



# Similarity Learning with (or without) Convolutional Neural Network

Moitreya Chatterjee, Yunan Luo



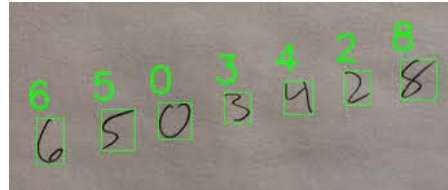
# Outline – This Section

- **Why do we need Similarity Measures**
- **Metric Learning as a measure of Similarity**
  - **Notion of a metric**
  - **Unsupervised Metric Learning**
  - **Supervised Metric Learning**
- **Traditional Approaches for Matching**
- **Challenges with Traditional Matching Techniques**
- Deep Learning as a Potential Solution
- Application of Siamese Network for different tasks



# Need for Similarity Measures

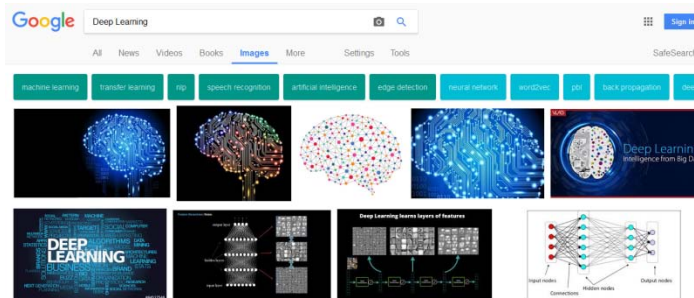
Several applications of Similarity Measures exists in today's world:



- Recognizing handwriting in checks.



- Automatic detection of faces in a camera image.



- Search Engines, such as Google, matching a **query** (could be text, image, etc.) with a set of **indexed documents** on the web.



# Notion of a Metric

- A **Metric** is a function that quantifies a “distance” between every pair of elements in a set, thus inducing a measure of similarity.
- A metric  $f(x,y)$  must satisfy the following properties for all  $x, y, z$  belonging to the set:
  - *Non-negativity*:  $f(x, y) \geq 0$
  - *Identity of Discernible*:  $f(x, y) = 0 \iff x = y$
  - *Symmetry*:  $f(x, y) = f(y, x)$
  - *Triangle Inequality*:  $f(x, z) \leq f(x, y) + f(y, z)$



# Types of Metrics

In broad strokes metrics are of two kinds:

- **Pre-defined Metrics:** Metrics which are fully specified without the knowledge of data.

E.g. Euclidian Distance:  $f(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T(\mathbf{x} - \mathbf{y})$

- **Learned Metrics:** Metrics which can only be defined with the **knowledge** of the **data**.

E.g. Mahalanobis Distance:  $f(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T \mathbf{M}(\mathbf{x} - \mathbf{y})$  ;  
where **M** is a matrix that is estimated from the data.

Learned Metrics are of two types:

- **Unsupervised** : Use unlabeled data
- **Supervised** : Use labeled data



# **UNSUPERVISED METRIC LEARNING**



# Mahalanobis Distance

- Mahalanobis Distance weighs the Euclidian distance between two points, by the standard deviation of the data.
  - $f(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T \Sigma^{-1} (\mathbf{x} - \mathbf{y})$ ; where  $\Sigma$  is the mean-subtracted covariance matrix of all data points.

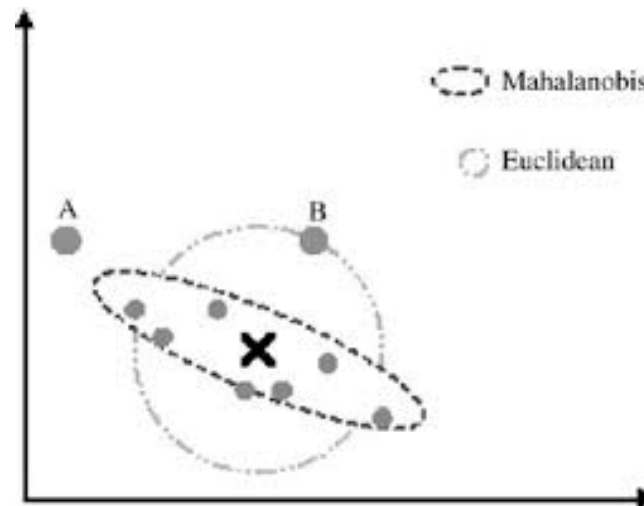


Image Source:  
Google

Chandra, M.P., 1936. On the generalised distance in statistics. In *Proceedings of the National Institute of Sciences of India* (Vol. 2, No. 1, pp. 49-55).



# **SUPERVISED METRIC LEARNING**





# Supervised Metric Learning

- In this setting, we have access to **labeled** data samples ( $z = \{x, y\}$ ).
- The typical strategy is to use a 2-step procedure:
  - Apply some **supervised** domain transform.
  - Then use one of the unsupervised metrics for performing the mapping.

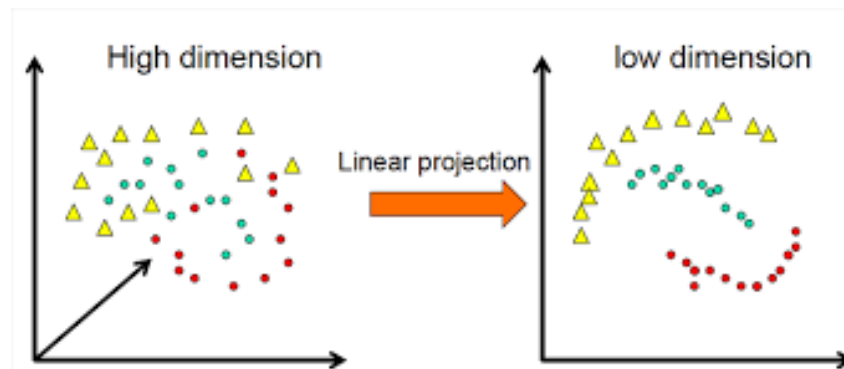


Image Source:  
Google

Bellet, A., Habrard, A. and Sebban, M., 2013. A survey on metric learning for feature vectors and structured data. *arXiv preprint arXiv:1306.6709*.



# Linear Discriminant Analysis (LDA)

- In Fisher-LDA, the goal is to project the data to a space such that the ratio of “**between class covariance**” to “**within class covariance**” is maximized.
- This is given by:  $J(w) = \max_w (w^T S_B w) / (w^T S_W w)$

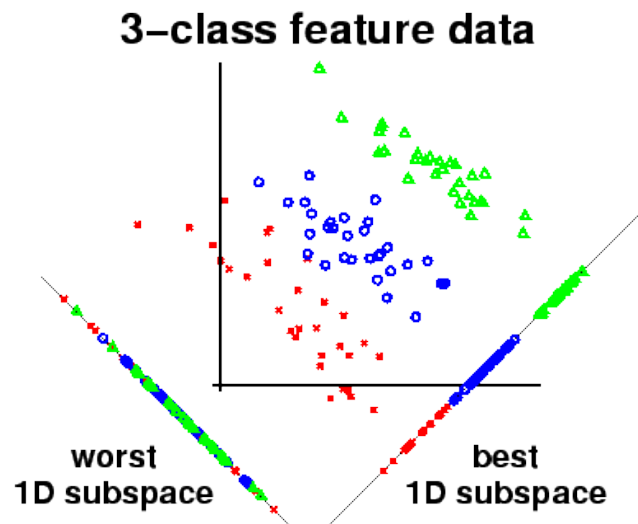


Image Source:  
Google

Fisher, R.A., 1936. The use of multiple measurements in taxonomic problems. *Annals of eugenics*, 7(2), pp.179-188.



# **TRADITIONAL MATCHING TECHNIQUES**



# Traditional Approaches for Matching

The traditional approach for matching images, relies on the following pipeline:

- 1. Extract Features:** For instance, color histograms of the input images.
- 2. Learn Similarity:** Use  $L_1$ -norm on the features.



# Challenges with Traditional Methods for Matching

The principal shortcoming of traditional metric learning based methods is that the **feature representation** of the data and the **metric** are **not learned jointly**.



# Outline – This Section

- Why do we need Similarity Measures
- Metric Learning as a measure of Similarity
- Traditional Approaches for Similarity Learning
- Challenges with Traditional Similarity Measures
- **Deep Learning as a Potential Solution**
  - **Siamese Networks**
    - **Architectures**
    - **Loss Function**
    - **Training Techniques**
- Application of Siamese Network to different tasks



# Deep Learning to the Rescue!

CNNs can **jointly optimize** the representation of the input data conditioned on the “similarity” measure being used, aka end-to-end learning.



Image Source:  
Google



# Revisit the Problem

- **Input:** Given a pair of input images, we want to know how “similar” they are to each other.
- **Output:** The output can take a variety of forms:
  - Either a binary label, i.e. 0 (same) or 1 (different).
  - A **Real** number indicating how similar a pair of images are.





# Typical Siamese CNN

- **Input:** A pair of input signatures.
- **Output (Target):** A label, **0** for similar, **1** else.

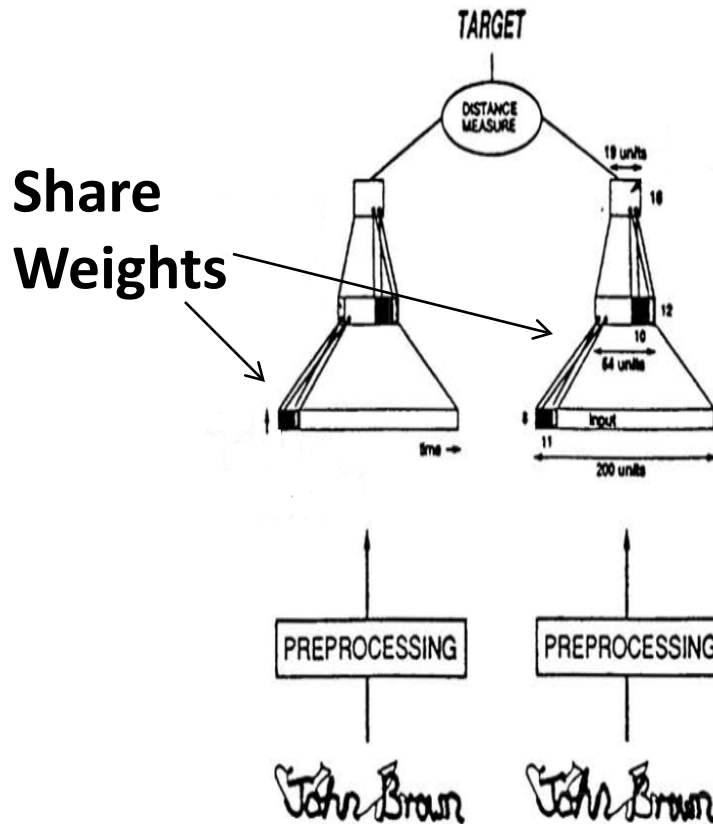


Image Source:  
Google

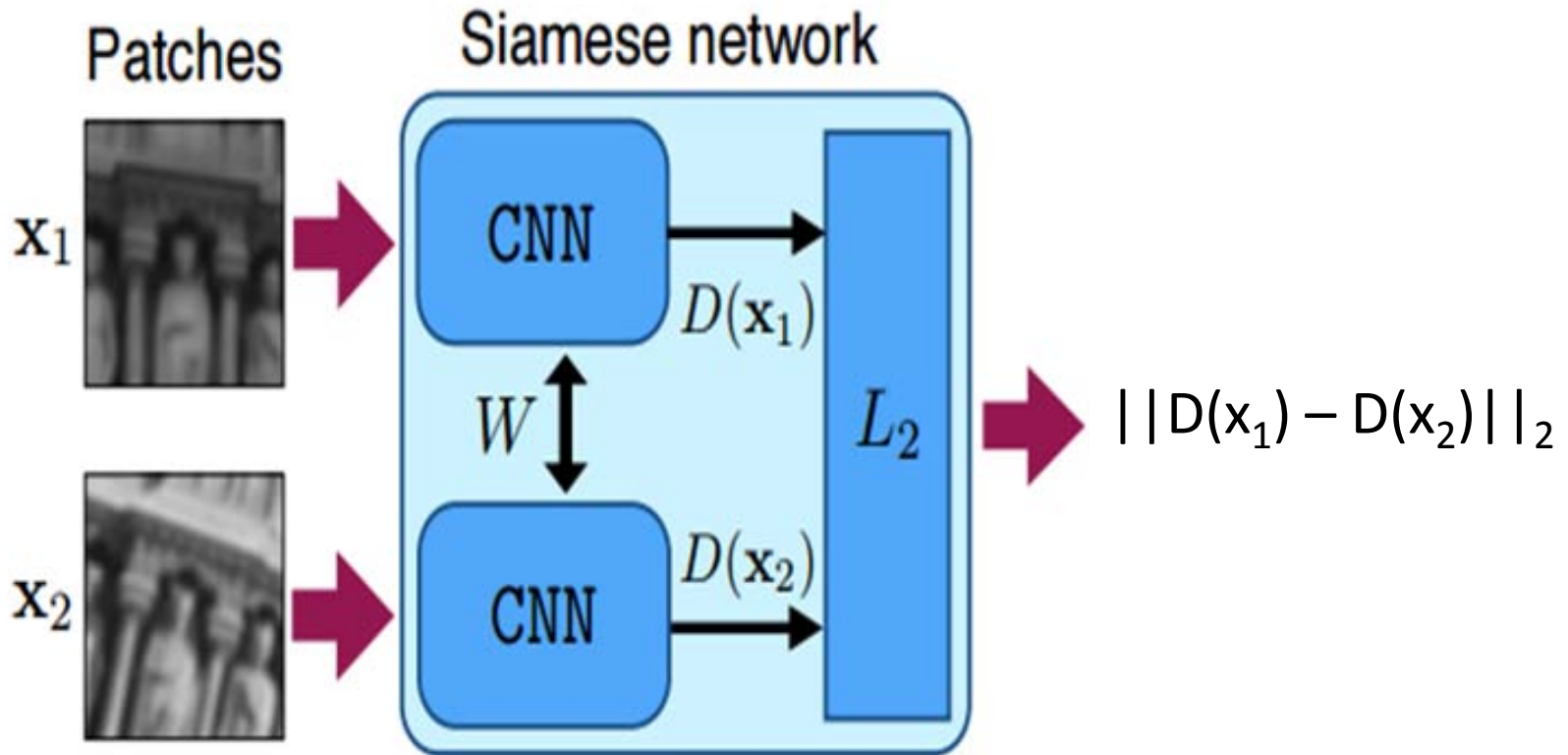
Bromley, J., Bentz, J.W., Bottou, L., Guyon, I., LeCun, Y., Moore, C., Säckinger, E. and Shah, R., 1993. Signature Verification Using A "Siamese" Time Delay Neural Network. *IJPRAI*, 7(4), pp.669-688.



