

# SALAD 🥗 – Semantics-Aware Logical Anomaly Detection

Matic Fučka

Vitjan Zavrtanik

Danijel Skočaj

University of Ljubljana, Faculty of Computer and Information Science, Slovenia

{matic.fucka, vitjan.zavrtanik, danijel.skocaj}@fri.uni-lj.si

## Abstract

Recent surface anomaly detection methods excel at identifying structural anomalies, such as dents and scratches, but struggle with logical anomalies, such as irregular or missing object components. The best-performing logical anomaly detection approaches rely on aggregated pretrained features or handcrafted descriptors (most often derived from composition maps), which discard spatial and semantic information, leading to suboptimal performance. We propose SALAD, a semantics-aware discriminative logical anomaly detection method that incorporates a newly proposed composition branch to explicitly model the distribution of object composition maps, consequently learning important semantic relationships. Additionally, we introduce a novel procedure for extracting composition maps that requires no hand-made labels or category-specific information, in contrast to previous methods. By effectively modelling the composition map distribution, SALAD significantly improves upon state-of-the-art methods on the standard benchmark for logical anomaly detection, MVTec LOCO, achieving an impressive image-level AUROC of 96.1%. Code: <https://github.com/MaticFuc/SALAD>

## 1. Introduction

Surface anomaly detection aims to detect and localise abnormal regions in the image while training only on anomaly-free images. It is commonly used in the industrial inspection domain [2–5, 46] where the limited availability and considerable diversity of abnormal images make training supervised models impractical. The problem of surface anomaly detection can be split into structural and logical anomaly detection. Structural anomalies are irregularities in the local appearance distribution, e.g., dents or scratches. They can be detected by modelling the object’s anomaly-free appearance and detecting local texture deviation. Logical anomalies break the semantic constraints of the image, e.g. incorrect number or misplacement of object components. For such anomalies, a model of local appearance does not suffice since object components may fit the anomaly-free appearance distribution

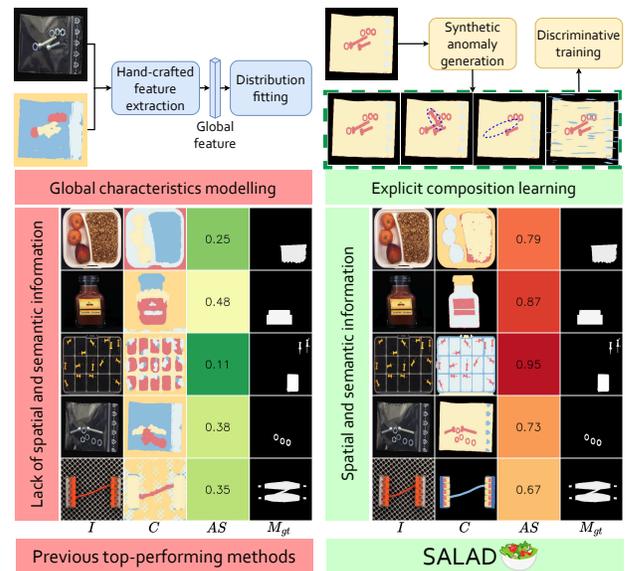


Figure 1. While previous approaches rely on handcrafted features, SALAD explicitly models the composition distribution by training directly on the composition maps  $C$ . SALAD is trained on simulated anomalous examples that are enclosed by a green dashed rectangle, whereas individual synthetic logical anomalies are highlighted with blue ellipses. SALAD improves the estimated anomaly scores (AS) on near-in-distribution logical anomalies.

locally. Recent top-performing surface anomaly detection methods [9, 12, 13, 23, 25, 33, 42, 44] are focused mostly on detecting structural anomalies and fail to model the semantic information required for logical anomaly detection.

Recent logical anomaly detection methods consist of a base structural branch that is typically a time-tested surface anomaly detection method, and a logical anomaly branch focusing on detecting semantic anomalies. The output scores of both branches are fused to produce a final anomaly score. Best performing logical anomaly methods attempt to detect deviations in appearance caused by logical anomalies by either modelling the global image appearance [36] or by utilising object composition maps [16, 19, 24]. Global appearance-based approaches build a global image descriptor

by aggregating pretrained features, disregarding the information contained in the object components’ position, orientation, and frequency. This results in missed detections in harder, near-in-distribution anomalous examples (Figure 1, left). Additionally, such approaches do not decouple logical and structural anomalies because they use pretrained features that focus on modelling appearance, leading to a considerable emphasis on structural anomalies. Recent composition map-based logical anomaly detection methods extract composition maps from the input RGB images. They then rely on handcrafted features extracted from composition maps, such as the class frequency, to better model the global distribution. Similar to the global appearance approaches, such representations do not sufficiently model the available semantic and spatial information. Additionally, the best-performing composition-based methods require either hand-labelled training examples [19] or category-specific procedures [16] to extract composition maps, making them infeasible to apply to new datasets.

We hypothesise that training a model (in our case, a composition branch) to model the composition map distribution would also capture the critical spatial and semantic information unobtainable with global representations, which would improve logical anomaly detection performance. Discriminative methods, which rely on synthetic anomalies to learn an anomaly-free distribution, present a possible solution. However, synthetic anomalies are currently defined only for RGB images [38, 42]. Therefore, we adapted them for composition maps. Additionally, as composition maps contain a compressed representation of object class, shape, and position while discarding appearance information, it is simpler to simulate more expressive anomalies appropriate for detecting logical anomalies. We propose a composition branch defined as a discriminative anomaly detection model operating with composition maps (Figure 1, right).

The contributions of this work are twofold. (i) As our main contribution, we propose SALAD, a **S**emantics-**A**ware discriminative **L**ogical **A**nomaly **D**etection method that extends the recent appearance and global branch framework with a new compositional branch that explicitly learns the anomaly-free object composition distribution. (ii) As an additional contribution, we propose a novel object composition map generation process. The proposed approach produces maps of high quality (Figure 1, right, columns *C*) without requiring hand-labelled training data or category-specific procedures in contrast to previous approaches. We showcase its robustness by applying it to several objects and datasets.

To emphasise the value of the proposed contributions, extensive experiments are performed on the standard logical anomaly detection benchmark, MVTec LOCO [3]. SALAD achieves a new state-of-the-art result on MVTec LOCO (AUROC of 96.1%), outperforming competing methods by a significant margin of 3.0 percentage points. Addition-

ally, SALAD is evaluated on standard MVTec AD [2] and VisA [46] datasets, achieving excellent results (AUROC of 98.9% and 97.9%).

## 2. Related work

**Surface anomaly detection** has been extensively researched, with methods categorised into three main paradigms: reconstructive, embedding-based, and discriminative.

*Reconstructive methods* train either an autoencoder-like [23, 43] network, a generative adversarial network [21], a diffusion model [37] or a transformer [30, 40] on anomaly-free images and assume the resulting model will not generalise well to anomalous regions, making them distinguishable by reconstruction error. *Embedding-based methods* leverage pretrained models to extract features and fit a normality model on top of them. The normality methods are often a coreset [33], a student-teacher network [1, 9, 34] or a normalising flow network [13, 41]. *Discriminative methods* focus on learning the boundary between normal and abnormal samples. Methods in this paradigm are learned to segment synthetic anomalies [12, 25, 32, 38, 42, 44] and learn a normality model to generalise to real-world scenarios. These methods fail on logical anomalies as they focus on local characteristics and do not consider global semantics.

**Logical anomaly detection** is a new surface anomaly detection research direction. The methods can be divided into three paradigms: local-global reconstruction, global distribution approaches and composition-based.

*Local-global reconstruction methods* use a two-branch neural network [1, 3, 8, 14, 39, 45]. These approaches contain a local and global appearance branch and assume that structural anomalies occur as local deviations and logical anomalies impact a large part of the image, which does not always hold. EfficientAD [1] uses a convolutional network with a small receptive field as the local branch and an autoencoder as the global branch. These methods have difficulties with categories without a constant object layout.

*Global distribution approaches* extract global appearance descriptors from images and use a descriptor distribution model [7, 28] to detect anomalies. LogicAD [17] uses a large Vision Language Model to extract a global description. This is converted into a formal representation and evaluated by an automatic theorem prover. PUAD [36] estimates the global distribution by fitting a Gaussian with mean feature vectors extracted from feature maps produced by EfficientAD’s student. During inference, the anomaly is detected using Mahalanobis distance. Simply using the mean of extracted pretrained features as a global appearance representation disregards considerable spatial information, leading to poor performance on spatially dependent logical anomalies. Introducing spatial information might improve the logical anomaly detection performance.

*Composition-based approaches* model the semantic infor-

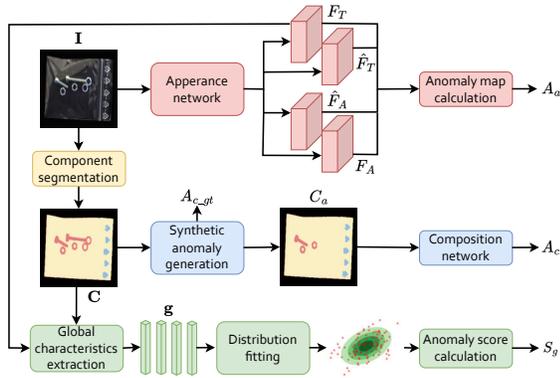


Figure 2. SALAD is constructed from a **local appearance branch**, a **composition branch**, and a **global branch**. Each branch focuses on a different level of image semantics. Synthetic anomalies are generated to train the composition branch. Composition maps are segmented using a **component segmentation network**.

mation by using composition maps, i.e. object component segmentation maps of the image. This paradigm was introduced with ComAD [24], which extracts features from a pretrained network and clusters them to create a rough semantic segmentation. The obtained segmentation maps are used only to extract handcrafted features and store them in a memory bank. PSAD [19] uses hand-labelled segmentation maps to finetune a pretrained feature extractor with an attached segmentation head. After the fine-tuning, PSAD creates a memory bank with global statistics and extracted feature vectors. PSAD’s manual annotation requirement is impractical in real-world scenarios. CSAD [16] obtains patch histograms of composition maps and stores them in a memory bank. CSAD extracts object composition maps automatically, but the parameters for each object category must be tuned, which is a drawback in practical use.

