



# GLAD: Towards Better Reconstruction with Global and Local Adaptive Diffusion Models for Unsupervised Anomaly Detection

Hang Yao<sup>1</sup>, Ming Liu<sup>1,2</sup>(✉), Haolin Wang<sup>1</sup>, Zhicun Yin<sup>1</sup>,  
Zifei Yan<sup>1,2</sup>, Xiaopeng Hong<sup>1</sup>, and Wangmeng Zuo<sup>1,2</sup>

[yaohang\\_1@outlook.com](mailto:yaohang_1@outlook.com), [csmliu@outlook.com](mailto:csmliu@outlook.com), [cswmzuo@gmail.com](mailto:cswmzuo@gmail.com)

<sup>1</sup>Harbin Institute of Technology, Harbin, China

<sup>2</sup>Pazhou Lab Huangpu, Guangzhou, China

ECCV 2024

Anomaly detection (AD) aims to detect and locate abnormal patterns of objects:

- it challenging to collect **enough abnormal** samples for all anomaly types in situations.
- ever-changing product design and production processes, it is impossible to collect **all anomalies** in advance.

**unsupervised anomaly detection** (UAD) has drawn much attention with only normal samples required.



## Related works

---

**Embedding-based methods:** extract feature of images to evaluate abnormal.

1. Knowledge distillation-based methods: train **student** network with **normal** samples, features from the **pre-trained teacher** network are compared with features from the student network to detect and locate anomalies.
2. PaDiM builds **multivariate Gaussian distributions** for **patch features of normal samples** and uses Mahalanobis distance as the anomaly score.
3. PatchCore proposes a **memory bank** to save **features of normal images**, which are **compared** with feature maps of **test images** to distinguish the difference between normal and abnormal features.

**Reconstruction-based methods:** detected and located via the comparison between the given sample and its normal counterpart.

Based on the hypothesis that models trained on normal samples only can reconstruct normal images well. Anomalies can be detected by comparing the samples before and after reconstruction.

AE(early), GAN, transformer, UNet architecture

# Motivation

$$\mathbf{x}^a \xrightarrow{diff} \mathbf{x}_t^a \xrightarrow{gen} \hat{\mathbf{x}}^a$$

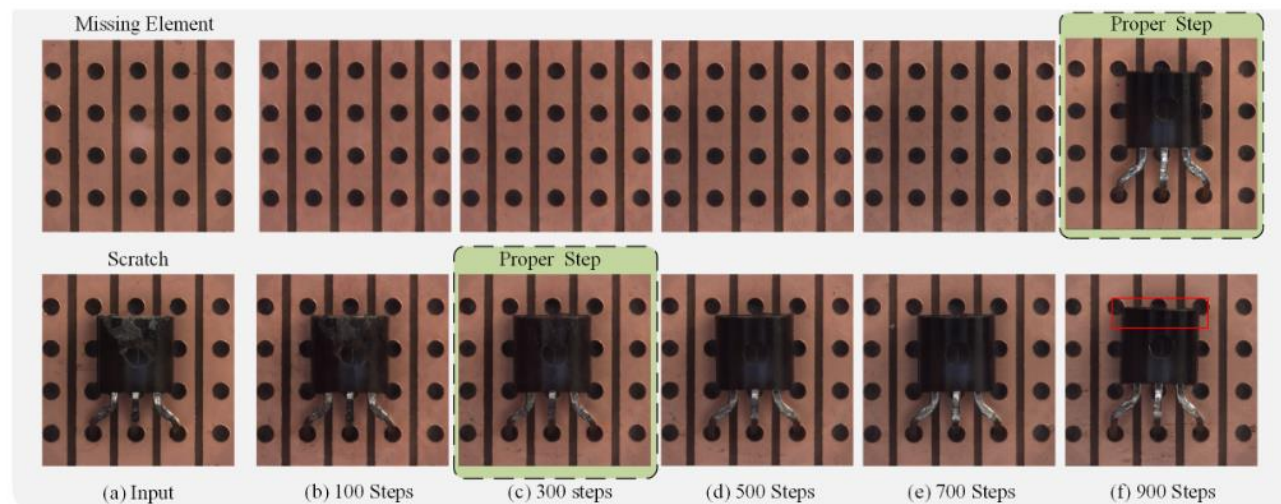
Diffusion models have prominent modeling ability.

During the training process , the diffusion model captures the distribution of normal samples only.

The same denoising step:

- Different anomalies is uneven.
- Less preserved details of the original

The anomaly noise inevitably deviates from the standard Gaussian distribution.



**Fig. 1:** Illustration of adaptive denoising process. For severe anomalies like missing elements, it requires a large number of denoising steps (900) to add the element back, while for small anomalies like scratch, 300 steps are already enough. Besides, setting a large enough denoising step (*e.g.*, 900) for all samples will affect the detail preservation. For example, in the area bounded by red lines, the position of the element is changed, which will be marked as anomalies during the comparison process.

