# Supplementary Materials of "Anomaly Heterogeneity Learning for Open-set Supervised Anomaly Detection"

Jiawen Zhu[1], Choubo Ding[2], Yu Tian[3], and Guansong Pang[1*]

[1]School of Computing and Information Systems, Singapore Management University
[2]Australian Institute for Machine Learning, University of Adelaide
[3]Harvard Ophthalmology AI Lab, Harvard University

## A. Dataset Details

### A.1. Key Data Statistics

We conduct extensive experiments on nine real-world Anomaly Detection (AD) datasets. Table 1 provides key data statistics for all the datasets used in this study. We follow exactly the same settings as prior Open-set Supervised Anomaly Detection (OSAD) studies [4, 7]. Particularly, we follow the original settings of MVTec AD and split the normal samples into training and test sets; for the other eight datasets, the normal samples are randomly split into training and test sets using a 3:1 proportion.

| Dataset | | | Original Training | Original Test | |
|---|---|---|---|---|---|
| | $|C|$ | Type | Normal | Normal | Anomaly |
| Carpet | 5 | Texture | 280 | 28 | 89 |
| Grid | 5 | Texture | 264 | 21 | 57 |
| Leather | 5 | Texture | 245 | 32 | 92 |
| Tile | 5 | Texture | 230 | 33 | 83 |
| Wood | 5 | Texture | 247 | 19 | 60 |
| Bottle | 3 | Object | 209 | 20 | 63 |
| Capsule | 5 | Object | 219 | 23 | 109 |
| Pill | 7 | Object | 267 | 26 | 141 |
| Transistor | 4 | Object | 213 | 60 | 40 |
| Zipper | 7 | Object | 240 | 32 | 119 |
| Cable | 8 | Object | 224 | 58 | 92 |
| Hazelnut | 4 | Object | 391 | 40 | 70 |
| Metal_nut | 4 | Object | 220 | 22 | 93 |
| Screw | 5 | Object | 320 | 41 | 119 |
| Toothbrush | 1 | Object | 60 | 12 | 30 |
| MVTec AD | 73 | - | 3,629 | 467 | 1,258 |
| AITEX | 12 | Texture | 1,692 | 564 | 183 |
| SDD | 1 | Texture | 594 | 286 | 54 |
| ELPV | 2 | Texture | 1,131 | 377 | 715 |
| Optical | 1 | Object | 10,500 | 3,500 | 2,100 |
| Mastcam | 11 | Object | 9,302 | 426 | 451 |
| BrainMRI | 1 | Medical | 73 | 25 | 155 |
| HeadCT | 1 | Medical | 75 | 25 | 100 |
| Hyper-Kvasir | 4 | Medical | 2,021 | 674 | 757 |

Table 1. Statistical details for nine real-world AD datasets, with the first 15 rows displaying detailed information for subsets of the MVTec AD dataset.

**MVTec AD** [1] is a widely-used dataset that enables researchers to benchmark the performance of anomaly detection methods in the context of industrial inspection applications. The dataset includes over 5,000 images that are divided into 15 object and texture categories. Each category contains a training set of anomaly-free images, as well as a test set that includes images with both defects and defect-free images.

**AITEX** [9] is a textile fabric database that comprises 245 images of 7 different fabrics, including 140 defect-free images (20 for each type of fabric) and 105 images with various types of defects.

**SDD** [10] is a collection of images captured in a controlled industrial environment, using defective production items as the subject. The dataset includes 52 images with visible defects and 347 product images without any defects.

**ELPV** [3] is a collection of 2,624 high-resolution grayscale images of solar cells extracted from photovoltaic modules. These images were extracted from 44 different solar modules, and include both intrinsic and extrinsic defects known to reduce the power efficiency of solar modules.

**Optical** [11] is a synthetic dataset created to simulate real-world industrial inspection tasks for defect detection. The dataset comprises ten individual subsets, with the first six subsets (referred to as development datasets) intended for algorithm development purposes. The remaining four subsets (known as competition datasets) can be used to evaluate algorithm performance.

**Mastcam** [5] is a novelty detection dataset constructed from geological images captured by a multispectral imaging system installed on Mars exploration rovers. The dataset comprises typical images as well as images of 11 novel geologic classes. Each image includes a shorter wavelength (color) channel and a longer wavelengths (grayscale) channel.

**BrainMRI** [8] is a dataset for brain tumor detection obtained from magnetic resonance imaging (MRI) of the brain.

**HeadCT** [8] is a dataset consisting of 100 normal head CT

*Corresponding author: G. Pang (gspang@smu.edu.sg)

slices and 100 slices with brain hemorrhage, without distinction between the types of hemorrhage. Each slice is from a different person, providing a diverse set of images for researchers to develop and test algorithms for hemorrhage detection and classification in medical imaging applications.

**Hyper-Kvasir** [2] is a large-scale open gastrointestinal dataset which is collected during real gastro- and colonoscopy procedures. It is comprised of four distinct parts, including labeled image data, unlabeled image data, segmented image data, and annotated video data.

## B. Implementation Details

### B.1. Generating Heterogeneous Anomaly Distribution Datasets

Our proposed approach creates a diverse collection of data subsets by randomly selecting a subset of normal clusters and incorporating labeled anomaly examples to form the support and query sets for learning heterogeneous anomaly distributions. Specifically, we generate each data subset as follows.

**Normal Samples in Each Data Subset.** We employ the K-means algorithm to cluster normal samples into three groups. In each data subset, we randomly select two clusters to form the support and query sets. We follow this way to generate six such data subsets. Furthermore, we include an additional subset consisting of all normal samples to capture a holistic view of the distribution of normal instances, in which normal samples are randomly divided into two parts to form the support and query sets.

**Abnormal Samples in Each Data Subset.** To simulate open-set environments in heterogeneous anomaly distribution generation component, each anomaly distribution $D_i$ is splited into two disjoint subsets, *i.e.*, $\mathcal{D}_i = \{\mathcal{D}_i^s, \mathcal{D}_i^q\}$, which correspond to support and query sets respectively, with the support set $\mathcal{D}_i^s = \mathcal{X}_{n,i}^s \cup \mathcal{X}_{a,i}^s$ used to train our base model $\phi_i$ and the query set $\mathcal{D}_i^q = \mathcal{X}_{n,i}^q \cup \mathcal{X}_{a,i}^q$ used to validate its open-set performance. For the setting with $M = 10$, we randomly select 50% samples from seen anomaly set and one normal cluster to form support set, while the remaining 50% seen anomalies are used in the query set. Under the protocol of having only one seen anomaly example, the example is included in both sets.

