



LONG BEACH  
CALIFORNIA  
June 16-20, 2019

# Deep Metric Learning Beyond Binary Supervision

Sungyeon Kim   Minkyoo Seo   Ivan Laptev   Minsu Cho   Suha Kwak

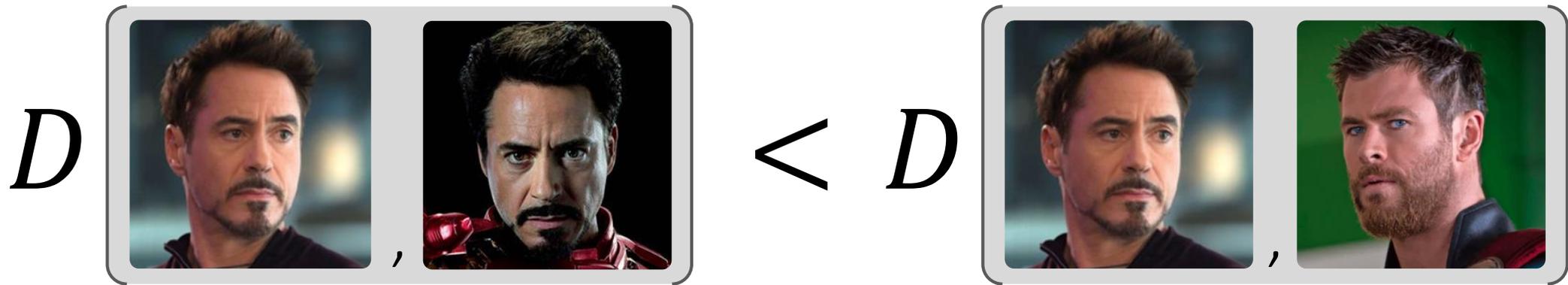
{tjddus9597, mkseo, mscho, suha.kwak}@postech.ac.kr, ivan.laptev@inria.fr

**POSTECH**

*Inria*

# Metric Learning

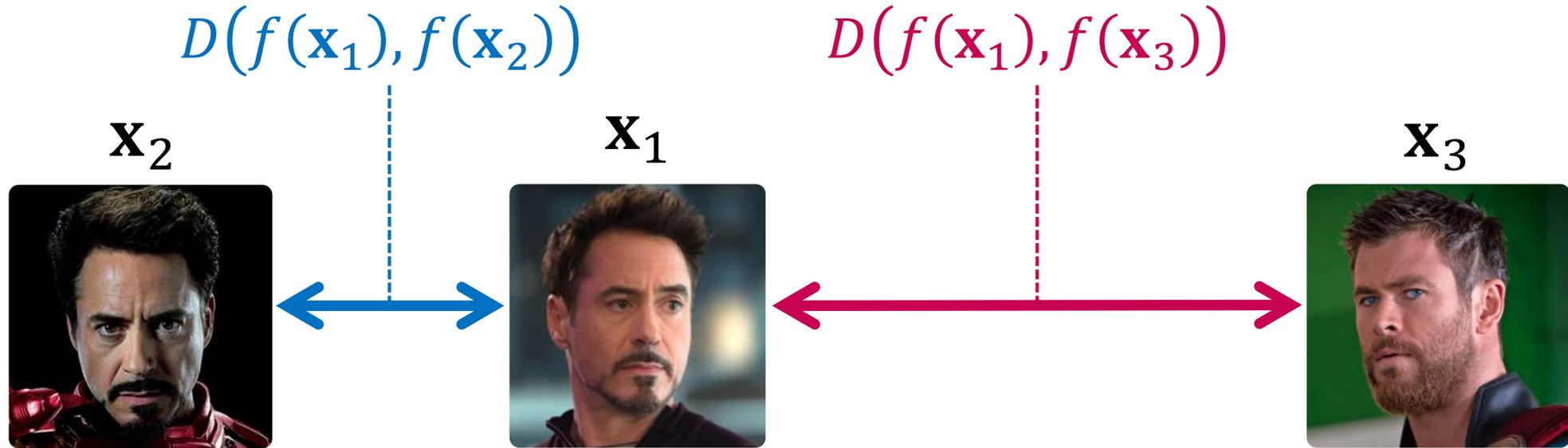
How much similar/dissimilar?



**Metric:** Function that quantifies a **distance**

**Metric Learning:** Learning a metric from a set of data

# Deep Metric Learning



Pairwise relation  
 $D(f_1, f_2) \downarrow, D(f_1, f_3) \uparrow$

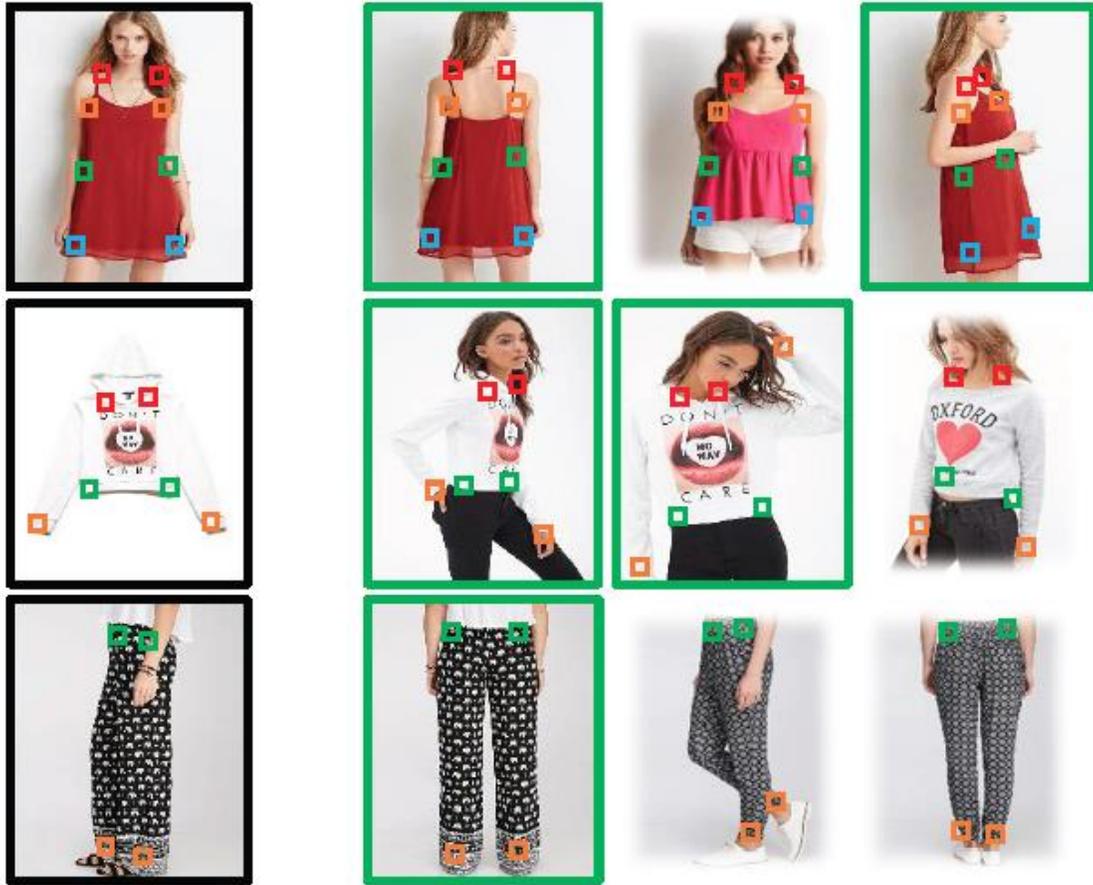
Triplet relation  
 $D(f_1, f_2) < D(f_1, f_3)$

...

