

where  $i$  represents the index of channel,  $u$  is the index of spatial position (height together with width for simplicity). *Anomaly localization* aims to localize anomalous regions, producing an anomaly score map,  $\mathbf{s}(u)$ , which assigns an anomaly score for each pixel,  $u$ .  $\mathbf{s}(u)$  is calculated as the  $L2$  norm of the feature difference vector,  $\mathbf{d}(:, u)$ .

$$\mathbf{s}(u) = \|\mathbf{d}(:, u)\|_2. \quad (3)$$

*Anomaly detection* aims to detect whether an image contains anomalous regions. We intuitively take the maximum value of the averagely pooled  $\mathbf{s}(u)$  as the anomaly score of the whole image.

### 3.2 Preventing “Identical Mapping” with Transformer

We suspect that, compared with CNN, the query embedding in attention layer makes transformer difficult to learn an “identical mapping”. We denote the features in a normal image as  $\mathbf{x}^+ \in \mathbb{R}^{K \times C}$ , where  $K$  is the feature number,  $C$  is the channel dimension. The features in an anomalous image are denoted as  $\mathbf{x}^- \in \mathbb{R}^{K \times C}$ . We take a 1-layer network as the reconstruction net, which is trained on  $\mathbf{x}^+$  with the MSE loss and tested to detect anomalous regions in  $\mathbf{x}^-$ .

**Convolutional layer in CNN.** We first visit a fully-connected layer, whose weights and bias are denoted as  $\mathbf{w} \in \mathbb{R}^{C \times C}$ ,  $\mathbf{b} \in \mathbb{R}^C$ , respectively. When using this layer as the reconstruction model of normal samples, it can be written as,

$$\hat{\mathbf{x}} = \mathbf{x}^+ \mathbf{w} + \mathbf{b} \in \mathbb{R}^{K \times C}. \quad (4)$$

With the MSE loss pushing  $\hat{\mathbf{x}}$  to  $\mathbf{x}^+$ , the model may take shortcut to regress  $\mathbf{w} \rightarrow \mathbf{I}$  (identity matrix),  $\mathbf{b} \rightarrow \mathbf{0}$ . Ultimately, this model could also reconstruct  $\mathbf{x}^-$  well, failing in anomaly detection. A convolutional layer with  $1 \times 1$  kernel is equivalent to a fully-connected layer. Besides, An  $n \times n$  ( $n > 1$ ) kernel has more parameters and larger capacity, and can complete whatever  $1 \times 1$  kernel can. Thus, the convolutional layer also has the chance to learn a shortcut.

**Transformer with query embedding** contains an attention layer with a learnable query embedding,  $\mathbf{q} \in \mathbb{R}^{K \times C}$ . This attention layer can be denoted as,

$$\hat{\mathbf{x}} = \text{softmax}(\mathbf{q}(\mathbf{x}^+)^T / \sqrt{C}) \mathbf{x}^+ \in \mathbb{R}^{K \times C}. \quad (5)$$

To push  $\hat{\mathbf{x}}$  to  $\mathbf{x}^+$ , the attention map,  $\text{softmax}(\mathbf{q}(\mathbf{x}^+)^T / \sqrt{C})$ , should approximate  $\mathbf{I}$  (identity matrix), so  $\mathbf{q}$  must be highly related to  $\mathbf{x}^+$ . Considering that  $\mathbf{q}$  in the trained model is relevant to normal samples, the model could not reconstruct  $\mathbf{x}^-$  well. The ablation study in Sec. 4.4 shows that without the attention layer or the query embedding, the performance of transformer respectively drops by 2.4% or 3%, which is almost the same as CNN. This reflects that the query embedding in attention layer helps prevent transformer from learning an “identical shortcut”.

### 3.3 Adaptation with Anomaly-available Case

In practice, anomalies gradually increase with the runs of production lines, which brings the demands of compatibility with these increasing anomalies. Thus we adapt ADTR to ADTR+ for compatibility with the anomaly-available case.