

---

# 내부 논문 Review Study 발표 : DeSTSeg

---

2023.08.04

Youngjae Kim

# 1. About the Paper

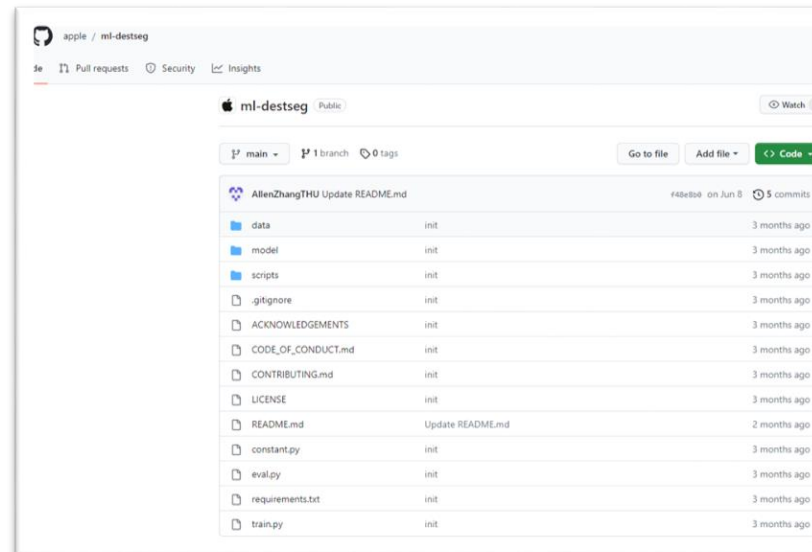
## DeSTSeg: Segmentation Guided Denoising Student-Teacher for Anomaly Detection

Xuan Zhang<sup>1</sup>, Shiyu Li<sup>2</sup>, Xi Li<sup>2</sup>, Ping Huang<sup>2</sup>, Jiulong Shan<sup>2</sup>, Ting Chen<sup>1</sup>  
<sup>1</sup>Tsinghua University <sup>2</sup>Apple

x-zhang18@mails.tsinghua.edu.cn, {shiyu.li, weston.li, huang.ping, j1shan}@apple.com,  
tingchen@tsinghua.edu.cn

- Accepted by CVPR 2023
- Tsinghua Univ, Apple

<https://github.com/apple/ml-destseg>



## 2. Introduction & Related Works

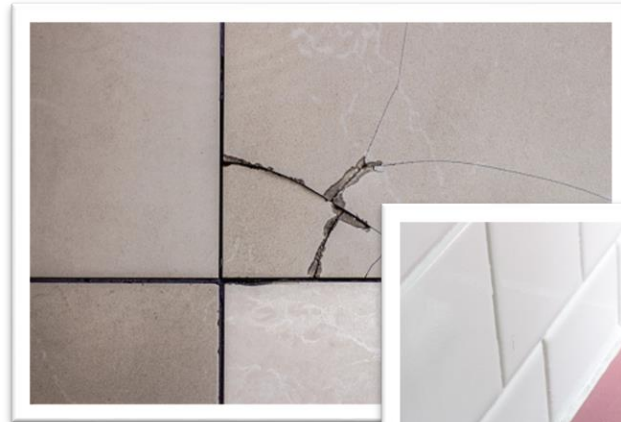
- Importance of Visual Anomaly Detection (AD)

- 1) Industrial inspection
- 2) Medical disease screening
- 3) Video surveillance (like CCTV)

- Properties of Anomalous samples

- 1) Types are enormous
- 2) Occurs rarely
- 3) Impossible to get all possible cases

-> usually, use only "normal" samples to train



## 2. Introduction & Related Works

- **Student – Teacher Framework (knowledge distillation)**
  - Proven to be effective in AD



- Pretrained on large-scale dataset (ex. ResNet50 on ImageNet)
- Large Models

**Pretrained  
Teacher  
Network**

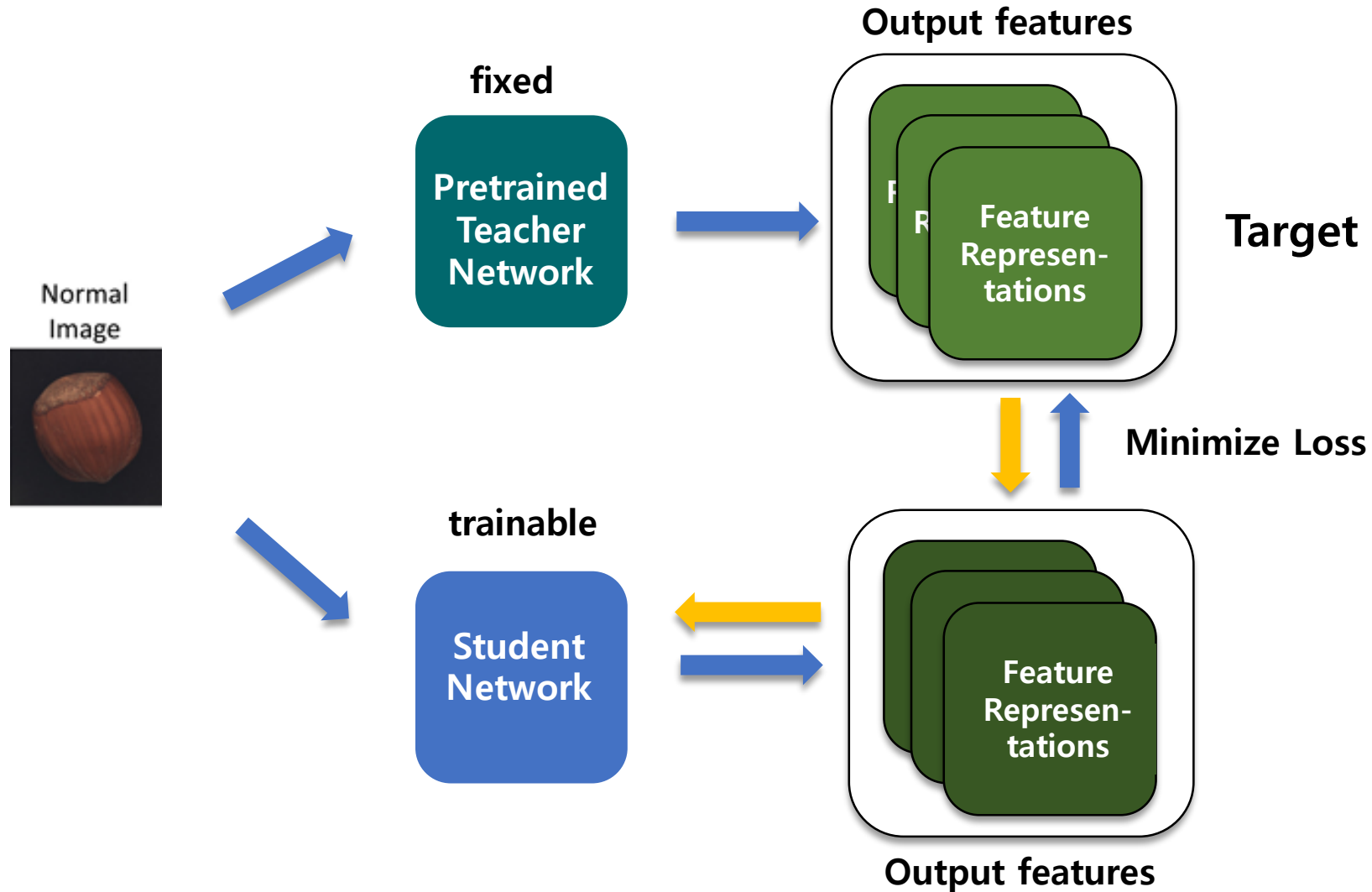


**Trainable  
Student  
Network**

- Mimic feature representations of teacher network
- Small size / but try to keep the performance same