Unlike recent logical anomaly detection methods, SALAD does not focus on feature averaging or handcrafted features but explicitly models the composition map distribution, thus learning important composition information. We also propose a highly accurate composition map generation procedure that does not require hand-labelled data or category-specific information.

### 3. SALAD

Recent logical anomaly detection methods are composed of a base local appearance model and a global model. The global model typically constructs the global distribution by either using aggregated pretrained features or handcrafted descriptors. The best-performing methods use composition maps to create better global representations. Due to such a construction, a significant amount of semantic information is not captured. In the proposed method, SALAD (Figure 2),

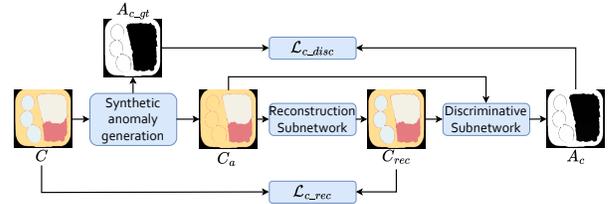


Figure 3. Composition branch architecture.

we follow the initial framework of a base structural anomaly detection model with a global appearance branch but propose a novel composition branch that explicitly learns the distribution of composition maps, consequently learning critical spatial and semantic information. Additionally, SALAD contains an improved global appearance model.

In the composition branch, the anomaly-free object composition distribution is learned in a discriminative fashion through a novel anomaly simulation process. It uses the information-dense composition maps to simulate near-in-distribution anomalies that are difficult to simulate using traditional simulation processes [42]. Such simulated anomalies facilitate the learning of a tight decision boundary around the anomaly-free semantic structure of the images, leading to an improvement in logical anomaly detection performance.

At inference time, structural anomalies are detected by the local appearance branch, while the global appearance and the composition branch focus on logical anomalies. Individual branch outputs are then combined using a score fusion model. In the rest of this section, we describe SALAD in detail.

#### 3.1. Local appearance model

Following recent advancements in logical anomaly detection, a powerful surface anomaly detection model is used as the base structural anomaly detection branch. In SALAD, the initial structural branch follows the EfficientAD architecture [1], a top-performing surface anomaly detection framework. EfficientAD takes an RGB image  $I$  as input and works as an embedding reconstruction model. The input image  $I$  is mapped to a feature representation  $F_T$  by a teacher encoder trained on natural images [10]. An autoencoder network takes  $I$  as input and is tasked with reconstructing the teacher output  $F_T$ , outputting  $F_A$ . The student encoder is tasked with reconstructing both  $F_T$  and  $F_A$ , outputting  $\hat{F}_T$  and  $\hat{F}_A$ . The anomaly map  $A_a$  output by EfficientAD is based on the difference between  $F_T$  and  $\hat{F}_T$  and between  $F_A$  and  $\hat{F}_A$ .

#### 3.2. Object composition model

The composition branch is formulated as a discriminative anomaly detection method that operates on object composition maps  $C$ . Training such a model requires both an automated way of obtaining object composition maps (Section 3.3) as well as a well-defined anomaly simulation process

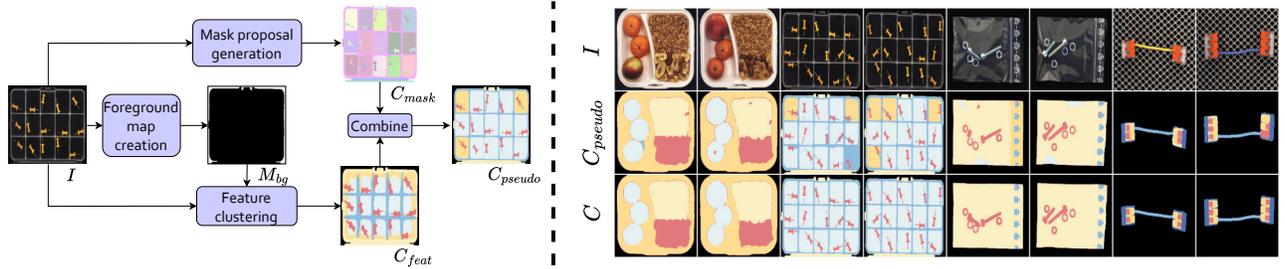


Figure 4. Left: Pseudo label  $C_{pseudo}$  generation procedure. The main idea is to combine precise mask generation from SAM-HQ [18] and the discriminative power from DINO [6]. Right: Comparisons of the final object composition maps  $C$  (generated by the *composition segmentation model*) and the pseudo labels  $C_{pseudo}$ . Input images  $I$  are also depicted for comparison. Final object composition maps  $C$  contain less misclassified components than the pseudo labels  $C_{pseudo}$ .

specifically designed for composition maps (Section 3.4).

The *composition branch* follows a general discriminative anomaly detection architecture [42] composed of a composition reconstruction network and a composition discriminative network (Figure 3). The composition reconstruction network first restores the (synthetically) anomalous parts of the composition map to their anomaly-free appearance. Then, the input composition map and its anomaly-free reconstruction are passed to the composition discriminative network to output an anomaly mask. The composition branch operates solely in the space of object composition maps  $C$ . During training, the composition reconstruction network takes as input a composition map  $C_a$ , which has been augmented to include simulated anomalies. The network is then trained to restore the original anomaly-free composition  $C$  by outputting an anomaly-free composition reconstruction  $C_{rec}$ . During inference,  $C$  is used as the input instead of  $C_a$ .

Since the composition map  $C$  values belong to individual classes, segmentation losses are used to train the composition reconstruction network. Namely, the focal loss  $\mathcal{L}_{foc}$  [22] and dice loss  $\mathcal{L}_{dice}$  [35] are used:

$$\mathcal{L}_{c_{rec}} = \mathcal{L}_{foc}(C, C_{rec}) + \mathcal{L}_{dice}(C, C_{rec}) , \quad (1)$$

where  $C$  is the original object composition map,  $C_{rec}$  is the reconstructed object composition map. The composition reconstruction network generalises to real anomalous examples, successfully reconstructing them to be anomaly-free.

After obtaining the anomaly-free composition reconstruction  $C_{rec}$ ,  $C_{rec}$  and the augmented composition map  $C_a$  (during inference,  $C_{rec}$  and  $C$  are used) are concatenated and used as input for the composition discriminative network. The discriminative network is trained to predict the difference between  $C_{rec}$  and  $C_a$  to output an anomaly segmentation map  $A_c$ .

Following recent literature [38], the loss for the composition discriminative network is defined as:

$$\mathcal{L}_{c_{disc}} = \alpha \mathcal{L}_{foc}(A_{c_{gt}}, A_c) + \mathcal{L}_1(A_{c_{gt}}, A_c) , \quad (2)$$

where  $A_{c_{gt}}$  is the ground truth anomaly map corresponding to synthetic anomalies,  $A_c$  is the predicted anomaly map, and  $\alpha$  is the weighting parameter (set to 5 in all our experiments).

### 3.3. Object composition maps

A semantically meaningful representation of the composition must first be extracted to model the object composition distribution accurately. Segmentation of object parts concisely represents part frequency, shapes, sizes, and positions without additional appearance information that is redundant for object composition and increases the complexity of the representation. Accurate component-level segmentation maps (dubbed composition maps) are used in SALAD.

A two-step process is used for composition map extraction. First, *pseudo-labels*  $C_{pseudo}$  for the training set are created by clustering DINO [6] features to obtain rough segmentation maps  $C_{feat}$ , which are then used to classify highly accurate mask proposals  $C_{mask}$  generated by SAM-HQ [18]. Finally, the pseudo-labels are used to train a *component segmentation model* to predict the final object composition maps  $C$ . Pseudo-label creation is illustrated in Figure 4.

For the initial pseudo-labels, background maps for each training image are generated by querying SAM-HQ on image corners and combining the resulting masks to produce a background mask. This mask is inverted to create the foreground mask  $M_{fg}$ . DINO [6] feature maps are extracted and resized to  $256 \times 256$ . Features outside  $M_{fg}$ , that is, the background features, are set to 0 to reduce noise; the rest are then subsampled and clustered into  $K$  clusters (in our case  $K=6$ ) to produce a rough object composition map  $C_{feat}$ .

SAM-HQ is queried on a grid over the input image  $I$  to obtain mask proposals  $C_{mask}$ . Each mask proposal is classified as the class of the corresponding majority cluster in  $C_{feat}$ , aligning high-quality masks with component labels to create high-quality pseudo-labels  $C_{pseudo}$ . Due to the computational intensity of SAM-HQ and DINO, a *component segmentation model* is trained with  $I$  and corresponding  $C_{pseudo}$  pairs. Specifically, we use a simple UNet trained

with a cross-entropy loss. Even though some  $C_{pseudo}$  contain mistakes, the components are correctly classified on average across the dataset. Due to that, the component segmentation model generalises and outputs composition maps  $C$  without incorrectly labelled components. The trained component segmentation model infers the desired composition map  $C$  directly from  $I$ , enabling efficient composition map extraction. Figure 4 shows examples of  $C_{pseudo}$  and object composition maps  $C$  produced by the component segmentation model.

### 3.4. Synthetic anomaly generation

An appropriate synthetic anomaly generation procedure is required to facilitate the training of the discriminative composition branch. Due to differences between structural and logical anomalies, different anomaly generation strategies are required. To simulate structural anomalies, we extend the synthetic anomaly generation proposed by DRÆM [42] by pasting a random class on top of the composition map according to an anomaly map generated using Perlin noise [29]. To generate near-in-distribution logical anomalies, a random component is either inpainted (from another image) or erased. When a component is erased, the corresponding region is inpainted with a component class randomly selected from the neighbouring components. In this case, the anomaly map marks both the erased component and the neighbouring component class used for inpainting. For cases where a component is inpainted, identifying the exact anomaly location is ill-posed (e.g., an extra screw added to a screw bag). Therefore, the anomaly map in these cases includes all regions containing the inpainted component class. Several examples of synthetic anomalies can be seen in Figure 1.

### 3.5. Global appearance model

Using a strong global appearance model can also improve the detection of structural anomalies. The global appearance branch utilises features extracted from the input image  $I$  and its corresponding object composition map  $C$ . In  $C$ , pixels belonging to individual image components are marked with their corresponding class labels  $c$ .

