



## Person Re-Id: Recent Advances and Challenges

# Session 2: Attribute Learning and Scalability



Shiliang Zhang  
Peking University  
Beijing, China



Jingdong Wang  
Microsoft Research  
Beijing, China



Qi Tian  
University of Texas at  
San Antonio, USA



Wen Gao  
Peking University  
Beijing, China



Longhui Wei  
Peking University  
Beijing, China

Institute of Digital Media

Peking University

2018.5.15



# Outline

---

- Attribute Learning
  - Attribute embedding in multi-task learning
  - Attribute learning with deep models
- Scalable Person Re-ID
  - Deep feature extraction and compression
  - Off-line index optimization



# Attribute

- Semantic description to an image
- a stable and robust mid-level description for pedestrian
  - “wear shorts”, “male”, “long hair”, “wear grey T-shirt”
- a compact binary feature
  - ‘1’ means presence, ‘0’ means absence



cam1



cam2



cam3



cam4



cam5





# Issues in attribute learning

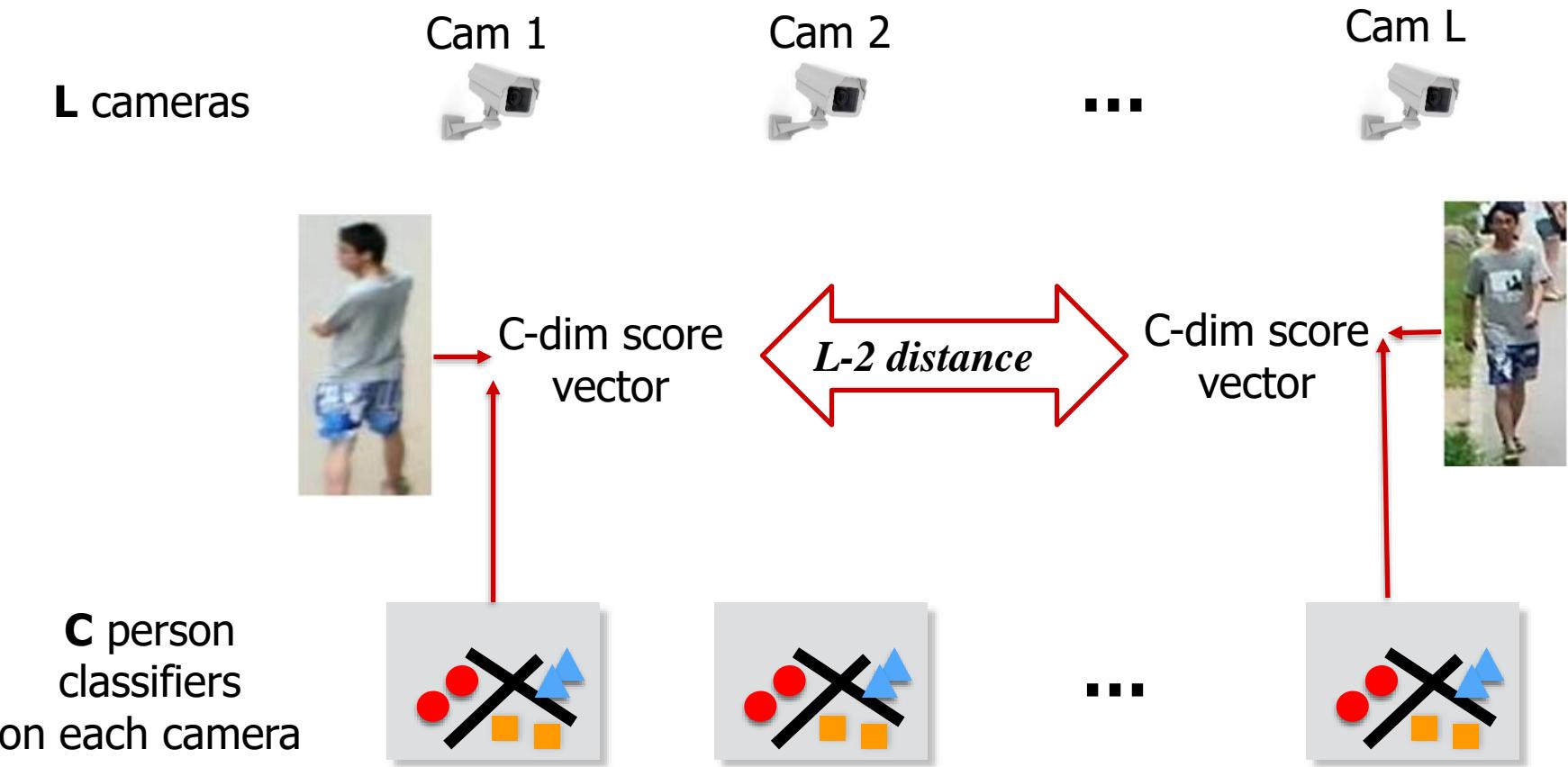
---

- Training data is hard to annotate
  - Expensive and time consuming for human annotation
- There are errors in attribute prediction
  - E.g., both ‘male’ and ‘female’ are ‘1’
- Strong correlations among attributes
  - E.g., ‘short’ and ‘bare leg’
  - ‘female’ and ‘long hair’



# Attribute embedding in multi-task learning

## □ Basic idea





# How to train C person classifiers

---

First issue:

- the classifier feature should be robust
- consider attributes and attribute correlations as features

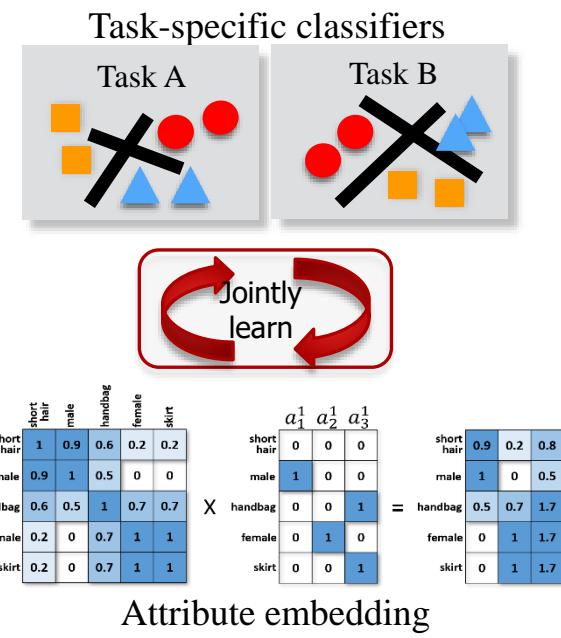
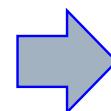
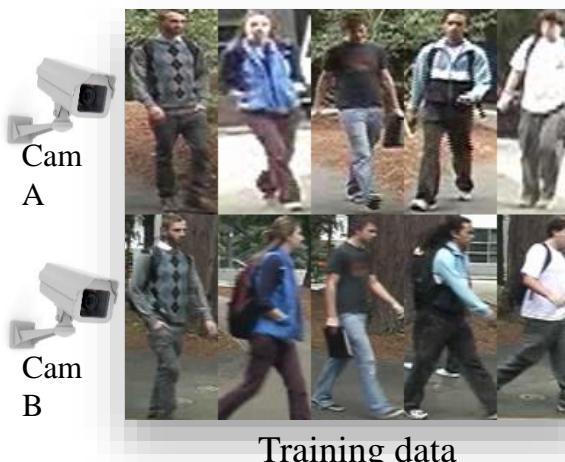
Second issue:

- different cameras corresponds to different classification models, due to different illumination, viewpoints, etc.
- the amount of training data on each camera is limited
- we treat classifier learning under different cameras as related tasks
  - Different cameras record the same set of persons



# Overview of our approach

- Multiple cameras are treated as related tasks to exploit their relationships
- Attributes correlations are embedded to improve the accuracy of predicted attributes and reduce noise
- Classifiers of multiple tasks and attribute embedding are jointly learned by alternating optimization





# attribute feature embedding

- Predict initial attributes with SVMs
- Learn attribute correlations to improve the accuracy

	short hair	male	handbag	female	skirt
short hair	1	0.9	0.6	0.2	0.2
male	0.9	1	0.5	0	0
handbag	0.6	0.5	1	0.7	0.7
female	0.2	0	0.7	1	1
skirt	0.2	0	0.7	1	1

$$\begin{array}{c}
 \text{Correlations Matrix: } \mathbf{Z} \\
 \times \quad \begin{matrix}
 a_1^1 & a_2^1 & a_3^1 \\
 \hline
 \text{short hair} & 0 & 0 & 0 \\
 \text{male} & 1 & 0 & 0 \\
 \text{handbag} & 0 & 0 & 1 \\
 \text{female} & 0 & 1 & 0 \\
 \text{skirt} & 0 & 0 & 1
 \end{matrix} = \begin{matrix}
 \text{short hair} & 0.9 & 0.2 & 0.8 \\
 \text{male} & 1 & 0 & 0.5 \\
 \text{handbag} & 0.5 & 0.7 & 1.7 \\
 \text{female} & 0 & 1 & 1.7 \\
 \text{skirt} & 0 & 1 & 1.7
 \end{matrix}
 \end{array}$$

Initial attributes:  $a_i^l$    Updated attributes:

$$\phi_{\mathbf{Z}}(\mathbf{a}_i^l) = \mathbf{Z}^\top \mathbf{a}_i^l$$

s.t.    $\text{rank}(\mathbf{Z}) \leq r,$

- Final feature:

$$\tilde{\mathbf{x}}_i^l = [\mathbf{x}_i^l; \phi_{\mathbf{Z}}(\mathbf{a}_i^l)] \in \mathbb{R}^{d+k}$$



# Formulation of classifier learning

- Loss function for person classification under camera  $l$ :

$$\ell(y_i^l, \mathbf{a}_i^l, \tilde{\mathbf{x}}_i^l, \mathbf{Z}) = \frac{1}{2} (\|y_i^l - \mathbf{w}^{l\top} \tilde{\mathbf{x}}_i^l\|^2 + \gamma \|\mathbf{a}_i^l - \mathbf{Z}^\top \mathbf{a}_i^l\|^2)$$

Person classification error

$$\tilde{\mathbf{x}}_i^l = [\mathbf{x}_i^l; \phi_{\mathbf{Z}}(\mathbf{a}_i^l)] \in \mathbb{R}^{d+k}$$

:concatenation of low level features  
and the updated attributes

attribute embedding  
regularization

- Goal: learn classifier  $\mathbf{w}$  and embedding matrix  $\mathbf{Z}$



# Formulation of multi-task learning

## □ Low-rank embedding

- the classifiers to be learned on  $L$  cameras:

$$\mathbf{W} = [\mathbf{w}^1, \mathbf{w}^2, \dots, \mathbf{w}^L] \in \mathbb{R}^{(d+k) \times L}$$

- We define:  $\mathbf{W} = \mathbf{R} + \mathbf{S}$
- $\mathbf{R}$ : a low-rank matrix shared by all tasks
- $\mathbf{S}$ : task-specific sparse component

$$\begin{aligned} \min_{\mathbf{R}, \mathbf{S}, \mathbf{Z}} \quad & \sum_{l=1}^L \sum_{i=1}^{n_l} \ell(y_i^l, \mathbf{a}_i^l, \tilde{\mathbf{x}}_i^l, \mathbf{Z}) + \lambda \|\mathbf{S}\|_0 \\ \text{s.t.} \quad & \mathbf{W} = \mathbf{R} + \mathbf{S}, \text{rank}(\mathbf{R}) \leq r_1, \text{rank}(\mathbf{Z}) \leq r_2 \end{aligned}$$



# Formulation of multi-task learning

- Relax the previous formulation into a smooth one

$$\min_{\mathbf{R}, \mathbf{S}, \mathbf{Z}} \quad \sum_{l=1}^L \sum_{i=1}^{n_l} \ell(y_i^l, \mathbf{a}_i^l, \tilde{\mathbf{x}}_i^l, \mathbf{Z}) + \lambda \|\mathbf{S}\|_1$$

$$\text{s.t.} \quad \mathbf{W} = \mathbf{R} + \mathbf{S}, \|\mathbf{R}\|_* \leq r_1, \|\mathbf{Z}\|_* \leq r_2,$$

- This formulation can be optimized by an alternating optimization
  - Fix  $\mathbf{Z}$  to update the classifier parameter  $\mathbf{w}$  by **MixedNorm**
  - Fix  $\mathbf{w}$  update  $\mathbf{Z}$  by **proximal gradient method**

For details:

Multi-Task Learning with Low Rank Attribute Embedding for Person Re-identification, *ICCV*, 2015  
Multi-Task learning with low rank attribute embedding for multi-camera person re-identification.  
*T-PAMI*, 2018



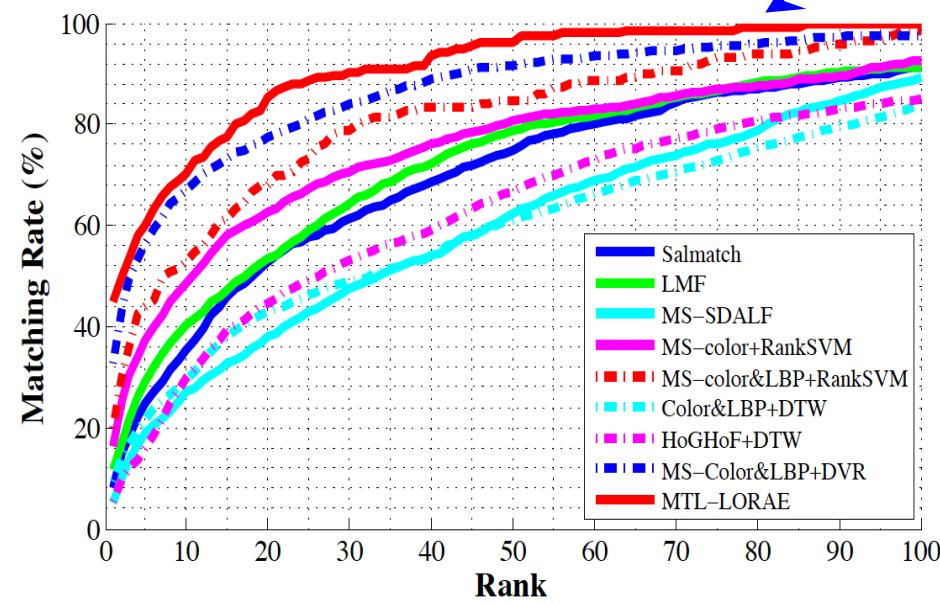
# Datasets

- **iLIDS-VID**
  - two camera, multi-shot
  - 600 image sets of 300 persons
- **PRID**
  - two camera, multi-shot
  - 385 and 749 persons in cameras A and B, respectively. 200 persons appear in both cameras
- **VIPeR**
  - two camera, single-shot
  - 632 persons in two cameras
- Use a 2784-dimensional color and texture descriptor as low level feature
- learn binary SVMs to predict 20-bit attributes for PRID, 90-bit attributes for VIPeR. Learn 32-bit attributes for iLIDS-VID.

