

Objectives								
Problem definition	 Person Re-Identification is the task of identifying a particular person across multiple images captured from the same/different cameras, from different viewpoints, at the same/different points in time. Two Deep Convolutional Neural Network (DCNN) architectures: 							
	 1 A Siamese network with novel Normalized correlation (Inexact matching) layer and Wider search space 2 A Fused model using Normalized correlation layer and a state of the art exact matching technique (Ahmed et al. CVPR-2015) 							

Challenges in Person Re-Identification





Viewpoint Change



Illumination Variation





Partial Occlusion

Prior approaches

Deep learning based approaches:

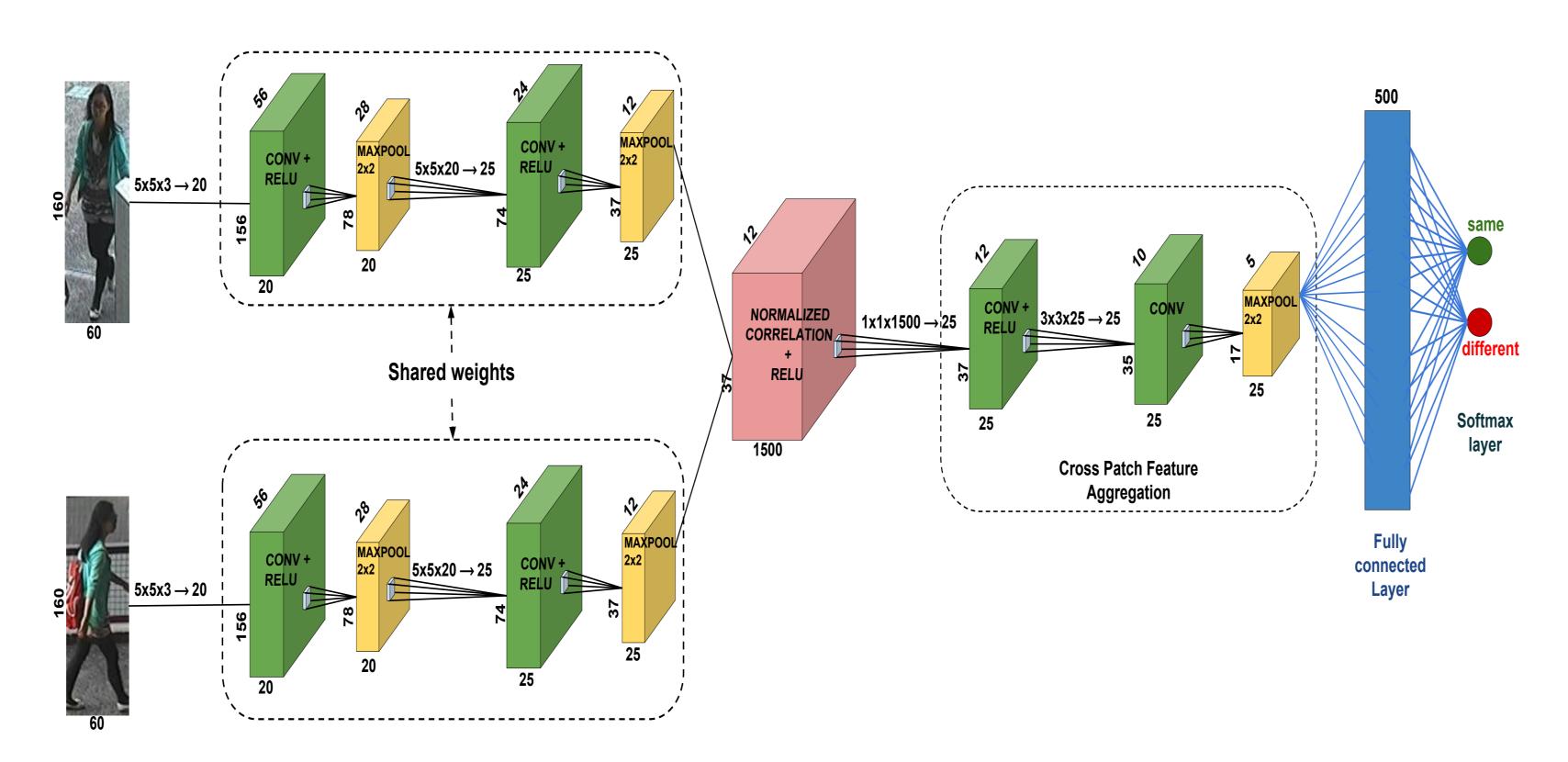
Yi et al.(ICPR-2014)	: Two-input network $\rightarrow 3$ stages of convolution (shared
	weights) \rightarrow Cosine similarity score
Li et al.(CVPR-2014)	: Two-input network $\rightarrow 1$ stage convolution (non-shared
	weights) \rightarrow take the product of the responses obtained from
	first set of convolutions corresponding to the two inputs \rightarrow
	Maxout grouping $\rightarrow 1$ stage convolution (shared weights)
	\rightarrow Softmax classifier (same / different)
Ahmed et al.(CVPR-2015)	: Siamese network $\rightarrow 2$ stages of convolution (shared weights)
	\rightarrow Cross-input neighborhood matching \rightarrow 2 stages of con-
	volution (non-shared weights) \rightarrow Softmax classifier (same /
	different)