# Method

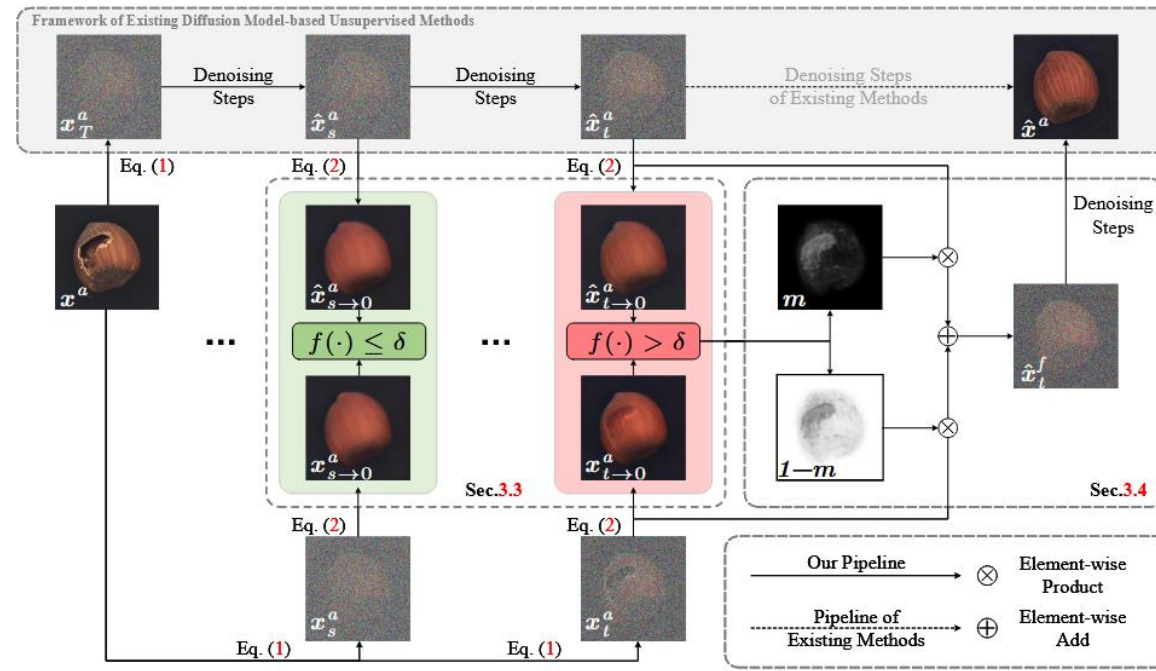
Inference:

Adaptive Denoising Step (ADS) : achieves a better trade-off between reconstruction quality and detail preservation ability.

Spatial-Adaptive Feature Fusion (SAFF): avoid reconstruction of normal regions.

Training:

Anomaly-oriented Training Paradigm (ATP): allow diffusion model to predict non-Gaussian noise at abnormal regions.



**Fig. 2:** The reconstruction pipeline of the proposed GLAD, including the Adaptive Denoising Steps (Sec. 3.3) and the Spatial-Adaptive Feature Fusion Scheme (Sec. 3.4).



## Method / Preliminary

Diffusion Process  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \bar{\alpha}_t = \prod_{i=1}^t \alpha_i$

Intermediate Result  $\mathbf{x}_{t \rightarrow 0} = \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(\mathbf{x}_t, t)),$

Generation Process  $\hat{\mathbf{x}}_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{\mathbf{x}}_{t \rightarrow 0} + \sqrt{1 - \bar{\alpha}_{t-1}} \epsilon_\theta(\hat{\mathbf{x}}_t, t).$

$$\mathbf{x}^a \xrightarrow{diff} \mathbf{x}_t^a \xrightarrow{gen} \hat{\mathbf{x}}^a$$

$$\|\hat{\mathbf{x}}^a - \mathbf{x}\|_\infty < \tau$$

detected and located by comparison between  $\hat{\mathbf{x}}^a$  and  $\mathbf{x}^a$

$$\mathbf{x}^a = \mathbf{x} + \mathbf{n}$$

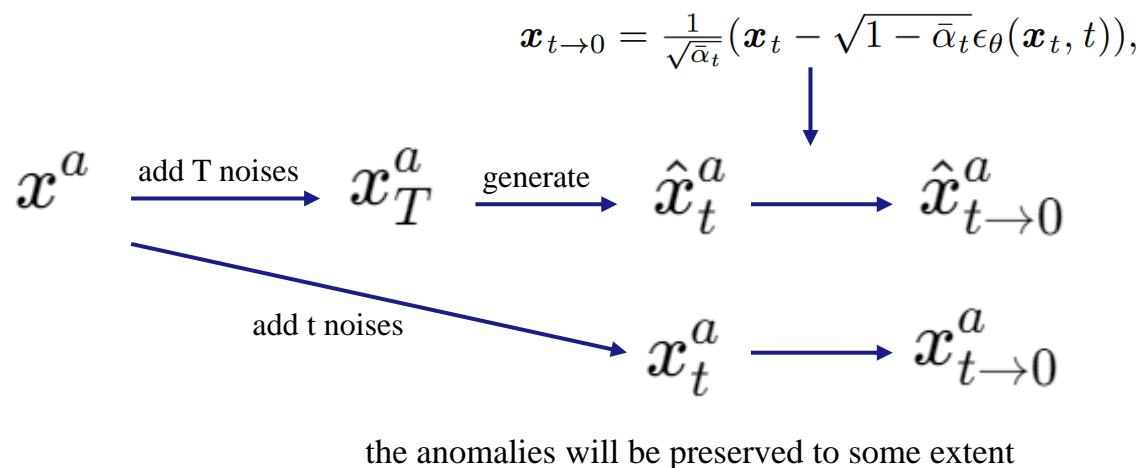
$$\begin{aligned} \mathbf{x}_t^a &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0^a + \sqrt{1 - \bar{\alpha}_t} \epsilon^a \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon^a + \sqrt{\bar{\alpha}_t} \mathbf{n}, \end{aligned}$$

$$\begin{aligned} \hat{\mathbf{x}}^a - \mathbf{x} &= g_t(\mathbf{x}_t^a) - g_t(\mathbf{x}_t) \\ &= g_t(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon^a + \sqrt{\bar{\alpha}_t} \mathbf{n}) - g_t(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon) \\ &\approx \sqrt{1 - \bar{\alpha}_t} (\epsilon^a - \epsilon) + \sqrt{\bar{\alpha}_t} \mathbf{n}, \end{aligned}$$