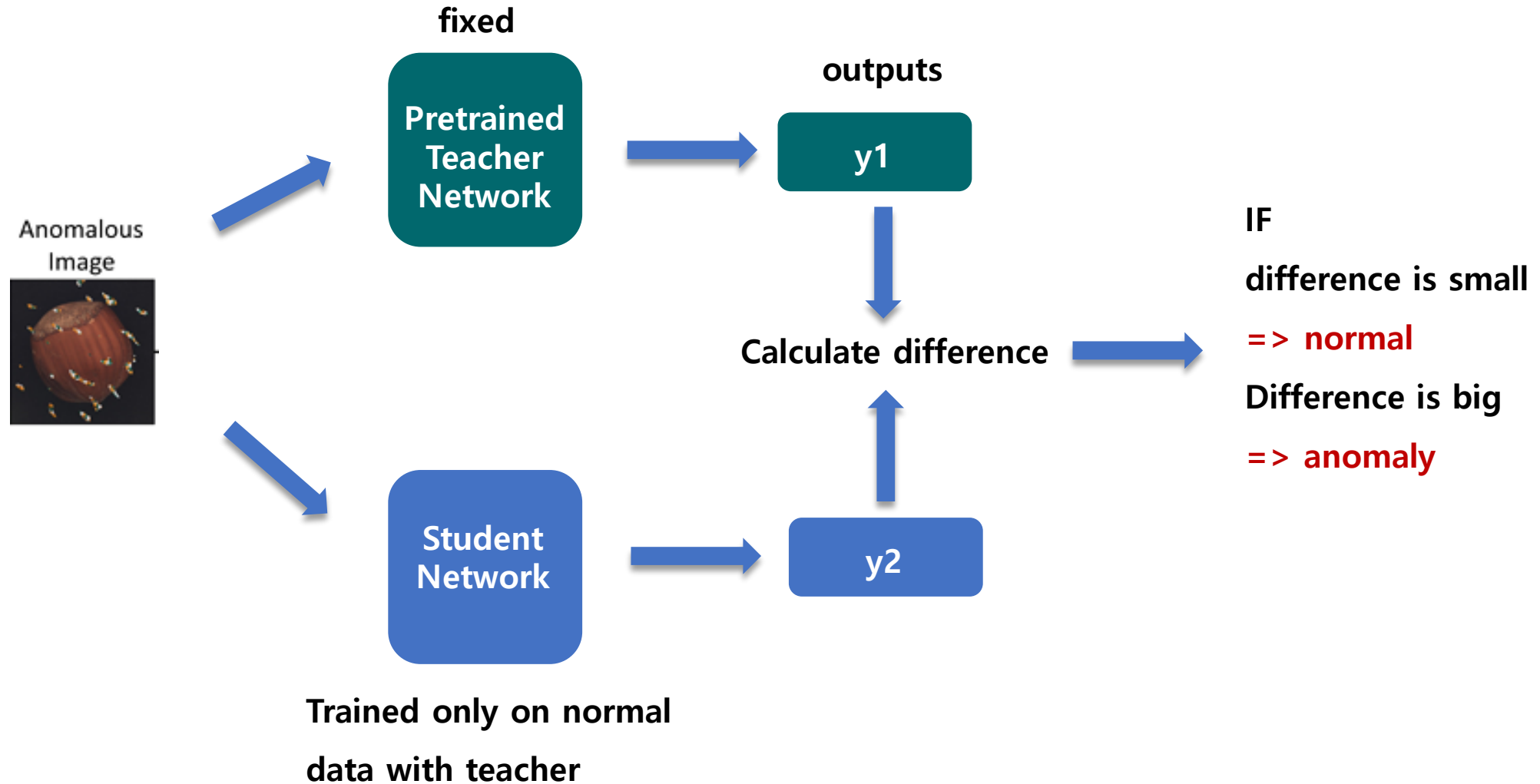
## 2. Introduction & Related Works

### TRAINING



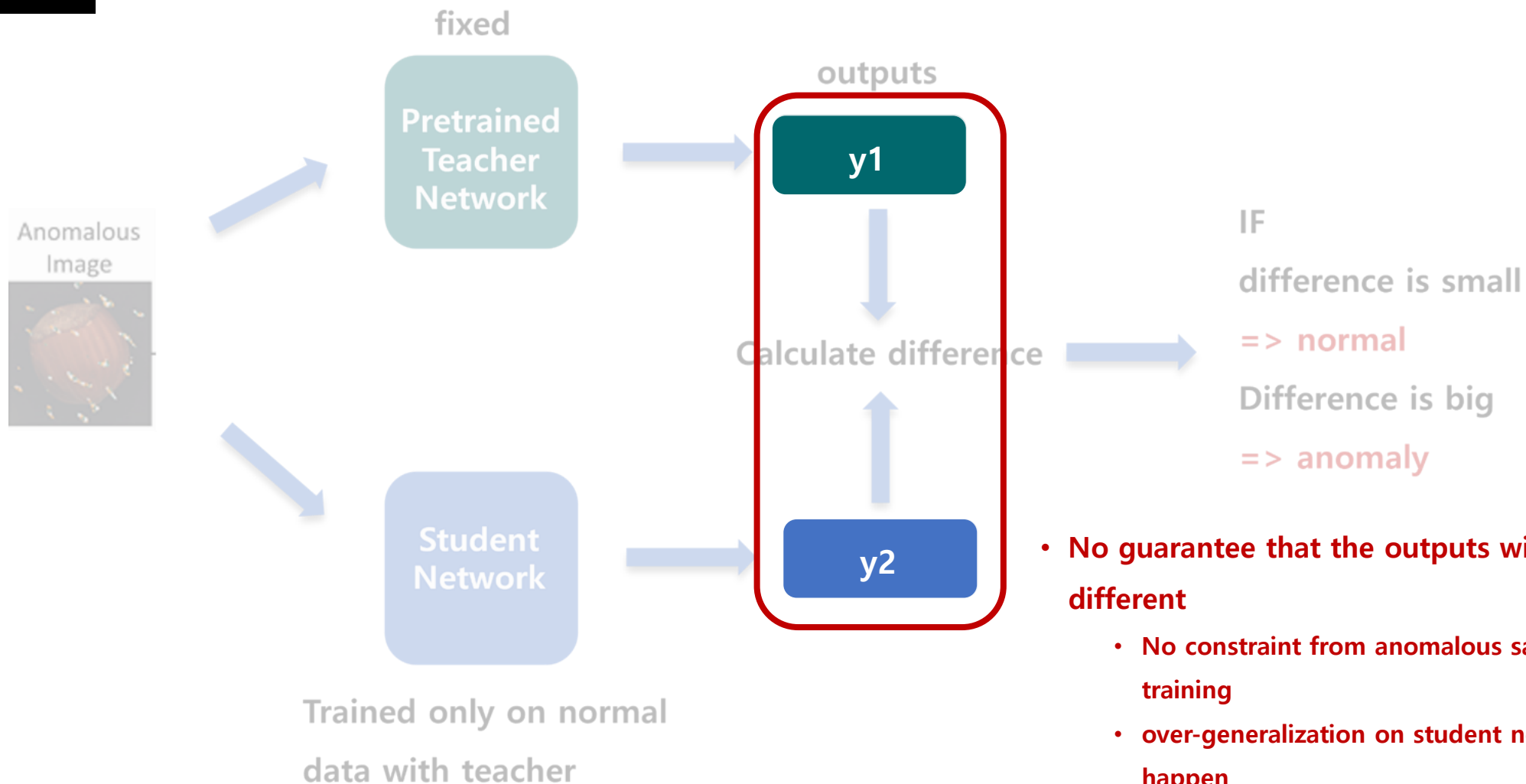
## 2. Introduction & Related Works

### INFERENCE



## 2. Introduction & Related Works

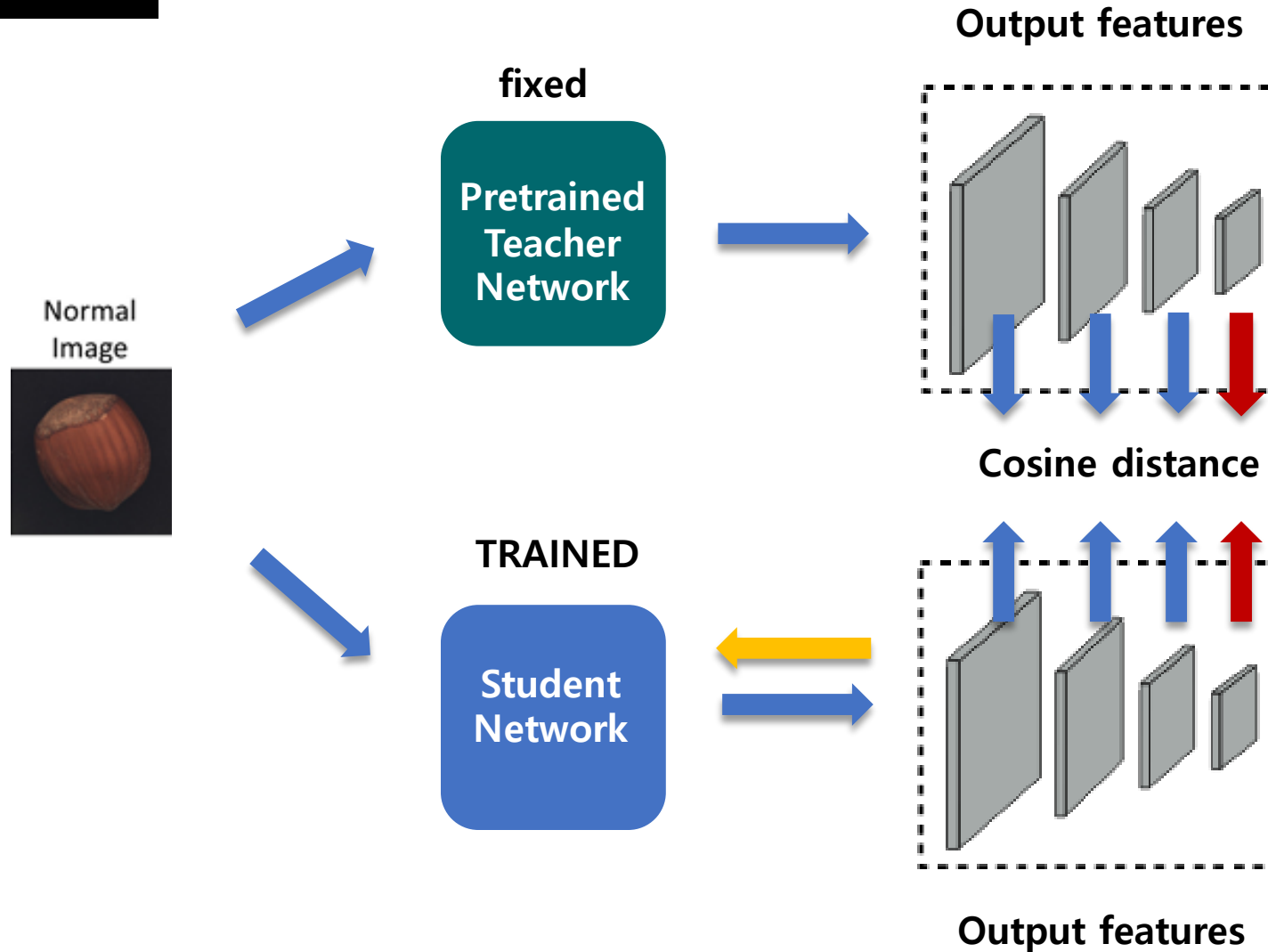
### INFERENCE



- No guarantee that the outputs will be always different
