

Siamese U-Net for Image Anomaly Detection and Segmentation with Contrastive Learning

Chia-Ying Lin¹, Shang-Hong Lai^{1,2}

¹Institute of Information Systems and Applications, National Tsing Hua University, Hsinchu, Taiwan ²Department of Computer Science, National Tsing Hua University, Hsinchu, Taiwan

Motivation

Mainstream approaches for anomaly detection and segmentation task in smart manufacturing usually suffer from two main issues:

- High FPR/FNR: Models trained with normal data only in an unsupervised manner are more likely to produce high false positive/ negative rates as they have no access to true anomalies.
- Indiscriminative image-level anomaly score: Segmentation-dominated models are prone to produce indistinguishable imagelevel anomaly scores for subtle anomalies. Slight discrepancy between pixel anomaly scores for normal and anomalous features often results in indiscriminative image-level anomaly scores, leading to degraded anomaly detection performance.

To tackle the aforementioned issues, we propose a novel Siamese U-Net model trained with contrastive learning and deviation finetuning, incorporating a few anomalous samples, either real-world anomalies or synthetic samples, into model training.

Proposed Method

Model Architecture



(a) Overview of the T-S U-Net block embedded with channel-positional attention module (CPAM). Laverwise cosine similarity maps are aggregated as anomaly map. (b) CPAM comprises two attention submodules, applying attention from channel and spatial aspects.

Model Flow: Two-Stage Training

(a) Stage 1: Contrastive Learning for Siamese U-Net





(a) Stage 1: Contrastive Learning (b) Stage 2: Deviation- based Detection Finetuning

Experimental Results



Fig 1: Qualitative comparison between our method and the baseline Reverse Distillation on MVTecAD (left) and MVTec3D-AD (right). For challenging case "tire" on the top row of the right half, our model can still generate a reliable anomaly map due to the aid of the deviation finetuning. Without deviation finetuning, the baseline tends to falsely predict the background with a high anomaly score.

Tab 1: Comparison on MVTecAD with full-shot training setting. AUROC% is reported in the format of (image-level, pixel-level). Our results include two variants: trained with a few (10-shot) real anomalies and with synthetic data, denoted as Ours(r) and Ours(s), respectively.

Tab 2: Quantitative comparison with SOTA methods on MVTec3D-AD under 2D setup. Our results include trained with a few (10-shot) real-world anomalies and with synthetic data, denoted as Ours(r) and Ours(s), respectively.