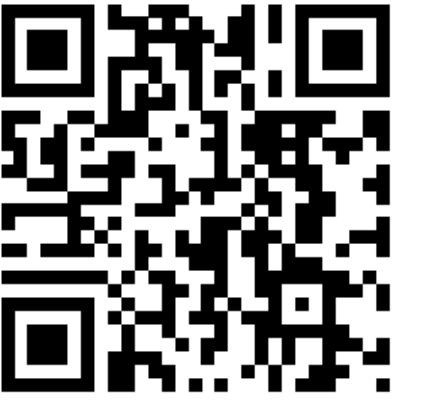


Regional Attention Based Deep Feature for Image Retrieval

Jaeyoon Kim Sung-Eui Yoon

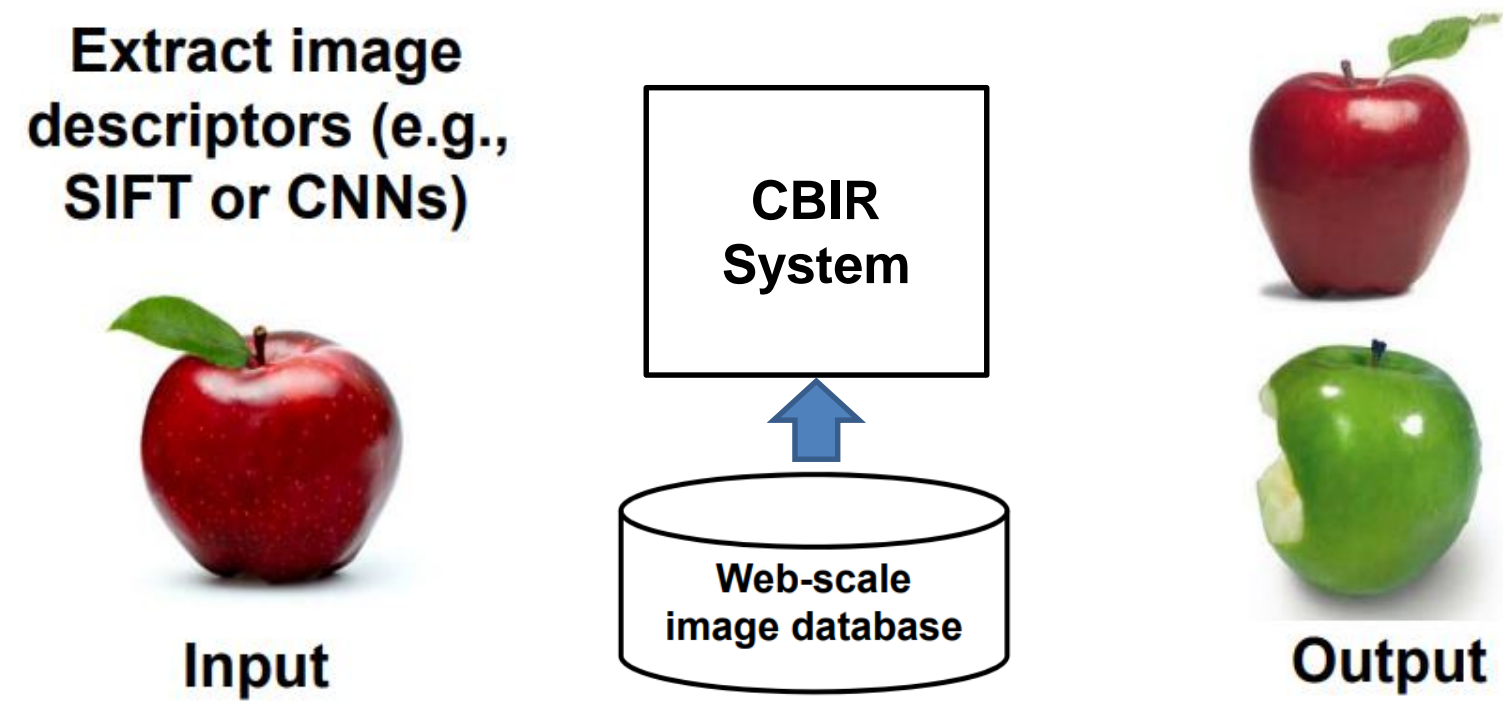
KAIST