## Deep Metric Learning

Learning a deep neural net  $f$  that satisfies the relations

# Applications



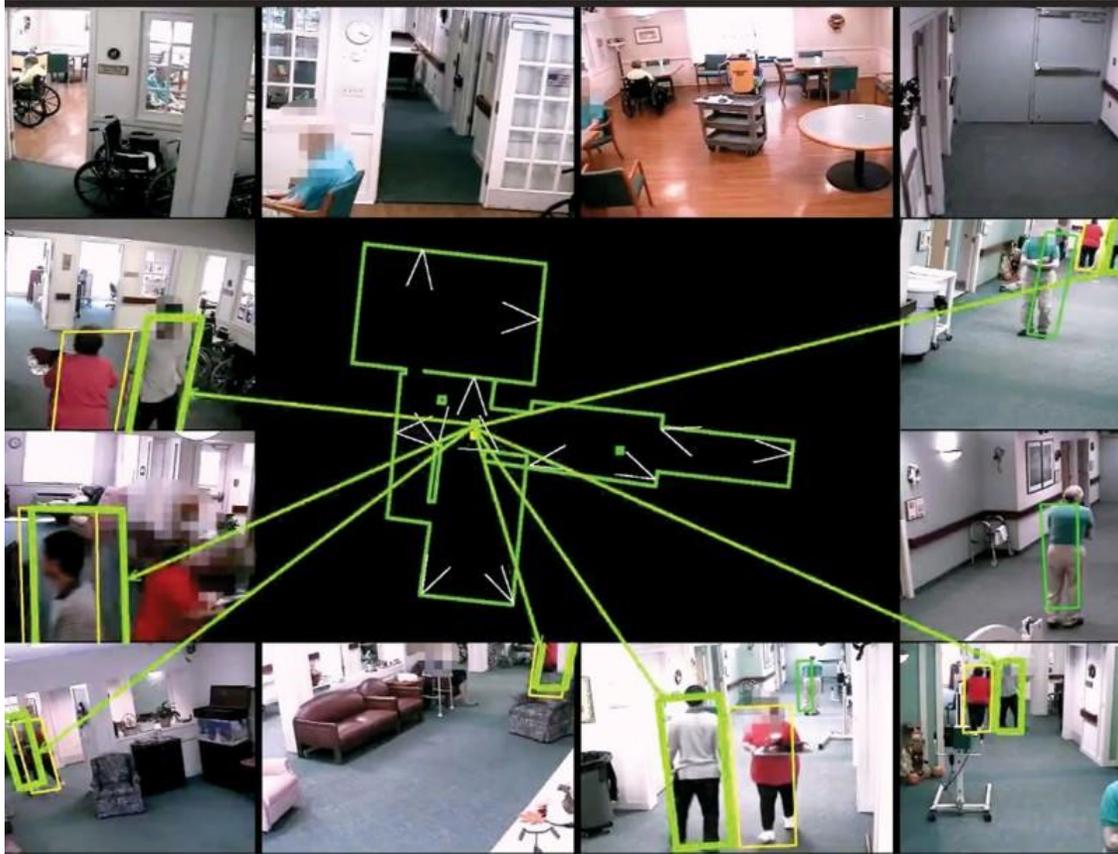
Content-based image retrieval



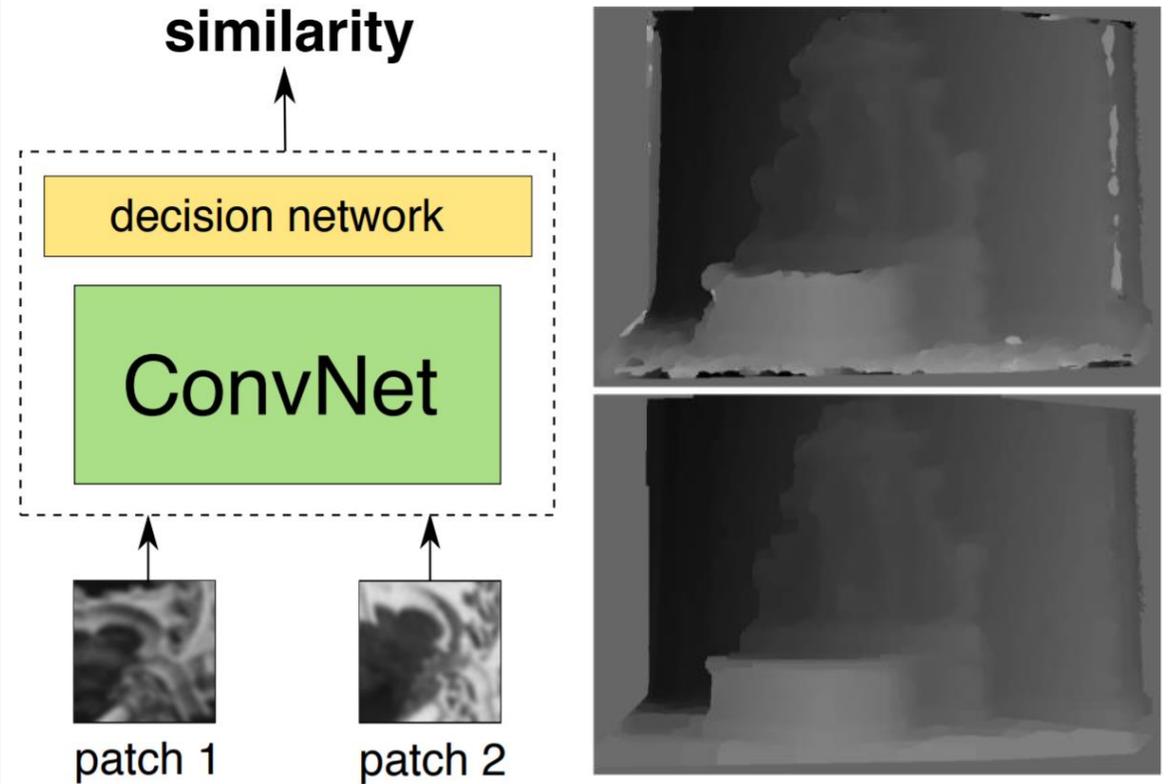
Face verification/identification<sup>[1]</sup>

[1] FaceNet: A unified embedding for face recognition and clustering, CVPR 2015

# Applications



Person re-identification<sup>[2]</sup>



Patch matching/stereo imaging<sup>[3]</sup>

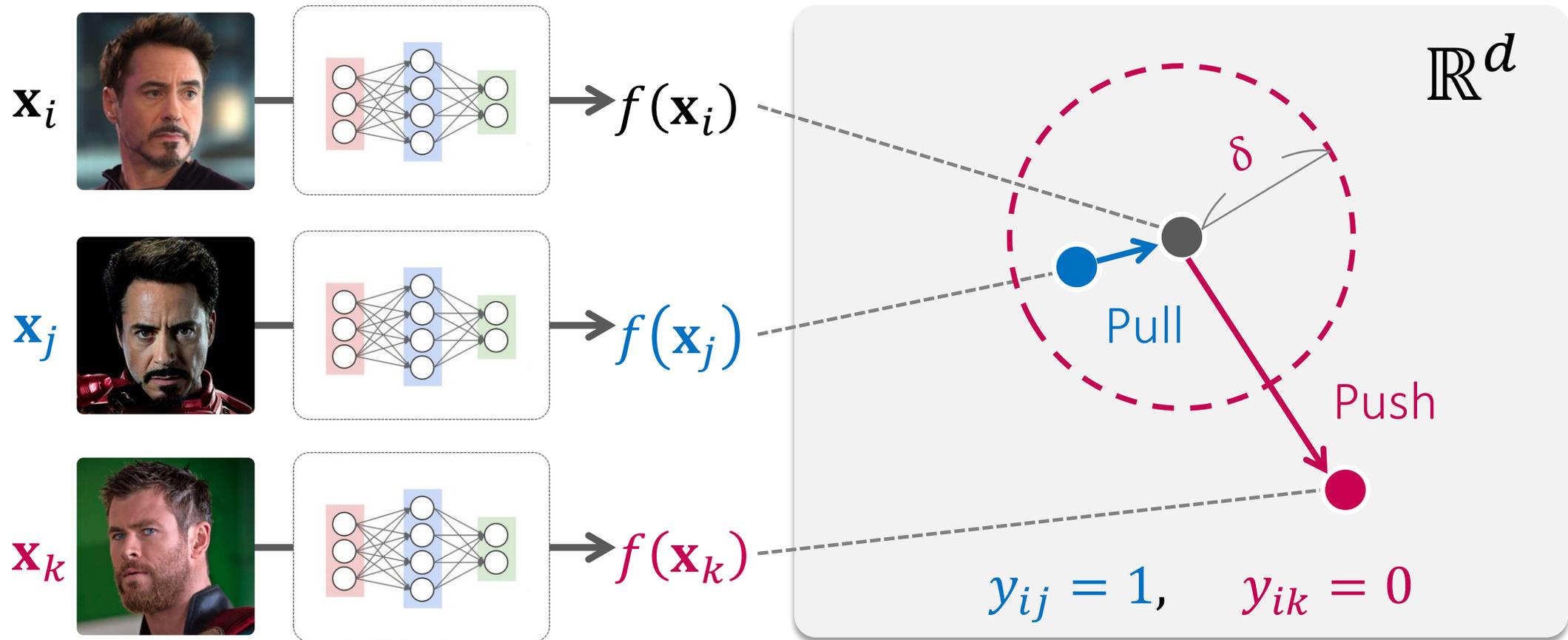
[2] Beyond triplet loss: a deep quadruplet network for person re-identification, CVPR 2017

[3] Learning to compare image patches via convolutional neural networks, CVPR 2015

# Existing Approaches

- Contrastive loss for Siamese networks<sup>[4]</sup>

$$\ell_{\text{ctr}}(i, j) = y_{ij}D(f_i, f_j)^2 + (1 - y_{ij})[\delta - D(f_i, f_j)]_+^2$$

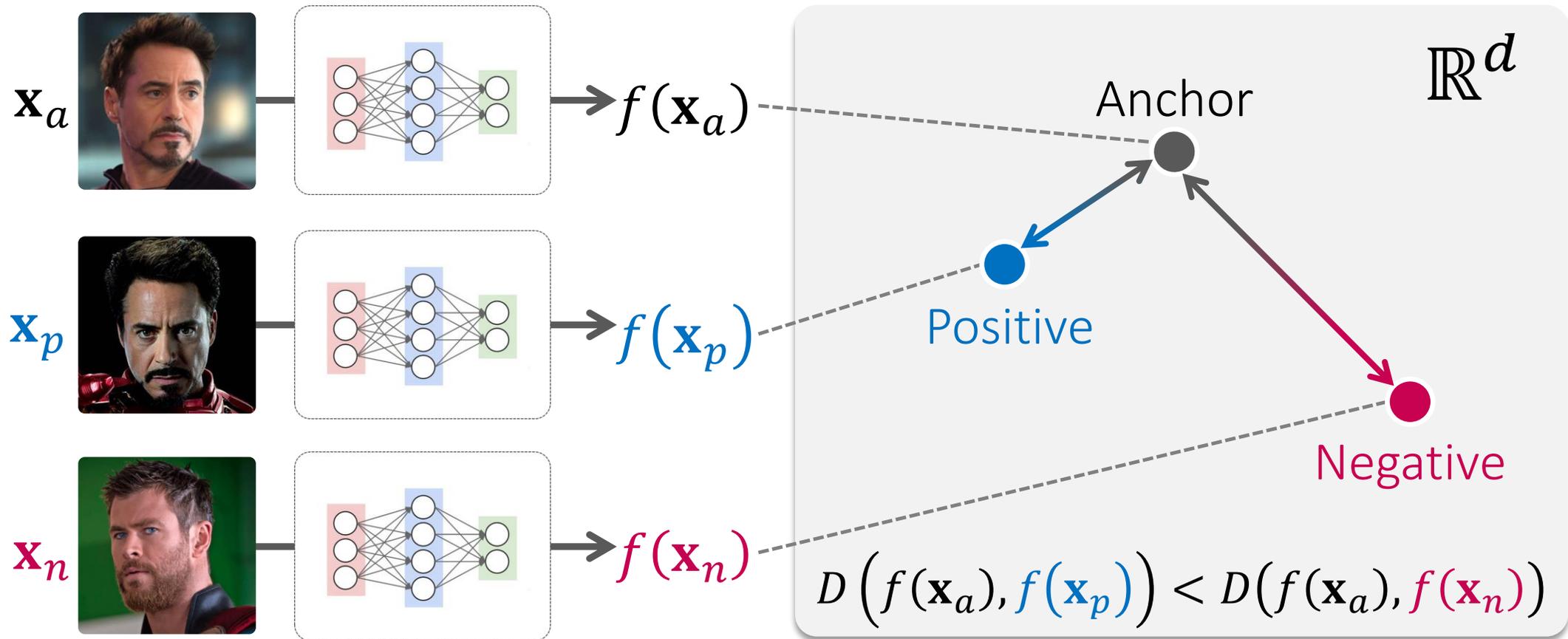


[4] Learning a similarity metric discriminatively with application to face verification, CVPR 2005

