

Poster Session:  
THU-PM-367



# Vector Quantization with Self-Attention for Quality-Independent Representation Learning

Zhou Yang<sup>1</sup>

Weisheng Dong<sup>1\*</sup>

Xin Li<sup>2</sup>

Menglun Huang<sup>1</sup>

Yulin Sun<sup>1</sup>

Guangming Shi<sup>1</sup>

<sup>1</sup>School of Artificial Intelligence  
Xidian University



西安电子科技大学  
XIDIAN UNIVERSITY

<sup>2</sup>Lane Department. of CSEE  
West Virginia University

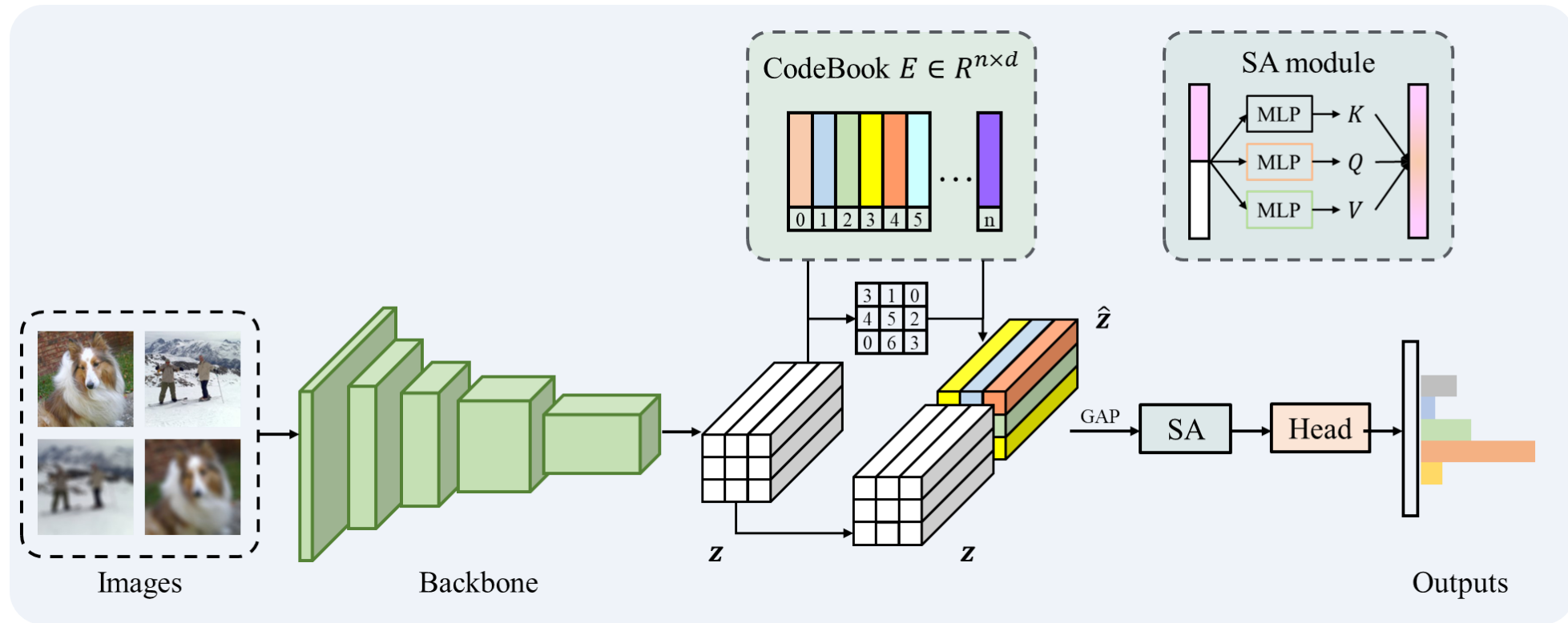


# Overview

- Problems:  
The robustness of deep models on the potential data distribution shift
- Motivation:  
Sparse representation can remove redundancy in signals and works well in image restoration
- Method:  
Vector quantization for quality-independent feature representation learning.
- Results:  
Better recognition performance on several benchmark datasets.