# SIAMESE CNN - ARCHITECTURE



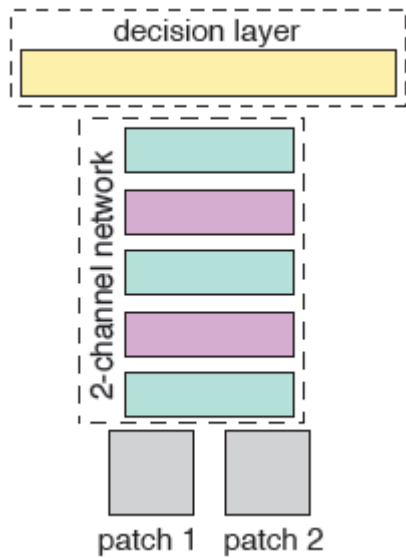
# Standard architecture of Siamese CNN



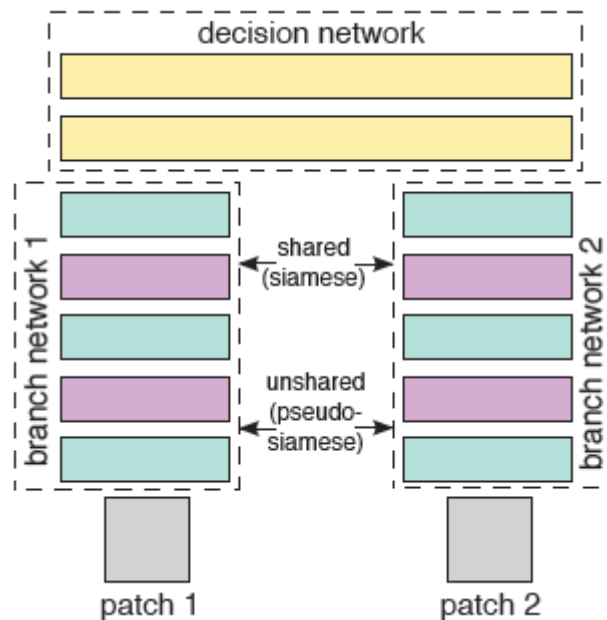


# Popular Architecture Varieties

- No one “architecture” fits all!
- Design largely governed by what performs well empirically on the task at hand.



Inputs are merged right at the onset

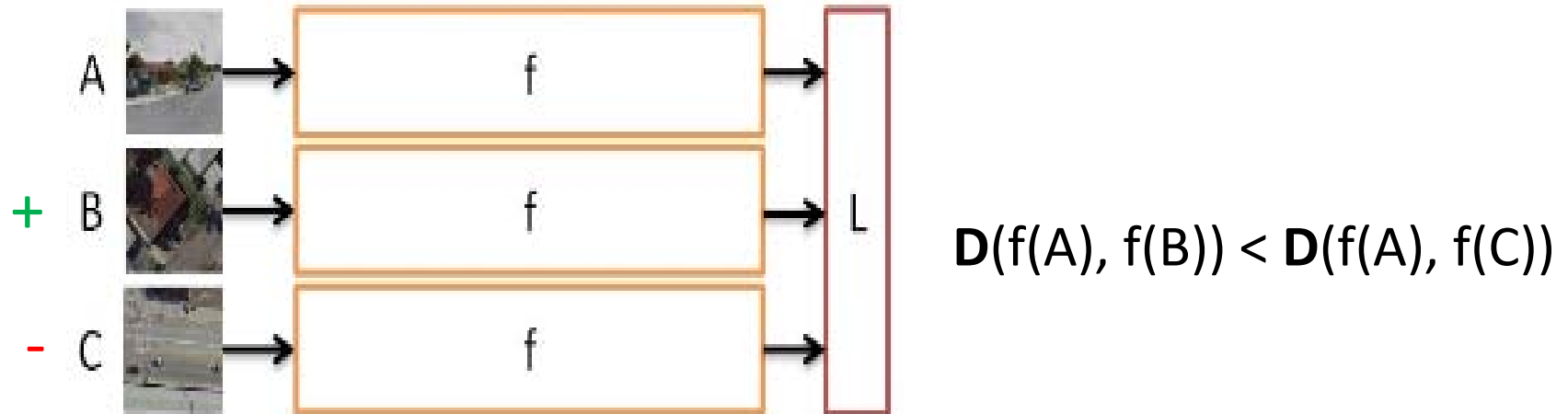


Inputs are first embedded independently, then merged.



# Siamese CNN – Variants

## TRIPLET NETWORK



- Compare triplets in one go.
- Check if the sample in the **topmost** channel, is more similar to the one in the **middle** or the one in the **bottom**.
- Allows us to learn ranking between samples.

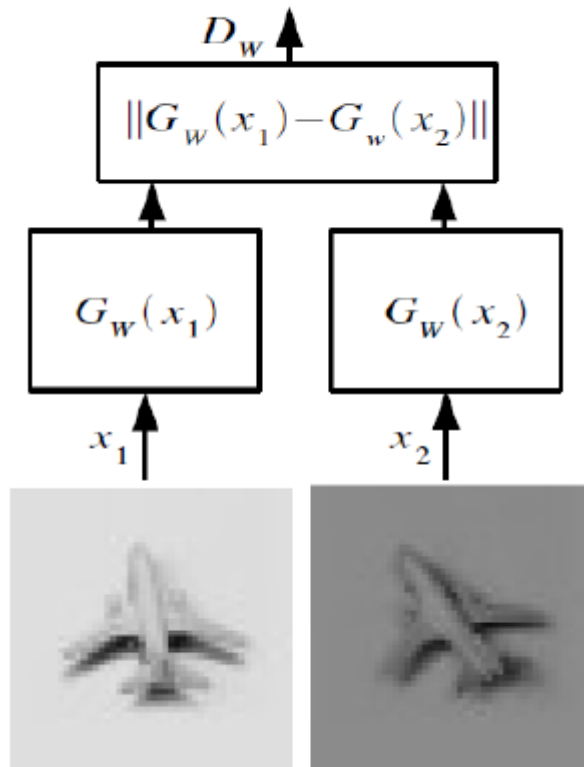


# SIAMESE CNN – LOSS FUNCTION



# Siamese CNN – Loss Function

Make this small



Similar images

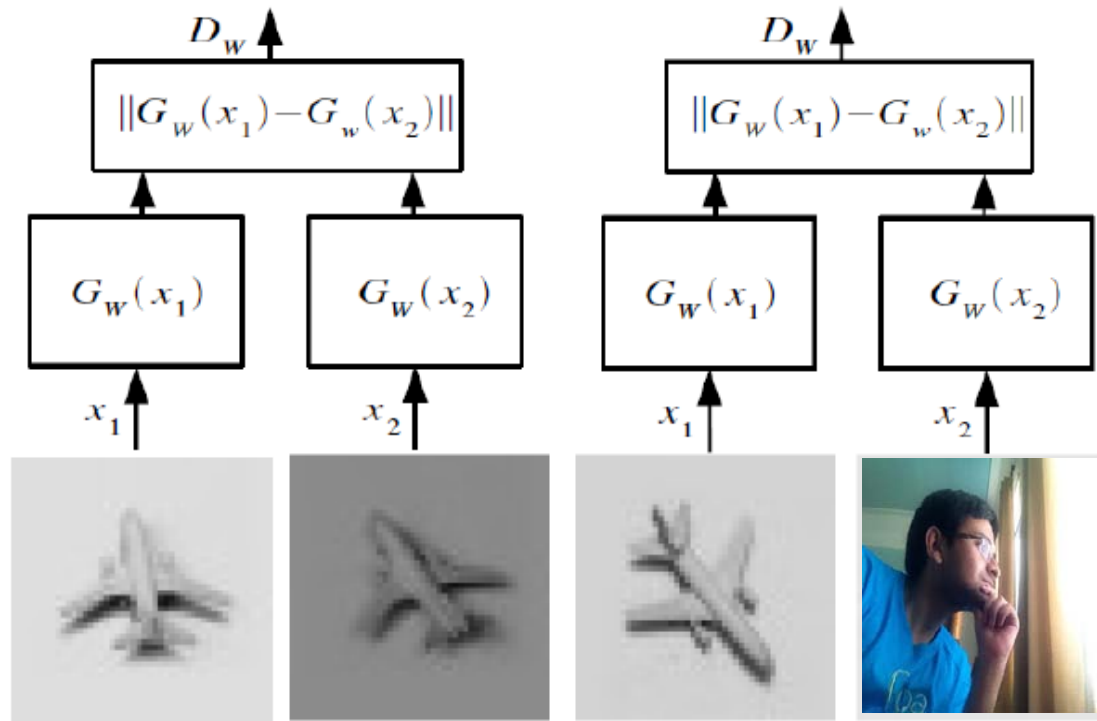
- Is there a problem with this formulation?
  - Yes.
  - The model could learn to embed every input to the same point, i.e. predict a **constant** as **output**.
  - In such a case, every pair of input would be categorized as a positive pair.



# Siamese CNN – Loss Function

Make this small

Make this large



Similar images

Dissimilar images

The final loss is defined as :

$$L = \sum \text{loss of positive pairs} + \sum \text{loss of negative pairs}$$

Chopra, S., Hadsell, R. and LeCun, Y., 2005, June. Learning a similarity metric discriminatively, with application to face verification. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on* (Vol. 1, pp. 539-546). IEEE.

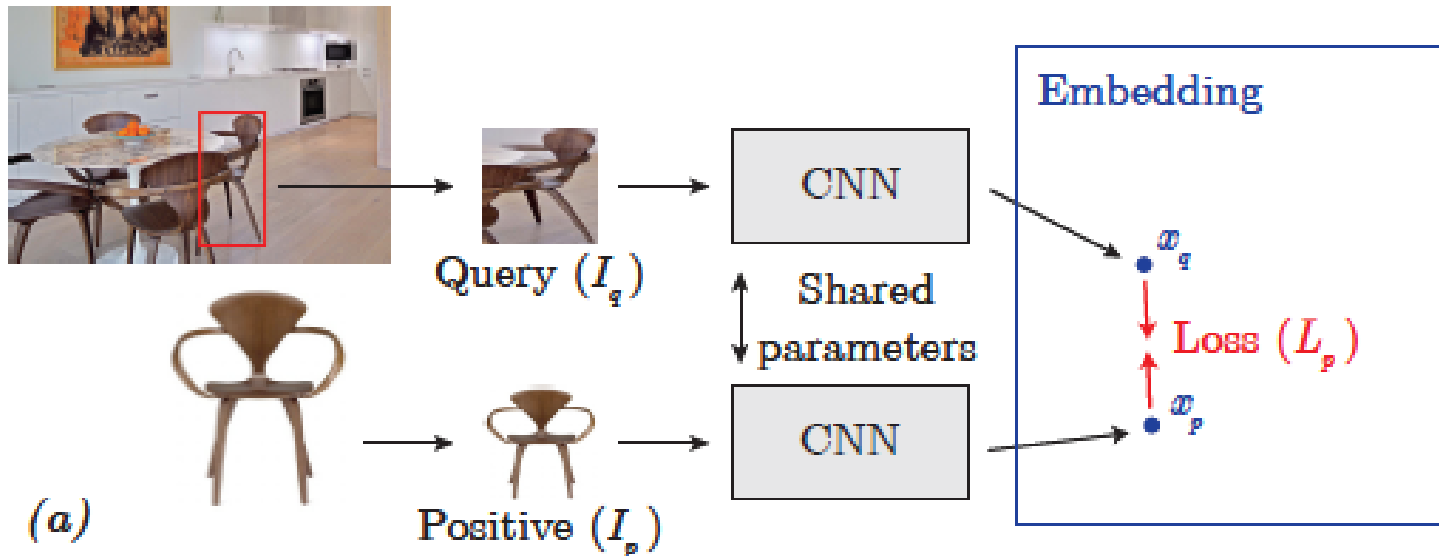




# Siamese CNN – Loss Function

We can use different loss functions for the two types of input pairs.