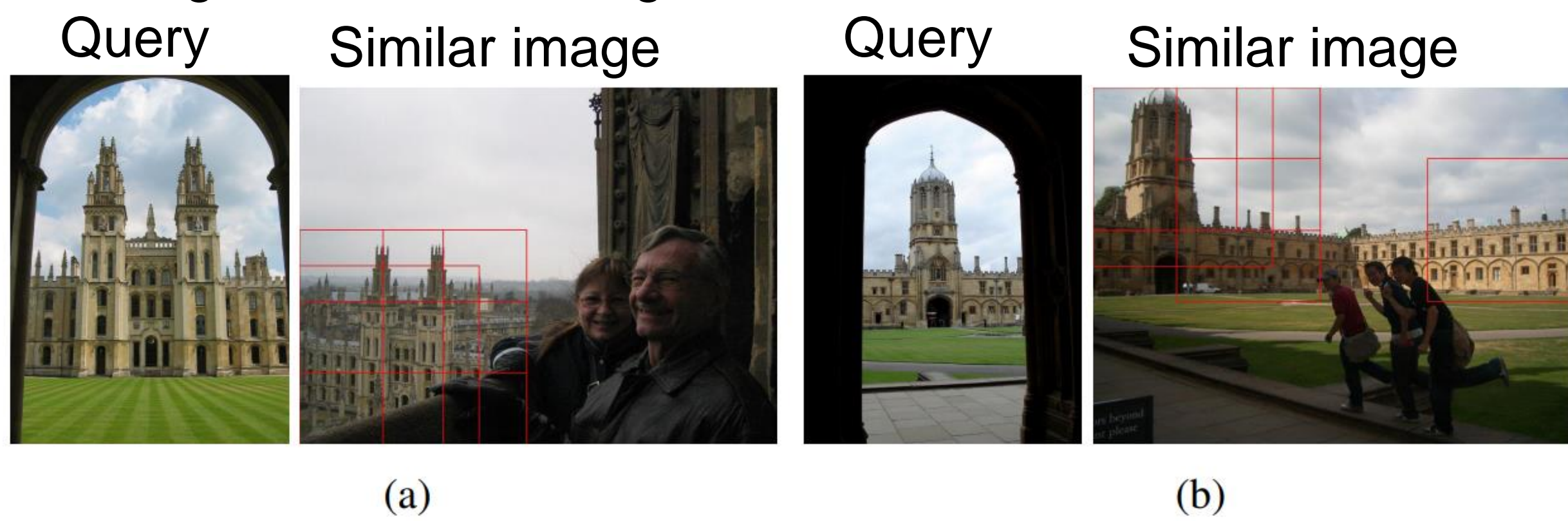
Trained model and source codes are available
<https://sglab.kaist.ac.kr/RegionalAttention/>

1. Intro. & Motivation

Content-Based Image Retrieval(CBIR)

