



Innovative R&D by NTT

# Weighted Generalized Mean Pooling for Deep Image Retrieval

**Xiaomeng Wu**, Go Irie, Hiramatsu Kaoru, and Kunio Kashino

NTT Communication Science Laboratories

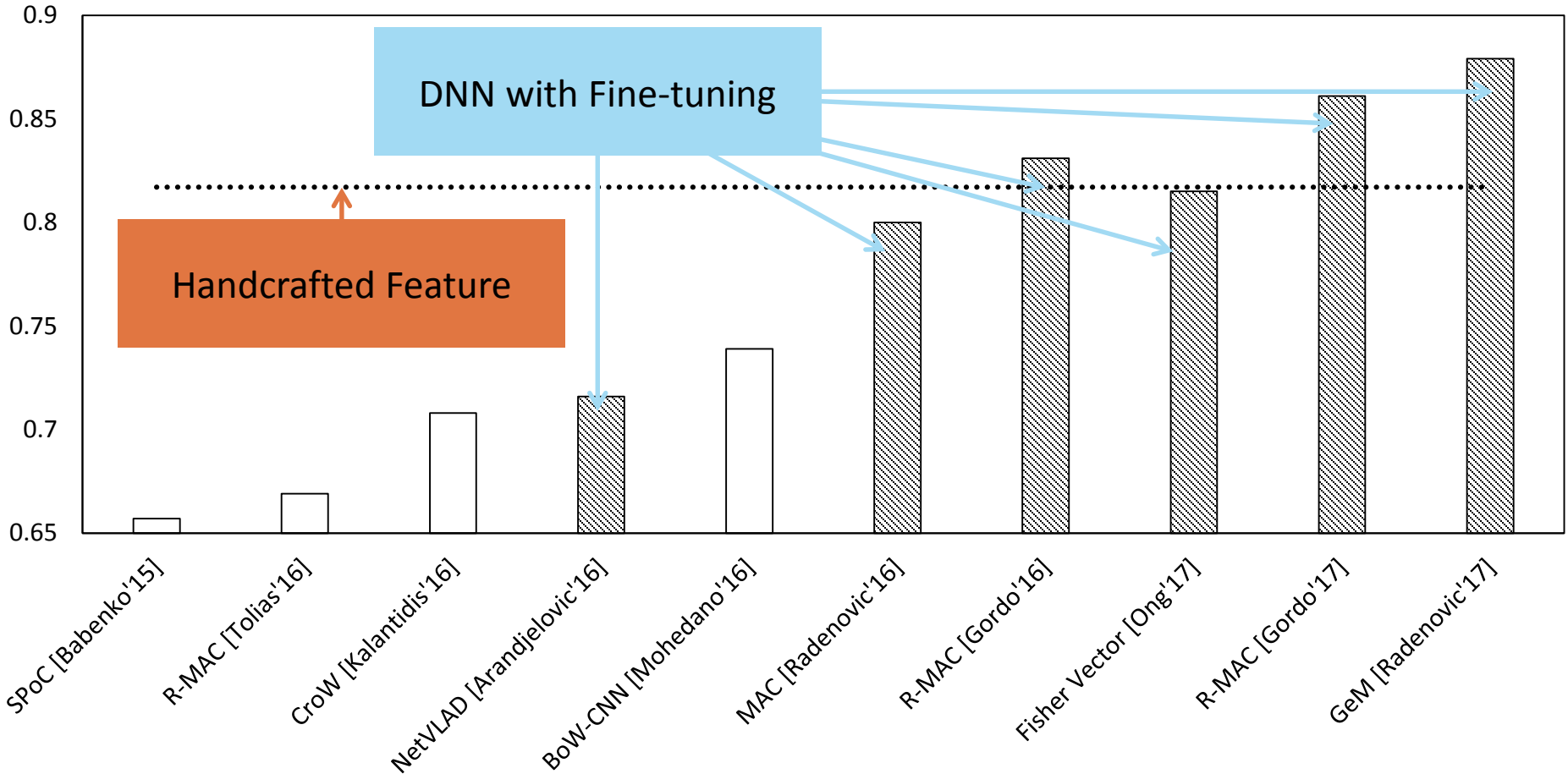
NTT Corporation, Japan

# Fine-tuning for Deep Image Retrieval



MAPs on Oxford5K w/o Query Expansion

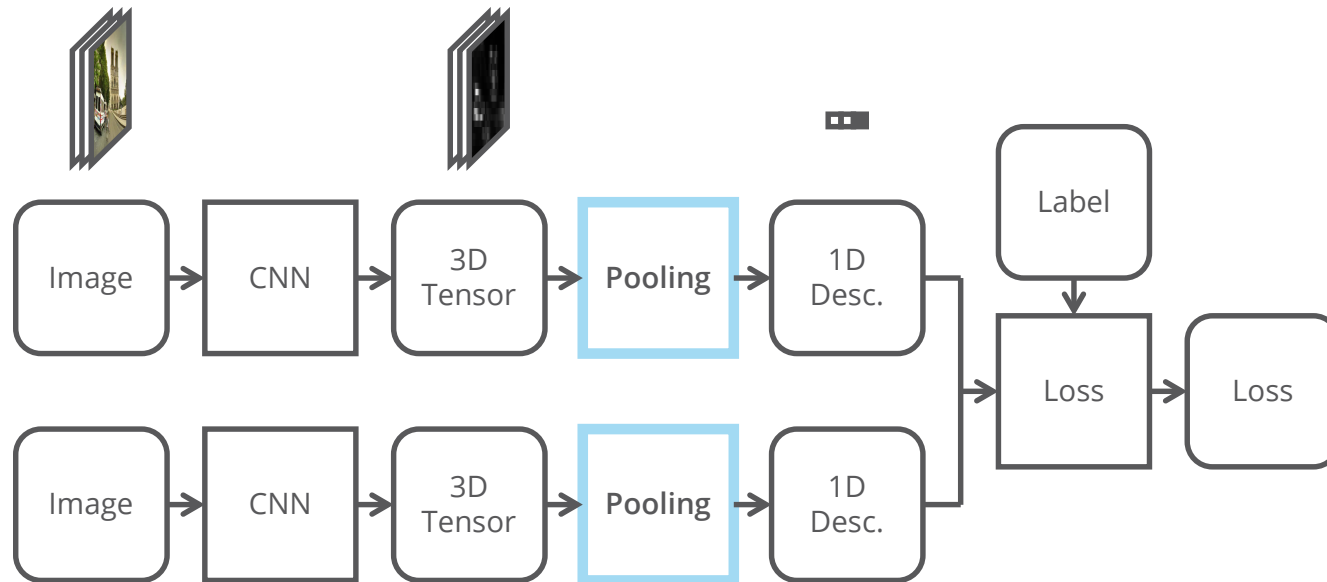
□ Deep Image Retrieval    ..... ASMK [Tolias'13]



# Siamese Network



- Crucial components of deep image retrieval include
  - A good pre-trained convolutional network
  - A good pooling method
  - A ranking loss





- **FC Layers** [Babenko'14][Gong'14]
- **Global Pooling**
  - Sum [Babenko'15]
  - Maximum [Radenovic'16][Tolias'16][Razavian'16]
  - Generalized Mean (GeM) [Radenovic'17]
- **Semi-local Regional Pooling**
  - R-MAC [Tolias'16]
  - Region Proposal Network (RPN) [Gordo'17]
- **Widely-used Encoding Techniques**
  - Bag of Visual Words [Mohedano'16]
  - VLAD [Arandjelovic'16]
  - Fisher Vector [Ong'17]

## UNIFORM POOLING

Each activation contributes equivalently to the construction of a global representation.