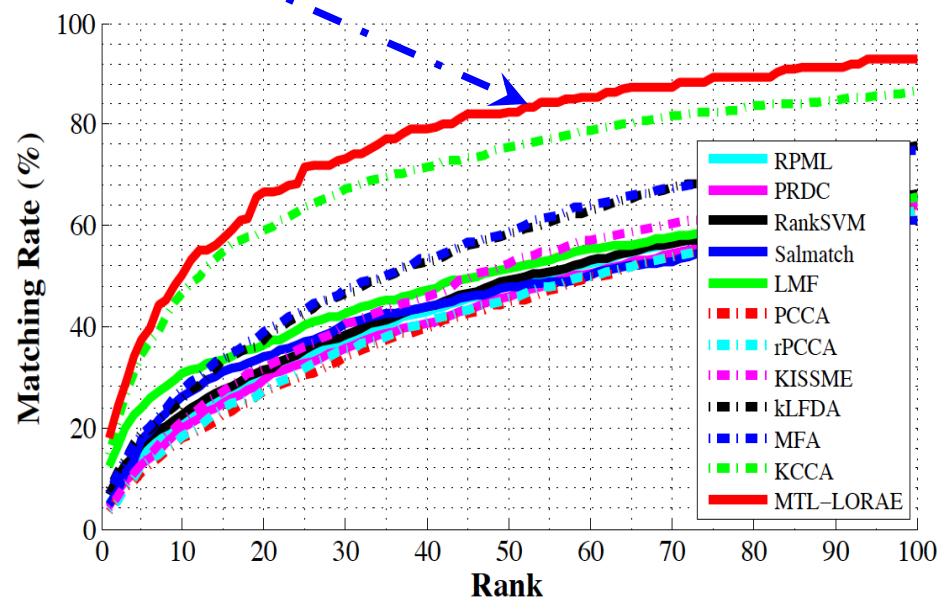


# Performance of two camera multi-shot Re-ID task

Our approach gets very competitive performance



iLIDS-VID



PRID





# Performance of two camera single-shot Re-ID task

## Viper dataset

TABLE 4  
CMC Scores of Ranks from 1 to 20 on the VIPeR Dataset

Rank	1	5	10	20
KISSME [11]	19.6	47.5	62.2	77.0
kLFDA [12]	32.2	65.8	79.7	90.9
KCCA [77]	37.3	71.4	<b>84.6</b>	92.3
LOMO + XQDA [78]	40.0	68.9	81.5	91.1
TSR [79]	31.6	68.6	82.8	<b>94.6</b>
EPKFM [80]	36.8	70.4	83.7	91.7
MLAPG [81]	40.7	69.9	82.3	92.4
MTL-LORAE	<b>42.3</b>	<b>72.2</b>	81.6	89.6

Our approach gets best performance for Rank -1 and 5



# Performance of combining deep feature

- VGG-FC7: use FC7 layer output of VGG-16
- Pre-trained on ImageNet, fine-tuned on cameras #1, #2, #4, #6 and #7 on SAIVT-SoftBio

	Methods	Rank 1	Rank 5	Rank 10	Rank 20
<i>iLIDS-VID</i>	VGG-FC7	24.1	43.6	52.8	65.6
	MTL-LOREA	43.0	60.0	70.2	85.3
	MTL-LOREA-FC7	56.4	69.0	78.4	87.4
<i>PRID</i>	VGG-FC7	19.8	28.5	42.4	53.9
	MTL-LOREA	18.0	37.4	50.1	66.6
	MTL-LOREA-FC7	21.0	44.0	55.9	68.7
<i>VIPeR</i>	VGG-FC7	25.1	39.8	48.5	60.6
	MTL-LOREA	42.3	72.2	81.6	89.6
	MTL-LOREA-FC7	45.4	76.6	85.3	91.7

- Our approach outperforms deep features
- The performance further boosted by combining deep features



# Examples of attribute correlations

blueshirt	<b>0.60±0.05</b>	-0.07±0.06	-0.06±0.06	-0.01±0.02	0.03±0.03
lightshirt	-0.42±0.06	<b>0.77±0.06</b>	-0.14±0.07	0.02±0.02	0.05±0.02
darkshirt	-0.40±0.06	-0.37±0.07	<b>0.80±0.10</b>	0.04±0.03	0.03±0.02
darkbottoms	-0.01±0.02	0.02±0.02	0.04±0.03	<b>0.76±0.02</b>	-0.19±0.03
lightbottoms	0.04±0.03	0.20±0.02	0.03±0.02	-0.43±0.03	<b>0.69±0.08</b>
	blueshirt	lightshirt	darkshirt	darkbottoms	lightbottoms

barelegs	<b>0.52±0.05</b>	0.14±0.02	-0.02±0.01	-0.02±0.02	-0.02±0.02
shorts	<b>0.62±0.02</b>	0.19±0.04	0.01±0.02	0.03±0.01	0.04±0.02
male	-0.03±0.01	0.01±0.02	<b>0.91±0.02</b>	-0.01±0.01	0.03±0.01
darkhair	-0.02±0.02	0.05±0.01	-0.01±0.01	<b>0.83±0.03</b>	-0.12±0.02
bald	-0.03±0.02	0.08±0.02	0.03±0.01	-0.29±0.02	<b>0.72±0.07</b>
	barelegs	shorts	male	darkhair	bald

hairBlack	<b>0.72±0.05</b>	-0.31±0.05	-0.52±0.02	0.03±0.02	0.04±0.03
hairBrown	-0.31±0.05	<b>0.78±0.04</b>	-0.51±0.01	0.06±0.03	0.01±0.03
hairWhite	-0.52±0.02	-0.51±0.01	<b>0.73±0.09</b>	0.00±0.02	0.01±0.02
hairLong	0.03±0.02	0.06±0.03	0.00±0.02	<b>0.80±0.03</b>	-0.33±0.04
hairShort	0.04±0.03	0.01±0.03	0.01±0.02	-0.33±0.04	<b>0.83±0.07</b>
	hairBlack	hairBrown	hairWhite	hairLong	hairShort

personalLess30	<b>0.88±0.03</b>	-0.31±0.03	0.01±0.03	0.00±0.03	-0.01±0.03
personalLess45	-0.31±0.03	<b>0.88±0.03</b>	0.02±0.03	0.03±0.03	0.01±0.04
carryingMessengerBag	0.01±0.03	0.02±0.03	<b>0.84±0.05</b>	-0.37±0.03	-0.31±0.04
carryingNothing	0.00±0.03	0.03±0.03	-0.37±0.03	<b>0.83±0.05</b>	-0.34±0.02
carryingOther	-0.01±0.03	0.01±0.04	-0.31±0.04	-0.34±0.03	<b>0.84±0.08</b>
	personalLess30	personalLess45	carryingMessengerBag	carryingNothing	carryingOther



北京大学



# Validity of MTL and attribute embedding

- STL: single task learning
- MTL-Att: without attribute embedding

Rank	<i>iLIDS-VID</i>					
	1	5	10	20	30	50
STL	14.7	42.7	41.8	58.5	83.5	91.7
MTL-FR	37.7	54.0	47.4	64.9	85.3	92.5
MTL-Att	40.5	54.9	47.5	64.2	84.2	91.2
MTL-LOREA	43.0	60.0	70.2	85.3	90.2	96.3

Rank	<i>PRID</i>					
	1	5	10	20	30	50
STL	11.3	27.9	41.8	53.0	68.5	74.6
MTL-FR	11.3	34.1	47.4	61.1	69.8	79.0
MTL-Att	12.2	34.7	47.5	61.7	70.9	79.8
MTL-LOREA	18.0	37.4	50.1	66.6	73.1	82.3

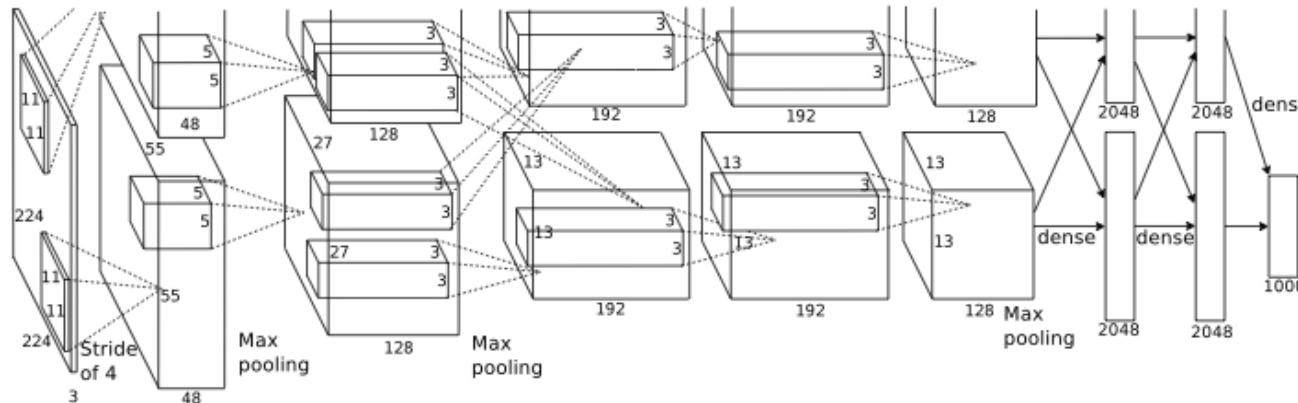
Rank	<i>VIPeR</i>					
	1	5	10	20	30	50
STL	13.3	27.4	32.8	42.7	56.2	68.3
MTL-FR	35.3	63.3	75.6	83.8	89.9	94.4
MTL-Att	37.2	64.2	76.3	84.9	91.4	95.3
MTL-LOREA	42.3	72.2	81.6	89.6	93.1	97.4

- Multi-task learning boosts the performance
- Attribute embedding boost the performance



# Attribute learning with deep models

- Deep learning (deep Convolutional Neural Network, dCNN):
  - Powerful learning ability and good performance in pattern recognition tasks
  - Good generalization ability



The structure of AlexNet (NIPS'12)



# Issues & Solution

## □ Issues

- dCNN needs a large amount of training data
- The available labeled attribute data is limited

## □ Solution: Semi-supervised Deep Attribute Learning (SSDAL)

- Train dCNN with a partially-labeled dataset
  - Attribute annotation: PETA
  - Another dataset with person ID: MOT Challenge
- Ensure its discriminative power and generalization ability in the person ReID tasks

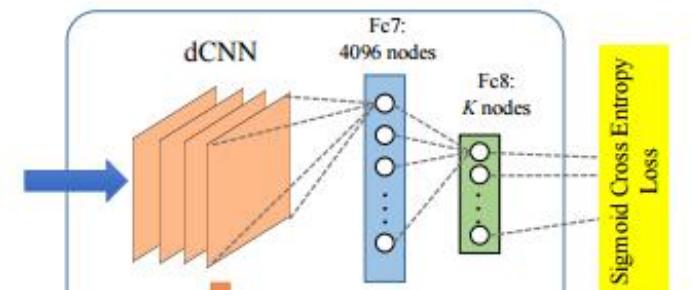


# Framework

Stage-1

Stage 1: Fully-supervised dCNN training

Independent dataset  
with attribute labels



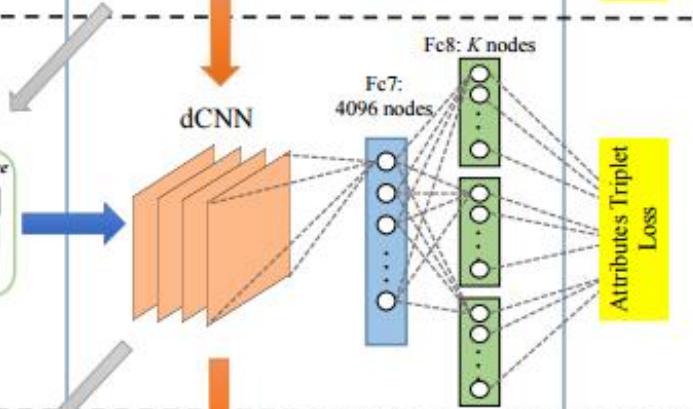
Stage-2

Stage 2: Fine-tuning using attributes triplet loss

Dataset with person ID labels

Anchor	Positive	Negative
⋮	⋮	⋮

Anchor	Positive	Negative
[110...]	[101...]	[11...]
[011...]	[101...]	[001...]
⋮	⋮	⋮