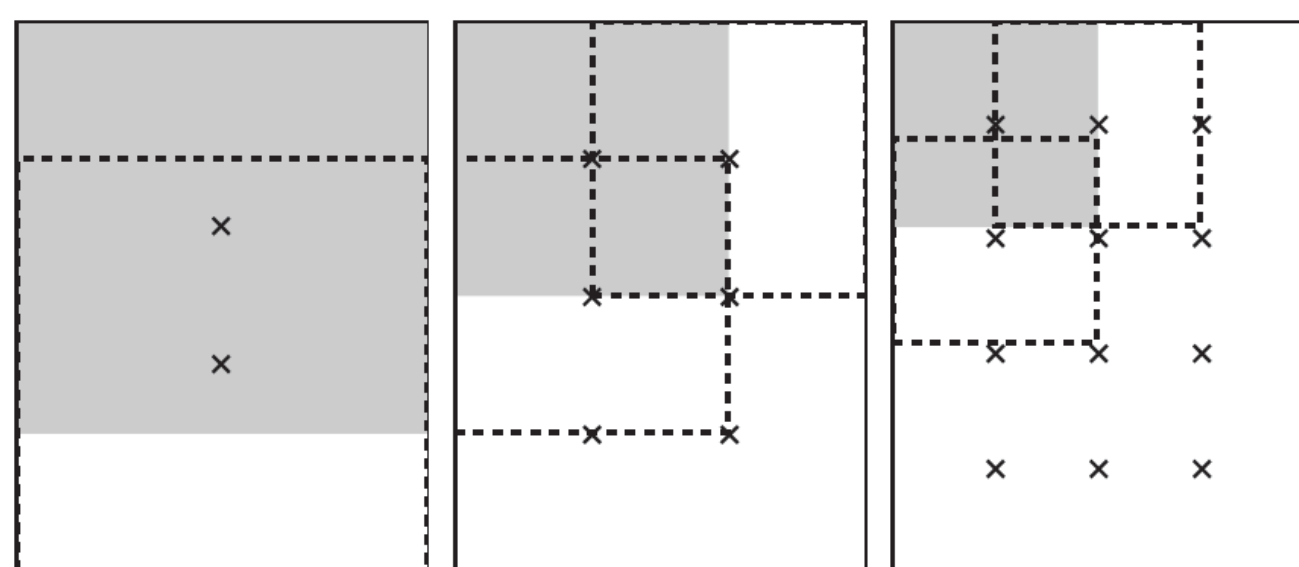


Challenge in CBIR – Backgrounds and Clutters



Regional Maximum Activation of Convolutions(R-MAC)[1]

R-MAC is sometimes suffered from backgrounds and clutters since it uniformly samples regions of an image.



[1] Tolias, G., Sivic, R., Jegou, H.: Particular object retrieval with integral max-pooling of cnn activations. In: ICLR. (2016)

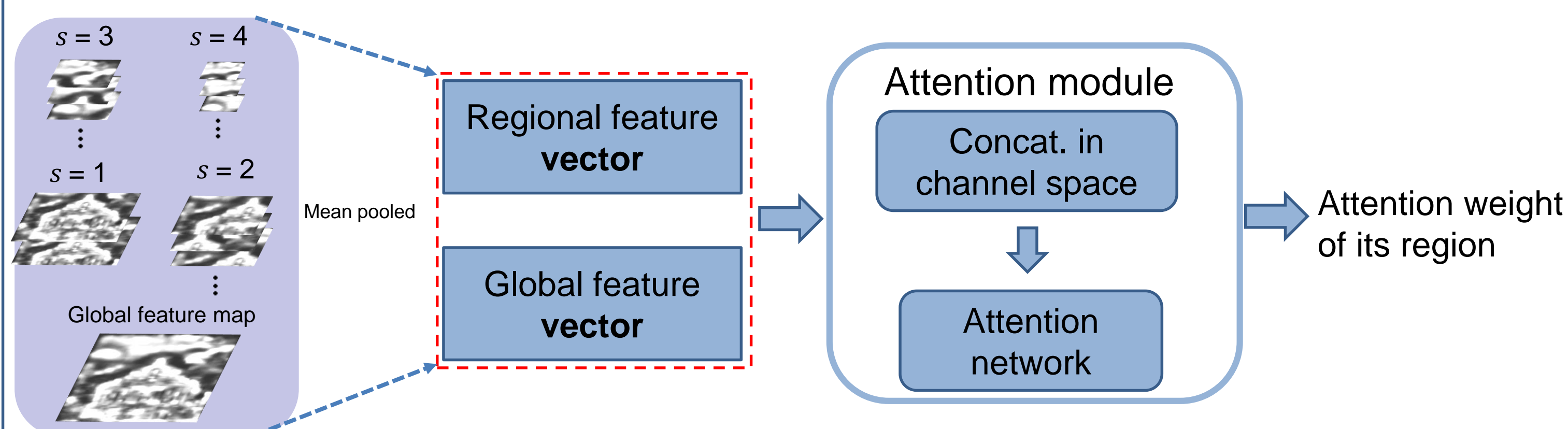
Categories in Image Retrieval

- "Fine-tuned": Fine-tuned CNN for specific dataset(data category).
- "Pre-trained": Off-the-shelf CNN from ImageNet

3. Context-Aware Regional Attention

Proposed Regional Attention Network

- Context awareness: Consider both of **local** and **global** feature of its region.



