

# EasyNet: An Easy Network for 3D Industrial Anomaly Detection

Ruitao Chen\*  
Southern University of Science and  
Technology  
Shenzhen, China  
chenrt2022@mail.sustech.edu.cn

Guoyang Xie\*  
Southern University of Science and  
Technology  
Shenzhen, China  
University of Surrey  
Guildford GU2 7XH, UK  
guoyang.xie@surrey.ac.uk

Jiaqi Liu\*  
Southern University of Science and  
Technology  
Shenzhen, China  
liujq32021@mail.sustech.edu.cn

Jinbao Wang†  
Southern University of Science and  
Technology  
Shenzhen, China  
linkingring@163.com

Ziqi Luo  
Southern University of Science and  
Technology  
Shenzhen, China  
luozq2022@mail.sustech.edu.cn

Jinfan Wang  
Southern University of Science and  
Technology  
Shenzhen, China  
Linkinsense  
Shenzhen, China  
wangjf@sustech.edu.cn

Feng Zheng†  
CSE and RITAS, Southern University  
of Science and Technology  
Shenzhen, China  
f.zheng@ieee.org

## ABSTRACT

3D anomaly detection is an emerging and vital computer vision task in industrial manufacturing (IM). Recently many advanced algorithms have been published, but most of them cannot meet the needs of IM. There are several disadvantages: i) difficult to deploy on production lines since their algorithms heavily rely on large pretrained models; ii) hugely increase storage overhead due to overuse of memory banks; iii) the inference speed cannot be achieved in real-time. To overcome these issues, we propose an easy and deployment-friendly network (called EasyNet) without using pretrained models and memory banks: firstly, we design a multi-scale multi-modality feature encoder-decoder to accurately reconstruct the segmentation maps of anomalous regions and encourage the interaction between RGB images and depth images; secondly, we adopt a multi-modality anomaly segmentation network to achieve a precise anomaly map; thirdly, we propose an attention-based information entropy fusion module for feature fusion during inference, making it suitable for real-time deployment. Extensive experiments show that EasyNet achieves an anomaly

detection AUROC of 92.6% without using pretrained models and memory banks. In addition, EasyNet is faster than existing methods, with a high frame rate of 94.55 FPS on a Tesla V100 GPU.

## CCS CONCEPTS

• **Computing methodologies** → **Computer vision tasks.**

## KEYWORDS

3D anomaly detection, multi-modality fusion, unsupervised learning, industrial manufacturing

## ACM Reference Format:

Ruitao Chen, Guoyang Xie, Jiaqi Liu, Jinbao Wang, Ziqi Luo, Jinfan Wang, and Feng Zheng. 2023. EasyNet: An Easy Network for 3D Industrial Anomaly Detection. In *Proceedings of the 31st ACM International Conference on Multimedia (MM '23)*, October 29–November 3, 2023, Ottawa, ON, Canada. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3581783.3611876>

## 1 INTRODUCTION

There is a strong need to propose a deployment-friendly 3D unsupervised anomaly detection (3D-AD) model to tap the gap, which brings 3D-AD's capabilities into the factory floor. Currently, most of anomaly detection methods [21, 23, 35, 36] are based on 2D images. But in the quality inspection of industrial products, human inspectors utilize both color (RGB) characteristics and depth information to determine whether it is a defective product, where depth information is essential for anomaly detection. As shown in Figure 2, for foam and peach, it is difficult to identify the anomalies from the RGB image alone. Though 3D-AD algorithms [7, 29, 33] are attracting interest from the academy, most of them are far from

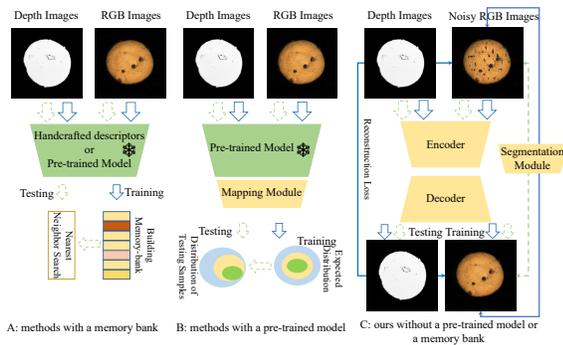
\*Equally contribute to this work

†Corresponding author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
MM '23, October 29–November 3, 2023, Ottawa, ON, Canada

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 979-8-4007-0108-5/23/10...\$15.00  
<https://doi.org/10.1145/3581783.3611876>

satisfactory for industrial manufacturing (IM). According to Figure 1, there are several issues: i) The cutting-edge 3D-AD methods steadily rely on the representational abilities of large pretrained models, leading to slow inference speed and huge storage overhead. ii) Feature embedding-based 3D-AD methods excessively use memory banks, leading to huge memory bank costs in real-world applications. Because of this, it is important and urgent to build an application-oriented 3D-AD model to meet the demands of IM.



**Figure 1: Illustration of 3D anomaly detection paradigms, including (A) feature embedding-based memory bank, (B) pretrained network, and (C) encoder-decoder (ours).**

To avoid using large pretrained models and memory banks, we propose an easy but effective multi-modality anomaly detection and localization network, called **EasyNet**. Specifically, EasyNet consists of two parts, the Multi-modality Reconstruction Network (MRN) and the Multi-modality Segmentation Network (MSN). First, instead of using pretrained features directly, we generate synthesized anomalies on RGB images and depth images, reconstruct the original images with semantically reasonable free content, and obtain multi-scale features. At the same time, in order to simplify the anomaly detection process based on memory bank, we input the abnormal and reconstructed multi-scale features into a simple MSN to obtain an anomaly map. As shown in Figure 3, the entire architecture, including MRN and MSN, significantly encourages interaction between RGB and depth features.

To reduce the disturbance between RGB and depth images, we propose an attention-based information entropy fusion module. We find that some 3D-AD methods, like AST [30] and BTF [17] cannot fully utilize the advantage of multi-modality fusion, i.e., RGB-D performance is not competitive than RGB performance. The main reason is that there are no uniform abnormal patterns in RGB or depth images. For example, some anomalies can be detected by pure RGB images and depth information works as the noise and may degrade the overall anomaly detection performance. Hence, we propose a dynamic multi-modality fusion scheme to make use of RGB and depth features. The architecture of the fusion scheme is shown in Figure 5. Moreover, as shown in Table 5, our proposed fusion scheme is simple and much more computationally efficient than the aforementioned 3D-AD models [17, 30, 33]. Our proposed attention-based information entropy fusion module is easy to train and apply, with outstanding performance and inference speed. As

a result, EasyNet can achieve 92.6% on MVTEC 3D-AD and 86.9% on Eyescandies in I-AUROC while running at 94.55 FPS, surpassing the previous best-published 3D-AD methods on accuracy and efficiency.

Our contributions can be summarized as follows:

- EasyNet is easy to implement and deploy for 3D unsupervised anomaly detection, i.e., eliminating the usage of pretrained models and memory banks, and achieves the fastest inference speed than the existing methods, with a high frame rate of 94.55 FPS on a Tesla V100 GPU.
- We propose an Attention-based Information Entropy Fusion Module to integrate the image features of the multi-modal characteristics well.
- We propose a Multi-modality Reconstruction Network (MRN) to accurately reconstruct the anomalous region and encourage the interaction of RGB and depth.
- We propose a Multi-modality Segmentation Network (MSN) to output the anomaly map precisely.
- EasyNet obtains the state-of-the-art result in Pure RGB. Note that EasyNet obtains the best anomaly detection I-AUROC of 92.6% in RGBD.

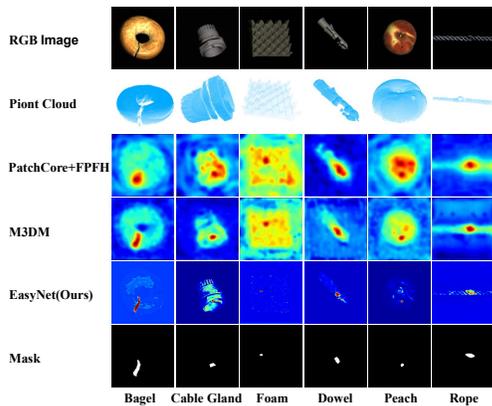
## 2 RELATED WORK

Anomaly detection (AD) is a classical topic, which aims to distinguish normal samples and abnormal samples. Existing experimental settings usually only take normal samples as the training set, and evaluate the ability of the model to distinguish abnormal samples in the test set. The current unsupervised AD can be mainly divided into feature extraction-based methods and image reconstruction-based methods. The former is restricted by the pretrained model, while the latter is free from this limitation. Based on this idea, we design a reconstructive AD algorithm for RGB-D data, removing the restrictions on pretrained models and memory banks.

### 2.1 2D Anomaly Detection

Since the emergence of the MVTEC AD dataset [4], research on AD in industrial 2D images has received more attention. Most existing research is based on this set for unsupervised AD tasks.