To enhance the variety of anomaly samples in our approach, we introduce three distinct anomaly augmentation techniques to generate pseudo anomaly samples: CutMix [12], CutPaste [6], DREAM Mask [13]. These techniques are randomly applied to each heterogeneous anomaly distribution subset to introduce diverse types of pseudo anomalies.

---

**Algorithm 1** Anomaly Heterogeneity Learning (`AHL`)

**Input:** Input $\mathcal{D} = \{\mathbf{x}, y\}, \{\phi\}_1^T, \psi$
**Output:** Output $g$
1: /* Heterogeneous Anomaly Distribution Generation */
2: Construct $\mathcal{D}_i$ through grouping training set $\mathcal{D}$ into $T$ groups
3: /* Collaborative Differentiable Learning */
4: **for** $epoch = 1$ to $N$ **do**
5:     Update parameter of base model $\phi_i$ for $\mathcal{D}_i$ based on Eq.(1)
6:     /* Learning Importance Scores of Individual Anomaly Distributions */
7:     **if** $epoch >= 5$ **then**
8:         Compute the generalization error $r_i$ and importance score $w_i$ for $\phi_i$ with the help of sequential model $\psi$ via Eq.(6) and Eq.(7), respectively
9:     **else**
10:         Treat all $\phi$ equally, $w_i = \frac{1}{T}$
11:     **end if**
12:     Update parameter of $g$ based on Eq.(4)
13:     Set the parameters of $\phi_i$ as the new weight parameters of $g$
14: **end for**

---

### B.2. The Algorithm of `AHL`

The overall objective of our `AHL` framework is to achieve a unified and robust AD model via synthesizing anomaly heterogeneities learned from various heterogeneous anomaly distributions. We summarize the Anomaly Heterogeneity Learning (`AHL`) procedure in Algorithm 1. Specifically, our framework first generates $T$ heterogeneous anomaly datasets, with each subset is sampled from the training set and contains a mixture of normal samples and (pseudo) anomaly samples, denoted as $\{\mathcal{D}_i\}_{i=1}^T$. In doing so, each subset is characterized by different set of normality/abnormality patterns, embodying heterogeneous anomaly distributions. We then employ a set of base models, denoted as $\{\phi_i\}_{i=1}^T$, to learn the underlying anomalous heterogeneity from these heterogeneous anomaly distributions. Moreover, a self-supervised sequential modeling approach is introduced to estimate the generalization errors $r_i$ and importance scores $w_i$ for each base model. Finally, we incorporate knowledge learned from heterogeneous anomaly distributions into a unified heterogeneous abnormality detection model $g$ to capture richer anomaly heterogeneity.

`AHL` is a generic framework, in which off-the-shelf open-set anomaly detectors can be easily plugged and gain significantly improved generalization and accuracy in detecting both seen and unseen anomalies. Once we choose a base anomaly detector, the training strategy and objective