Attention Network

- Two linear layers and Two non-linear layers

$$\Phi(\mathbf{k}) = \text{softplus}(\mathbf{W}_c \pi(\mathbf{k}) + \mathbf{b}_c),$$

$$\pi(\mathbf{k}) = \tanh(\mathbf{W}_r \mathbf{k} + \mathbf{b}_r).$$

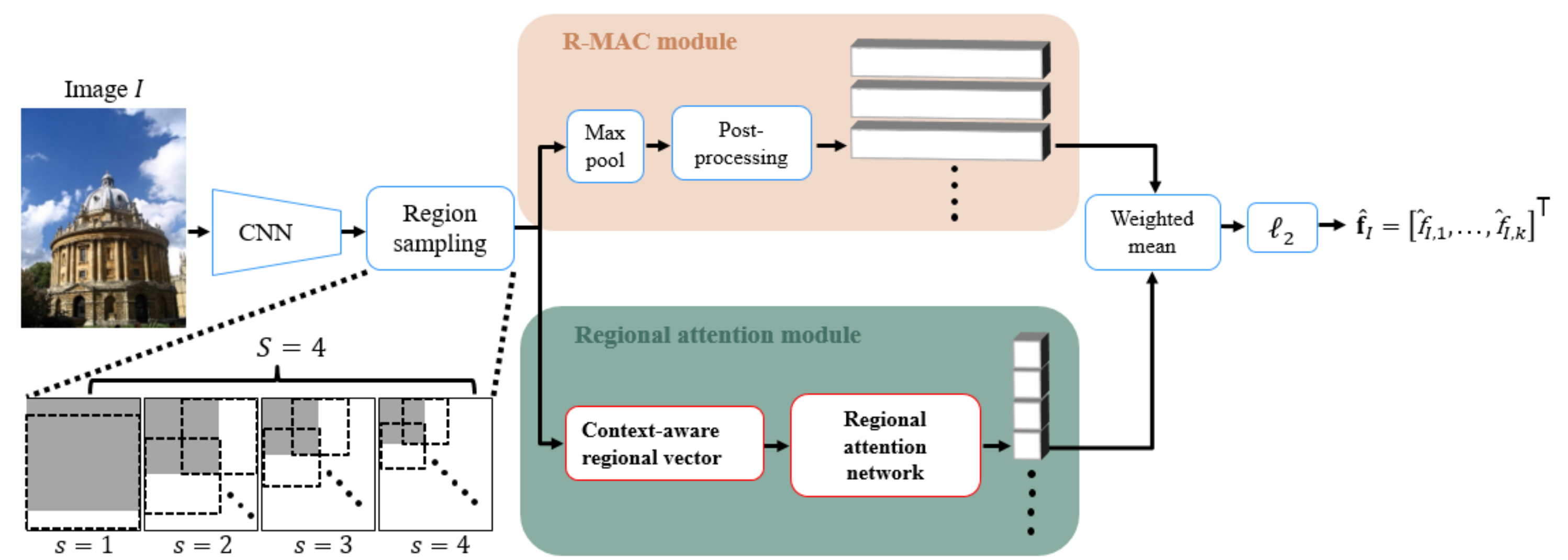
\mathbf{k} : combined feature vector(local, global), $\Phi(\mathbf{k})$: Attention weight of \mathbf{k}