There is more research on feature embedding-based methods than reconstruction-based methods. The most basic idea is to regard AD as a one-class classification problem and turn the AD problem into a problem of finding boundaries for classification. CutPaste [20] and SimpleNet [24] are representative methods. They make abnormal samples and change unsupervised AD datasets into supervised datasets. Teacher-student architecture is another useful approach. The teacher network distills knowledge to the student network by extracting features from normal samples. While the teacher network and the student network perform differently when producing abnormal samples and they detect anomalies through this characteristic [5, 14]. Normalizing flow methods map samples into a Gaussian distribution, while abnormal samples deviate from this distribution [15, 28]. Methods based on memory banks are simple but effective, whose ideas come from the k-nearest neighbors (KNN) algorithm. They store features of normal samples and calculate the distance between the features of test samples and the features of normal samples during testing to determine whether the samples are abnormal [13, 26]. As for reconstruction-based



**Figure 2: Simple examples of 3D anomaly detection on MVTEC 3D-AD [6]. The first and second rows are RGB images and point clouds, and the third to fifth rows show predicted results of PatchCore+FPFH [18], M3DM [33], and EasyNet (ours), respectively. Ground Truth refers to the anomaly regions.**

methods, most of them are similar in structure. They synthesize abnormal samples and restore abnormal samples to normal samples. For example, DRAEM [39] and NSA [31] synthesize abnormal samples in image level, while DSR [41] and UniAD [37] synthesize abnormal samples in feature level.

Generally, most of the 2D-AD methods use the pretrained model of natural images to extract RGB’s features while they don’t process depth information, so it is difficult to apply to 3D-AD directly. There is a certain gap between these two, and our method tries to get rid of this dependence so that 2D-AD can smoothly transition to 3D-AD.

## 2.2 3D Anomaly Detection

Different from 2D-AD, 3D-AD is a new research topic since the publication of MVTEC 3D-AD [6]. As shown in Figure 2, 3D-AD is a more challenging but also more promising research direction. The effective use of depth information can greatly improve detection accuracy in specific scenarios. On the other hand, how to integrate depth information and prevent it from interfering with RGB information is the current difficulty.

Bergmann *et al.* [7] introduce a point-cloud feature extraction network of the teacher-student model. During training, the features extracted by the student network and the teacher network are forced to be consistent. During the test, the differences between the features extracted by the teacher-student model are compared to locate anomalies. Horwitz *et al.* [17] combine hand-crafted 3D descriptors with the KNN framework, a classic AD approach. These two methods are efficient, but with poor performance. AST [29] gets a better result in MVTEC 3D-AD. However, it only uses depth information to remove the background and still uses the 2D-AD method to detect anomalies and the depth information about items is ignored. Similar to BTF, but M3DM [33] extracts features from point clouds and RGB images, respectively, and fuses them to make a decision, which has a better performance than treating RGB and

depth as six-channel images as BTF. The visualization effect of M3DM is shown in the fourth row of Figure 2. CPMF [10] also adopts the KNN paradigm, but the difference lies in the fact that the authors project the point cloud from different angles into 2D images and fuse the 2D image information obtained for detection.

In summary, existing 3D-AD models either suffer from poor performance or reliance on pretrained models and memory banks. In contrast, EasyNet is simple, effective, and without relying on pretrained models or memory banks. It achieves SOTA performance outperforming all previous methods without pre-training.

## 3 APPROACH

### 3.1 Problem Definition and Challenges

Our 3D-AD setting is similar to M3DM [33] and AST [30] and can be formally stated as follows. Given a set of training examples  $\mathcal{T} = \{t_i\}_{i=1}^N$ , in which  $\{t_1, t_2, \dots, t_N\}$  are the normal samples and each of them consists of paired images, RGB image  $I_{rgb}$  and depth image  $I_{depth}$ . In addition,  $\mathcal{T}_n$  belongs to a certain category,  $c_j$ ,  $c_j \in \mathcal{C}$ , where  $\mathcal{C}$  denotes the set of all categories. During testing, given a normal or abnormal sample from a target category  $c_j$ , the AD model should predict whether or not the test 3D object is anomalous and localize the anomaly region if the anomaly is detected.

The following are the main challenges. (1) Information on normal samples is limited, each category’s training dataset only contains normal samples, i.e., no pixel-level annotations of  $I_{rgb}$  and  $I_{depth}$ . (2) It is difficult to find an effective multi-modality fusion way for anomalies that may appear in RGB, depth, or both. Simply fusing these features may negatively impact overall AD performance. (3) Real-world applications have limited storage space, so it is impractical to build a model that uses large pretrained models and memory banks.

### 3.2 EasyNet

This section provides a complete description of EasyNet. As illustrated in Figure 3, the proposed model comprises a multi-scale Multi-modality Reconstruction Network (MRN), a multi-scale Multi-modality Segmentation Network (MSN) and an attention-based information entropy fusion module, with the fusion network being exclusively applied during reasoning stages. The following sections elaborate on the design and functionality of each module.

**3.2.1 Multi-modality Reconstruction Network (MRN).** The multi-modality reconstruction network establishes a task of image reconstruction. In this task, the network reconstructs the original image from an artificially corrupted image obtained from the simulator. The network is designed as an encoder-decoder structure to transform the local features of the input image into a mode that more closely resembles the normal sample distribution.

The framework of the simulator is depicted in Figure 4. We generate a foreground mask on the original depth image and apply a mask operation on the randomly generated Berlin noise figure. Our empirical evaluation reveals that only adding foreground noise exclusively assists the network in recognizing the noise on the foreground object rapidly. Then, the Berlin noise map undergoes binarization to produce positive and negative mask maps. Both normal and original RGB images undergo weighting, along with

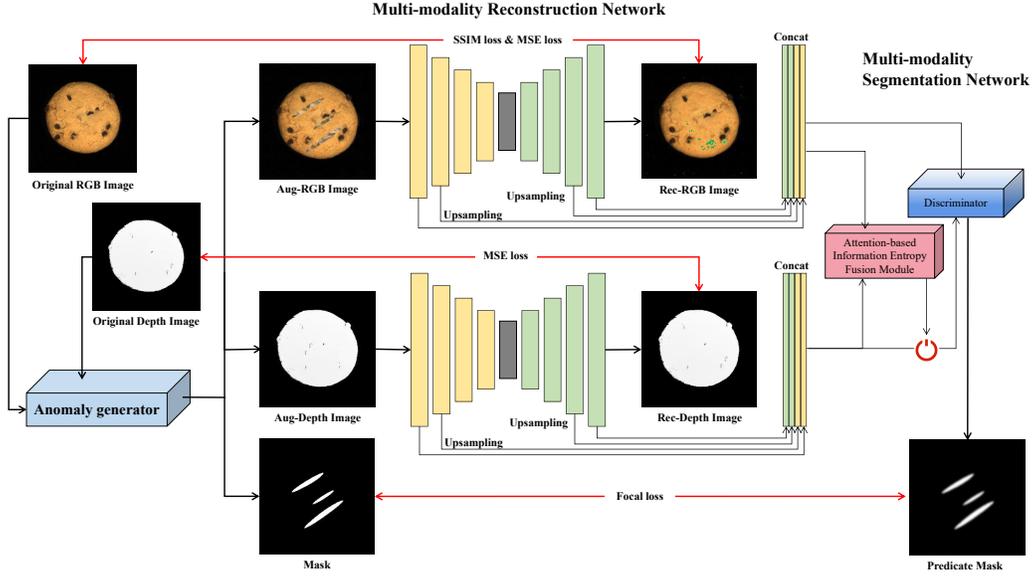


Figure 3: Total architecture of EasyNet. EasyNet consists of three main components. (1) Anomaly Generator adds Berlin noise to the original multi-modal images to simulate synthetic abnormal and normal images (RGB and depth). (2) The Multi-modality Reconstruction Network (MRN) constructs a reconstruction task to restore the enhanced synthetic anomaly images to RGB and depth images without anomalies while obtaining multi-modal feature information from multiple layers, where two layers are used as an example. (3) The Multi-modality Segmentation Network (MSN) fuses the extracted multi-modal features by utilizing an attention-based information entropy fusion module, which is fully open during training and takes control of feature flow by calculating the self-attention information entropy score for feature information from multiple modes during inference. For the reconstruction task, we use SSIM loss [34] and MSE loss to calculate the reconstruction loss for RGB images, while only using MSE loss for depth images. Moreover, Focal loss [22] is used to calculate the pixel classification task loss.

the Berlin noise map and depth image. Finally, the resulting outputs include RGB and depth images with anomalies and masks.

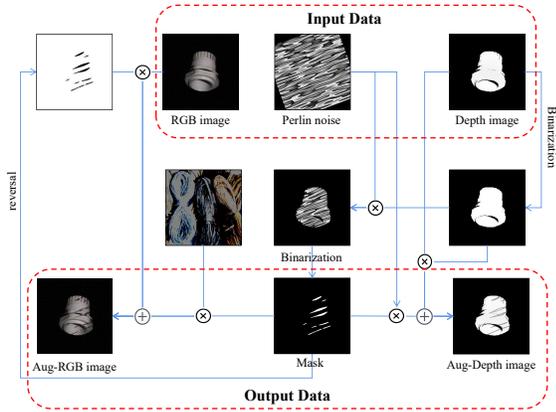


Figure 4: Illustration of anomaly generation processing.

