

# Part-based Pseudo Label Refinement for Unsupervised Person Re-identification

윤성의

SGVR Lab, KAIST

# Project Guidelines: Project Topics

---

- **Any topics related to the course theme are okay**
  - **You can find topics by browsing recent papers**

# Expectations

---

- **Mid-term project presentation**
  - **Introduce problems and explain why it is important**
  - **Give an overall idea on the related work**
  - **Explain what problems those existing techniques have**
  - **(Optional) explain how you can address those problems**
  - **Explain roles of each member**

# Expectations

---

- **Final-term project presentation**
  - **Cover all the materials that you talked for your mid-term project**
  - **Present your ideas that can address problems of those state-of-the-art techniques**
  - **Give your qualitatively (or intuitive) reasons how your ideas address them**
  - **Also, explain expected benefits and drawbacks of your approach**
  - **(Optional) backup your claims with quantitative results collected by some implementations**
  - **Explain roles of each members**

# A few more comments

---

- **Start to implement a paper, if you don't have any clear ideas**
  - **While you implement it, you may get ideas about improving it**

Speaker	Novelty of the project and idea (1 ~ 5)	Practical benefits of the method (1 ~ 5)	Completeness level of the project (1 ~ 5)	Total score (3 ~ 15)	Role of each student is clear and well balanced? (Yes or No)
XXX					
YYY					

# Class Objectives

---

- Person Re-identification
- Unsupervised Approaches
- Part-based Pseudo Label Refinement for Unsupervised Person Re-ID (CVPR 2022)

# Person Re-identification (Person Re-ID)

- Person re-ID aims to **retrieve a person corresponding to a given query** across disjoint camera views or different time stamps.
- Applications: Surveillance system, Finding a missing person, etc.

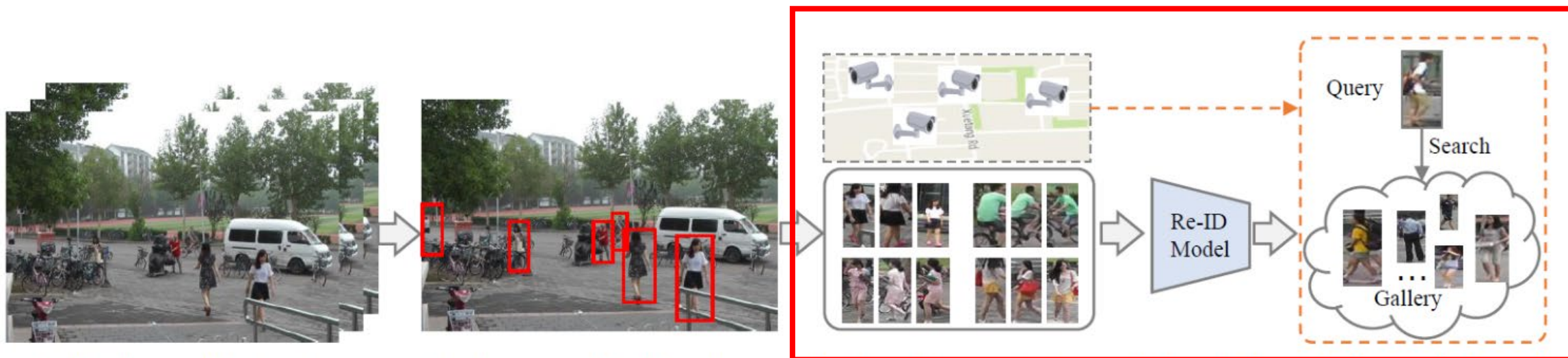


Fig. 1: The flow of designing a practical person Re-ID system, including five main steps: 1) *Raw Data Collection*, (2) *Bounding Box Generation*, 3) *Training Data Annotation*, 4) *Model Training* and 5) *Pedestrian Retrieval*.

# Person Re-identification (Person Re-ID)

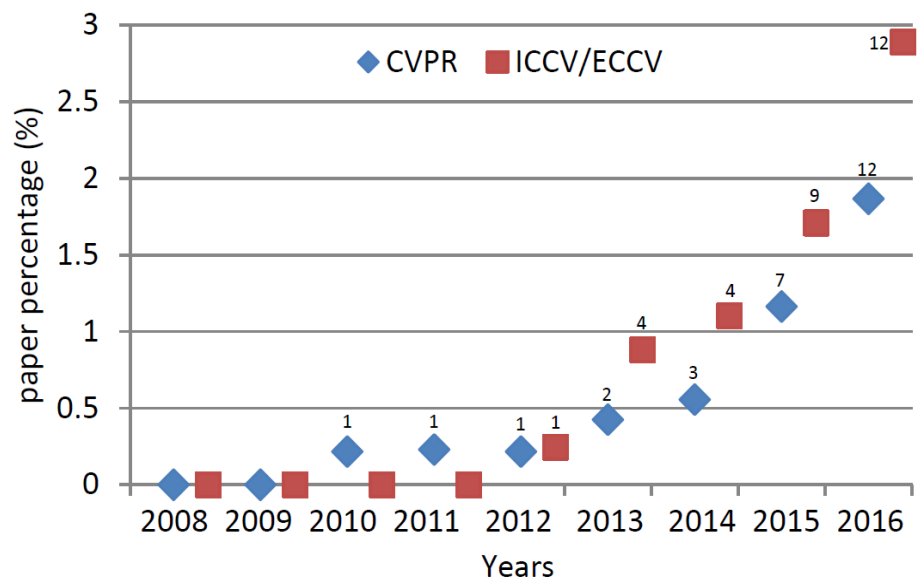


Fig. Percentage of person re-ID papers on top conferences over the years. Numbers above the markers indicate the number of re-ID papers.

Conference	Link	#Total	Person Re-ID
ICCV2021	<a href="#">click</a>	34	24
CVPR2021	<a href="#">click</a>	32	25
ECCV2020	<a href="#">Click</a>	30	23
CVPR2020	<a href="#">Click</a>	34	24
ICCV2019	<a href="#">Click</a>	39	33
CVPR2019	<a href="#">Click</a>	29	21
ECCV2018	<a href="#">Click</a>	19	15
CVPR2018	<a href="#">Click</a>	31	30
ICCV2017	<a href="#">Click</a>	16	14
CVPR2017	<a href="#">Click</a>	16	14

### 14. CVPR2023

- Person re-identification
  - 1) "Diverse Embedding Expansion Network and Low-Light Cross-Modality Benchmark for Visible-Infrared Person Re-Identification" [[paper](#)]
  - 2) "PHA: Patch-Wise High-Frequency Augmentation for Transformer-Based Person Re-Identification" [[paper](#)]
  - 3) "Shape-Erased Feature Learning for Visible-Infrared Person Re-Identification" [[paper](#)]
  - 4) "TranSG: Transformer-Based Skeleton Graph Prototype Contrastive Learning With Structure-Trajectory Prompted Reconstruction for Person Re-Identification" [[paper](#)]
  - 5) "PartMix: Regularization Strategy To Learn Part Discovery for Visible-Infrared Person Re-Identification" [[paper](#)]
  - 6) "Event-Guided Person Re-Identification via Sparse-Dense Complementary Learning" [[paper](#)]
  - 7) "Clothing-Change Feature Augmentation for Person Re-Identification" [[paper](#)]

**Awesome Person Re-identification (Person ReID), github**

Zheng et al. Person Re-identification: Past, Present and Future. In arXiv 2016.  
<https://github.com/bismex/Awesome-person-re-identification>.