  - No constraint from anomalous sample during training
  - over-generalization on student network could happen

## 2. Introduction & Related Works

### INFERENCE

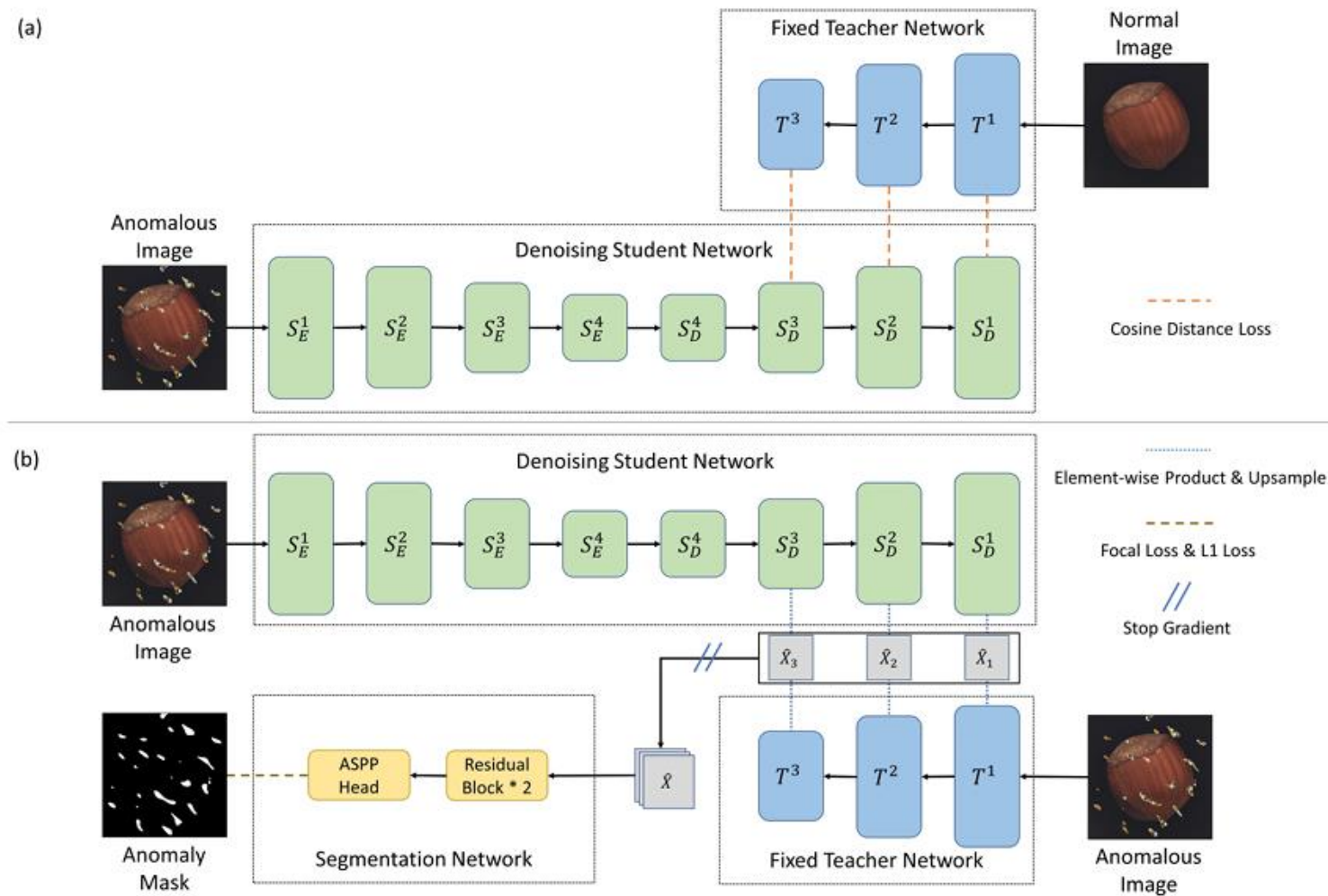


- Using multi level feature could be suboptimal
  - MVTec AD dataset (category of transistor)
  - 88.4% using only last layer feature representation
  - 81.9% on multi-level features