For each class label  $c$  in  $C$ , the mean feature vector  $g_c$  is computed from the feature vectors in  $F_T$  corresponding to pixels belonging to  $c$  in  $C$ . The set of  $g_c$  values for all classes  $C$  represents the global appearance descriptor  $g$ .

The procedure is repeated for each sample  $i$  in the training set to obtain global appearance descriptors  $g^{(i)}$  upon which the global distribution is estimated by fitting a Gaussian distribution [31]. The mean  $\mu_c$  and covariance  $\Sigma_c$  for each class in  $C$  are calculated from all samples  $g_c^{(i)}$  in the training set. During inference, the anomaly score  $S_g$  is calculated using the average Mahalanobis distance [31] for each class,

that is:

$$S_g = \frac{1}{K} \sum_{c=1}^K \sqrt{(g_c - \mu_c)^\top \Sigma_c^{-1} (g_c - \mu_c)}, \quad (3)$$

where  $K$  is the total number of classes in  $C$ .

### 3.6. Anomaly score calculation

Each model branch outputs an anomaly score at inference:  $AS_a$ ,  $AS_c$  and  $AS_g$  for the appearance, compositional and global branches, respectively. Individual scores are calculated as follows:

$$AS_a = \max(A_a), \quad AS_c = \max(A_c), \quad AS_g = S_g, \quad (4)$$

where  $A_a$  is the output of the appearance branch,  $A_c$  is the output of the composition branch, and  $S_g$  is the output of the global branch. The outputs are normalised using the means  $\mu_a$ ,  $\mu_c$  and  $\mu_g$  and standard deviations  $\sigma_a$ ,  $\sigma_c$  and  $\sigma_g$  of the anomaly scores on the validation set. The final anomaly score is then defined as:

$$AS = \frac{AS_a - \mu_a}{\sigma_a} + \frac{AS_c - \mu_c}{\sigma_c} + \frac{AS_g - \mu_g}{\sigma_g}. \quad (5)$$

## 4. Experiments

### 4.1. Datasets

Experiments are performed on the standard anomaly detection dataset for logical anomalies: the MVTec LOCO [3] and two standard anomaly detection datasets for structural anomalies: MVTec AD [2] and VisA [46]. The MVTec LOCO dataset comprises 3,644 images distributed across five object categories, the MVTec AD dataset comprises 5,354 images distributed across ten object categories and five texture categories, and the VisA dataset comprises 10,821 images distributed across twelve categories. All three datasets provide pixel-level annotations for the test images, enabling accurate evaluation and analysis.

### 4.2. Implementation Details

For the composition map generation, SAM-HQ-h [18] and a DINO [6] pretrained ViT-b\8 [11] are used. A UNet is used as the component segmentation model. The UNet was trained for 15 epochs with AdamW [26] using cross-entropy loss with a learning rate of  $5 \cdot 10^{-4}$  and a batch size of 8.

SALAD follows the training regime from EfficientAD-70000 iterations with the Adam [20] optimizer. The learning rate was set to  $10^{-4}$  for the appearance branch and  $10^{-5}$  for the composition branch. Both learning rates were multiplied by 0.1 after 90% (66500) of the iterations. Synthetic anomalies were added to the object composition maps with a 50% probability. All of the images were resized to  $256 \times 256$  pixels. Following the standard protocol, a separate model was

Method	Venue	Supervised Masks	Breakfast box	Juice bottle	Pushpins	Screw bag	Splicing conn.	Average
DRÆM [42]	ICCV'21		80.2	94.3	68.6	70.6	85.4	79.8
TransFusion [12]	ECCV'24		82.4	<b>99.7</b>	63.8	71.5	83.7	80.2
DSR [44]	ECCV'22		85.8	99.2	76.5	64.9	85.5	82.6
THFR [14]	ICCV'23		77.3	80.1	80.8	79.7	81.0	83.3
LogicAD [17]	AAAI'25		<b>92.1</b>	81.6	<b>98.1</b>	83.8	73.4	86.0
Sinbad [7]	ArXiv'23		<b>91.8</b>	94.4	83.9	<b>86.8</b>	84.5	88.3
ComAD + Patchcore [24]	AEI'23		86.4	96.6	93.4	80.2	94.1	90.1
SLSG [39]	PR'24		88.9	99.1	95.5	79.4	88.5	90.3
SAM-LAD [28]	KBS'25		<b>91.0</b>	97.6	88.2	<b>86.6</b>	90.0	<b>90.7</b>
EfficientAD [1]	WACV'24		88.5	99.0	93.6	73.6	<b>97.1</b>	<b>90.7</b>
PUAD [36]	ICIP'24		87.1	<b>99.7</b>	<b>98.0</b>	81.1	<b>96.8</b>	<b>93.1</b>
<b>SALAD</b>			89.3	<b>99.7</b>	<b>99.4</b>	<b>95.0</b>	<b>97.3</b>	<b>96.1</b>
PSAD [19]	AAAI'24	✓	<b>92.5</b>	<b>98.7</b>	<b>94.9</b>	<b>97.5</b>	<b>90.6</b>	<b>94.9</b>
CSAD [16]	BMVC'24	✓	<b>92.8</b>	<b>95.3</b>	<b>98.7</b>	<b>96.5</b>	<b>93.5</b>	<b>95.3</b>
<b>SALAD<sup>†</sup></b>		✓	<b>94.2</b>	<b>99.3</b>	<b>99.1</b>	96.6	<b>97.0</b>	<b>97.2</b>

Table 1. Anomaly detection (AUROC) on MVTec LOCO [3]. **First**, **second** and **third** place are marked. **SALAD<sup>†</sup>** is trained using composition maps from PSAD [19].

Category	Venue	Supervised Masks	Breakfast box		Juice bottle		Pushpins		Screw bag		Splicing conn.		Average	
			Log.	Str.	Log.	Str.	Log.	Str.	Log.	Str.	Log.	Str.	Log.	Str.
DRÆM [42]	ICCV'21		75.1	85.4	97.8	90.8	55.7	81.5	56.2	85.0	75.2	95.5	72.0	87.6
TransFusion [12]	ECCV'24		78.8	86.0	<b>99.8</b>	<b>99.6</b>	56.4	71.3	54.8	88.2	69.2	<b>98.2</b>	71.8	88.7
DSR [44]	ECCV'22		83.6	<b>88.0</b>	<b>99.5</b>	98.9	69.4	83.6	54.4	75.4	75.9	94.9	75.0	90.2
Sinbad [7]	ArXiv'23		<b>97.7</b>	85.9	97.1	91.7	88.9	78.9	81.1	<b>92.4</b>	91.5	78.3	<b>91.2</b>	85.5
ComAD + Patchcore [24]	AEI'23		81.6	<b>91.1</b>	98.2	95.0	91.1	<b>95.7</b>	<b>88.5</b>	71.9	<b>94.9</b>	93.3	89.4	90.9
SLSG [39]	PR'24		<b>93.7</b>	84.5	99.2	98.8	<b>97.4</b>	93.4	69.4	<b>91.6</b>	88.4	88.5	89.6	<b>91.4</b>
SAM-LAD [28]	KBS'25		<b>96.7</b>	85.2	98.7	96.5	<b>97.2</b>	79.2	<b>95.2</b>	77.9	91.4	88.6	<b>95.8</b>	85.5
EfficientAD [1]	WACV'24		87.4	<b>89.5</b>	98.8	<b>99.1</b>	93.5	<b>93.6</b>	58.1	89.1	<b>96.0</b>	<b>98.2</b>	86.8	<b>94.7</b>
<b>SALAD</b>			92.9	85.7	<b>99.7</b>	<b>99.6</b>	<b>100.0</b>	<b>98.8</b>	<b>93.9</b>	<b>96.0</b>	<b>96.0</b>	<b>98.6</b>	<b>96.5</b>	<b>95.7</b>
PSAD [19]	AAAI'24	✓	<b>100.0</b>	<b>84.9</b>	99.1	98.2	<b>100.0</b>	<b>89.8</b>	<b>99.3</b>	<b>95.7</b>	<b>91.9</b>	<b>89.3</b>	98.1	<b>91.6</b>
CSAD [16]	BMVC'24	✓	<b>94.4</b>	<b>91.1</b>	<b>94.9</b>	<b>95.6</b>	<b>99.5</b>	<b>97.8</b>	<b>99.9</b>	<b>93.2</b>	<b>94.8</b>	<b>92.2</b>	<b>96.7</b>	<b>94.0</b>
<b>SALAD<sup>†</sup></b>		✓	<b>99.6</b>	<b>88.8</b>	<b>99.6</b>	<b>98.9</b>	99.9	<b>98.3</b>	<b>98.6</b>	<b>94.7</b>	<b>95.8</b>	<b>98.6</b>	<b>98.7</b>	<b>95.8</b>

Table 2. Anomaly detection (AUROC) split by type (Logical/Structural) on MVTec LOCO [3].

Category	Logical	MVTec AD	VisA	Average
SimpleNet [25]		<b>99.6</b>	87.9	93.8
DRÆM [42]		98.0	88.7	93.4
TransFusion [12]		<b>99.2</b>	<b>98.5</b>	<b>98.9</b>
DSR [44]		98.2	91.6	94.9
RD4AD [9]		98.5	96.0	97.3
Patchcore [33]		<b>99.1</b>	94.3	96.7
EfficientAD [1]	✓	<b>99.1</b>	<b>98.1</b>	<b>98.6</b>
PUAD [36]	✓	98.5	96.9	97.7
CSAD [16]	✓	96.2	89.5	92.6
PSAD [19]	✓	98.0	90.3	94.2
<b>SALAD</b>	✓	98.8	<b>97.9</b>	<b>98.3</b>

Table 3. Anomaly detection (AUROC) on MVTec AD [2] and VisA [46].

trained with a predefined train-test split for each category, and the same hyperparameters were set across all datasets and all categories. Anomaly Scores are normalized using the anomaly scores from the validation set. As MVTec AD and VisA do not have a validation set, we create it by taking a part of the training set (more specifically, 10%).

### 4.3. Experimental results

Following recent literature [1, 33], anomaly detection performance is evaluated using the Area Under the Receiver Operator Curve (AUROC). Most concurrent works [7, 16, 19, 36] do not report localisation results due to the ambiguity of the ground-truth masks regarding logical anomalies. Due to that, we omitted them from the main paper. However, they are reported in the supplementary material for completeness.