# Person Search

- Task to detect the person of interest from the entire image
- We need to detect for the target person from a gallery of whole scene images before doing a re-ID



(b) Person search: finding from whole scene images



(a) Person re-id: matching with manually cropped pedestrians

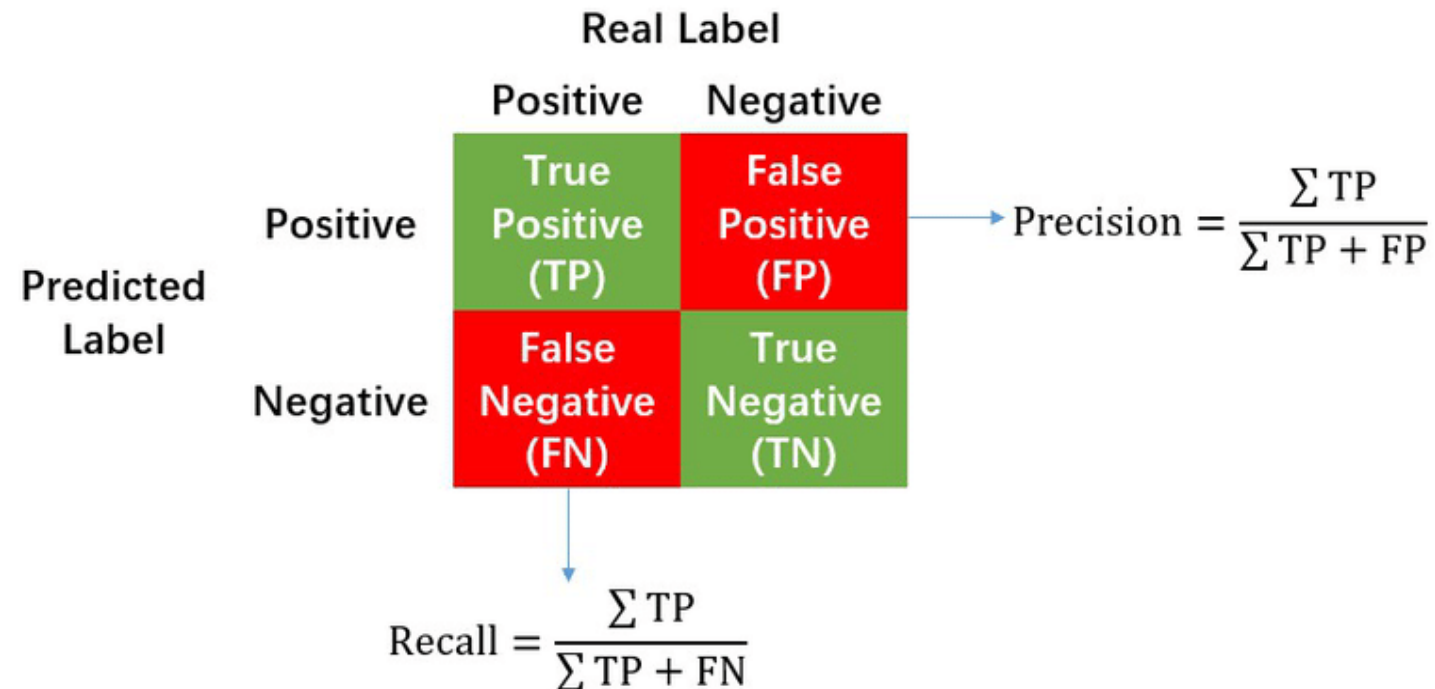
# Datasets

- The dataset scale (both #image and #ID) has increased rapidly.
- The camera number is greatly increased to approximate the large-scale camera network in practical scenarios.

Dataset	<i>Image datasets</i>						
	Time	#ID	#image	#cam.	Label	Res.	Eval.
VIPeR	2007	632	1,264	2	hand	fixed	CMC
iLIDS	2009	119	476	2	hand	vary	CMC
GRID	2009	250	1,275	8	hand	vary	CMC
PRID2011	2011	200	1,134	2	hand	fixed	CMC
CUHK01	2012	971	3,884	2	hand	fixed	CMC
CUHK02	2013	1,816	7,264	10	hand	fixed	CMC
CUHK03	2014	1,467	13,164	2	both	vary	CMC
Market-1501	2015	1,501	32,668	6	both	fixed	C&M
DukeMTMC	2017	1,404	36,411	8	both	fixed	C&M
Airport	2017	9,651	39,902	6	auto	fixed	C&M
MSMT17	2018	4,101	126,441	15	auto	vary	C&M
Dataset	<i>Video datasets</i>						
	time	#ID	#track(#bbox)	#cam.	label	Res.	Eval
PRID-2011	2011	200	400 (40k)	2	hand	fixed	CMC
iLIDS-VID	2014	300	600 (44k)	2	hand	vary	CMC
MARS	2016	1261	20,715 (1M)	6	auto	fixed	C&M
Duke-Video	2018	1,812	4,832 (-)	8	auto	fixed	C&M
Duke-Tracklet	2018	1,788	12,647 (-)	8	auto	C&M	
LPW	2018	2,731	7,694(590K)	4	auto	fixed	C&M
LS-VID	2019	3,772	14,943 (3M)	15	auto	fixed	C&M

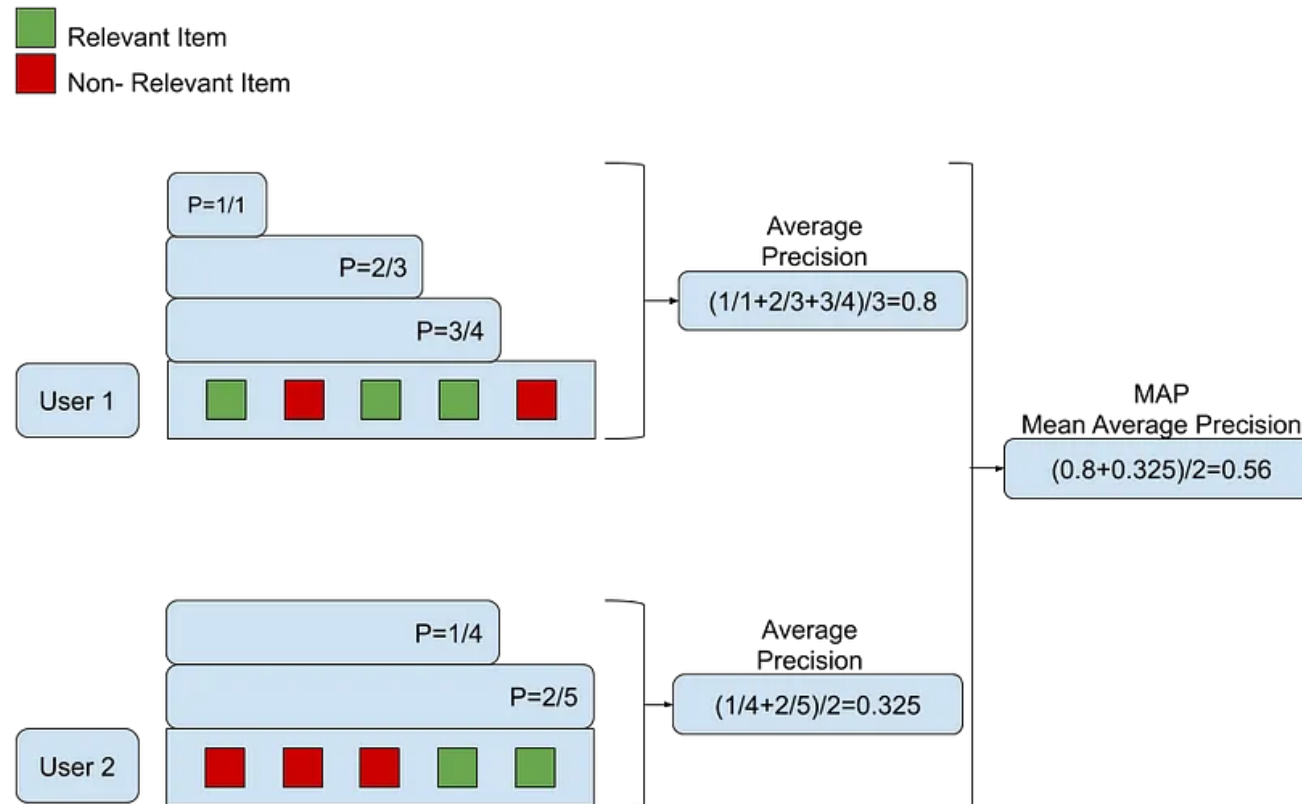
# Evaluation Metrics

- mean Average Precision (mAP)
- Cumulative Matching Characteristics (CMC- $k$ , Rank- $k$  matching accuracy); the probability that a correct match appears in the top- $k$  retrieved results.



# Evaluation Metrics

- mean Average Precision (mAP)
- Cumulative Matching Characteristics (CMC- $k$ , Rank- $k$  matching accuracy); the probability that a correct match appears in the top- $k$  retrieved results.