| Dataset | One Anomaly Examples (Random) | | | | | | | Ten Anomaly Examples (Random) | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | SAOE | MLEP | FLOS | DevNet | DRA | AHL(DevNet) | AHL(DRA) | SAOE | MLEP | FLOS | DevNet | DRA | AHL(DevNet) | AHL(DRA) |
| Carpet | 0.766±0.098 | 0.701±0.091 | 0.755±0.026 | 0.778±0.055 | **0.873±0.035** | 0.802±0.018 | **0.877±0.004** | 0.755±0.136 | 0.781±0.049 | 0.780±0.009 | 0.864±0.012 | **0.945±0.014** | 0.867±0.006 | **0.953±0.001** |
| Grid | 0.921±0.032 | 0.839±0.028 | 0.871±0.076 | 0.868±0.031 | **0.972±0.016** | 0.872±0.032 | **0.975±0.001** | 0.952±0.011 | 0.980±0.009 | 0.966±0.005 | 0.901±0.016 | **0.990±0.008** | 0.914±0.003 | **0.992±0.002** |
| Leather | 0.996±0.007 | 0.781±0.020 | 0.791±0.057 | 0.874±0.016 | **0.988±0.003** | 0.880±0.005 | **0.988±0.001** | 1.000±0.000 | 0.813±0.158 | 0.993±0.004 | 0.986±0.033 | **1.000±0.000** | 0.996±0.008 | **1.000±0.000** |
| Tile | 0.935±0.034 | 0.927±0.036 | 0.787±0.038 | 0.872±0.035 | **0.966±0.014** | 0.909±0.007 | **0.968±0.001** | 0.944±0.013 | **0.988±0.009** | 0.952±0.010 | **1.000±0.000** | **1.000±0.000** | **1.000±0.000** | **1.000±0.000** |
| Wood | 0.948±0.009 | 0.660±0.142 | 0.927±0.065 | 0.917±0.029 | **0.987±0.012** | 0.947±0.020 | **0.987±0.003** | 0.976±0.031 | **0.999±0.002** | **1.000±0.000** | **0.999±0.000** | 0.998±0.013 | **1.000±0.000** | 0.998±0.000 |
| Bottle | 0.989±0.019 | 0.927±0.090 | 0.975±0.023 | 0.986±0.012 | **1.000±0.000** | 0.994±0.004 | **1.000±0.000** | 0.998±0.003 | 0.981±0.004 | 0.995±0.002 | 0.996±0.005 | **1.000±0.000** | 0.998±0.001 | **1.000±0.000** |
| Capsule | 0.611±0.109 | 0.558±0.075 | 0.666±0.020 | 0.567±0.042 | **0.646±0.029** | 0.581±0.202 | **0.665±0.030** | 0.850±0.054 | 0.818±0.063 | 0.902±0.017 | 0.872±0.017 | **0.928±0.011** | 0.885±0.012 | **0.930±0.001** |
| Pill | 0.652±0.078 | 0.656±0.061 | 0.745±0.064 | 0.779±0.018 | **0.831±0.026** | 0.781±0.087 | **0.840±0.003** | 0.872±0.049 | 0.845±0.048 | **0.929±0.012** | 0.882±0.008 | **0.918±0.009** | 0.900±0.004 | **0.918±0.001** |
| Transistor | 0.680±0.182 | 0.695±0.124 | 0.709±0.041 | 0.732±0.075 | 0.727±0.105 | **0.737±0.098** | **0.796±0.002** | 0.860±0.053 | 0.927±0.043 | 0.862±0.037 | 0.907±0.004 | **0.919±0.003** | 0.912±0.002 | **0.926±0.009** |
| Zipper | 0.970±0.033 | 0.856±0.086 | 0.885±0.033 | 0.914±0.027 | **0.983±0.008** | 0.928±0.006 | **0.986±0.000** | 0.995±0.004 | 0.965±0.002 | 0.990±0.008 | 0.992±0.008 | **1.000±0.000** | 0.995±0.002 | **1.000±0.000** |
| Cable | 0.819±0.060 | 0.688±0.017 | 0.790±0.039 | 0.790±0.086 | **0.855±0.007** | 0.793±0.091 | **0.858±0.011** | 0.862±0.022 | 0.857±0.062 | 0.890±0.063 | 0.901±0.006 | **0.914±0.006** | 0.907±0.004 | **0.921±0.001** |
| Hazelnut | 0.961±0.042 | 0.704±0.090 | 0.976±0.021 | 0.970±0.005 | **0.982±0.005** | 0.978±0.003 | **0.989±0.004** | **1.000±0.000** | **1.000±0.000** | **1.000±0.000** | **1.000±0.000** | **1.000±0.000** | **1.000±0.000** | **1.000±0.000** |
| Metal_nut | 0.922±0.033 | 0.878±0.038 | 0.930±0.022 | 0.874±0.014 | **0.950±0.011** | 0.880±0.002 | **0.952±0.003** | 0.976±0.013 | 0.974±0.009 | 0.984±0.004 | 0.992±0.003 | **0.997±0.005** | 0.997±0.003 | **0.998±0.000** |
| Screw | 0.653±0.074 | 0.675±0.294 | 0.337±0.091 | 0.766±0.045 | **0.902±0.032** | 0.769±0.031 | **0.927±0.009** | 0.975±0.023 | 0.899±0.039 | 0.940±0.017 | 0.981±0.013 | 0.978±0.012 | **0.984±0.005** | **0.985±0.002** |
| Toothbrush | 0.686±0.110 | 0.617±0.058 | 0.731±0.028 | **0.790±0.029** | 0.675±0.019 | **0.794±0.016** | 0.710±0.007 | 0.865±0.062 | 0.783±0.048 | 0.900±0.008 | **0.950±0.025** | 0.908±0.007 | **0.959±0.002** | 0.921±0.007 |
| MVTec AD | 0.834±0.007 | 0.744±0.019 | 0.792±0.014 | 0.832±0.016 | **0.889±0.013** | 0.843±0.021 | **0.901±0.003** | 0.926±0.010 | 0.907±0.005 | 0.939±0.007 | 0.948±0.005 | **0.966±0.002** | 0.954±0.003 | **0.970±0.002** |
| AITEX | 0.675±0.094 | 0.564±0.055 | 0.538±0.073 | 0.609±0.054 | 0.693±0.031 | **0.704±0.004** | **0.734±0.007** | 0.874±0.024 | 0.867±0.037 | 0.841±0.049 | 0.889±0.007 | 0.892±0.007 | **0.903±0.011** | **0.925±0.013** |
| SDD | 0.781±0.009 | 0.811±0.045 | 0.840±0.043 | 0.851±0.003 | **0.907±0.002** | 0.864±0.001 | **0.909±0.001** | 0.955±0.020 | 0.783±0.013 | 0.967±0.018 | 0.985±0.004 | **0.990±0.000** | **0.991±0.001** | **0.991±0.000** |
| ELPV | 0.635±0.092 | 0.578±0.062 | 0.457±0.056 | **0.810±0.024** | 0.676±0.003 | **0.828±0.005** | 0.723±0.008 | 0.793±0.047 | 0.794±0.047 | 0.818±0.032 | 0.843±0.001 | 0.843±0.002 | **0.849±0.003** | **0.850±0.004** |
| Optical | 0.815±0.014 | 0.516±0.009 | 0.518±0.003 | 0.513±0.001 | **0.880±0.002** | 0.547±0.009 | **0.888±0.007** | 0.941±0.013 | 0.740±0.039 | 0.720±0.055 | 0.785±0.012 | **0.966±0.002** | 0.841±0.010 | **0.976±0.004** |
| Mastcam | 0.662±0.018 | 0.625±0.045 | 0.542±0.017 | 0.627±0.049 | **0.709±0.011** | 0.644±0.013 | **0.743±0.003** | 0.810±0.029 | 0.798±0.026 | 0.703±0.029 | 0.797±0.021 | **0.849±0.003** | 0.825±0.020 | **0.855±0.005** |
| BrainMRI | 0.531±0.060 | 0.632±0.017 | 0.693±0.036 | **0.853±0.045** | 0.747±0.001 | **0.866±0.004** | 0.760±0.013 | 0.900±0.041 | 0.959±0.011 | 0.955±0.011 | 0.951±0.007 | **0.971±0.001** | 0.959±0.008 | **0.977±0.001** |
| HeadCT | 0.597±0.022 | 0.758±0.038 | 0.698±0.092 | 0.755±0.029 | **0.804±0.010** | 0.781±0.007 | **0.825±0.014** | 0.935±0.021 | 0.972±0.014 | 0.971±0.004 | **0.997±0.002** | 0.988±0.001 | **0.999±0.003** | 0.993±0.002 |
| Hyper-Kvasir | 0.498±0.100 | 0.445±0.040 | 0.668±0.004 | 0.734±0.020 | 0.712±0.010 | **0.768±0.015** | 0.742±0.015 | 0.666±0.050 | 0.600±0.069 | 0.773±0.029 | 0.822±0.031 | 0.844±0.001 | **0.873±0.009** | **0.880±0.003** |

Table 2. AUC results(mean±std) on nine real-world AD datasets under the general setting. Best results and the second-best results are respectively highlighted in **Red** and **Bold**.

function should be consistent with it. Following the proposed loss of base models (*i.e.*, DRA and DevNet), we adopt the deviation loss [7] to evaluate the loss between predicted anomaly scores and ground truths in the whole training phase:

$$\ell_{dev}(\mathbf{x}, y; h) = \mathbb{I}(y = 0)|(dev(\mathbf{x}; h))| $$
$$+ \mathbb{I}(y = 1) \max(0, m - dev(\mathbf{x}; h)),$$

where $\mathbb{I}(.)$ is an indicator function that is equal to one when the condition is true, and zero otherwise; $h(\cdot)$ denotes the anomaly detection model. $dev(\mathbf{x}) = \frac{h(\mathbf{x}) - \mu_r}{\sigma_r}$ with $\mu_r$ and $\sigma_r$ representing the mean and standard deviation of a set of sampled anomaly scores from the Gaussian prior distribution $\mathcal{N}(0, 1)$. $m$ is a confidence margin which defines a radius around the deviation.