- Typical **positive pair** ( $x_p, x_q$ ) loss:  $L(x_p, x_q) = ||x_p - x_q||^2$   
(Euclidian Loss)

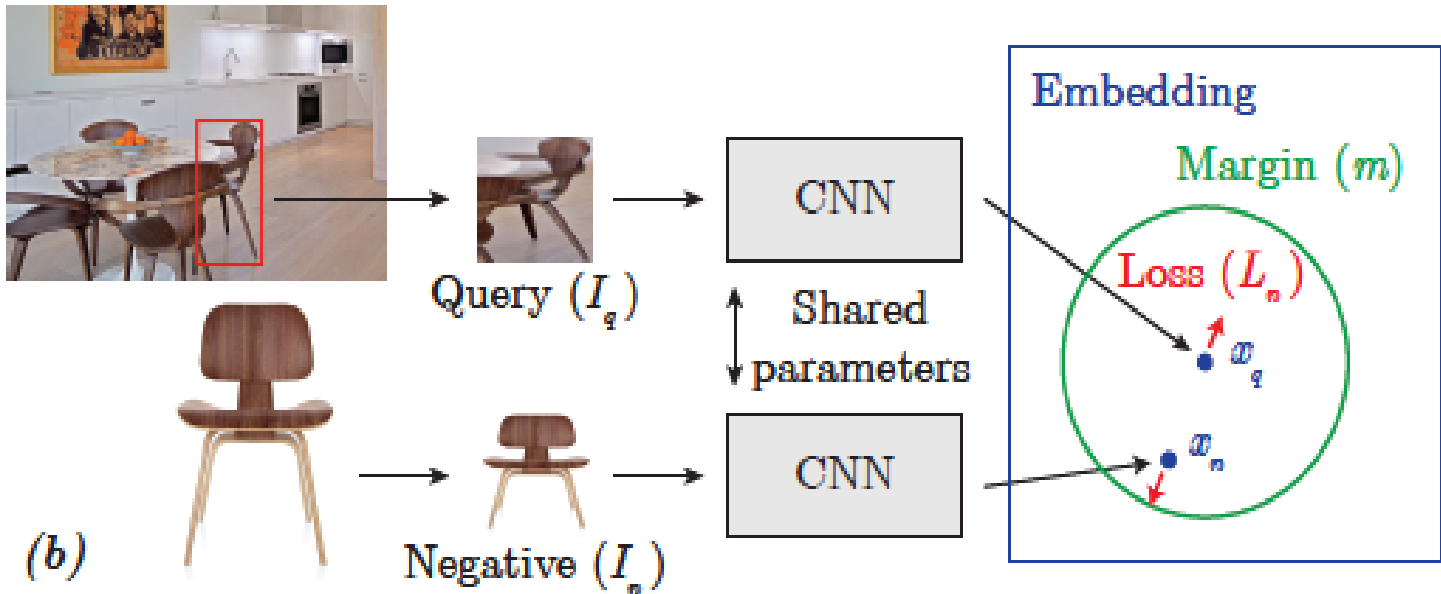




# Siamese CNN – Loss Function

- Typical **negative pair**  $(x_n, x_q)$  loss :

$$L(x_n, x_q) = \max(0, m^2 - ||x_n - x_q||^2) \text{ (Hinge Loss)}$$





# Choices of Loss Function

- Several choices for the Loss Functions are available. Choice depends on the task at hand.
- Loss Functions for **2-Stream Networks**:
  - **Margin Based**:
    - **Contrastive Loss**:  $\text{Loss}(x_p, x_q, y) = y * ||x_p - x_q||^2 + (1 - y) * \max(0, m^2 - ||x_p - x_q||^2)$ 
      - Allows us to learn a margin of separation.
      - Extensible for Triplet Networks
  - **Non-Margin Based**:
    - **Distance-Based Logistic Loss**:  
 $P(x_p, x_q) = (1 + \exp(-m)) / (1 + \exp(||x_p - x_q|| - m))$   
 $\text{Loss}(x_p, x_q, y) = \text{LogLoss}(P(x_p, x_q), y)$ 
      - Good for quicker convergence.

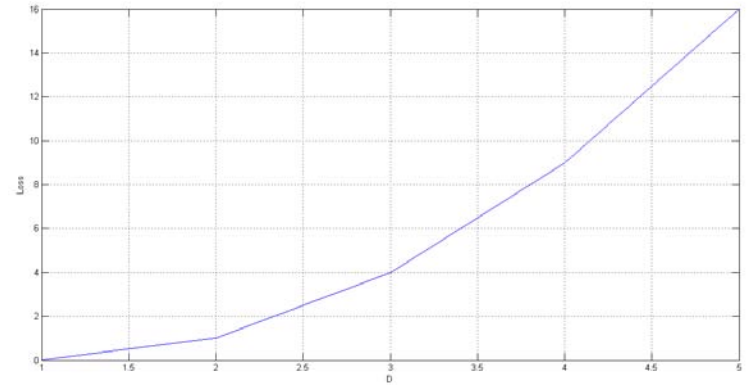


# Choices of Loss Function

- Contrastive Loss:

For similar samples:

$$\text{Loss}(x_p, x_q) = ||x_p - x_q||^2$$

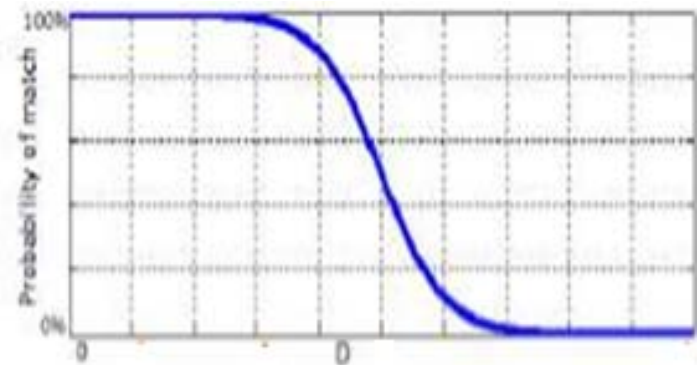


- Distance-Based Logistic Loss:

For similar pairs

$$P(x_p, x_q) = (1 + \exp(-m)) / (1 + \exp(||x_p - x_q|| - m)) \rightarrow 1 \text{ quickly}$$

$$\text{Loss}(x_p, x_q, y) = \text{LogLoss}(P(x_p, x_q), y)$$



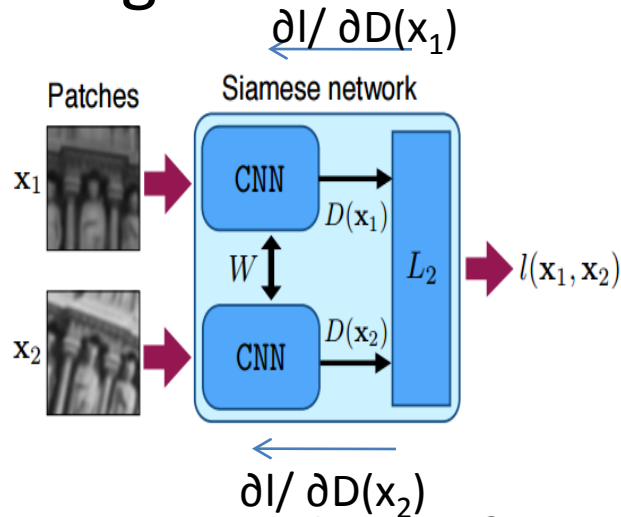


# SIAMESE CNN – TRAINING



# Siamese CNN – Training

- Update each of the two streams independently and then average the weights.



- Does this technique remind us of anything?
  - Training in RNNs.
- Data augmentation may be used for more effective training.
  - Typically we hallucinate more examples by performing random crops, image flipping, etc.



# Outline – This Section

- Why do we need Similarity Measures
- Metric Learning as a measure of Similarity
- Traditional Approaches for Similarity Learning
- Challenges with Traditional Similarity Measures
- Deep Learning as a Potential Solution
- **Application of Siamese Network to different tasks**
  - **Generating invariant and robust descriptors**
  - **Person Re-Identification**
  - **Rendering a street from Different Viewpoints**
  - **Newer nets for Person Re-Id, Viewpoint Invariance and Multimodal Data.**
  - **Use of Siamese Networks for Sentence Matching**

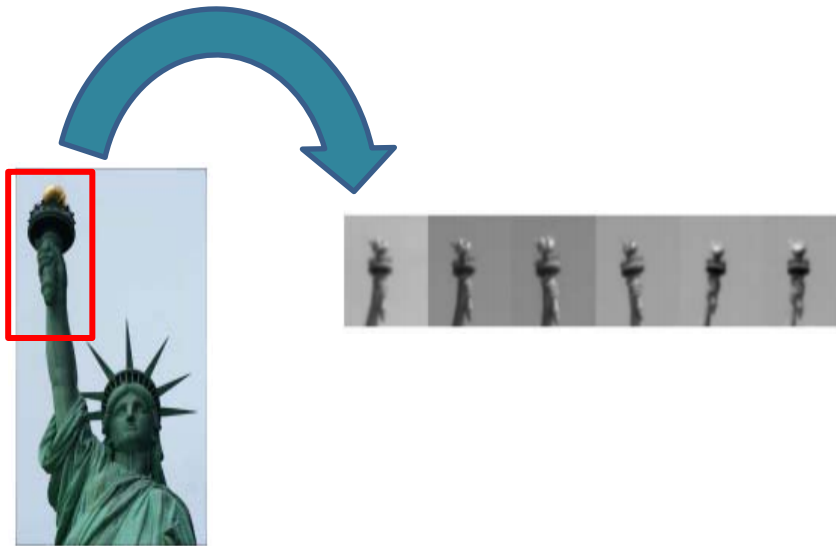


# APPLICATIONS





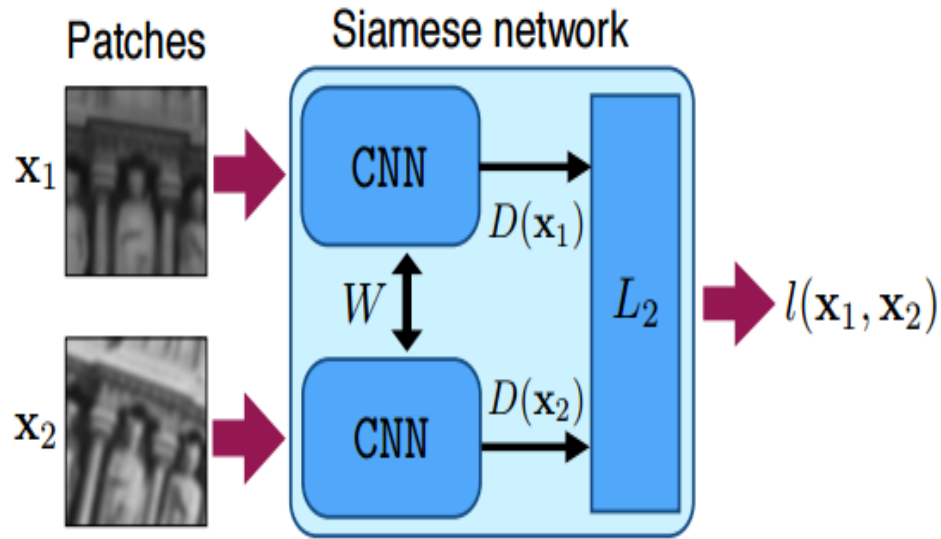
# Discriminative Descriptors for Local Patches



Learn a discriminative representation of patches from different views of 3D points



# Deep Descriptor



$$l(x_1, x_2) = \begin{cases} \|D(x_1) - D(x_2)\|_2, & p_1 = p_2 \\ \max(0, C - \|D(x_1) - D(x_2)\|_2), & p_1 \neq p_2 \end{cases}$$

Use the CNN outputs of our Siamese networks as descriptor



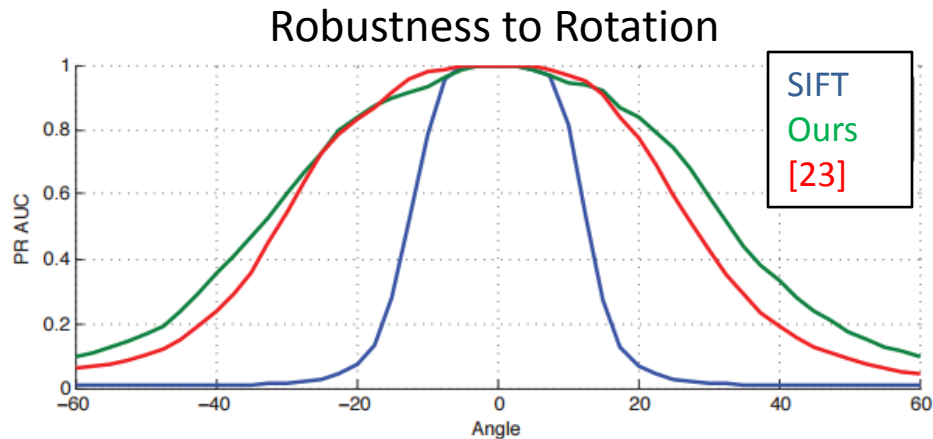
# Evaluation

Comparison of area under precision-recall curve

| Dataset | SIFT (Non-deep) | [23](Non-deep) | Ours         |
|---------|-----------------|----------------|--------------|
| ND      | 0.346           | 0.663          | <b>0.667</b> |
| TO      | 0.425           | <b>0.709</b>   | 0.545        |
| LY      | 0.226           | 0.558          | <b>0.608</b> |
| All     | 0.370           | 0.693          | <b>0.756</b> |

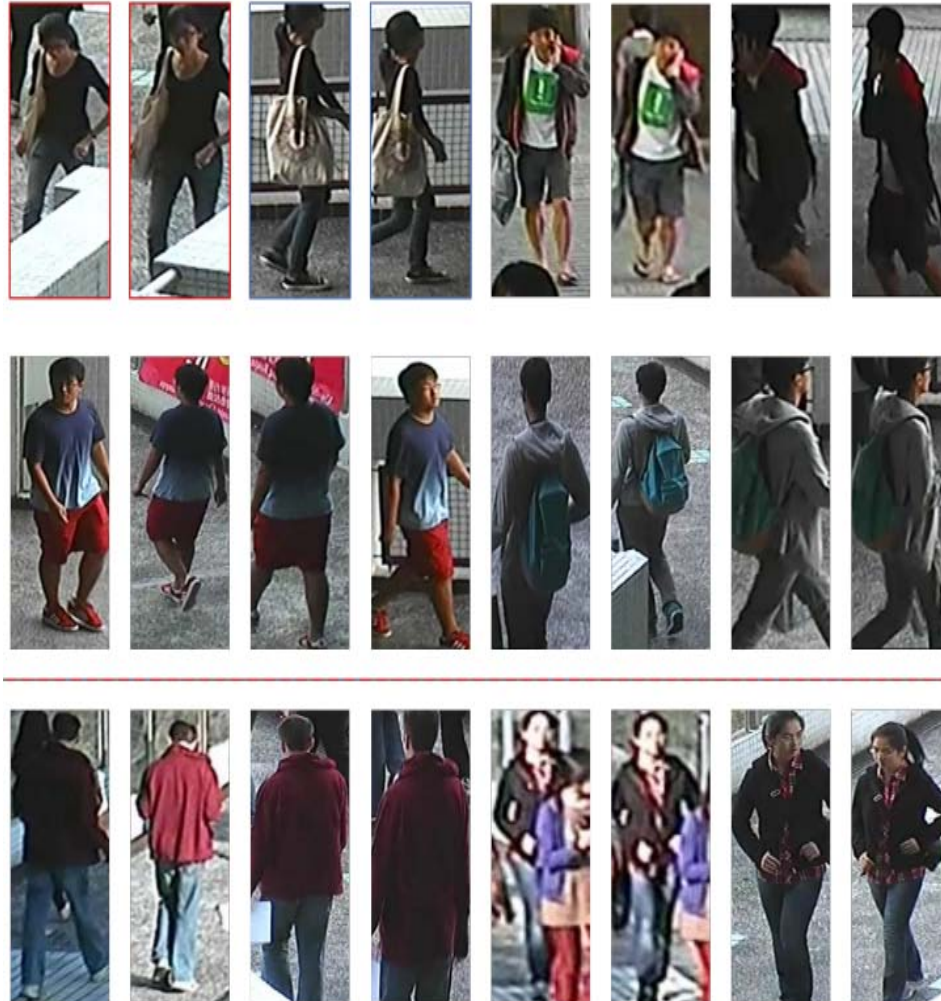
SIFT: hand-crafted features

[23]: descriptor via convex optimization





# Person Re-Identification



CUHK03 Dataset



# Quick Test

Are they the same person?