# Challenges in Person Re-ID

- Challenges by different camera views and time stamps.
  - Variance of viewpoints, illumination, pose, etc.
  - Occlusions.
  - Low resolutions.
- Large intra-variation & Small inter-variation



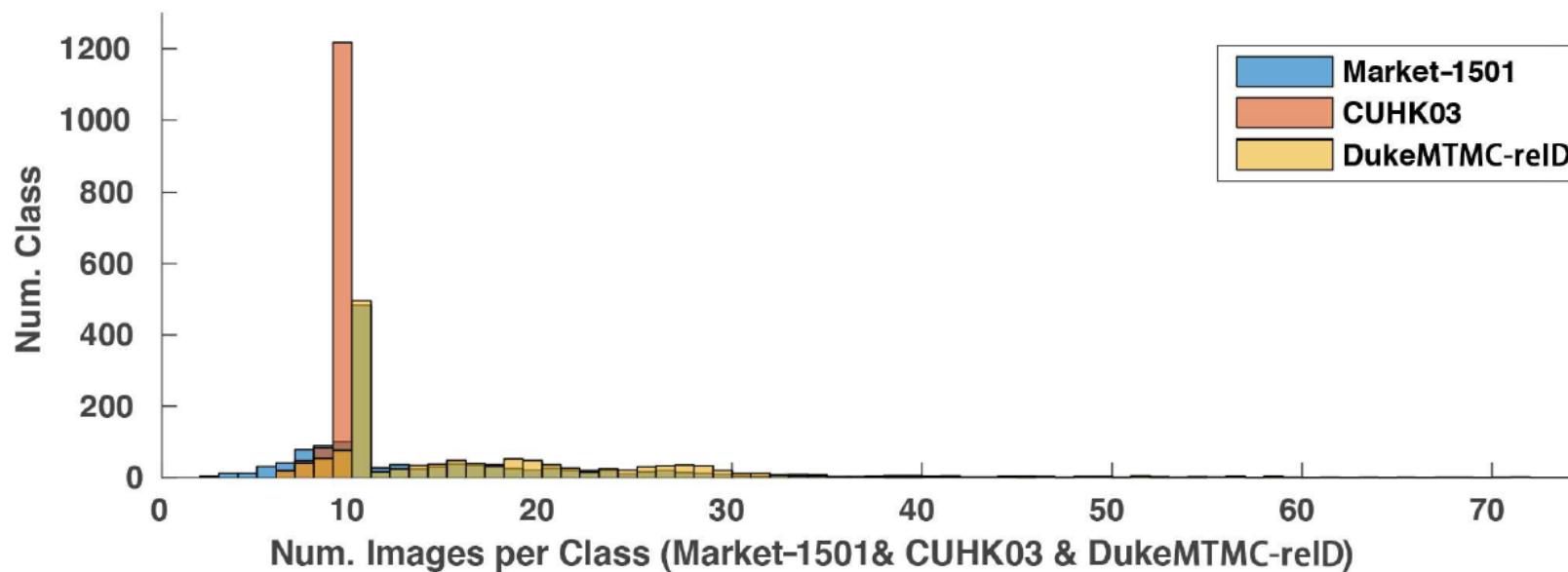
True match



False match

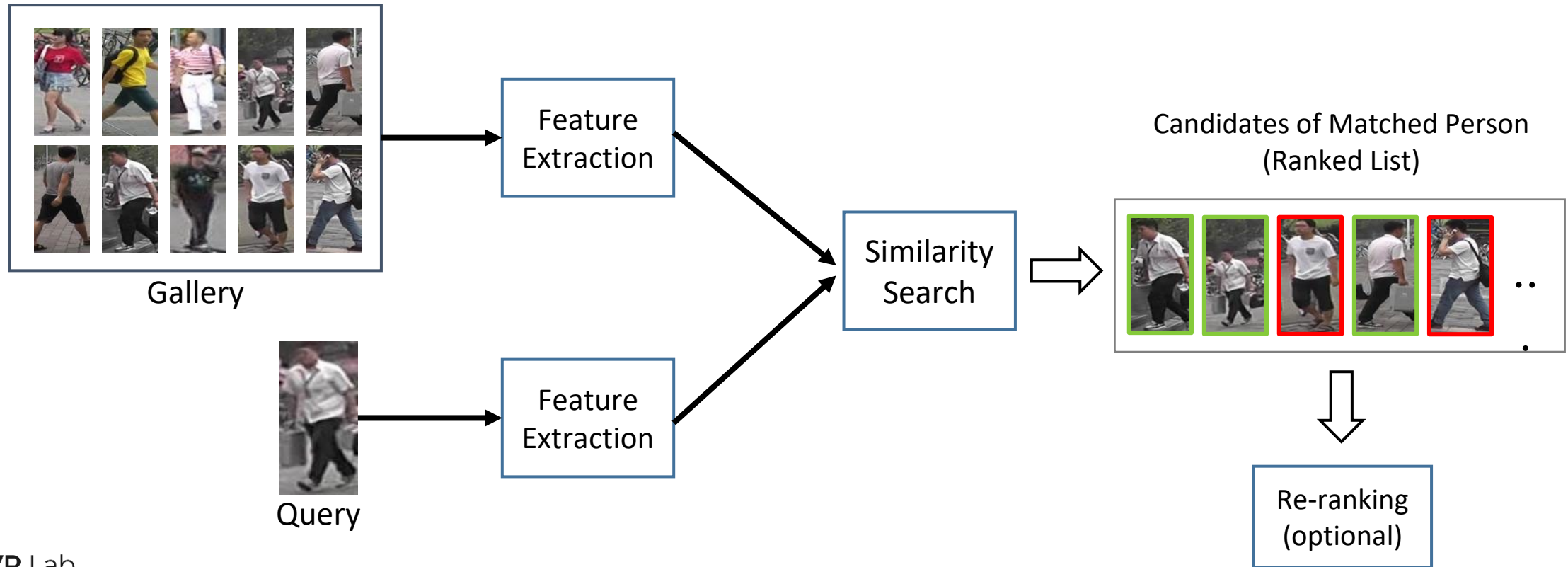
# Challenges in Person Re-ID

- Long-tail problem.
  - In person re-ID, all datasets suffer from the insufficient training set.
  - Long-tail distribution training sets can yield unstable convergence and overfitting to head distributions.
  - MNIST 10 class/ 5000 per class, CIFAR 100 class/500 per class.



# General Protocol of Person Re-ID

- Person re-identification pipeline.



# Feature Representation Learning for Person Re-ID

- Most studies focus on **learning discriminative representation for person retrieval**.
- Recently, deep neural networks (DNN) have provided powerful descriptors.

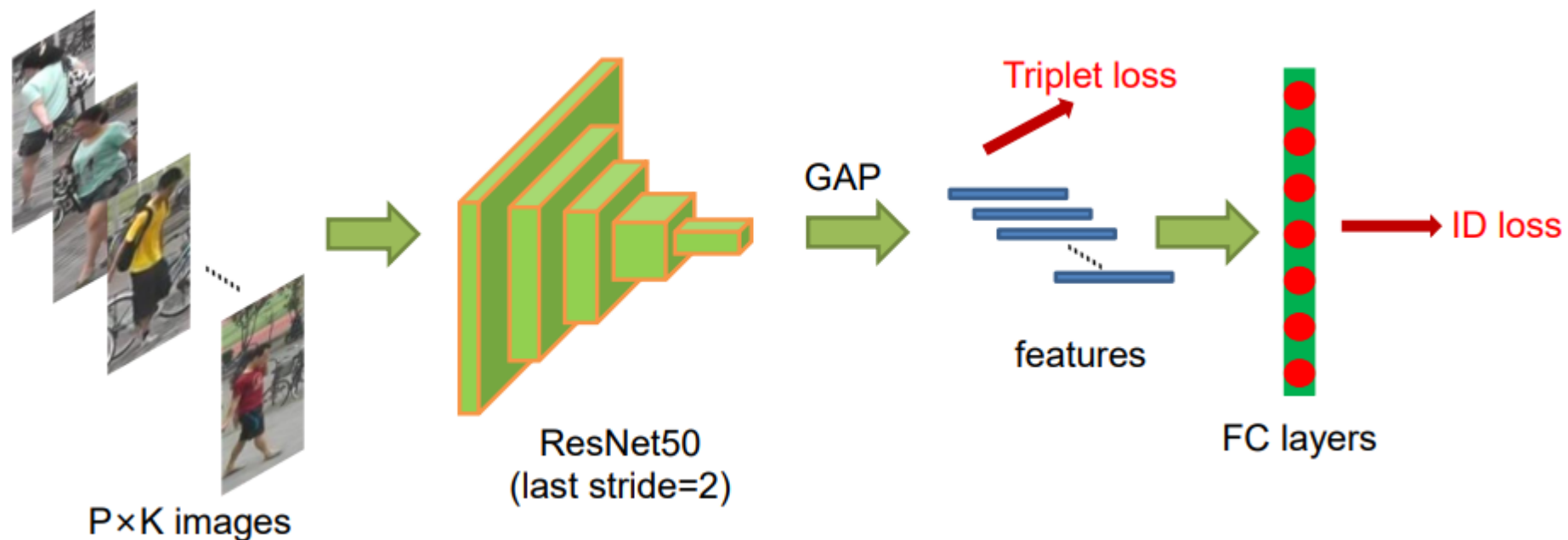


Fig. 2: Four different feature learning strategies. a) Global Feature, learning a global representation for each person image in § 2.1.1; b) Local Feature, learning part-aggregated local features in § 2.1.2; c) Auxiliary Feature, learning the feature representation using auxiliary information, *e.g.*, attributes [62], [63] in § 2.1.3 and d) Video Feature, learning the video representation using multiple image frames and temporal information [64], [65] in § 2.1.4.