The results for anomaly detection on MVTec LOCO are shown in Table 1. SALAD achieves the best score with a mean average AUROC of 96.1%, beating the best previous method by a significant margin of 3.0 percentage points. To enable future comparison, SALAD was also trained with composition maps from PSAD [19], which were obtained in a supervised manner. In this scenario, marked SALAD<sup>†</sup>, it outperforms all methods with supervised (or category-tuned) composition maps by a significant margin of 1.9 percentage points. The results split by the anomaly type can be seen in Table 2. SALAD achieves both the highest logical anomaly detection and the highest structural anomaly detection result,

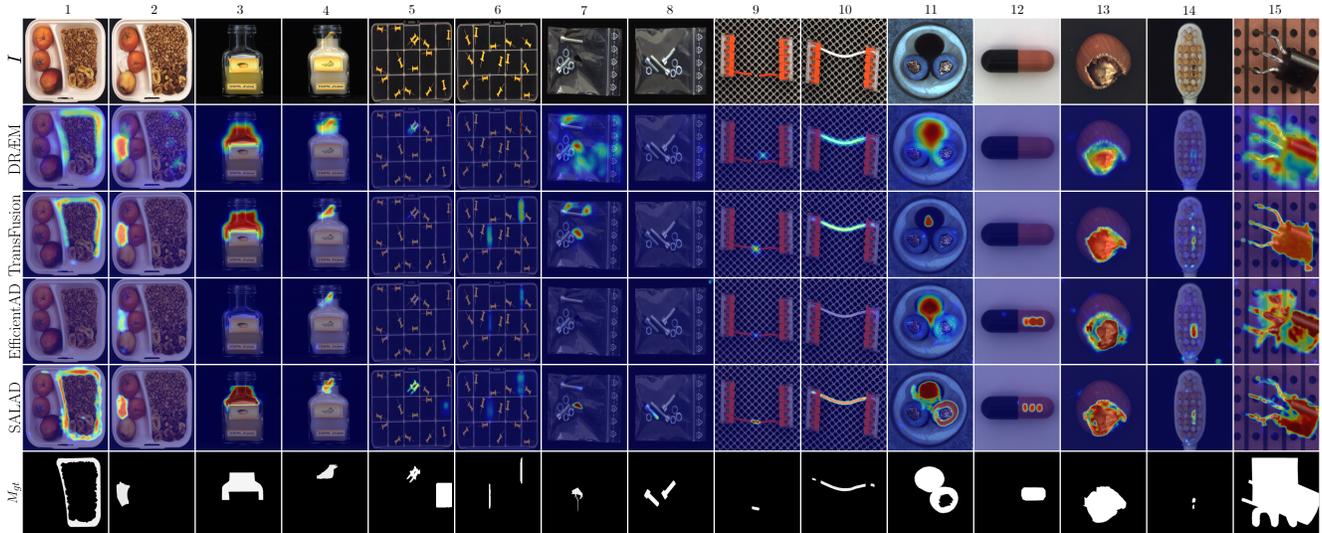


Figure 5. Qualitative comparison of the anomaly segmentation masks produced by SALAD and three other state-of-the-art methods. In the first row, the image is shown. In the next four rows, the anomaly segmentations produced by DRÆM [42], TransFusion [12], EfficientAD [1] and SALAD are depicted, and in the last row, the ground truth mask is depicted. For SALAD, we visualised the sum of  $A_a$  and  $A_c$  (the outputs of the appearance and the composition branch).

suggesting that introducing the composition branch improves anomaly detection efficiency.

Table 3 shows the results on the structural anomaly datasets MVTEC AD [2] and Visa [46]. SALAD achieves a state-of-the-art result with a mean average AUROC of 98.3% over both datasets. SALAD outperforms the vast majority of logical anomaly detection methods, showing superior performance in scenarios with only structural anomalies despite not being specialised for this task.

Qualitative examples can be seen in Figure 5. SALAD produces accurate anomaly localisation even in hard near-distribution cases, with which previously proposed methods struggled (Columns 5, 8, 11 and 12). The extra successful detections are primarily due to the composition branch, and the presented anomalies mainly concern the image’s composition. SALAD also detects all of the regions containing an anomaly as opposed to previous methods (Columns 1, 5 and 11). The better coverage comes from the composition branch, which detects different parts of the anomaly.

## 5. Ablation study

Ablation experiments validating the contributions of SALAD are performed. Results are shown in Tables 4 and 5.

**Branch performance** To show each branch’s overall importance and especially the composition branch’s importance, we evaluated the model by excluding one branch at a time. Dropping the appearance branch leads to a 0.9 p. p. drop in logical anomalies and a 4.5 p. p. drop in structural anomalies. Dropping the composition branch results in a 3.5 p. p. drop in performance with logical anomalies and a 1.0 p. p. drop

with structural anomalies. Dropping the stat branch results in the lowest overall performance drop with an overall drop of 1.8 p. p. While not using the composition branch would still achieve SOTA results, it wouldn’t improve them significantly. This confirms that modelling the composition map distribution will improve logical anomaly detection.

**Choice of the architecture** To show the generality of the proposed framework, we exchanged the appearance branch with three other state-of-the-art models: DSR [44], TransFusion [12] and DRÆM [25]. To maintain a unified evaluation process, we disabled the centre-cropping for TransFusion. No model selection strategy or parameter tuning was performed for each model. Using all three architectures for the appearance branch produced state-of-the-art results. TransFusion performs better than the other two due to better overall performance with structural anomalies, where the composition and the global branch struggle. The results still show robustness to the choice of the appearance branch.

**Different synthetic anomaly generation strategies** Three synthetic anomaly generation techniques are used during training - DRÆM anomaly, component inpainting and component removal. Each strategy was removed from training to verify its contribution to the overall performance. Removing each strategy resulted in a drop in performance. The highest drop is seen by removing the component removal strategy (0.9 p. p.) and the lowest when we remove the component inpainting strategy (0.4 p. p.). The results show the contribution and necessity of each strategy.

**Different global representation** To investigate the improvement of the global representation, we exchanged our global

Group	Condition	Det. Logical	Det. Struct.	Det. Avg
<i>Full model</i>	w/o Appearance branch	95.6 (-0.9)	91.2 (-4.5)	93.4 (-2.7)
	w/o Composition branch	93.0 (-3.5)	94.7 (-1.0)	93.8 (-2.3)
	w/o Stat branch	93.6 (-2.9)	94.9 (-0.8)	94.2 (-1.9)
<i>Appearance branch</i>	DSR [44]	95.2 (-1.3)	93.7 (-2.0)	94.4 (-1.7)
	TransFusion [12]	95.3 (-1.2)	95.5 (-0.2)	95.4 (-0.7)
	DRÆM [42]	94.8 (-1.7)	93.1 (-2.6)	94.0 (-2.1)
<i>Composition branch</i>	w/o DRÆM anomalies	96.0 (-0.5)	95.2 (-0.5)	95.6 (-0.5)
	w/o component inpainting anomalies	95.9 (-0.6)	95.6 (-0.1)	95.7 (-0.4)
	w/o component removal anomalies	95.5 (-1.0)	95.1 (-0.6)	95.2 (-0.9)
<i>Global branch</i>	Only Global Vector	95.6 (-0.9)	95.6 (-0.1)	95.6 (-0.5)
	$g_{\text{DINOv2}}$	95.4 (-0.5)	93.5 (-2.4)	94.4 (-1.7)
	$g_{\text{DINO}}$	96.2 (-0.3)	94.3 (-1.4)	95.3 (-0.8)
	$g_{\text{ResNet50}}$	92.7 (-3.8)	92.8 (-3.0)	92.7 (-3.4)
<i>Object composition map generation</i>	DINOv2	95.6 (-0.9)	95.4 (-0.3)	95.5 (-0.6)
	4 clusters	95.4 (-0.3)	95.4 (-0.8)	95.4 (-0.7)
	8 clusters	96.0 (-0.5)	95.8 (+0.1)	95.9 (-0.2)
<i>SALAD</i>	EfficientAD, DINO, 6 clusters	96.5	95.7	96.1

Table 4. Ablation study results. Results are reported for MVTEC LOCO [3] in AUROC and separated by the type of anomalies. In the last column, the average for both types is reported. The difference from the base model is shown in blue.

representation with the one from PUAD [36], that is, using only the global mean vector. The performance drops by 0.9 p. p. in logical anomalies and by 0.1 p. p. in structural anomalies. The drop is due to the lack of spatial information inside the global representation. The results indicate that our representation does indeed improve the results.

**Different feature representation for the global representation** To verify the effectiveness of EfficientAD’s feature extractor for the global representation, we exchanged it with a few different high-performing feature extractors – DINOv2 [27], DINO [6] and ResNet50 [15]. The performance drops the least (0.8 p. p.) when using DINO and the most using ResNet50 (3.4 p. p.). We hypothesise this is due to the subpar representation of ResNet50 features. However, it shows the importance of choosing the right feature extractor for the global representation.

**Different feature extractor for the composition map generation** To investigate the importance of the feature extractor in the object composition map generation, we exchanged DINO with DINOv2 [27]. The performance drops by 0.9 p. p. in logical anomalies and by 0.3 p. p. in structural anomalies. The drop is due to consistent small mistakes in the generated composition maps made by DINOv2. Some examples are in the Supplementary material. Nevertheless, the results suggest that the choice of feature extractor for the composition map generation is robust.

**Different number of clusters** To show the robustness of the cluster number parameter, we also evaluated our model for  $K = 4$  and  $K = 8$ . Having 4 clusters results in a 0.6 p. p. drop in overall performance, and having 8 clusters results in a 0.2 p. p. drop in overall performance. These results suggest that the results are robust when  $K$  gets high enough. If  $K$  is

Method	DRÆM [42]	Patchcore [33]	EfficientAD [1]	<i>SALAD</i>
Inference [ms]	52.6	224.4	6.2	64.6

Table 5. Results for average inference time of a single sample with NVIDIA A100 GPU. Inference times are reported in milliseconds.

too low, the results are lower. Qualitative examples and the results for other values are in the Supplementary material.