# Existing Approaches

- Triplet rank loss for triplet networks<sup>[1]</sup>

$$\ell_{\text{tri}}(a, p, n) = [D(f_a, f_p) - D(f_a, f_n) + \delta]_+$$



[1] FaceNet: A unified embedding for face recognition and clustering, CVPR 2015

# Existing Approaches

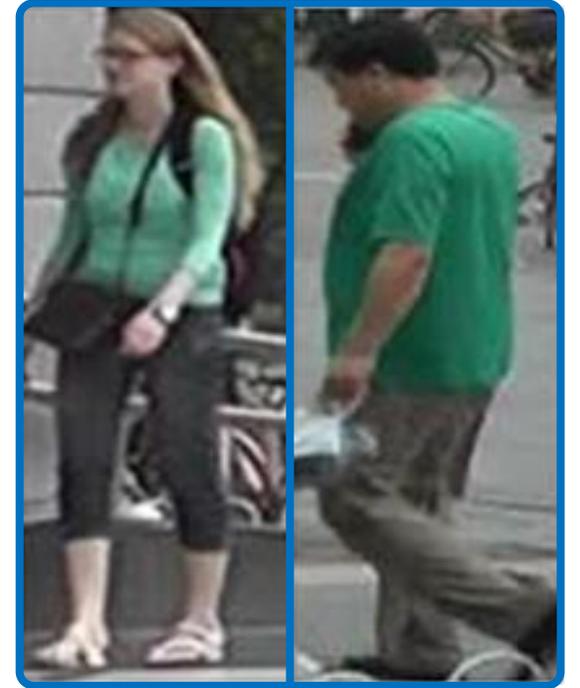
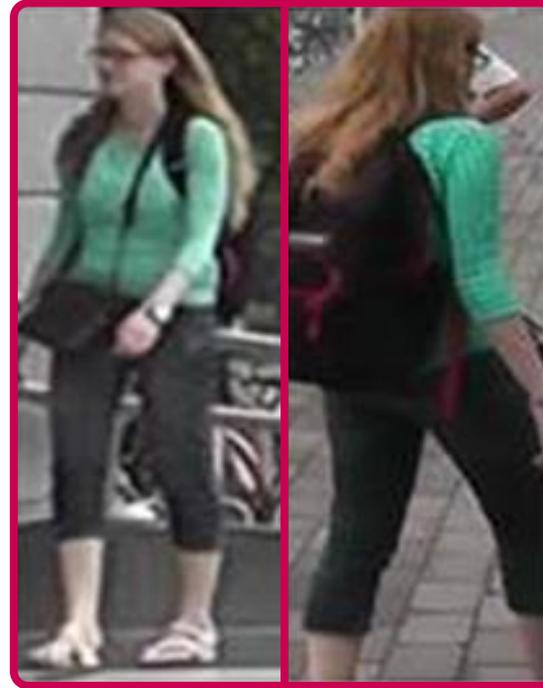
- A common issue
  - Existing (deep) metric learning approaches rely on binary relations between images: “*same*” or “*not*”.



Face verification



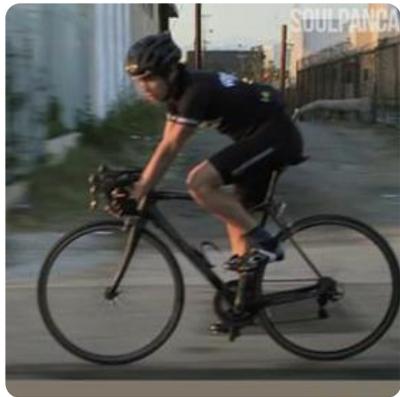
Content-based image retrieval



Person re-identification

# Existing Approaches

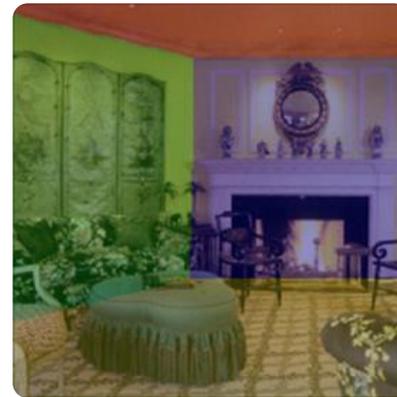
- A common issue
  - However, relations between real world images are *not binary* but often represented as *continuous similarities*.



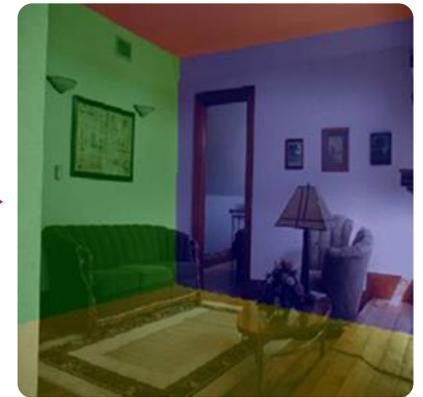
1.65



1.47



0.26



0.41



2.86

3.41

0.34

0.29

# Existing Approaches

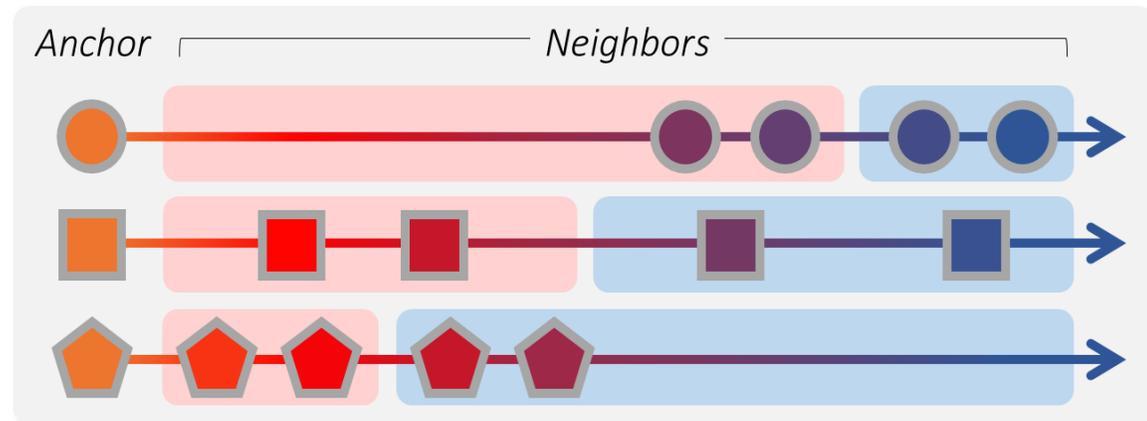
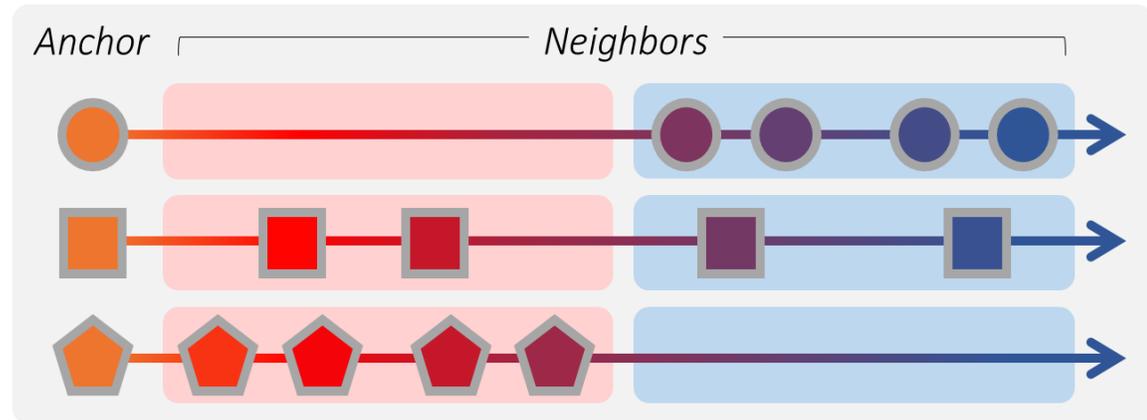
- Conventional approaches to handle the issue
  - Existing metric learning loss + *similarity quantization*

## *Binary thresholding*<sup>[5]</sup>

Populations of positive and negative examples would be significantly imbalanced.

## *Nearest neighbor search*<sup>[6]</sup>

Positive neighbors of a rare example would be dissimilar and negative neighbors of a common example would be too similar.



[5] Pose embeddings: A deep architecture for learning to match human poses, arXiv 2015

[6] Thin-slicing for pose: Learning to understand pose without explicit pose estimation, CVPR 2016