# Person Re-Identification

**True  
positive**

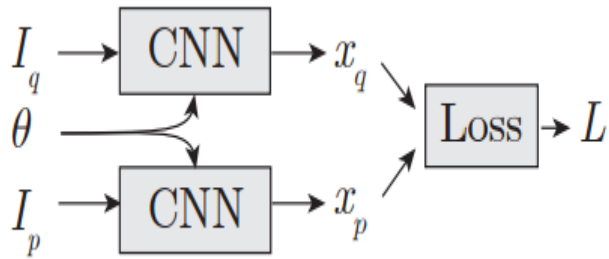


**True  
negative**





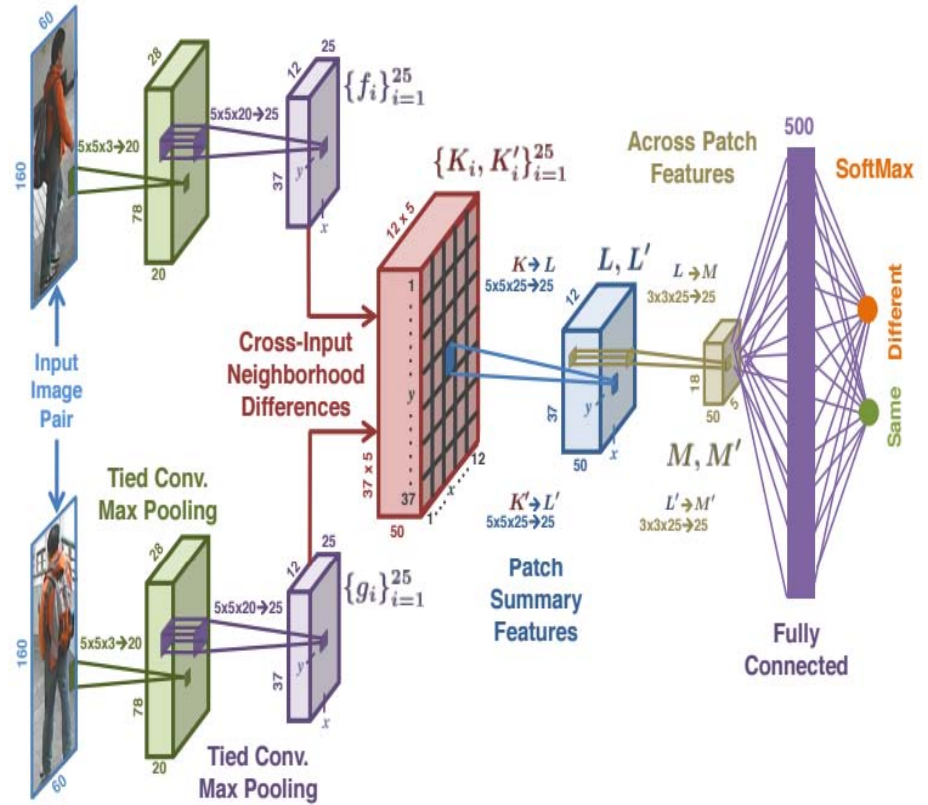
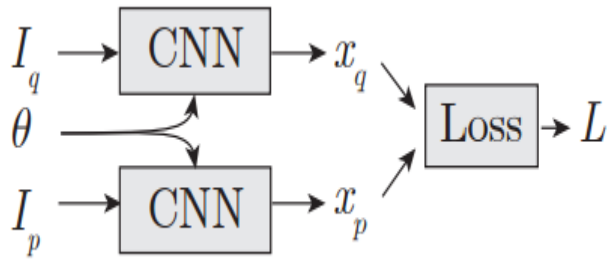
# Proposed Architecture



Ahmed, E., Jones, M. and Marks, T.K., 2015. An improved deep learning architecture for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 3908-3916).



# Proposed Architecture

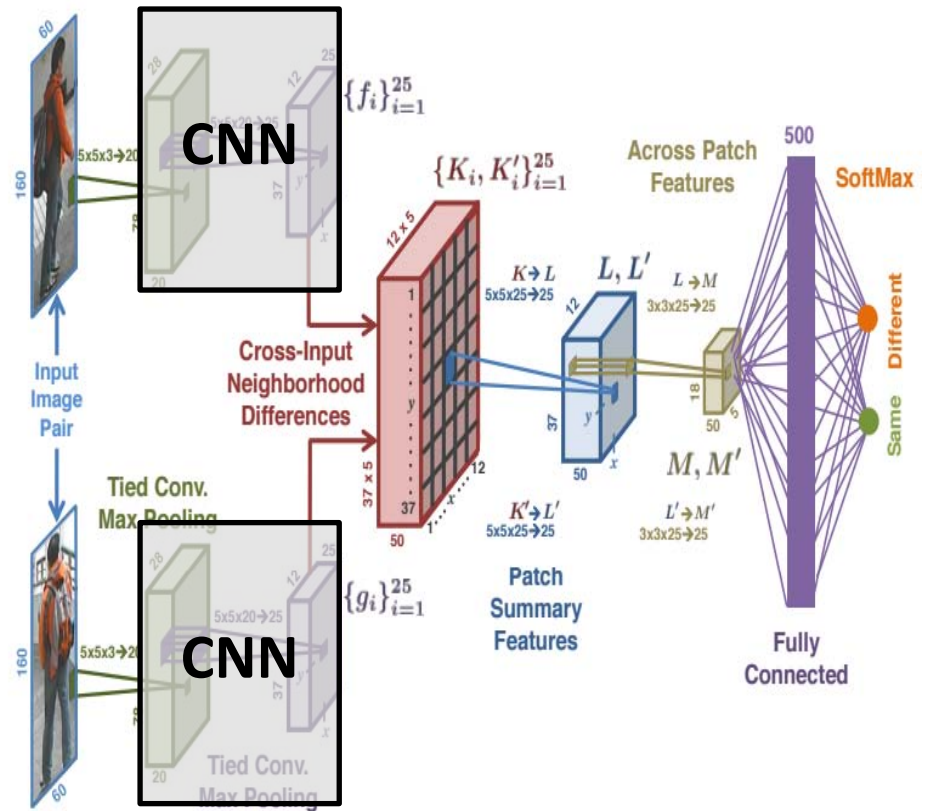
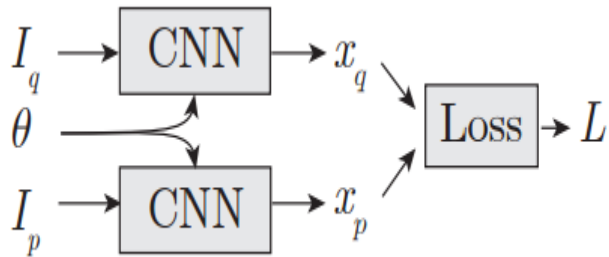


Ahmed, E., Jones, M. and Marks, T.K., 2015. An improved deep learning architecture for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 3908-3916).





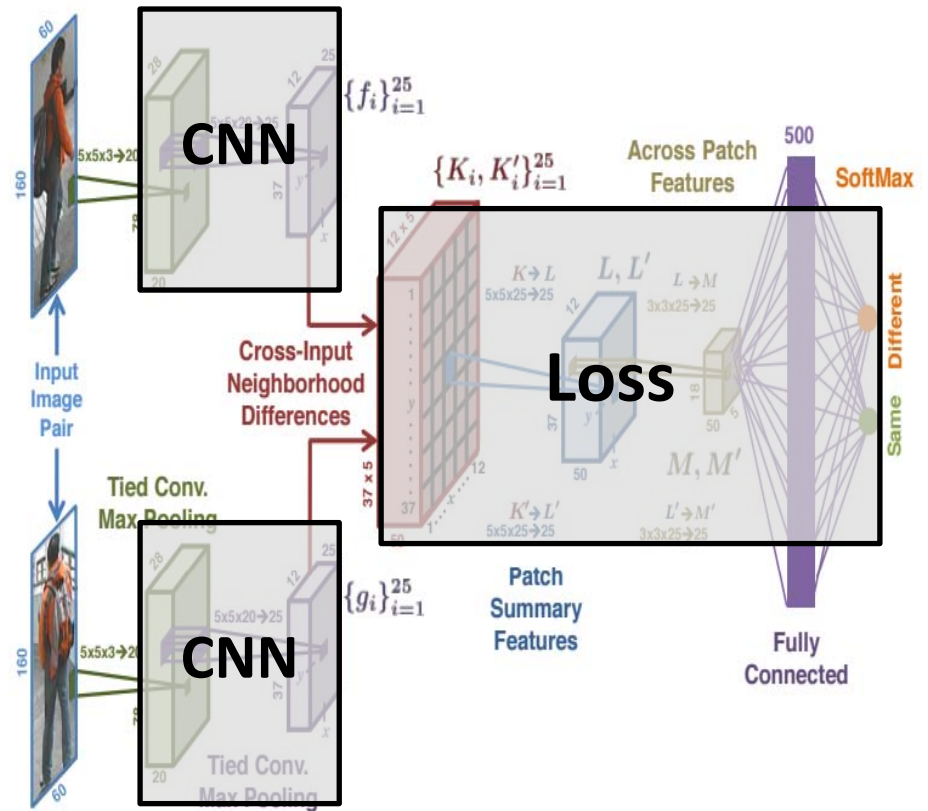
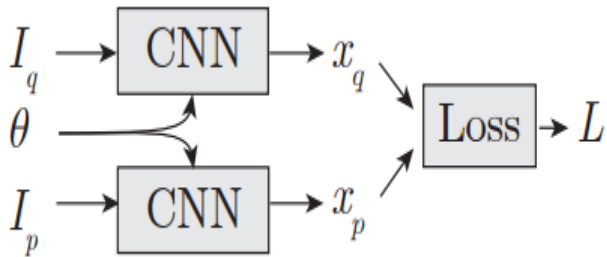
# Proposed Architecture



Ahmed, E., Jones, M. and Marks, T.K., 2015. An improved deep learning architecture for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 3908-3916).



# Proposed Architecture

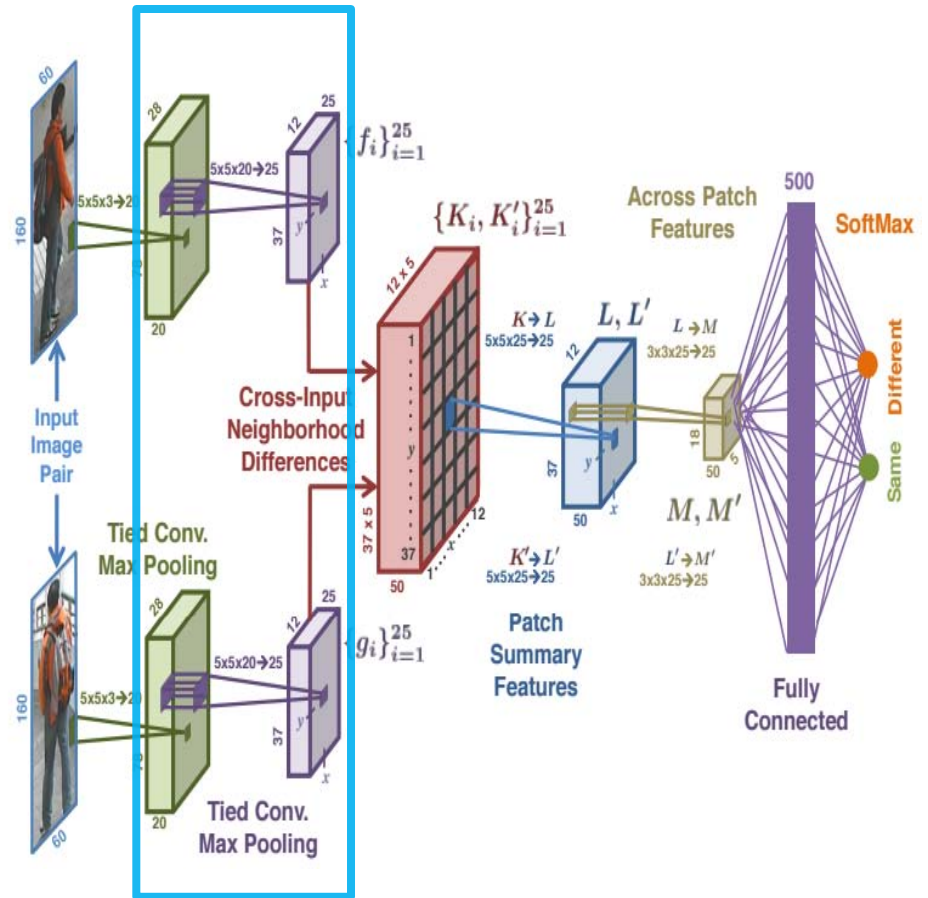


Ahmed, E., Jones, M. and Marks, T.K., 2015. An improved deep learning architecture for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 3908-3916).