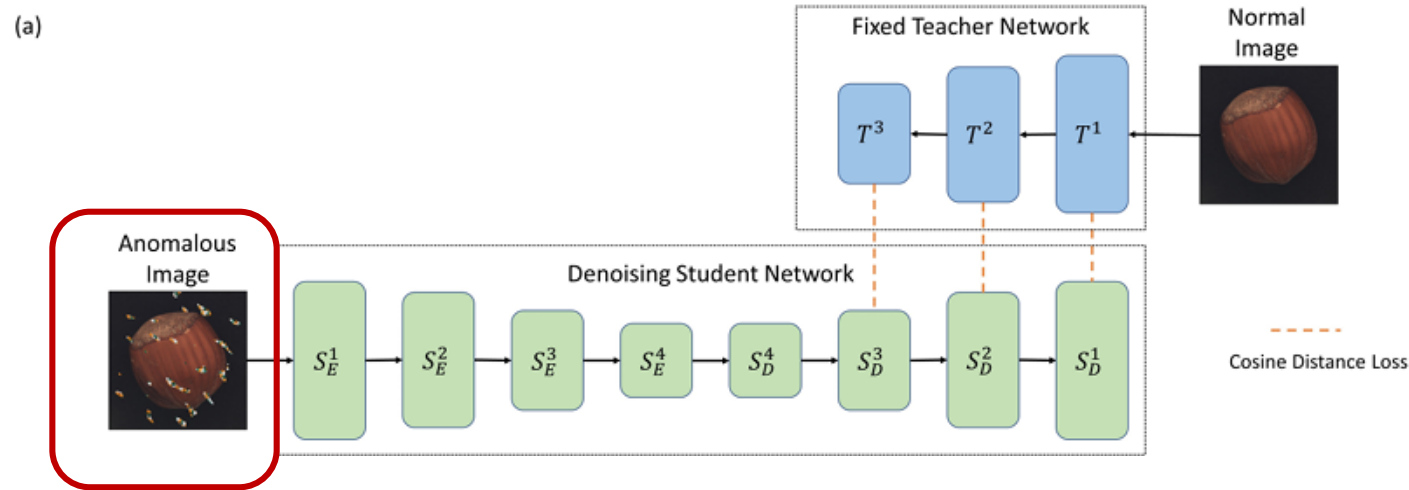


### 3. Method (DeSTSeg)

DeSTSeg = Denoising S + T + Segmentation



### 3. Method (DeSTSeg)



- Generated synthetic anomaly image
  - 1) Generate random two-dimensional Perlin noise
  - 2) Binarize to obtain anomaly mask ( $M$ )
  - 3) Replace mask region with anomaly-free image ( $I_a$ ) & arbitrary image from external data source ( $A$ )
  - 4) Apply opacity ( $\beta$ ) [0.15,1]

$$I_a = \beta(M \odot A) + (1 - \beta)(M \odot I_n) + (1 - M) \odot I_n \quad (1)$$

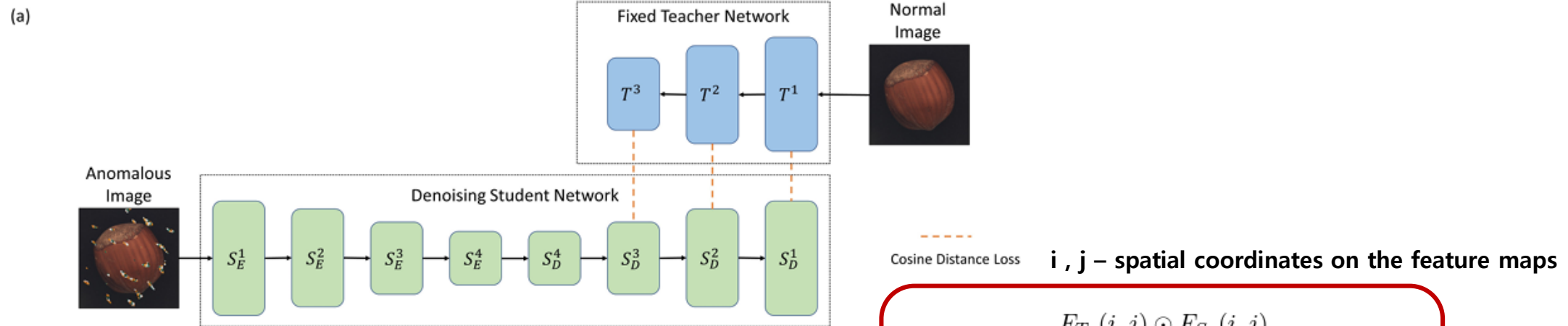
$\odot$  means the element-wise multiplication operation.

**DR/EM – A discriminatively trained reconstruction embedding for surface anomaly detection**

Vitjan Zavrtanik      Matej Kristan      Danijel Skočaj  
University of Ljubljana, Faculty of Computer and Information Science  
{vitjan.zavrtanik, matej.kristan, danijel.skocaj}@fri.uni-lj.si

### 3. Method (DeSTSeg)

- ImageNet pretrained ResNet18 with final block removed(conv5\_x)



- Encoder – randomly initialized ResNet18
- Decoder – reversed ResNet18
- Train student network to remove noise from anomalous image