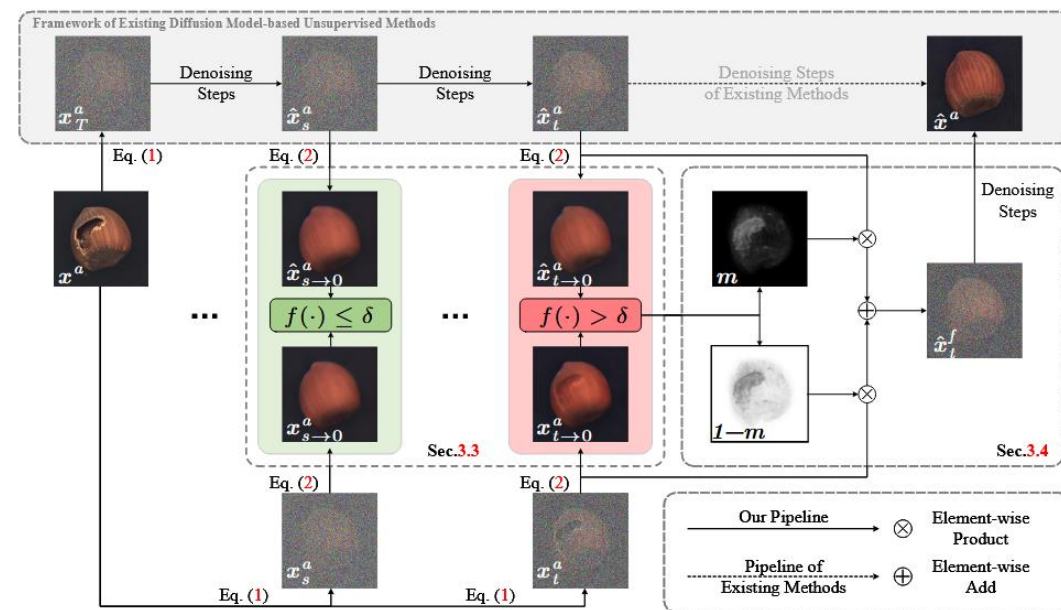
we make our efforts to reduce the errors for a better reconstruction quality.



## Method / Adaptive Denoising Steps (ADS)



(t + n) step of denoising,  
 normal regions are best preserved;  
 the anomalies can be reconstructed



**Fig. 2:** The reconstruction pipeline of the proposed GLAD, including the Adaptive Denoising Steps (Sec. 3.3) and the Spatial-Adaptive Feature Fusion Scheme (Sec. 3.4).

# Method / Spatial-Adaptive Feature Fusion (SAFF)

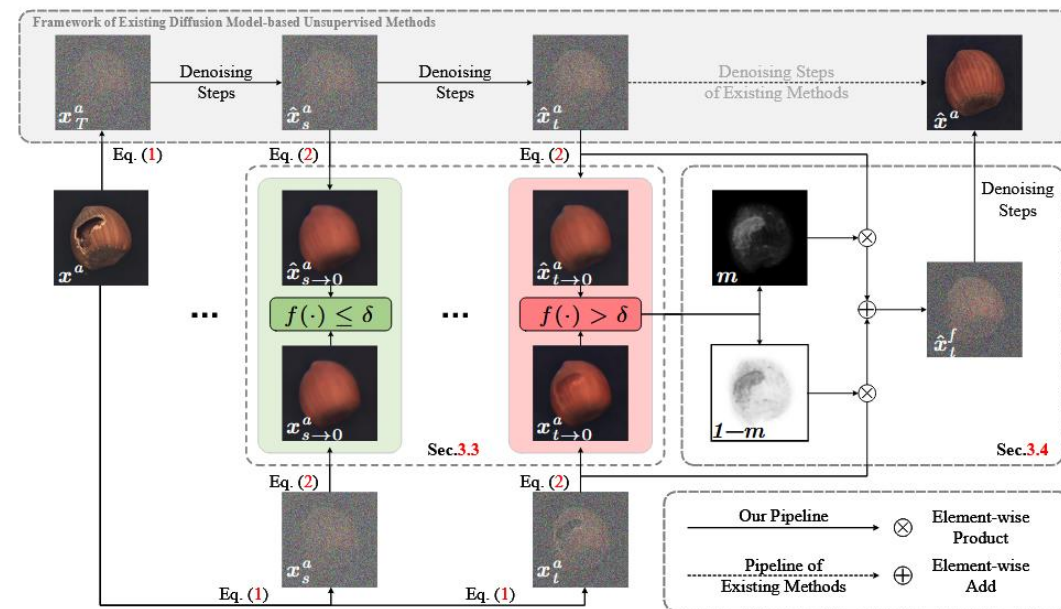
the whole image is sub-optimal

we need only to set a larger step for the abnormal regions.

mask  $m$ , which means the possibility for the pixels to be part of the anomalies

$$\hat{x}_t^f = m \cdot \hat{x}_t^a + (1 - m) \cdot x_t^a.$$

$$\begin{aligned} \hat{x}_t^f &= \sqrt{\bar{\alpha}_t} \hat{x}_{t \rightarrow 0}^f + \sqrt{1 - \bar{\alpha}_t} \epsilon \\ &= \sqrt{\bar{\alpha}_t} (m \cdot \hat{x}_{t \rightarrow 0}^a + (1 - m) \cdot x_{t \rightarrow 0}^a) + \sqrt{1 - \bar{\alpha}_t} \epsilon. \end{aligned}$$



**Fig. 2:** The reconstruction pipeline of the proposed GLAD, including the Adaptive Denoising Steps (Sec. 3.3) and the Spatial-Adaptive Feature Fusion Scheme (Sec. 3.4).