Stage-3

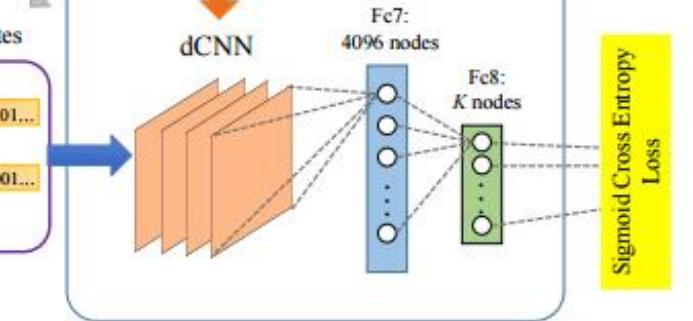
Stage 3: Final fine-tuning on the combined dataset

Independent dataset



Dataset with refined attributes

	[110...]		[110...]		[101...]
	[101...]		[101...]		[001...]
⋮	⋮	⋮	⋮	⋮	⋮



北京大学



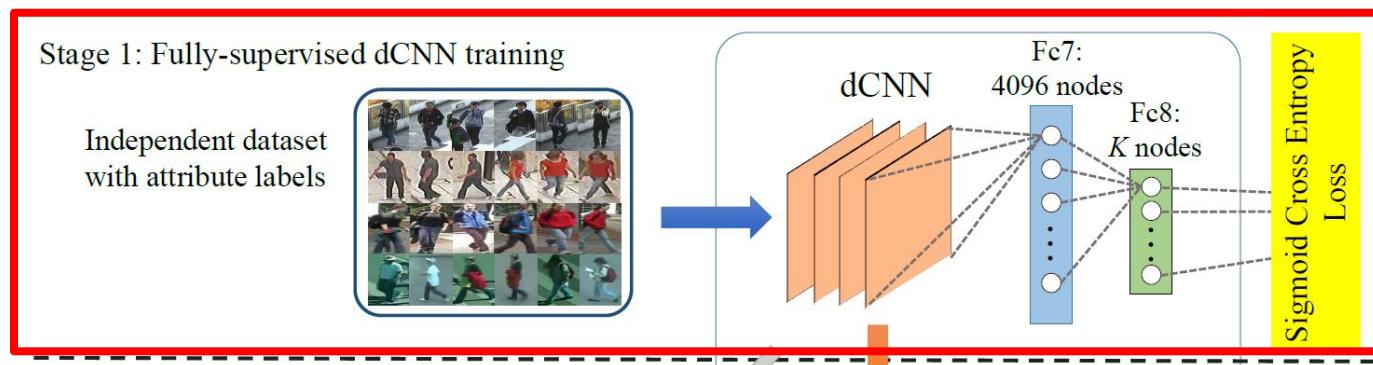
# Our Approach

---

- Solution:
  - Three Stage:
    - Stage 1: a fully supervised training on an independent dataset labeled with human attributes
    - Stage 2: fine-tune on another target dataset labeled with person IDs using *Attributes Triplet Loss*
    - Stage 3: The updated dCNN predicts attribute labels for the target dataset, which is finally combined with the independent dataset for the second round of fine-tuning.



# Stage-1



## Fully-Supervised dCNN Training

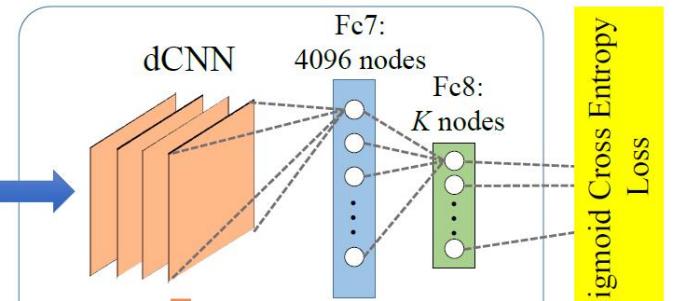
- Train with a small **independent dataset labeled with:**  
**image + attributes**
- Use a sigmoid cross-entropy loss layer for multi-label  
prediction

# Stage-2

Fine-Tune with  
Attributes  
Triplet Loss

Stage 1: Fully-supervised dCNN training

Independent dataset  
with attribute labels

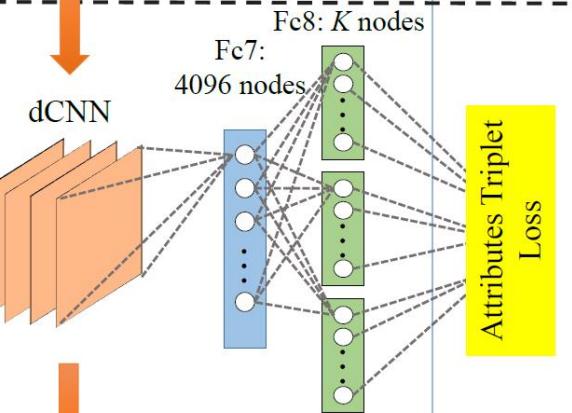


Stage 2: Fine-tuning using attributes triplet loss

Dataset with person ID labels      Predicted attributes

Anchor	Positive	Negative
$\vdots$	$\vdots$	$\vdots$

Anchor	Positive	Negative
110...	101...	111...
011...	101...	001...
$\vdots$	$\vdots$	$\vdots$



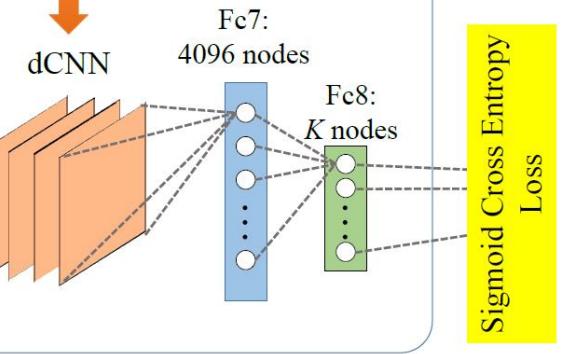
Stage 3: Final fine-tuning on the combined dataset

Independent dataset



Dataset with refined attributes

	110...		110...		101...
	101...		101...		001...
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$





# Fine-Tune with Attributes Triplet Loss

- Goal: refine the predicted attributes to make same person share more similar attributes
- Network predicts the attributes  $[A_{(a)}, A_{(p)}, A_{(n)}]$  of three images labeled with person ID
- Network loss:

$$\mathcal{L} = \sum_e^E \left\{ \max \left( 0, \boxed{\mathbf{D} \left( A_{(a)}^{(e)}, A_{(p)}^{(e)} \right)} + \theta - \right. \right.$$

distance between the **same** person

$$\left. \left. \boxed{\mathbf{D} \left( A_{(a)}^{(e)}, A_{(n)}^{(e)} \right)} \right) + \gamma \times \mathcal{E} \right\}$$

distance between **different** person

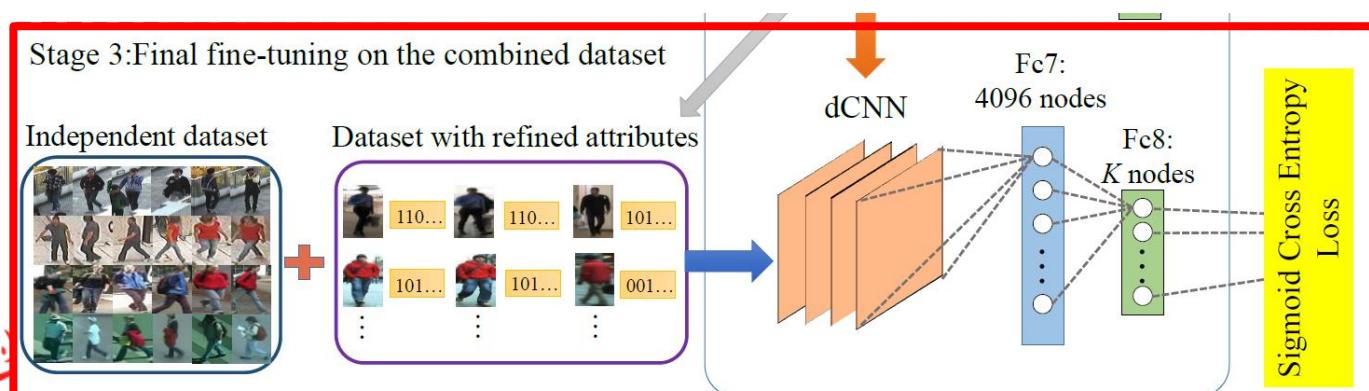
$$\mathcal{E} = \mathbf{D} \left( A_{(a)}^{(e)}, \tilde{A}_{(a)}^{(e)} \right) + \mathbf{D} \left( A_{(p)}^{(e)}, \tilde{A}_{(p)}^{(e)} \right) + \mathbf{D} \left( A_{(n)}^{(e)}, \tilde{A}_{(n)}^{(e)} \right)$$



# Stage-3

## Fine-Tuning on the Combined Dataset

- The labeled dataset from Stage-2 + the original independent dataset
- The attributes predicted by the final dCNN model are named as *deep attributes*
- *deep attributes* + Cosine Distance for Person Re-ID



北京大



# Examples of Predicted Attributes



lowerBodyBlack  
lowerBodyTrousers  
hairShort  
personalFemale  
carryingMessengerBag  
  
accessoryNothing  
hairBlack



lowerBodyTrousers hairShort  
hairGrey  
upperBodyWhite  
  
personalMale  
carryingNothing  
footwearLeatherShoes



upperBodyGrey  
accessoryHat  
lowerBodyTrousers  
hairShort  
  
upperBodyOther



upperBodyRed  
lowerBodyJeans  
personalLess30  
lowerBodyBlue  
accessoryNothing  
  
footwearSneakers  
hairShort



upperBodyWhite  
hairShort  
accessoryNothing  
lowerBodyWhite  
carryingNothing  
lowerBodyGrey  
  
lowerBodyWhite



lowerBodyTrousers  
upperBodyBlue  
lowerBodyGrey  
carryingMessengerBag  
  
accessoryNothing  
hairBlack  
personalMale



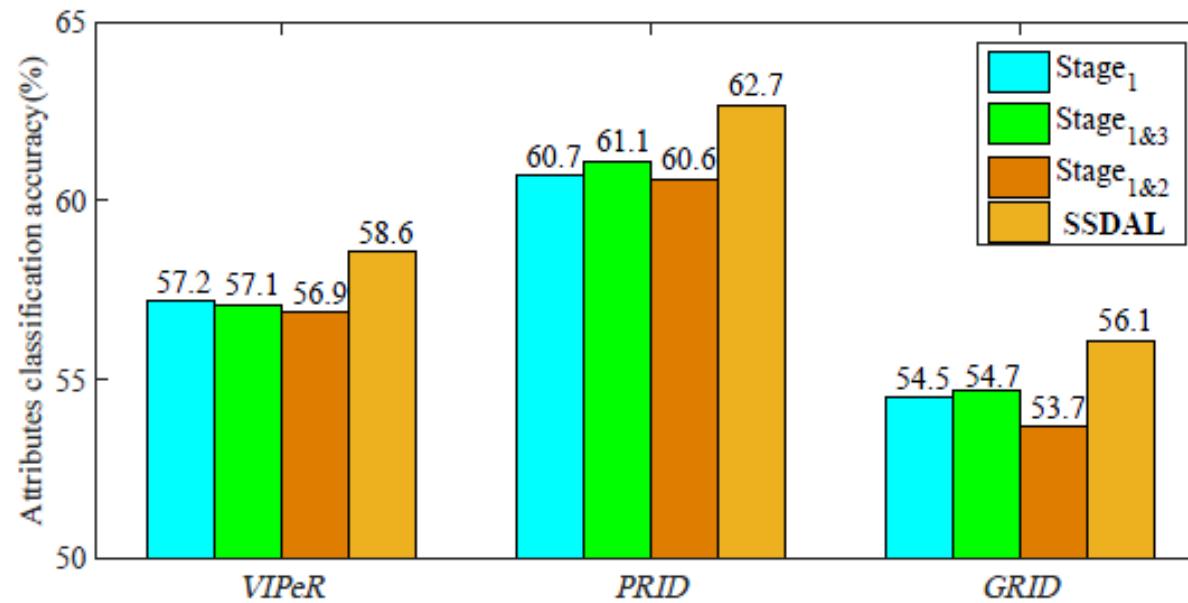
accessoryNothing  
hairLong  
carryingMessengerBag  
footwearBlack  
lowerBodyTrousers  
  
upperBodyWhite



lowerBodyShorts  
upperBodyTshirt  
hairShort  
personalLess30  
accessoryNothing  
footwearWhite  
  
upperBodyShortSleeve  
footwearSandals



# Why using three stages



The three stage training gets the best performance



# Datasets for experiments

---

- Training dataset:
  - PETA: annotated with 105 attributes
  - MOT challenge: annotated with person ID
  