## PROBLEM

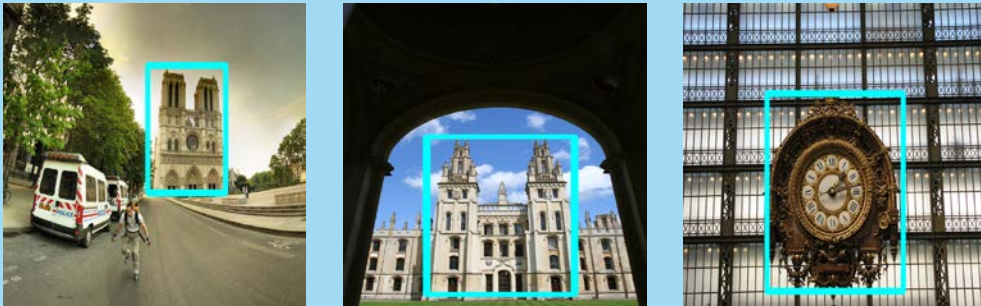
Uniform pooling suffers from the presence of activations, e.g., from background clutter, that play a negative role as regards matching.

# CNN Activations vs. Objectness



- CNN Firing to Background Clutter

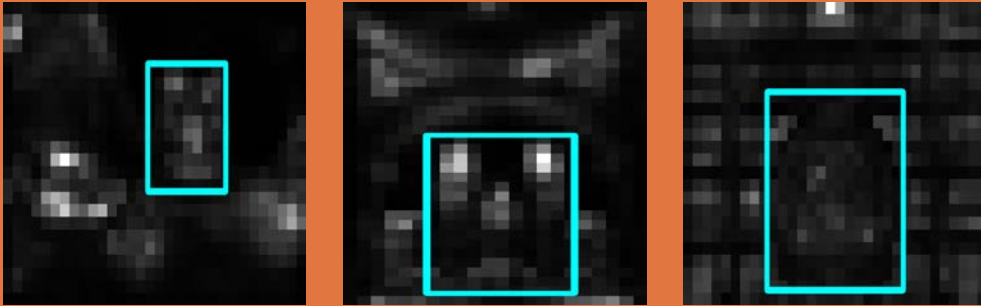
Images from Oxford5K



Objects of interest (landmarks) predefined in Oxford5K

The CNN has been fine-tuned on landmark images.

Activation Strengths



# Goal of the Study

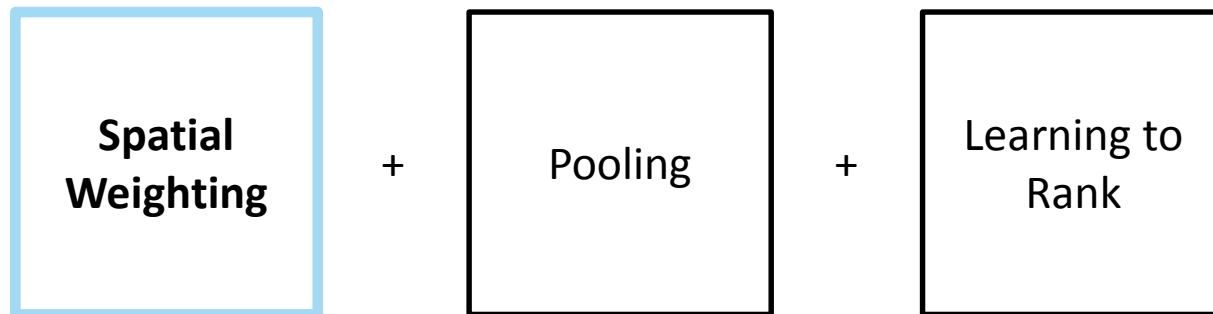


- **Problem of Uniform Pooling**

- Each activation contributes equivalently to the construction of a global representation.