In the reconstruction task, we used the classic  $L_2$  to reduce perceptual differences in RGB image reconstruction, we used the SSIM loss function in the RGB image reconstruction task. In our experiment, it is also found that spatial variation and multi-scale features

have limited and even negative effects on depth images. Therefore, the final image reconstruction loss function should be:

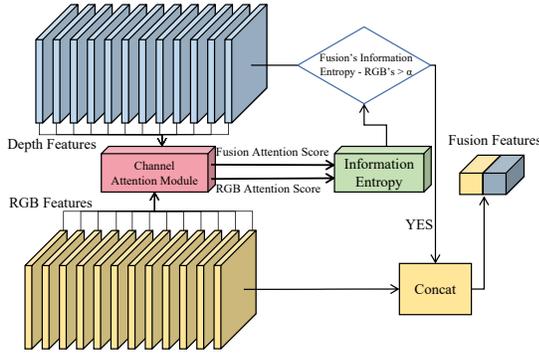
$$\begin{aligned} L_{rec}(I, I_r) &= L_{rec}^{RGB}(I, I_r) + L_{rec}^{depth}(I, I_r) \\ &= \lambda_1 L_{SSIM}^{RGB}(I, I_r) + \lambda_2 l_2^{RGB}(I, I_r) + \lambda_3 l_2^{depth}(I, I_r), \end{aligned} \quad (1)$$

where  $\lambda_1, \lambda_2, \lambda_3$  are loss balancing hyper-parameters, and all are set to 1 in our experimental setting.

**3.2.2 Multi-modality Segmentation Network (MSN).** The MSN evaluates the normality of each time slot ( $H, W$ ). Similar to DRAEM [39], the training set samples are processed by the simulator and the discriminator performs mask identification by identifying the input of enhanced images and reconstructed images. The difference between DRAEM and EasyNet refers to the *supplementary materials*. EasyNet extracts multi-layer features evaluated by discriminators through an MRN. MSN utilizes multi-layer reconstruction features and enhanced image features, which come from our assumption that some features that deviate from the normal distribution will be removed gradually with the deepening of the multi-modality reconstruction network. By comparing the difference of eigenvalues before and after removal, the locations of anomalies can be obtained.

When extracting reconstruction features and enhancing image features of multiple layers, we mainly adopt the first three layers of shallow networks of MRN and the last three layers of reconstructed

features and carry out up-sampling operations to adapt for features of multiple layers. Moreover, we conduct ablation experiments. As shown in Section 4.2.4, experiments show that when two-layer features are adopted, both the accuracy and computing cost of the network are optimized. We use a two-layer multi-layer perceptron (MLP) to process multi-layer scale features extracted from RGB and depth images respectively. Finally, we use another two-layer MLP structure to combine the features of the two modes and perform positive and negative discriminations for each pixel in the image. As shown in Section 4.2, the proposed straightforward strategy is successful in reaching its goal.



**Figure 5: Illustration of Attention-based Information Entropy Fusion Module.**

**3.2.3 Attention-based Information Entropy Fusion Module.** As noted in Section 3.1, an anomaly may occur solely in pure RGB or depth images, or both. The direct combination of both features may diminish the overall performance of AD and lead to an inverse outcome. So we generate multi-channel self-attention scores from input features in the input layer of MSN’s multi-layer perceptron module, as shown in Figure 5. We then compare the information entropy of the channel that integrates RGB and depth features with that of the channel integrating only pure RGB features. We hypothesize that the greater the information entropy of the channel attention score, the richer the feature knowledge it contains. If fusion features enhance the information gain beyond RGB features, it could positively affect the performance of the results. The experimental results presented in Section 4.2.3 provide support for our theory. The mathematical representation of this process is shown in Formula 2.

$$F_{fusion} = \begin{cases} F_{RGB} + F_{depth}, & f_{IE}(F_{RGB} + F_{depth}) > f_{IE}(F_{RGB}) + \alpha \\ F_{RGB}, & f_{IE}(F_{RGB} + F_{depth}) \leq f_{IE}(F_{RGB}) + \alpha \end{cases} \quad (2)$$

where  $F_{fusion}$  represents the features after fusion,  $F_{RGB}$  represents the features of RGB,  $F_{depth}$  represents the features of depth,  $f_{IE}(\cdot)$  represents the function of calculating information entropy, and  $\alpha$  represents the threshold adjustment factor.

When calculating the loss between the predicted mask and the ground truth mask, we use the Focal Loss [22] function (Formula 3), which could well solve the problem of sample imbalance in the single-class classification of pixels.

$$L_{focal}(M, M_{out}) = -\alpha_t (1 - p_t)^\gamma \log(p_t), \quad (3)$$

where  $\alpha_t$  is a scaling factor related to class  $t$ ,  $\gamma$  is an adjustable parameter,  $p_t$  corresponds to the predicted classification of pixel points, the abnormal category is 1, and the normal category is 0.

To sum up, EasyNet optimization objectives and tasks are reconstruction loss and classification loss. Finally, the overall loss of the network during training is as follows:

$$\begin{aligned} L_{all}(I, I_r) &= L_{rec}^{RGB}(I, I_r) + L_{rec}^{depth}(I, I_r) + L_{focal}(M, M_{out}) \\ &= \lambda_1 L_{SSIM}^{RGB}(I, I_r) + \lambda_2 L_2^{RGB}(I, I_r) \\ &\quad + \lambda_3 L_2^{depth}(I, I_r) + \lambda_4 L_{focal}(M, M_{out}), \end{aligned} \quad (4)$$

where  $\lambda_1, \lambda_2, \lambda_3$ , and  $\lambda_4$  are loss balancing hyper-parameters. EasyNet aims to meet the objectives of optimizing anomaly detection and reconstruction tasks while training, so we optimize the above objectives by assigning weights to different losses. All four  $\lambda$  are set to 1 in our experimental setting.

**3.2.4 Algorithms.** The EasyNet is implemented as Algorithm 1. when training, images  $I_{rgb}$  and  $I_{depth}$  are enhanced by the anomaly generator  $\Phi_{ag}$  to produce augmented images  $A_{rgb}$  and  $A_{depth}$  respectively. The Multi-modality Reconstruction Network  $\Phi_{rec}$  extracts multi-scale features ( $F_{rgb}, F_{depth}$ ) and generate reconstructed images ( $R_{rgb}, R_{depth}$ ) from these augmented images and origin images. The Multi-modality Segmentation Network  $\Phi_{seg}$  generates an anomaly score maps  $M$  and  $M_{rgb}$  by fusion and pure RGB features. When inferring, the function  $\Phi_{ai}$  generates corresponding self-attention information entropy scores from both RGB and RGB-D channels to combine RGB and depth features.

## 4 EXPERIMENTS

### 4.1 Experimental Details

**4.1.1 Datasets.** We mainly used MVTEC 3D-AD [6] and Eyescandies [8] data sets in the experiment. MVTEC 3D-AD dataset is the data set of the real scene, and Eyescandies is the data set of the virtual synthesis. More detailed introduction to these two datasets please refer to the *supplementary materials*.

**4.1.2 Evaluation Metrics.** Due to the unsupervised experimental setting, the common evaluation metrics we used include Area Under the Receiver Operating Characteristic Curve (AUROC) (I-AUROC and P-AUROC) and the Area Under the Precision-Recall curve (AUPR/AP), the explanation of I-AUROC and P-AUROC please refer to the *supplementary materials*.