## Method / Anomaly-oriented Training Paradigm (ATP)

$$\hat{\mathbf{x}}^a - \mathbf{x} \stackrel{\infty}{\sim} \sqrt{1 - \bar{\alpha}_t}(\epsilon^a - \epsilon) + \sqrt{\bar{\alpha}_t} \mathbf{n} \rightarrow 0.$$

$$\epsilon^a \rightarrow \epsilon - \frac{\sqrt{\bar{\alpha}_t}}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{n}.$$

$$\begin{aligned} L_{ATP} &= \mathbb{E}_{(\mathbf{x}, \mathbf{x}^a) \sim p_{data}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t} [\|(\epsilon - \frac{\sqrt{\bar{\alpha}_t}}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{n}) - \epsilon^a\|_2] \\ &= \mathbb{E}_{(\mathbf{x}, \mathbf{x}^a) \sim p_{data}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t} [\|(\epsilon - \frac{\sqrt{\bar{\alpha}_t}}{\sqrt{1 - \bar{\alpha}_t}} (\mathbf{x}^a - \mathbf{x})) - \epsilon_\theta(\mathbf{x}_t^a, t)\|_2]. \end{aligned}$$

MemSeg [29] to synthesize abnormal samples with normal ones, which enables the training to proceed in an unsupervised manner



## Method / Anomaly Scoring and Map Construction

pre-trained model: DINO

$F_t \in \mathbb{R}^{c \times u \times v}$  features of test images

$F_r \in \mathbb{R}^{c \times u \times v}$  features of reconstructed images

$$M_l^{(i,j)}(F_t^l, F_r^l) = \min(1 - \langle F_t^{l(i,j)}, F_r^l \rangle),$$

$$M = \sum_l M_l(F_t^l, F_r^l),$$

Anomaly score: the whole image is the average of the top K maximum values of M .

We use the pre-trained latent diffusion model (LDM) [21] and fine-tune the UNet to adapt data for reconstruction.

DINO [5] with ViT-B/8 architecture is utilized as a feature extraction model



# Experiment

**Table 1:** Comparison with SOTA methods on MVTec-AD dataset. I-AUROC and P-AUROC are displayed in each entry. The best results among all methods are shown in bold, and the underlined results denote the best results among reconstruction-based methods.

Category	Embedding-based methods			Reconstruction-based methods					
	PatchCore [22]	RD4AD [8]	SimpleNet [17]	DRAEM [32]	OCR-GAN [14]	Lu <i>et al.</i> [18]	DiffAD [33]	DDAD [19]	Ours
Carpet	98.7/ <b>99.0</b>	98.9/98.8	<b>99.7</b> /98.2	97.0/95.5	<u>99.4</u> /-	-/97.7	98.3/98.1	99.3/ <u>98.7</u>	99.0/98.5
Grid	98.2/98.7	<b>100</b> /97.0	99.7/98.8	99.9/ <b>99.7</b>	99.6/-	-/95.6	<b>100</b> / <b>99.7</b>	<b>100</b> /99.4	<b>100</b> /99.6
Leather	<b>100</b> /99.3	<b>100</b> /98.6	<b>100</b> /99.2	<b>100</b> /98.5	97.1/-	-/97.5	<b>100</b> /99.1	<b>100</b> /99.4	<b>100</b> / <b>99.8</b>
Tile	98.7/95.6	99.3/98.9	99.8/97.0	99.6/99.2	95.5/-	-/98.9	<b>100</b> / <b>99.4</b>	<b>100</b> /98.2	<b>100</b> /98.7
Wood	99.2/95.0	99.2/ <b>99.3</b>	<b>100</b> /94.5	99.1/96.4	95.7/-	-/99.1	<b>100</b> /96.7	<b>100</b> /95.0	99.4/98.4
Bottle	<b>100</b> /98.6	<b>100</b> /99.0	<b>100</b> /98.0	99.2/ <b>99.1</b>	99.6/-	-/97.3	<b>100</b> /98.8	<b>100</b> /98.7	<b>100</b> /98.9
Cable	99.5/98.4	95.0/99.4	<b>99.9</b> /97.6	91.8/94.7	99.1/-	-/99.5	94.6/96.8	99.4/98.1	<b>99.9</b> /98.1
Capsule	98.1/98.8	96.3/97.3	97.7/ <b>98.9</b>	98.5/94.3	96.2/-	-/96.8	97.5/98.2	99.4/95.7	<b>99.5</b> / <u>98.5</u>
Hazelnut	<b>100</b> /98.7	99.9/98.2	<b>100</b> /97.9	<b>100</b> / <b>99.7</b>	98.5/-	-/92.5	<b>100</b> /99.4	<b>100</b> /98.4	<b>100</b> /99.5
Metal nut	<b>100</b> /98.4	<b>100</b> / <b>99.6</b>	<b>100</b> /98.8	98.7/ <u>99.5</u>	99.5/-	-/99.0	<b>100</b> /99.4	<b>100</b> /99.0	<b>100</b> /98.8
Pill	96.6/97.4	<b>99.6</b> /95.7	99.0/ <b>98.6</b>	98.9/97.6	98.3/-	-/92.1	97.7/97.7	<b>100</b> /99.1	98.1/ <u>97.9</u>
Screw	98.1/ <b>99.4</b>	97.0/99.1	98.2/99.3	93.9/97.6	<b>100</b> /-	-/98.6	97.2/99.0	99.0/ <u>99.3</u>	96.9/99.1
Toothbrush	<b>100</b> /98.7	99.5/93.0	99.7/98.5	<b>100</b> /98.1	98.7/-	-/93.1	<b>100</b> /99.2	<b>100</b> /98.7	<b>100</b> / <b>99.4</b>
Transistor	<b>100</b> /96.3	96.7/95.4	<b>100</b> / <b>97.6</b>	93.1/90.9	<u>98.3</u> /-	-/94.5	96.1/93.7	<b>100</b> /95.3	98.3/ <u>96.2</u>
Zipper	99.4/98.8	98.5/98.2	99.9/98.9	<b>100</b> /98.8	99.0/-	-/97.6	<b>100</b> / <b>99.0</b>	<b>100</b> /98.2	98.5/97.9
Average	99.1/98.1	98.5/97.8	99.6/98.1	98.0/97.3	98.3/-	-/96.7	98.7/98.3	<b>99.8</b> /98.1	99.3/ <b>98.6</b>

**Table 3:** Comparison with SOTA methods on VisA dataset. I-AUROC and P-AUROC are displayed in each entry. The best results among all methods are shown in bold, and the underlined results denote the best results among reconstruction-based methods.