- **Proposal**

- Exploit a spatial weighting mechanism for pooling
  - Predict a weight that describes how discriminating each activation at each location is as regards image matching.
- Lead to the end-to-end learning based on a weighted generalized mean (wGeM) pooling method





Innovative R&D by NTT

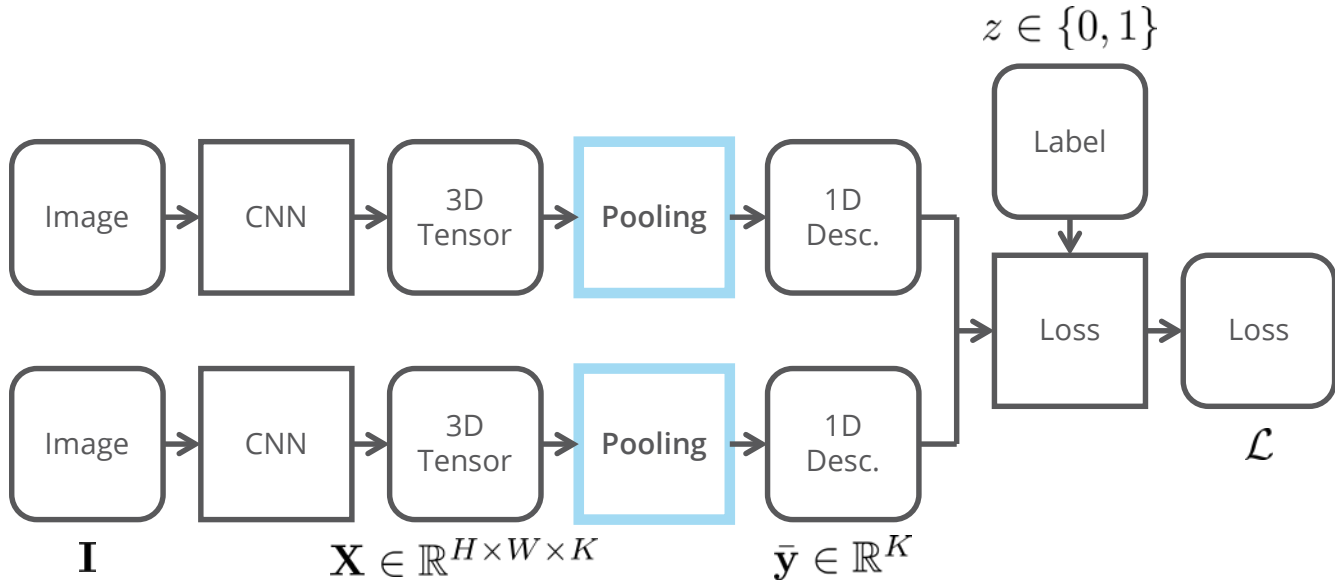
# Proposed Method

# Network Architecture and Learning



- **Contrastive Loss** [Chopra'05][Radenovic'16]

$$\mathcal{L}_{i,j} = \begin{cases} \frac{1}{2} \|\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_j\|_2^2 & \text{if } z_{i,j} = 1 \\ \frac{1}{2} [\max(0, \tau - \|\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_j\|_2)]^2 & \text{otherwise} \end{cases}$$