- Independent testing dataset:
  - VIPeR: 632 persons, each has two 48\*128 images
  - PRID: 385 and 749 persons captured by camera A and B
  - GRID: 1025 persons taken by 8 non-adjacent cameras
  - Market: 25,000 images of 1501 persons taken by 6 cameras

The training data and testing data are  
totally independent



# Performance-1

## □ Two cams datasets:

### ■ VIPeR:

Methods		Rank 1	Rank 5	Rank 10	Rank 20
Metric Learning based ReID	RPML [10]	27.0	57.0	69.0	83.0
	Salmatch [48]	30.2	52.4	65.5	79.1
	LMF [49]	29.1	52.3	65.9	80.0
	KISSME [13]	19.6	47.5	62.2	77.0
	KCCA [50]	37.3	71.4	<b>84.6</b>	92.3
	kLFDA [14]	32.2	65.8	79.7	90.9
	LOMO + XQDA [20]	40.0	68.9	81.5	91.1
	CSL [22]	34.8	68.7	82.3	91.8
	MLAPG [23]	40.7	69.9	82.3	92.4
	TSR [51]	31.6	68.6	82.8	<b>94.6</b>
	EPKFM [19]	36.8	70.4	83.7	91.7
Traditional Attributes Learning based ReID	AIR [29]	18.0	38.8	51.1	71.2
	OAR [31]	21.4	41.5	55.2	71.5
	LORAE [34]	42.3	<b>72.2</b>	81.6	89.6
Deep Learning based ReID	IDLA [42]	34.8	54.3	76.5	87.6
	DML [41]	28.2	59.3	73.5	86.4
	Deep-RDC [24]	40.5	60.8	70.4	84.4
Proposed	<i>Stage<sub>1</sub></i>	34.5	63.9	73.1	87.0
	<b>SSDAL</b>	37.9	65.5	75.6	88.4
	<b>SSDAL + XQDA</b>	<b>43.5</b>	71.8	81.5	89.0



# Performance-2

## □ Two cams datasets:

### ■ PRID:

Methods	Rank 1	Rank 5	Rank 10	Rank 20
RPML [10]	4.8	14.3	21.6	30.2
PRDC [17]	4.5	12.6	19.7	29.5
RSVM [52]	6.8	16.5	22.7	31.5
Salmatch [48]	4.9	17.5	26.1	33.9
LMF [49]	12.5	23.9	30.7	36.5
PCCA [9]	3.5	10.9	17.9	27.1
KISSME [13]	4.1	12.8	21.1	31.8
kLFDA [14]	7.6	18.9	25.6	37.4
KCCA [50]	14.5	34.3	46.7	59.1
LOREA [34]	18.0	37.4	50.1	66.6
LOMO + XQDA [20]	15.3	35.7	41.2	53.8
MLAPG [23]	16.6	33.1	41.4	52.5
<i>Stage<sub>1</sub></i>	18.7	46.9	55.0	65.8
<b>SSDAL</b>	20.1	47.4	55.7	68.6
<b>SSDAL + XQDA</b>	<b>22.6</b>	<b>48.7</b>	<b>57.8</b>	<b>69.2</b>



# Performance-3

- Eight cams datasets:
  - GRID:

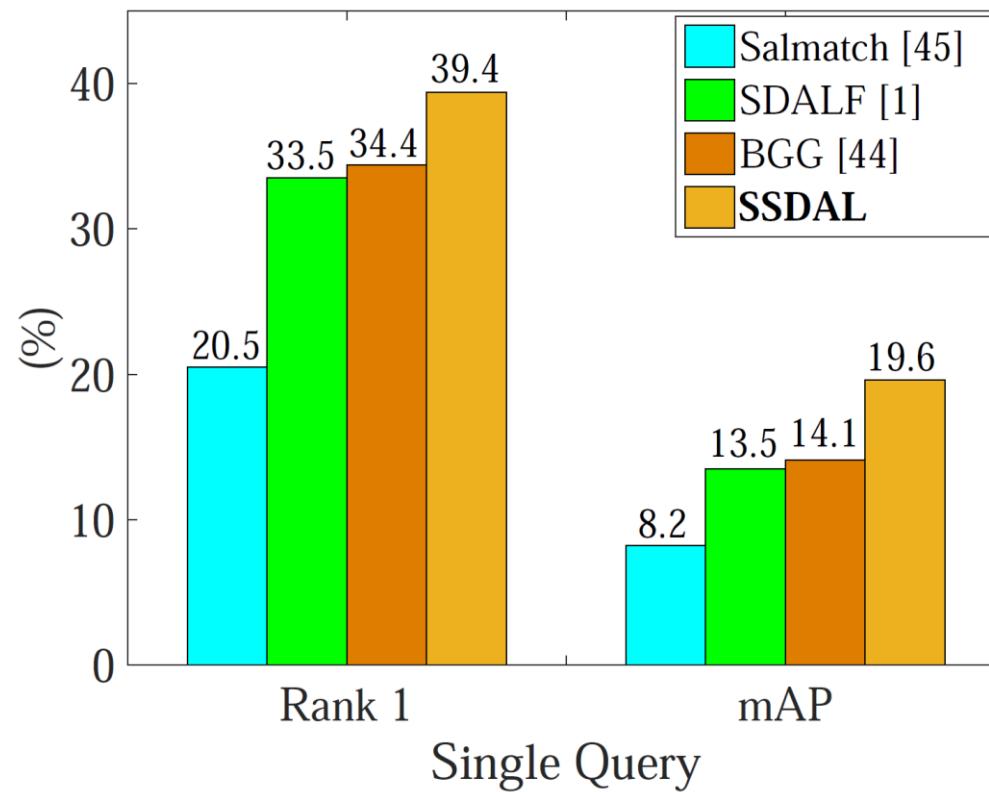
Methods	Rank 1	Rank 5	Rank 10	Rank 20
PRDC [17]	9.7	22.0	33.0	44.3
RSVM [52]	10.2	24.6	33.3	43.7
MRank-PRDC [28]	11.1	26.1	35.8	46.6
MRank-RSVM [28]	12.2	27.8	36.3	49.3
RQDA [53]	15.2	30.1	39.2	49.3
EPKFM [19]	16.3	35.8	46.0	57.6
LOMO + XQDA [20]	16.6	35.4	41.8	52.4
<i>Stage1</i>	16.9	30.1	40.7	50.2
<b>SSDAL</b>	19.1	35.6	45.8	58.1
<b>SSDAL + XQDA</b>	<b>22.4</b>	<b>39.2</b>	<b>48.0</b>	<b>58.4</b>



# Performance-4

- Large-scale multi-cam datasets:

- Market:





# Issues of SSDAL

---

- Shallow network, 5 conv layers
- Has thresholds to select positive attributes
- Does not consider the correlations between attributes



# Solution and Extension

---

- Propose a Semi-supervised Multi-Type Attribute Learning (SMTAL)
- Uses a deeper network: VGG-16, 13 conv layers
- Takes the correlations between attributes into consideration
  - Divide the attributes into sub-groups
  - in each sub-group only one attribute could be positive, e.g., a group contains: male and female

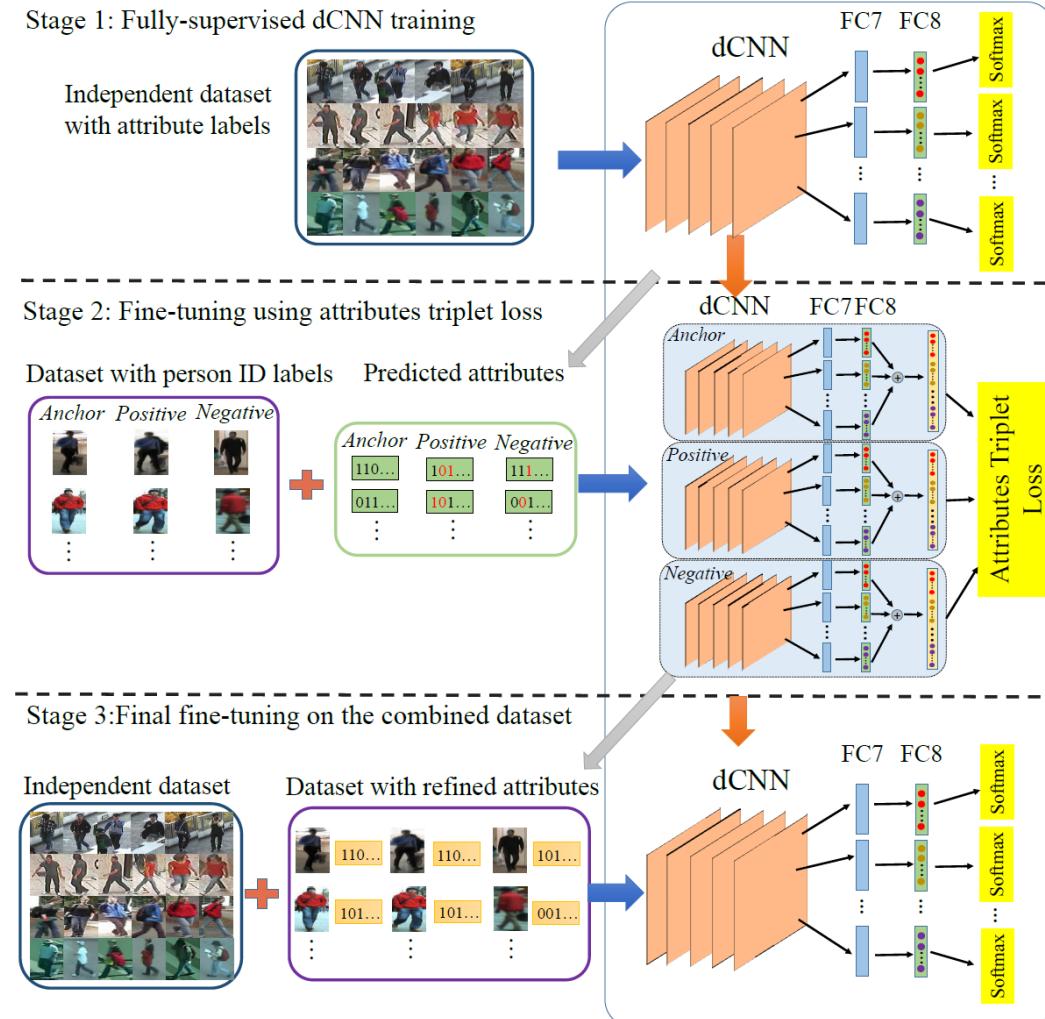


# Attribute groups

Type	Attribute Label
<b>Age</b>	personalLess15, personalLess30, personalLess45, personalLess60, personalLarger60
<b>Gender</b>	personalFemale, personalMale
<b>CarryObject</b>	carryingBabyBuggy, carryingBackpack, carryingOther, carryingShoppingTro, carryingUmbrella, carryingFolder, carryingLuggageCase, carryingMessengerBag, carryingNothing, carryingPlasticBags, carryingSuitcase
<b>AccessoryObject</b>	accessoryHeadphone, accessoryHairBand, accessoryHat, accessoryKerchief, accessoryMuffler, accessoryNothing, accessorySunglasses
<b>SleeveStyle</b>	upperBodyNoSleeve, upperBodyShortSleeve, upperBodyLongSleeve
<b>UpperStyle</b>	upperBodyCasual, upperBodyFormal
<b>UpperType</b>	upperBodyJacket, upperBodyLogo, upperBodyPlaid, upperBodyThinStripes, upperBodySuit, upperBodySweater, upperBodyThickStripes, upperBodyTshirt, upperBodyOther, upperBodyVNeck
<b>LowerStyle</b>	lowerBodyCasual, lowerBodyFormal
<b>LowerType</b>	lowerBodyCapri, lowerBodyHotPants, lowerBodyJeans, lowerBodyLongSkirt, lowerBodyPlaid, lowerBodyThinStripes, lowerBodyShorts, lowerBodyShortSkirt, lowerBodySuits, lowerBodyTrousers
<b>HairStyle</b>	hairBald, hairShort, hairLong
<b>FootStyle</b>	footwearBoots, footwearLeatherShoes, footwearSandals, footwearShoes, footwearSneakers, footwearStocking
<b>UpperColor</b>	upperBodyBlack, upperBodyBlue, upperBodyBrown, upperBodyGreen, upperBodyGrey, upperBodyOrange, upperBodyPink, upperBodyPurple, upperBodyRed, upperBodyWhite, upperBodyYellow
<b>LowerColor</b>	lowerBodyBlack, lowerBodyBlue, lowerBodyBrown, lowerBodyGreen, lowerBodyGrey, lowerBodyOrange, lowerBodyPink, lowerBodyPurple, lowerBodyRed, lowerBodyWhite, lowerBodyYellow
<b>HairColor</b>	hairBlack, hairBlue, hairBrown, hairGreen, hairGrey, hairOrange, hairPink, hairPurple, hairRed, hairWhite, hairYellow
<b>FootColor</b>	footwearBlack, footwearBlue, footwearBrown, footwearGreen, footwearGrey, footwearOrange, footwearPink, footwearPurple, footwearRed, footwearWhite, footwearYellow