# Tied Convolution

- Use convolutional layers to compute higher-order features
- Shared weights



Ahmed, E., Jones, M. and Marks, T.K., 2015. An improved deep learning architecture for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 3908-3916).

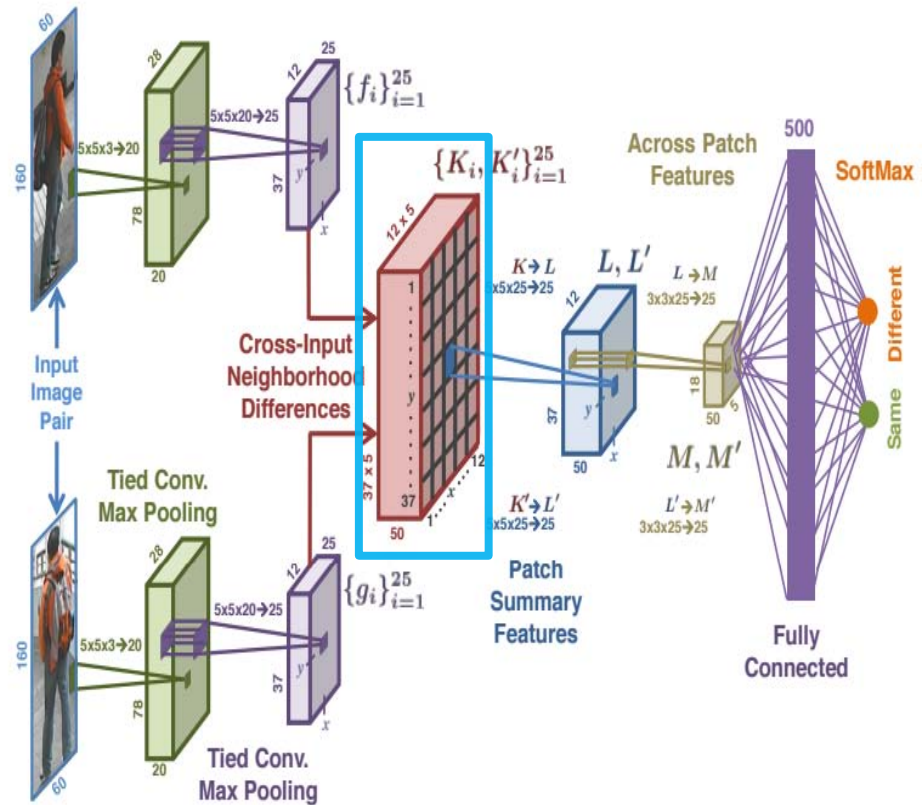


# Cross-Input Neighborhood Differences

- Compute *neighborhood difference* of two *feature maps*, instead of elementwise difference.

Example:  $f, g$  are feature maps of two input images

|     |   |   |   |     |   |   |   |
|-----|---|---|---|-----|---|---|---|
| $f$ | 5 | 7 | 2 | $g$ | 1 | 4 | 1 |
|     | 1 | 4 | 2 |     | 2 | 3 | 5 |
|     | 3 | 4 | 4 |     | 1 | 2 | 3 |





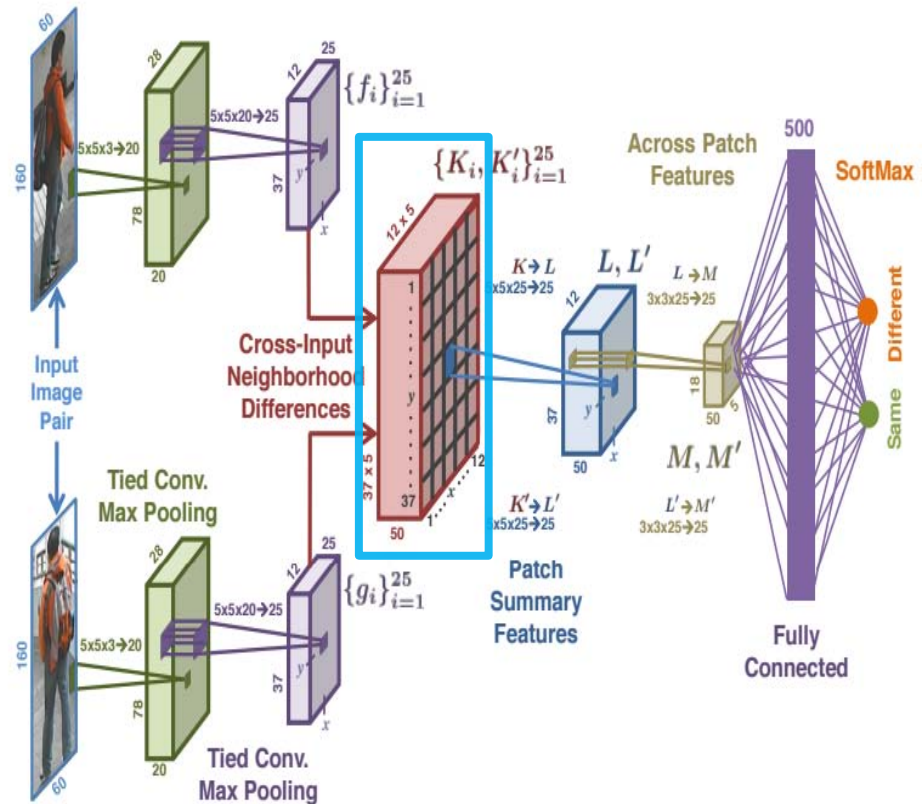
# Cross-Input Neighborhood Differences

- Compute *neighborhood difference* of two *feature maps*, instead of elementwise difference.

Example:  $f, g$  are feature maps of two input images

|     |   |   |   |     |   |   |   |
|-----|---|---|---|-----|---|---|---|
| $f$ | 5 | 7 | 2 | $g$ | 1 | 4 | 1 |
|     | 1 | 4 | 2 |     | 2 | 3 | 5 |
|     | 3 | 4 | 4 |     | 1 | 2 | 3 |

$$K(1,1) = \begin{bmatrix} 5 & 5 \\ 5 & 5 \end{bmatrix} - \begin{bmatrix} 1 & 4 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 4 & 4 \\ 3 & 2 \end{bmatrix}$$





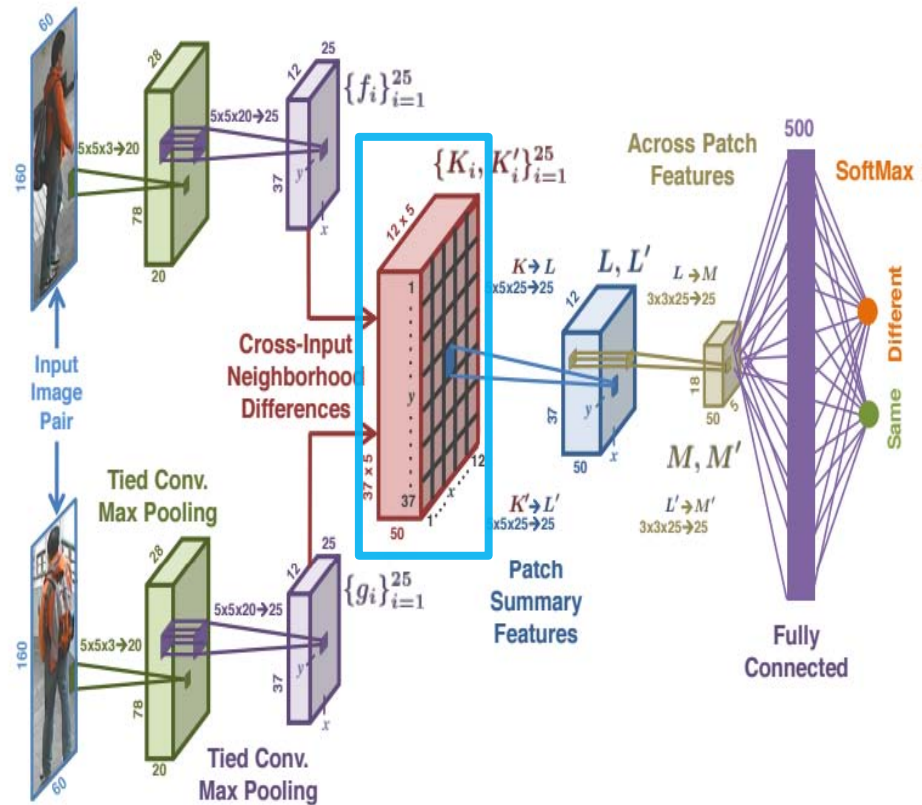
# Cross-Input Neighborhood Differences

- Compute *neighborhood difference* of two *feature maps*, instead of elementwise difference.
- A neighborhood-patch size of 5 was used in the paper:

$$K_i(x,y) = f_i(x,y)I(5,5) - N[g_i(x,y)]$$

where  $I(5,5)$  is a 5x5 matrix of 1s,  $N[g_i(x,y)]$  is the 5x5 neighborhood of  $g_i$  centered at  $(x,y)$

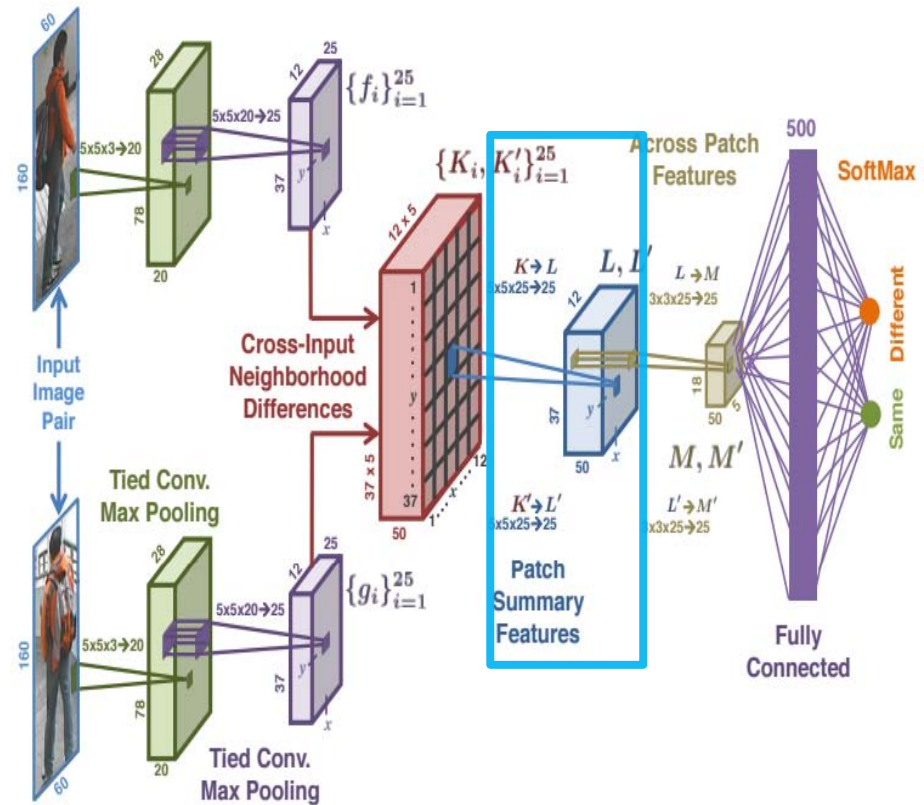
- Another neighborhood difference map  $K'$  was also computed where  $f$  and  $g$  were revised.





# Patch Summary Features

- Convolutional layers with 5x5 filters and stride 5 (the size of neighborhood patch).
- Provides a high-level summary of the cross-input differences in a neighborhood patch.

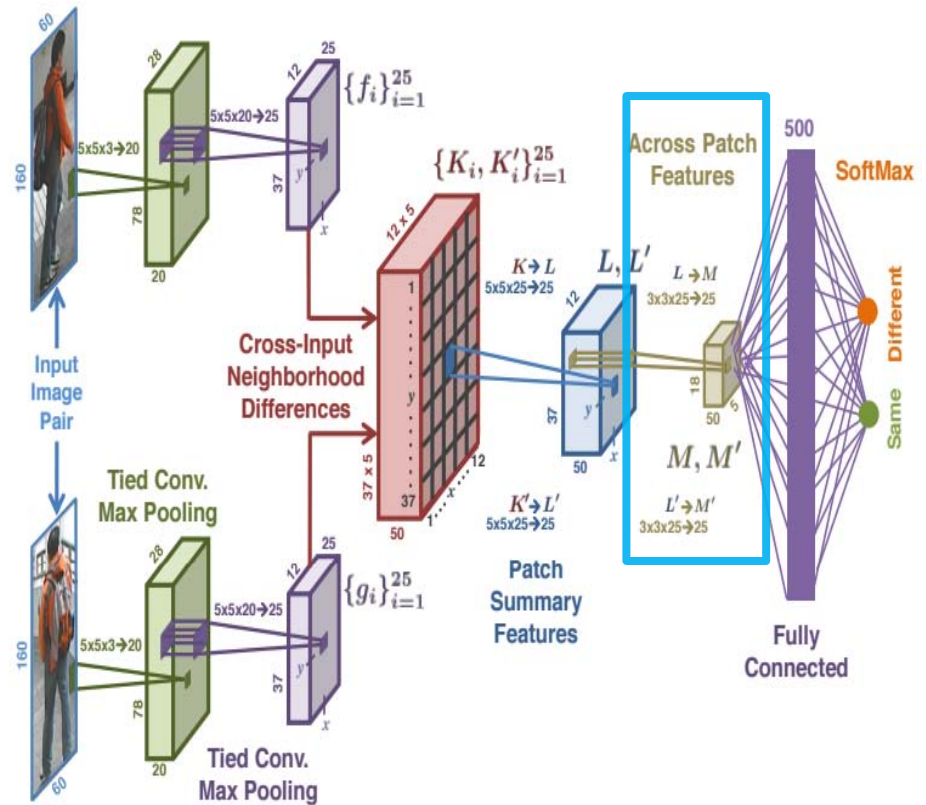


Ahmed, E., Jones, M. and Marks, T.K., 2015. An improved deep learning architecture for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 3908-3916).