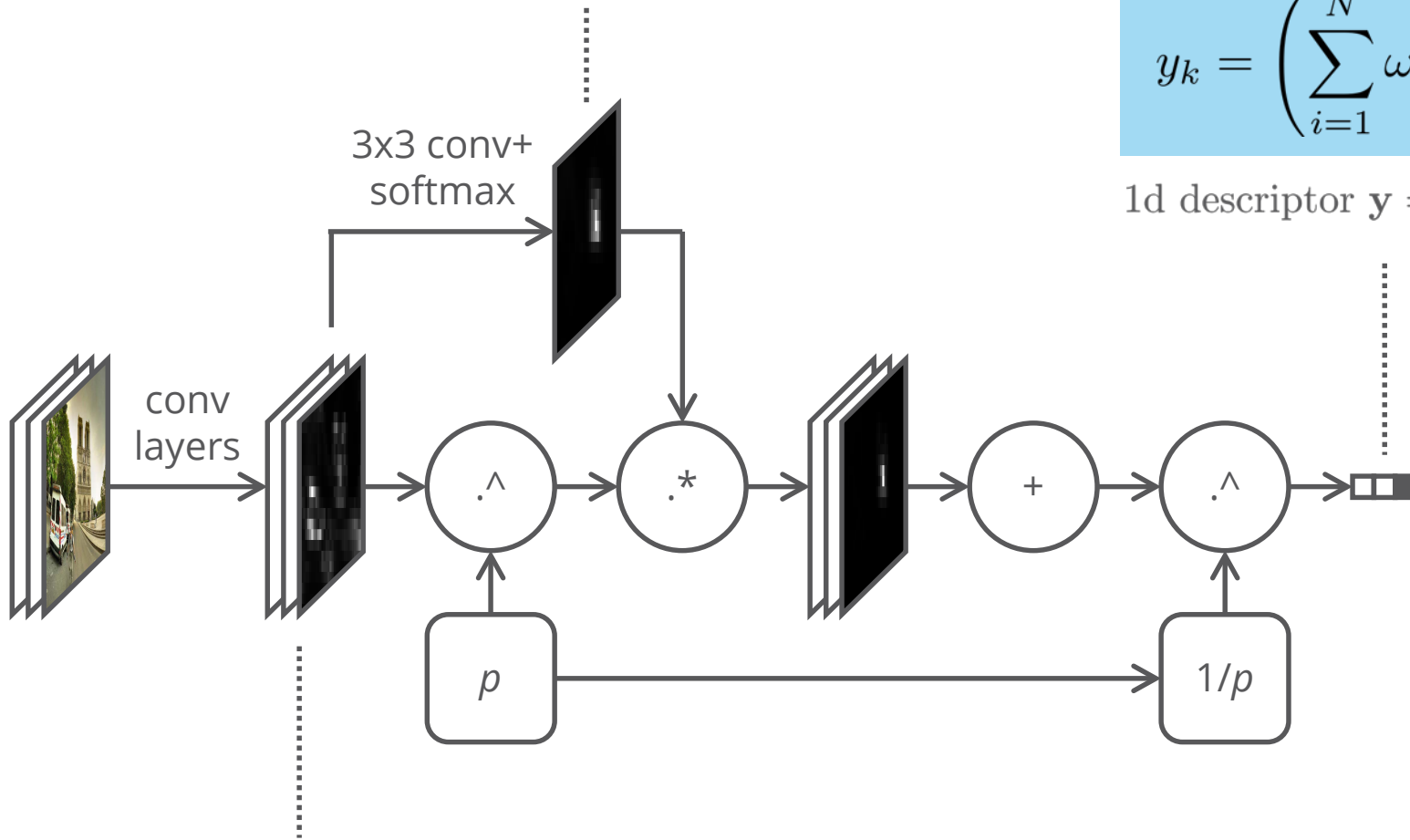




# Weighted Generalized Mean (wGeM) Pooling



$$2d \text{ mask } \Omega = [\omega_1 \dots \omega_N]^\top$$



$$3d \text{ tensor } \mathbf{X} = \{\mathbf{x}_1 \dots \mathbf{x}_K\}$$

$$2d \text{ activation map } \mathbf{x}_k = [x_{1,k} \dots x_{N,k}]^\top$$

## Forward Propagation

$$y_k = \left( \sum_{i=1}^N \omega_i x_{i,k}^p \right)^{1/p}$$

$$1d \text{ descriptor } \mathbf{y} = [y_1 \dots y_K]^\top$$

# Derivatives and Property



- **Derivatives**

$$\frac{\partial y_k}{\partial x_{i,k}} = \sum_{j=1}^N \frac{\partial y_k}{\partial \omega_j} \frac{\partial \omega_j}{\partial x_{i,k}} + \omega_i \left( \frac{x_{i,k}}{y_k} \right)^{p-1}$$

$$\frac{\partial y_k}{\partial \omega_i} = \frac{x_{i,k}}{p} \left( \frac{x_{i,k}}{y_k} \right)^{p-1}$$

$$\frac{\partial y_k}{\partial p} = \frac{y_k}{p} \left( \frac{\sum_i \omega_i x_{i,k}^p \log x_{i,k}}{y_k^p} - \log y_k \right)$$