# Overview

- The overall architecture of our proposed method:



# background



clean



fog

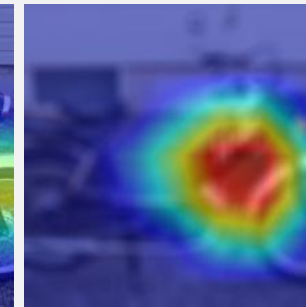
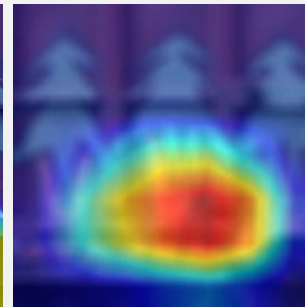
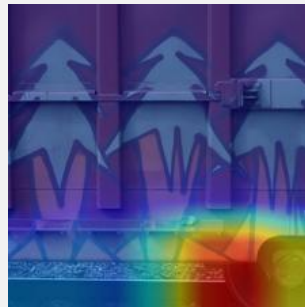


blur



noise

Deep models are typically trained on high-quality images but have poor recognition accuracy for low-quality ones

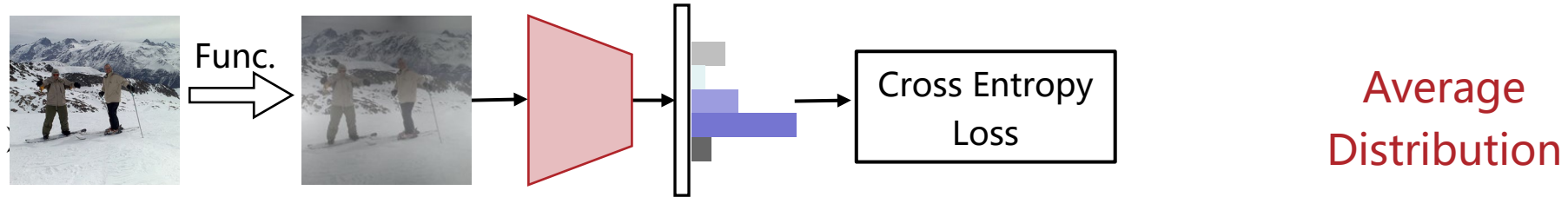


Visualization of CAM Class Activation Maps for Deep Models on Clean and Defocus-Blur Images

Deep features extracted from low-quality images are interfered, affecting the recognition

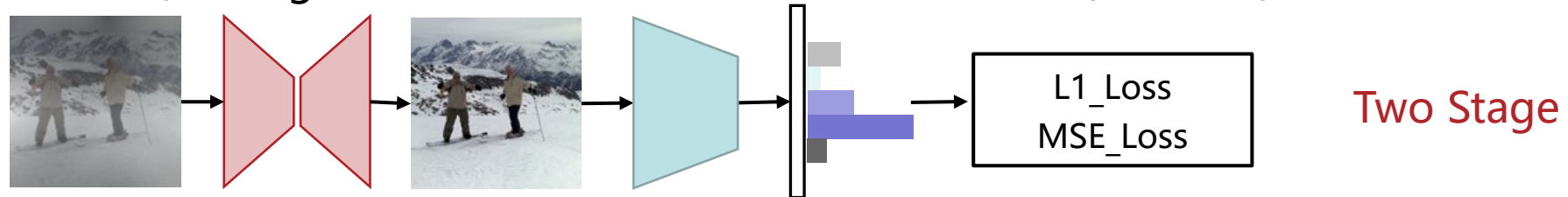
# Related Work

- Existing recognition methods for degraded images
  - Fine-tuning the models by mimicking real-world corruption

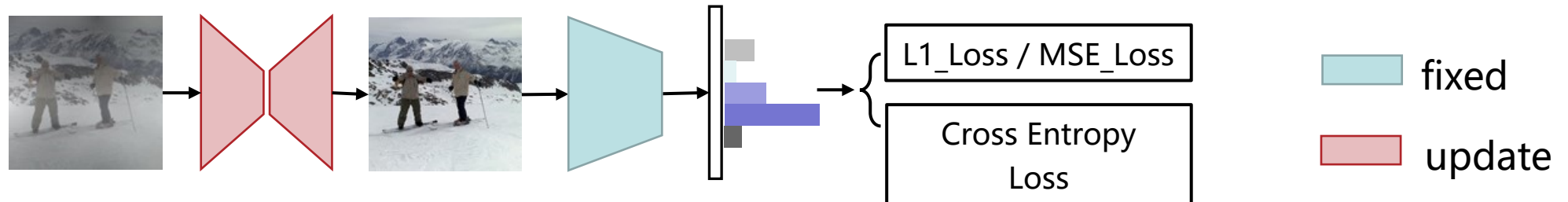


- Based on image restoration:

- a): recognition after restoration, Haze removal (ECCV'18)

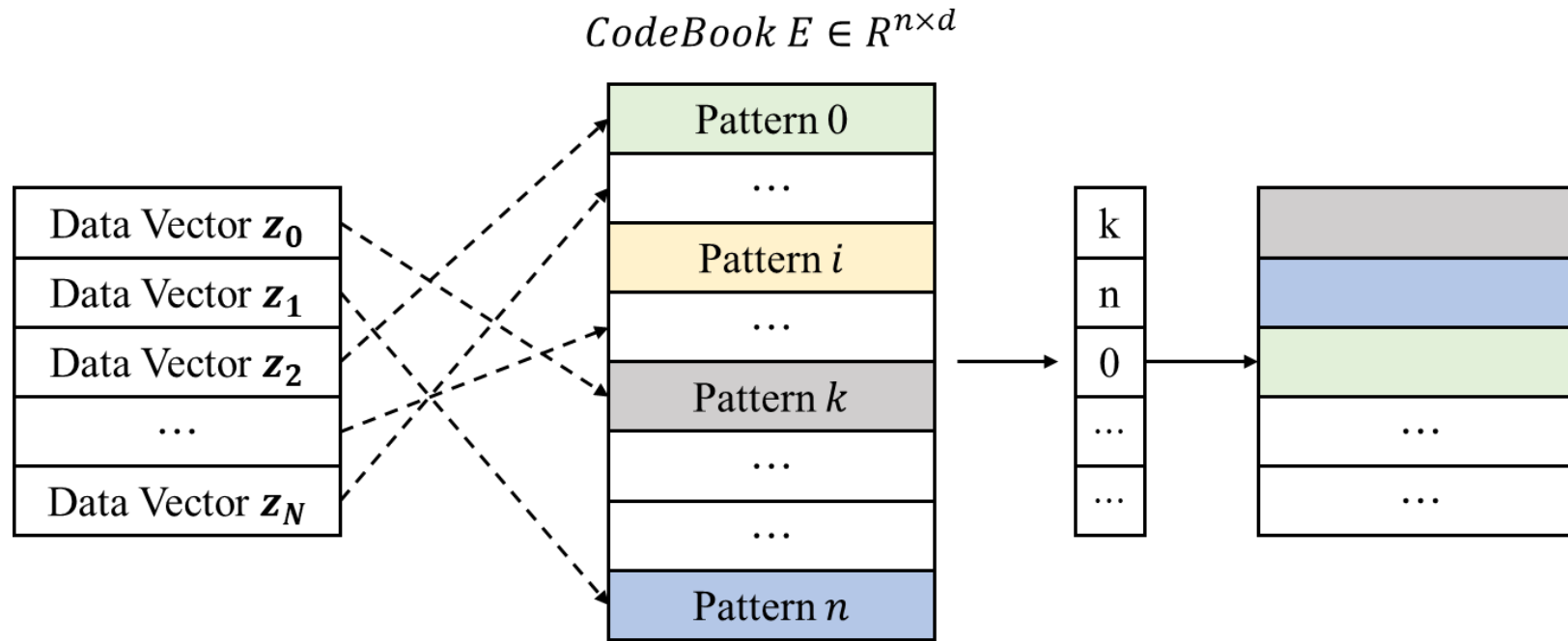


- b): Recognition friendly: Denoise&Recog (IJCAI'19), URIE (ECCV'20)



# Related Work

- Vector quantization:



$$\hat{\mathbf{z}} = E(\mathbf{z}) = \mathbf{e}_k, \quad \text{where } k = \arg \min_i \|\mathbf{z} - \mathbf{e}_i\|_2^2$$

# Method

➤ Quality-independent feature representation learning:

- Assuming that the quality-independent feature vector  $\hat{\mathbf{z}}$  of an image is a linear combination of a series of features (atoms):

$$\hat{\mathbf{z}} = \sum_i \alpha_i * \mathbf{e}_i = \alpha_0 * \mathbf{e}_0 + \alpha_1 * \mathbf{e}_1 + \dots + \alpha_n * \mathbf{e}_n, \mathbf{e}_i \in E, E \in R^{n \times d} \quad (1)$$

- We have this sparse representation, which need to optimize  $\alpha$  and  $E$  alternately:

$$\hat{\mathbf{z}} = E * \hat{\alpha}, \hat{\alpha} = \underset{\alpha}{\operatorname{argmin}} \|\mathbf{z} - E * \alpha\|_2^2 + \lambda * \|\alpha\|_0 \quad (2)$$

- Simplify  $\alpha$  as an one-hot vector, we have:

$$\hat{\mathbf{z}} = \mathbf{e}_k \quad (3)$$

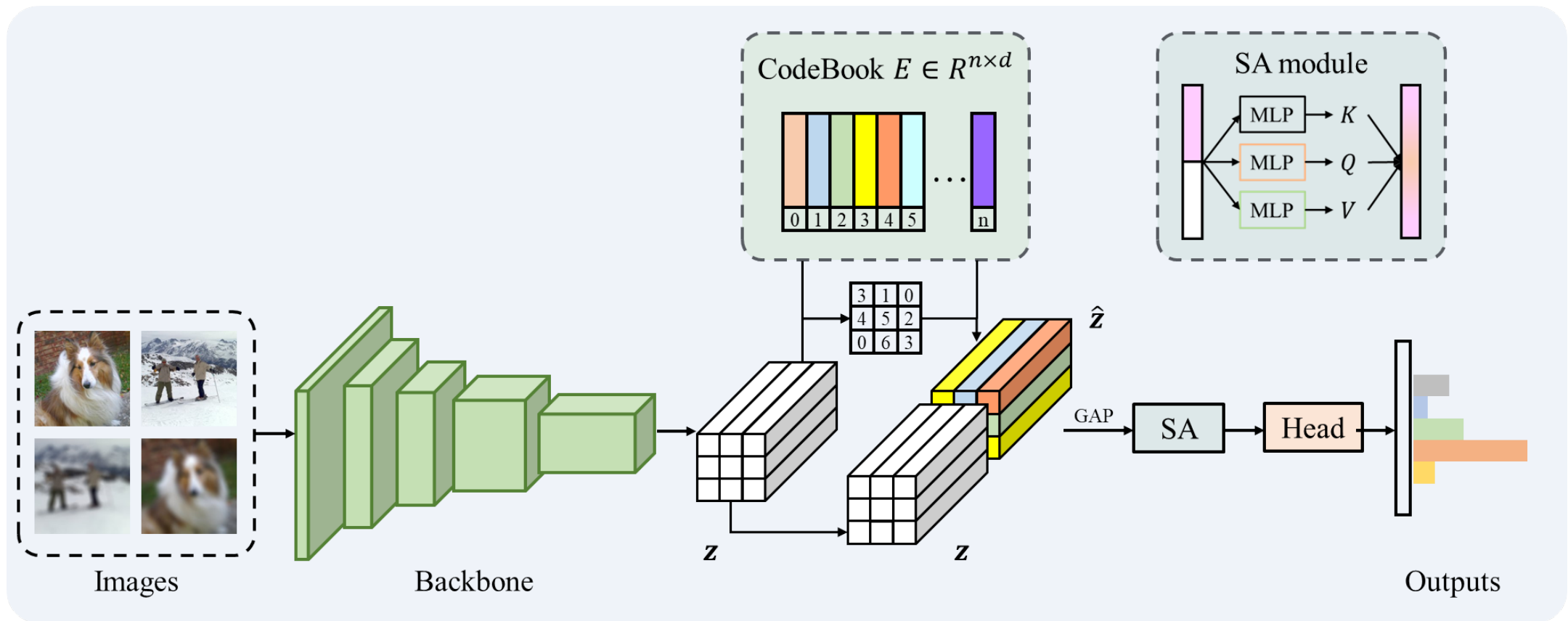
- Vector quantization as VQ-VAE:

$$L_{vq} = \|\operatorname{sg}(\mathbf{z}) - \hat{\mathbf{z}}\|_2^2 \quad (4)$$

$$L_{cmt} = \|\mathbf{z} - \operatorname{sg}(\hat{\mathbf{z}})\|_2^2 \quad (5)$$

# Method

➤ Overall architecture of our proposed approach:





# Method

- Training Loss:

$$L_{total} = L_{ce} + \lambda * (L_{vq} + \beta * L_{cmt}) \quad (6)$$

- Concatenate & self-attention:

$$f = \text{Cat}(\mathbf{z}, \hat{\mathbf{z}}) \quad (7)$$

$$f_{sa} = \text{softmax} \left( K * \frac{Q^T}{\sqrt{d^n}} \right) * V \quad (8)$$

- Experiments have shown that these skills can further improve performance.

# Ablation Study

- The impact of codebook size  $n$ :

Size	# Params	clean $\uparrow$	mCE $\downarrow$
$n = 1k$	$4.0 \times 10^7$	76.1	45.7
$n = 10k$	$5.8 \times 10^7$	76.6	43.1
$n = 100k$	$2.4 \times 10^8$	76.6	42.9

- The choice of fusion mode:

CodeBook	Fusion mode	SA	clean $\uparrow$	mCE $\downarrow$
-	-	-	73.1	53.7
✓	replace	-	74.3	50.1
✓	add	-	74.7	48.9
✓	concat	-	76.2	45.7
✓	concat	✓	<b>76.6</b>	<b>43.1</b>

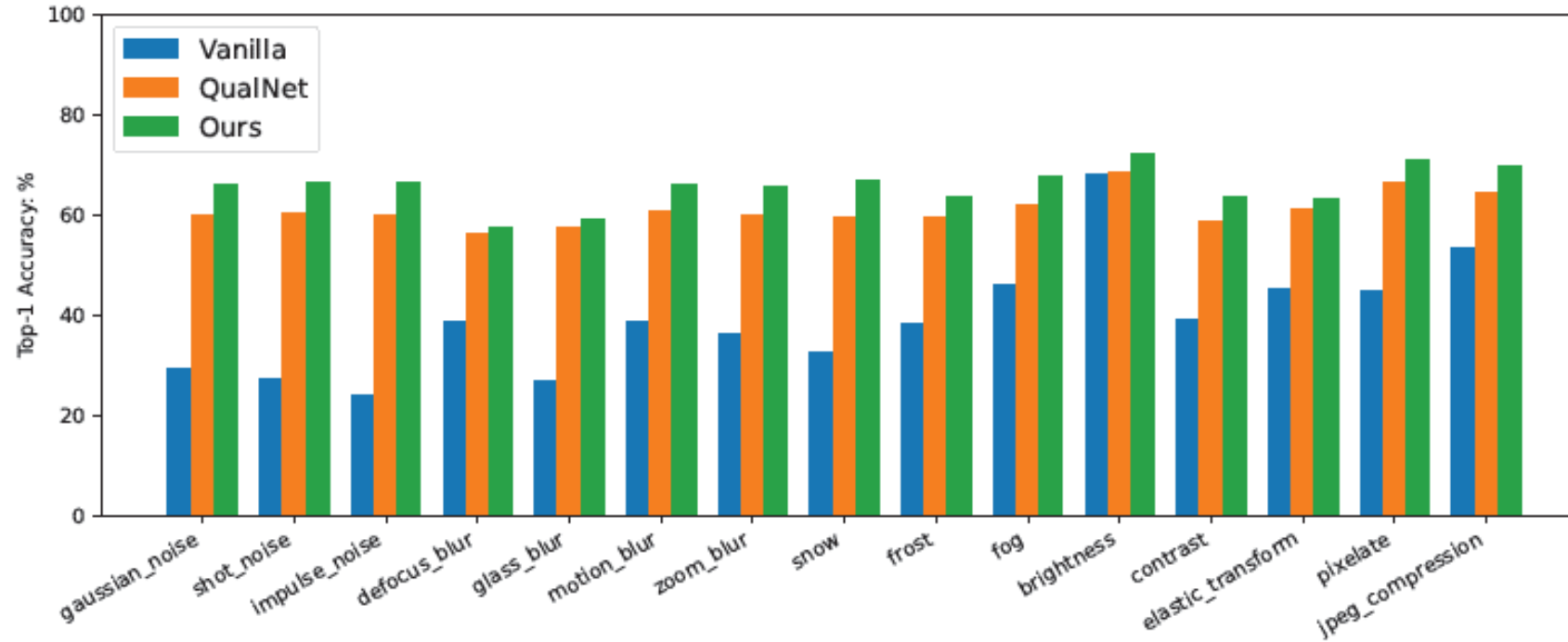
# Results

Method	Backbone	Clean	Known	UnKnown	mCE ↓
Vanilla [20]	ResNet50	76.1	39.1	46.7	76.7
DDP [55]		72.1	48.2	50.7	62.78
URIE [47]		73.8	55.1	56.5	55.7
QualNet [29]		75.4	61.1	58.1	50.3
Ours		<b>76.6</b>	<b>65.6</b>	<b>60.2</b>	<b>43.1</b>
Vanilla [20]	ResNeXt101	79.6	47.1	55.5	69.7
QualNet [29]		77.8	65.5	63.3	42.6
Ours		<b>80.3</b>	<b>68.6</b>	<b>64.5</b>	<b>37.9</b>

Method	Clean	ImageNet-C ↓	ImageNet-A	ImageNet-R
Vanilla [20]	76.1	76.7	0.0	36.2