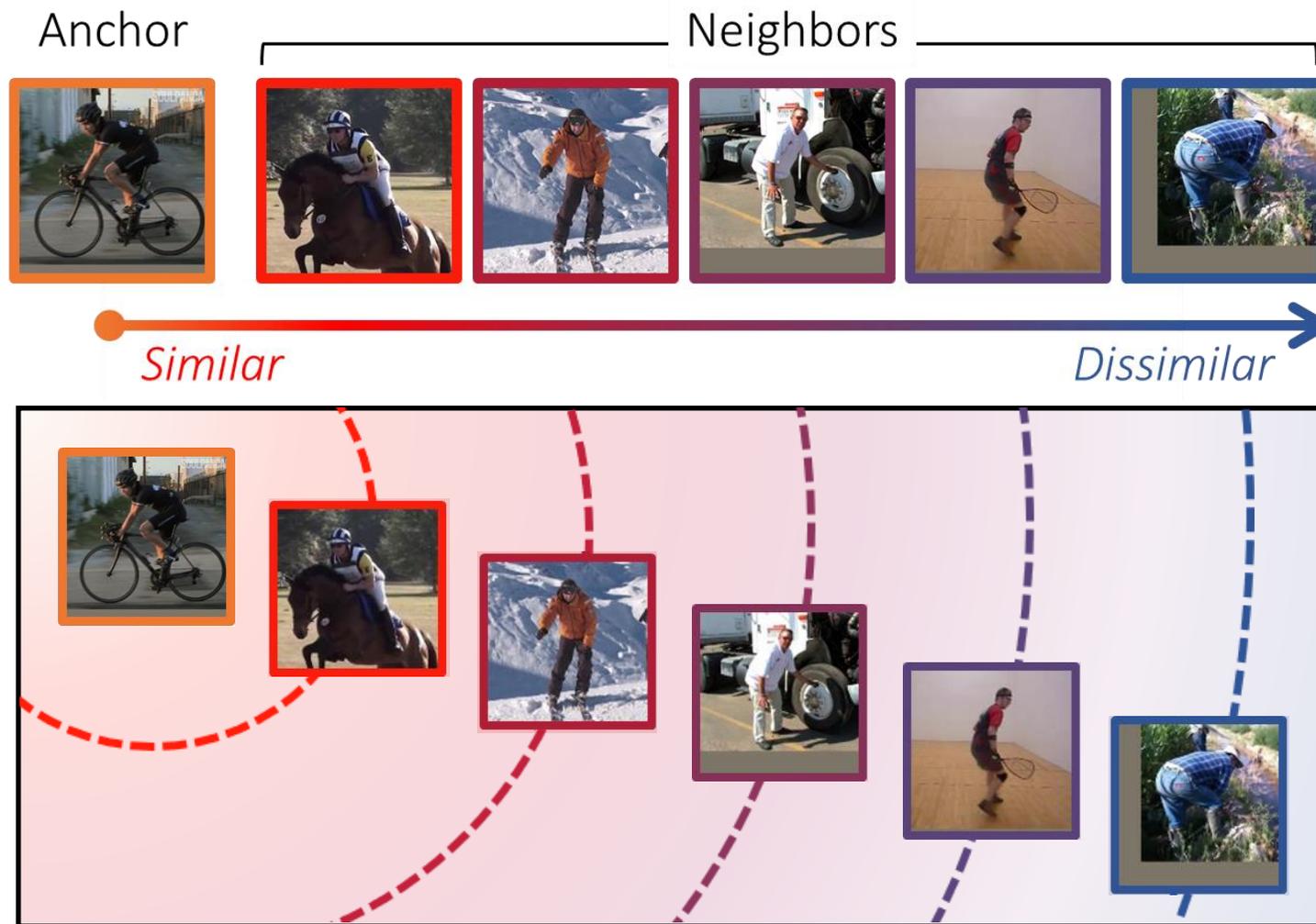
# Existing Approaches

- Conventional approaches to handle the issue
  - Degree of similarity is ignored in the learned embedding space.



# Our Approach

- Our goal
  - Learning a metric space that reflects the degree of similarity directly



# Our Approach

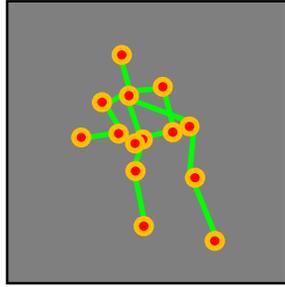
- Our goal
  - Learning a metric space that reflects the degree of similarity directly
- Contributions
  - A new triplet loss: *Log-ratio loss*
  - A new triplet sampling technique: *Dense triplet sampling*
  - Various applications
    - Human pose retrieval
    - Room layout retrieval
    - Caption-aware image retrieval
    - Representation learning for image captioning

# Log-ratio Loss

- Definition



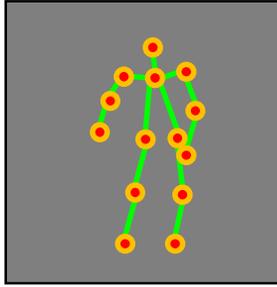
$\mathbf{x}_a$



$\mathbf{y}_a$



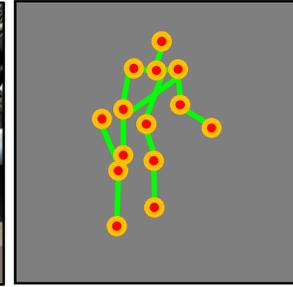
$\mathbf{x}_i$



$\mathbf{y}_i$



$\mathbf{x}_j$



$\mathbf{y}_j$

$$\ell_{\text{lr}}(a, i, j) = \left\{ \log \frac{D(f_a, f_i)}{D(f_a, f_j)} - \log \frac{D_y(\mathbf{y}_a, \mathbf{y}_i)}{D_y(\mathbf{y}_a, \mathbf{y}_j)} \right\}^2$$

where  $f_i := f(\mathbf{x}_i)$  is the embedding vector of image  $i$ ,  
and  $D(\cdot)$  denotes the squared Euclidean distance.

The distance between two images in the learned metric space  
will be proportional to **their distance in the label space**.

# Log-ratio Loss

- Analysis on its gradients

$$\frac{\partial \ell_{\text{lr}}(a, i, j)}{\partial f_a} = - \frac{\partial \ell_{\text{lr}}(a, i, j)}{\partial f_i} - \frac{\partial \ell_{\text{lr}}(a, i, j)}{\partial f_j}$$

$$\frac{\partial \ell_{\text{lr}}(a, i, j)}{\partial f_i} = \frac{(f_i - f_a)}{D(f_a, f_i)} \cdot \ell'_{\text{lr}}(a, i, j)$$

$$\frac{\partial \ell_{\text{lr}}(a, i, j)}{\partial f_j} = \frac{(f_a - f_j)}{D(f_a, f_j)} \cdot \ell'_{\text{lr}}(a, i, j)$$

Direction between  
the anchor and neighbors

Discrepancy between  
the label distance ratio and  
the embedding distance ratio

$$4 \left\{ \log \frac{D(f_a, f_i)}{D(f_a, f_j)} - \log \frac{D_{\mathbf{y}}(\mathbf{y}_a, \mathbf{y}_i)}{D_{\mathbf{y}}(\mathbf{y}_a, \mathbf{y}_j)} \right\}$$

# Log-ratio Loss

- Comparison to the triplet rank loss

*Log-ratio loss*

$$\ell_{\text{lr}}(a, i, j) = \left\{ \log \frac{D(f_a, f_i)}{D(f_a, f_j)} - \log \frac{D(y_a, y_i)}{D(y_a, y_j)} \right\}^2$$

$$\frac{\partial \ell_{\text{lr}}(a, i, j)}{\partial f_a} = - \frac{\partial \ell_{\text{lr}}(a, i, j)}{\partial f_i} - \frac{\partial \ell_{\text{lr}}(a, i, j)}{\partial f_j}$$

$$\frac{\partial \ell_{\text{lr}}(a, i, j)}{\partial f_i} = \frac{(f_i - f_a)}{D(f_a, f_i)} \cdot \ell'_{\text{lr}}(a, i, j)$$

$$\frac{\partial \ell_{\text{lr}}(a, i, j)}{\partial f_j} = \frac{(f_a - f_j)}{D(f_a, f_j)} \cdot \ell'_{\text{lr}}(a, i, j)$$

Although the rank constraint holds, the gradients' magnitudes could be significant if  $\ell'_{\text{lr}}(a, i, j)$  is large.

*Triplet rank loss*

$$\ell_{\text{tri}}(a, i, j) = [D(f_a, f_i) - D(f_a, f_j) + \delta]_+$$

$$\frac{\partial \ell_{\text{tri}}(a, i, j)}{\partial f_a} = - \frac{\partial \ell_{\text{tri}}(a, i, j)}{\partial f_i} - \frac{\partial \ell_{\text{tri}}(a, i, j)}{\partial f_j}$$

$$\frac{\partial \ell_{\text{tri}}(a, i, j)}{\partial f_i} = 2(f_i - f_a) \cdot \mathbb{I}(\ell_{\text{tri}}(a, i, j) > 0)$$

$$\frac{\partial \ell_{\text{tri}}(a, i, j)}{\partial f_j} = 2(f_a - f_j) \cdot \mathbb{I}(\ell_{\text{tri}}(a, i, j) > 0)$$

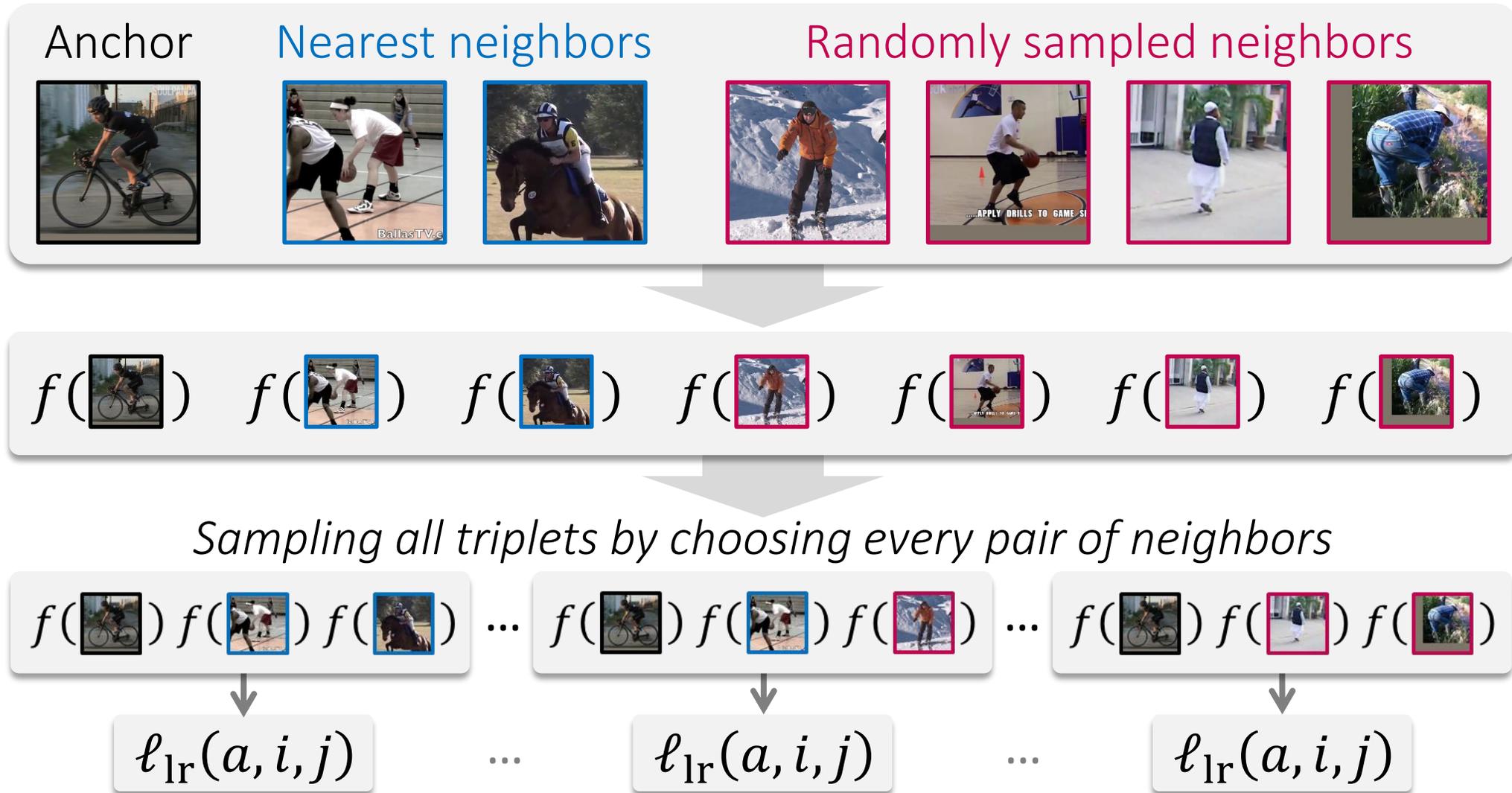
The gradients are zero if the triplet satisfies the rank constraint due to the indicator  $\mathbb{I}(\ell_{\text{tri}}(a, i, j) > 0)$ .