# Across-Patch Features

- Convolutional layers with 3x3 filters and stride 1.
- Learn spatial relationships across neighborhood differences

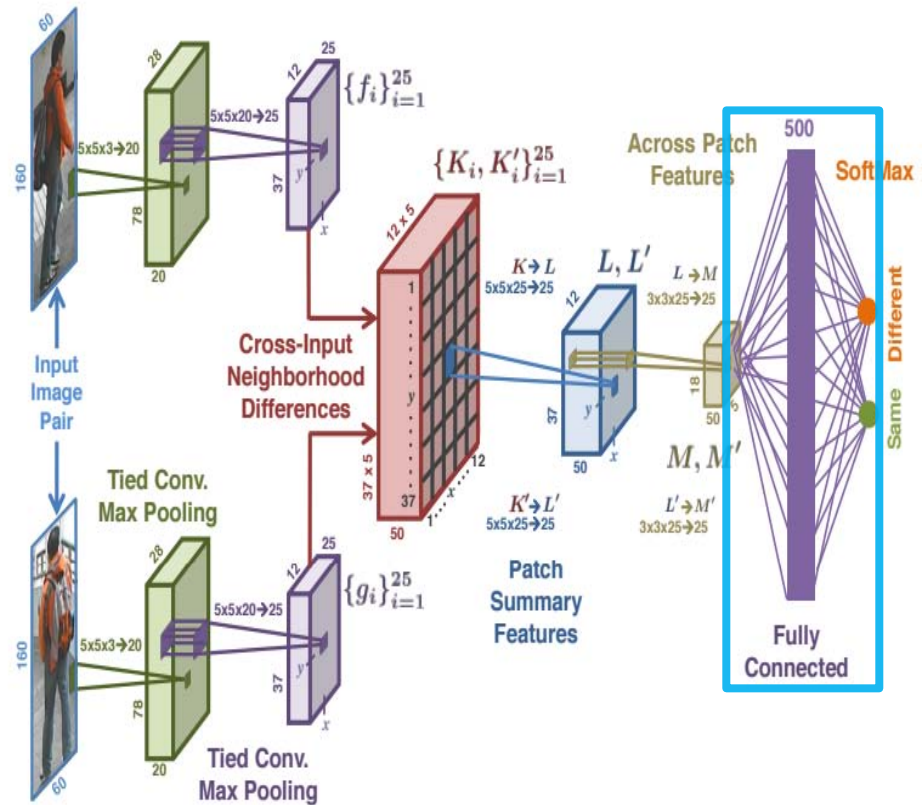






# Across-Patch Features

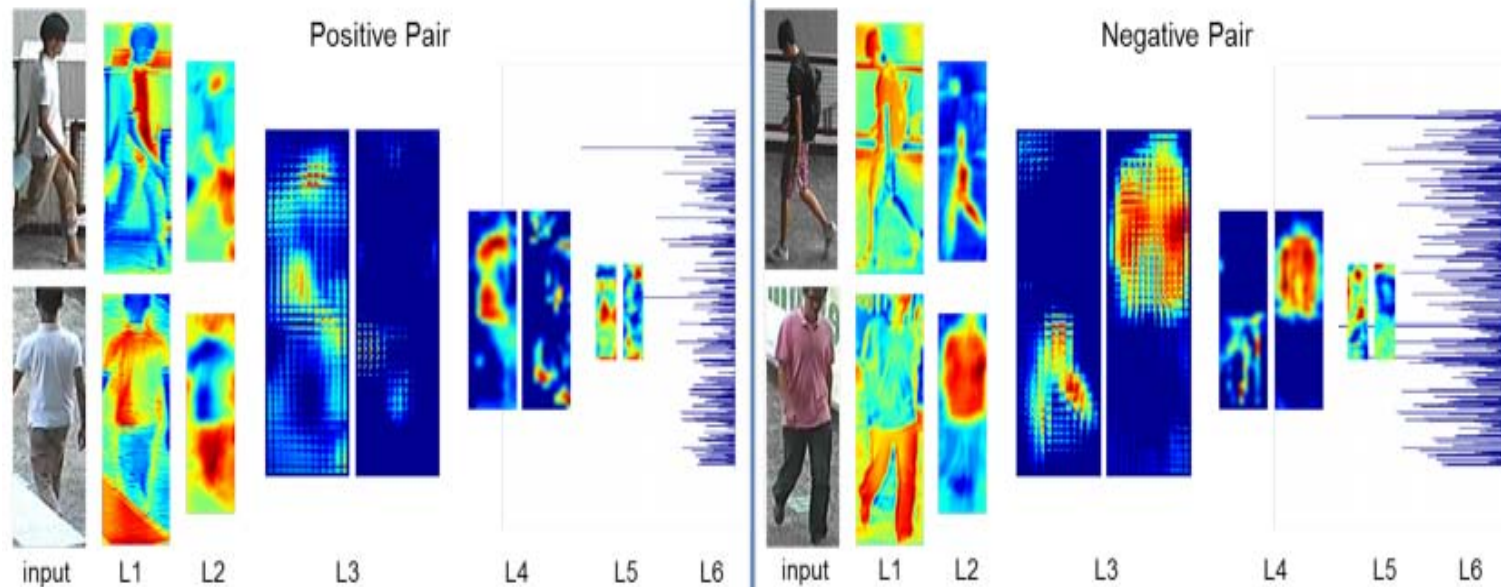
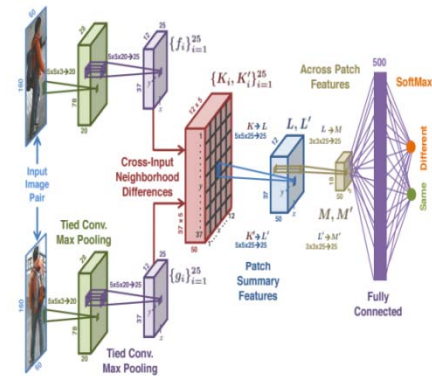
- Fully connected layer.
- Combine information from patches that are far from each other.
- Output: 2 softmax units



Ahmed, E., Jones, M. and Marks, T.K., 2015. An improved deep learning architecture for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 3908-3916).



# Visualization of Learned Features



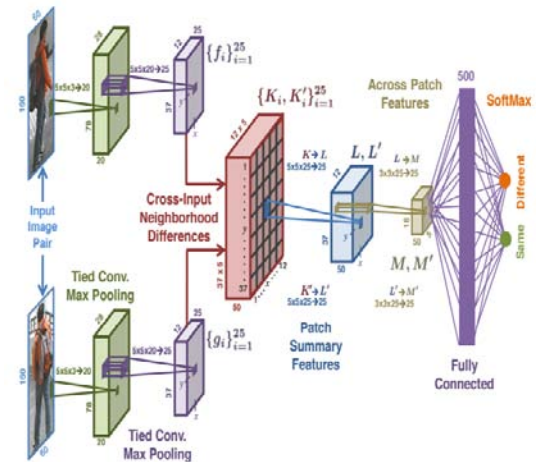
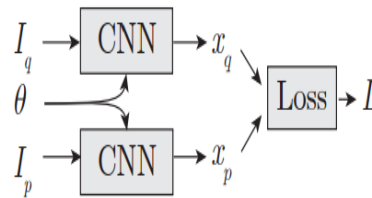
Ahmed, E., Jones, M. and Marks, T.K., 2015. An improved deep learning architecture for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 3908-3916).



# Evaluation

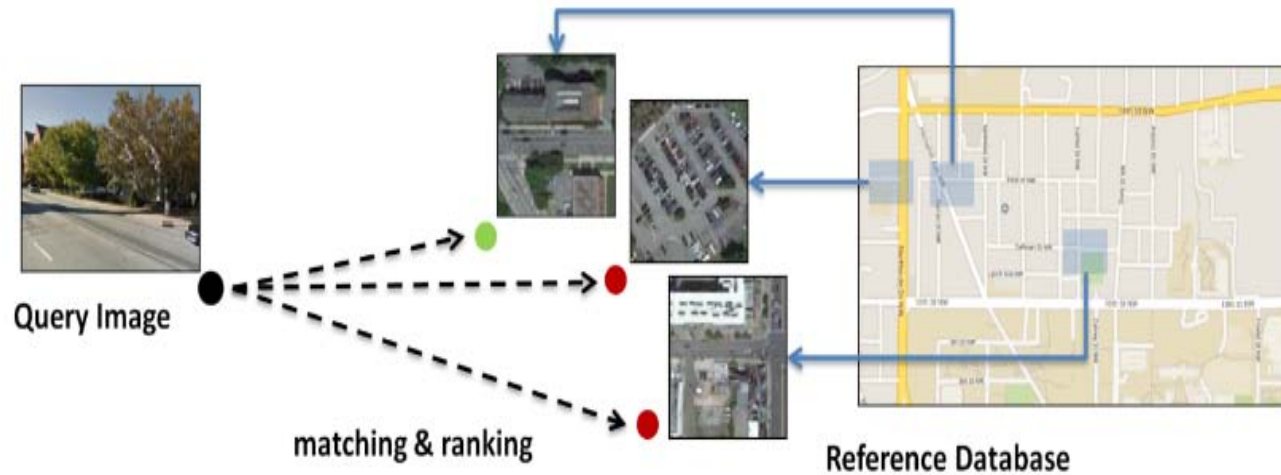
| Method              | Elementwise Difference | Neighborhood Difference |
|---------------------|------------------------|-------------------------|
| Identification rate | 27.66%                 | 54.74%                  |

| Method              | Regular Siamese Network | This work |
|---------------------|-------------------------|-----------|
| Identification rate | 42.19%                  | 54.74%    |





# Street-View to Overhead-View Image Matching





# Street-View to Overhead-View Image Matching

Query:



Matching  
Image:



# Quick Test

Which one is the correct match?

Query Image



A



B



C



D



E





# CNN Architectures

Classification CNN:



$$L(A, B, l) = \text{LogLossSoftMax}(f(I), l)$$

$I = \text{concatenation}(A, B)$

$f = \text{AlexNet}$

$l = \{0, 1\}, \text{label}$



# CNN Architectures

## Classification CNN:



$$L(A, B, l) = \text{LogLossSoftMax}(f(I), l)$$

$I = \text{concatenation}(A, B)$

$f = \text{AlexNet}$

$l = \{0, 1\}, \text{label}$

## Siamese-like CNN:



$$L(A, B, l) = l * D + (1 - l) * \max(0, m - D)$$

$D = \|f(A) - f(B)\|_2$

$m = \text{margin parameter}$





# CNN Architectures

## Classification CNN:



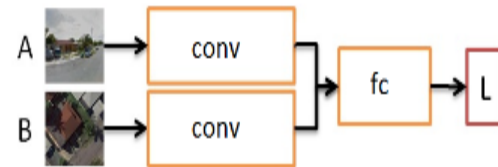
$$L(A, B, l) = \text{LogLossSoftMax}(f(I), l)$$

$I = \text{concatenation}(A, B)$

$f = \text{AlexNet}$

$l = \{0, 1\}, \text{label}$

## Siamese-classification hybrid network:



$$L(A, B, l) = \text{LogLossSoftMax}(f_{fc}(I_{conv}), l)$$

$I_{conv} = \text{concatenation}(f_{conv}(A), f_{conv}(B))$

## Siamese-like CNN:



$$L(A, B, l) = l * D + (1 - l) * \max(0, m - D)$$

$D = \|f(A) - f(B)\|_2$

$m = \text{margin parameter}$



# CNN Architectures

## Classification CNN:



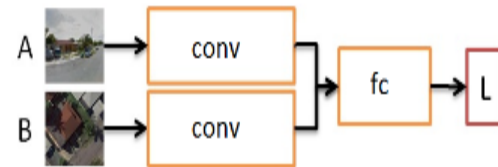
$$L(A, B, l) = \text{LogLossSoftMax}(f(I), l)$$