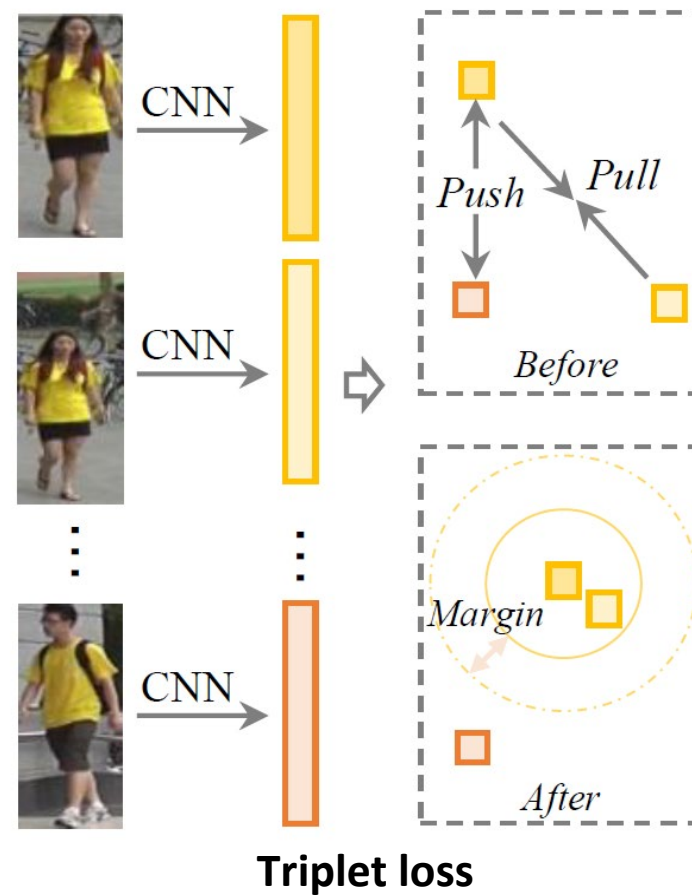
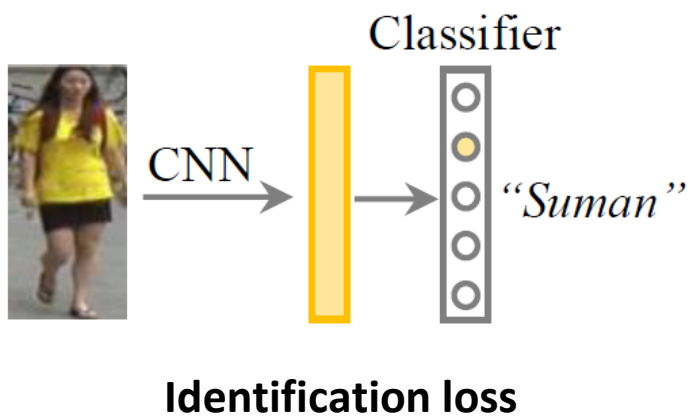
# Standard Approaches

- Recent approaches **utilize both identification and triplet loss.**



# Standard Approaches

- Recent approaches **utilize both identification and triplet loss.**

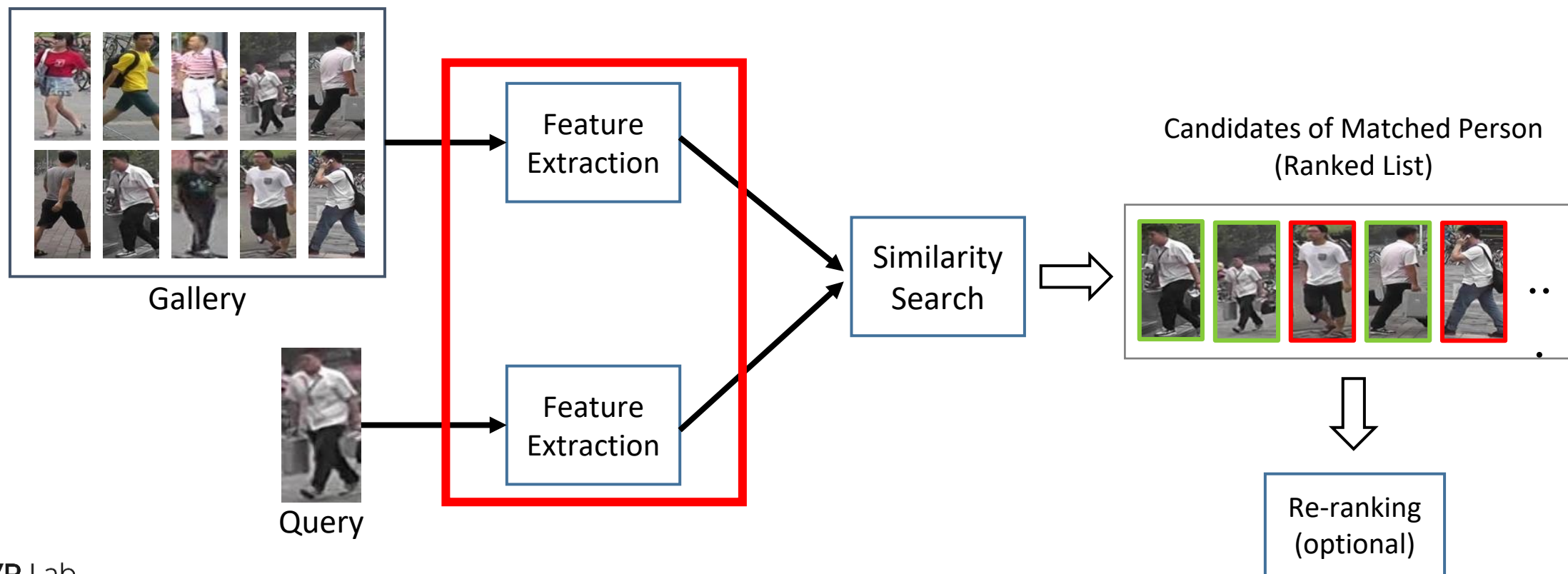


# Unsupervised Person Re-identification

Recent techniques for unsupervised approaches

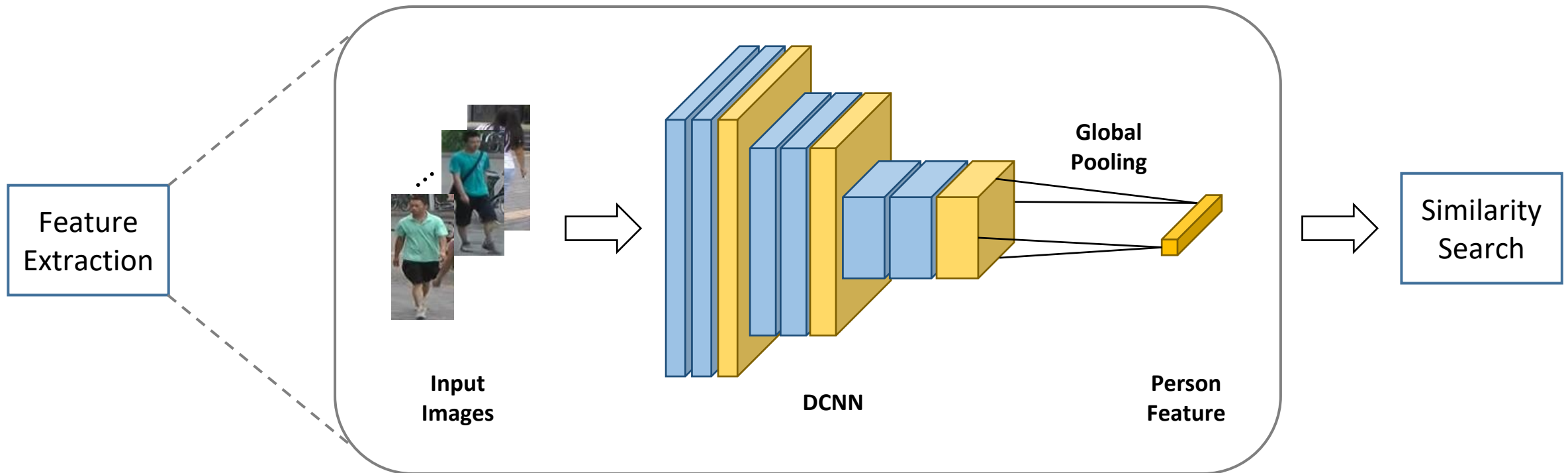
# General Protocol of Person Re-ID

- Person re-identification pipeline.