Category	Embedding-based methods			Reconstruction-based methods			
	PatchCore [22]	RD4AD [8]	SimpleNet [17]	DRAEM [32]	OCR-GAN [14]	DDAD [19]	Ours
Candle	98.7/ <b>99.2</b>	96.2/98.9	96.9/98.6	89.6/91.0	98.9/-	<b>99.9</b> /98.7	<b>99.9</b> /94.8
Capsules	68.8/96.5	91.8/99.4	89.5/99.2	89.2/99.0	98.8/-	<b>100</b> /99.5	99.1/ <b>99.6</b>
Cashew	97.7/99.2	<b>98.7</b> /94.4	94.8/ <b>99.0</b>	88.3/85.0	97.4/-	94.5/ <u>97.4</u>	98.4/97.0
Chewinggum	99.1/98.9	99.3/97.6	<b>100</b> /98.5	96.4/97.7	99.4/-	98.1/96.5	99.6/ <b>99.1</b>
Fryum	91.6/95.9	96.9/96.4	96.6/94.5	94.7/82.5	96.3/-	99.0/ <b>96.9</b>	<b>99.4</b> / <b>96.9</b>
Macaroni1	90.1/98.5	98.7/99.3	97.6/99.6	93.9/99.4	97.2/-	99.2/98.7	<b>99.9</b> / <b>99.8</b>
Macaroni2	63.4/93.5	91.4/99.1	83.4/96.4	88.3/99.7	95.1/-	<b>99.2</b> /98.2	98.9/ <b>99.8</b>
Pcb1	96.0/ <b>99.8</b>	96.7/99.6	99.2/ <b>99.8</b>	84.7/98.4	96.1/-	<b>100</b> /93.4	99.6/ <u>99.6</u>
Pcb2	95.1/98.4	97.2/98.3	99.2/ <b>98.8</b>	96.2/94.0	98.3/-	99.7/97.4	<b>100</b> /98.6
Pcb3	93.0/98.9	96.5/ <b>99.3</b>	98.6/99.2	97.4/94.3	98.1/-	97.2/96.3	<b>99.9</b> /98.9
Pcb4	99.5/98.3	99.4/98.2	98.9/98.6	98.9/97.6	99.7/-	<b>100</b> /98.5	99.9/ <b>99.5</b>
Pipe fryum	99.0/99.3	99.6/99.1	99.2/99.3	94.7/65.8	99.7/-	<b>100</b> / <b>99.5</b>	98.9/99.4
Average	91.0/98.1	96.9/98.3	96.2/98.5	92.4/92.0	97.9/-	98.9/97.6	<b>99.5</b> / <b>98.6</b>

**Table 2:** Comparison with SOTA on MPDD dataset. I-AUROC and P-AUROC are displayed in each entry. The best results among all methods are shown in bold, and the underlined results denote the best results among reconstruction-based methods.

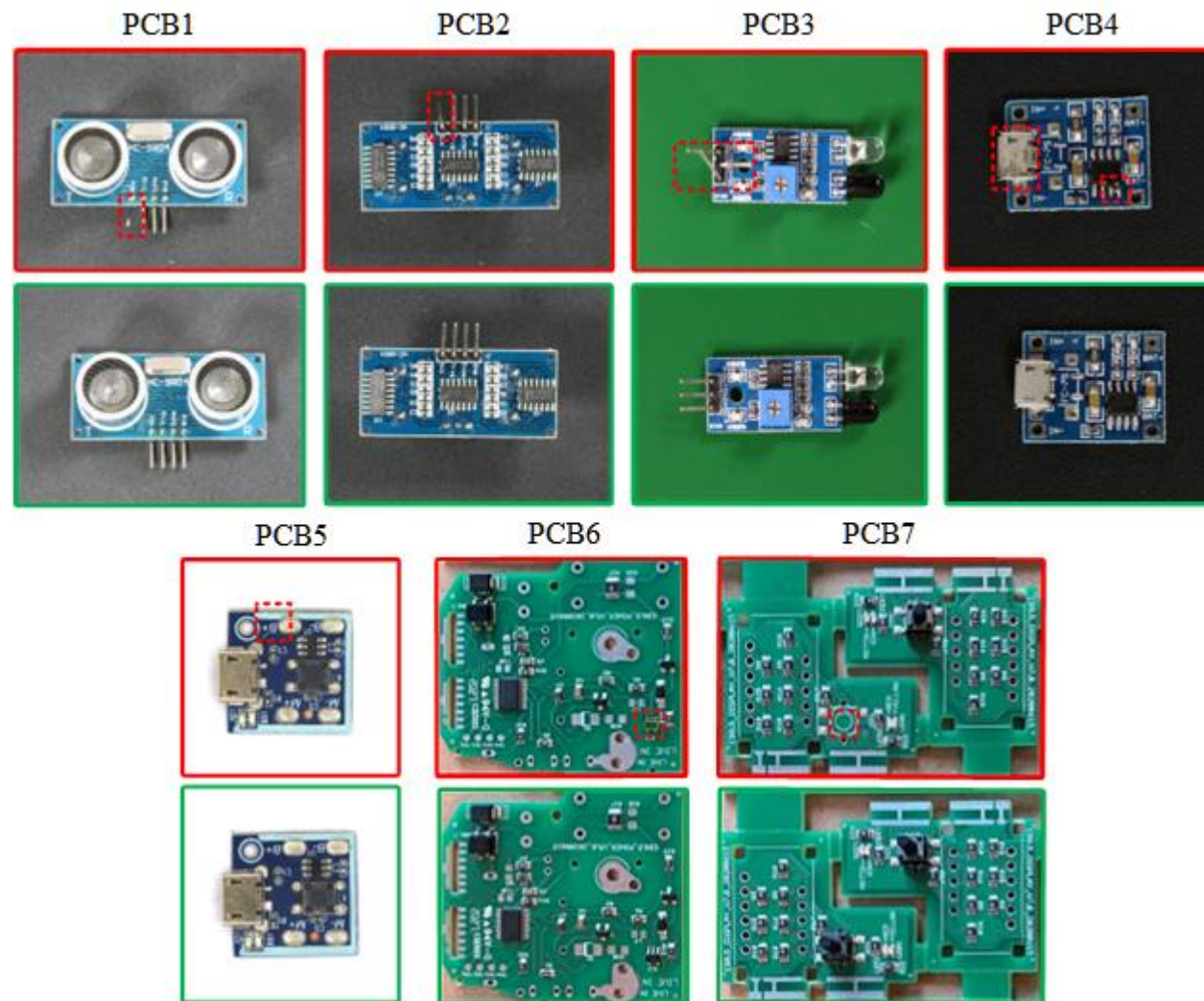
Category	Embedding-based methods			Reconstruction-based methods			
	PatchCore [22]	RD4AD [8]	SimpleNet [17]	DRAEM [32]	OCR-GAN [14]	DDAD [19]	Ours
Bracket Black	85.3/97.6	90.2/98.0	85.1/96.0	91.8/98.2	99.9/-	<b>98.7</b> /96.7	98.0/ <b>99.4</b>
Bracket Brown	92.5/ <b>98.1</b>	94.2/97.0	<b>98.3</b> /94.4	90.3/63.7	89.4/-	92.7/97.2	90.7/97.5
Bracket White	92.3/99.7	90.1/99.3	98.0/96.7	88.8/98.9	88.1/-	96.6/91.8	<b>98.3</b> / <b>99.7</b>
Connector	100/99.4	99.5/99.4	<b>100</b> / <b>99.5</b>	<b>100</b> /91.2	100/-	96.2/98.6	<b>100</b> /98.2
Metal Plate	<b>100</b> /98.8	<b>100</b> /99.1	<b>100</b> /98.5	<b>100</b> /96.6	100/-	100/98.1	99.9/ <b>99.4</b>
Tubes	77.4/97.2	97.6/99.1	97.9/ <b>99.2</b>	94.7/95.9	98.1/-	99.2/99.0	<b>98.1</b> / <u>97.8</u>
Average	91.3/98.5	95.3/ <b>98.7</b>	96.6/97.4	94.3/90.7	95.9/-	97.2/96.9	<b>97.5</b> / <b>98.7</b>