# Log-ratio Loss

- Compared to the triplet rank loss, our loss
  - Captures continuous similarities between images better, (the triplet rank loss focuses only on partial ranks of similarities.)
  - Does not require any hyperparameter, (for the triplet rank loss the margin should be tuned carefully.)
  - Does not demand  $L_2$  normalization of the embedding vectors, (such a normalization is essential for the triplet rank loss.)
  - Performs much better with a low embedding dimension.

# Dense Triplet Sampling

- Main idea: Using all triplets within a minibatch



# Dense Triplet Sampling

- Why not using existing sampling techniques<sup>[1,7]</sup>
  - They rely on binary relations between images.
  - They are designed to be combined with conventional triplet losses.
  - The notion of hardness is not clear in our setting.
- Our sampling strategy is well matched with the log-ratio loss.
  - The log-ratio loss enables every triplet to well contribute to training.

$$\frac{\partial \ell_{\text{lr}}(a, i, j)}{\partial f_i} = \frac{(f_i - f_a)}{D(f_a, f_i)} \cdot 4 \left\{ \log \frac{D(f_a, f_i)}{D(f_a, f_j)} - \log \frac{D_{\mathbf{y}}(\mathbf{y}_a, \mathbf{y}_i)}{D_{\mathbf{y}}(\mathbf{y}_a, \mathbf{y}_j)} \right\}$$

Non-trivial even if the triplet complies the rank constraint

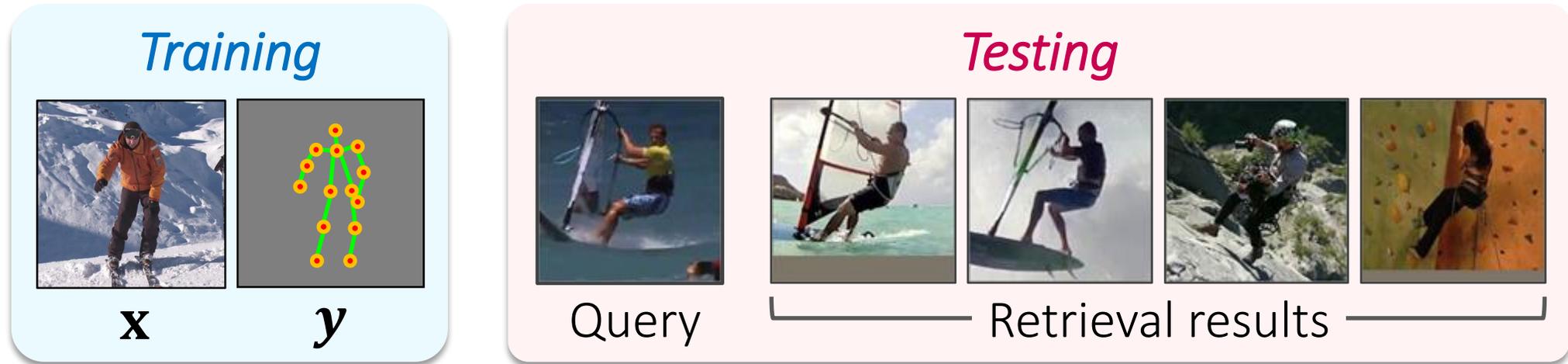
- *Exploiting all triplets improves embedding performance.*

[1] FaceNet: A unified embedding for face recognition and clustering, CVPR 2015

[7] Sampling matters in deep embedding learning, ICCV 2017

# Experiments – Three Retrieval Tasks

- Human pose retrieval

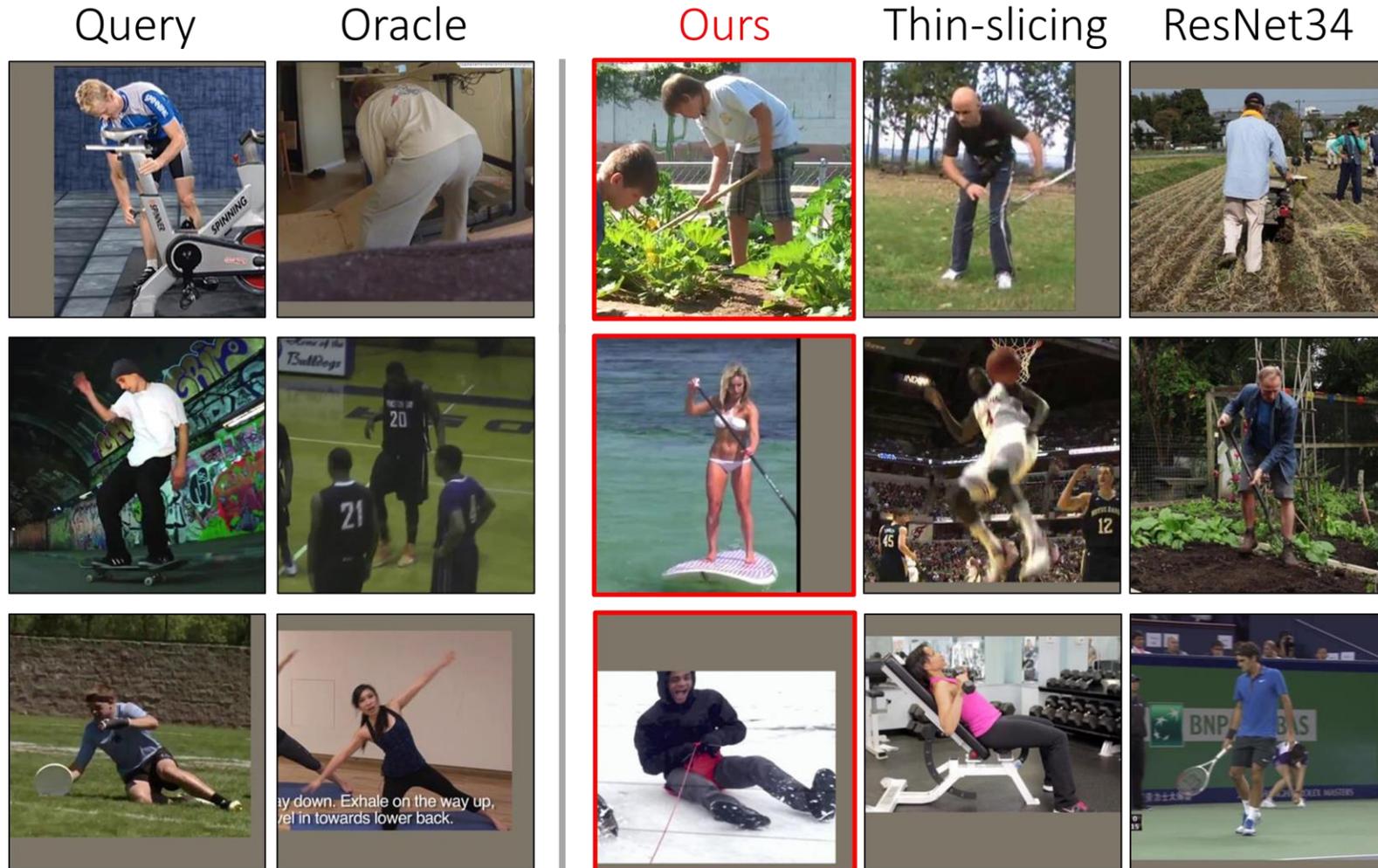


- Conducted on the *MPII human pose dataset*
- Similarity between images: *inverse pose distances*
- Application: *pose-aware representation for action recognition*
- Label distance between images:

$$D_y(\mathbf{y}_i, \mathbf{y}_j) = \|\mathbf{y}_i - \mathbf{y}_j\|_2^2,$$

# Experiments – Three Retrieval Tasks

- Human pose retrieval



ResNet34: ImageNet pre-trained network

Typically focuses on objects or background other than human poses.

Thin-slicing<sup>[6]</sup>: A previous work on pose embedding

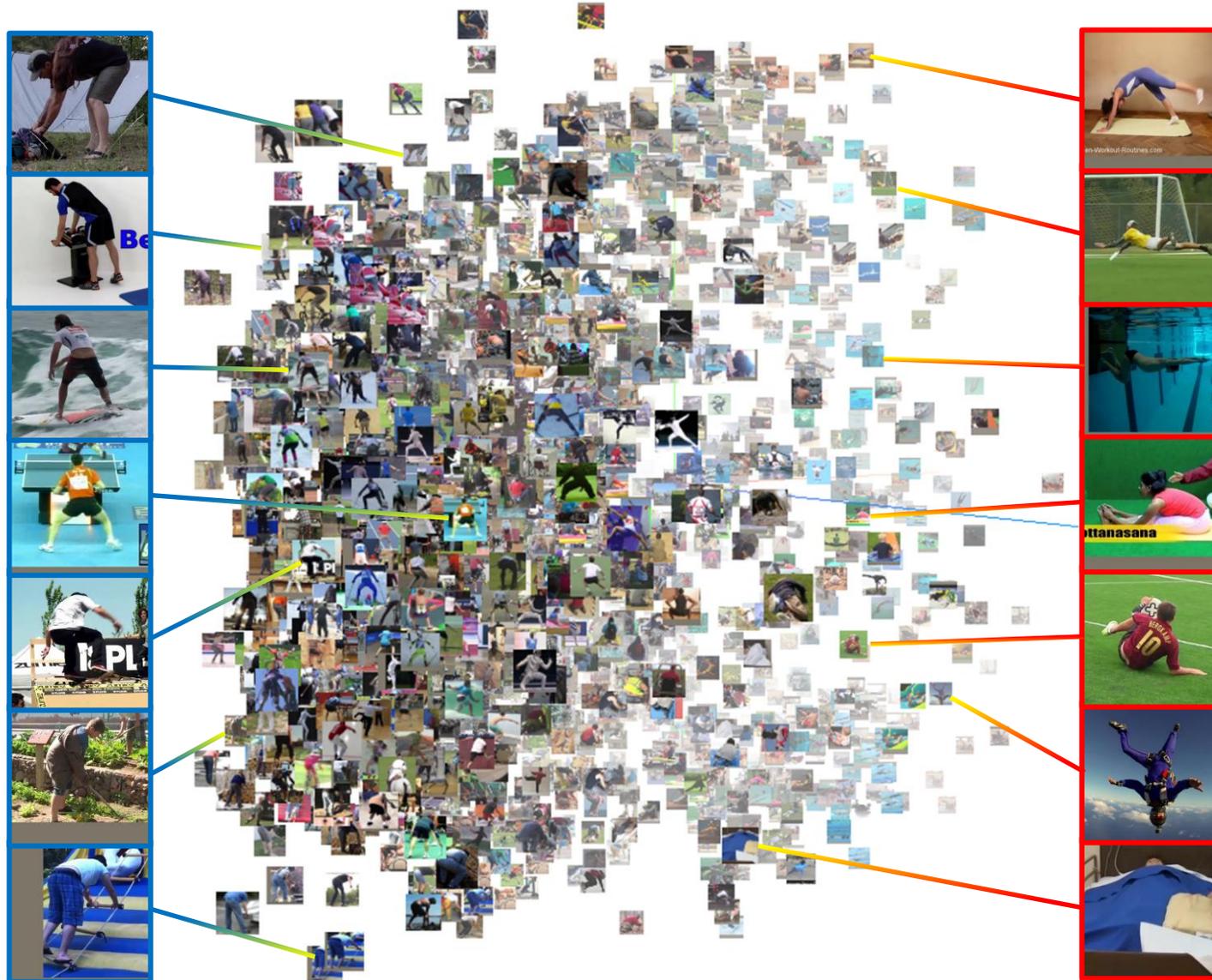
Often fails to address rare human poses.