## C. Detailed Empirical Results

### C.1. Full Results under General Setting

Table 2 shows the detailed comparison results of AHL and SOTA competing methods under the general setting. It includes the performance metrics of each category of MVTec AD dataset. Overall, our proposed AHL model consistently outperforms the baseline methods in both ten-shot and one-shot settings across all three application scenarios. AHL (DRA) achieves the best performance in terms of AUC. On average, AHL improves the AUC of DRA and DevNet by up to 4% and 9%, respectively.

### C.2. Full Results under Hard Setting

To investigate the detection performance of AHL framework on novel anomaly classes, we evaluate its performance under the hard setting, and present the detailed results for six multi-subset datasets, including each anomaly class-level

performance, in Table 3. Overall, our models – AHL (DRA) and AHL (DevNet) – achieve the best AUC results in both $M = 1$ and $M = 10$ setting protocols. Specifically, AHL improves the performances of DRA and DevNet by up to 3.2% and 3% AUC, respectively. The results here are consistent with the superiority performance of AHL in the general setting.

### C.3. More Ablation Study Results

**Class-level Results under Hard Setting.** To evaluate the effectiveness of each module in our AHL approach (**+ HADG + CDL$^+$**), we compare it with the base model (**DRA**), the base model using randomly sampled data distribution subsets and initial CDL component (**+ CDL**), and the base model with HADG and initial CDL component (**+ HADG + CDL**). Table 4 shows the detailed results of class-level anomaly detection in the hard setting. The results show that all the modules in the AHL framework contribute to improving the detection performance on unseen anomaly classes, demonstrating the importance of anomaly heterogeneity.

**Importance of Pseudo Anomalies.** Moreover, since DevNet [7] does not use pseudo anomalies, we also evaluate the impact of removing the pseudo anomalies on the AHL framework when using DevNet as the base model. As shown in Figure 1, AHL (DevNet) remains substantially better than the original DevNet model in such cases.

**Effectiveness of Data Overlapping in HADG Module.** Table 5 shows the results of AHL (DRA) using three data overlapping types: overlap occurring exclusively in (1) the query set ($\mathcal{D}^s$->$\mathcal{D}^q$), (2) the support set ($\mathcal{D}^q$->$\mathcal{D}^s$), and (3) in both sets ($\mathcal{D}^s <> \mathcal{D}^q$). It is clear that the overlapping between $\mathcal{D}^q$ and $\mathcal{D}^s$ typically does not lead to performance improvement.

Table 3 columns header:

| Dataset | | One Example from One Anomaly Class | | | | | | | Ten Example from One Anomaly Class | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | SAOE | MLEP | FLOS | DevNet | DRA | AHL(DevNet) | AHL(DRA) | SAOE | MLEP | FLOS | DevNet | DRA | AHL(DevNet) | AHL(DRA) |
| Carpet | Color | $0.763_{\pm0.100}$ | $0.547_{\pm0.056}$ | $0.467_{\pm0.278}$ | $0.701_{\pm0.046}$ | $0.890_{\pm0.011}$ | $0.718_{\pm0.009}$ | $0.894_{\pm0.004}$ | $0.467_{\pm0.067}$ | $0.698_{\pm0.025}$ | $0.760_{\pm0.005}$ | $0.774_{\pm0.009}$ | $0.899_{\pm0.019}$ | $0.778_{\pm0.004}$ | $0.929_{\pm0.007}$ |
| | Cut | $0.664_{\pm0.165}$ | $0.658_{\pm0.056}$ | $0.685_{\pm0.007}$ | $0.679_{\pm0.018}$ | $0.890_{\pm0.024}$ | $0.684_{\pm0.014}$ | $0.934_{\pm0.003}$ | $0.793_{\pm0.175}$ | $0.653_{\pm0.120}$ | $0.688_{\pm0.059}$ | $0.817_{\pm0.021}$ | $0.942_{\pm0.012}$ | $0.825_{\pm0.006}$ | $0.943_{\pm0.002}$ |
| | Hole | $0.772_{\pm0.071}$ | $0.653_{\pm0.065}$ | $0.594_{\pm0.041}$ | $0.729_{\pm0.032}$ | $0.915_{\pm0.045}$ | $0.736_{\pm0.062}$ | $0.935_{\pm0.014}$ | $0.831_{\pm0.125}$ | $0.674_{\pm0.076}$ | $0.733_{\pm0.014}$ | $0.808_{\pm0.016}$ | $0.958_{\pm0.031}$ | $0.815_{\pm0.036}$ | $0.960_{\pm0.003}$ |
| | Metal | $0.780_{\pm0.172}$ | $0.706_{\pm0.047}$ | $0.701_{\pm0.028}$ | $0.822_{\pm0.016}$ | $0.877_{\pm0.013}$ | $0.846_{\pm0.012}$ | $0.931_{\pm0.007}$ | $0.883_{\pm0.043}$ | $0.764_{\pm0.061}$ | $0.678_{\pm0.083}$ | $0.885_{\pm0.012}$ | $0.916_{\pm0.017}$ | $0.899_{\pm0.026}$ | $0.921_{\pm0.003}$ |
| | Thread | $0.787_{\pm0.204}$ | $0.831_{\pm0.117}$ | $0.941_{\pm0.005}$ | $0.937_{\pm0.017}$ | $0.954_{\pm0.010}$ | $0.941_{\pm0.006}$ | $0.966_{\pm0.005}$ | $0.834_{\pm0.297}$ | $0.967_{\pm0.006}$ | $0.946_{\pm0.005}$ | $0.981_{\pm0.005}$ | $0.985_{\pm0.005}$ | $0.984_{\pm0.013}$ | $0.991_{\pm0.001}$ |
| | Mean | $0.753_{\pm0.029}$ | $0.679_{\pm0.029}$ | $0.678_{\pm0.040}$ | $0.774_{\pm0.007}$ | $0.905_{\pm0.006}$ | $0.785_{\pm0.015}$ | $0.932_{\pm0.003}$ | $0.762_{\pm0.073}$ | $0.751_{\pm0.023}$ | $0.761_{\pm0.001}$ | $0.853_{\pm0.005}$ | $0.940_{\pm0.006}$ | $0.860_{\pm0.013}$ | $0.949_{\pm0.002}$ |
| Metal_nut | Bent | $0.864_{\pm0.032}$ | $0.743_{\pm0.013}$ | $0.851_{\pm0.046}$ | $0.817_{\pm0.033}$ | $0.952_{\pm0.015}$ | $0.831_{\pm0.020}$ | $0.954_{\pm0.003}$ | $0.901_{\pm0.023}$ | $0.956_{\pm0.013}$ | $0.827_{\pm0.075}$ | $0.907_{\pm0.018}$ | $0.987_{\pm0.003}$ | $0.909_{\pm0.016}$ | $0.989_{\pm0.000}$ |
| | Color | $0.857_{\pm0.037}$ | $0.835_{\pm0.075}$ | $0.821_{\pm0.059}$ | $0.903_{\pm0.019}$ | $0.930_{\pm0.021}$ | $0.910_{\pm0.008}$ | $0.933_{\pm0.008}$ | $0.879_{\pm0.018}$ | $0.945_{\pm0.039}$ | $0.978_{\pm0.008}$ | $0.992_{\pm0.015}$ | $0.956_{\pm0.009}$ | $0.995_{\pm0.002}$ | $0.958_{\pm0.000}$ |
| | Flip | $0.751_{\pm0.090}$ | $0.813_{\pm0.031}$ | $0.799_{\pm0.058}$ | $0.751_{\pm0.039}$ | $0.931_{\pm0.017}$ | $0.755_{\pm0.022}$ | $0.931_{\pm0.001}$ | $0.795_{\pm0.062}$ | $0.805_{\pm0.057}$ | $0.942_{\pm0.009}$ | $0.982_{\pm0.010}$ | $0.931_{\pm0.010}$ | $0.987_{\pm0.003}$ | $0.937_{\pm0.003}$ |
| | Scratch | $0.792_{\pm0.075}$ | $0.907_{\pm0.085}$ | $0.947_{\pm0.027}$ | $0.974_{\pm0.061}$ | $0.929_{\pm0.009}$ | $0.981_{\pm0.035}$ | $0.934_{\pm0.005}$ | $0.845_{\pm0.041}$ | $0.805_{\pm0.153}$ | $0.943_{\pm0.002}$ | $0.998_{\pm0.005}$ | $0.998_{\pm0.006}$ | $0.998_{\pm0.001}$ | $0.999_{\pm0.00}$ |
| | Mean | $0.816_{\pm0.020}$ | $0.825_{\pm0.023}$ | $0.855_{\pm0.024}$ | $0.861_{\pm0.019}$ | $0.936_{\pm0.041}$ | $0.869_{\pm0.004}$ | $0.939_{\pm0.004}$ | $0.855_{\pm0.016}$ | $0.878_{\pm0.058}$ | $0.922_{\pm0.014}$ | $0.970_{\pm0.009}$ | $0.968_{\pm0.006}$ | $0.972_{\pm0.002}$ | $0.971_{\pm0.001}$ |
| AITEX | Broken end | $0.778_{\pm0.068}$ | $0.441_{\pm0.111}$ | $0.645_{\pm0.030}$ | $0.702_{\pm0.037}$ | $0.696_{\pm0.057}$ | $0.716_{\pm0.014}$ | $0.704_{\pm0.005}$ | $0.712_{\pm0.068}$ | $0.732_{\pm0.065}$ | $0.585_{\pm0.037}$ | $0.658_{\pm0.062}$ | $0.708_{\pm0.062}$ | $0.688_{\pm0.013}$ | $0.735_{\pm0.010}$ |
| | Broken pick | $0.644_{\pm0.039}$ | $0.476_{\pm0.070}$ | $0.598_{\pm0.023}$ | $0.567_{\pm0.016}$ | $0.719_{\pm0.004}$ | $0.575_{\pm0.005}$ | $0.727_{\pm0.003}$ | $0.629_{\pm0.012}$ | $0.555_{\pm0.027}$ | $0.548_{\pm0.054}$ | $0.595_{\pm0.017}$ | $0.671_{\pm0.034}$ | $0.612_{\pm0.005}$ | $0.683_{\pm0.002}$ |
| | Cut selvage | $0.681_{\pm0.077}$ | $0.434_{\pm0.149}$ | $0.694_{\pm0.036}$ | $0.674_{\pm0.021}$ | $0.751_{\pm0.006}$ | $0.680_{\pm0.017}$ | $0.753_{\pm0.007}$ | $0.770_{\pm0.014}$ | $0.682_{\pm0.025}$ | $0.745_{\pm0.035}$ | $0.703_{\pm0.062}$ | $0.777_{\pm0.021}$ | $0.737_{\pm0.012}$ | $0.781_{\pm0.006}$ |
| | Fuzzyball | $0.650_{\pm0.064}$ | $0.525_{\pm0.157}$ | $0.525_{\pm0.043}$ | $0.629_{\pm0.103}$ | $0.631_{\pm0.018}$ | $0.644_{\pm0.031}$ | $0.647_{\pm0.040}$ | $0.842_{\pm0.026}$ | $0.677_{\pm0.023}$ | $0.550_{\pm0.082}$ | $0.736_{\pm0.101}$ | $0.749_{\pm0.033}$ | $0.755_{\pm0.002}$ | $0.775_{\pm0.024}$ |
| | Nep | $0.710_{\pm0.044}$ | $0.517_{\pm0.059}$ | $0.734_{\pm0.038}$ | $0.741_{\pm0.011}$ | $0.685_{\pm0.010}$ | $0.754_{\pm0.012}$ | $0.703_{\pm0.005}$ | $0.771_{\pm0.032}$ | $0.740_{\pm0.052}$ | $0.746_{\pm0.060}$ | $0.806_{\pm0.039}$ | $0.784_{\pm0.025}$ | $0.836_{\pm0.007}$ | $0.792_{\pm0.007}$ |
| | Weft crack | $0.582_{\pm0.108}$ | $0.400_{\pm0.029}$ | $0.546_{\pm0.114}$ | $0.561_{\pm0.085}$ | $0.693_{\pm0.002}$ | $0.588_{\pm0.018}$ | $0.706_{\pm0.009}$ | $0.618_{\pm0.172}$ | $0.370_{\pm0.037}$ | $0.636_{\pm0.051}$ | $0.614_{\pm0.097}$ | $0.710_{\pm0.016}$ | $0.624_{\pm0.005}$ | $0.713_{\pm0.003}$ |
| | Mean | $0.674_{\pm0.034}$ | $0.466_{\pm0.030}$ | $0.624_{\pm0.024}$ | $0.646_{\pm0.014}$ | $0.696_{\pm0.011}$ | $0.660_{\pm0.007}$ | $0.707_{\pm0.007}$ | $0.724_{\pm0.032}$ | $0.626_{\pm0.041}$ | $0.635_{\pm0.043}$ | $0.685_{\pm0.016}$ | $0.733_{\pm0.011}$ | $0.709_{\pm0.006}$ | $0.747_{\pm0.002}$ |