Solutions proposed in our work							
Challenge	Solution						
Illumination variation	Normalized correlation between pixel's patch neigh-						
	borhood						
Pose / viewpoint variations	Wider search space						
Partial occlusion	Normalized correlation (Inexact matching) +						
	Wider search space						

Deep Neural Networks with Inexact Matching for Person Re-Identification

Arulkumar Subramaniam, Moitreya Chatterjee, Anurag Mittal Indian Institute of Technology Madras

NormXcorr model ($\sim 1.12M$ parameters)



Normalized correlation (Inexact Matching) layer

Given two corresponding input feature maps $X(12 \times 37)$ and $Y(12 \times 37)$ after first two convolution layers, we compute the normalized correlation as follows.

- We start with every pixel of X located at (x, y), where x is along the width and y along the height (denoted as X(x, y)).
- We then create two matrices. The first is a 5×5 matrix representing the 5×5 neighborhood of X(x, y), while the second is the corresponding 5×5 neighborhood of Y centered at (a, b), where $1 \le a \le 12$ and $y - 2 \le b \le y + 2$.
- We perform **inexact matching** over a **wider search space**, by computing a Normalized Correlation between the two patch matrices.

Given two matrices, E and F, whose elements are arranged as two N-dimensional vectors, Normalized Correlation is given by:

$$normxcorr(E, F) = \frac{\sum_{i=1}^{N} (E_i - \mu_E) * (F_i - \mu_F)}{(N-1) * \sigma_E * \sigma_F},$$

where μ_E, μ_F denotes the means of the elements of the 2 matrices E and F respectively, while σ_E, σ_F denotes their respective unbiased standard deviation (a small $\epsilon = 0.01$ is added to the unbiased standard deviation to avoid division by 0).

- The mean of a N-dimensional vector E, $\mu_E = \frac{\sum_{i=1}^{N} E_i}{N}$
- The unbiased standard deviation of a N-dimensional vector E, $\sigma_E = \sqrt{\frac{\sum_{i=1}^{N} (E_i \mu_E)^2}{N-1}}$

For every pair of feature maps X and Y, this gives us 60 feature maps of dimension 12×37 each. For all 25 pairs of maps that are input to the Normalized Correlation layer, we obtain an output of 1500, 12×37 maps.