# General Protocol of Person Re-ID

- Deep convolutional neural network (DCNN) brings impressive improvements in person re-ID fields.

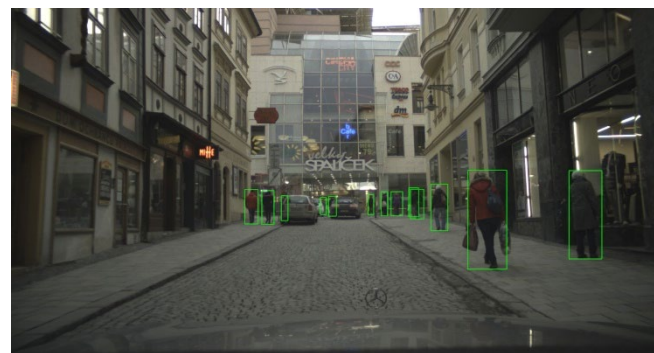


# Problems in DCNN

- Require many **training data with labels**.
- Challenges in identity annotation.
  - illumination changes.
  - Low resolution.
  - **Occlusions**.



Camera view



Person Detection

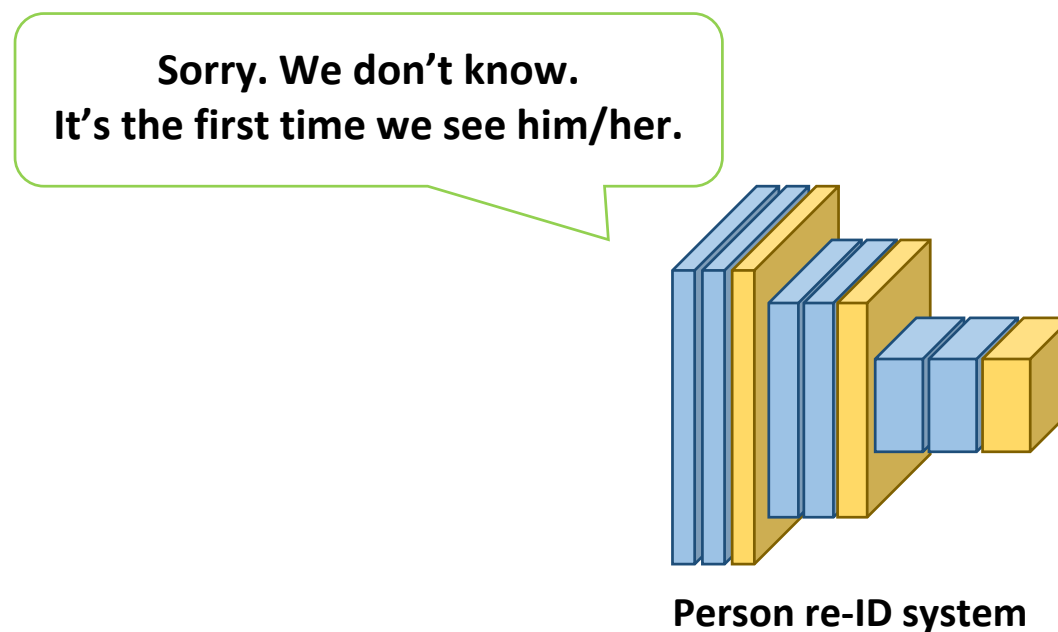


Identity Annotation

# Problems in DCNN

---

- The real-world scenario of person re-ID is an **open set** problem.
- New people (= new class) will appear from the camera views.



# Problem Setting

---

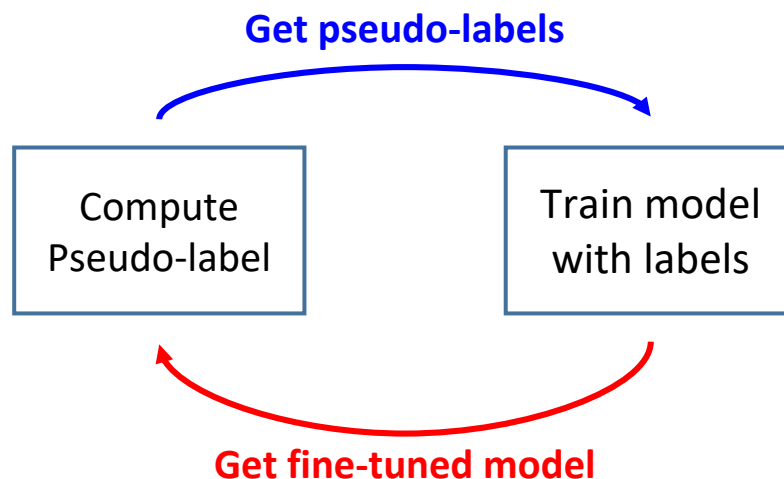
- Goal: Learn discriminative features for person retrieval **without ID labels**.
- Protocol: Training on **target** domain w/o labels → Testing on **target** domain.
- Challenges: Poor pseudo-supervision from unlabeled data.



# Pseudo Label-based Approaches

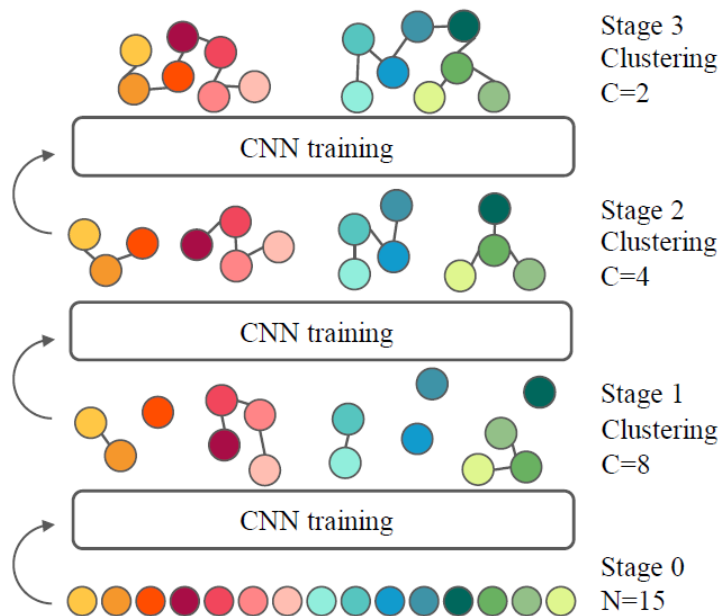
---

- Most recent studies **utilize pseudo-labels to train a re-ID model**.
  - K-nearest neighbor search; regard k-NN as the same class.
  - Clustering; regard each cluster as a class
- Clustering-based approaches have shown state-of-the-art results.