| ELPV | Mono | $0.563_{\pm0.102}$ | $0.649_{\pm0.027}$ | $0.717_{\pm0.025}$ | $0.620_{\pm0.057}$ | $0.762_{\pm0.017}$ | $0.638_{\pm0.019}$ | $0.774_{\pm0.013}$ | $0.569_{\pm0.035}$ | $0.756_{\pm0.045}$ | $0.629_{\pm0.072}$ | $0.639_{\pm0.067}$ | $0.735_{\pm0.008}$ | $0.663_{\pm0.007}$ | $0.745_{\pm0.004}$ |
| | Poly | $0.665_{\pm0.173}$ | $0.483_{\pm0.247}$ | $0.665_{\pm0.021}$ | $0.705_{\pm0.011}$ | $0.681_{\pm0.026}$ | $0.717_{\pm0.007}$ | $0.705_{\pm0.006}$ | $0.796_{\pm0.084}$ | $0.734_{\pm0.078}$ | $0.662_{\pm0.042}$ | $0.806_{\pm0.004}$ | $0.806_{\pm0.004}$ | $0.842_{\pm0.003}$ | $0.831_{\pm0.011}$ |
| | Mean | $0.614_{\pm0.048}$ | $0.566_{\pm0.111}$ | $0.691_{\pm0.008}$ | $0.663_{\pm0.008}$ | $0.722_{\pm0.006}$ | $0.678_{\pm0.006}$ | $0.740_{\pm0.003}$ | $0.683_{\pm0.047}$ | $0.745_{\pm0.020}$ | $0.646_{\pm0.032}$ | $0.722_{\pm0.018}$ | $0.771_{\pm0.005}$ | $0.752_{\pm0.005}$ | $0.788_{\pm0.003}$ |
| Mastcam | Bedrock | $0.636_{\pm0.072}$ | $0.532_{\pm0.036}$ | $0.499_{\pm0.056}$ | $0.508_{\pm0.107}$ | $0.653_{\pm0.019}$ | $0.533_{\pm0.065}$ | $0.679_{\pm0.012}$ | $0.636_{\pm0.068}$ | $0.512_{\pm0.062}$ | $0.499_{\pm0.098}$ | $0.586_{\pm0.012}$ | $0.654_{\pm0.013}$ | $0.589_{\pm0.010}$ | $0.673_{\pm0.006}$ |
| | Broken-rock | $0.699_{\pm0.058}$ | $0.544_{\pm0.088}$ | $0.569_{\pm0.025}$ | $0.558_{\pm0.016}$ | $0.640_{\pm0.023}$ | $0.572_{\pm0.014}$ | $0.661_{\pm0.009}$ | $0.712_{\pm0.062}$ | $0.651_{\pm0.063}$ | $0.608_{\pm0.085}$ | $0.562_{\pm0.033}$ | $0.704_{\pm0.007}$ | $0.572_{\pm0.024}$ | $0.722_{\pm0.004}$ |
| | Drill-hole | $0.697_{\pm0.074}$ | $0.636_{\pm0.066}$ | $0.539_{\pm0.077}$ | $0.555_{\pm0.026}$ | $0.642_{\pm0.035}$ | $0.563_{\pm0.012}$ | $0.654_{\pm0.004}$ | $0.682_{\pm0.042}$ | $0.660_{\pm0.002}$ | $0.601_{\pm0.009}$ | $0.590_{\pm0.074}$ | $0.757_{\pm0.008}$ | $0.610_{\pm0.075}$ | $0.760_{\pm0.003}$ |
| | Drt | $0.735_{\pm0.020}$ | $0.624_{\pm0.042}$ | $0.591_{\pm0.042}$ | $0.570_{\pm0.048}$ | $0.733_{\pm0.027}$ | $0.581_{\pm0.023}$ | $0.724_{\pm0.006}$ | $0.761_{\pm0.062}$ | $0.616_{\pm0.048}$ | $0.652_{\pm0.024}$ | $0.620_{\pm0.031}$ | $0.757_{\pm0.006}$ | $0.629_{\pm0.016}$ | $0.772_{\pm0.004}$ |
| | Dump-pile | $0.682_{\pm0.022}$ | $0.545_{\pm0.127}$ | $0.508_{\pm0.021}$ | $0.510_{\pm0.008}$ | $0.741_{\pm0.022}$ | $0.519_{\pm0.004}$ | $0.756_{\pm0.011}$ | $0.750_{\pm0.037}$ | $0.696_{\pm0.047}$ | $0.700_{\pm0.070}$ | $0.689_{\pm0.070}$ | $0.757_{\pm0.008}$ | $0.695_{\pm0.021}$ | $0.802_{\pm0.005}$ |
| | Float | $0.711_{\pm0.041}$ | $0.530_{\pm0.075}$ | $0.551_{\pm0.030}$ | $0.507_{\pm0.039}$ | $0.688_{\pm0.031}$ | $0.524_{\pm0.017}$ | $0.702_{\pm0.005}$ | $0.718_{\pm0.046}$ | $0.671_{\pm0.032}$ | $0.736_{\pm0.041}$ | $0.640_{\pm0.012}$ | $0.749_{\pm0.009}$ | $0.647_{\pm0.008}$ | $0.765_{\pm0.002}$ |
| | Meteorite | $0.669_{\pm0.037}$ | $0.476_{\pm0.014}$ | $0.462_{\pm0.017}$ | $0.436_{\pm0.033}$ | $0.604_{\pm0.020}$ | $0.463_{\pm0.008}$ | $0.616_{\pm0.013}$ | $0.647_{\pm0.030}$ | $0.473_{\pm0.047}$ | $0.568_{\pm0.053}$ | $0.561_{\pm0.053}$ | $0.689_{\pm0.010}$ | $0.572_{\pm0.015}$ | $0.691_{\pm0.001}$ |
| | Scuff | $0.679_{\pm0.048}$ | $0.492_{\pm0.037}$ | $0.508_{\pm0.070}$ | $0.496_{\pm0.121}$ | $0.573_{\pm0.017}$ | $0.515_{\pm0.013}$ | $0.581_{\pm0.020}$ | $0.676_{\pm0.019}$ | $0.504_{\pm0.052}$ | $0.575_{\pm0.042}$ | $0.447_{\pm0.043}$ | $0.626_{\pm0.016}$ | $0.506_{\pm0.104}$ | $0.656_{\pm0.009}$ |
| | Veins | $0.688_{\pm0.068}$ | $0.489_{\pm0.028}$ | $0.493_{\pm0.052}$ | $0.530_{\pm0.006}$ | $0.650_{\pm0.012}$ | $0.548_{\pm0.010}$ | $0.687_{\pm0.017}$ | $0.686_{\pm0.053}$ | $0.510_{\pm0.090}$ | $0.608_{\pm0.044}$ | $0.577_{\pm0.029}$ | $0.644_{\pm0.007}$ | $0.598_{\pm0.037}$ | $0.650_{\pm0.003}$ |
| | Mean | $0.689_{\pm0.037}$ | $0.541_{\pm0.007}$ | $0.524_{\pm0.013}$ | $0.519_{\pm0.014}$ | $0.658_{\pm0.021}$ | $0.535_{\pm0.003}$ | $0.673_{\pm0.010}$ | $0.697_{\pm0.014}$ | $0.588_{\pm0.016}$ | $0.616_{\pm0.021}$ | $0.588_{\pm0.025}$ | $0.704_{\pm0.007}$ | $0.602_{\pm0.008}$ | $0.721_{\pm0.003}$ |
| Hyper-Kvasir | Barretts | $0.382_{\pm0.117}$ | $0.438_{\pm0.111}$ | $0.703_{\pm0.040}$ | $0.682_{\pm0.007}$ | $0.788_{\pm0.008}$ | $0.701_{\pm0.013}$ | $0.792_{\pm0.007}$ | $0.698_{\pm0.037}$ | $0.540_{\pm0.014}$ | $0.764_{\pm0.066}$ | $0.837_{\pm0.014}$ | $0.820_{\pm0.005}$ | $0.850_{\pm0.002}$ | $0.829_{\pm0.002}$ |
| | Barretts-short-seg | $0.367_{\pm0.050}$ | $0.532_{\pm0.075}$ | $0.538_{\pm0.033}$ | $0.608_{\pm0.077}$ | $0.643_{\pm0.013}$ | $0.629_{\pm0.027}$ | $0.651_{\pm0.006}$ | $0.661_{\pm0.034}$ | $0.480_{\pm0.107}$ | $0.810_{\pm0.034}$ | $0.790_{\pm0.017}$ | $0.829_{\pm0.006}$ | $0.812_{\pm0.005}$ | $0.895_{\pm0.003}$ |
| | Esophagitis-a | $0.518_{\pm0.063}$ | $0.491_{\pm0.084}$ | $0.536_{\pm0.040}$ | $0.567_{\pm0.034}$ | $0.759_{\pm0.015}$ | $0.583_{\pm0.015}$ | $0.760_{\pm0.006}$ | $0.820_{\pm0.034}$ | $0.646_{\pm0.036}$ | $0.815_{\pm0.022}$ | $0.867_{\pm0.004}$ | $0.854_{\pm0.006}$ | $0.876_{\pm0.002}$ | $0.878_{\pm0.021}$ |
| | Esophagitis-b-d | $0.358_{\pm0.030}$ | $0.457_{\pm0.086}$ | $0.505_{\pm0.039}$ | $0.535_{\pm0.025}$ | $0.604_{\pm0.022}$ | $0.562_{\pm0.010}$ | $0.622_{\pm0.014}$ | $0.611_{\pm0.017}$ | $0.621_{\pm0.042}$ | $0.754_{\pm0.045}$ | $0.812_{\pm0.025}$ | $0.785_{\pm0.006}$ | $0.842_{\pm0.010}$ | $0.815_{\pm0.010}$ |
| | Mean | $0.406_{\pm0.018}$ | $0.480_{\pm0.044}$ | $0.571_{\pm0.004}$ | $0.598_{\pm0.006}$ | $0.699_{\pm0.009}$ | $0.619_{\pm0.005}$ | $0.706_{\pm0.007}$ | $0.698_{\pm0.021}$ | $0.571_{\pm0.014}$ | $0.786_{\pm0.024}$ | $0.827_{\pm0.008}$ | $0.822_{\pm0.013}$ | $0.845_{\pm0.003}$ | $0.854_{\pm0.004}$ |