+ Ours	<b>76.6</b>	<b>71.1</b>	<b>3.7</b>	<b>38.6</b>
DeepAugment [21]	76.6	60.4	3.5	42.2
+ AugMix [23]	75.8	53.5	3.9	46.8
+ DAT [35]	77.1	50.8	<b>6.8</b>	47.8
DAu+AM+Ours	<b>77.4</b>	<b>48.7</b>	5.9	<b>49.3</b>

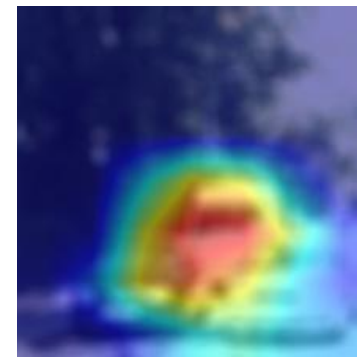
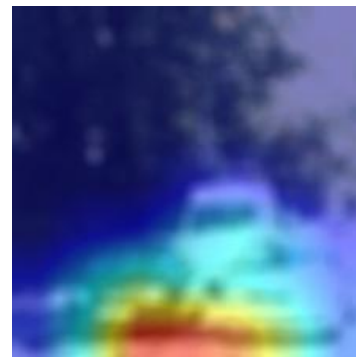
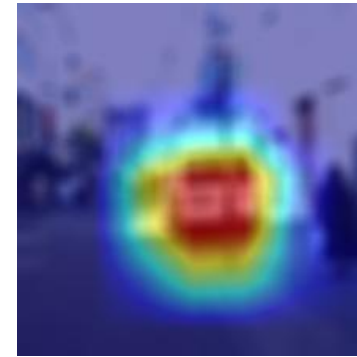
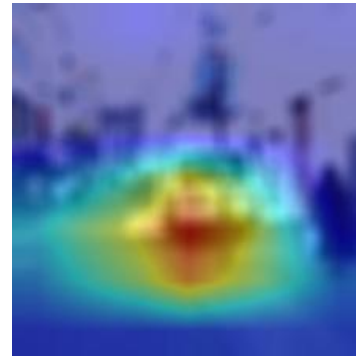
The top-1 accuracy of each method on several benchmark datasets.

# Results



The detailed top-1 accuracy results of the different methods for each corruption type in the benchmark dataset ImageNet-C.

# Results



Clean image

CAM on clean

CAM on blur

After our method

# Summary

- We propose to introduce vector quantization into the recognition model and improve the models' robustness on common corruptions.
- We concatenate the quantized feature vector with the original one and use the self-attention module to enhance the quality-independent feature representation instead of direct replacement in the standard vector quantization method.
- Extensive experimental results show that our method has achieved higher accuracy on benchmark low-quality datasets than several current sota methods.

# Thank you for watching!

Code is available at: <https://github.com/yangzhou321/VQSA>