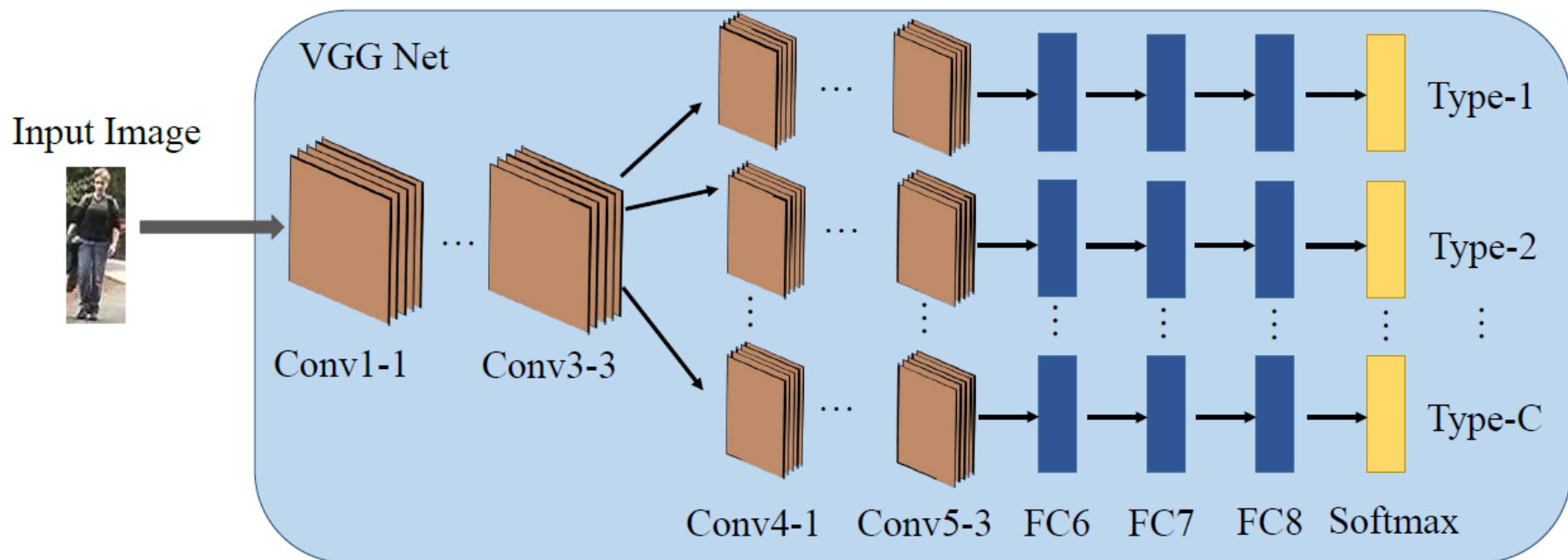
# framework





# Network structure

- Each group shares several conv layers, has independent classification layer





# Attribute prediction accuracy

Types	Number	Validation(%)		ViPeR(%)		PRID(%)		GRID(%)		
		Top-1	Top-2	Top-1	Top-2	Top-1	Top-2	Top-1	Top-2	
Age	5	<i>Stage<sub>1</sub></i>	86.41	97.03	44.34	80.42	48.13	77.35	42.56	78.35
		<i>Stage<sub>1&amp;2</sub></i>	84.56	95.18	41.09	77.85	46.03	75.45	39.87	76.42
		<i>Stage<sub>1&amp;3</sub></i>	<b>86.47</b>	<b>97.51</b>	45.02	80.33	48.23	78.24	42.45	77.99
		SMTAL	86.02	96.87	<b>46.99</b>	<b>81.84</b>	<b>50.43</b>	<b>79.13</b>	<b>45.13</b>	<b>81.84</b>
Sex	2	<i>Stage<sub>1</sub></i>	<b>96.53</b>	<b>100.00</b>	70.05	<b>100.00</b>	67.50	<b>100.00</b>	67.82	<b>100.00</b>
		<i>Stage<sub>1&amp;2</sub></i>	95.47	<b>100.00</b>	69.12	<b>100.00</b>	66.43	<b>100.00</b>	66.34	<b>100.00</b>
		<i>Stage<sub>1&amp;3</sub></i>	96.49	<b>100.00</b>	70.13	<b>100.00</b>	67.66	<b>100.00</b>	68.03	<b>100.00</b>
		SMTAL	95.83	<b>100.00</b>	<b>71.20</b>	<b>100.00</b>	<b>69.64</b>	<b>100.00</b>	<b>69.73</b>	<b>100.00</b>
CarryObject	11	<i>Stage<sub>1</sub></i>	<b>86.35</b>	<b>96.31</b>	26.46	46.82	48.13	77.35	26.96	41.28
		<i>Stage<sub>1&amp;2</sub></i>	84.75	95.92	25.12	47.03	46.03	75.45	25.47	41.96
		<i>Stage<sub>1&amp;3</sub></i>	86.32	95.69	27.19	46.71	48.23	78.24	26.88	41.63
		SMTAL	85.08	94.24	<b>29.27</b>	<b>49.01</b>	<b>31.71</b>	<b>54.82</b>	<b>28.74</b>	<b>44.53</b>
AccessoryObject	7	<i>Stage<sub>1</sub></i>	<b>91.86</b>	<b>97.54</b>	44.32	66.33	57.66	75.58	52.45	83.77
		<i>Stage<sub>1&amp;2</sub></i>	90.25	96.44	43.21	65.74	56.41	74.93	51.68	83.40
		<i>Stage<sub>1&amp;3</sub></i>	91.55	97.19	44.83	66.85	57.38	75.69	51.85	83.79
		SMTAL	90.62	97.40	<b>46.17</b>	<b>69.10</b>	<b>60.14</b>	<b>78.85</b>	<b>54.19</b>	<b>86.06</b>
SleeveStyle	3	<i>Stage<sub>1</sub></i>	99.14	<b>100.00</b>	80.17	93.73	85.37	95.62	42.71	77.96
		<i>Stage<sub>1&amp;2</sub></i>	98.75	<b>100.00</b>	79.64	93.46	84.33	94.23	42.22	77.68
		<i>Stage<sub>1&amp;3</sub></i>	<b>99.93</b>	<b>100.00</b>	79.38	93.63	84.96	95.73	43.09	76.93
		SMTAL	98.96	<b>100.00</b>	<b>85.76</b>	<b>98.73</b>	<b>87.09</b>	<b>97.88</b>	<b>45.13</b>	<b>81.84</b>
UpperStyle	2	<i>Stage<sub>1</sub></i>	<b>99.01</b>	<b>100.00</b>	90.17	<b>100.00</b>	82.64	<b>100.00</b>	80.53	<b>100.00</b>
		<i>Stage<sub>1&amp;2</sub></i>	97.85	<b>100.00</b>	89.25	<b>100.00</b>	82.01	<b>100.00</b>	80.11	<b>100.00</b>
		<i>Stage<sub>1&amp;3</sub></i>	98.69	<b>100.00</b>	90.33	<b>100.00</b>	82.49	<b>100.00</b>	81.36	<b>100.00</b>
		SMTAL	98.96	<b>100.00</b>	<b>92.96</b>	<b>100.00</b>	<b>84.10</b>	<b>100.00</b>	<b>83.63</b>	<b>100.00</b>
UpperType	10	<i>Stage<sub>1</sub></i>	85.31	99.36	62.57	83.71	57.48	79.13	55.97	76.53
		<i>Stage<sub>1&amp;2</sub></i>	84.27	99.17	61.48	83.42	57.39	79.11	55.48	75.68
		<i>Stage<sub>1&amp;3</sub></i>	<b>86.33</b>	<b>100.00</b>	61.34	83.66	57.24	78.54	54.76	76.42
		SMTAL	84.90	98.44	<b>64.34</b>	<b>85.93</b>	<b>60.81</b>	<b>81.56</b>	<b>58.99</b>	<b>79.83</b>

Types	Number	Validation(%)		ViPeR(%)		PRID(%)		GRID(%)		
		Top-1	Top-2	Top-1	Top-2	Top-1	Top-2	Top-1	Top-2	
LowerStyle	2	<i>Stage<sub>1</sub></i>	96.27	<b>100.00</b>	90.23	<b>100.00</b>	81.97	<b>100.00</b>	79.19	<b>100.00</b>
		<i>Stage<sub>1&amp;2</sub></i>	95.69	<b>100.00</b>	89.96	<b>100.00</b>	81.05	<b>100.00</b>	78.56	<b>100.00</b>
		<i>Stage<sub>1&amp;3</sub></i>	<b>96.42</b>	<b>100.00</b>	90.15	<b>100.00</b>	82.01	<b>100.00</b>	79.08	<b>100.00</b>
		SMTAL	96.33	<b>100.00</b>	93.24	<b>100.00</b>	<b>84.44</b>	<b>100.00</b>	<b>81.70</b>	<b>100.00</b>
LowerType	10	<i>Stage<sub>1</sub></i>	<b>93.55</b>	<b>100.00</b>	70.39	90.14	64.32	89.44	41.75	64.35
		<i>Stage<sub>1&amp;2</sub></i>	92.29	99.42	70.11	90.27	63.57	89.28	40.19	63.70
		<i>Stage<sub>1&amp;3</sub></i>	93.27	99.33	69.74	89.36	64.68	88.93	41.38	63.77
		SMTAL	91.67	98.96	<b>74.92</b>	<b>95.94</b>	<b>66.20</b>	<b>91.66</b>	<b>44.04</b>	<b>67.03</b>
HairStyle	3	<i>Stage<sub>1</sub></i>	<b>96.79</b>	<b>100.00</b>	68.59	94.19	70.13	94.53	66.23	95.43
		<i>Stage<sub>1&amp;2</sub></i>	96.23	<b>100.00</b>	67.13	94.00	68.47	94.53	66.12	94.13
		<i>Stage<sub>1&amp;3</sub></i>	<b>96.93</b>	<b>100.00</b>	68.48	94.32	69.88	94.89	67.05	95.00
		SMTAL	96.88	<b>100.00</b>	<b>71.42</b>	<b>95.91</b>	<b>71.51</b>	<b>96.81</b>	<b>68.82</b>	<b>96.03</b>
FootStyle	6	<i>Stage<sub>1</sub></i>	<b>71.53</b>	95.06	35.38	64.63	45.34	72.53	30.15	55.91
		<i>Stage<sub>1&amp;2</sub></i>	71.59	94.69	33.62	63.92	44.43	72.04	29.96	56.82
		<i>Stage<sub>1&amp;3</sub></i>	<b>71.66</b>	<b>95.88</b>	34.96	64.55	46.02	73.10	30.05	56.39
		SMTAL	70.83	94.79	<b>38.36</b>	<b>66.25</b>	<b>47.75</b>	<b>74.14</b>	<b>33.80</b>	<b>60.13</b>
UpperColor	11	<i>Stage<sub>1</sub></i>	79.35	93.23	44.33	67.11	28.21	51.03	30.25	52.47
		<i>Stage<sub>1&amp;2</sub></i>	78.02	91.68	42.96	66.49	27.63	50.41	30.23	51.39
		<i>Stage<sub>1&amp;3</sub></i>	<b>79.44</b>	91.99	45.36	66.52	28.04	50.88	31.18	53.00
		SMTAL	78.12	91.15	<b>47.92</b>	<b>69.13</b>	<b>28.97</b>	<b>52.31</b>	<b>32.90</b>	<b>55.04</b>
LowerColor	11	<i>Stage<sub>1</sub></i>	<b>95.62</b>	98.96	47.76	72.36	38.06	70.53	50.91	76.35
		<i>Stage<sub>1&amp;2</sub></i>	94.18	98.78	46.33	<b>71.45</b>	37.82	69.65	50.13	76.42
		<i>Stage<sub>1&amp;3</sub></i>	94.37	<b>100.00</b>	48.25	71.88	38.27	70.48	49.97	75.99
		SMTAL	94.23	98.44	<b>51.75</b>	<b>75.69</b>	<b>40.99</b>	<b>72.82</b>	<b>52.18</b>	<b>77.89</b>
HairColor	11	<i>Stage<sub>1</sub></i>	96.39	<b>99.13</b>	55.37	76.14	48.11	80.37	44.35	64.50
		<i>Stage<sub>1&amp;2</sub></i>	95.44	98.27	54.23	75.38	46.27	81.04	43.19	65.39
		<i>Stage<sub>1&amp;3</sub></i>	<b>96.50</b>	98.69	56.30	<b>76.23</b>	48.09	82.10	43.96	66.42
		SMTAL	95.84	97.92	<b>57.74</b>	<b>79.30</b>	<b>50.80</b>	<b>84.49</b>	<b>47.09</b>	<b>70.71</b>
FootColor	11	<i>Stage<sub>1</sub></i>	<b>93.68</b>	<b>98.66</b>	44.28	68.13	62.45	83.69	30.30	60.13
		<i>Stage<sub>1&amp;2</sub></i>	92.71	96.75	43.19	67.48	60.97	82.47	28.56	59.48
		<i>Stage<sub>1&amp;3</sub></i>	93.07	97.58	45.07	68.25	62.37	83.02	30.66	59.88
		SMTAL	92.19	96.88	<b>46.03</b>	<b>70.84</b>	<b>67.14</b>	<b>84.81</b>	<b>34.48</b>	<b>63.80</b>