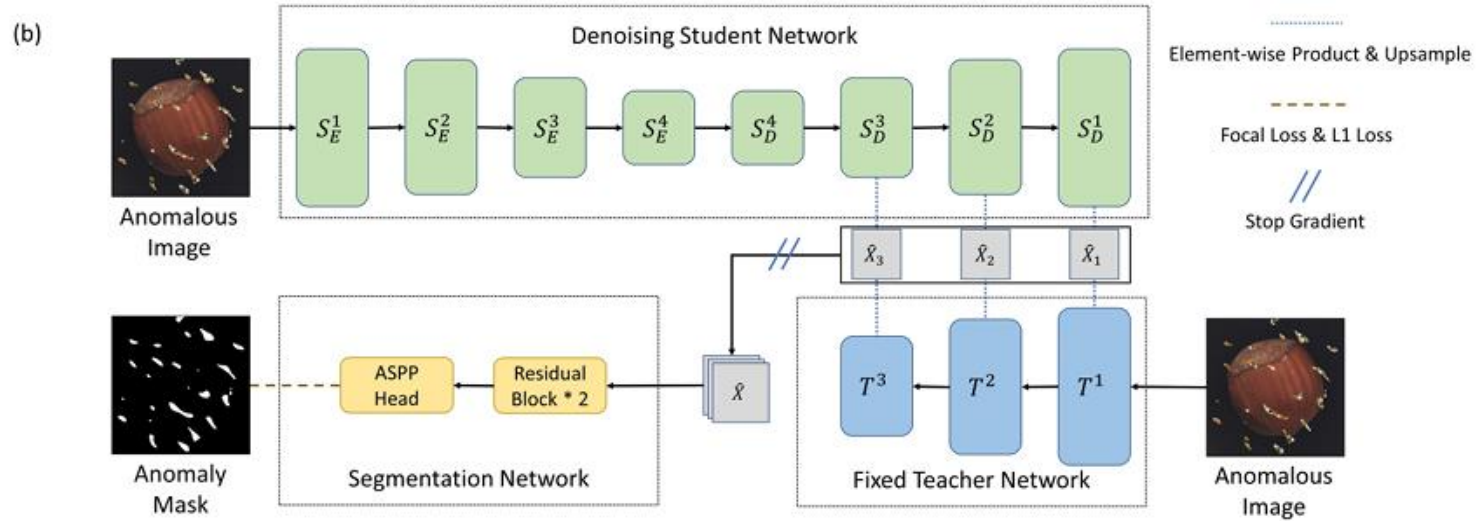
$$X_k(i, j) = \frac{F_{T_k}(i, j) \odot F_{S_k}(i, j)}{\|F_{T_k}(i, j)\|_2 \|F_{S_k}(i, j)\|_2} \quad (2)$$

Element-wise product between feature maps of S-T

$$D_k(i, j) = 1 - \sum_{c=1}^{C_k} X_k(i, j)_c \quad (3)$$

$$L_{cos} = \sum_{k=1}^3 \left( \frac{1}{H_k W_k} \sum_{i,j=1}^{H_k, W_k} D_k(i, j) \right) \quad (4)$$

### 3. Method

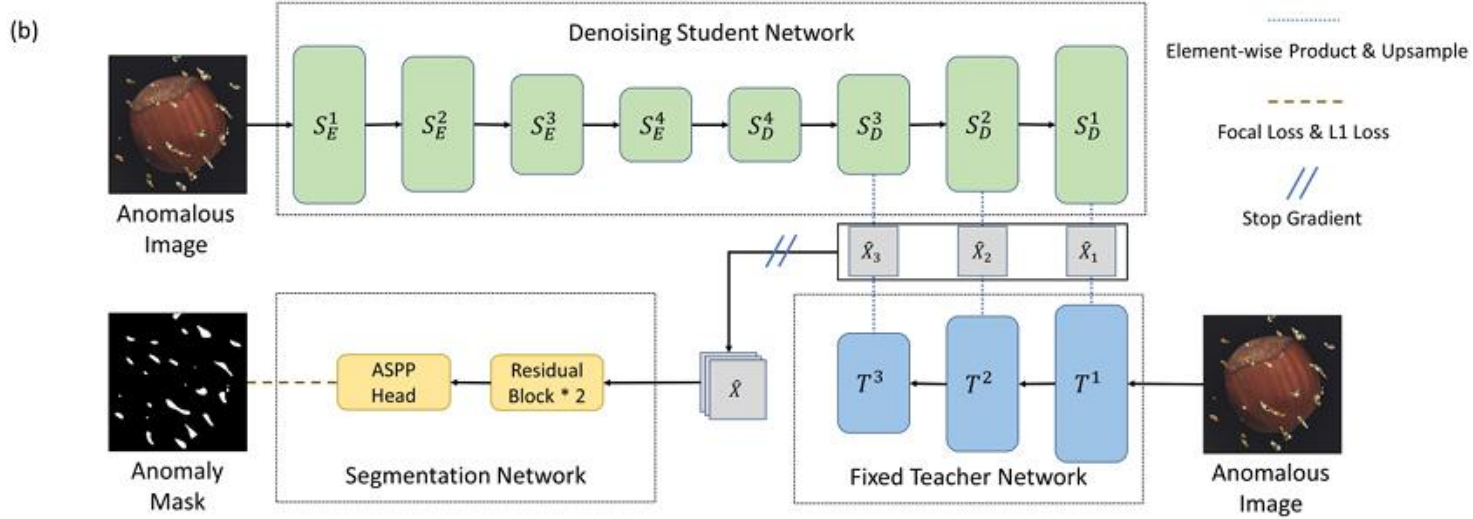


- Freeze Student & Teacher Network
- Anomalous images are used as input for both S & T networks
- Binary anomaly mask (M) which is generated before -> Ground Truth
- $\hat{X}_{1\sim3}$  : similarities of the paired feature maps -> calculated with (2)
- Then, concatenated as  $\hat{X}$  -> fed into segmentation network

- Similarities of the feature maps calculated
- Element-wise product between feature maps of S-T
- + Upsample =  $\hat{X}_{1\sim3}$

$$X_k(i, j) = \frac{F_{T_k}(i, j) \odot F_{S_k}(i, j)}{\|F_{T_k}(i, j)\|_2 \|F_{S_k}(i, j)\|_2} \quad (2)$$

### 3. Method



- Loss for segmentation

$$L_{focal} = -\frac{1}{H_1 W_1} \sum_{i,j=1}^{H_1, W_1} (1 - p_{ij})^\gamma \log(p_{ij}) \quad (5)$$

$$L_{l1} = \frac{1}{H_1 W_1} \sum_{i,j=1}^{H_1, W_1} |M_{ij} - \hat{Y}_{ij}| \quad (6)$$

$$L_{seg} = L_{focal} + L_{l1} \quad (7)$$

- Focal loss (5)

- 1) Even in anomalous image, majorities are background
- 2) Focus on the minority category

- L1 loss (6)

- 1) Improve sparsity of the output
- 2) Segmentation mask's boundaries are more distinct