[6] Thin-slicing for pose: Learning to understand pose without explicit pose estimation, CVPR 2016

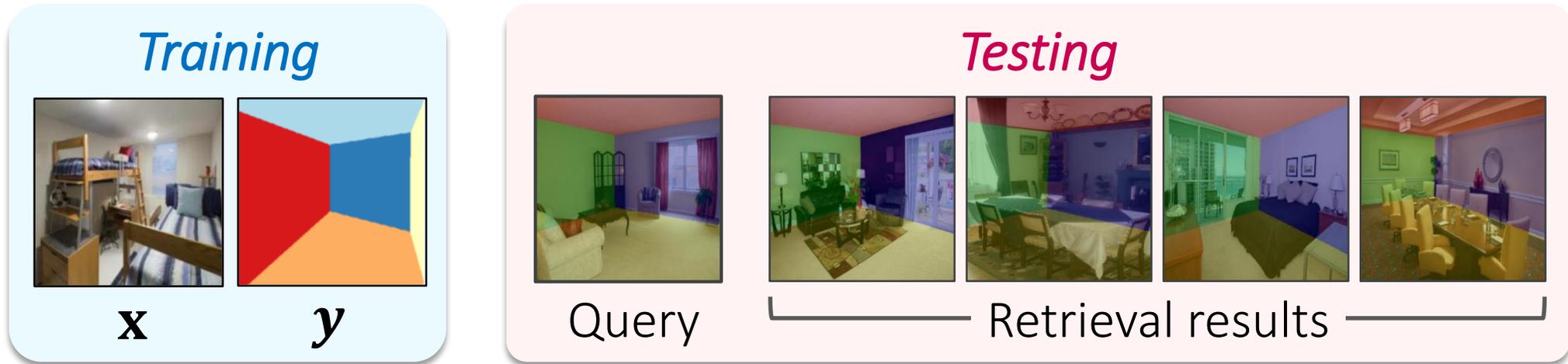
# Experiments – Three Retrieval Tasks

- Human pose retrieval



# Experiments – Three Retrieval Tasks

- Room layout retrieval



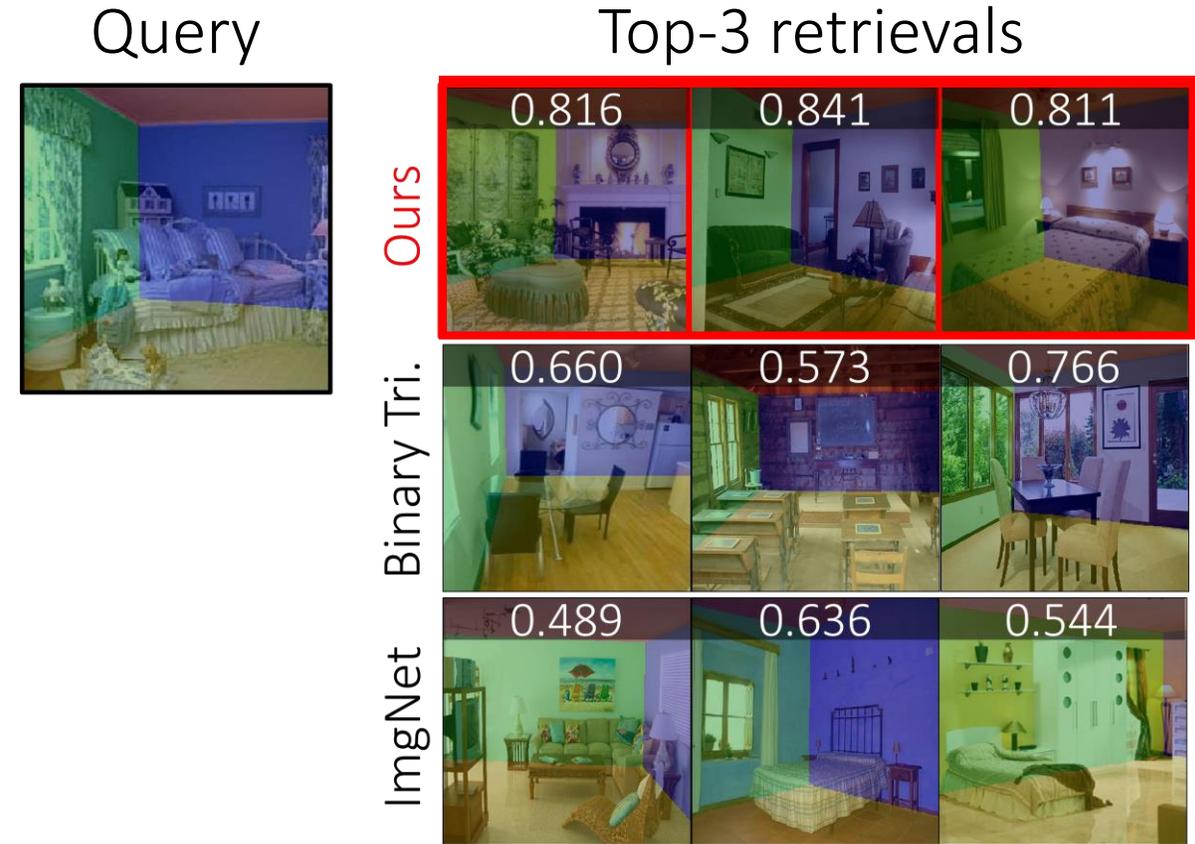
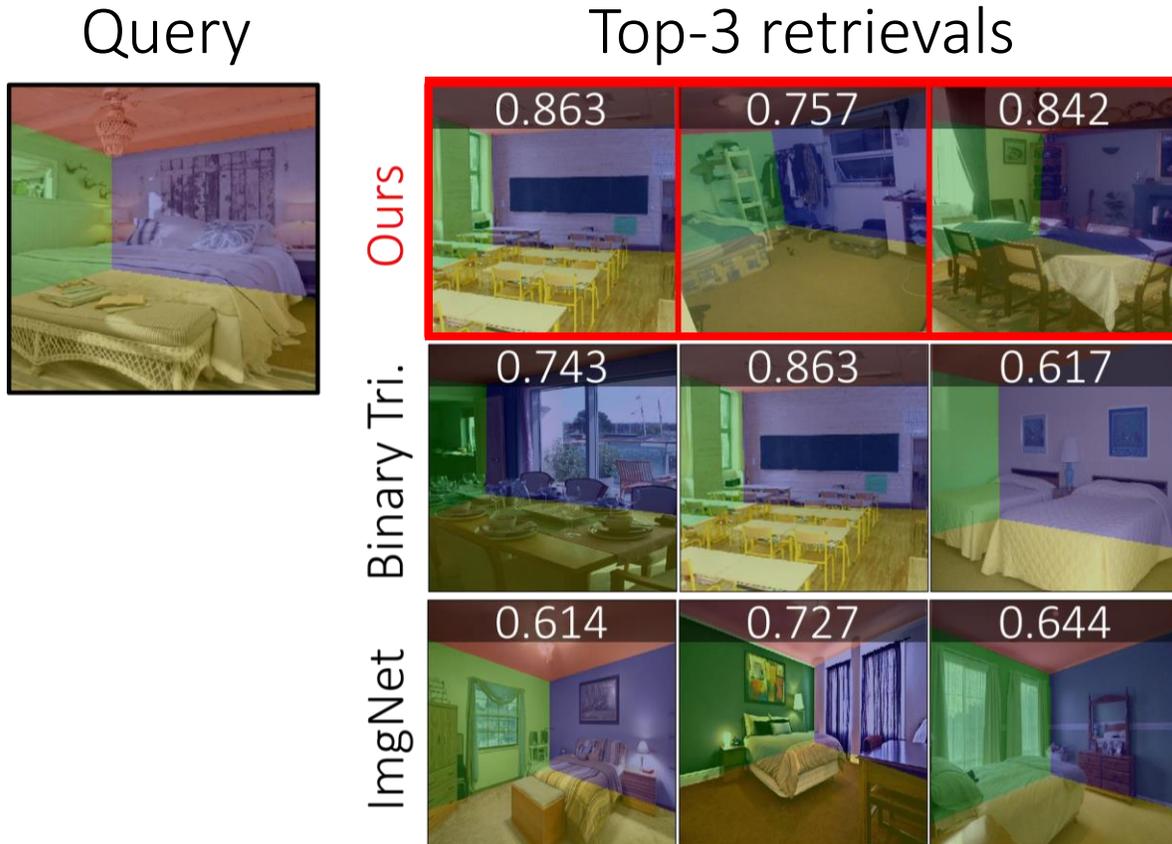
- Conducted on the *LSUN room layout dataset*
- Label distance between images:

$$D_{\mathbf{y}}(\mathbf{y}_i, \mathbf{y}_j) = 1 - \text{mIoU}(\mathbf{y}_i, \mathbf{y}_j),$$

where  $\mathbf{y}_i$  and  $\mathbf{y}_j$  denote groundtruth room segmentations