### 4.2 Experimental Results and Analysis

**4.2.1 RGB+Depth on MVTEC 3D-AD and Eyescandies.** Table 1 and Table 2 clearly demonstrate that EasyNet achieves the state-of-the-art performance on MVTEC 3D-AD and Eyescandies without using a pretrained model and memory bank. Specifically, EasyNet outperforms AutoEncoder [9] and PatchCore+FPFH [18] in RGB+Depth setting of Eyescandies by a large margin, 20.7% and 6.9%, respectively. For MVTEC 3D-AD, AST [30] and M3DM [33] are the cutting-edge models in MVTEC 3D-AD. However, both of them use large pretrained models or memory banks. In specific, M3DM uses two pretrained models (Point Transformer and Vision Transformer) to extract the features from depth images and RGB images. In addition, M3DM employs two large memory banks (average 6.098 GB) to

**Table 1: I-AUROC score for anomaly detection of MVTec 3D-AD. The best is in red and the second best is in blue.**

Method	Bagel	Cable Gland	Carrot	Cookie	Dowel	Foam	Peach	Potato	Rope	Tire	Mean	Memory Bank Usage	pretrained Model Usage	
Pure Depth	Depth GAN [6]	0.530	0.376	0.607	0.603	0.497	0.484	0.595	0.489	0.536	0.521	0.523		
	Depth AE [6]	0.468	<b>0.731</b>	0.497	0.673	0.534	0.417	0.485	0.549	0.564	0.546	0.546		
	Depth VM [6]	0.510	0.542	0.469	0.576	0.609	0.699	0.450	0.419	0.668	0.520	0.546		
	Voxel GAN [6]	0.383	0.623	0.474	0.639	0.564	0.409	0.617	0.427	0.663	0.577	0.537		
	Voxel AE [6]	0.693	0.425	0.515	0.790	0.494	0.558	0.537	0.484	0.639	0.583	0.571		
	Voxel VM [6]	0.750	<b>0.747</b>	0.613	0.738	0.823	0.693	0.679	0.652	0.609	<b>0.690</b>	0.699		
	3D-ST [7]	0.862	0.484	0.832	0.894	0.848	0.663	0.763	0.687	<b>0.958</b>	0.486	0.748	✓	
	PatchCore+FPFH [18]	0.825	0.551	<b>0.952</b>	0.797	<b>0.883</b>	0.582	0.758	0.889	0.929	0.653	0.782	✓	✓
	AST [30]	<b>0.881</b>	0.576	<b>0.965</b>	<b>0.957</b>	0.679	<b>0.797</b>	<b>0.990</b>	<b>0.915</b>	<b>0.956</b>	0.611	<b>0.833</b>	✓	✓
	M3DM [33]	<b>0.941</b>	0.651	<b>0.965</b>	<b>0.969</b>	<b>0.905</b>	<b>0.760</b>	<b>0.880</b>	<b>0.974</b>	0.926	<b>0.765</b>	<b>0.874</b>	✓	✓
<b>EasyNet(ours)</b>	0.735	0.678	0.747	0.864	0.719	0.716	0.713	0.725	0.885	0.687	0.747			
Pure RGB	DifferNet [27]	0.859	0.703	0.643	0.435	0.797	0.790	0.787	0.643	0.715	0.590	0.696	✓	
	PADiM [12]	<b>0.975</b>	0.775	0.698	0.582	0.959	0.663	0.858	0.535	0.832	0.760	0.764	✓	
	PatchCore [25]	0.876	0.880	0.791	0.682	0.912	0.701	0.695	0.618	0.841	0.702	0.770	✓	
	STEPM [32]	0.930	0.847	0.890	0.575	0.947	0.766	0.710	0.598	0.965	0.701	0.793	✓	
	CS-Flow [16]	0.941	0.930	0.827	0.795	<b>0.990</b>	0.886	0.731	0.471	<b>0.986</b>	0.745	0.830	✓	
	AST [30]	0.947	0.928	0.851	<b>0.825</b>	0.981	<b>0.951</b>	<b>0.895</b>	0.613	<b>0.992</b>	<b>0.821</b>	<b>0.880</b>	✓	
	M3DM [33]	0.944	0.918	0.896	0.749	0.959	0.767	<b>0.919</b>	0.648	0.938	0.767	0.850	✓	
	SPADE [11]	0.771	0.793	0.760	0.531	0.848	0.683	0.646	0.460	0.879	0.502	0.687	✓	
	FastFlow [38]	0.624	0.472	0.654	0.694	0.501	0.667	0.595	0.632	0.816	0.731	0.639	✓	
	RD4AD [14]	<b>0.975</b>	<b>0.987</b>	<b>0.943</b>	0.575	<b>0.999</b>	0.830	0.863	0.618	0.984	<b>0.899</b>	0.867	✓	
	STPM [32]	0.899	0.706	0.796	0.486	0.512	0.678	0.502	<b>0.666</b>	0.962	0.581	0.679	✓	
	<b>EasyNet(ours)</b>	<b>0.982</b>	<b>0.992</b>	<b>0.917</b>	<b>0.953</b>	0.919	<b>0.923</b>	0.840	<b>0.785</b>	<b>0.986</b>	0.742	<b>0.904</b>		
	RGB+ Depth	Depth GAN [6]	0.538	0.372	0.580	0.603	0.430	0.534	0.642	0.601	0.443	0.577	0.532	
Depth AE [6]		0.648	0.502	0.650	0.488	0.805	0.522	0.712	0.529	0.540	0.552	0.595		
Depth VM [6]		0.513	0.551	0.477	0.581	0.617	0.716	0.450	0.421	0.598	0.623	0.555		
Voxel GAN [6]		0.680	0.324	0.565	0.399	0.497	0.482	0.566	0.579	0.601	0.482	0.517		
Voxel AE [6]		0.510	0.540	0.384	0.693	0.446	0.632	0.550	0.494	0.721	0.413	0.538		
Voxel VM [6]		0.553	0.772	0.484	0.701	0.751	0.578	0.480	0.466	0.689	0.611	0.609		
3D-ST [7]		0.950	0.483	<b>0.986</b>	0.921	0.905	0.632	0.945	<b>0.988</b>	0.976	0.542	0.833	✓	
PatchCore+FPFH [18]		0.918	0.748	0.967	0.883	0.932	0.582	0.896	0.912	0.921	<b>0.886</b>	0.865	✓	
AST [30]		0.983	0.873	<b>0.976</b>	<b>0.971</b>	0.932	0.885	<b>0.974</b>	<b>0.981</b>	<b>1.000</b>	0.797	<b>0.937</b>	✓	
M3DM [33]		<b>0.994</b>	<b>0.909</b>	0.972	<b>0.976</b>	<b>0.960</b>	<b>0.942</b>	<b>0.973</b>	0.899	0.972	<b>0.850</b>	<b>0.945</b>	✓	
<b>EasyNet(ours)</b>	<b>0.991</b>	<b>0.998</b>	0.918	0.968	<b>0.945</b>	<b>0.945</b>	0.905	0.807	<b>0.994</b>	0.793	0.926			

**Table 2: I-AUROC score for anomaly detection of all categories of Eyescandies. The best is in red and the second best is in blue.**

Method	Candy Cane	Chocolate Cookie	Chocolate Cookie	Confetto	Gummy Bear	Hazelnut Truffle	Licorice Sandwich	Lollipop	Marsh-mallow	Peppermint Candy	Mean	Memory Bank usage	pretrained Model Usage
Pure Depth	Raw [18]	<b>0.654</b>	0.510	0.563	0.451	0.433	0.454	0.472	0.515	0.626	0.366	0.504	✓
	HoG [18]	0.653	0.510	0.470	0.723	<b>0.728</b>	<b>0.520</b>	0.717	0.667	0.699	0.742	0.643	✓
	SIFT [18]	0.589	0.582	0.683	<b>0.885</b>	0.663	0.480	<b>0.778</b>	0.702	<b>0.746</b>	<b>0.790</b>	0.690	✓
	FPFH [18]	<b>0.670</b>	<b>0.710</b>	<b>0.805</b>	<b>0.806</b>	<b>0.748</b>	0.515	<b>0.794</b>	<b>0.757</b>	<b>0.765</b>	<b>0.757</b>	<b>0.733</b>	✓
	<b>EasyNet(ours)</b>	0.629	<b>0.716</b>	<b>0.768</b>	0.731	0.660	<b>0.710</b>	0.712	<b>0.711</b>	0.688	0.731	<b>0.706</b>	
Pure RGB	GANomaly [3]	0.485	0.512	0.532	0.504	0.558	0.486	0.467	0.511	0.481	0.528	0.507	
	DFKDE [2]	0.539	0.577	0.482	0.548	0.541	0.492	0.524	0.602	0.658	0.591	0.555	✓
	DFM [1]	0.532	0.776	0.624	0.675	0.681	0.596	0.685	0.618	0.964	0.770	0.692	✓
	STEPM [32]	<b>0.551</b>	0.654	0.576	0.784	<b>0.737</b>	<b>0.790</b>	0.778	0.620	0.840	0.749	0.708	✓
	PaDiM [12]	0.531	0.816	<b>0.821</b>	<b>0.856</b>	<b>0.826</b>	0.727	<b>0.784</b>	0.665	<b>0.987</b>	<b>0.924</b>	<b>0.794</b>	✓
	AutoEncoder [9]	0.527	<b>0.848</b>	0.772	0.734	0.590	0.508	0.693	<b>0.760</b>	0.851	0.730	0.701	
<b>EasyNet(ours)</b>	<b>0.723</b>	<b>0.925</b>	<b>0.849</b>	<b>0.966</b>	0.705	<b>0.815</b>	<b>0.806</b>	<b>0.851</b>	<b>0.975</b>	<b>0.960</b>	<b>0.858</b>		
RGB+ Depth	AutoEncoder [9]	0.529	0.861	0.739	0.752	0.594	0.498	0.679	0.651	0.838	0.750	0.689	
	PatchCore+FPFH [18]	<b>0.606</b>	<b>0.904</b>	<b>0.792</b>	<b>0.939</b>	<b>0.720</b>	<b>0.563</b>	<b>0.867</b>	<b>0.860</b>	<b>0.992</b>	<b>0.842</b>	<b>0.809</b>	✓
<b>EasyNet(ours)</b>	<b>0.737</b>	<b>0.934</b>	<b>0.866</b>	<b>0.966</b>	<b>0.717</b>	<b>0.822</b>	<b>0.847</b>	<b>0.863</b>	<b>0.977</b>	<b>0.960</b>	<b>0.869</b>		✓