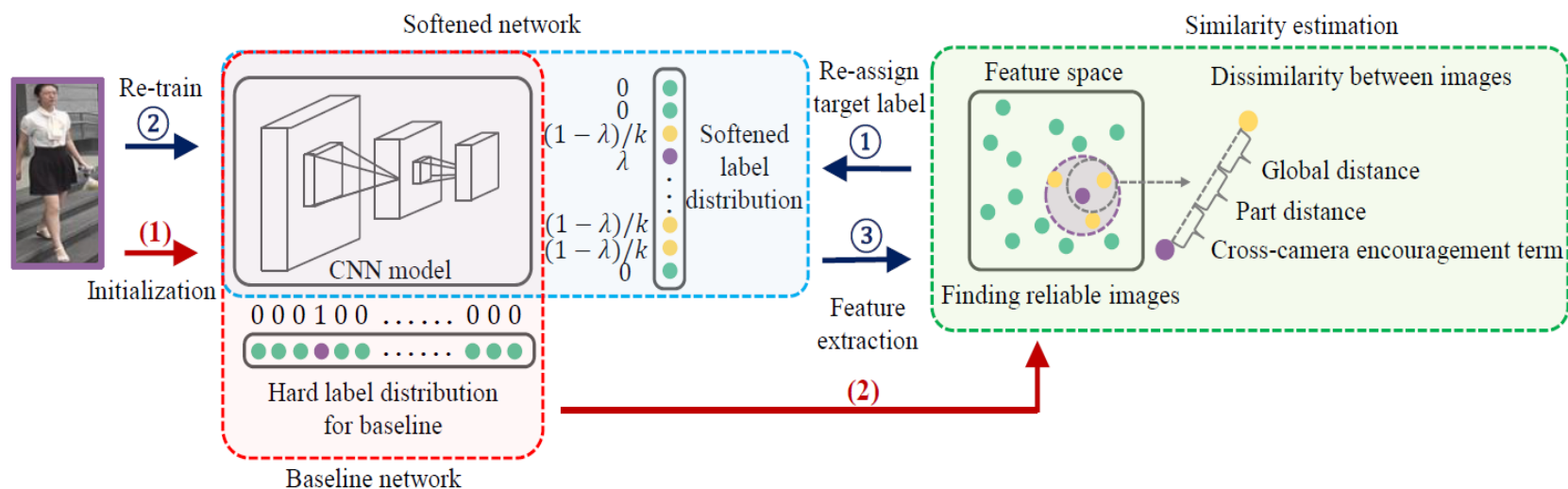
# Pseudo Label-based Approaches

- In early studies of this field focus on **how to obtain pseudo-labels**.
- Nowadays, most methods utilize **DBSCAN** clustering with **re-ranked distances** (performed by k-NN neighboring, etc.).



Bottom-Up Clustering (AAAI 19)

Different clustering

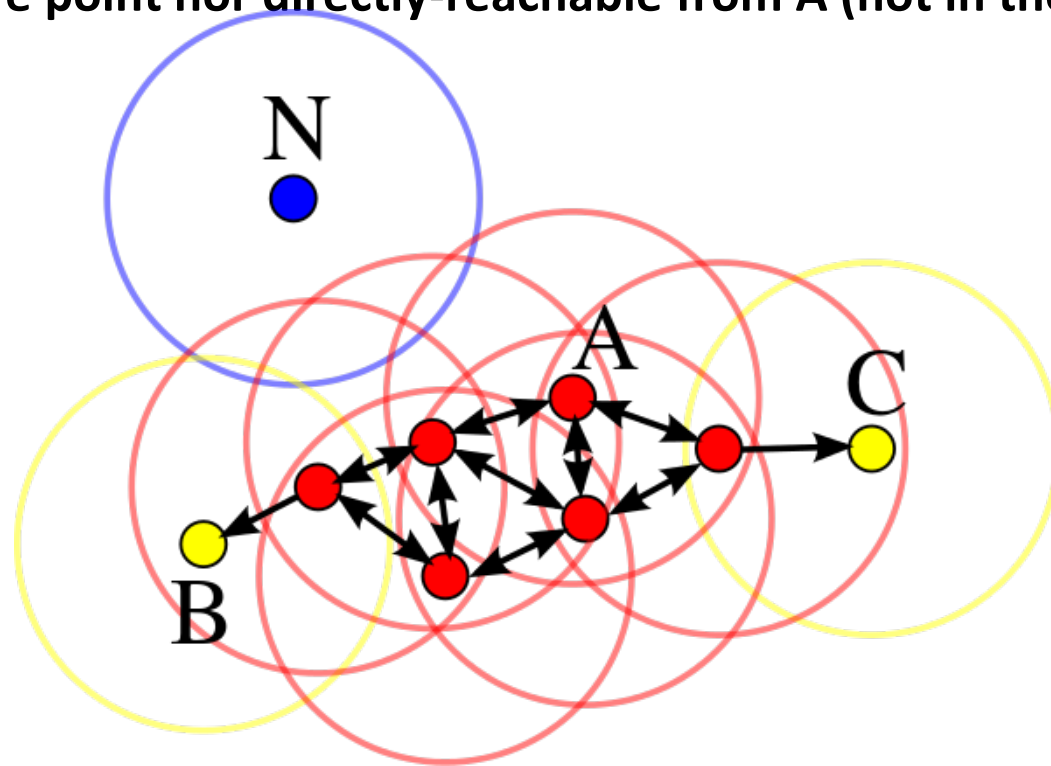


Softened Similarity Learning (CVPR 20)  
Different distance measure

Lin et al. A Bottom-Up Clustering Approach to Unsupervised Person Re-Identification. In AAAI 2019.  
Lin et al. Unsupervised Person Re-identification via Softened Similarity Learning. In CVPR 2020.

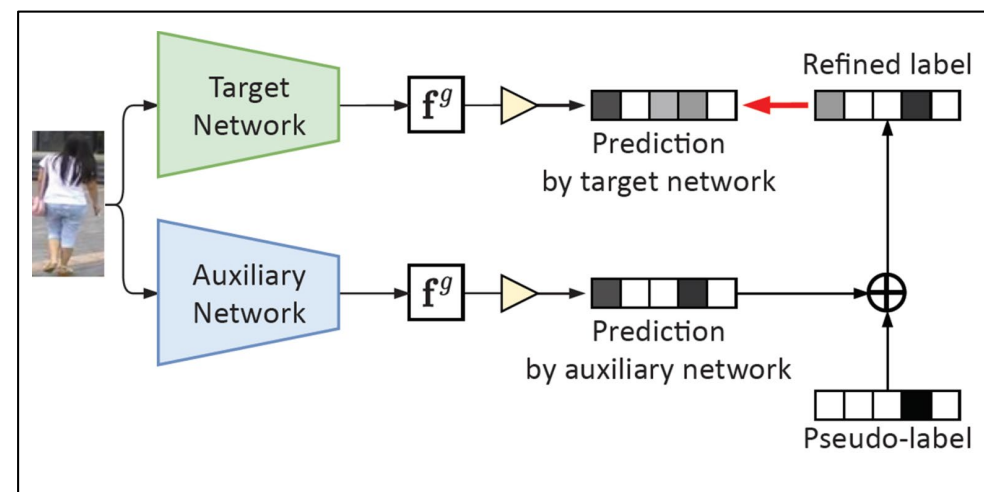
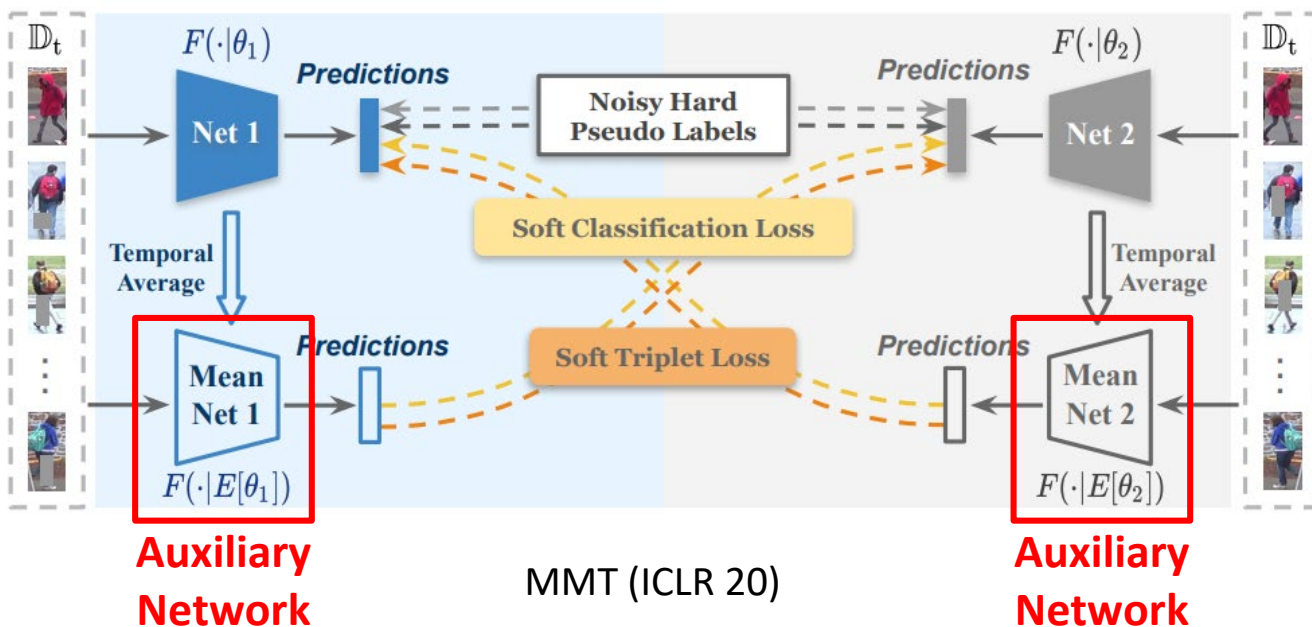
# Pseudo Label-based Approaches