# Experiments – Three Retrieval Tasks

- Room layout retrieval

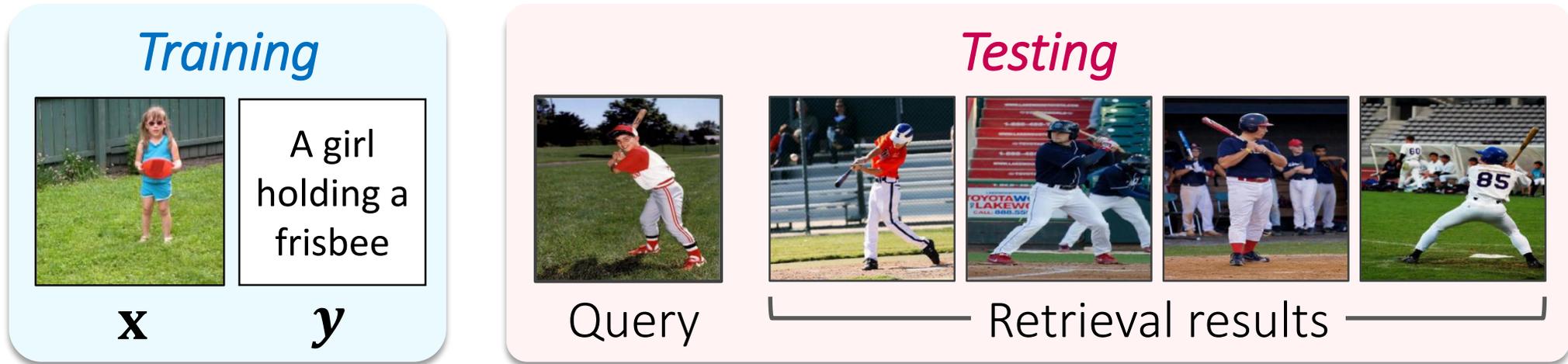


***Binary Tri.***: Triplet rank loss + Binary thresholding

***ImgNet***: ImageNet pre-trained ResNet101

# Experiments – Three Retrieval Tasks

- Caption-aware image retrieval



- Conducted on the *MS-COCO 2014 caption dataset*
- Label distance between images:

$$D_y(\mathbf{y}_i, \mathbf{y}_j) = \sum_{c_i \in \mathbf{y}_i} \min_{c_j \in \mathbf{y}_j} W(c_i, c_j) + \sum_{c_j \in \mathbf{y}_j} \min_{c_i \in \mathbf{y}_i} W(c_i, c_j),$$

where  $\mathbf{y}_i$  and  $\mathbf{y}_j$  are sets of 5 captions and  $W(\cdot)$  is the WMD<sup>[8]</sup> between two captions

[8] From word embeddings to document distances, ICML 2015

# Experiments – Three Retrieval Tasks

- Caption-aware image retrieval

Query



Top-3 retrievals

Ours



Binary Tri.



ImgNet



Query



Top-3 retrievals

Ours



Binary Tri.



ImgNet



Binary Tri.: Triplet rank loss + Binary thresholding

ImgNet: ImageNet pre-trained ResNet101

# Experiments – Three Retrieval Tasks

- Caption-aware image retrieval

Query



Top-3 retrievals

Ours



Binary Tri.



ImgNet



Query



Top-3 retrievals

Ours



Binary Tri.



ImgNet

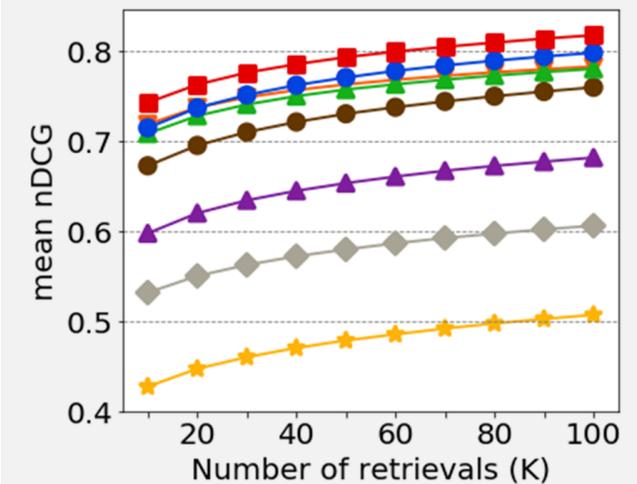
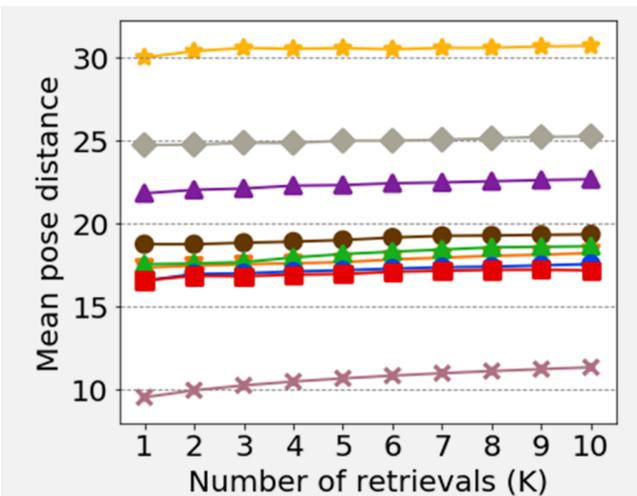


Binary Tri.: Triplet rank loss + Binary thresholding

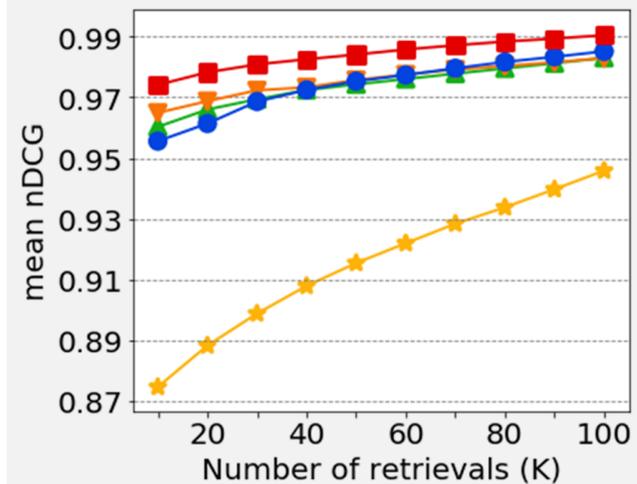
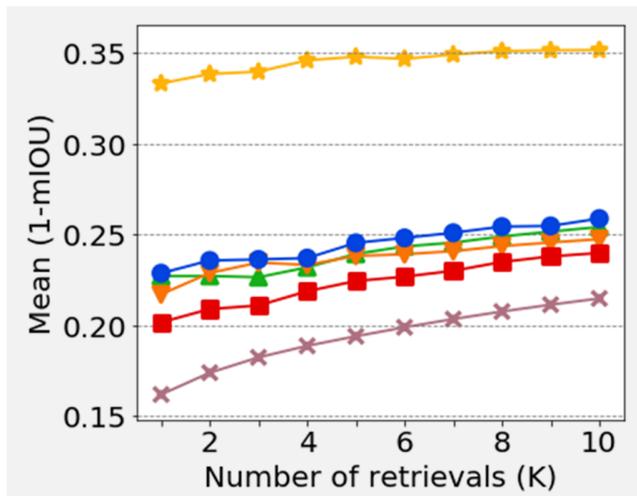
ImgNet: ImageNet pre-trained ResNet101

# Experiments – Three Retrieval Tasks

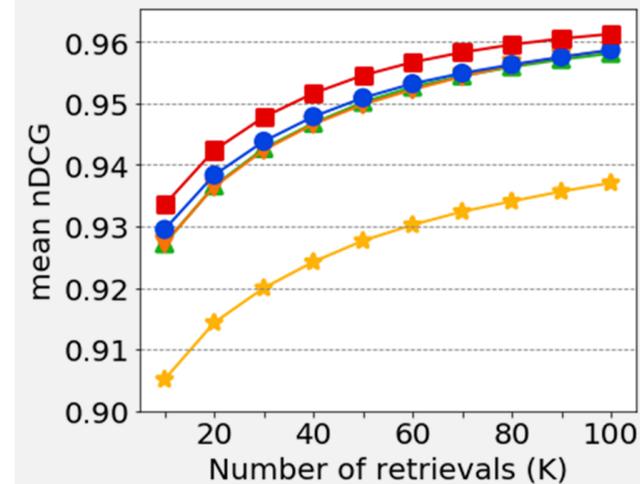
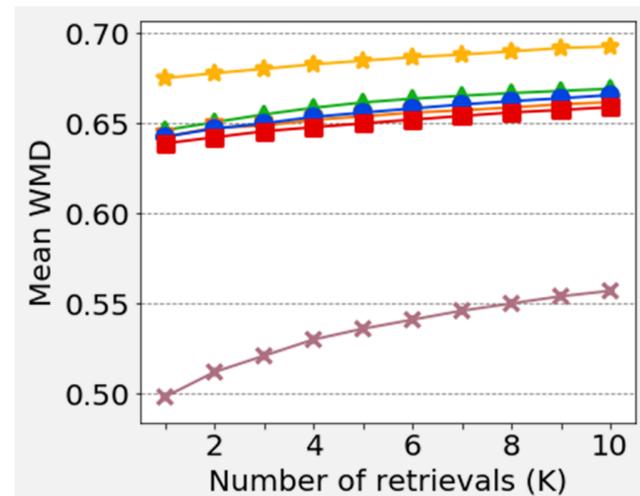
- Quantitative performance analysis