**Inference Speed and Computational Complexity** The inference speed can be seen in Table 5. SALAD is faster than Patchcore [33] and lags slightly behind DRÆM [42] and EfficientAD [1]. SALAD could be further optimised for speed by successfully parallelising each branch. SALAD requires approximately 1.5 hours to train on a single A100 GPU and has 65.1 million parameters.

## 6. Conclusion

A novel model for logical anomaly detection, SALAD, is proposed. Unlike recent methods, SALAD explicitly models object composition information by introducing a novel discriminatively trained composition branch. For this purpose, it introduces a novel automatic composition map generation strategy and an anomaly simulation process, facilitating discriminative training. SALAD achieves a new state-of-the-art of 96.1% AUROC on the MVTEC LOCO Dataset, outperforming all previous methods by a significant margin of 3.0 percentage points. Furthermore, SALAD also performs very well on datasets with only structural anomalies, achieving 98.9% on MVTEC AD and 97.9% on VisA. Further interaction between branches in the architecture may improve performance and is a good avenue for future research. The results indicate that explicit composition distribution modelling is also a viable future research direction.

**Acknowledgements** This work was in part supported by the ARIS research project MUXAD (J2-60055), research programme P2-0214 and the supercomputing network SLING (ARNES, EuroHPC Vega).

## References

- [1] Kilian Batzner, Lars Heckler, and Rebecca König. EfficientAD: Accurate Visual Anomaly Detection at Millisecond-Level Latencies. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 128–138, 2024. 2, 3, 6, 7, 8
- [2] Paul Bergmann, Kilian Batzner, Michael Fauser, David Sattler, and Carsten Steger. The MVTec Anomaly Detection Dataset: A Comprehensive Real-World Dataset for Unsupervised Anomaly Detection. *International Journal of Computer Vision*, 129(4):1038–1059, 2021. 1, 2, 5, 6, 7
- [3] Paul Bergmann, Kilian Batzner, Michael Fauser, David Sattler, and Carsten Steger. Beyond Dents and Scratches: Logical Constraints in Unsupervised Anomaly Detection and Localization. *International Journal of Computer Vision*, 130(4):947–969, 2022. 2, 5, 6, 8, 1, 4
- [4] Jakob Božič, Domen Tabernik, and Danijel Skočaj. Mixed supervision for surface-defect detection: From weakly to fully supervised learning. *Computers in Industry*, 129:103459, 2021.
- [5] Jakob Božič, Matic Fučka, Vitjan Zavrtanik, and Danijel Skočaj. Robustness of unsupervised methods for image surface-anomaly detection. *Pattern Analysis and Applications*, 28(2):99, 2025. 1
- [6] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 4, 5, 8, 1, 2, 3
- [7] Niv Cohen, Issar Tzachor, and Yedid Hoshen. Set Features for Anomaly Detection. *arXiv preprint arXiv:2311.14773*, 2023. 2, 6
- [8] Songmin Dai, Yifan Wu, Xiaoqiang Li, and Xiangyang Xue. Generating and reweighting dense contrastive patterns for unsupervised anomaly detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1454–1462, 2024. 2
- [9] Hanqiu Deng and Xingyu Li. Anomaly Detection via Reverse Distillation From One-Class Embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9737–9746, 2022. 1, 2, 6
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 3
- [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 5
- [12] Matic Fučka, Vitjan Zavrtanik, and Danijel Skočaj. TransFusion – A Transparency-Based Diffusion Model for Anomaly Detection. In *European conference on computer vision*, pages 91–108. Springer, 2025. 1, 2, 6, 7, 8
- [13] Denis Gudovskiy, Shun Ishizaka, and Kazuki Kozuka. CFLOW-AD: Real-Time Unsupervised Anomaly Detection with Localization via Conditional Normalizing Flows. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 98–107, 2022. 1, 2
- [14] Hwei Guo, Liping Ren, Jingjing Fu, Yuwang Wang, Zhizheng Zhang, Cuiling Lan, Haoqian Wang, and Xinwen Hou. Template-guided hierarchical feature restoration for anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6447–6458, 2023. 2, 6
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 8, 1, 2, 3
- [16] Yu-Hsuan Hsieh and Shang-Hong Lai. CSAD: Unsupervised component segmentation for logical anomaly detection. In *In Proceedings of the British Machine Vision Conference (BMVC)*, 2024. 1, 2, 3, 6
- [17] Er Jin, Qihui Feng, Yongli Mou, Stefan Decker, Gerhard Lakemeyer, Oliver Simons, and Johannes Stegmaier. LogicaAD: Explainable Anomaly Detection via VLM-based Text Feature Extraction. 2025. 2, 6
- [18] Lei Ke, Mingqiao Ye, Martin Danelljan, Yu-Wing Tai, Chi-Keung Tang, Fisher Yu, et al. Segment anything in high quality. *Advances in Neural Information Processing Systems*, 36, 2024. 4, 5, 1, 2, 3
- [19] Soopil Kim, Sion An, Philip Chikontwe, Myeongkyun Kang, Ehsan Adeli, Kilian M Pohl, and Sang Hyun Park. Few shot part segmentation reveals compositional logic for industrial anomaly detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 8591–8599, 2024. 1, 2, 3, 6
- [20] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. 5
- [21] Yufei Liang, Jiangning Zhang, Shiwei Zhao, Runze Wu, Yong Liu, and Shuwen Pan. Omni-Frequency Channel-Selection Representations for Unsupervised Anomaly Detection. *IEEE Transactions on Image Processing*, 32:4327–4340, 2023. 2
- [22] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollar. Focal loss for dense object detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017. 4
- [23] Tongkun Liu, Bing Li, Xiao Du, Bingke Jiang, Leqi Geng, Feiyang Wang, and Zhuo Zhao. FAIR: Frequency-Aware Image Restoration for Industrial Visual Anomaly Detection. *arXiv preprint arXiv:2309.07068*, 2023. 1, 2
- [24] Tongkun Liu, Bing Li, Xiao Du, Bingke Jiang, Xiao Jin, Liuyi Jin, and Zhuo Zhao. Component-aware anomaly detection framework for adjustable and logical industrial visual inspection. *Advanced Engineering Informatics*, 58:102161, 2023. 1, 3, 6, 2

- [25] Zhikang Liu, Yiming Zhou, Yuansheng Xu, and Zilei Wang. SimpleNet: A Simple Network for Image Anomaly Detection and Localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20402–20411, 2023. 1, 2, 6, 7
- [26] Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019. 5
- [27] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafranec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2024. 8, 1
- [28] Yun Peng, Xiao Lin, Nachuan Ma, Jiayuan Du, Chuangwei Liu, Chengju Liu, and Qijun Chen. SAM-LAD: Segment Anything Model meets zero-shot logic anomaly detection. *Knowledge-Based Systems*, page 113176, 2025. 2, 6
- [29] Ken Perlin. An image synthesizer. *ACM Siggraph Computer Graphics*, 19(3):287–296, 1985. 5
- [30] Jonathan Pirnay and Keng Chai. Inpainting transformer for anomaly detection. In *Image Analysis and Processing – ICIAP 2022*, pages 394–406, Cham, 2022. Springer International Publishing. 2
- [31] Oliver Rippel, Patrick Mertens, and Dorit Merhof. Modeling the distribution of normal data in pre-trained deep features for anomaly detection. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages 6726–6733. IEEE, 2021. 5
- [32] Blaž Rolih, Matic Fučka, and Danijel Skočaj. SuperSimpleNet: Unifying Unsupervised and Supervised Learning for Fast and Reliable Surface Defect Detection. In *International Conference on Pattern Recognition*, 2024. 2
- [33] Karsten Roth, Latha Pemula, Joaquin Zepeda, Bernhard Schölkopf, Thomas Brox, and Peter Gehler. Towards Total Recall in Industrial Anomaly Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14318–14328, 2022. 1, 2, 6, 8
- [34] Marco Rudolph, Tom Wehrbein, Bodo Rosenhahn, and Bastian Wandt. Asymmetric student-teacher networks for industrial anomaly detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2592–2602, 2023. 2
- [35] Carole H Sudre, Wenqi Li, Tom Vercauteren, Sebastien Ourselin, and M Jorge Cardoso. Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: Third International Workshop, DLMIA 2017, and 7th International Workshop, ML-CDS 2017, Held in Conjunction with MICCAI 2017, Québec City, QC, Canada, September 14, Proceedings 3*, pages 240–248. Springer, 2017. 4
- [36] Shota Sugawara and Ryuji Imamura. PUAD: Frustratingly simple method for robust anomaly detection. In *2024 IEEE international conference on image processing (ICIP)*, 2024. 1, 2, 6, 8
- [37] Julian Wyatt, Adam Leach, Sebastian M. Schmon, and Chris G. Willcocks. AnoDDPM: Anomaly Detection With Denoising Diffusion Probabilistic Models Using Simplex Noise. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 650–656, 2022. 2
- [38] Minghui Yang, Peng Wu, and Hui Feng. MemSeg: A semi-supervised method for image surface defect detection using differences and commonalities. *Engineering Applications of Artificial Intelligence*, 119:105835, 2023. 2, 4
- [39] Minghui Yang, Jing Liu, Zhiwei Yang, and Zhaoyang Wu. SLGS: Industrial Image Anomaly Detection by Learning Better Feature Embeddings and One-Class Classification. *Pattern Recognition*, 156:110862, 2024. 2, 6
- [40] Xincheng Yao, Ruoqi Li, Zefeng Qian, Yan Luo, and Chongyang Zhang. Focus the discrepancy: Intra-and inter-correlation learning for image anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6803–6813, 2023. 2
- [41] Jiawei Yu, Ye Zheng, Xiang Wang, Wei Li, Yushuang Wu, Rui Zhao, and Liwei Wu. FastFlow: Unsupervised Anomaly Detection and Localization via 2D Normalizing Flows. *arXiv preprint arXiv:2111.07677*, 2021. 2
- [42] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. DRAEM-A discriminatively trained reconstruction embedding for surface anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8330–8339, 2021. 1, 2, 3, 4, 5, 6, 7, 8
- [43] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. Reconstruction by inpainting for visual anomaly detection. *Pattern Recognition*, 112:107706, 2021. 2
- [44] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. DSR—A dual subspace re-projection network for surface anomaly detection. In *European Conference on Computer Vision*, pages 539–554. Springer, 2022. 1, 2, 6, 7, 8
- [45] Jie Zhang, Masanori Suganuma, and Takayuki Okatani. Contextual affinity distillation for image anomaly detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 149–158, 2024. 2
- [46] Yang Zou, Jongheon Jeong, Latha Pemula, Dongqing Zhang, and Onkar Dabeer. SPot-the-Difference Self-supervised Pre-training for Anomaly Detection and Segmentation. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXX*, pages 392–408. Springer, 2022. 1, 2, 5, 6, 7

# SALAD 🥗 – Semantics-Aware Logical Anomaly Detection

## Supplementary Material

This supplementary material includes additional information and visualisations. More specifically, we ablate the object composition map generation and add a further experiment to verify the importance of each branch. Ultimately, we add localisation results for MVTEC LOCO and more qualitative examples.