**Table 4:** Comparison with SOTA methods on PCB-Bank dataset. I-AUROC and P-AUROC are displayed in each entry. The best results among all methods are shown in bold, and the underlined results denote the best results among reconstruction-based methods.

Category	Embedding-based methods			Reconstruction-based methods			
	PatchCore [22]	RD4AD [8]	SimpleNet [17]	DRAEM [32]	OCR-GAN [14]	DDAD [19]	Ours
Pcb1	96.0/ <b>99.8</b>	96.7/99.6	99.2/ <b>99.8</b>	84.7/98.4	96.1/-	<b>100</b> /93.4	99.6/ <u>99.6</u>
Pcb2	95.1/98.4	97.2/98.3	99.2/ <b>98.8</b>	96.2/94.0	98.3/-	99.7/97.4	<b>100</b> / <u>98.6</u>
Pcb3	93.0/98.9	96.5/ <b>99.3</b>	98.6/99.2	97.4/94.3	98.1/-	97.2/96.3	<b>99.9</b> /98.9
Pcb4	99.5/98.3	99.4/98.2	98.9/98.6	98.9/97.6	99.7/-	<b>100</b> /98.5	99.9/ <b>99.5</b>
Pcb5	94.6/ <b>99.8</b>	94.1/99.5	94.5/99.4	97.2/97.8	85.9/-	<b>99.7</b> /96.0	99.6/ <u>99.1</u>
Pcb6	82.2/98.9	89.4/98.9	91.7/97.5	72.4/94.6	75.1/-	87.8/98.5	<b>92.2</b> / <b>99.7</b>
Pcb7	93.7/98.8	99.0/99.6	<b>100</b> / <b>99.9</b>	97.7/98.3	85.7/-	94.4/98.7	99.6/99.8
Average	94.2/99.1	96.0/99.1	96.2/98.5	91.5/96.4	91.3/-	97.4/96.5	<b>98.7</b> / <b>99.3</b>