- **DBSCAN (Density-Based Spatial Clustering of Applications with Noise)**
  - **Red**; core point (when  $\text{minPts} = 4$ ) that are closed to other nearby points
  - **Yellow**; not core points, but are reachable from A (belong to the same cluster)
  - **Blue**; neither a core point nor directly-reachable from A (not in the cluster)



# Pseudo Label Refinement

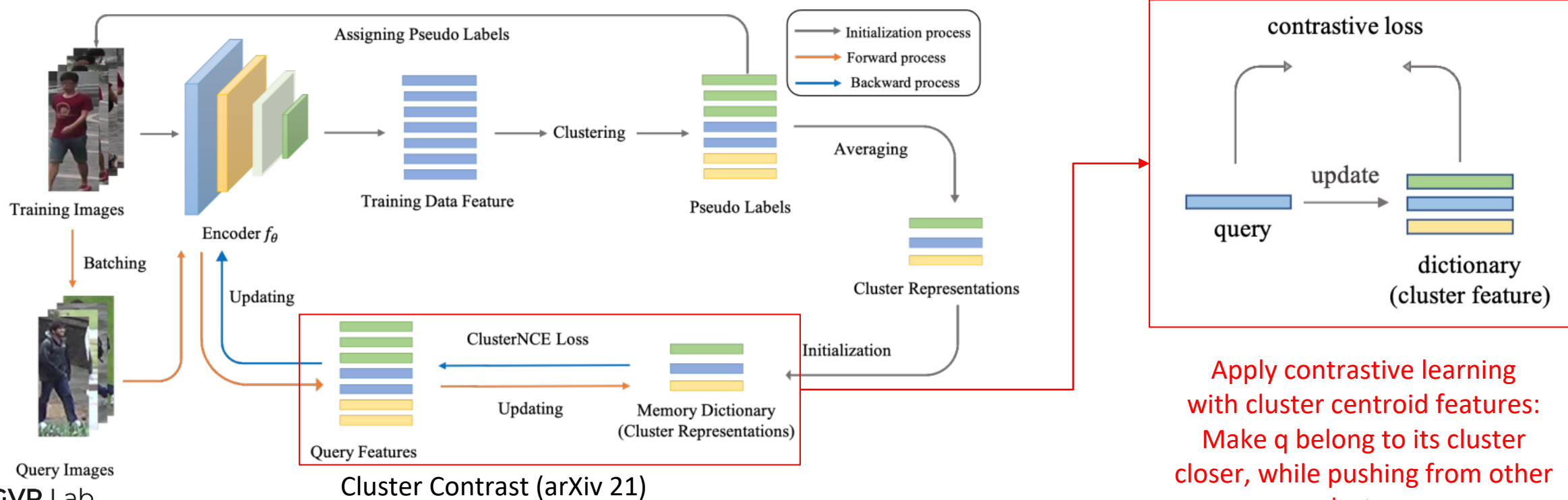
- Key idea: Re-ID performance  $\propto$  Quality of pseudo-labels.
- There are inevitable noises in pseudo-labels (noisy label problem), and some studies utilize predictions of an auxiliary network to refine labels (to avoid bias due to noise)



Label Refinement Process

# Cluster-based Contrastive Learning

- Key idea: Utilize clustering results for contrastive learning which is demonstrated its effectiveness in various unsupervised (self-supervised) tasks.
- Apply a contrastive learning in cluster-level



Apply contrastive learning with cluster centroid features: Make  $q$  belong to its cluster closer, while pushing from other clusters

# SimCLR

[Chen et al, A simple framework for contrastive learning of visual representations, ICML'20]

Maximizing the agreement of representations under data transformation, using a contrastive loss in the latent/feature space.

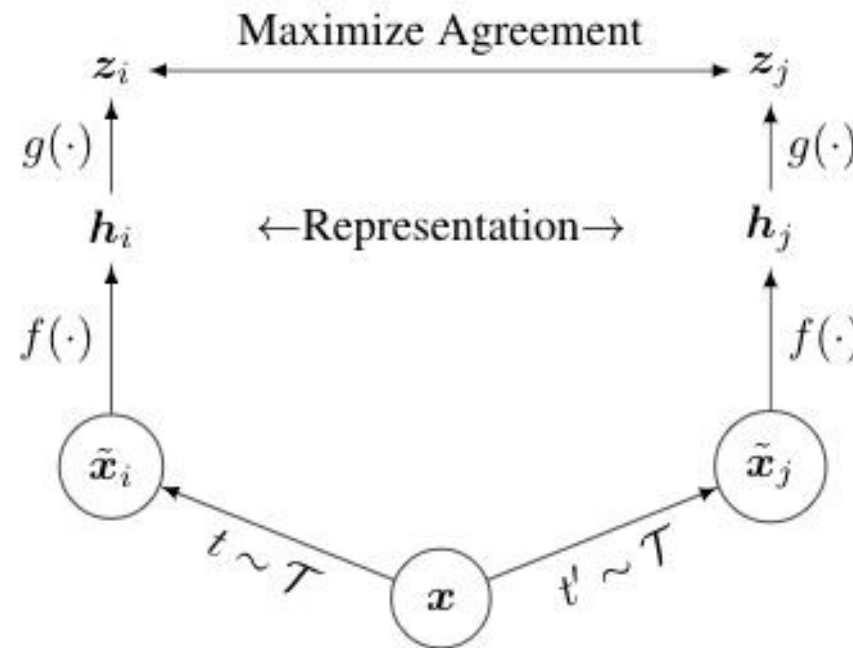


Figure 2. A framework for contrastive representation learning. Two separate stochastic data augmentations  $t, t' \sim \mathcal{T}$  are applied to each example to obtain two correlated views. A base encoder network  $f(\cdot)$  with a projection head  $g(\cdot)$  is trained to maximize agreement in *latent representations* via a contrastive loss.

# Semi-supervised learning

SimCLR as an example: strong semi-supervised learners, outperforms AlexNet with 100X fewer labels.

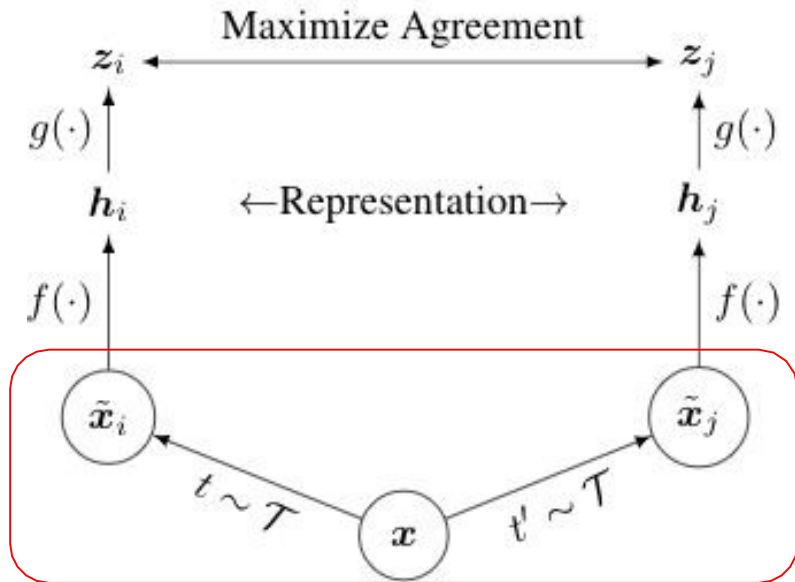
Method	Architecture	Label fraction	
		1%	10%
		Top 5	
<i>Methods using other label-propagation:</i>			
Pseudo-label	ResNet50	51.6	82.4
VAT+Entropy Min.	ResNet50	47.0	83.4
UDA (w. RandAug)	ResNet50	-	88.5
FixMatch (w. RandAug)	ResNet50	-	89.1
S4L (Rot+VAT+En. M.)	ResNet50 (4×)	-	91.2
<i>Methods using representation learning only:</i>			
InstDisc	ResNet50	39.2	77.4
BigBiGAN	RevNet-50 (4×)	55.2	78.8
PIRL	ResNet-50	57.2	83.8
CPC v2	ResNet-161(*)	77.9	91.2
Ours	ResNet-50	75.5	87.8
Ours	ResNet-50 (2×)	83.0	91.2
Ours	ResNet-50 (4×)	<b>85.8</b>	<b>92.6</b>

Table 7. ImageNet accuracy of models trained with few labels.