- **Behavior when  $p \rightarrow \infty$**

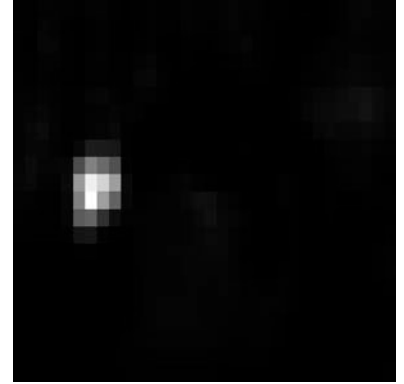
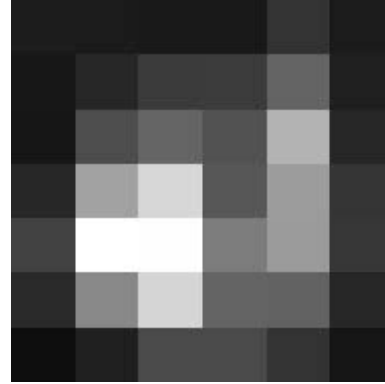
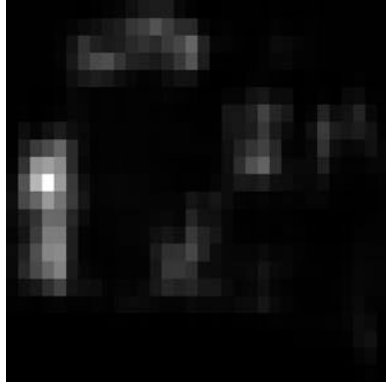
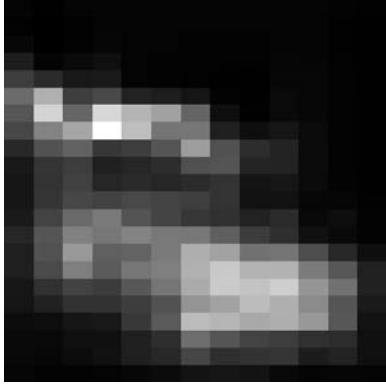
- Forward

$$\lim_{p \rightarrow \infty} y_k = \lim_{p \rightarrow \infty} \left( \sum_{i=1}^N \omega_i x_{i,k}^p \right)^{1/p} = \max_i x_{i,k}$$

- Back-propagation

$$\lim_{p \rightarrow \infty} \frac{\partial \mathcal{L}}{\partial \omega_i} = \lim_{p \rightarrow \infty} \frac{1}{p} \sum_{k=1}^K \frac{\partial \mathcal{L}}{\partial y_k} \left( \frac{x_{i,k}}{y_k} \right)^{p-1} x_{i,k} = 0$$

# Matching Images and Spatial Weights





Innovative R&D by NTT

# Experiments

# Experimental Setup



<b>Dataset for Training</b>	36K pairs out of 163K+ landmark images [Radenovic'16]
<b>Dataset for Testing</b>	Oxford5K, Paris6K, & Oxford105K
<b>Backbone Network</b>	ResNet101 [He'16]
<b>#epoch</b>	30
<b>Batch Size</b>	30 pairs (#pos./#neg.: 5/25)
<b>Learning Rate</b>	$10^{-6}e^{-0.1i}$ over epoch $i$
<b>Margin of Contrastive Loss</b>	0.85
<b>Learning Method</b>	Adam [Kingma'14]
<b>Pre- and Post-processing</b>	Multi-scale representation & supervised whitening [Radenovic'18]
<b>Performance Measure</b>	Mean Average Precision (MAP)