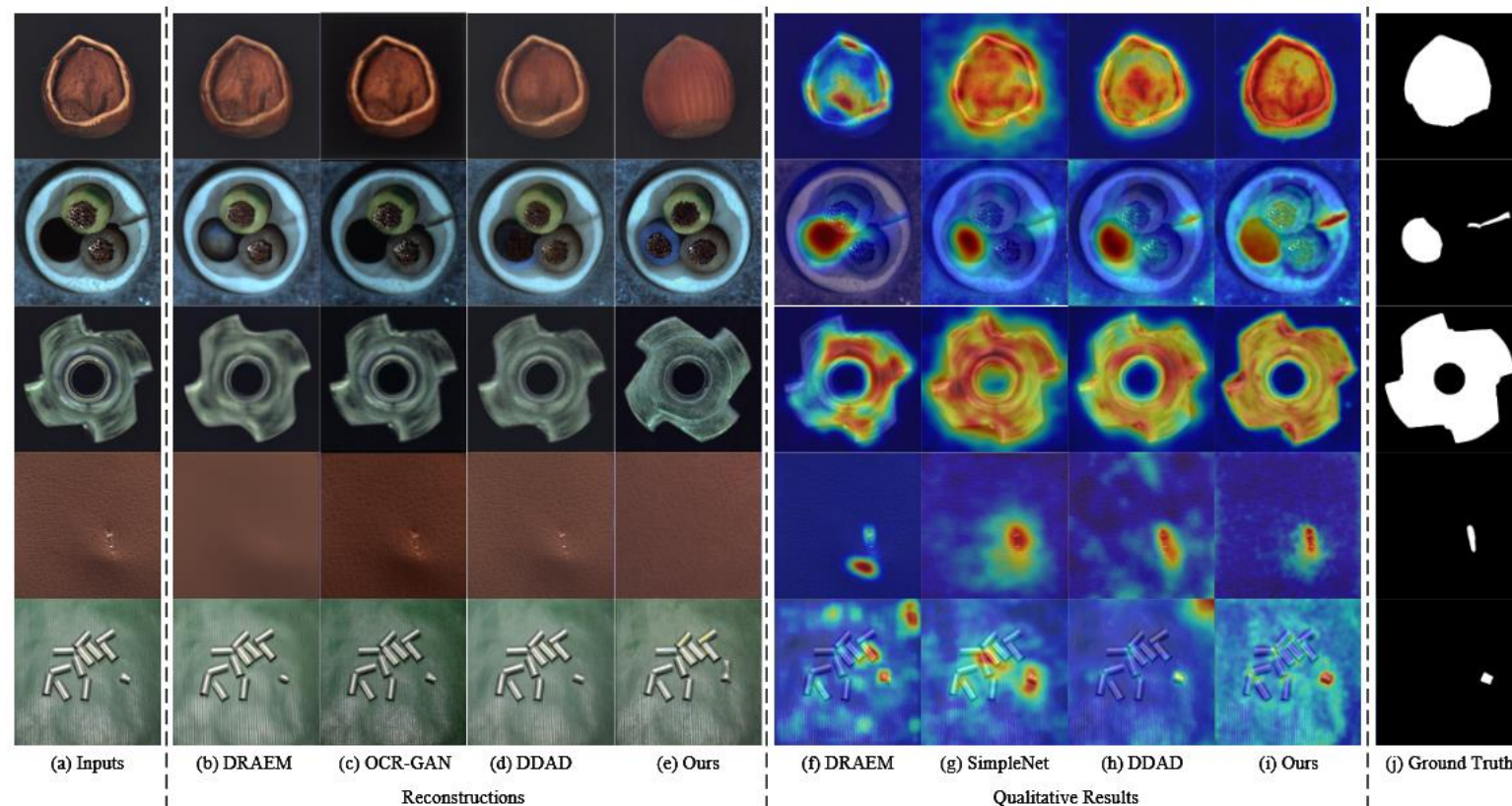
# Experiment





# Experiment

other methods usually fail to reconstruct **large-scale** anomalies into normal regions



**Fig. 3:** Reconstructions and qualitative comparisons with other methods. The first four rows display examples of the MVTec-AD dataset, and the last row is for the MPDD dataset. OCR-GAN only produces anomaly scores, and there is no anomaly map. SimpleNet is the embedding-based method.





# Experiment / Ablation Study

Adaptive Denoising Step (ADS)  
Anomaly-oriented Training Paradigm (ATP)  
Spatial-Adaptive Feature Fusion (SAFF)

Compare to fix steps  
ADS: ensures anomaly-free  
reconstruction and preserves as much  
information as possible about normal  
regions.

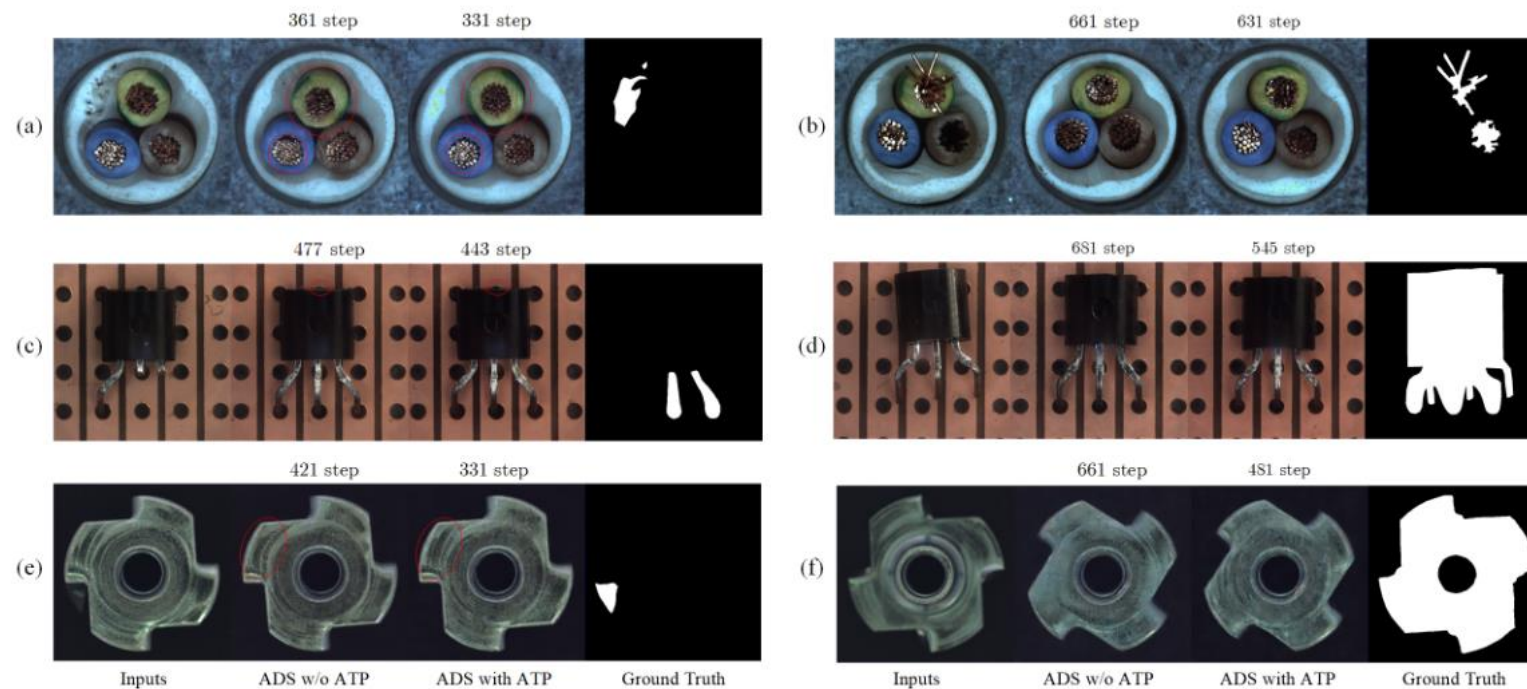
**Table 6:** Performance of each component on MVTec-AD dataset. The best results are shown in bold.