Table 3. AUC results(mean±std) on nine real-world AD datasets under the hard setting. Best results and the second-best results are respectively highlighted in **Red** and **Bold**. Carpet and Meta_nut are two subsets of MVTec AD. The same set of datasets is used as in [4].

| Datsset | Subset | DRA | + CDL | + HADG + CDL | + HADG+ CDL$^+$ |
|---|---|---|---|---|---|
| Carpet | Color | $0.899_{\pm0.019}$ | $0.917_{\pm0.004}$ | $0.919_{\pm0.003}$ | $0.929_{\pm0.007}$ |
| | Cut | $0.942_{\pm0.012}$ | $0.938_{\pm0.004}$ | $0.943_{\pm0.001}$ | $0.943_{\pm0.002}$ |
| | Hole | $0.958_{\pm0.031}$ | $0.954_{\pm0.010}$ | $0.952_{\pm0.002}$ | $0.960_{\pm0.003}$ |
| | Metal | $0.916_{\pm0.017}$ | $0.919_{\pm0.008}$ | $0.914_{\pm0.006}$ | $0.921_{\pm0.003}$ |
| | Thread | $0.985_{\pm0.005}$ | $0.985_{\pm0.003}$ | $0.988_{\pm0.002}$ | $0.991_{\pm0.001}$ |
| | Mean | $0.940_{\pm0.006}$ | $0.943_{\pm0.002}$ | $0.943_{\pm0.003}$ | $0.949_{\pm0.002}$ |
| AITEX | Broken end | $0.708_{\pm0.062}$ | $0.714_{\pm0.011}$ | $0.719_{\pm0.005}$ | $0.735_{\pm0.010}$ |
| | Broken pick | $0.671_{\pm0.034}$ | $0.670_{\pm0.005}$ | $0.678_{\pm0.002}$ | $0.683_{\pm0.002}$ |
| | Cut selvage | $0.777_{\pm0.021}$ | $0.777_{\pm0.009}$ | $0.779_{\pm0.018}$ | $0.781_{\pm0.006}$ |
| | Fuzzyball | $0.749_{\pm0.033}$ | $0.742_{\pm0.010}$ | $0.756_{\pm0.021}$ | $0.775_{\pm0.024}$ |
| | Nep | $0.784_{\pm0.025}$ | $0.786_{\pm0.007}$ | $0.788_{\pm0.005}$ | $0.792_{\pm0.007}$ |
| | Weft crack | $0.710_{\pm0.016}$ | $0.708_{\pm0.003}$ | $0.711_{\pm0.003}$ | $0.713_{\pm0.003}$ |
| | Mean | $0.733_{\pm0.011}$ | $0.733_{\pm0.005}$ | $0.739_{\pm0.007}$ | $0.747_{\pm0.002}$ |
| elpv | Mono | $0.735_{\pm0.008}$ | $0.738_{\pm0.008}$ | $0.739_{\pm0.006}$ | $0.745_{\pm0.004}$ |
| | Poly | $0.806_{\pm0.004}$ | $0.809_{\pm0.006}$ | $0.817_{\pm0.012}$ | $0.831_{\pm0.011}$ |
| | Mean | $0.771_{\pm0.005}$ | $0.774_{\pm0.004}$ | $0.784_{\pm0.004}$ | $0.788_{\pm0.003}$ |
| Hyper-Kvasir | Barretts | $0.820_{\pm0.005}$ | $0.819_{\pm0.006}$ | $0.822_{\pm0.009}$ | $0.829_{\pm0.002}$ |
| | Barretts-short-seg | $0.829_{\pm0.006}$ | $0.864_{\pm0.012}$ | $0.887_{\pm0.016}$ | $0.895_{\pm0.003}$ |
| | Esophagitis-a | $0.854_{\pm0.006}$ | $0.863_{\pm0.006}$ | $0.871_{\pm0.005}$ | $0.878_{\pm0.021}$ |
| | Esophagitis-b-d | $0.785_{\pm0.006}$ | $0.795_{\pm0.005}$ | $0.806_{\pm0.004}$ | $0.815_{\pm0.010}$ |
| | Mean | $0.822_{\pm0.013}$ | $0.835_{\pm0.004}$ | $0.847_{\pm0.008}$ | $0.854_{\pm0.004}$ |