**Table 3: Ablation studies on an attention-based information entropy fusion module. The best is in red and the second best is in blue.**

MVTEC 3D-AD	Evaluation	Bagel	Cable gland	Carrot	Cookie	Dowel	Foam	Peach	Potato	Rope	Tire	Mean
	Gate Close	<b>0.982</b>	0.992	<b>0.917</b>	<b>0.953</b>	0.919	0.923	0.840	<b>0.785</b>	<b>0.986</b>	<b>0.742</b>	0.904
	Gate Open	0.974	<b>0.996</b>	0.914	<b>0.968</b>	<b>0.941</b>	<b>0.938</b>	<b>0.882</b>	0.781	0.982	<b>0.793</b>	<b>0.917</b>
	Gate Control	<b>0.991</b>	<b>0.998</b>	<b>0.918</b>	<b>0.968</b>	<b>0.945</b>	<b>0.945</b>	<b>0.905</b>	<b>0.807</b>	<b>0.994</b>	<b>0.793</b>	<b>0.926</b>
Eyescandies	Evaluation	Candy Cane	Chocolate Cookie	Chocolate Cookie	Confetto	Gummy Bear	Hazelnut Truffle	Licorice Sandwich	Lollipop	Marsh-mallow	Peppermint Candy	Mean
	Gate Close	<b>0.723</b>	<b>0.925</b>	<b>0.849</b>	<b>0.966</b>	<b>0.705</b>	<b>0.815</b>	0.806	<b>0.851</b>	<b>0.975</b>	<b>0.960</b>	<b>0.857</b>
	Gate Open	0.722	0.919	0.827	<b>0.945</b>	0.685	0.813	<b>0.846</b>	0.850	<b>0.975</b>	<b>0.959</b>	0.854
	Gate Control	<b>0.737</b>	<b>0.934</b>	<b>0.866</b>	<b>0.966</b>	<b>0.717</b>	<b>0.822</b>	<b>0.847</b>	<b>0.863</b>	<b>0.977</b>	<b>0.960</b>	<b>0.869</b>

**Table 4: The ablation study for the number of fusion layers. The best is in red and the second best is in blue.**

RGB+Depth	Evaluation	Bagel	Cable gland	Carrot	Cookie	Dowel	Foam	Peach	Potato	Rope	Tire	Mean
Image AUC	1 layer	0.967	0.981	0.912	0.890	0.901	0.908	0.763	0.717	1.000	0.731	0.877
	2 layers	0.974	0.996	0.914	0.968	0.941	0.938	0.882	0.781	0.982	0.793	0.917
	3 layers	0.964	1.000	0.884	0.889	0.967	0.932	0.734	0.697	1.000	0.866	0.893
Image AP	1 layer	0.992	0.995	0.982	0.969	0.976	0.975	0.915	0.873	1.000	0.906	0.958
	2 layers	0.994	0.999	0.981	0.991	0.984	0.984	0.966	0.930	0.992	0.927	0.975
	3 layers	0.991	1.000	0.975	0.966	0.992	0.982	0.903	0.916	1.000	0.963	0.969
Pixel AUC	1 layer	0.875	0.861	0.963	0.678	0.716	0.998	0.969	0.921	0.928	0.861	0.877
	2 layers	0.935	0.941	0.971	0.897	0.885	0.997	0.992	0.888	0.955	0.728	0.919
	3 layers	0.904	0.717	0.836	0.651	0.809	0.997	0.914	0.942	0.986	0.970	0.873
Pixel AP	1 layer	0.020	0.025	0.097	0.015	0.026	0.803	0.081	0.032	0.133	0.163	0.139
	2 layers	0.039	0.062	0.188	0.025	0.034	0.562	0.298	0.034	0.144	0.031	0.142
	3 layers	0.041	0.043	0.084	0.008	0.121	0.648	0.017	0.014	0.249	0.134	0.136

**Algorithm 1** EasyNet pseudo-code

```

1: Input: train dataloader  $D_{train}$ , test dataloader  $D_{test}$ , epochs
2: Output: trained  $\Phi_{rec}$  and  $\Phi_{seg}$ ,  $M$ 
3: Initialization randomly:  $\Phi_{rec}$  and  $\Phi_{seg}$ 
4: /*Training time*/
5: for  $i = 0$  to epochs do
6:   for  $I_{rgb}, I_{depth}, M_{gt} \leftarrow D_{train}$  do
7:      $A_{rgb}, A_{depth} = \Phi_{ag}(I_{rgb}, I_{depth})$ 
8:      $F_{rgb}, F_{depth}, R_{rgb}, R_{depth} = \Phi_{rec}(A_{rgb}, A_{depth})$ 
9:      $F_{fusion} = \text{Concat}(F_{rgb}, F_{depth})$ 
10:     $M_{rgb} = \Phi_{seg}(F_{rgb})$ 
11:     $M = \Phi_{seg}(F_{fusion})$ 
12:     $L_{rgb} = \Phi_{loss}(R_{rgb}, I_{rgb}, M_{rgb}, M_{gt})$ 
13:     $L_{total} = \Phi_{loss}(R_{rgb}, R_{depth}, I_{rgb}, I_{depth}, M, M_{gt})$ 
14:     $L_{rgb.backward}, L_{total.backward}$ 
15:   end for
16: end for
17: /*Inference time*/
18: for  $I_{rgb}, I_{depth}, M_{gt} \leftarrow D_{test}$  do
19:    $A_{rgb}, A_{depth} = \Phi_{ag}(I_{rgb}, I_{depth})$ 
20:    $F_{rgb}, F_{depth}, R_{rgb}, R_{depth} = \Phi_{rec}(A_{rgb}, A_{depth})$ 
21:    $S_{rgb} = \Phi_{ai}(\text{Feature}_{rgb})$ 
22:    $S_{fusion} = \Phi_{ai}(F_{rgb}, F_{depth})$ 
23:   if  $S_{fusion} - S_{rgb} > \alpha$  then
24:      $F_{fusion} = \text{Concat}(F_{rgb}, F_{depth})$ 
25:   else
26:      $F_{fusion} = F_{rgb}$ 
27:   end if
28:    $M = \Phi_{seg}(F_{fusion})$ 
29: end for

```

store the features from depth images and RGB images. Due to strict storage limitations in practice, M3DM cannot be perfectly fit in real-world applications. Although PatchCore+FPFH [18] has a relatively small memory footprint (249.260 MB on average), the actual performance is not as good as EasyNet. The memory bank size of M3DM and PatchCore+FPFH on the MVTEC 3D-AD can be found in *supplementary materials*. AST also adopts two EfficientNet-B5 as the feature extractors for depth images and RGB images, which violate the storage limitation in IM. Moreover, according to Table 5, massive usage of pretrained models will slow down the inference

speed, which cannot meet the requirement of IM. Furthermore, the performance gap among EasyNet, AST and M3DM is very small, 2.1% and 1.2%. Hence, EasyNet is the best 3D-AD model to meet all the demands of IM.

**4.2.2 Pure RGB Performance.** In real-world applications, the limitation of the depth sensor is very large since the effective distance range of the depth sensor is 3 meters. In addition, most depth sensors are easily affected by the lighting condition. Hence, as for simulating the failure of the depth sensor, we conducted the experiment using only RGB images as the input. In the pure RGB branch of Table 1 and Table 2, EasyNet achieves state-of-the-art performance in I-AUROC in the pure RGB track. In particular, EasyNet outperforms in pure RGB of Eyescandies by a large margin, 15.7% to AutoEncoder [9] and 6.4% to PaDiM [12]. Moreover, EasyNet gets 5.97% better I-AUROC score than M3DM and 2.65% better I-AUROC score than AST. In total, The performance of EasyNet is robust even though the depth sensor is a failure.

**4.2.3 Attention-based Information Entropy Fusion Module.** Table 3 clearly illustrates the effectiveness of our proposed attention-based information entropy fusion module in EasyNet. The gate network is the key to control multi-feature fusion. We conduct the ablation studies on three options, Gate Close, Gate Open and Gate Control, respectively. Gate Close means that EasyNet only utilizes RGB images as the input and ignores depth information. Gate Open denotes that EasyNet uses both the RGB images and the depth images and combines their features for evaluation. Gate Control means that EasyNet adopts an attention-based information entropy fusion module to select depth features during the inference case. In Gate Open, we discover that depth information may work as the noise and degrade the total performance if we select both RGB features and depth features during inference, so we design an attention-based information entropy fusion module to select the feature for fusion, which can enhance the performance of all classes in MVTEC 3D-AD and Eyescandies.