$I = \text{concatenation}(A, B)$

$f = \text{AlexNet}$

$l = \{0, 1\}, \text{label}$

## Siamese-classification hybrid network:



$$L(A, B, l) = \text{LogLossSoftMax}(f_{fc}(I_{conv}), l)$$

$I_{conv} = \text{concatenation}(f_{conv}(A), f_{conv}(B))$

## Siamese-like CNN:

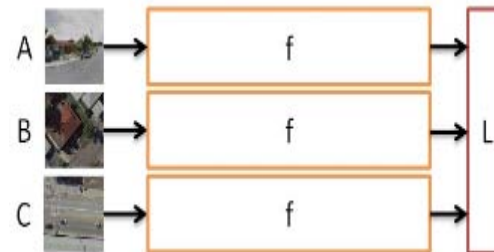


$$L(A, B, l) = l * D + (1 - l) * \max(0, m - D)$$

$D = \|f(A) - f(B)\|_2$

$m = \text{margin parameter}$

## Triplet network CNN:



$$L(A, B, C) = \max(0, m + D(A, B) - D(A, C))$$

$(A, B)$  is a match pair

$(A, C)$  is a non-match pair



# Distance-based Logistic Loss

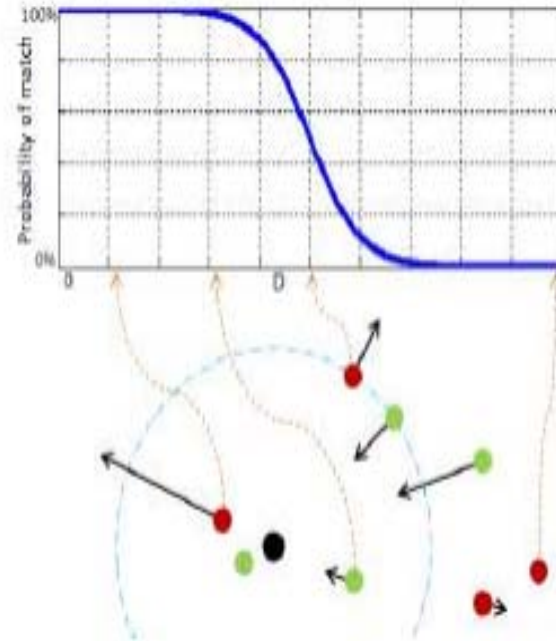
$$p(A, B) = \frac{1 + \exp(-m)}{1 + \exp(D - m)}$$

$$L(A, B, l) = \text{LogLoss}(p(A, B), l)$$

where

$$D = \|f(A) - f(B)\|_2$$

$m = \text{margin parameter}$



**Matched/Nonmatched**  
instances are pushed away from  
the “boundary” in the  
**inward/outward** direction.



# Performance of Different Networks

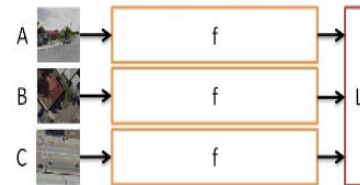
Matching accuracy

| Test set | Denver      | Detroit     | Seattle     |
|----------|-------------|-------------|-------------|
| Siamese  | 85.6        | 83.2        | 82.9        |
| Triplet  | <b>88.8</b> | <b>86.8</b> | <b>86.4</b> |

Siamese-like CNN:



Triplet network CNN:



## Observation 1:

- Triplet network outperforms the Siamese by a large margin



# Performance of Different Networks

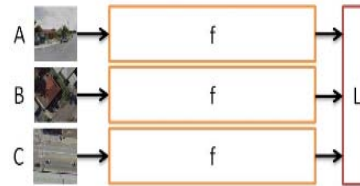
## Matching accuracy

| Test set    | Denver      | Detroit     | Seattle     |
|-------------|-------------|-------------|-------------|
| Siamese     | 85.6        | 83.2        | 82.9        |
| Siamese-DBL | <b>90.0</b> | <b>88.0</b> | <b>88</b>   |
| Triplet     | 88.8        | 86.8        | 86.4        |
| Triplet-DBL | <b>90.2</b> | <b>88.4</b> | <b>87.6</b> |

Siamese-like CNN:



Triplet network CNN:



Distance-based logistic (DBL) loss:

$$p(A, B) = \frac{1 + \exp(-m)}{1 + \exp(D - m)}$$
$$L(A, B, l) = \text{LogLoss}(p(A, B), l)$$

### Observation 2:

- Distance-based logistic (DBL) Nets significantly outperform the original network.



# Performance of Different Networks

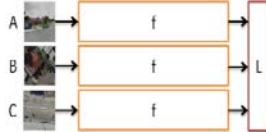
## Matching accuracy

| Test set           | Denver      | Detroit     | Seattle     |
|--------------------|-------------|-------------|-------------|
| Siamese Net        | 85.6        | 83.2        | 82.9        |
| Triplet Net        | 88.8        | 86.8        | 86.4        |
| Classification Net | <b>90.0</b> | <b>87.8</b> | <b>87.7</b> |
| Hybrid Net         | <b>91.5</b> | <b>88.7</b> | <b>89.4</b> |

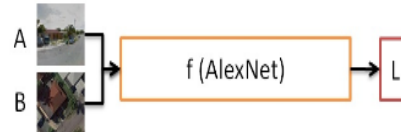
Siamese-like CNN:



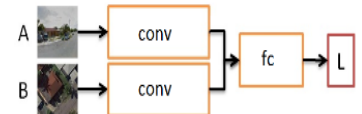
Triplet network CNN:



Classification CNN:



Classification-siamese hybrid:



### Observation 3:

- Classification networks achieved better accuracy than Siamese and triplet networks.
- Jointly extract and exchange information from both input images.

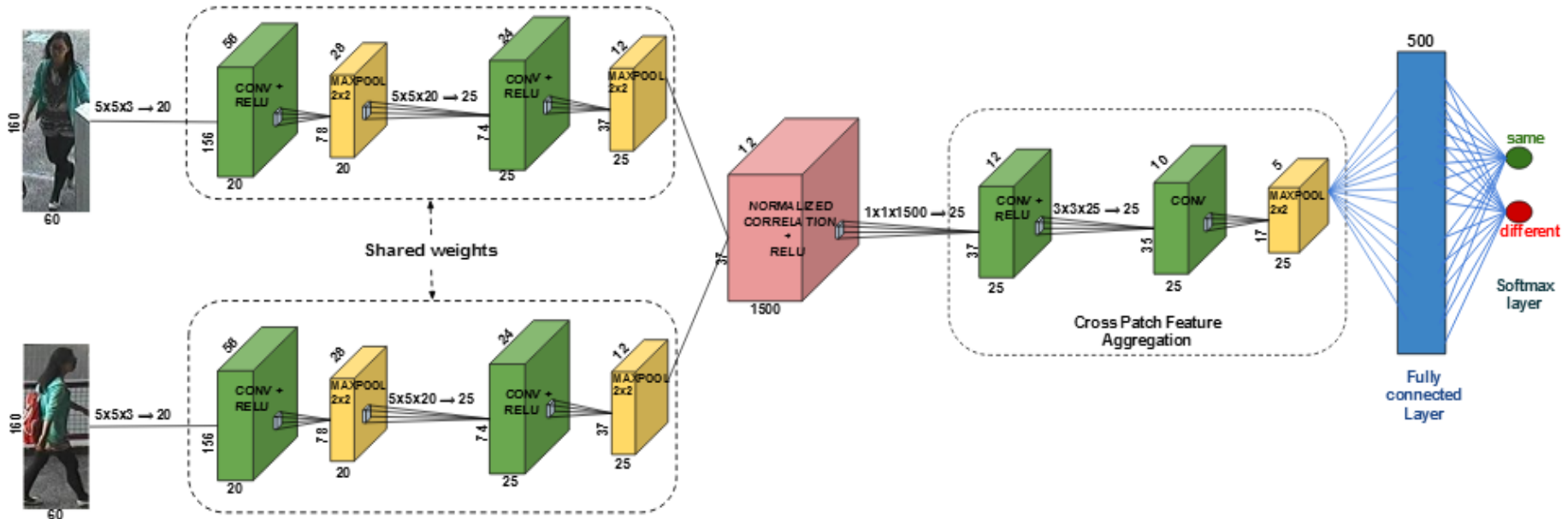


# **MORE VARIANTS OF SIAMESE CNNs**



# Siamese CNN – Variants

## SIAMESE CNN – INTERMEDIATE MERGING



- Combining at an **intermediate stage** allows us to capture patch-level variability.
- Performing inexact (soft) matching yields superior performance.  $Match(X, Y) = (X - \mu_X)(Y - \mu_Y) / \sigma_X \sigma_Y$





# Siamese CNN – Variants

## SIAMESE CNN – INTERMEDIATE MERGING

Results:

- Handling Partial Occlusion:



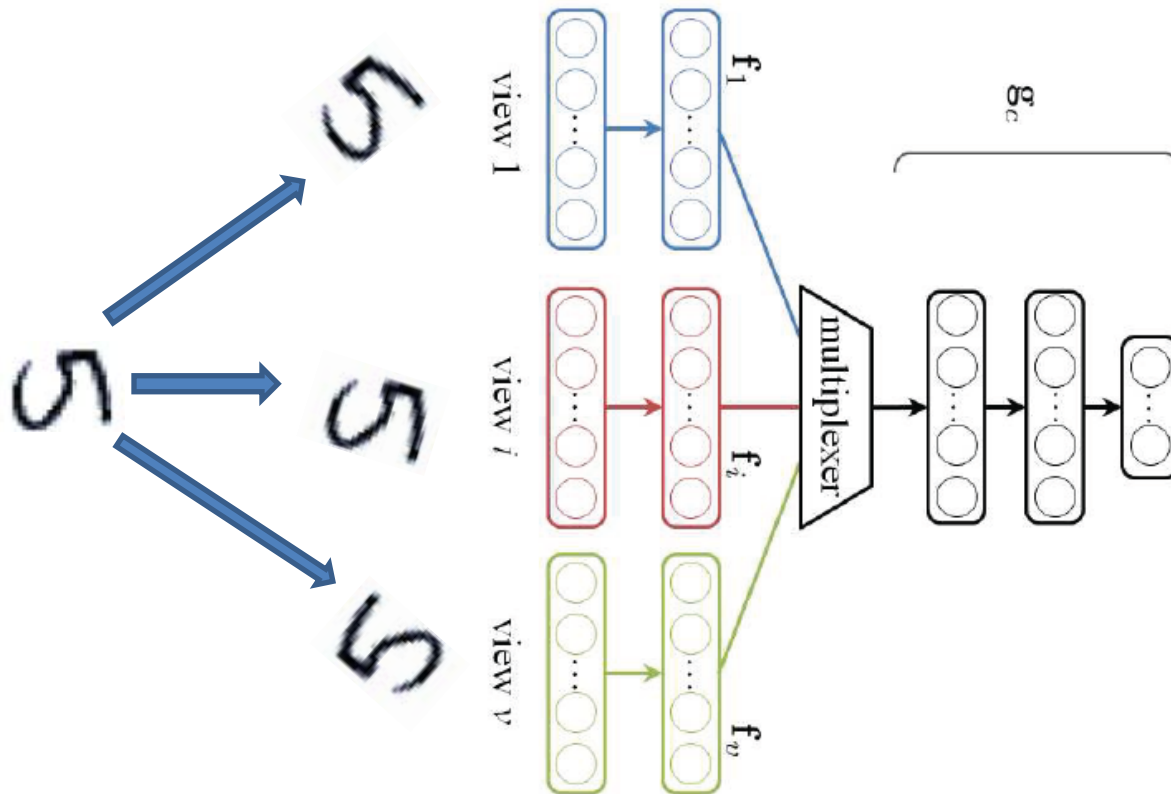
Baseline:

Proposed Method:



# Siamese CNN – Variants

## SIAMESE CNN – FOR VIEWPOINT INVARIANCE



**Viewpoint** invariance is incorporated by considering the similarity of response across the individual streams.



# Siamese CNN – Variants

## SIAMESE CNN – FOR VIEWPOINT INVARIANCE

Results on the CMU MultiPIE Dataset, for recognition across 7 poses.

| Methods         | -45 deg     | -30 deg     | -15 deg     | 15deg       | 30 deg      | 45 deg      |
|-----------------|-------------|-------------|-------------|-------------|-------------|-------------|
| CCA             | 0.73        | 0.96        | 1.00        | 0.99        | 0.96        | 0.69        |
| KCCA (RBF)      | 0.80        | 0.98        | 0.99        | 1.00        | 0.98        | 0.72        |
| FIP+LDA         | 0.93        | 0.96        | 1.00        | 0.99        | 0.96        | 0.90        |
| MVP+LDA         | 0.93        | <b>1.00</b> | 1.00        | 1.00        | 0.99        | 0.96        |
| <b>Proposed</b> | <b>0.99</b> | 0.99        | <b>1.00</b> | <b>1.00</b> | <b>0.99</b> | <b>0.98</b> |

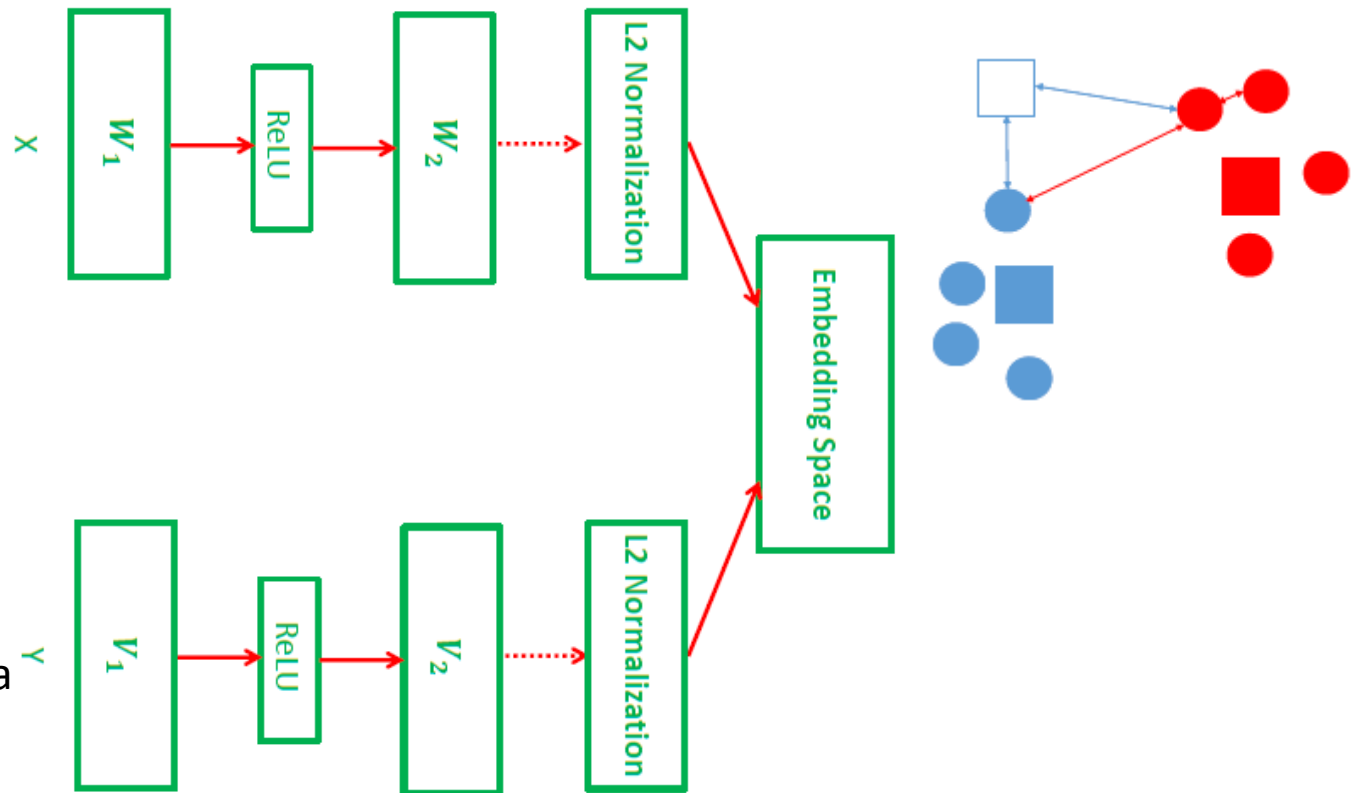


# Siamese CNN – Variants

## TWO STREAM CNN – FOR CROSS-MODAL EMBEDDING



Man in black shirt playing a guitar



Two stream networks have also been used for cross-modal embedding tasks. Here inputs from different modalities are mapped to a common space.