# Examples of predicted attributes



personalLess45  
personalMale  
carryingMessengerBag  
accessoryNothing  
upperBodyLongSleeve  
upperBodyCasual  
**upperBodyPlaid**  
lowerBodyCasual  
lowerBodyTrousers  
hairShort  
footwearLeatherShoes  
upperBodyGreen  
lowerBodyBlack  
**hairYellow**  
footwearBrown



personalLess30  
personalMale  
carryingBackpack  
accessoryNothing  
upperBodyLongSleeve  
upperBodyCasual  
upperBodyThickStripes  
lowerBodyCasual  
lowerBodyTrousers  
hairShort  
footwearSneakers  
**upperBodyRed**  
**lowerBodyGrey**  
hairGrey  
footwearBrown



personalLess30  
personalMale  
carryingBackpack  
accessoryNothing  
upperBodyLongSleeve  
upperBodyCasual  
**upperBodyPlaid**  
lowerBodyCasual  
lowerBodyTrousers  
hairShort  
footwearSneakers  
upperBodyBrown  
lowerBodyGrey  
**hairGrey**  
footwearBrown



personalLess30  
personalFemale  
carryingNothing  
accessoryNothing  
upperBodyLongSleeve  
upperBodyCasual  
upperBodyOther  
lowerBodyCasual  
lowerBodyJeans  
hairLong  
footwearShoes  
upperBodyRed  
**lowerBodyGrey**  
hairBlack  
footwearBlack



personalLess30  
personalMale  
carryingNothing  
accessoryHat  
**upperBodyNoSleeve**  
upperBodyCasual  
upperBodySweater  
lowerBodyCasual  
lowerBodyTrousers  
**hairLong**  
footwearSneakers  
upperBodyWhite  
**lowerBodyWhite**  
hairBlack  
footwearWhite



personalLarger60  
personalFemale  
carryingPlasticBags  
accessoryNothing  
upperBodyLongSleeve  
upperBodyCasual  
**upperBodyJacket**  
lowerBodyCasual  
lowerBodyTrousers  
hairShort  
footwearLeatherShoes  
**upperBodyGrey**  
lowerBodyBrown  
hairWhite  
footwearBlack



personalLess30  
personalFemale  
carryingLuggageCase  
accessoryNothing  
upperBodyLongSleeve  
upperBodyCasual  
**upperBodyThickStripes**  
lowerBodyCasual  
lowerBodyTrousers  
hairLong  
footwearSneakers  
**upperBodyGrey**  
lowerBodyBrown  
**hairYellow**  
footwearBlack



personalLess30  
personalMale  
carryingBackpack  
accessorySunglasses  
upperBodyShortSleeve  
upperBodyCasual  
upperBodyTshirt  
lowerBodyCasual  
lowerBodyShorts  
hairShort  
footwearSneakers  
upperBodyGreen  
lowerBodyWhite  
hairBlack  
footwearWhite



北京大学

C. Su, S. Zhang, J. Xing, Q. Tian, and W. Gao. Multi-type attributes driven multi-camera person re-identification, *Pattern Recognition*, 2018.

# Comparsions-1

## ■ VIPeR:

Methods		Rank 1	Rank 5	Rank 10	Rank 20
Metric Learning based ReID	RPML [24]	27.0	57.0	69.0	83.0
	Salmatch [63]	30.2	52.4	65.5	79.1
	LMF [46]	29.1	52.3	65.9	80.0
	KISSME [27]	19.6	47.5	62.2	77.0
	KCCA [79]	37.3	71.4	<b>84.6</b>	92.3
	kLFDA [28]	32.2	65.8	79.7	90.9
	LOMO + XQDA [34]	40.0	68.9	81.5	91.1
	CSL [36]	34.8	68.7	82.3	91.8
	MLAPG [37]	40.7	69.9	82.3	92.4
Traditional Attributes based ReID	TSR [80]	31.6	68.6	82.8	<b>94.6</b>
	EPKFM [33]	36.8	70.4	83.7	91.7
	AIR [39]	18.0	38.8	51.1	71.2
Deep Learning based ReID	OAR [40]	21.4	41.5	55.2	71.5
	LORAE [87]	42.3	72.2	81.6	89.6
	IDLA [52]	34.8	54.3	76.5	87.6
Proposed	DML [51]	28.2	59.3	73.5	86.4
	Deep-RDC [53]	40.5	60.8	70.4	84.4
	Deep-TCP [54]	<b>47.8</b>	<b>74.7</b>	<b>84.8</b>	<b>91.1</b>
	Stage <sub>1</sub>	36.2	63.9	73.5	84.7
Proposed	SSDAL	37.9	65.5	75.6	88.4
	SSDAL + XQDA	43.5	71.8	81.5	89.0
	SMTAL	39.7	66.9	76.5	86.6
	SMTAL + XQDA	47.1	71.5	80.3	88.2





# Comparsions-2

## PRID:

Methods	Rank 1	Rank 5	Rank 10	Rank 20
RPML [24]	4.8	14.3	21.6	30.2
PRDC [31]	4.5	12.6	19.7	29.5
RSVM [21]	6.8	16.5	22.7	31.5
Salmatch [63]	4.9	17.5	26.1	33.9
LMF [46]	12.5	23.9	30.7	36.5
PCCA [23]	3.5	10.9	17.9	27.1
KISSME [27]	4.1	12.8	21.1	31.8
kLFDA [28]	7.6	18.9	25.6	37.4
KCCA [79]	14.5	34.3	46.7	59.1
LOREA [87]	18.0	37.4	50.1	66.6
LOMO + XQDA [34]	15.3	35.7	41.2	53.8
MLAPG [37]	16.6	33.1	41.4	52.5
Deep-TCP [54]	22.0	-	47.0	57.0
<i>Stage<sub>1</sub></i>	19.6	46.7	55.1	66.4
SSDAL	20.1	47.4	55.7	68.6
SSDAL + XQDA	22.6	48.7	57.8	69.2
SMTAL	22.4	47.8	56.8	67.6
SMTAL + XQDA	<b>24.4</b>	<b>52.3</b>	<b>62.5</b>	<b>74.2</b>



# Comparsions-3

---

## ■ GRID:

Methods	Rank 1	Rank 5	Rank 10	Rank 20
PRDC [31]	9.7	22.0	33.0	44.3
RSVM [21]	10.2	24.6	33.3	43.7
MRank-PRDC [67]	11.1	26.1	35.8	46.6
MRank-RSVM [67]	12.2	27.8	36.3	49.3
RQDA [92]	15.2	30.1	39.2	49.3
EPKFM [33]	16.3	35.8	46.0	57.6
LOMO + XQDA [34]	16.6	35.4	41.8	52.4
<i>Stage<sub>1</sub></i>	17.5	34.5	42.8	55.3
SSDAL	19.1	35.6	45.8	58.1
SSDAL + XQDA	22.4	39.2	48.0	58.4
SMTAL	19.2	38.1	47.8	58.7
<b>SMTAL + XQDA</b>	<b>23.4</b>	<b>39.6</b>	<b>49.8</b>	<b>60.3</b>



# Comparsions-4

## ■ Market:

Single Query	Rank 1	mAP
Salmatch [63]	20.5	8.2
SDALF [13]	33.5	13.5
BGG [19]	34.4	14.1
KISSME [27]	40.5	19.0
MFA [28]	45.7	18.2
kLFDA [28]	<b>51.3</b>	24.4
LOMO + XQDA [34]	43.8	22.2
<b>SSDAL</b>	39.4	19.6
<b>SMTAL</b>	49.5	<b>29.2</b>

Multiple Query	Rank 1	MAP
BGG+MultiQ_max [19]	42.1	18.5
kLFDA [28]	52.7	27.4
LOMO + XQDA [34]	54.1	28.4
<b>SSDAL</b>	49.0	25.8
<b>SMTAL</b>	<b>56.6</b>	<b>31.2</b>



# Summary

---

- Attribute is robust, effective, and efficient for Person Re-ID
- Attribute correlation is important
- Our deep model effectively learns attributes with a limited number of labeled data
- It works well even when the training and testing data are independent

## Reference:

1. C. Su, **S. Zhang**, J. Xing, Q. Tian, and W. Gao. Multi-type attributes driven multi-camera person re-identification, *Elsevier Pattern Recognition (PR)*, 75: 77-89, 2018.
2. C. Su, F. Yang, **S. Zhang**, Q. Tian, L. S. Davis, and W. Gao. Multi-Task learning with low rank attribute embedding for multi-camera person re-identification. *IEEE Trans. on Pattern Analysis and Machine Intelligence (T-PAMI)*, 2018
3. C. Su, **S. Zhang**, J. Xing, W. Gao, and Qi Tian. Deep attributes driven multi-camera person re-identification. *European Conference on Computer Vision (ECCV)*, 2016.
4. C. Su, F. Yang, **S. Zhang**, Q. Tian, L. Davis, W. Gao. Multi-Task Learning with Low Rank Attribute Embedding for Person Re-identification, *IEEE International Conference on Computer Vision (ICCV)*, 2015.



# Outline

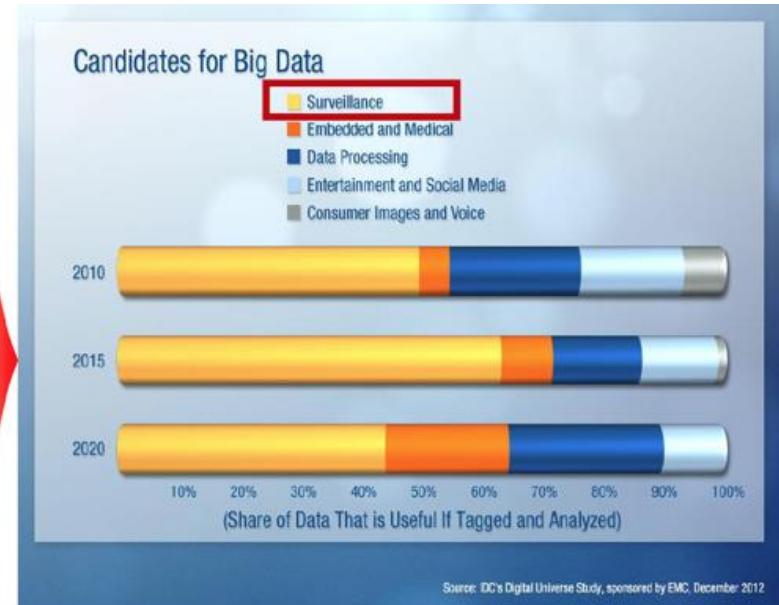
---

- Attribute Learning
  - Attribute embedding in multi-task learning
  - Attribute learning with deep models
- Scalable Person Re-ID
  - Deep feature extraction and compression
  - Off-line index optimization



# Scalability is required for person re-id

- Camera is everywhere in smart city



## Camera Networks in Smart City

- Big data, fast growth
  - 30 millions cameras have been deployed in China

## More than half of all big data

- Complicated content, rich Information
  - Cars, persons, others



北京大学



# Issues

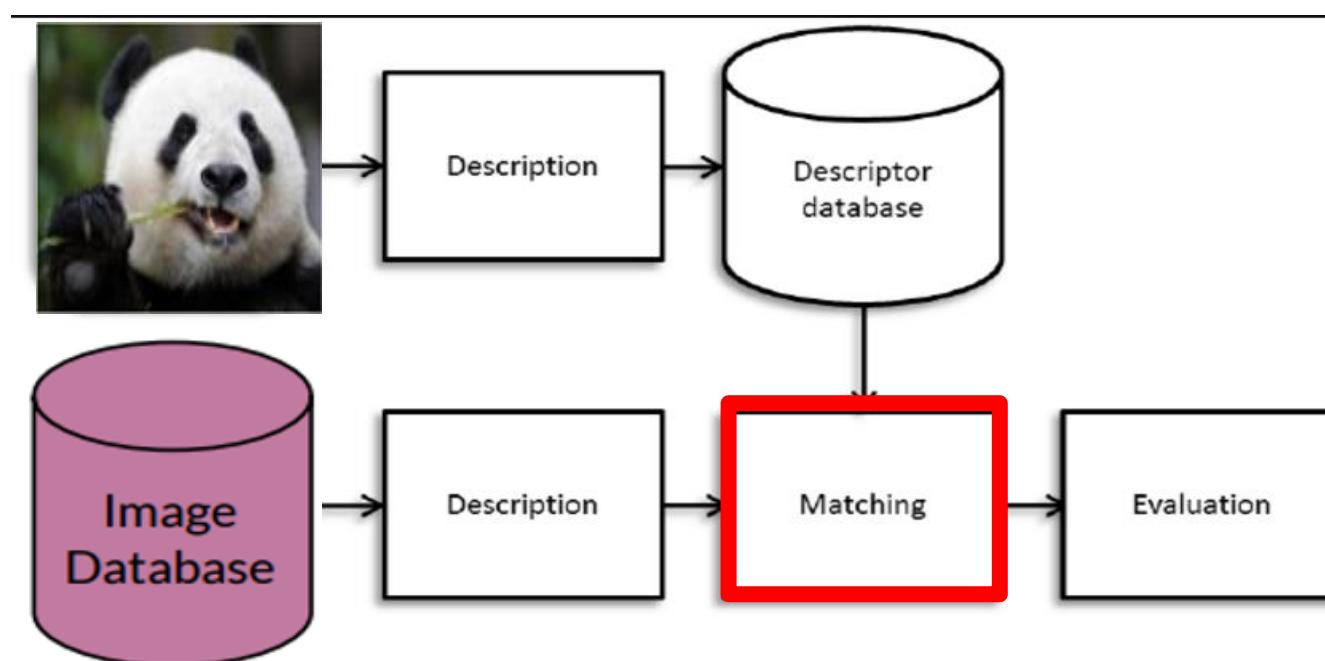
---

- Most of existing algorithms are designed and evaluated on small-scale datasets,
  - CHUK03 consists of 14,096 images, 1,467 identities,
  - Market1501 consists of 36,036 images, 1,501 identities
- Working on those datasets, existing methods focus on improving the accuracy, and pay less attention to the efficiency and generalization ability
- person Re-ID can be conducted by matching the query person image to database images and return only the images of the same person.
- This is similar to image retrieval