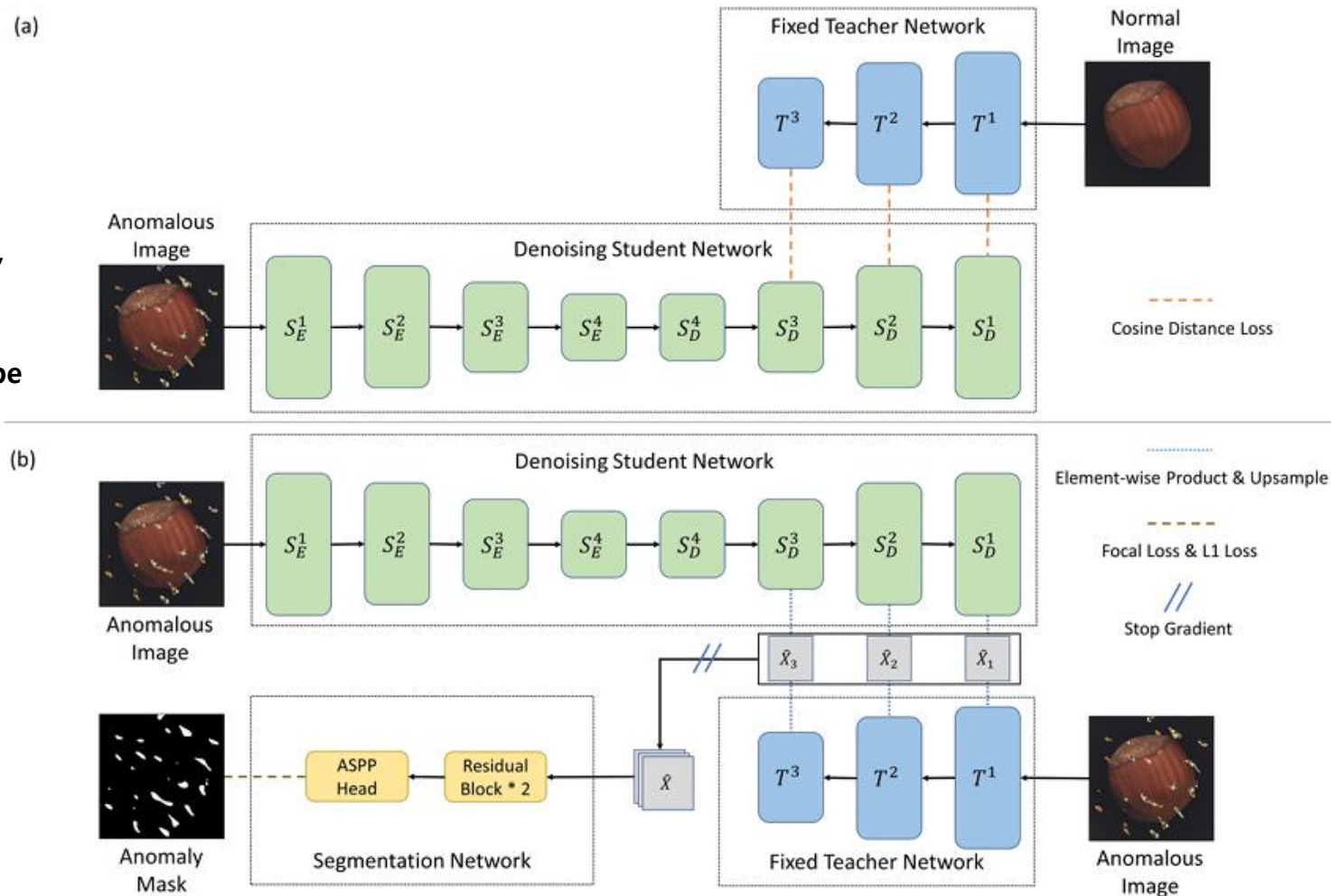
### 3. Method (DeSTSeg)

**Problem 1 : No guarantee that the outputs will be different**

- 1) Training "Denoising Student Network" with encoder-decoder structure  
=> Guaranteed the output features will be different

**Problem 2 : Using multi level feature could be suboptimal**

- 2) Training "Segmentation Network"  
=> Will use multi-level feature as input of segmentation training



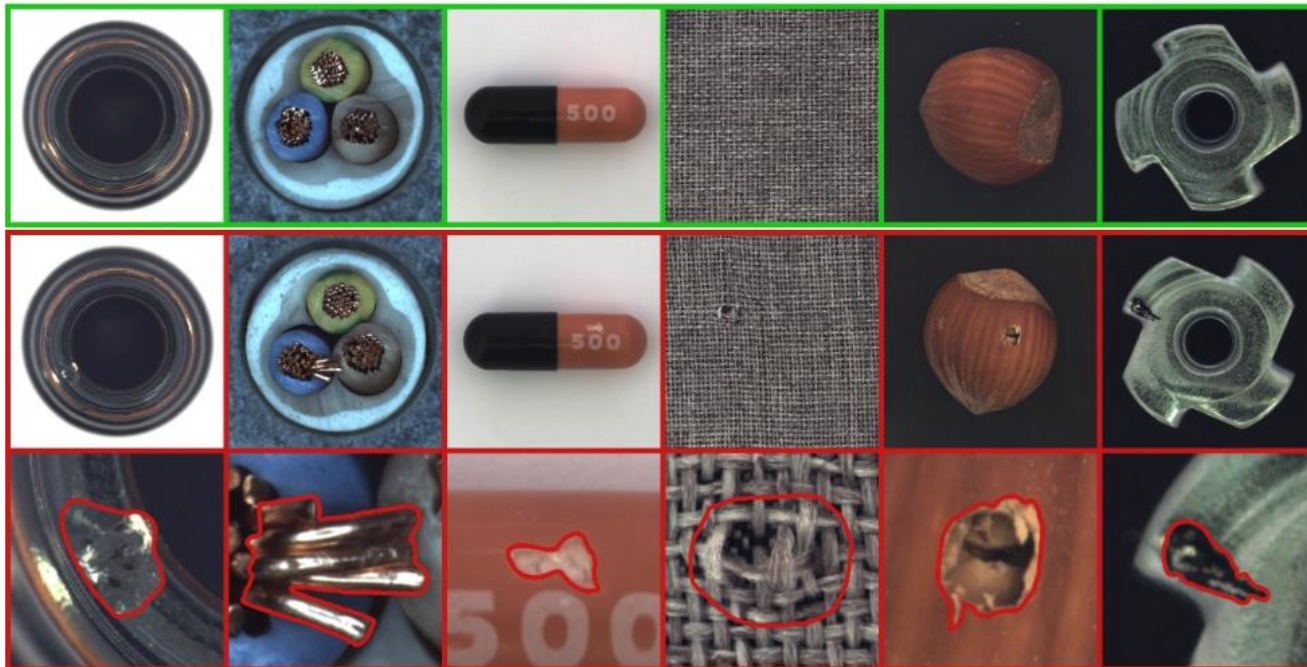


## 4. Experiments

### Dataset

MVTec AD – one of the most widely used benchmarks for anomaly detection and localization

- 15 categories
- Hundreds of normal images for training
- mixture of anomalous and normal images for evaluation



## 4. Experiments

US [3]	STPM [31]	CutPaste [16]	DRAEM [36]	DSR [37]	PatchCore [24]	Ours
87.7	95.1	95.2	98.0	98.2	98.5	<b>98.6<math>\pm</math>0.4</b>

Table 1. Image-level anomaly detection AUC (%) on MVTec AD dataset. Results are averaged over all categories.