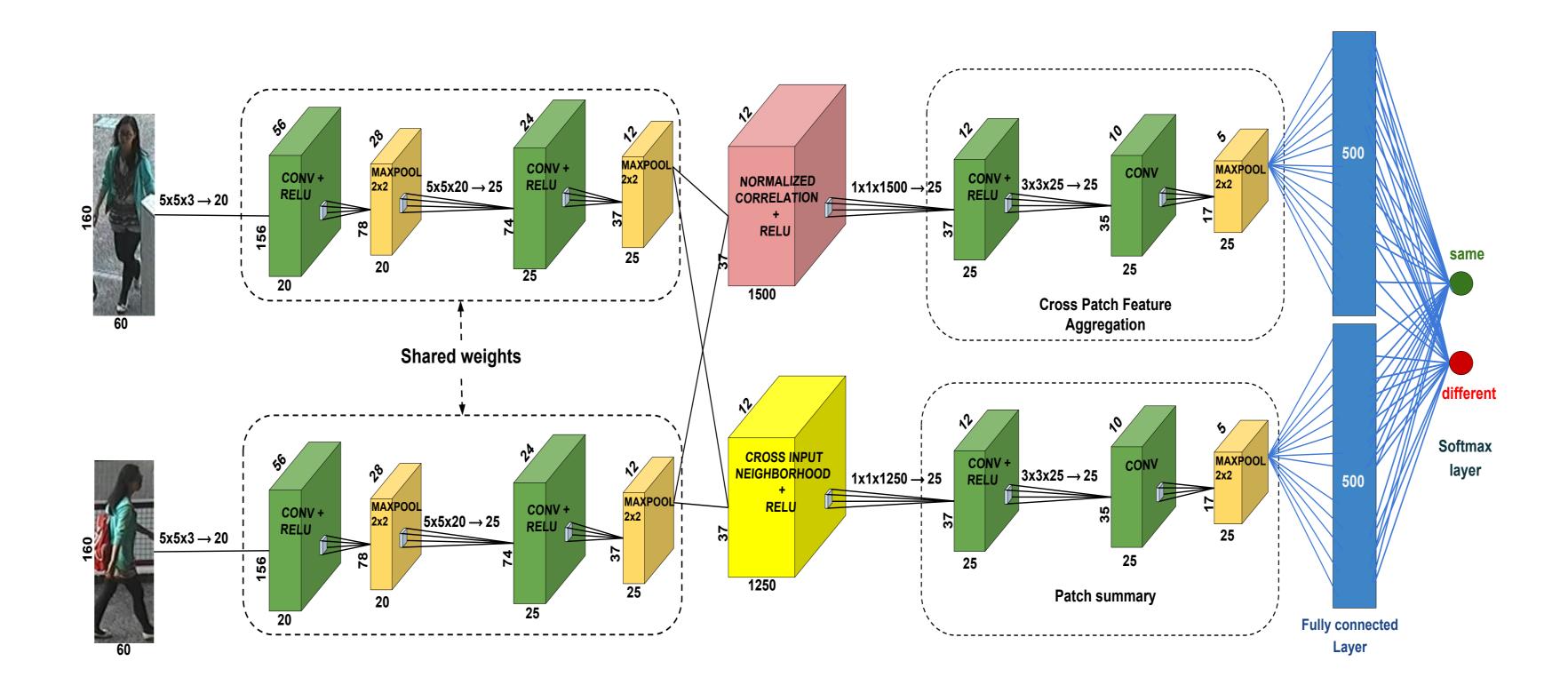
Gradient of Normalized correlation

The gradient for the Normalized Correlation layer is calculated as:

$$\frac{\partial normxcorr(E,F)}{\partial E_i} = \frac{1}{(N-1)*\sigma_E} * \left(\frac{F_i - \mu_F}{\sigma_F} - \frac{normxcorr(E,F)*(E_i - \mu_E)}{\sigma_E}\right)$$
(1)

$$\frac{\partial normxcorr(E,F)}{\partial F_i} = \frac{1}{(N-1)*\sigma_F} * \left(\frac{E_i - \mu_E}{\sigma_E} - \frac{normxcorr(E,F)*(F_i - \mu_F)}{\sigma_F}\right)$$
(2)

Fused model ($\sim 2.22M$ parameters)



Fused Model

We fuse our novel Normalized correlation layer with CrossInput Neighborhood(CIN) layer (Ahmed et al.[1]) to overcome occasional false matches due to wider search space.