# Large-Scale Person Re-ID as Retrieval

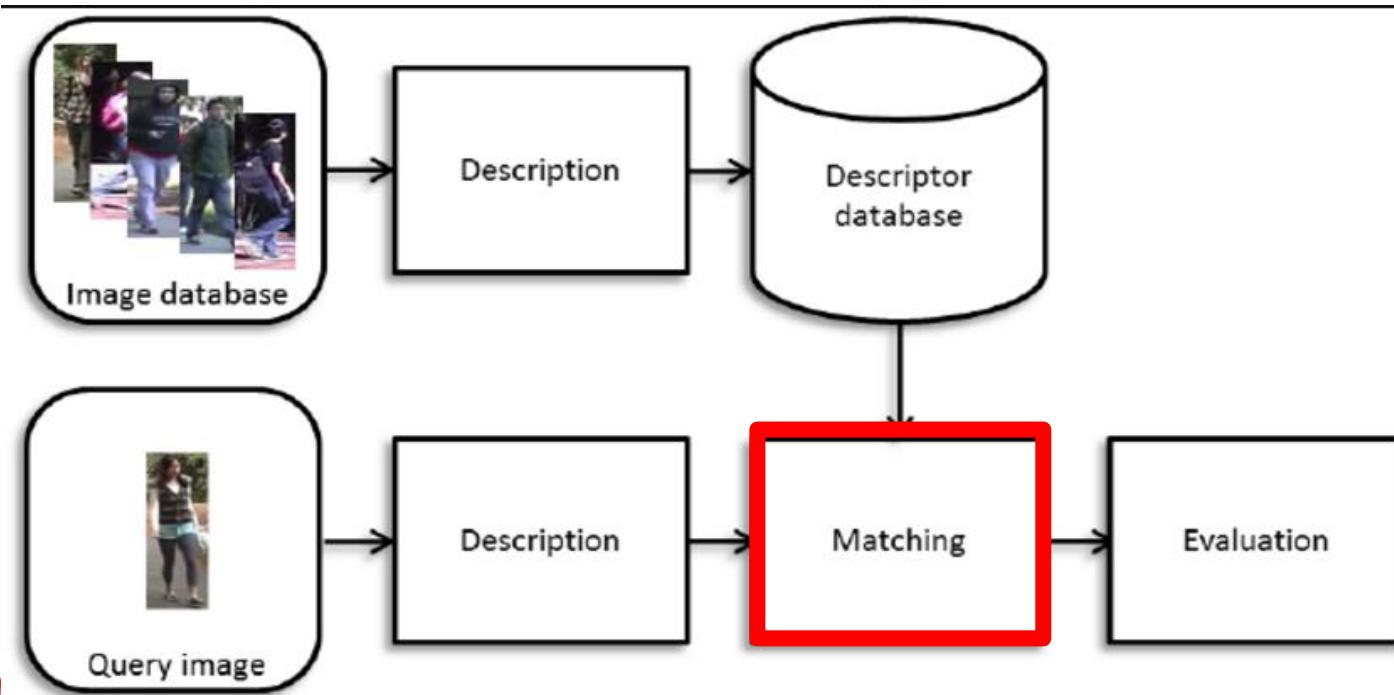
- Image Retrieval: Hash could be used for large-scale image retrieval to accelerate online matching.





# Large-Scale Person Re-ID as Retrieval

- Person Re-ID is commonly conducted with linear search
- Time-consuming for large-scale dataset
- research efforts on image retrieval could make the large-scale person Re-ID matching more efficient





# Challenges



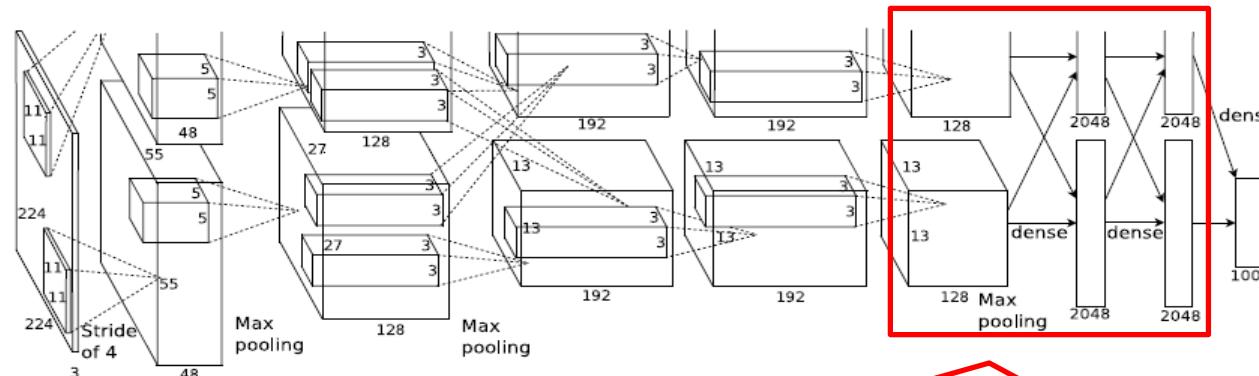
**Fig. 1.** Example images of four persons from Market1501 (first two) and CUHK03 (last two). The subtle differences among different persons and the large variance among images of the same person make person ReID challenging.

The key is image representation, which should be fast and discriminative for unseen persons



# Convolutional Neural Network

- AlexNet [Krizhevsky *et al*, NIPS,2012].
  - CNN = **Feature Extraction** + **Classifier**
  - Disadvantage of Fully-Connected Layer
    - Huge parameters make the network easy to overfitting.
    - Fully-connected operation makes the input for CNN is fixed size image.

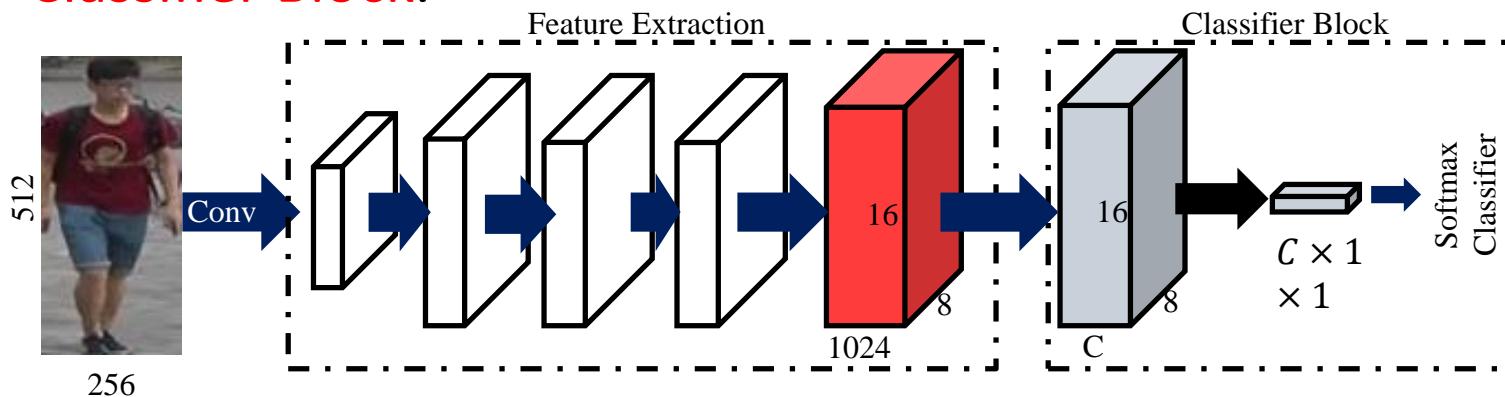


Pool5->Fc6:  $9126 \times 4096 = 35.6M$   
Fc6->Fc7:  $4096 \times 4096 = 16M$

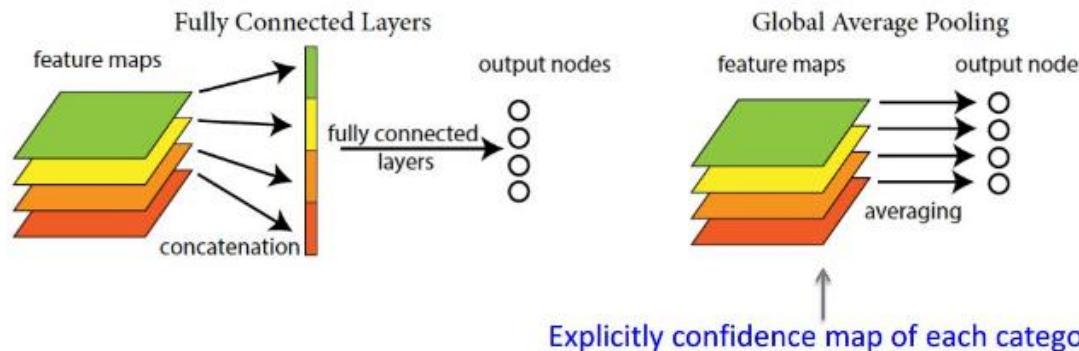


# Conv-Net for feature learning

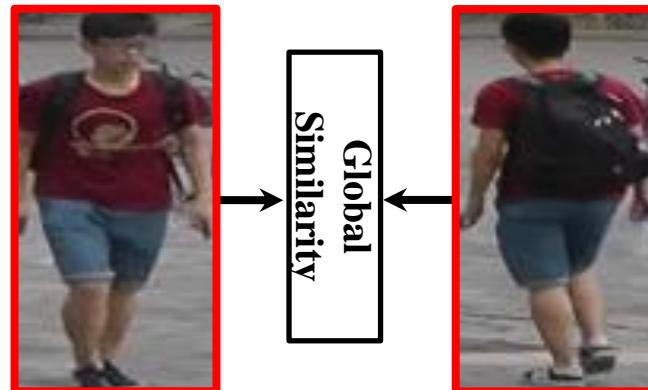
- Replace the Fully-Connected classifier with **Convolutional Classifier Block**.



- Convolutional Classifier Block= Convolutional Operation+ Global Average Pooling (GAP).



# Global Feature Extraction



- Person Descriptor = **Global Descriptor** + Local Descriptor
  - Global descriptor is used to describe the global appearance of the person.
  - Global Descriptor:  $f^g$ ,

$$f^g = [f_1, f_2, \dots, f_K], f_k = \frac{1}{W \times H} \sum_{h=1}^H \sum_{w=1}^W \chi_{k,h,w}$$

Where  $\chi$  denotes the output of the last layer of feature extraction block.



# Local Feature Extraction

- Person Descriptor = Global Descriptor + Local Descriptor
  - Global descriptor lacks the ability to describe the local parts.
  - For person Re-ID images, most person images are weakly aligned.
  - Therefore, we divide the image into four stripes to generate local descriptors.