## ▪ Implementation of wGeM

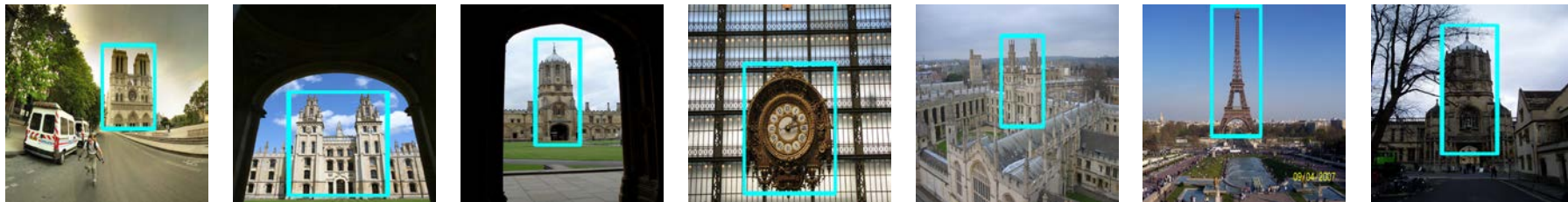
- Initialize the parameters of the  $3 \times 3$  conv layer with 0s such that  $\forall \omega_i \in \Omega, \omega_i = \frac{1}{N}$
- Use learning rates that are 10 times as large as those of the pre-trained ResNet

# MAPs for Different Initializations of $p$

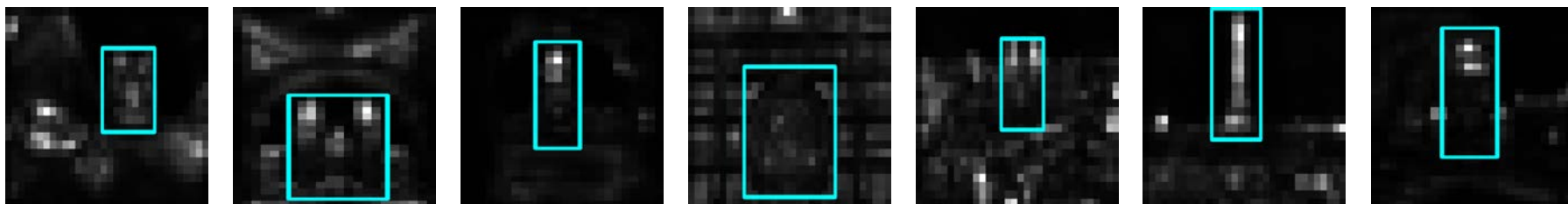


	Initial $p$	Oxford5K	Paris6K
<b>GeM</b> [Radenovic'18]	3	87.8	92.7
<b>wGeM</b>	2	88.6	92.2
	3	88.9	92.3
	4	88.8	92.5
	5	88.4	92.6

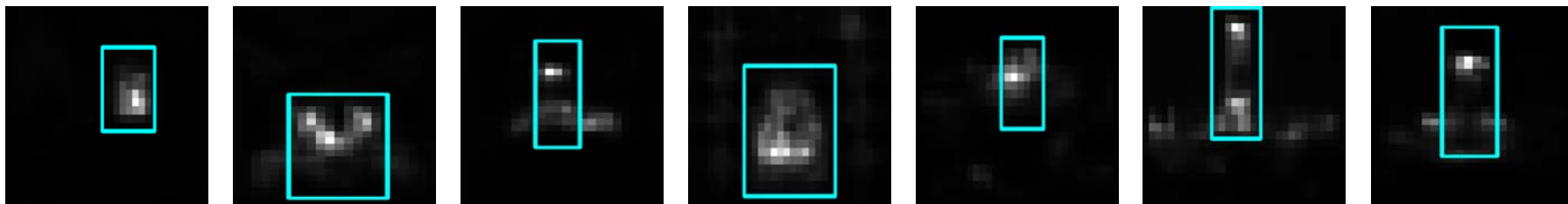
The MAP of the wGeM was slightly poorer for Paris6K because the training set is composed of more historical landmarks, while Paris6K contains more contemporary buildings.