Table 4. Ablation study class-level results of AHL and its three main variants under hard settings. Best results and the second-best results are respectively highlighted in **Red** and **Bold**.

| Dataset | $\mathcal{D}^s$->$\mathcal{D}^q$ | $\mathcal{D}^q$->$\mathcal{D}^s$ | $\mathcal{D}^s$ <> $\mathcal{D}^q$ | Ours |
|---|---|---|---|---|
| AITEX | 0.918 | 0.917 | 0.917 | **0.925** |
| SDD | 0.988 | 0.990 | **0.991** | 0.991 |
| ELPV | 0.853 | 0.844 | 0.842 | **0.850** |
| BrainMRI | 0.970 | 0.969 | 0.970 | **0.977** |
| HeadCT | 0.993 | 0.991 | 0.993 | **0.993** |
| Hyper-Kvasir | 0.872 | 0.875 | 0.876 | **0.880** |

Table 5. Results with various data overlapping in the support and query sets in HADG.



Figure 1. Comparison of the AUC performance between DevNet and AHL (DRA) without the pseudo anomaly augmentation generation module. Here *wo. aug* indicates the augmentation techniques are excluded during training.

# References

[1] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Mvtec ad–a comprehensive real-world dataset for unsupervised anomaly detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9592–9600, 2019. 1

[2] Hanna Borgli, Vajira Thambawita, Pia H Smedsrud, Steven Hicks, Debesh Jha, Sigrun L Eskeland, Kristin Ranheim Randel, Konstantin Pogorelov, Mathias Lux, Duc Tien Dang Nguyen, et al. Hyperkvasir, a comprehensive multi-class image and video dataset for gastrointestinal endoscopy. *Scientific data*, 7(1):283, 2020. 2

[3] Sergiu Deitsch, Vincent Christlein, Stephan Berger, Claudia Buerhop-Lutz, Andreas Maier, Florian Gallwitz, and Christian Riess. Automatic classification of defective photovoltaic module cells in electroluminescence images. *Solar Energy*, 185:455–468, 2019. 1

[4] Choubo Ding, Guansong Pang, and Chunhua Shen. Catching both gray and black swans: Open-set supervised anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7388–7398, 2022. 1, 4

[5] Hannah R Kerner, Kiri L Wagstaff, Brian D Bue, Danika F Wellington, Samantha Jacob, Paul Horton, James F Bell, Chiman Kwan, and Heni Ben Amor. Comparison of novelty detection methods for multispectral images in rover-based planetary exploration missions. *Data Mining and Knowledge Discovery*, 34:1642–1675, 2020. 1

[6] Chun-Liang Li, Kihyuk Sohn, Jinsung Yoon, and Tomas Pfister. Cutpaste: Self-supervised learning for anomaly detection and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9664–9674, 2021. 2

[7] Guansong Pang, Choubo Ding, Chunhua Shen, and Anton van den Hengel. Explainable deep few-shot anomaly detection with deviation networks. *arXiv preprint arXiv:2108.00462*, 2021. 1, 3

[8] Mohammadreza Salehi, Niousha Sadjadi, Soroosh Baselizadeh, Mohammad H Rohban, and Hamid R Rabiee. Multiresolution knowledge distillation for anomaly detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14902–14912, 2021. 1

[9] Javier Silvestre-Blanes, Teresa Albero-Albero, Ignacio Miralles, Rubén Pérez-Llorens, and Jorge Moreno. A public fabric database for defect detection methods and results. *Autex Research Journal*, 19(4):363–374, 2019. 1

[10] Domen Tabernik, Samo Šela, Jure Skvarč, and Danijel Skočaj. Segmentation-based deep-learning approach for surface-defect detection. *Journal of Intelligent Manufacturing*, 31(3):759–776, 2020. 1

[11] Matthias Wieler and Tobias Hahn. Weakly supervised learning for industrial optical inspection. In *DAGM symposium in*, 2007. 1

[12] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019. 2

[13] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. Draema discriminatively trained reconstruction embedding for surface anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8330–8339, 2021. 2