### A. Limitations and Failure Cases

Composition map creation depends on the performance of SAM-HQ and DINO, although they perform very well across diverse datasets. In the future, this can even be improved with stronger models (e.g. Perception Encoder). Additionally, there are also some cases (some are depicted in Figure 1) in which SALAD fails to detect anomalies. SALAD mostly fails on extremely near-distribution structural (Columns 1-6) and logical anomalies (Columns 7-10). Architectural improvements to the compositional and appearance branch might improve this.

### B. Differences from other methods utilising composition maps

Currently, there are three different methods using composition maps – ComAD [24], CSAD [16] and PSAD [19]. SALAD’s biggest difference from all three is the introduction of a specialised composition branch. This means SALAD is directly trained on the composition maps in contrast to the other three. Additionally, CSAD and PSAD require extra category-specific information, either via hand-labelled samples or via category-specific composition map procedures. SALAD performs all of this automatically without any additional information. ComAD produces composition maps of low quality, whilst SALAD produces high-quality maps.

### C. Object composition map generation ablation

This section compares the design choices for the object composition map generation. First, we examine the effect of using a different feature extractor than DINO [6]. Then, we examine the importance of the number of clusters parameter. **Different feature extractor** To evaluate the choice of feature extractor in component map generation, we replaced the original with other standard feature extractors: ResNet50 [15], ResNet50 DINO [6], SAM [18], ViT DINOv2 [27]. Their performance is qualitatively evaluated by comparing feature clusters  $C_{feat}$ , pseudo labels  $C_{pseudo}$ , and final composition maps  $C$ . Figure 2 depicts that ResNet50, ViT DINOv2, and ViT DINO cluster features effectively, discriminating similar objects (e.g., Columns 5 and 6). In contrast, ResNet, DINO

Condition	Det. Logical	Det. Struct.	Det. Avg
Only Appearance branch	87.5 (-9.0)	94.1 (-1.6)	90.8 (-5.3)
Only Composition branch	88.1 (-8.4)	82.8 (-12.9)	85.4 (-10.7)
Only Global branch	90.8 (-5.7)	87.3 (-8.4)	89.1 (-8.1)
SALAD	96.5	95.7	96.1

Table 1. Branch importance is evaluated with the downstream performance in Anomaly detection on the MVTEC LOCO dataset [3] (results are presented in AUROC). The importance is evaluated by using only one branch. The results are categorised by anomaly type, and the overall average detection rate is reported in the final column. The performance difference relative to the base model is highlighted in blue.

and SAM yield poor clusters, as seen in Columns 3 and 4. This pattern continues with pseudo labels in Figure 3, where ResNet DINO and SAM exhibit loss of detail and class mismatches (Columns 7 and 8). Due to noisy pseudo labels, the lightweight semantic segmentation model struggles with generalisation (Figure 4, Columns 7 and 8). Consequently, we evaluated downstream anomaly detection performance only for ViT DINO [6] and ViT DINOv2 [27], with results detailed in the main paper.

**Different number of clusters** To investigate the importance of the number of clusters during composition map generation, we qualitatively and quantitatively assessed the output composition maps. More specifically, we checked for different values of  $K$  ranging from 4 to 8. We qualitatively assessed the feature clusters, pseudo-labels and the generated composition maps. The results for different stages in the pipeline can be seen in Figure 5, Figure 6 and Figure 7. From the Figures, it can be seen that there are no significant differences, especially with the final composition maps. This would suggest that the choice of the number of clusters is robust (once it is high enough). In Figure 8, the effect of this parameter on downstream anomaly detection is depicted. All values are above the current state-of-the-art, suggesting that the parameter choice is robust.

### D. Branch importance

To further show the overall importance of each branch, we evaluated the model by using one branch at a time. The results can be seen in Table 1 and in Table 2. Using only the appearance branch leads to a drop in performance of 9.0 percentage points (p. p.) for logical anomalies and 1.6 p. p. for structural anomalies. Using only the composition branch leads to a drop of 8.4 p. p. on logical anomalies and a 12.9 p. p. drop for structural anomalies. By solely using the global

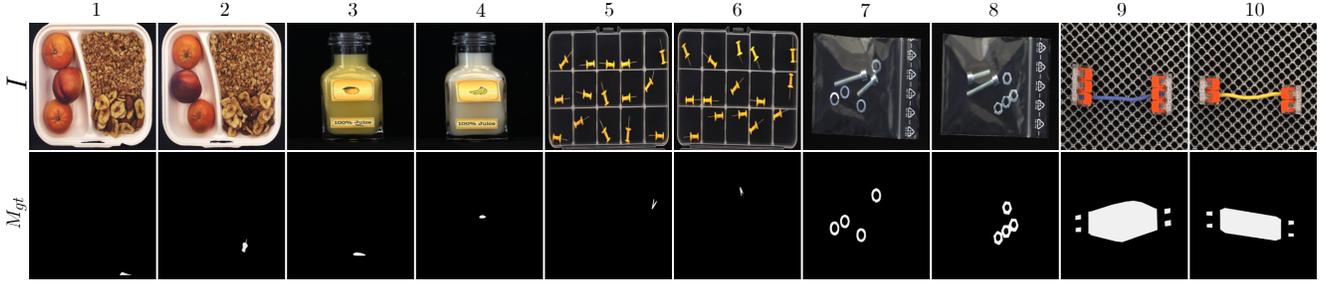


Figure 1. Failure case results. In all of the cases, SALAD produces a very low anomaly score. Most of the cases also represent near-distribution logical and structural anomalies.

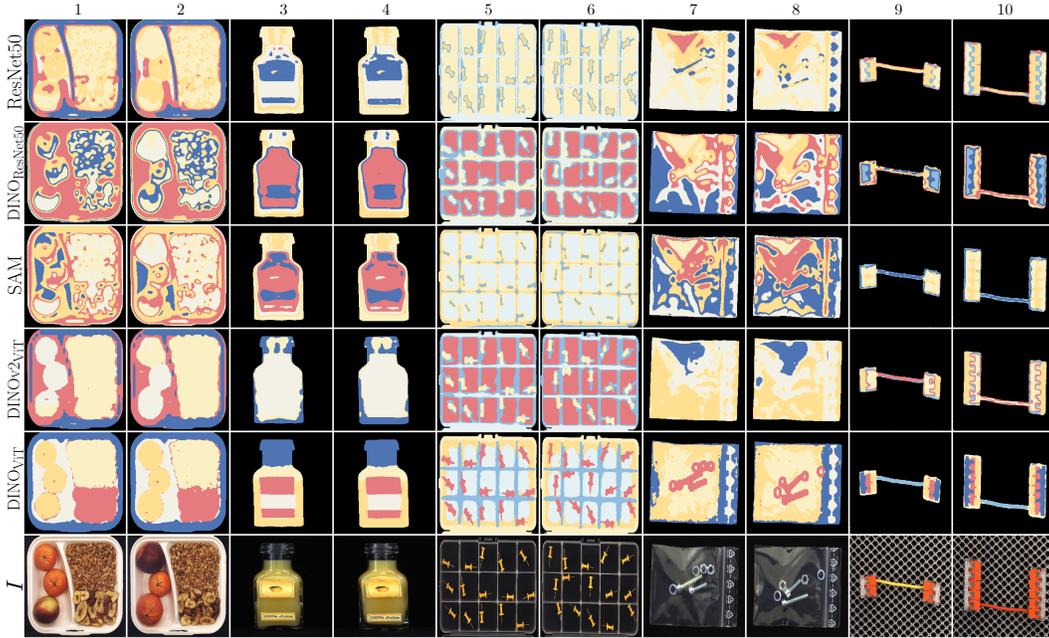


Figure 2. Qualitative comparison of the feature clusters  $C_{feat}$  produced by 5 different feature extractors: ResNet50 [15], ResNet50 DINO [6], SAM [18] and ViT DINO [6]. In the bottom row, the original image  $I$  is shown. It can be observed that both ViT DINOv2 and ViT DINO separate the objects effectively (e.g. Columns 5 and 6), while other feature extractors face problems (e.g. Columns 3 and 4).

branch, the performance drops 5.7 p. p. for logical anomalies and 8.4 p. p. for structural anomalies. The results show that the branches complement each other, especially with logical anomalies.

### E. Localisation results for MVTec LOCO

Following recent literature [1, 39], the AUsPRO Metric [3] is used to evaluate the localisation performance. Again, it is important to highlight that most concurrent works [7, 16, 19, 24, 36] strayed away from reporting these results due to the ambiguity in pixel-level ground truths in images containing a logical anomaly. Some such cases are depicted in Figure 9. The localisation results on MVTec LOCO are given in Table 3. SALAD achieves the second-highest result with an AUsPRO of 68.7%.

### F. Additional qualitative results

In this section, we provide additional qualitative mask comparisons to the state-of-the-art models DRÆM [42], TransFusion [12] and EfficientAD [1]. The comparisons can be seen in Figure 10 and Figure 11. SALAD can detect more near-distribution and harder anomalies compared to previous state-of-the-art methods.