	US [3]	STPM [31]	CutPaste [16]	DRAEM [36]	DSR [37]	PatchCore [24]	Ours
bottle	97.8 / 74.2	98.8 / 80.6	97.6 / -	<b>99.3</b> / 89.8	- / <b>91.5</b>	98.9 / 80.1	99.2 $\pm$ 0.2 / 90.3 $\pm$ 1.8
cable	91.9 / 48.2	94.8 / 58.0	90.0 / -	95.4 / 62.6	- / <b>70.4</b>	<b>98.8</b> / 70.0	97.3 $\pm$ 0.4 / 60.4 $\pm$ 2.3
capsule	96.8 / 25.9	98.2 / 35.9	97.4 / -	94.1 / 43.5	- / 53.3	<b>99.1</b> / 48.1	<b>99.1<math>\pm</math>0.0</b> / <b>56.3<math>\pm</math>1.1</b>
carpet	93.5 / 52.2	<b>99.1</b> / 65.3	98.3 / -	96.2 / 64.4	- / <b>78.2</b>	<b>99.1</b> / 66.7	96.1 $\pm$ 2.2 / 72.8 $\pm$ 5.8
grid	89.9 / 10.1	99.1 / 45.4	97.5 / -	<b>99.5</b> / 56.8	- / <b>68.0</b>	98.9 / 41.0	99.1 $\pm$ 0.1 / 61.5 $\pm$ 1.6
hazelnut	98.2 / 57.8	98.9 / 60.3	97.3 / -	99.5 / 88.1	- / 87.3	99.0 / 61.5	<b>99.6<math>\pm</math>0.2</b> / <b>88.4<math>\pm</math>2.2</b>
leather	97.8 / 40.9	99.2 / 42.9	99.5 / -	98.9 / 69.9	- / 62.5	99.4 / 51.0	<b>99.7<math>\pm</math>0.0</b> / <b>75.6<math>\pm</math>1.2</b>
metal_nut	97.2 / 83.5	97.2 / 79.3	93.1 / -	98.7 / 91.7	- / 67.5	<b>98.8</b> / 88.8	98.6 $\pm$ 0.4 / <b>93.5<math>\pm</math>1.1</b>
pill	96.5 / 62.0	94.7 / 63.3	95.7 / -	97.6 / 46.1	- / 65.7	98.2 / 78.7	<b>98.7<math>\pm</math>0.4</b> / <b>83.1<math>\pm</math>4.2</b>
screw	97.4 / 7.8	98.6 / 26.9	96.7 / -	<b>99.7</b> / <b>71.5</b>	- / 52.5	99.5 / 41.4	98.5 $\pm$ 0.3 / 58.7 $\pm$ 3.7
tile	92.5 / 65.3	96.6 / 61.7	90.5 / -	<b>99.5</b> / <b>96.9</b>	- / 93.9	96.6 / 59.3	98.0 $\pm$ 0.7 / 90.0 $\pm$ 2.5
toothbrush	97.9 / 37.7	98.9 / 48.8	98.1 / -	98.1 / 54.7	- / 74.2	98.9 / 51.6	<b>99.3<math>\pm</math>0.1</b> / <b>75.2<math>\pm</math>1.8</b>
transistor	73.7 / 27.1	81.9 / 44.4	93.0 / -	90.0 / 51.7	- / 41.1	<b>96.2</b> / 63.2	89.1 $\pm$ 3.4 / <b>64.8<math>\pm</math>4.0</b>
wood	92.1 / 53.3	95.2 / 47.0	95.5 / -	97.0 / 80.5	- / 68.4	95.1 / 52.3	<b>97.7<math>\pm</math>0.3</b> / <b>81.9<math>\pm</math>1.2</b>
zipper	95.6 / 36.1	98.0 / 54.9	99.3 / -	98.6 / 72.3	- / 78.5	99.0 / 64.0	<b>99.1<math>\pm</math>0.5</b> / <b>85.2<math>\pm</math>3.3</b>
average	93.9 / 45.5	96.6 / 54.3	96.0 / -	97.5 / 69.3	- / <b>70.2</b>	<b>98.4</b> / 61.2	97.9 $\pm$ 0.3 / <b>75.8<math>\pm</math>0.8</b>

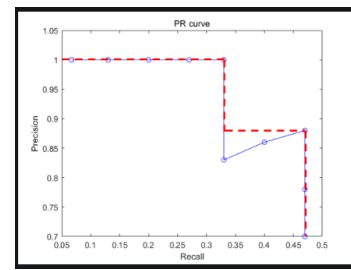
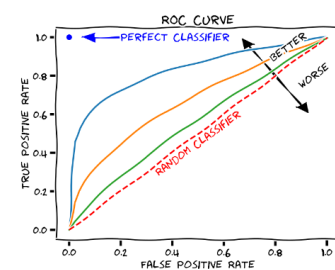
Table 2. Pixel-level anomaly localization AUC / AP (%) on MVTec AD dataset.

### Image-level anomaly detection

- > discriminate anomaly based on the whole image
- > if there is a anomaly region, the whole image is counted as an anomaly

### Pixel-level anomaly detection

- > each pixels are discriminated as normal or anomaly



- AUC (Area under ROC curve)
- AP (Area under PR curve)



## 4. Experiments

### Instance-level anomaly detection

-> Calculate overlap area between the object in ground truth mask and predicted mask

$$IAP = (TP + TN) / (TP + TN + FP + FN)$$

IAP@k% = if and only if overlapped area is over k%, considered as "detected"

Using different k thresholds, obtained average precision of this curve is called IAP.

	STPM [31]	DRAEM [36]	PatchCore [24]	Ours
bottle	83.2 / 73.3	90.3 / <b>84.8</b>	81.8 / 70.1	<b>90.5</b> $\pm$ 1.7 / 82.5 $\pm$ 4.1
cable	54.9 / 17.2	47.0 / 10.8	<b>69.2</b> / <b>50.6</b>	51.1 $\pm$ 2.5 / 26.7 $\pm$ 3.7
capsule	37.2 / 17.9	<b>50.7</b> / 21.4	44.2 / 26.9	49.4 $\pm$ 1.5 / <b>27.3</b> $\pm$ 3.3
carpet	68.4 / 52.2	76.8 / 32.3	64.4 / 43.7	<b>84.5</b> $\pm$ 4.9 / <b>58.6</b> $\pm$ 17.1
grid	45.7 / 21.0	55.5 / 42.3	39.1 / 15.6	<b>61.6</b> $\pm$ 1.8 / <b>47.4</b> $\pm$ 2.9
hazelnut	64.8 / 56.2	<b>95.7</b> / <b>89.0</b>	63.8 / 52.5	87.7 $\pm$ 1.8 / 77.6 $\pm$ 3.4
leather	46.2 / 24.9	<b>78.6</b> / 55.0	50.1 / 30.1	77.5 $\pm$ 1.8 / <b>65.3</b> $\pm$ 3.9
metal_nut	83.4 / 81.7	92.6 / 83.9	90.1 / 84.6	<b>93.6</b> $\pm$ 1.3 / <b>86.5</b> $\pm$ 2.7
pill	72.0 / 45.5	46.9 / 41.5	82.7 / <b>63.5</b>	<b>84.8</b> $\pm$ 3.8 / 61.1 $\pm$ 12.4
screw	24.4 / 4.2	<b>68.8</b> / <b>33.0</b>	38.4 / 16.3	53.6 $\pm$ 3.6 / 8.6 $\pm$ 2.3
tile	62.9 / 55.3	<b>98.9</b> / <b>98.2</b>	60.0 / 52.1	94.7 $\pm$ 1.8 / 86.5 $\pm$ 3.6
toothbrush	41.9 / 23.4	44.7 / 21.5	40.4 / 22.1	<b>59.8</b> $\pm$ 2.9 / <b>32.1</b> $\pm$ 5.1
transistor	53.4 / 8.5	59.3 / 22.8	69.9 / 36.8	<b>78.3</b> $\pm$ 2.5 / <b>49.6</b> $\pm$ 8.4
wood	56.0 / 35.4	<b>88.4</b> / 72.6	59.7 / 35.6	87.8 $\pm$ 2.8 / <b>76.4</b> $\pm$ 3.4
zipper	59.1 / 46.6	78.7 / 67.0	66.0 / 52.4	<b>90.6</b> $\pm$ 2.3 / <b>80.3</b> $\pm$ 4.9
average	56.9 / 37.5	71.5 / 51.7	61.3 / 43.5	<b>76.4</b> $\pm$ 1.0 / <b>57.8</b> $\pm$ 1.8

Table 3. Instance-level anomaly detection IAP / IAP@90 (%) on MVTec AD dataset.