Training the context-aware regional attention network

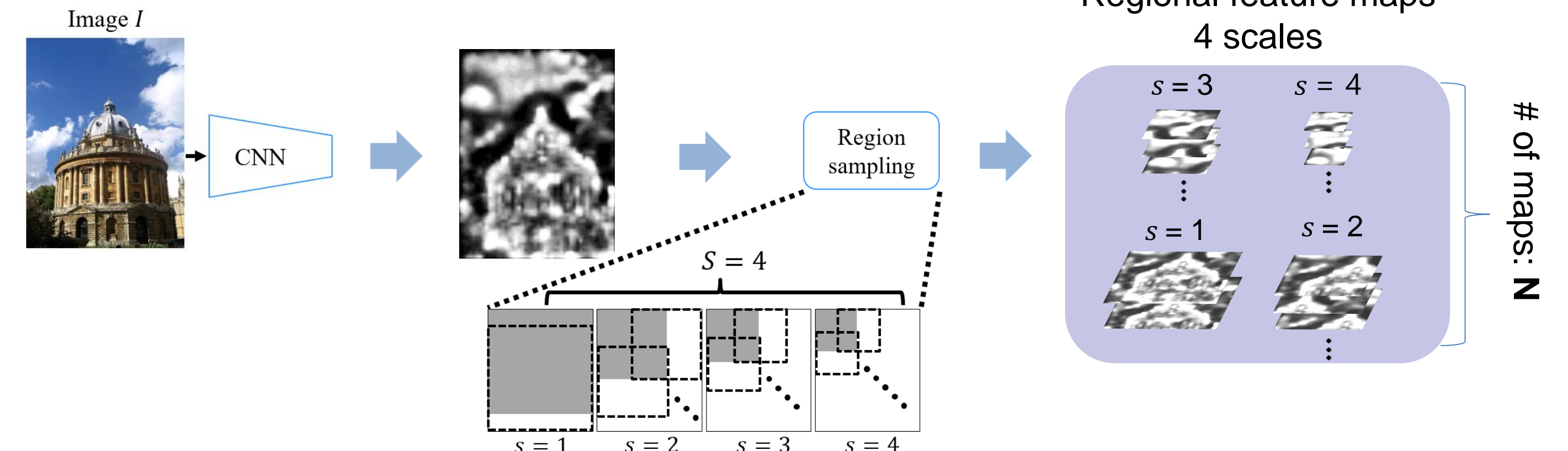
- Parameters to train: $\mathbf{W}_r, \mathbf{b}_r, \mathbf{W}_c, \mathbf{b}_c$
- Freezing the CNN(Resnet101) while training our attention network
- Dataset: ILSVRC2012-ImageNet for "Pre-trained" category
- Classification loss(Cross entropy)