**4.2.4 Ablation study on the number of fusion layers.** As described in Section 3.2.2, the MSN is utilized for segmenting anomalies by fusing enhanced multilayer image features and reconstructed ones. Table 4 presents that EasyNet performs optimally when incorporating features from only two layers. Specifically, the performance metrics of I-AUROC and P-AUROC are respectively improved by 4.36% and 4.57% with two layers compared to one layer, and they

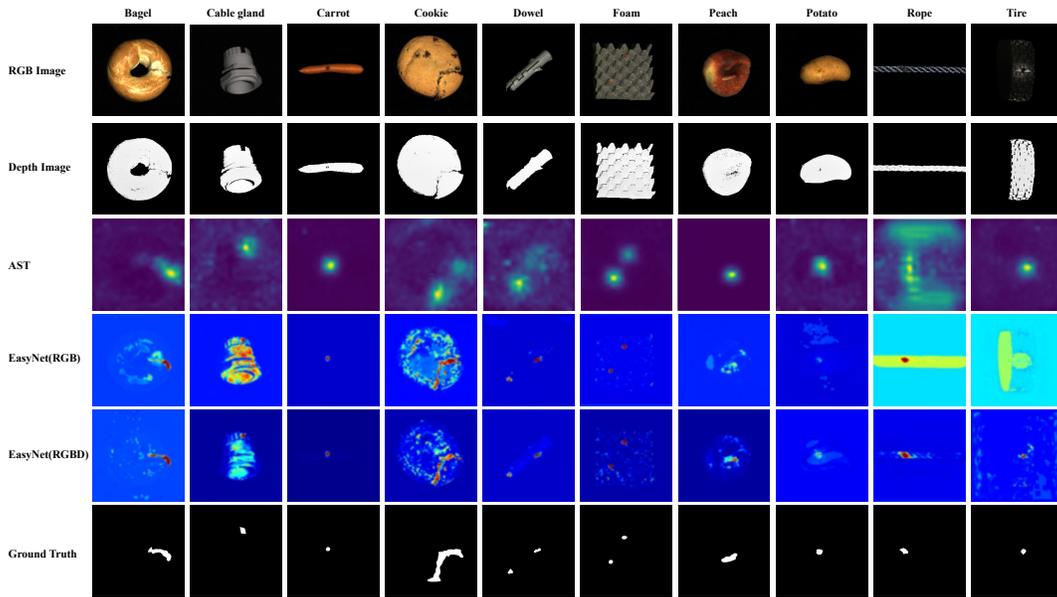


Figure 6: Visualizations on MVtec 3D-AD, which are obtained by AST [30], EasyNet (RGB) and EasyNet (RGB-D), respectively.

are further increased by 2.62% and 5.27% with three layers compared to one layer. These results indicate that the use of two feature layers can achieve the best performance in both anomaly detection and localization. In contrast, using three feature layers leads to performance degradation, which verifies our hypothesis. The deepening of the multi-modality reconstruction network can gradually eliminate some abnormal feature parts, whereas using three feature layers introduces more feature removal of abnormal parts, thereby leading to poor discriminator performance.

#### 4.2.5 Accuracy VS Inference Speed.

As we previously described in Section 1, inference speed is one of the important factors to be considered in IM. Since the real production lines need to check each product in real time. Table 5 shows EasyNet obtains the fastest speed among the cutting-edge anomaly detection methods without sacrificing performance. In particular, EasyNet gets 125% FPS better than AST and 93900% FPS better than M3DM. Regarding performance, the performance gap among EasyNet, AST and M3DM is very small, 2.1% and 1.2%. Therefore, EasyNet is the most deployment-friendly 3D-AD method for IM.

### 4.3 Visualization

Figure 6 visualizes the performance of EasyNet on MVtec 3D-AD, demonstrating the effectiveness of the proposed method. In Figure 6, compared to existing 3D anomaly detection methods (AST [30]), EasyNet can significantly reduce false positive rates and achieve higher segmentation accuracy. Furthermore, the fusion method reduces false positives for anomalies such as *cable gland*, compared

to using only RGB images. Anomalies in *peach* and *potato* are also more clearly visible on depth images, indicating the importance of using depth image information in industrial AD. In addition, EasyNet is not affected by the domain gap between natural and industrial images and has a higher inference speed than existing methods. Note that we put the visualization results of EasyNet on Eyescandies in *supplementary materials*.

## 5 CONCLUSIONS

This paper addresses a promising and challenging task, i.e., deployment-friendly 3D-AD and proposes an easy but effective neural network (termed as EasyNet) to achieve competitive performance without using large pretrained models and memory banks. Specifically, as for getting rid of large pretrained models and memory banks, EasyNet employs MRN to implicitly detect and reconstruct the anomalies with semantically plausible anomaly-free content, while keeping the non-anomalous regions of the input image unchanged. Meanwhile, EasyNet proposes an MSN to produce an accurate anomaly segmentation map from the concatenated reconstructed RGB images and depth images and their original appearances. In the test phase, EasyNet adopts a self-attention information entropy score in the early fusion stage to select the informative depth features before fusing with RGB features. To this end, EasyNet achieves the fastest inference speed without sacrificing performance.

## 6 ACKNOWLEDGMENTS

This work was supported by the National Key R&D Program of China (Grant NO. 2022YFF1202903), the National Natural Science Foundation of China (Grant NO. 62122035, 62206122), and the Key-Area Research and Development Program of Guangdong Province (2020B0101130003).

Table 5: Inference abilities.

Method	I-AUROC	FPS
BTF [17]	0.865	27.92
AST [30]	0.937	41.94
M3DM [33]	0.945	0.10
EasyNet(ours)	<b>0.926</b>	<b>94.55</b>