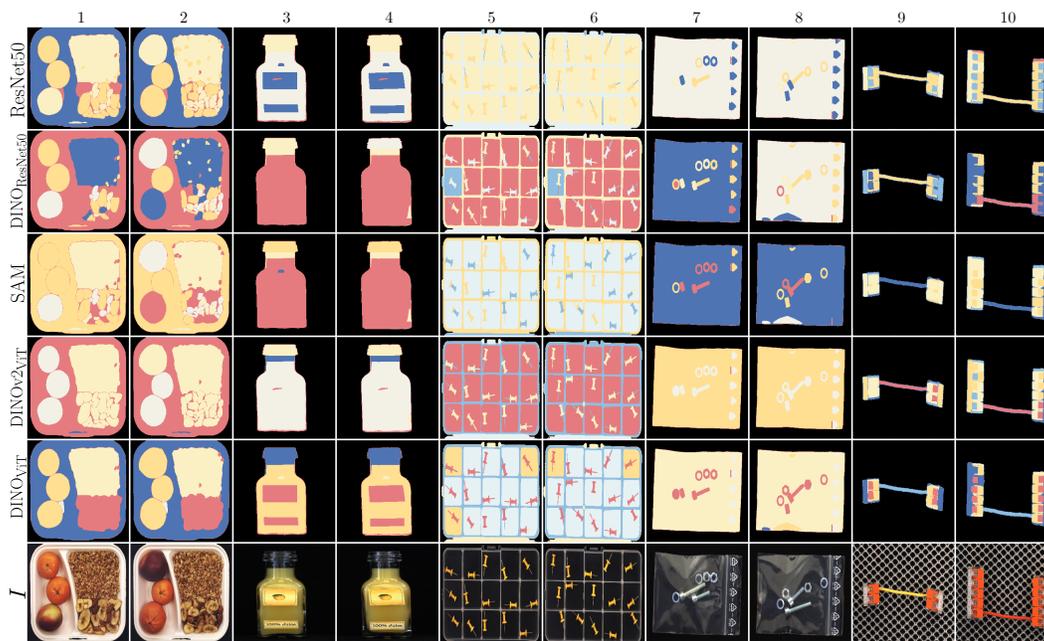


Figure 3. Qualitative comparison of the pseudo-labels  $C_{pseudo}$  produced by 5 different feature extractors: ResNet50 [15], ResNet50 DINO [6], SAM [18] and ViT DINO [6]. In the bottom row, the original image  $I$  is shown. Most methods do not face class mismatches and loss of detail except for ResNet50, DINO, and SAM (e.g. Columns 7 and 8).

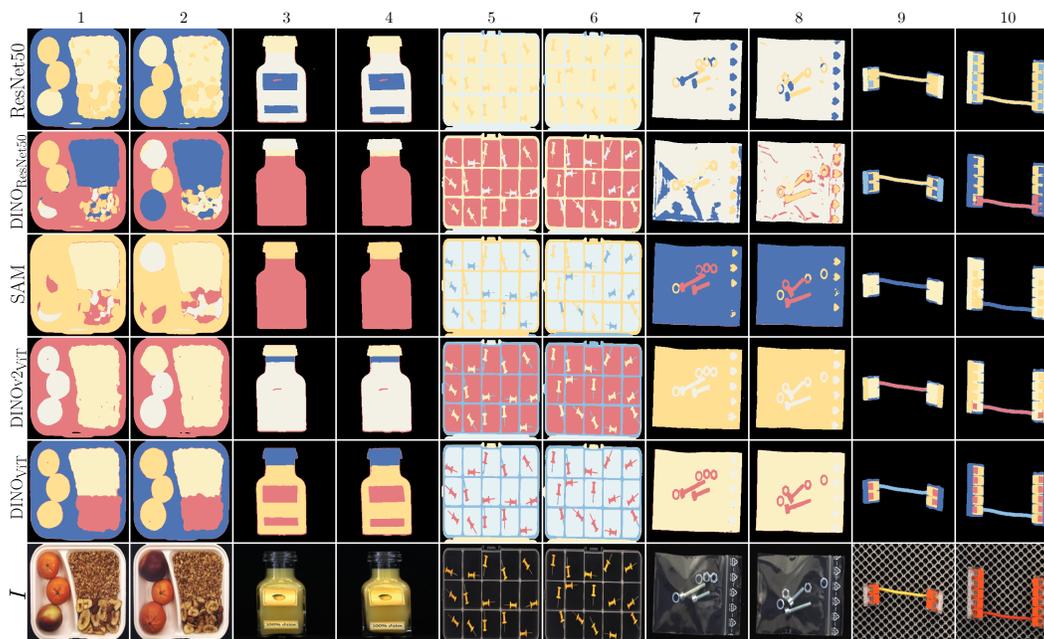


Figure 4. Qualitative comparison of the composition maps  $C$  produced by 5 different feature extractors: ResNet50 [15], ResNet50 DINO [6], SAM [18] and ViT DINO [6]. In the bottom row, the original image  $I$  is shown. While ViT DINO and ViT DINOv2 can generalise effectively, other methods face problems (e.g. Columns 7 and 8).

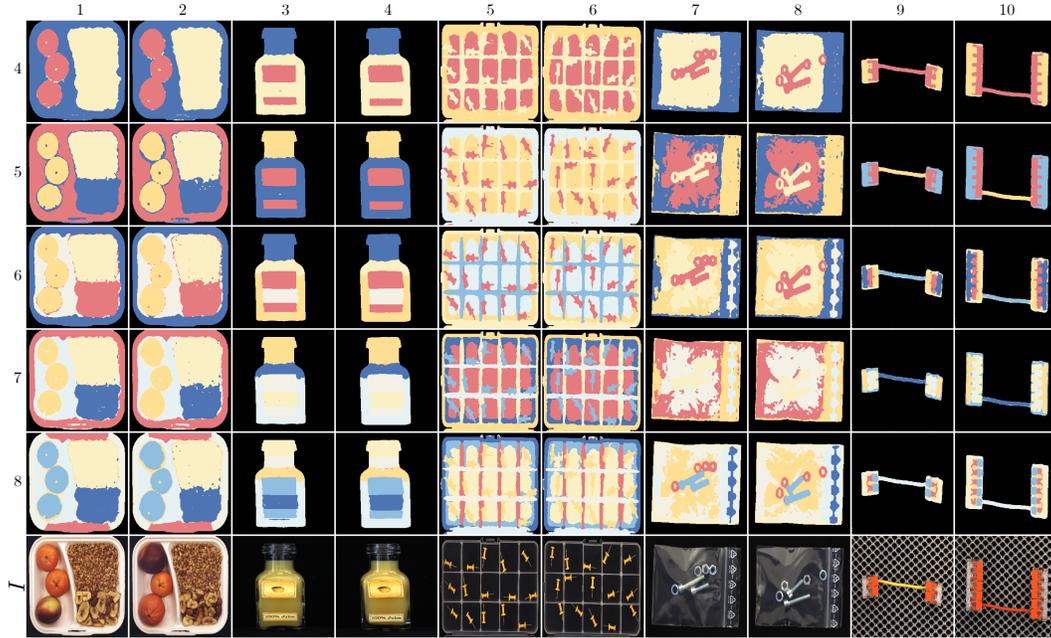


Figure 5. Qualitative comparison of the feature clusters  $C_{feat}$  produced by different numbers of clusters  $K$  (from 4 to 8). In the bottom row, the original image  $I$  is shown.

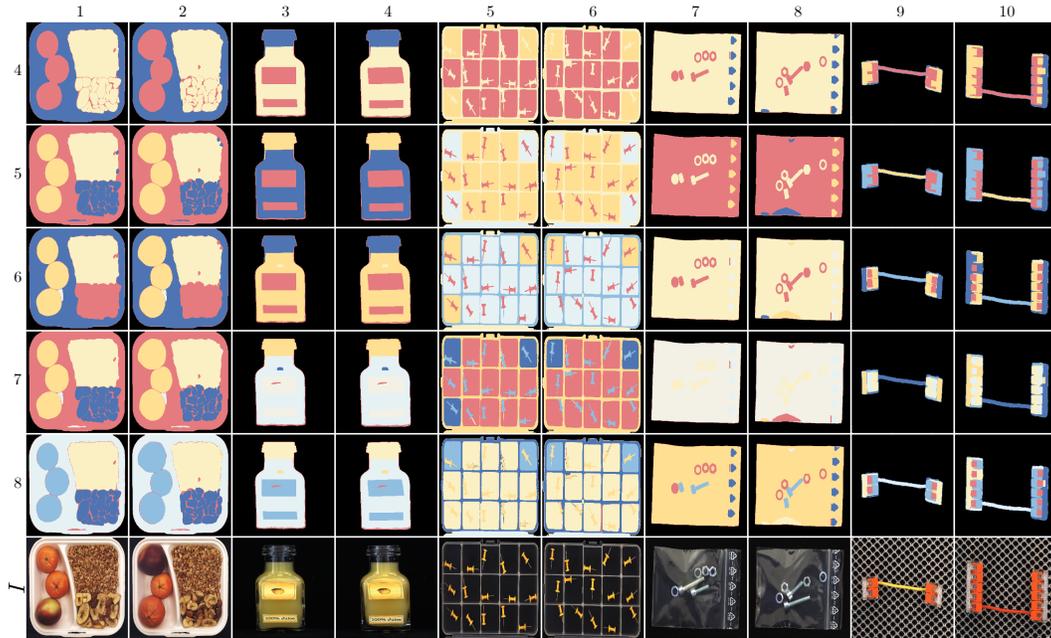


Figure 6. Qualitative comparison of the pseudo-labels  $C_{pseudo}$  produced by different numbers of clusters  $K$  (from 4 to 8).

Branch	Breakfast box	Juice bottle	Pushpins	Screw bag	Splicing conn.	Average
Only Appearance Branch	85.7	96.9	96.8	77.9	96.6	90.8
Only Composition Branch	77.1	87.0	87.7	88.2	86.2	85.4
Only Global Branch	82.2	97.7	91.8	86.3	87.3	89.1

Table 2. Anomaly detection (AUROC) for each branch on MVTec LOCO [3].