IAP@90 : when 90% of anomaly instances are detected, the value of pixel-level precision

## 4. Experiments

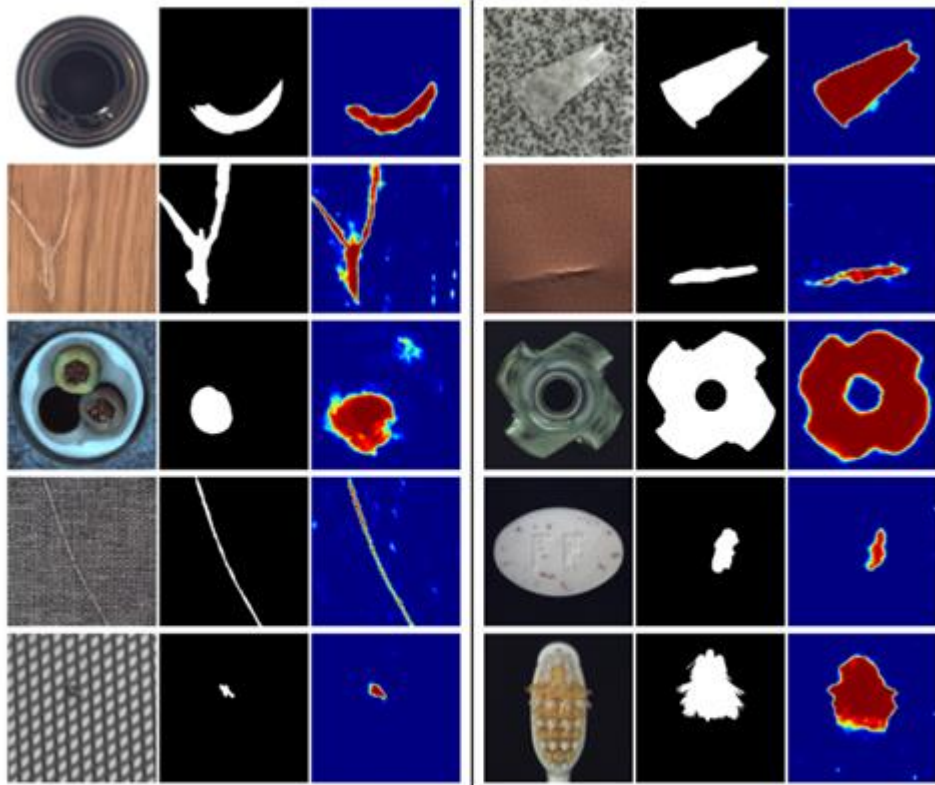


Figure 3. Visualization examples of our method. For each example, left: input image; middle: ground truth; right: prediction map.

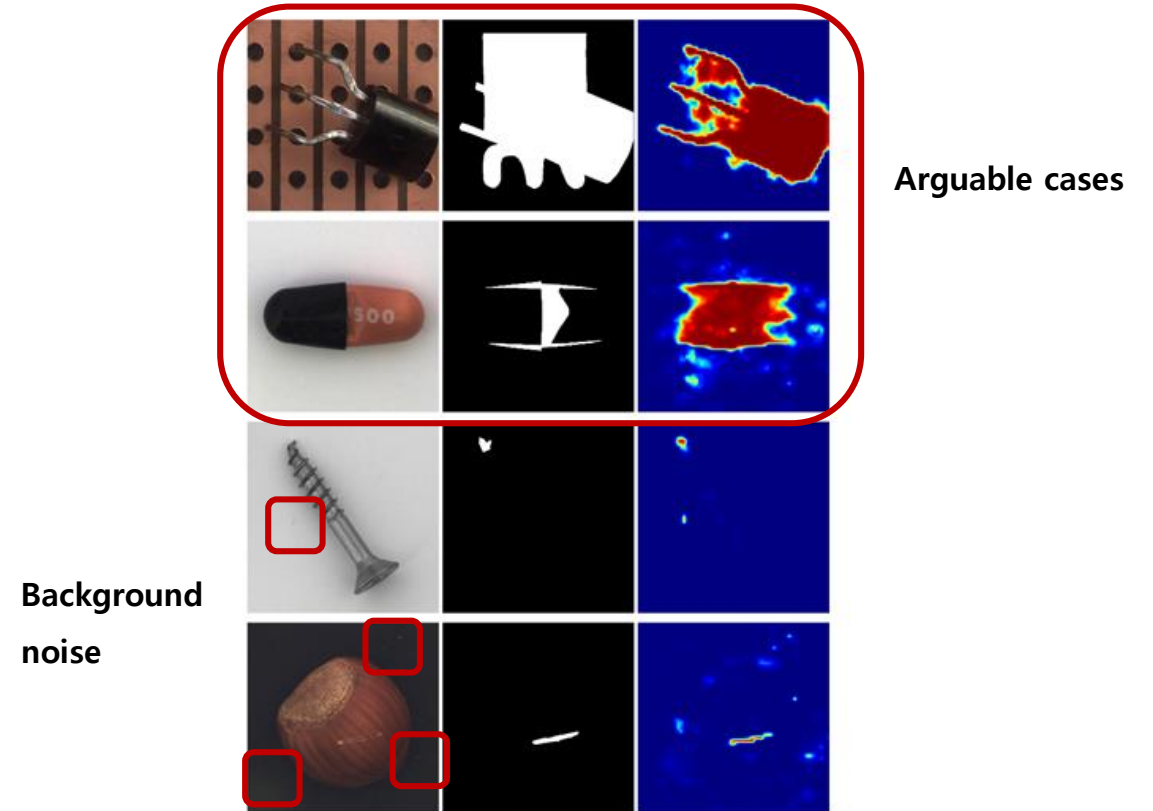


Figure 4. Failure cases of our method. The examples are chosen from transistor, capsule, screw, and hazelnut (from top to bottom). For each example, left: input image; middle: ground truth; right: prediction map.

## 4. Ablation studies

Exp.	den	ed	seg	img (AUC)	pix (AP)	ins (IAP)
1				94.8	52.9	55.8
2	✓			93.4	49.6	53.9
3		✓		95.4	53.3	57.7
4			✓	97.3	70.1	71.8
5	✓	✓		94.5	54.0	58.5
6	✓		✓	97.3	70.9	72.3
7		✓	✓	97.7	69.7	71.2
8	✓	✓	✓	<b>98.6</b>	<b>75.8</b>	<b>76.4</b>

Table 4. Ablation studies on our main designs: denoising training (den), the encoder-decoder architecture of student network (ed), and segmentation network (seg). AUC, AP, and IAP (%) are used to evaluate image-level, pixel-level, and instance-level detection, respectively. Exp. 1 uses the same architecture of [31], but different training settings to align with Exp. 2~8.

	img (AUC)	pix (AP)	ins (IAP)
w/o L1 loss	97.9	72.2	74.4
w/ L1 loss	<b>98.6</b>	<b>75.8</b>	<b>76.4</b>

Table 5. Ablation studies on the segmentation loss: AUC, AP, and IAP (%) are used to evaluate image-level, pixel-level, and instance-level detection, respectively.

## 4. Ablation studies

	img (AUC)	pix (AP)	ins (IAP)
concatenated-ST input	98.0	72.2	72.6
cosine-distance input	98.5	72.0	74.5
DeSTSeg	<b>98.6</b>	<b>75.8</b>	<b>76.4</b>

Table 6. Ablation studies on the **input of segmentation network**: AUC, AP, and IAP (%) are used to evaluate image-level, pixel-level, and instance-level detection, respectively.

Direct concatenation of feature maps of S-T networks

Computing cosine distance of S-T network's feature map

$$X_k(i, j) = \frac{F_{T_k}(i, j) \odot F_{S_k}(i, j)}{\|F_{T_k}(i, j)\|_2 \|F_{S_k}(i, j)\|_2} \quad (2)$$

$$D_k(i, j) = 1 - \sum_{c=1}^{C_k} X_k(i, j)_c \quad (3)$$

Element-wise product between feature maps of S-T

$$X_k(i, j) = \frac{F_{T_k}(i, j) \odot F_{S_k}(i, j)}{\|F_{T_k}(i, j)\|_2 \|F_{S_k}(i, j)\|_2} \quad (2)$$