## REFERENCES

- [1] Nilesh A. Ahuja, Ibrahim J. Ndiour, Trushant Kalyanpur, and Omesh Tickoo. 2019. Probabilistic Modeling of Deep Features for Out-of-Distribution and Adversarial Detection. *ArXiv abs/1909.11786* (2019).
- [2] Samet Akçay, Dick Ameln, Ashwin Vaidya, Barath Lakshmanan, Nilesh Ahuja, and Utku Genc. 2022. Anomalib: A Deep Learning Library for Anomaly Detection. *arXiv:2202.08341 [cs.CV]*
- [3] Samet Akçay, Amir Atapour-Abarghouei, and T. Breckon. 2018. GANomaly: Semi-Supervised Anomaly Detection via Adversarial Training. In *Asian Conference on Computer Vision*.
- [4] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. 2019. MVTec AD–A comprehensive real-world dataset for unsupervised anomaly detection. *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2019), 9592–9600.
- [5] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. 2020. Uninformed students: Student-teacher anomaly detection with discriminative latent embeddings. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2020), 4183–4192.
- [6] Paul Bergmann, Xin Jin, David Sattlegger, and Carsten Steger. 2021. The MVTec 3D-AD Dataset for Unsupervised 3D Anomaly Detection and Localization. *ArXiv abs/2112.09045* (2021).
- [7] Paul Bergmann and David Sattlegger. 2022. Anomaly Detection in 3D Point Clouds using Deep Geometric Descriptors. *2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)* (2022), 2612–2622.
- [8] Luca Bonfiglioli, Marco Toschi, Davide Silvestri, Nicola Fioraio, and Daniele De Gregorio. 2022. The Eyecandies Dataset for Unsupervised Multimodal Anomaly Detection and Localization. *Proceedings of the Asian Conference on Computer Vision* (2022), 3586–3602.
- [9] Luca Bonfiglioli, Marco Toschi, Davide Silvestri, Nicola Fioraio, and Daniele De Gregorio. 2022. The Eyecandies Dataset for Unsupervised Multimodal Anomaly Detection and Localization. In *Asian Conference on Computer Vision*.
- [10] Yunkang Cao, Xiaohao Xu, and Weiming Shen. 2023. Complementary Pseudo Multimodal Feature for Point Cloud Anomaly Detection. *arXiv preprint arXiv:2303.13194* (2023).
- [11] Niv Cohen and Yedid Hoshen. 2020. Sub-Image Anomaly Detection with Deep Pyramid Correspondences. *ArXiv abs/2005.02357* (2020).
- [12] Thomas Defard, Aleksandr Setkov, Angélique Loesch, and Romaric Audigier. 2020. PaDiM: a Patch Distribution Modeling Framework for Anomaly Detection and Localization. In *ICPR Workshops*.
- [13] Thomas Defard, Aleksandr Setkov, Angélique Loesch, and Romaric Audigier. 2021. Padim: a patch distribution modeling framework for anomaly detection and localization. *International Conference on Pattern Recognition* (2021), 475–489.
- [14] Hanqiu Deng and Xingyu Li. 2022. Anomaly Detection via Reverse Distillation from One-Class Embedding. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2022), 9727–9736.
- [15] Denis Gudovskiy, Shun Ishizaka, and Kazuki Kozuka. 2022. Cflow-ad: Real-time unsupervised anomaly detection with localization via conditional normalizing flows. *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision* (2022), 98–107.
- [16] Denis A. Gudovskiy, Shun Ishizaka, and Kazuki Kozuka. 2021. CFLOW-AD: Real-Time Unsupervised Anomaly Detection with Localization via Conditional Normalizing Flows. *2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)* (2021), 1819–1828.
- [17] Eliahu Horwitz and Yedid Hoshen. 2022. Back to the feature: classical 3d features are (almost) all you need for 3d anomaly detection. *arXiv preprint arXiv:2203.05550* (2022).
- [18] Eliahu Horwitz and Yedid Hoshen. 2022. An Empirical Investigation of 3D Anomaly Detection and Segmentation. *ArXiv abs/2203.05550* (2022).
- [19] Jie Hu, Li Shen, Samuel Albanie, Gang Sun, and Enhua Wu. 2017. Squeeze-and-Excitation Networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42 (2017), 2011–2023.
- [20] Chun-Liang Li, Kihyuk Sohn, Jinsung Yoon, and Tomas Pfister. 2021. Cutpaste: Self-supervised learning for anomaly detection and localization. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), 9664–9674.
- [21] Wujin Li, Jiawei Zhan, Jinbao Wang, Bizhong Xia, Bin-Bin Gao, Jun Liu, Chengjie Wang, and Feng Zheng. 2022. Towards Continual Adaptation in Industrial Anomaly Detection. *Proceedings of the 30th ACM International Conference on Multimedia* (2022), 2871–2880.
- [22] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. 2017. Focal Loss for Dense Object Detection. *2017 IEEE International Conference on Computer Vision (ICCV)* (2017), 2999–3007.
- [23] Jiaqi Liu, Guoyang Xie, Jingbao Wang, Shangnian Li, Chengjie Wang, Feng Zheng, and Yaochu Jin. 2023. Deep Industrial Image Anomaly Detection: A Survey. *arXiv:2301.11514 [cs.CV]*
- [24] Zhikang Liu, Yiming Zhou, Yuansheng Xu, and Zilei Wang. 2023. SimpleNet: A Simple Network for Image Anomaly Detection and Localization. *arXiv preprint arXiv:2303.15140* (2023).
- [25] Karsten Roth, Latha Pemula, Joaquin Zepeda, Bernhard Schölkopf, Thomas Brox, and Peter Gehler. 2021. Towards Total Recall in Industrial Anomaly Detection. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2021), 14298–14308.
- [26] Karsten Roth, Latha Pemula, Joaquin Zepeda, Bernhard Schölkopf, Thomas Brox, and Peter Gehler. 2022. Towards total recall in industrial anomaly detection. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), 14318–14328.
- [27] Marco Rudolph, Bastian Wandt, and Bodo Rosenhahn. 2020. Same Same But DifferNet: Semi-Supervised Defect Detection with Normalizing Flows. *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)* (2020), 1906–1915.
- [28] Marco Rudolph, Bastian Wandt, and Bodo Rosenhahn. 2021. Same same but different: Semi-supervised defect detection with normalizing flows. *Proceedings of the IEEE/CVF winter conference on applications of computer vision* (2021), 1907–1916.
- [29] Marco Rudolph, Tom Wehrbein, Bodo Rosenhahn, and Bastian Wandt. 2022. Asymmetric Student-Teacher Networks for Industrial Anomaly Detection. *2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)* (2022), 2591–2601.
- [30] Marco Rudolph, Tom Wehrbein, Bodo Rosenhahn, and Bastian Wandt. 2022. Asymmetric Student-Teacher Networks for Industrial Anomaly Detection. *arXiv preprint arXiv:2210.07829* (2022).
- [31] Hannah M Schlüter, Jeremy Tan, Benjamin Hou, and Bernhard Kainz. 2022. Natural Synthetic Anomalies for Self-supervised Anomaly Detection and Localization. *European Conference on Computer Vision* (2022), 474–489.
- [32] Guodong Wang, Shumin Han, Errui Ding, and Di Huang. 2021. Student-Teacher Feature Pyramid Matching for Anomaly Detection. In *British Machine Vision Conference*.
- [33] Yue Wang, Jinlong Peng, Jiangning Zhang, Ran Yi, Yabiao Wang, and Chengjie Wang. 2023. Multimodal Industrial Anomaly Detection via Hybrid Fusion. *ArXiv abs/2303.00601* (2023).
- [34] Zhou Wang, Alan Conrad Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing* 13 (2004), 600–612.
- [35] Guoyang Xie, Jinbao Wang, Jiaqi Liu, Yaochu Jin, and Feng Zheng. 2023. Pushing the Limits of Fewshot Anomaly Detection in Industry Vision: Graphcore. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=xzmqxHdZAwO>
- [36] Guoyang Xie, Jinbao Wang, Jiaqi Liu, Jiayi Lyu, Y. Liu, Chengjie Wang, Feng Zheng, and Yaochu Jin. 2023. IM-IAD: Industrial Image Anomaly Detection Benchmark in Manufacturing. *ArXiv abs/2301.13359* (2023).
- [37] Zhiyuan You, Lei Cui, Yujun Shen, Kai Yang, Xin Lu, Yu Zheng, and Xinyi Le. 2022. A Unified Model for Multi-class Anomaly Detection. *arXiv preprint arXiv:2206.03687* (2022).
- [38] Jiawei Yu1, Ye Zheng, Xiang Wang, Wei Li, Yushuang Wu, Rui Zhao, and Liwei Wu. 2021. FastFlow: Unsupervised Anomaly Detection and Localization via 2D Normalizing Flows. *ArXiv abs/2111.07677* (2021).
- [39] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. 2021. Draem—a discriminatively trained reconstruction embedding for surface anomaly detection. *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), 8330–8339.
- [40] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. 2021. DRAEM – A discriminatively trained reconstruction embedding for surface anomaly detection. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)* (2021), 8310–8319.
- [41] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. 2022. DSR–A dual subspace re-projection network for surface anomaly detection. *arXiv preprint arXiv:2208.01521* (2022).

## 7 SUPPLEMENTARY

### 7.1 Difference between DRAEM and EasyNet

As stated in EasyNet, we learn from DRAEM's reconstruction and abnormal image generation methods, but the architecture between EasyNet and DRAEM is quite different. Firstly, the reconstruction network of EasyNet employs the multi-layer and multi-scale feature information before and after reconstruction, which is also highlighted in the ablation studies, while the counterpart of DREAM can only be segmented by the images before and after reconstruction. Secondly, our experiments find that 2-layer MLP usage is sufficient to effectively segment the abnormal region without the need to use a large U-Net like DREAM. Finally, EasyNet pays more attention to the fusion and segmentation of multi-modal features, while DREAM only focuses on RGB features.

### 7.2 Datasets

**MVTec 3D-AD** [6] includes ten categories and a total of 2,656 training samples along with 1,137 testing samples. The 3D scans in this dataset were acquired via a structured-light-powered industrial scanner that captured the x, y, and z coordinates of the target object. Additionally, RGB data is also collected at the same time for each point in the cloud. To process the 3D data accurately, it is crucial to remove all the background noise. A RANSAC algorithm is employed to estimate the background plane, ensuring that points within 0.005 distances were eliminated without disturbing the RGB data. However, their corresponding pixels in the RGB image were set to zero. This step minimized disturbances while enhancing the accuracy of anomaly detection.

**Eyescandies** [8] is a novel synthetic dataset comprising ten different categories of candies rendered in a controlled environment. Bonfiglioli *et al.* [9] generated item instances through modeling software and collected relevant data. The dataset provides depth and RGB images in an industrial conveyor scenario. The ten categories of candies show different challenges, such as complex textures, self-occlusions, and specularities. By controlling the lighting conditions and parameters of a procedural rendering pipeline in the modeling software, the authors of the dataset produced datasets containing complex instances with varying conditions. Similar to MVTec 3D-AD, the training dataset only consists of normal samples, while the testing dataset consists of normal and abnormal samples.

### 7.3 I-AUROC and P-AUROC

I-AUROC notes Image-level AUROC and P-AUROC notes Pixel-level AUROC. Image-level AUROC is based on the area under the ROC curve of the entire image, where the horizontal axis represents a false positive rate and the vertical axis represents a true positive rate. Each point represents the performance of a model under different classification thresholds. Different points can be obtained by changing the image classification threshold to plot the entire ROC curve. Image-level AUROC is used to evaluate the classification quality of the overall image. Pixel-level AUROC reflects the segmentation accuracy of a model at the pixel level based on the area under the ROC curve of each pixel.

### 7.4 Implementation Details

This section presents the implementation details of our experiments.

MRN uses the "UNet-like" structure as the primary network with intermediate skip operations subtracted primarily from the original UNet. The input image is resized to  $256 \times 256$ , and the abnormal and normal images are allocated according to a 1:1 ratio. The abnormal images are applied with Berlin noise [40] added on top of normal images.

