ours)
AP	22.9	33.1	25.3	37.3	51.1	-	-	75.7
Dice	37.8	44.0	40.8	45.3	54.4	61.9	63.1	68.7

Table 2. Quantitative comparison of pixel-wise anomaly detection in AUROC on the MVTec-AD dataset.

Class\Method	US [9]	DRAEM [8]	LSGS(ours)
Bottle	67.9	87.6	92.5
Cable	78.3	71.3	91.2
Capsule	85.5	50.5	95.1
Hazelnut	93.7	96.9	94.0
Metal Nut	76.6	62.2	91.7
Pill	80.3	94.4	94.2
Screw	90.8	95.5	89.6
Toothbrush	86.9	97.7	93.7
Transistor	68.3	64.5	88.5
Zipper	84.2	98.3	87.8
Carpet	87.7	98.6	84.4
Grid	64.5	98.7	93.3
Leather	95.4	97.3	89.1
Tile	82.7	98.0	80.0
Wood	83.3	96.0	83.4
Avg.	81.8	87.2	89.9

3.2. Comparison to Existing Methods

We compare the supposed method to several state-of-the-art reconstruction-based UAD methods including f-AnoGAN [16] and DRAEM [8]. Note that the compared methods are trained and evaluated in a same dataset setting introduced in Sec. 3.1. The quantitative results as shown in Table 1 and Table 2 on datasets from two different domains indicate the generality and robustness of the proposed method.

**Fig. 3.** Visualized results on BraTS2018 and MVTec-AD datasets.

3.3. Ablation Studies

In this section, we evaluate the influence of several components of our framework in a controlled setting.

Aggregated codebook. As show in Table. 3, the supposed aggregated codebook better represents the distribution of discrete latent space with an increasing number of effective embeddings than the learned one. Next, we show lower reconstruction error (which is L1 loss in our work) of VQVAE with the aggregated codebook, which helps improve anomaly detection (evaluated in Dice on BraTS2018 and AUROC on MVTec-AD).

Attention. We compare the metrics on the BraTS2018 dataset of the following structures: (1) without the transformers (2) transformers of casual attention (used in [11]), and (3) transformers of full attention trained with the proposed self-supervised training strategy. Experiments shown in Table. 4 demonstrate that the proposed global-sensitive transformers achieves better anomaly detection than previous work.

Table 3. Effect of aggregated codebook on VQVAE reconstruction results. Size means the number of effective embeddings.

Dataset	Agg.	Size	Rec. Loss↓	Dice↑
BraTS2018	-	686	0.0115	67.4
	✓	4096	0.0084	68.7
MVTec-AD	-	472	0.0685	89.2
	✓	1024	0.0672	89.9

Table 4. Effect of global-sensitive transformer for anomaly detection on BraTS2018.

Transformer Type	AP	Dice
w/o	0.237	0.295
vanilla transformers	0.624	0.581
global-sensitive transformers	0.757	0.687

4. CONCLUSION

In this work, we propose a method for unsupervised anomaly detection, which builds on local-sensitive VQVAE and global-sensitive transformers. Two novel strategies are employed to improve model performance on anomaly detection. The experimental results show that the proposed model outperforms existing state-of-the-art methods.

5. REFERENCES

- [1] Christoph Baur, Stefan Denner, Benedikt Wiestler, Nassir Navab, and Shadi Albarqouni, “Autoencoders for unsupervised anomaly segmentation in brain mr images: a comparative study,” *Medical Image Analysis*, p. 101952, 2021.
- [2] Chenxin Li, Yunlong Zhang, Jiongcheng Li, Yue Huang, and Xinghao Ding, “Unsupervised anomaly segmentation using image-semantic cycle translation,” 2021.
- [3] Hans E Atlason, Askell Love, Sigurdur Sigurdsson, Vilmundur Gudnason, and Lotta M Ellingsen, “Unsupervised brain lesion segmentation from mri using a convolutional autoencoder,” in *Medical Imaging 2019: Image Processing*. International Society for Optics and Photonics, 2019, vol. 10949, p. 109491H.
- [4] Thomas Schlegl, Philipp Seeböck, Sebastian M Waldstein, Ursula Schmidt-Erfurth, and Georg Langs, “Unsupervised anomaly detection with generative adversarial networks to guide marker discovery,” in *International conference on information processing in medical imaging*. Springer, 2017, pp. 146–157.
- [5] Suhang You, Kerem C Tezcan, Xiaoran Chen, and Ender Konukoglu, “Unsupervised lesion detection via image restoration with a normative prior,” in *International Conference on Medical Imaging with Deep Learning*. PMLR, 2019, pp. 540–556.
- [6] Dong Gong, Lingqiao Liu, Vuong Le, Budhaditya Saha, Moussa Reda Mansour, Svetha Venkatesh, and Anton van den Hengel, “Memorizing normality to detect anomaly: Memory-augmented deep autoencoder for unsupervised anomaly detection,” in *IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [7] Jinlei Hou, Yingying Zhang, Qiaoyong Zhong, Di Xie, Shiliang Pu, and Hong Zhou, “Divide-and-assemble: Learning block-wise memory for unsupervised anomaly detection,” *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 8771–8780, 2021.
- [8] Vitjan Zavrtanik, Matej Kristan, and Danijel Skocaj, “Draem - a discriminatively trained reconstruction embedding for surface anomaly detection,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 8330–8339.
- [9] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger, “Uninformed students: Student-teacher anomaly detection with discriminative latent embeddings,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. jun 2020, IEEE.
- [10] Sergio Naval Marimont and Giacomo Tarroni, “Anomaly detection through latent space restoration using vector quantized variational autoencoders,” in *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2021, pp. 1764–1767.
- [11] Walter Hugo Lopez Pinaya, Petru-Daniel Tudosiu, Robert Gray, Geraint Rees, Parashkev Nachev, Sébastien Ourselin, and M Jorge Cardoso, “Unsupervised brain anomaly detection and segmentation with transformers,” in *Medical Imaging with Deep Learning*, 2021.
- [12] Aaron van den Oord, Nal Kalchbrenner, Oriol Vinyals, Lasse Espeholt, Alex Graves, and Koray Kavukcuoglu, “Conditional image generation with pixelcnn decoders,” *arXiv preprint arXiv:1606.05328*, 2016.
- [13] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al., “Language models are unsupervised multitask learners,” *OpenAI blog*, vol. 1, no. 8, pp. 9, 2019.
- [14] Yoshua Bengio, Nicholas Léonard, and Aaron Courville, “Estimating or propagating gradients through stochastic neurons for conditional computation,” *arXiv preprint arXiv:1308.3432*, 2013.
- [15] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu, “Neural discrete representation learning,” in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 6309–6318.
- [16] Thomas Schlegl, Philipp Seeböck, Sebastian M. Waldstein, Georg Langs, and Ursula Schmidt-Erfurth, “fanogan: Fast unsupervised anomaly detection with generative adversarial networks,” *Medical Image Analysis*, vol. 54, pp. 30–44, 2019.
- [17] Deng-Ping Fan, Ge-Peng Ji, Tao Zhou, Geng Chen, Huazhu Fu, Jianbing Shen, and Ling Shao, “Pranet: Parallel reverse attention network for polyp segmentation,” in *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*, Cham, 2020, pp. 263–273, Springer International Publishing.
- [18] Ange Lou, Shuyue Guan, Hanseok Ko, and Murray H. Loew, “CaraNet: context axial reverse attention network for segmentation of small medical objects,” in *Medical Imaging 2022: Image Processing*. International Society for Optics and Photonics, 2022, vol. 12032, pp. 81 – 92, SPIE.