Method	I-AUROC	P-AUROC
Baseline	98.3	98.0
Baseline + ADS	99.0	98.5
Baseline + ATP	98.7	98.5
Baseline + ADS + ATP w/o SAFF	99.2	98.3
Baseline + ADS + ATP with SAFF (Ours)	<b>99.3</b>	<b>98.6</b>

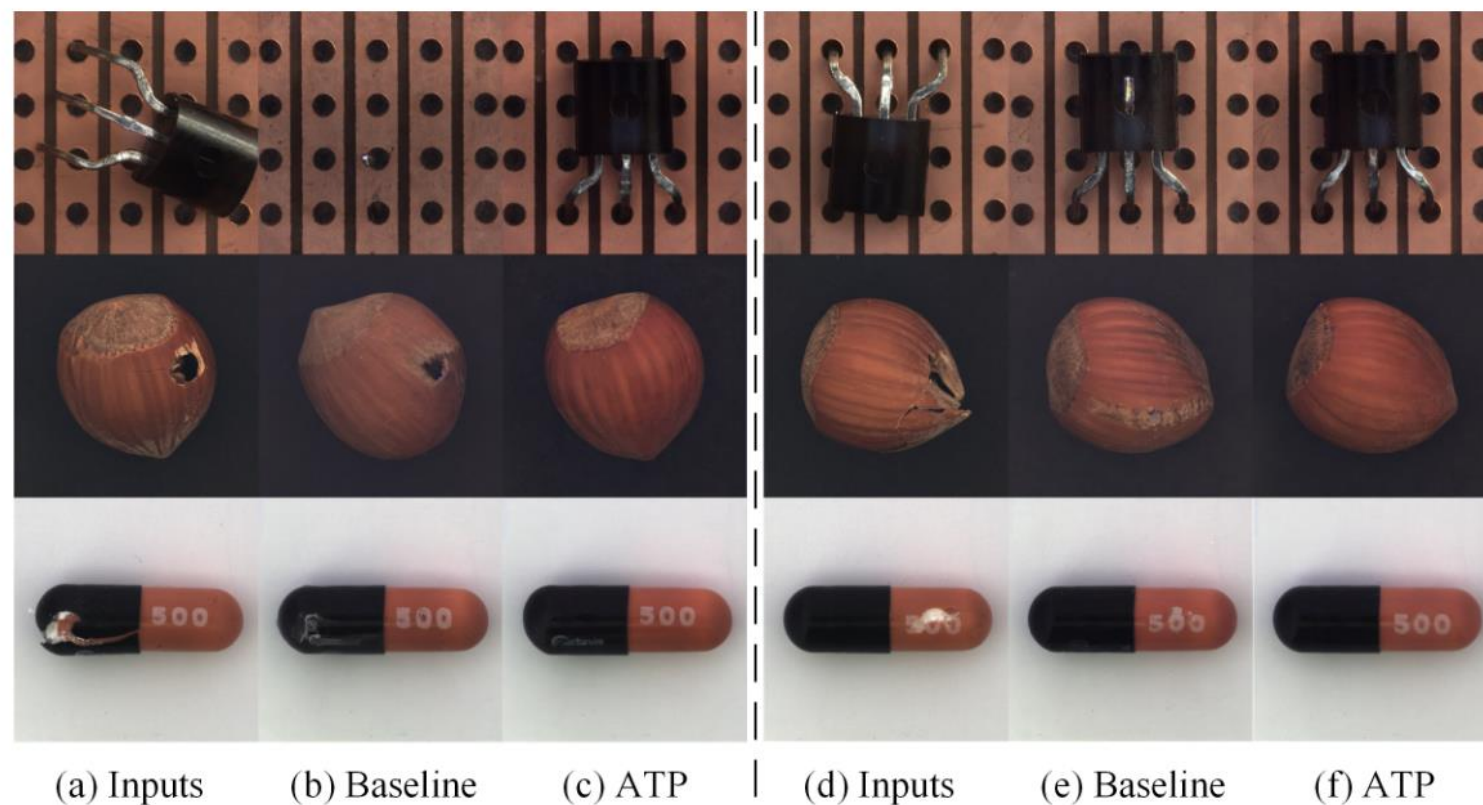
**Table 7:** Comparison of adaptive steps and different fixed steps on MVTec-AD dataset. The best results are shown in bold.

Denoising steps	fixed steps			Adaptive steps
	350 step	550 step	750 step	
I-AUROC	98.8	98.8	98.7	<b>99.3</b>
P-AUROC	97.6	98.1	98.5	<b>98.6</b>

# Experiment / Ablation Study



**Fig. 4:** Reconstructions of different types of anomaly and proper steps. Examples (a), (c), and (e) contain small-scale anomalies, and (b), (d), and (f) are large-scale anomalies. The numbers above the reconstructed images represent the proper steps. Differences in details of normal areas are marked in red circles.



**Fig. 5:** Qualitative comparisons between baseline and proposed ATP on MVTec-AD. The same denoising steps are used for the two methods.



南京航空航天大学  
NANJING UNIVERSITY OF AERONAUTICS AND ASTRONAUTICS

**Thank you**