CrossInput Neighborhood(CIN) layer: Given two corresponding input feature maps $X(12 \times 37)$ and $Y(12 \times 37)$ after first two convolution layers, for every pixel of X located at (x, y), where $1 \le x \le 12$ and $1 \le y \le 37$, the CIN is calculated as ,

 $CIN(x,y) = X(x,y)\mathbb{1}(5,5) - \mathcal{N}(Y(x,y))$

where $\mathcal{N}(Y(x,y))$ is the 5 \times 5 neighborhood of pixel Y(x,y). Unlike Normalized correlation, CIN is an asymmetric operation and needs to be computed in both directions $(X \to Y)$ and $(Y \to X).$

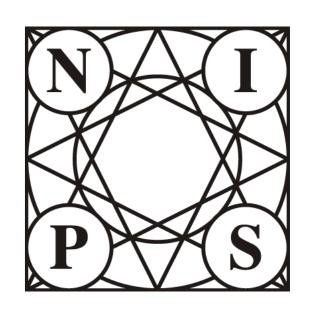
Doculto										
Results										
Table: CUHK03 (labeled & detected) datasets										
Method $r = 1$ $r = 10$ $r = 20$										
thod	r = 1	r = 10	r = 20							
del (ours)	72.04	96.00	98.26							
Corr (ours)	67.13	94.49	97.66							
APG	51.15	92.05	96.90							
	44.96	83.47	93.15							
DA	46.25	88.55	94.25							
	19.89	64.79	81.14							
	<u> </u>	<u> </u>								
Table: CUHKO1 tost100 & tost186 datasata										
$\begin{array}{llllllllllllllllllllllllllllllllllll$										
	r = 1	r = 10	m r=20							
	r = 1 65.04	r = 10 89.76	r = 20 94.49							
chod	65.04									
zhod del (ours)	65.04	89.76	94.49							
zhod del (ours)	65.04 60.17	89.76 86.26	94.49 91.47							
zhod del (ours)	65.04 60.17 59.5	89.76 86.26 89.70	94.49 91.47 93.10							
chod del (ours) Corr (ours)	65.04 60.17 59.5 51.9	89.76 86.26 89.70 83.00	94.49 91.47 93.10 89.40							
chod del (ours) Corr (ours)	65.04 60.17 59.5 51.9 47.50	89.76 86.26 89.70 83.00 80.00	94.49 91.47 93.10 89.40							
chod del (ours) Corr (ours)	65.04 60.17 59.5 51.9 47.50	89.76 86.26 89.70 83.00 80.00	94.49 91.47 93.10 89.40							
t Deep Lear:	65.04 60.17 59.5 51.9 47.50	89.76 86.26 89.70 83.00 80.00	94.49 91.47 93.10 89.40							
t Deep Lear Y	65.04 60.17 59.5 51.9 47.50	89.76 86.26 89.70 83.00 80.00	94.49 91.47 93.10 89.40							
et Deep Lear Y	65.04 60.17 59.5 51.9 47.50	89.76 86.26 89.70 83.00 80.00	94.49 91.47 93.10 89.40							
	hod del (ours) Corr (ours) APG DA	hod $r = 1$ del (ours)72.04corr (ours)67.13APG51.1544.96DA46.2519.89	hod $r = 1$ $r = 10$ del (ours)72.0496.00corr (ours)67.1394.49APG51.1592.0544.9683.47DA46.2588.5519.8964.79							

16.30 35.80 46.00 57.60

14.08 34.64 45.84 59.84

PolyMap

MtMCML



Ablation study

	NormXcorr(ours),			NormXcorr(ours),			Ahmed et al.[1],			Ahmed et al.[1],		
Dataset	5x5 search			5x12 search			5x5 search			5x12 search		
	r=1	r=10	r=50	r=1	r=10	r=50	r=1	r=10	r=50	r=1	r=10	r=50
CUHK03	62.43 9	02.22	00 60	61 73	09.77	99.60	54 74	93.30	00 70	57 60	00.63	00.17
labeled	02.40	94.44	99.00	04.70	92.11	99.00	04.14	30.00	33.10	51.00	90.00	33.11
CUHK03	63.12	2 92.76	99.20	67.13	94.49	99.73	44.96	83.47	99.40	54.31	90.24	99.18
detected												
CUHK01	72.3	79.2	2.3 95.80 99.60 77.43 96.6 '	06 67	00.20 65	65.00	94.00	99.90	60 70	05.02	00 12	
test100		2.0 90.00 99	99.00		90.07	99.29	00.00	94.00	33.3 0	09.70	<i>30.00</i>	33.10
CUHK01	56.79	5.79 84.43 95.95 60.17	86.96	06 14	17 50	20.25	06.20	10 21	01 /0	05.05		
test486			90.90	00.17	00.20	90.44	47.00	00.20	96.30	49.01	01.40	95.95