Wang, L., Li, Y. and Lazebnik, S., 2016. Learning deep structure-preserving image-text embeddings. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 5005-5013).



# Siamese CNN - Variants

Application: Sentence completion, response to tweet, paraphrase identification

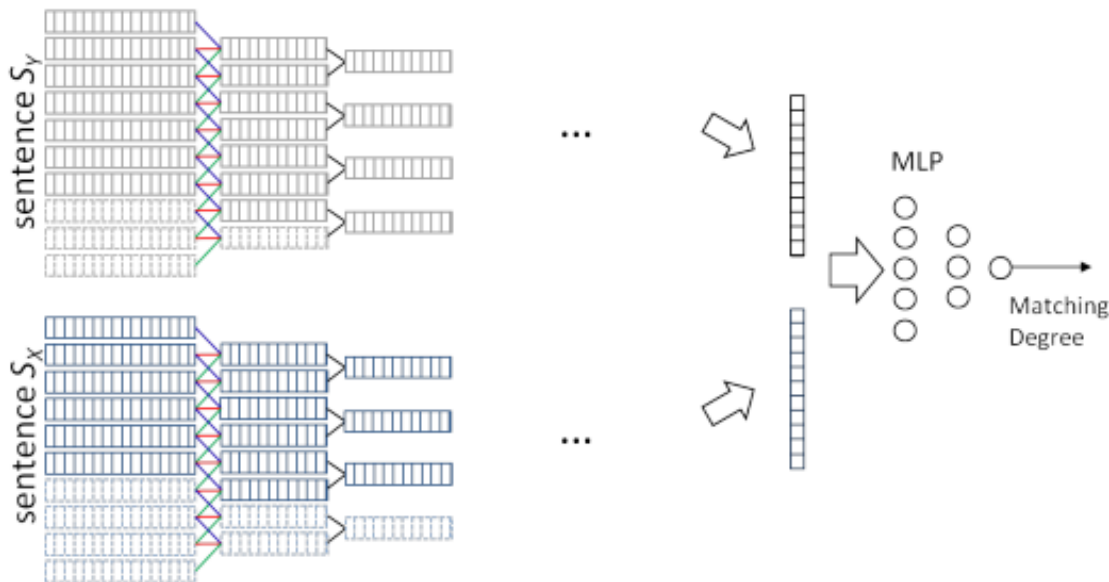
Example:

$x$  : Damn, I have to work overtime this weekend!

$y^+$ : Try to have some rest buddy.

$y^-$ : It is hard to find a job, better start polishing your resume.

word2vec



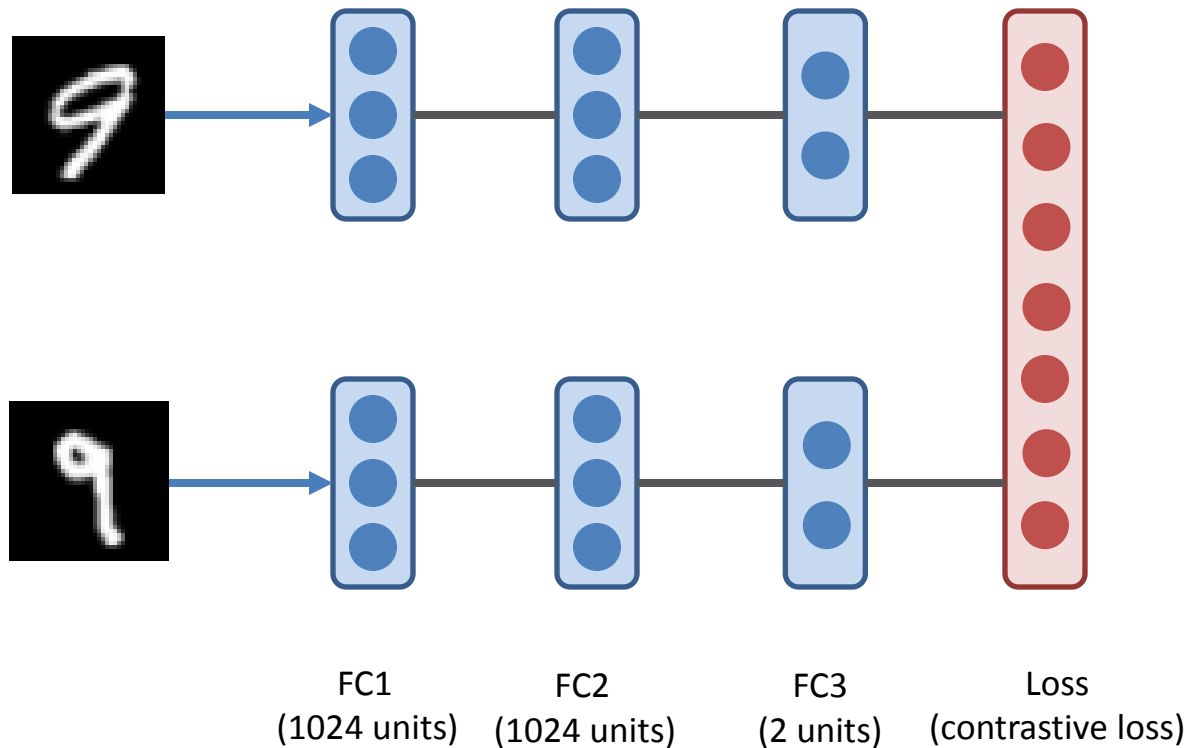


# DEMO OF SIAMESE NETWORK



# Demo: Architecture

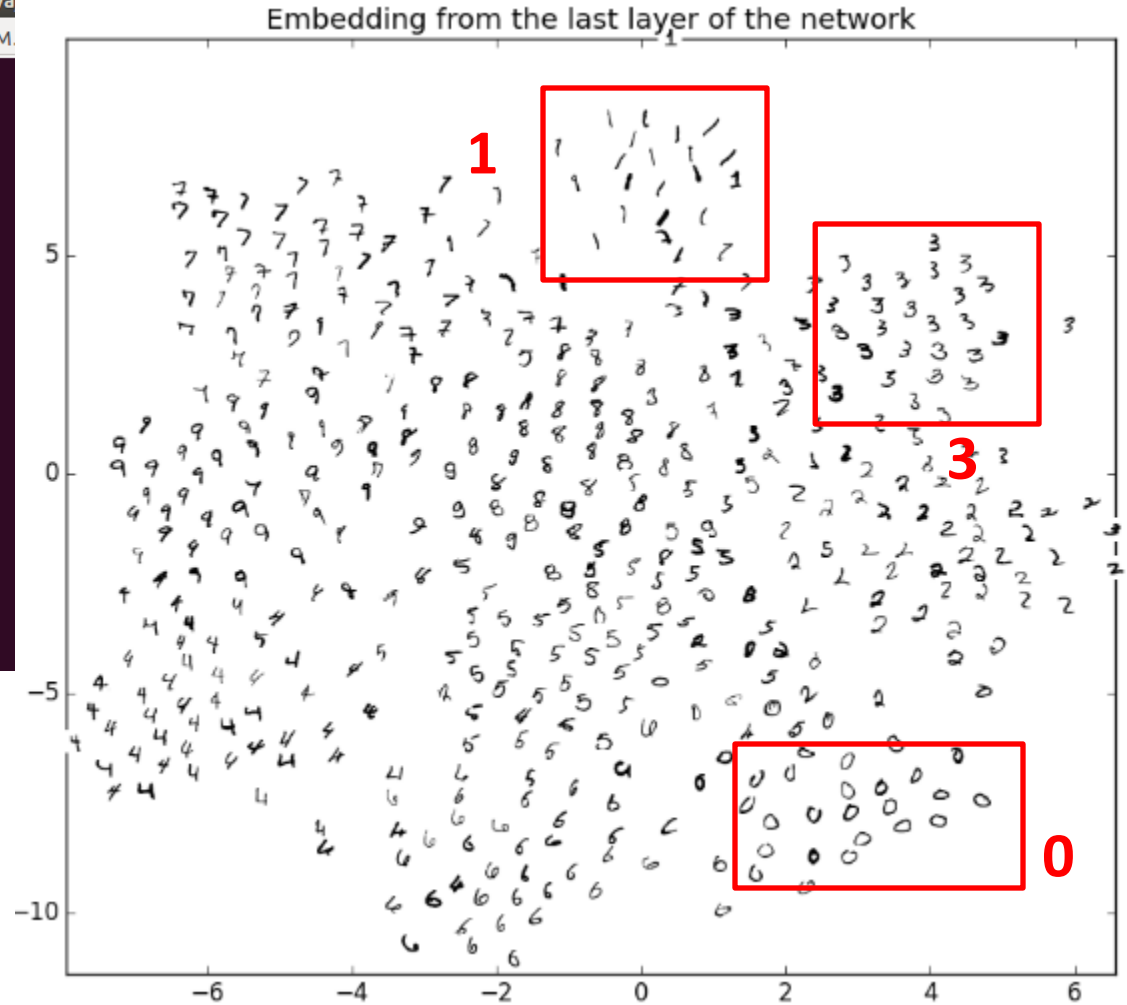
## MNIST Digit Similarity Assessment





# Demo: Results

```
moitreya@moitreya-pc: /media/moitreya/RAID_1T/Users/Moitreya
moitreya@moitreya-pc: /media/moitreya/RAID_1T/Users/M.
step 99390: loss 0.007
step 99400: loss 0.008
step 99410: loss 0.007
step 99420: loss 0.007
step 99430: loss 0.028
step 99440: loss 0.011
step 99450: loss 0.013
step 99460: loss 0.010
step 99470: loss 0.008
step 99480: loss 0.006
step 99490: loss 0.009
step 99500: loss 0.007
step 99510: loss 0.012
step 99520: loss 0.019
step 99530: loss 0.009
step 99540: loss 0.014
step 99550: loss 0.013
step 99560: loss 0.013
step 99570: loss 0.013
step 99580: loss 0.012
step 99590: loss 0.016
step 99600: loss 0.010
step 99610: loss 0.006
step 99620: loss 0.008
```







# Summary

- Quantifying “similarity” is an essential component of data analytics.
- Deep Learning approaches, such as “Siamese” Convolution Neural Nets, have shown promise recently.
- Several variants of Siamese CNN are available for making our life easier for a variety of tasks.



# Reading List

- Bell, Sean, and Kavita Bala, [Learning visual similarity for product design with convolutional neural networks](#), ACM Transactions on Graphics (TOG), 2015
- Chopra, Sumit, Raia Hadsell, and Yann LeCun, [Learning a similarity metric discriminatively, with application to face verification](#), CVPR 2005
- Zagoruyko, Sergey, and Nikos Komodakis, [Learning to compare image patches via convolutional neural networks](#), CVPR 2015
- Hoffer, Elad, and Nir Ailon, [Deep metric learning using triplet network](#), arXiv:1412.6622
- Simo-Serra, Edgar, et al., [Discriminative Learning of Deep Convolutional Feature Point Descriptors](#), ICCV 2015
- Vo, Nam N., and James Hays, [Localizing and Orienting Street Views Using Overhead Imagery](#), ECCV 2016
- Ahmed, Ejaz, Michael Jones, and Tim K. Marks, [An Improved Deep Learning Architecture for Person Re-Identification](#), CVPR 2015
- Hu, Baotian, et al., [Convolutional neural network architectures for matching natural language sentences](#), NIPS 2014
- Kulis, Brian, [Metric learning: A survey](#), Foundations and Trends in Machine Learning, 2013
- Su, Hang, et al., [Multi-view convolutional neural networks for 3d shape recognition](#), ICCV 2015
- Zheng, Yi, et al., [Time Series Classification Using Multi-Channels Deep Convolutional Neural Networks](#), WAIM 2014
- Yi, Kwang Moo, et al., [LIFT: Learned Invariant Feature Transform](#), arXiv:1603.09114
- Stricker, M.A. and Orengo, M. [Similarity of color images](#). In *IS&T/SPIE's Symposium on Electronic Imaging: Science & Technology* (pp. 381-392), 1995.



**Appreciate your kind attention!**