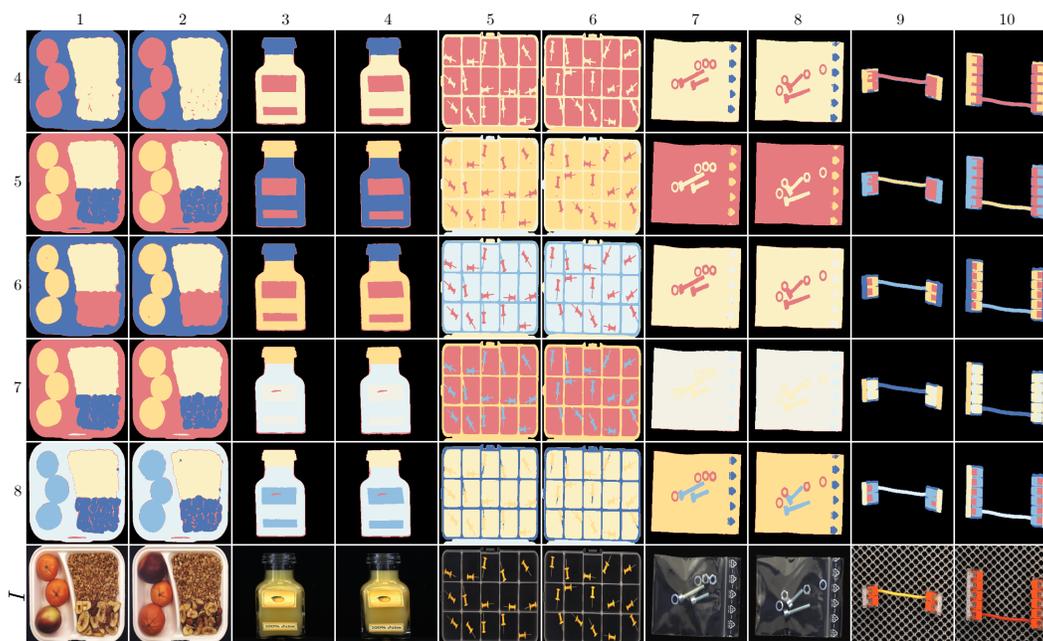


Figure 7. Qualitative comparison of the composition maps  $C$  produced by different numbers of clusters  $K$  (from 4 to 8).

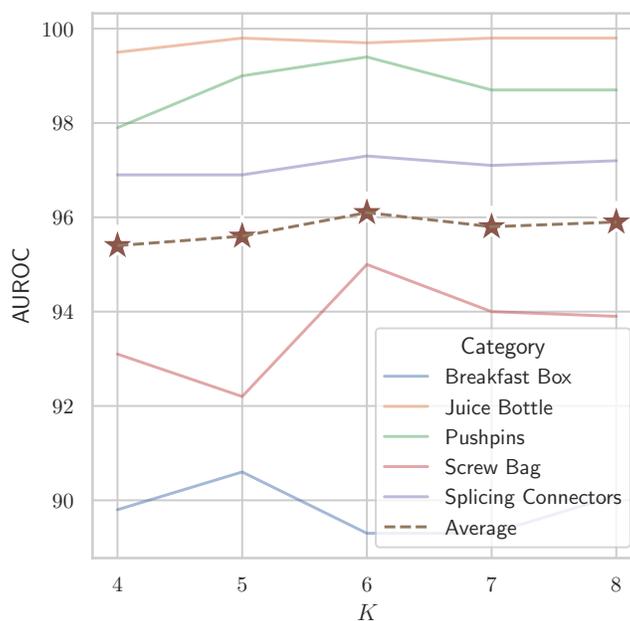


Figure 8. Anomaly detection performance on MVTec LOCO under different values for  $K$  in the object composition map generation. The default settings for  $K$  is 6.

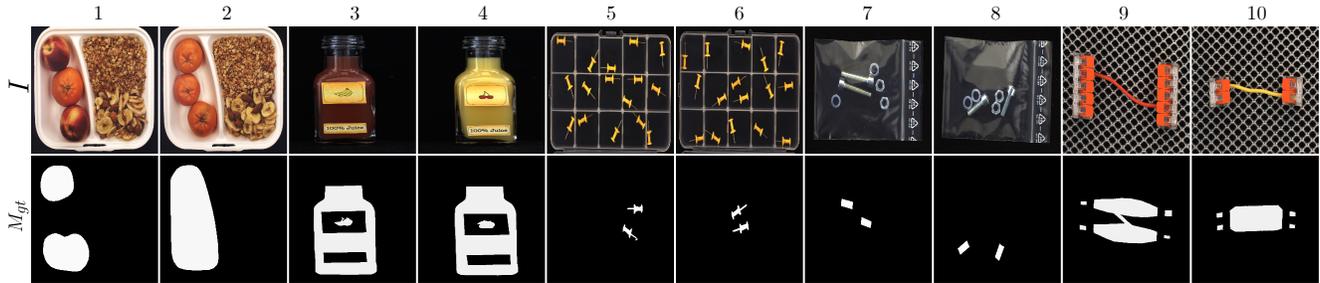


Figure 9. Examples of problematic pixel-level ground truths ( $M_{gt}$ ) and their corresponding images ( $I$ ) in MVTec LOCO [3] show issues with how annotations are done. The ground truths are designed to include all possible solutions, which causes ambiguity. For example, in Column 7, there are two long screws instead of one long screw and one short screw as expected. The annotation requires marking both long screws, even though marking one would still be a correct interpretation of the anomaly. This approach unfairly lowers the scores of methods that label only one screw, even if their prediction makes sense.

Category	SimpleNet [25]	DRÆM [42]	TransFusion [12]	DSR [44]	Patchcore [33]	SLSG [39]	EfficientAD [1]	SALAD
Breakfast box	38.8	49.9	53.5	49.9	46.6	65.9	60.4	49.1
Juice bottle	43.9	80.0	90.1	86.8	41.2	82.0	93.4	81.5
Pushpins	27.2	49.3	51.9	59.1	31.4	74.4	62.3	73.5
Screw bag	66.0	49.0	39.3	37.9	48.1	47.2	64.4	58.4
Splicing connectors	36.9	67.3	67.0	58.6	31.3	66.9	73.3	81.2
Average	36.3	59.1	60.4	58.5	39.7	67.3	69.4	68.7

Table 3. Anomaly localization (AUsPRO) on MVTec LOCO [3].



Figure 10. Qualitative comparison of the anomaly segmentation masks produced by SALAD and three other state-of-the-art methods on MVTEC LOCO. In the first row, the image is shown. In the next four rows, the anomaly segmentations produced by DRÆM [42], TransFusion [12], EfficientAD [1] and SALAD are depicted, and in the last row, the ground truth mask is shown. For SALAD, we visualised the sum of  $A_a$  and  $A_c$ .

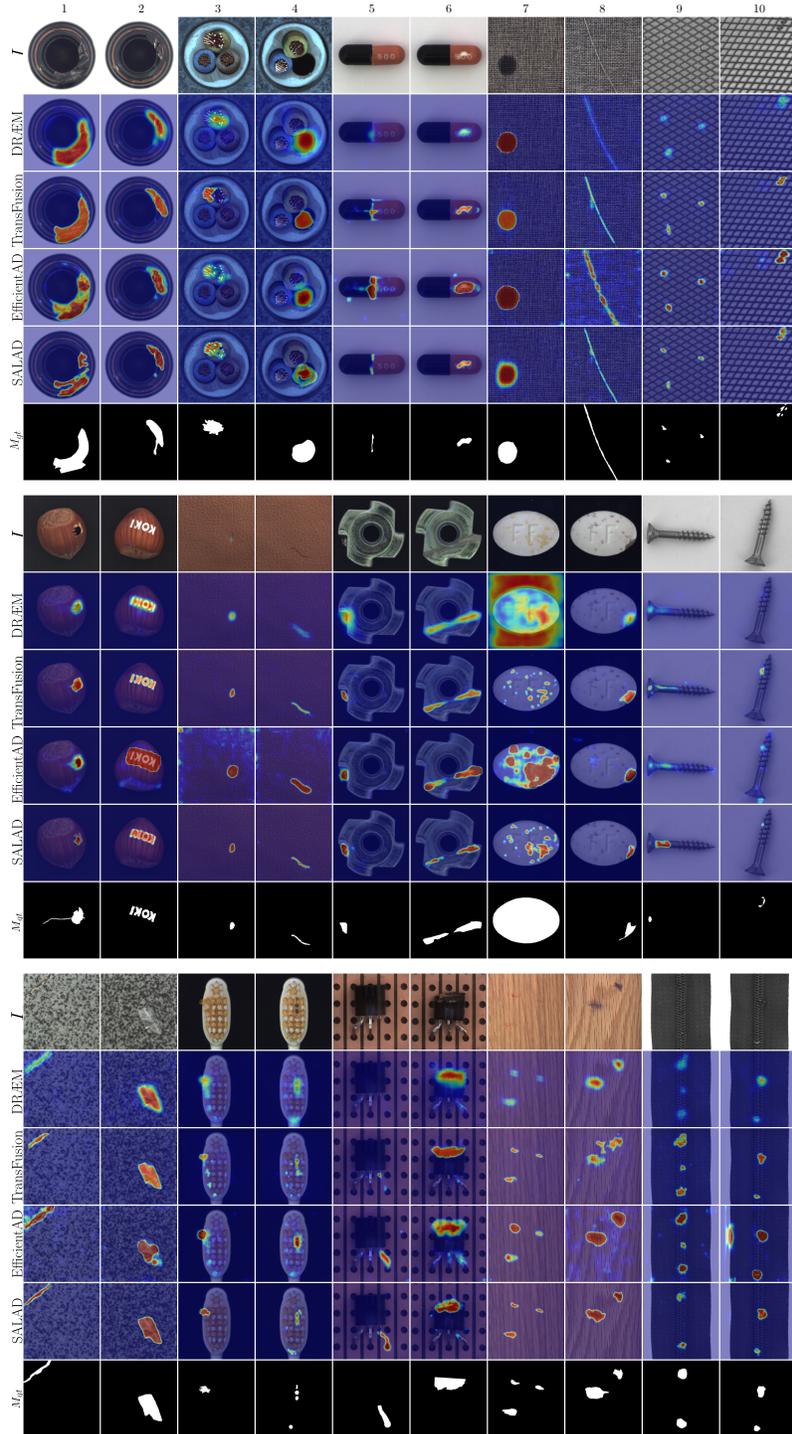


Figure 11. Qualitative comparison of the anomaly segmentation masks produced by SALAD and three other state-of-the-art methods on MVTEC AD. In the first row, the image is shown. In the next four rows, the anomaly maps produced by DRÆM [42], TransFusion [12], EfficientAD [1] and SALAD are depicted, and in the last row, the ground truth mask is shown.