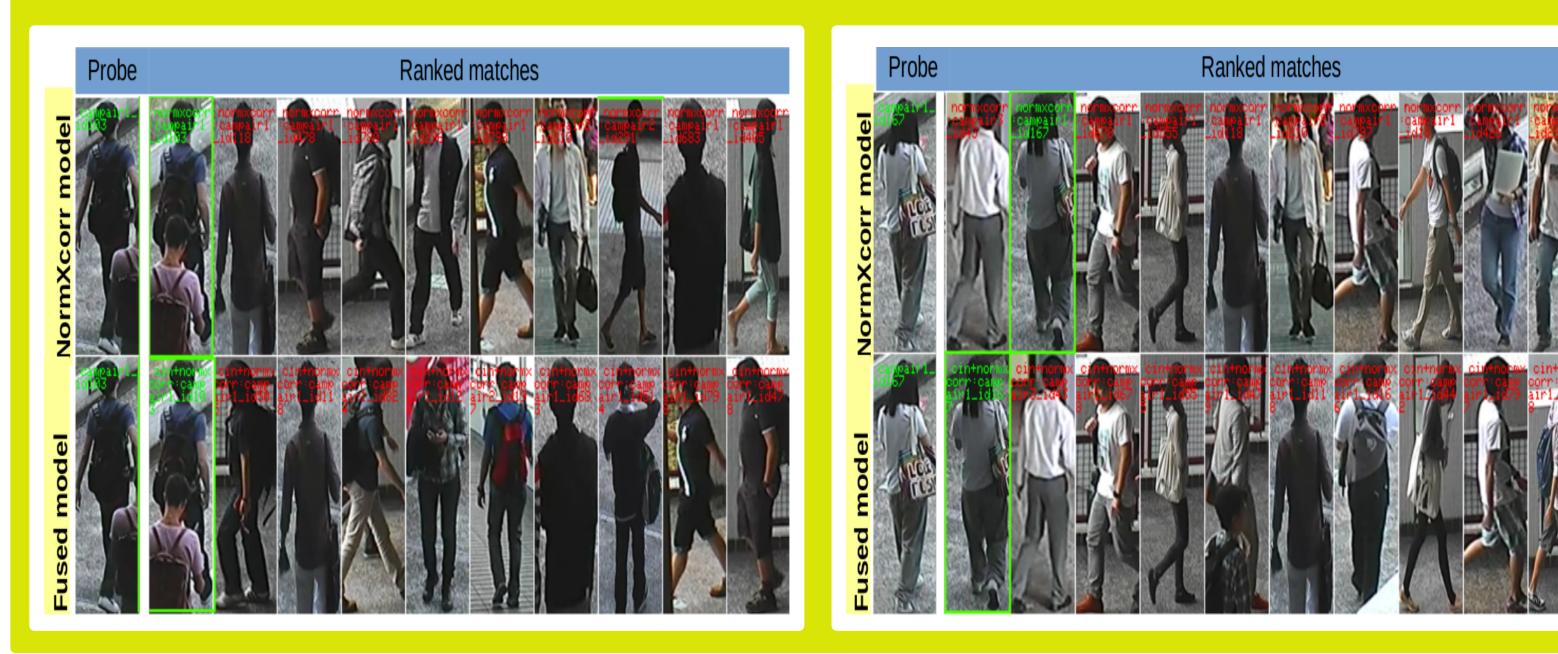
Illumination invariance

Viewpoint change



Partial occlusion

Occasional false matches



References

[1] E. Ahmed, M. Jones, and T. K. Marks, "An improved deep learning architecture for person re-identification," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3908–3916, 2015.

Acknowledgements

This work is partly supported by travel grant from Google to Arulkumar Subramaniam.