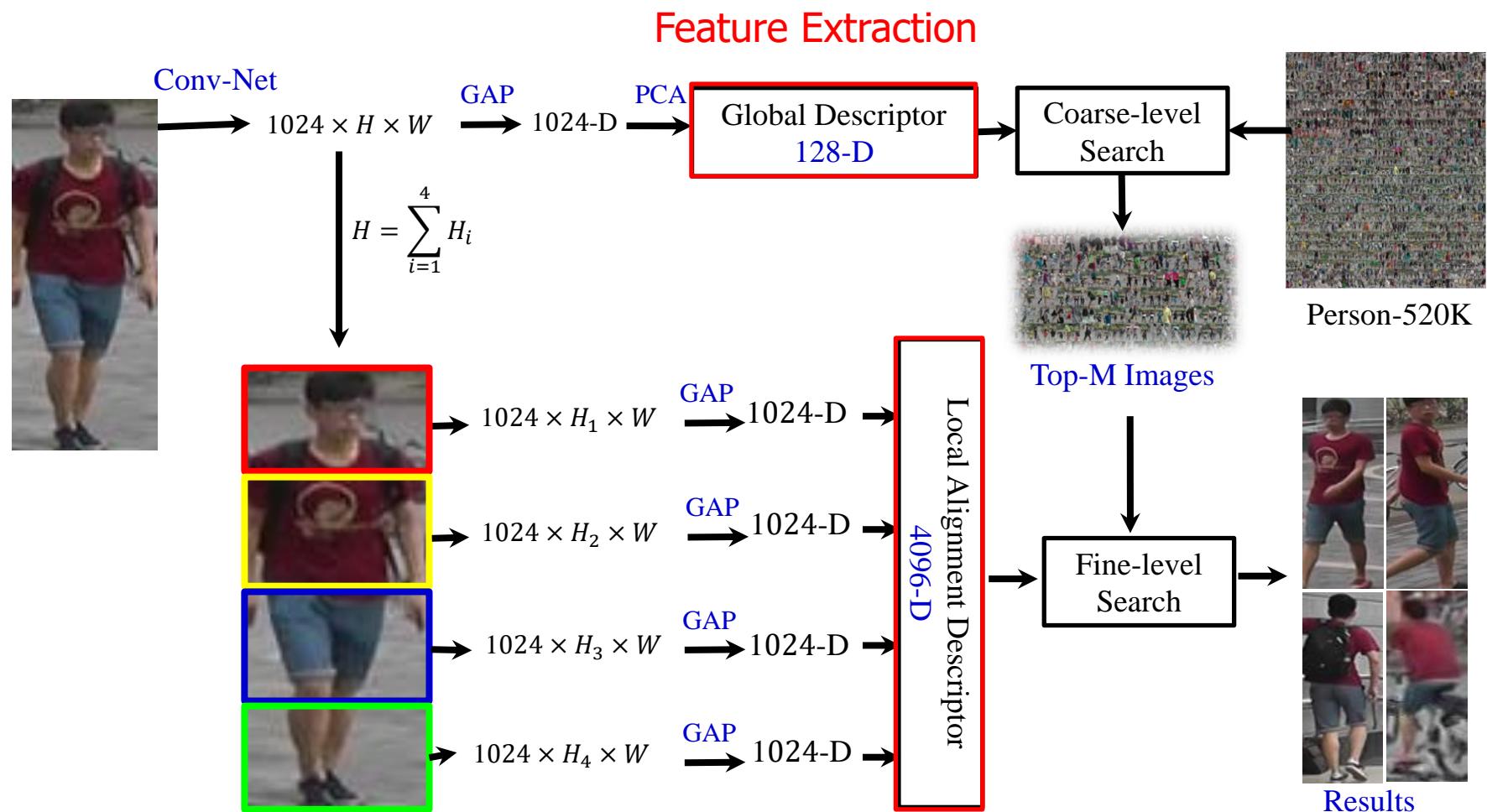
$$\begin{aligned} \mathbf{f}^l &= [\mathbf{f}^{l_1}; \mathbf{f}^{l_2}; \mathbf{f}^{l_3}; \mathbf{f}^{l_4}], \\ \mathbf{f}^{l_i} &= [f_1, f_2, \dots, f_K], \end{aligned}$$



$$f_k = \frac{1}{W \times \mathcal{H}} \sum_{h=(i-1) \times \mathcal{H}}^{i \times \mathcal{H}} \sum_{w=1}^W \chi_{k,h,w}, \quad i \in [1,4] \text{ and } \mathcal{H} = \frac{H}{4}$$



# Proposed Approach





# Online Retrieval

---

- Coarse-to-Fine Retrieval
  - Global Descriptor used for coarse retrieval.
    - Feature dim is reduced to 128-D
  - Local Descriptor used for fine re-ranking.
    - Use 4096-D feature for fine retrieval



# Dataset Construction

- A large-scale Person-520K dataset is constructed for evaluation.
- Combines CUHK03 and Market1501

**Table 1.** Details of Person-520K dataset.

	Training		Testing		Query	
	Identity	Image	Identity	Image	Identity	Image
CUHK03	736	7,029	731	5,606	731	1,461
Market1501	751	12,936	750	19,732	750	3,368
Distractors	-	-	-	500K	-	-
Total	1,487	19,965	1,481	520K	1,481	4,829

The code and Person-520K could be  
downloaded from:  
<http://www.yahantao.com/>

58



# Performance of Conv-Net

- Input size 224 \* 224
- Convolutional Classifier block is better than fully-connected classifier.

Comparison on Market1501

Models	mAP(%)	Rank-1 (%)
AlexNet [13]	26.79	50.89
VGG16Net [14]	38.27	65.02
GoogleNet [16]	48.24	70.27
Res50Net [15]	51.48	73.69
Conv-Net	<b>54.86</b>	<b>76.36</b>



# Effect of network input size

---

Larger input size corresponds to better performance

**Table 3.** Effect of various input sizes on Market1501.

Scales	mAP(%)	Rank-1 (%)
$224 \times 224$	54.86	76.36
$128 \times 64$	44.57	68.27
$256 \times 128$	53.50	75.89
$384 \times 192$	56.04	77.64
$512 \times 256$	61.9	81.5

# Performance on Person-520K

**Table 6.** The performance on Person-520K. “Time” denotes the retrieval time. Feature extraction takes about 50ms for Conv-Net on GPU.

Methods	Dim	Time(ms)	mAP(%)	Rank-1(%)
AlexNet	4,096	3932	17.13	33.46
GoogleNet	1,024	960	36.38	56.05
Conv-Net_f <sup>g</sup>	1,024	961	43.42	60.84
Conv-Net_f <sup>g</sup>	512	500	41.06	61.79
Conv-Net_f <sup>g</sup>	256	262	40.84	61.21
Conv-Net_f <sup>g</sup>	128	148	39.99	59.3
Conv-Net_f	5,120	4746	46.95	64.60
coarse-level	128	150	39.99	59.3
C2F	-	180	46.74	64.58

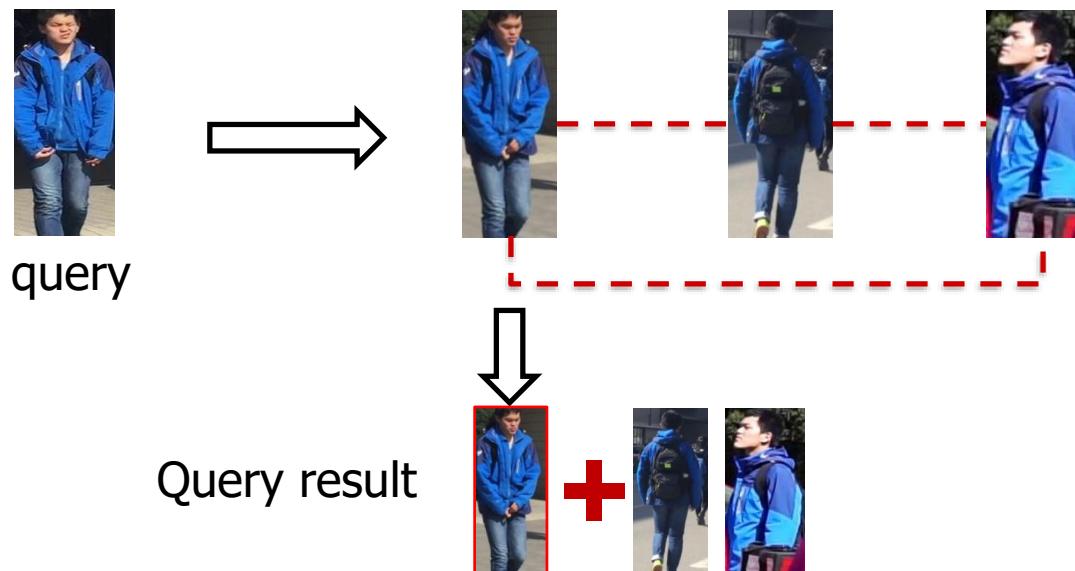
- Reduce dimensionality with PCA drops the accuracy
- C2F achieves a good balance between accuracy and efficiency



# Introduce indexing to ReID

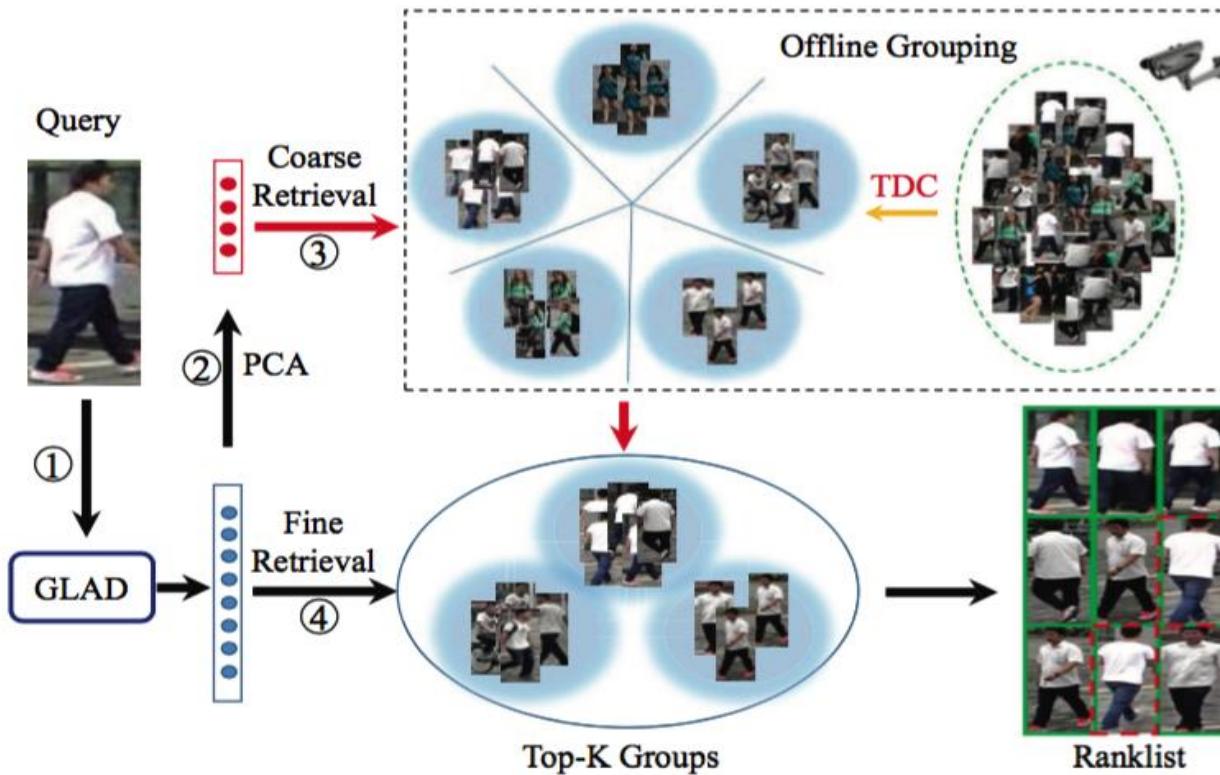
## □ Basic Idea

- each person has multiple samples in the gallery set
- Can we link those images during offline indexing?
- The re-id accuracy can be improved, if the linkage is accurate and smart



# Indexing & Retrieval Framework

- Generate image groups for indexing,
- Two retrieval stages, to retrieve groups and images, respectively.
  - Offline group generation can use more expensive features
  - Group number smaller than image number





# Group Generation

- Compute the image distance with Jaccard distance

$$dis(I_i^a, I_j^b) = 1 - \frac{|\overline{KNN}(I_i^a, K) \cap \overline{KNN}(I_j^b, K)|}{|\overline{KNN}(I_i^a, K) \cup \overline{KNN}(I_j^b, K)|}, \rightarrow \text{K-reciprocal neighbor set}$$

- Propose a Two-fold Divisive Clustering (TDC)
  - Hierarchically divides image sets into two sets and ensure the dissimilarity within each group smaller than a threshold
  - Only need one parameter  $\Theta$ , need not to specify the group number
- Use average pooling for group feature generation:

$$\mathbf{f}^G(i) = \frac{1}{N} \sum_{j=1}^N \mathbf{f}_j^{GLAD}(i),$$

# Efficiency and Accuracy Boost

TABLE IX

RE-ID PERFORMANCE OF OUR RETRIEVAL FRAMEWORK ON MARKET1501 WITH DIFFERENT  $\theta$  AND  $K$ . THE SECOND ROW DENOTES THE PERFORMANCE OF LINEAR SEARCH WITH THE ORIGINAL GLAD.

$\theta$	$K$	Group Number	Dim	mAP	Rank-1	Times(ms)
-	-	19732	4096	73.9	89.9	352.0
0.2	10	16692	128	73.8	89.7	9.7
0.4	10	12785	128	74.8	89.8	7.7
0.6	10	8302	128	76.0	89.7	5.3
0.8	10	4911	128	76.4	89.8	4.5
0.8	20	3726	128	75.9	89.7	5.2
0.8	30	3532	128	75.4	89.6	6.2
0.8	40	3493	128	75.0	89.6	6.8
0.8	50	3433	128	74.7	89.6	6.5

100X speed up

Rank-1 Acc. improves 2.5%



# Performance on Person520K

---

- Improved accuracy
- Improved efficiency

TABLE XI  
RE-ID PERFORMANCE OF OUR RETRIEVAL FRAMEWORK ON  
PERSON-520K.

Method	Group Number	Dim	mAP	Rank-1	Speedup Ratio
[56]	525K	128	46.7	64.6	26X
GLAD+Linear	525K	4096	56.5	72.7	1X
GLAD+C2F	<b>199K</b>	128	<b>58.2</b>	<b>72.9</b>	<b>61X</b>



# Conclusions

---

- Person ReID and image retrieval are essentially the same
- A coarse-to-fine framework is proposed for large-scale person ReID.
- Off-line indexing boosts both efficiency and accuracy

## Reference:

1. H. Yao, **S. Zhang**, D. Zhang, Y. Zhang, J. Li, Y. Wang, Q. Tian. Large-Scale Person Re-Identification as Retrieval. ICME, 2017.
2. L. Wei, **S. Zhang**, H. Yao, W. Gao, and Q. Tian. GLAD: Global-Local-Alignment Descriptor for Pedestrian Retrieval. ACM MM, 2017.



---

# Thanks!

Homepage



[www.pkuvmc.com](http://www.pkuvmc.com)