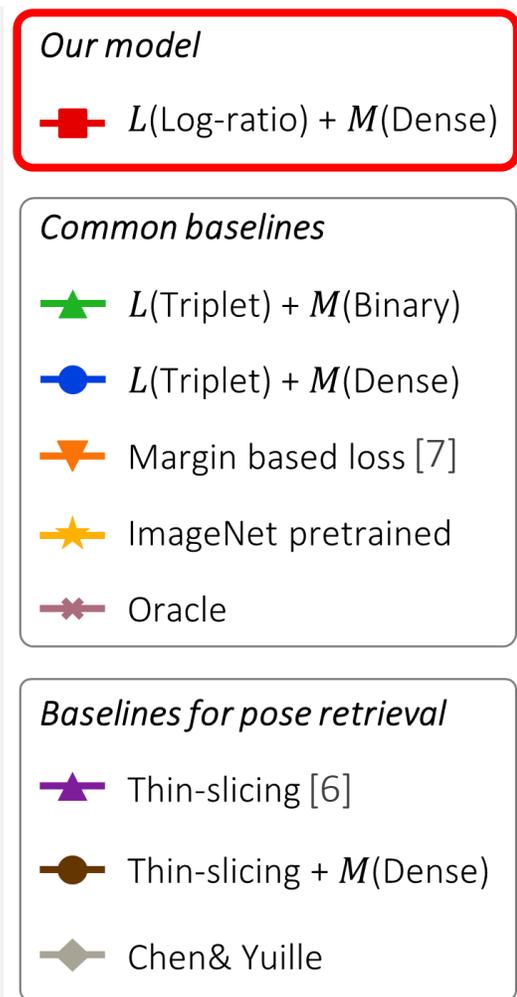
Human pose retrieval



Room layout retrieval

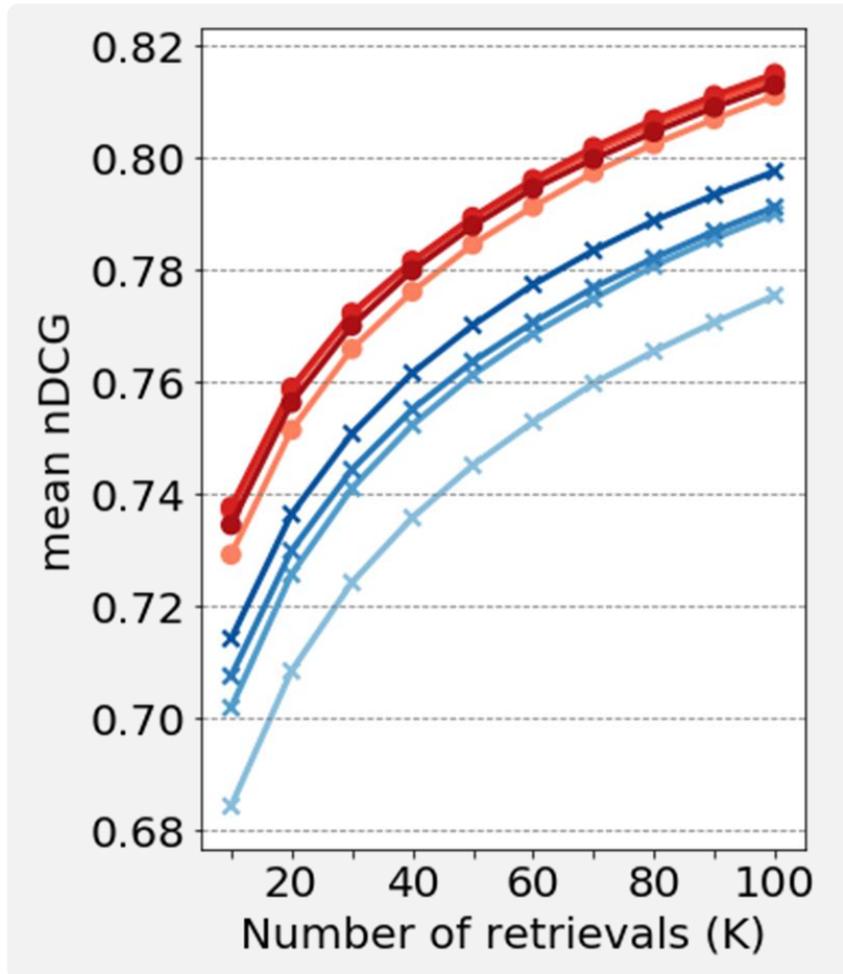


Caption-aware image retrieval



# Experiments – Three Retrieval Tasks

- Embedding dimension vs. retrieval performance



## Our models

- $L(\text{Log-ratio}) + M(\text{Dense})$  128-D
- $L(\text{Log-ratio}) + M(\text{Dense})$  64-D
- $L(\text{Log-ratio}) + M(\text{Dense})$  32-D
- $L(\text{Log-ratio}) + M(\text{Dense})$  16-D

## Baselines

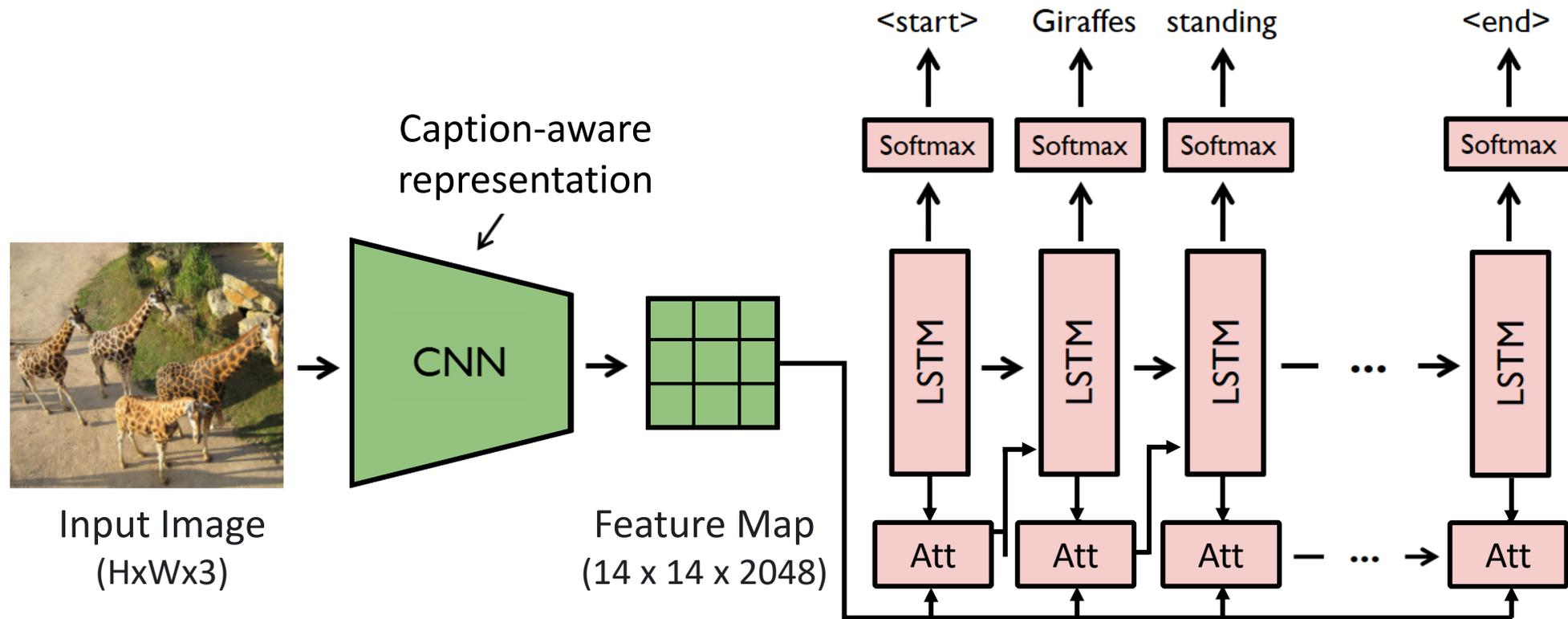
- $L(\text{Triplet}) + M(\text{Dense})$  128-D
- $L(\text{Triplet}) + M(\text{Dense})$  64-D
- $L(\text{Triplet}) + M(\text{Dense})$  32-D
- $L(\text{Triplet}) + M(\text{Dense})$  16-D

$L(\text{Log-ratio}) + M(\text{Dense})$ : Log-ratio loss + Dense triplet sampling

$L(\text{Triplet}) + M(\text{Dense})$ : Triplet rank loss + Dense triplet sampling

# Experiments – Representation Learning

- Representation learning for image captioning



## *Our approach*

Using the caption embedding network trained with caption similarities as an initial visual representation for image captioning

# Experiments – Representation Learning

- Quantitative results

**115.9** in CIDEr  
Caption-aware feature + RL

**34.65** in BLEU-4  
Caption-aware feature + RL



**+2.5%**

**+3.5%**

**113.1** in CIDEr  
ImageNet pretrained feature + RL

**33.48** in BLEU-4  
ImageNet pretrained feature + RL

[9] Self-critical sequence training for image captioning, CVPR 2017

[10] Bottom-up and top-down attention for image captioning and visual question answering, CVPR 2018

# Experiments – Representation Learning

- Qualitative results obtained by the top-down attention model



GT1	There are some zebras <b>standing</b> in a grassy field
GT2	A field with tall grass, bushes and trees, that has zebra <b>standing</b> in the field
Img XE	A group of zebras <b>grazing</b> in a field
Cap XE	Two zebras are <b>standing</b> in a grassy field
Img RL	A group of zebras are <b>grazing</b> in a field
Cap RL	A couple of zebras and a zebra <b>standing</b> in a field



GT1	A baseball batter <b>swinging</b> a bat over home plate
GT2	A baseball player <b>swings</b> a bat at a game
Img XE	A baseball player <b>holding</b> a bat on a field
Cap XE	A baseball player <b>swinging</b> a bat on top of a field
Img RL	A baseball player <b>holding</b> a bat on a field
Cap RL	A baseball player <b>swinging</b> a bat at a ball

# Experiments – Representation Learning

- Visualization of attentions drawn by the Att2all2 model



Img RL A baseball player **holding** a bat on a field

Cap RL A baseball player **swinging** a bat at a ball

# Conclusion

- Summary
  - A new framework for metric learning with continuous labels
  - Various applications including visual representation learning
  - Performance boost over existing approaches
- Future directions
  - A better distance metric for continuous and structured labels
  - A hard triplet mining technique for continuous metric learning
  - More applications of semantic nearest neighbor search
  - A new benchmark for continuous metric learning

# References

- [1] FaceNet: A unified embedding for face recognition and clustering, CVPR 2015
- [2] Beyond triplet loss: a deep quadruplet network for person re-identification, CVPR 2017
- [3] Learning to compare image patches via convolutional neural networks, CVPR 2015
- [4] Learning a similarity metric discriminatively with application to face verification, CVPR 2005
- [5] Pose embeddings: A deep architecture for learning to match human poses, arXiv 2015
- [6] Thin-slicing for pose: Learning to understand pose without explicit pose estimation, CVPR 2016
- [7] Sampling matters in deep embedding learning, ICCV 2017
- [8] From word embeddings to document distances, ICML 2015
- [9] Self-critical sequence training for image captioning, CVPR 2017
- [10] Bottom-up and top-down attention for image captioning and visual question answering, CVPR 2018