For MSN, the two-layer MLP network is used to fuse different scale features of RGB and depth features. In the experiment of Section 4.2.3, the two layers MLPs network are employed to fuse different modal features. The input and output features of all the MLPs have the same size of  $256 \times 256$ . And a SE block [19] is utilized for Attention-based Information Entropy Fusion Module to score channel attention for both modes.

The training process adopted the Adam optimizer with a learning rate of 0.002, which is dynamically adjusted twice, at  $0.8 \times$  epochs and  $0.9 \times$  epochs, with a multiplier factor of 0.2. The batch size is set to 8. Finally, we report the best anomaly detection results obtained after 800 training steps of MRN.

### 7.5 Memory of Mainstream Methods

Memory-based methods are also deployed in real-world applications, so the use of memory banks should not be restricted. Therefore, we calculate in Table 1 the memory size required by the current mainstream methods to use the RGBD method in the MVTec 3D-AD dataset.

The M3DM method occupies a large amount of memory, resulting in a large amount of data transmission and loading time in the actual reasoning process. Although PatchCore+FPFH occupies a smaller memory, its performance is worse than EasyNet. The size of the memory bank is mainly affected by the training set, which indicates that they are using large memory sizes and hindering their practical application.

### 7.6 AUPRO Score of MVTec 3D-AD

In addition to I-AUROC, we also calculated the EasyNet model's AUPRO performance on the MvTec 3D-AD dataset, and Table ?? clearly shows that EasyNet achieves state-of-the-art performance on MvTec 3D-AD without the use of pre-trained models and memory banks and is slightly worse performance than 3D-ST using pre-trained models. In the experimental settings of Pure RGB and Pure Depth, our model does not take priority, but it also proves to a certain extent that our Attention-based Information Entropy Fusion Module plays a role, which blends the information of the two modes well.

### 7.7 Visualizations on Eyescandies

Figure 7 visualizes the performance of EasyNet on Eyescandies dataset, proving the effectiveness of the proposed method. In Figure 7, in the EasyNet dataset, although RGB images are interfered with by various lighting conditions, the anomaly heat map result of RGB image generated by EasyNet has a high coincidence degree with the ground truth. Compared with only RGB images, the contour of the anomaly heat map generated by EasyNet fusion method is clearer.

Method	Bagel	Cable Gland	Carrot	Cookie	Dowel	Foam	Peach	Potato	Rope	Tire	Mean
PatchCore+FPFH [18](MB)	228.984	209.281	268.403	197.083	270.282	221.479	338.790	281.547	279.667	197.083	249.260
M3DM [33](GB)	5.580	5.122	6.569	4.824	6.615	5.421	8.292	6.891	6.845	4.824	6.098

Table 6: The size of the memory of mainstream methods using memory bank on the MVTec 3D-AD.

Method	Bagel	Cable Gland	Carrot	Cookie	Dowel	Foam	Peach	Potato	Rope	Tire	Mean	memory bank use	pretrain model use	
Pure Depth	Depth GAN [6]	0.111	0.072	0.212	0.174	0.160	0.128	0.003	0.042	0.446	0.075	0.143		
	Depth AE [6]	0.147	0.069	0.293	0.217	0.207	0.181	0.164	0.066	0.545	0.142	0.203		
	Depth VM [6]	0.280	0.374	0.243	0.526	0.485	0.314	0.199	0.388	0.543	0.385	0.374		
	Voxel GAN [6]	0.440	0.453	0.875	0.755	0.782	0.378	0.392	0.639	0.775	0.389	0.583		
	Voxel AE [6]	0.260	0.341	0.581	0.351	0.502	0.234	0.351	0.658	0.015	0.185	0.348		
	Voxel VM [6]	0.453	0.343	0.521	0.697	0.680	0.284	0.349	0.634	0.616	0.346	0.492		
	FPFH [18]	<b>0.973</b>	<b>0.879</b>	<b>0.982</b>	<b>0.906</b>	<b>0.892</b>	<b>0.735</b>	<b>0.977</b>	<b>0.982</b>	<b>0.956</b>	<b>0.961</b>	<b>0.924</b>	✓	
	M3DM [33]	<b>0.943</b>	<b>0.818</b>	<b>0.977</b>	<b>0.882</b>	<b>0.881</b>	<b>0.743</b>	<b>0.958</b>	<b>0.974</b>	<b>0.950</b>	<b>0.929</b>	<b>0.906</b>	✓	✓
<b>EasyNet(ours)</b>	0.160	0.030	0.680	0.759	0.758	0.069	0.225	0.734	0.797	0.509	0.472			
Pure RGB	PatchCore [25]	<b>0.901</b>	<b>0.949</b>	<b>0.928</b>	<b>0.877</b>	<b>0.892</b>	0.563	0.904	<b>0.932</b>	<b>0.908</b>	<b>0.906</b>	<b>0.876</b>	✓	✓
	M3DM [33]	<b>0.944</b>	<b>0.918</b>	0.896	<b>0.749</b>	<b>0.959</b>	<b>0.767</b>	<b>0.919</b>	0.648	<b>0.938</b>	<b>0.767</b>	<b>0.850</b>	✓	✓
	<b>EasyNet(ours)</b>	0.751	0.825	<b>0.916</b>	0.599	0.698	<b>0.699</b>	<b>0.917</b>	<b>0.827</b>	0.887	0.636	0.776		
RGB+ Depth	Depth GAN [6]	0.421	0.422	0.778	0.696	0.494	0.252	0.285	0.362	0.402	<b>0.631</b>	0.474		
	Depth AE [6]	0.432	0.158	0.808	0.491	<b>0.841</b>	0.406	0.262	0.216	0.716	0.478	0.481		
	Depth VM [6]	0.388	0.321	0.194	0.570	0.408	0.282	0.244	0.349	0.268	0.331	0.335		
	Voxel GAN [6]	0.664	0.620	0.766	<b>0.740</b>	0.783	0.332	0.582	0.790	0.633	0.483	0.639		
	Voxel AE [6]	0.467	<b>0.750</b>	0.808	0.550	0.765	0.473	0.721	<b>0.918</b>	0.019	0.170	0.564		
	Voxel VM [6]	0.510	0.331	0.413	0.715	0.680	0.279	0.300	0.507	0.611	0.366	0.471		
	3D-ST [7]	<b>0.950</b>	0.483	<b>0.986</b>	<b>0.921</b>	<b>0.905</b>	<b>0.632</b>	<b>0.945</b>	<b>0.988</b>	<b>0.976</b>	<b>0.542</b>	<b>0.833</b>		✓
	<b>EasyNet(ours)</b>	<b>0.839</b>	<b>0.864</b>	<b>0.951</b>	0.618	0.828	<b>0.836</b>	<b>0.942</b>	0.889	<b>0.911</b>	0.528	<b>0.821</b>		

Table 7: AUPRO score for anomaly detection of all categories of MVTec 3D-AD.

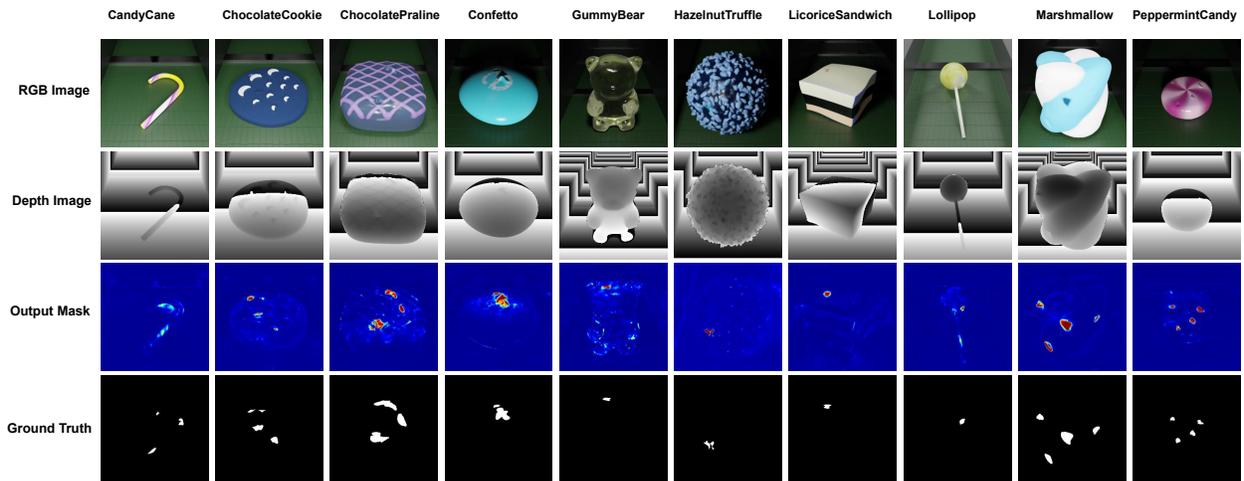


Figure 7: Visualizations on Eyescandies, which are obtained by EasyNet (RGB) and EasyNet (RGB-D).