2. Image Encoding Pipeline

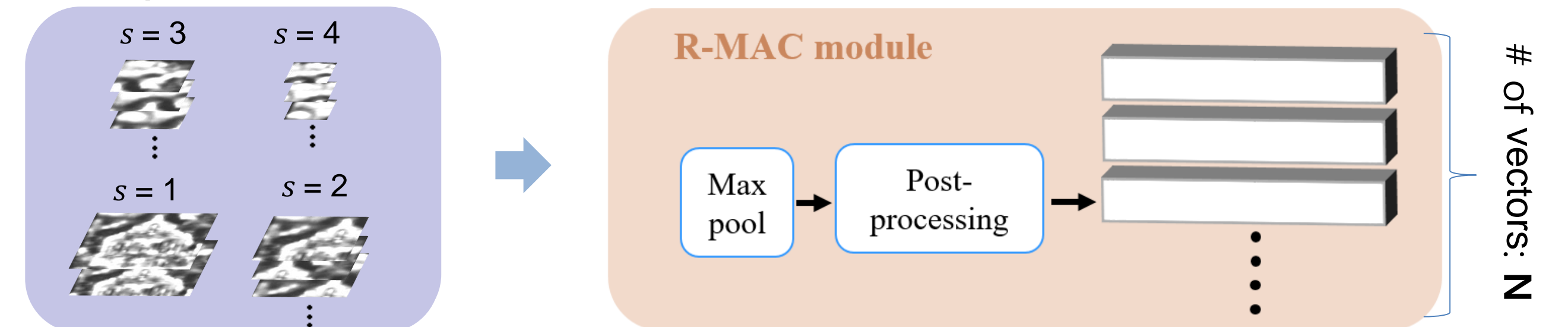
Overall Sequence of encoding pipeline



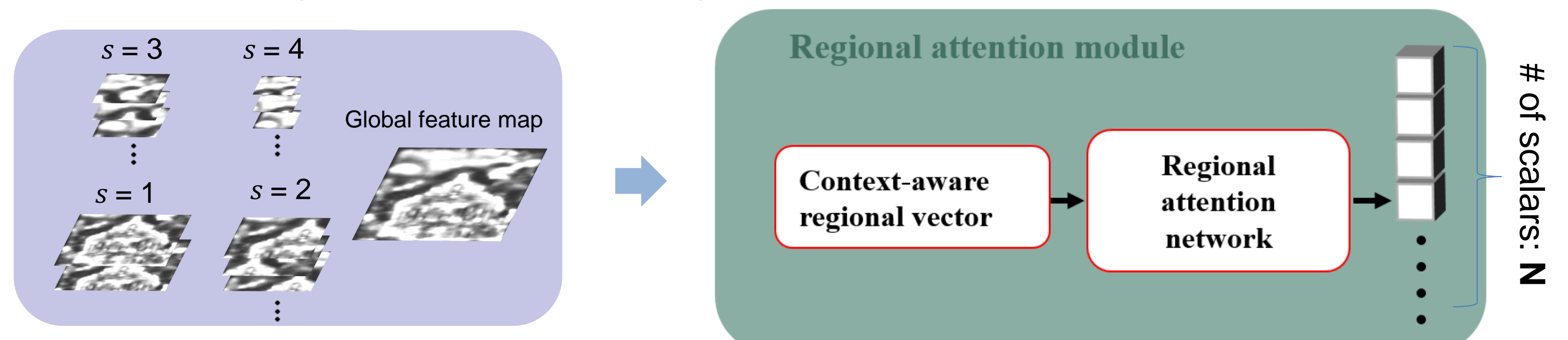
1. Extract a CNN feature map and sample regional feature maps in a R-MAC manner



2.1 Produce R-MAC feature vectors with the regional feature maps.



2.2 Calculate regional attention weights $\Phi(\mathbf{k})$



3. Obtain a global feature vector, $\hat{\mathbf{f}}_I$, through combining R-MAC features with regional attention weights

$$\text{Weighted mean} \rightarrow \ell_2 \rightarrow \hat{\mathbf{f}}_I = [\hat{f}_{I,1}, \dots, \hat{f}_{I,k}]^T$$

4. Experiments

Method	Scale (S)				Method	Oxford5k	Paris6k
	S=3	S=4	S=5	S=6			
Baseline	69.9	70.7	70.1	69.0	RPN + PCA Landmark	64.7	75.5
Ours	75.1	76.7	76.8	76.4	+ Regional attention	66.6	75.8
					+ Context awareness	67.9	76.4

Ablation study

Method	(a)			(b)		
	Oxford5k	Paris6k	Time (s)	Oxford5k	Paris6k	Time (s)
Baseline + PCA Landmark	70.1	85.4	0.095	70.1	85.4	0.095
+ Regional attention	74.9	86.0	0.115	74.9	86.0	0.115
+ Context awareness	76.8	87.5	0.123	76.8	87.5	0.123
- PCA Landmark	77.6	88.3	-	77.6	88.3	-
+ PCA Paris, Oxford						

Comparison with state-of-the-arts

Method	Dim.	Query expansion (QE)			
		Oxford5k	Paris6k	Oxford105k	Paris106k
Resnet101	SDCF [10]	69.1	81.7	65.4	74.3
	CroW [11]	68.7	82.8	62.7	75.1
	R-MAC [12]	70.1	85.4	66.9	80.8
	CAM [13]	69.9	84.3	64.3	77.1
	Ours	2048	76.8	87.5	73.6
Resnet101	SDCF+QE [10]	68.5	84.9	66.8	79.4
	CroW+QE [11]	69.5	85.1	66.7	79.9
	R-MAC+QE [12]	73.8	86.4	71.8	82.6
	CAM+QE [13]	71.3	86.1	68.7	80.8
	Ours+QE	2048	81.8	89.3	80.4