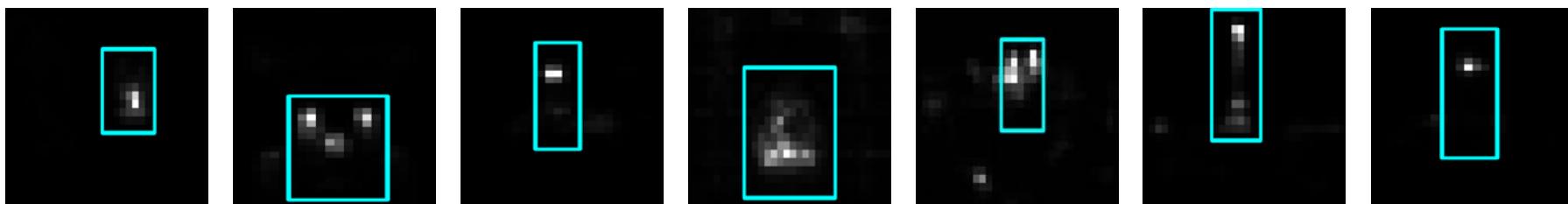
Images from Oxford5K



Activation Strengths w/o Spatial Weighting



Spatial Weights



Activation Strengths w/ Spatial Weighting

# Comparison with State of the Art

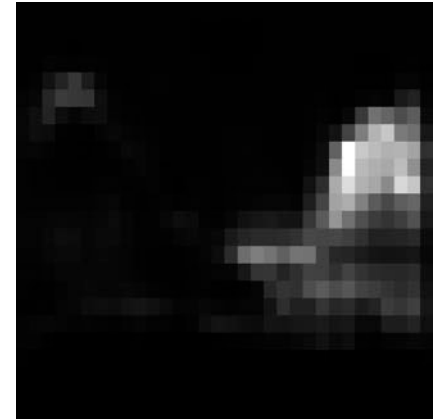
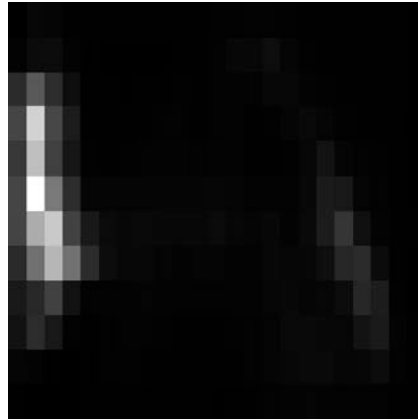


	CNN	QE	Oxford5K	Oxford105K	Paris6K
<b>NetVLAD</b> [Arandjelovic'16]	VGG	n	71.6	n/a	79.7
<b>MAC</b> [Radenovic'16]	VGG	n	80.0	75.1	82.9
<b>Fisher Vector</b> [Ong'17]	VGG	n	81.5	76.6	82.4
<b>R-MAC</b> [Gordo'17]	ResNet	n	86.1	82.8	<b>94.5</b>
<b>GeM</b> [Radenovic'18]	ResNet	n	87.8	84.6	92.7
<b>* wGeM</b>	ResNet	n	<b>88.8</b>	<b>85.6</b>	92.5
<b>MAC+QE</b> [Radenovic'16]	VGG	y	85.4	82.3	87.0
<b>R-MAC+QE</b> [Gordo'17]	ResNet	y	90.6	89.4	<b>96.0</b>
<b>GeM+<math>\alpha</math>QE</b> [Radenovic'18]	ResNet	y	91.0	89.5	95.5
<b>* wGeM+QE</b>	ResNet	y	<b>91.7</b>	<b>89.7</b>	<b>96.0</b>

The wGeM outperforms or is on a par with the state of the art based on fine-tuned deep networks.



# Spatial Weighting Failure



Here, the object of interest, the Louvre Pyramid, was largely ignored by the wGeM because regular patterns are potentially less discriminating when matching the landmarks in the training set. Therefore, the wGeM was trained to fire less for such regions.



Innovative R&D by NTT

# Conclusions

## ▪ **Characteristics of Proposed Method**

- Generalize sum, max, GeM pooling, and CroW [Kalantidis'16].
- Trainable
- Require no bounding box annotations for training

## ▪ **Future Directions**

- Complementarity with state-of-the-art QE techniques, e.g., diffusion [Iscen'17] [Wu'18].
- How do different wGeM block structures affect the learning of deep representations?



Innovative R&D by NTT

# Thank you for your attention.

**Xiaomeng Wu**, Go Irie, Hiramatsu Kaoru, and Kunio Kashino

NTT Communication Science Laboratories

NTT Corporation, Japan