# SimCLR component: data augmentation

We use random crop and color distortion for augmentation.

Examples of augmentation applied to the left most images:

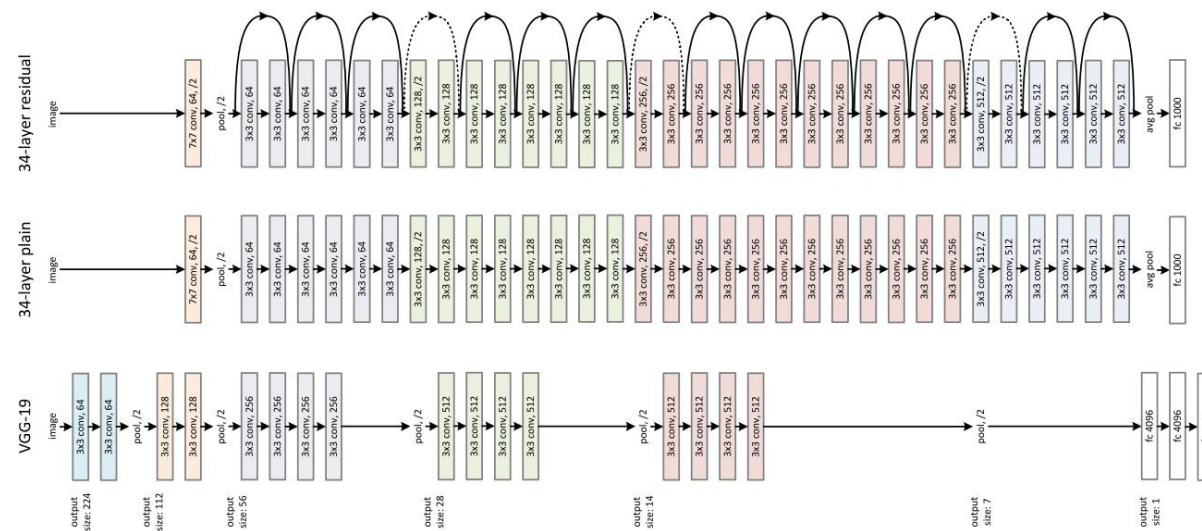
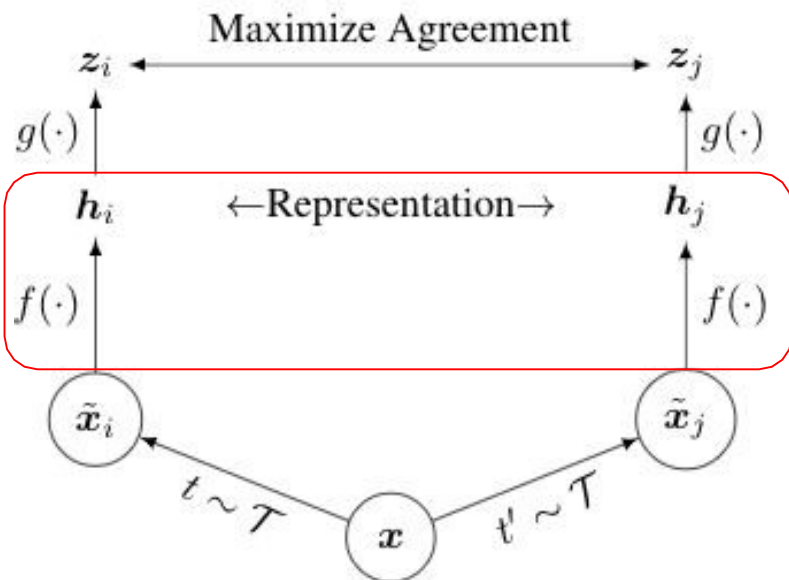




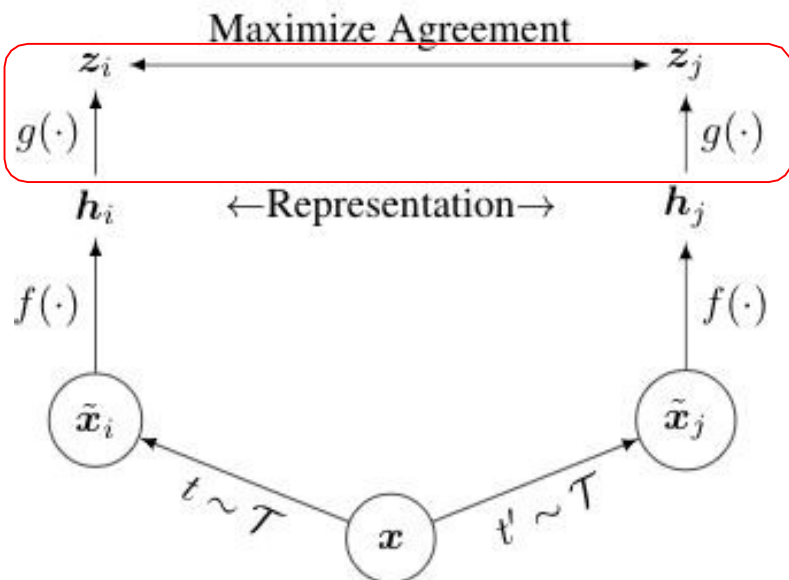
# SimCLR component: encoder

$f(x)$  is the base network that computes internal representation.

We can use (unconstrained) ResNet in this work. However, it can be other networks.

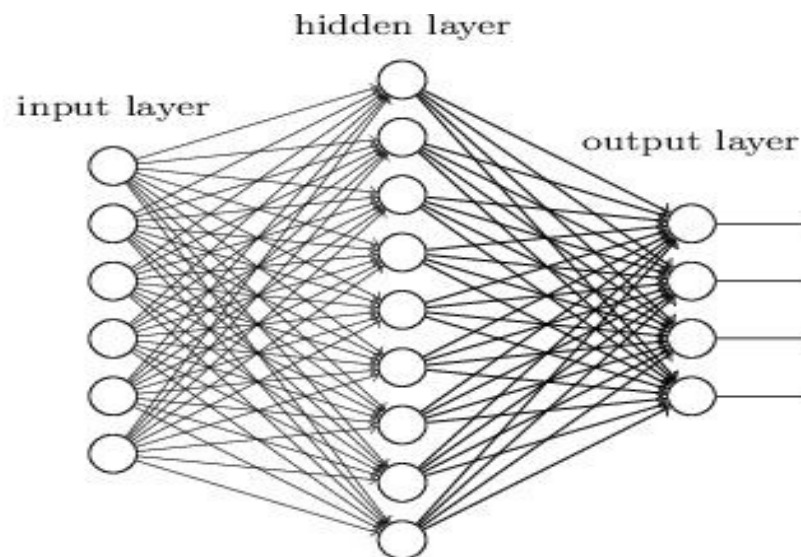


# SimCLR component: projection head



$g(h)$  is a projection network that project representation to a latent space.

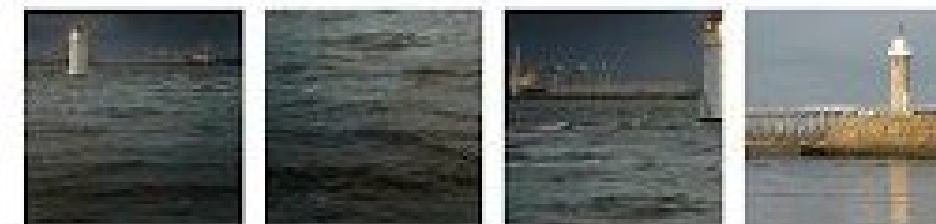
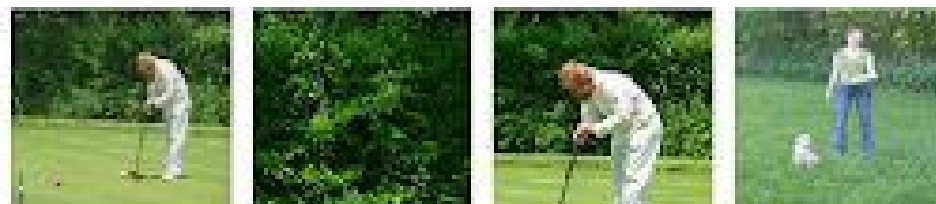
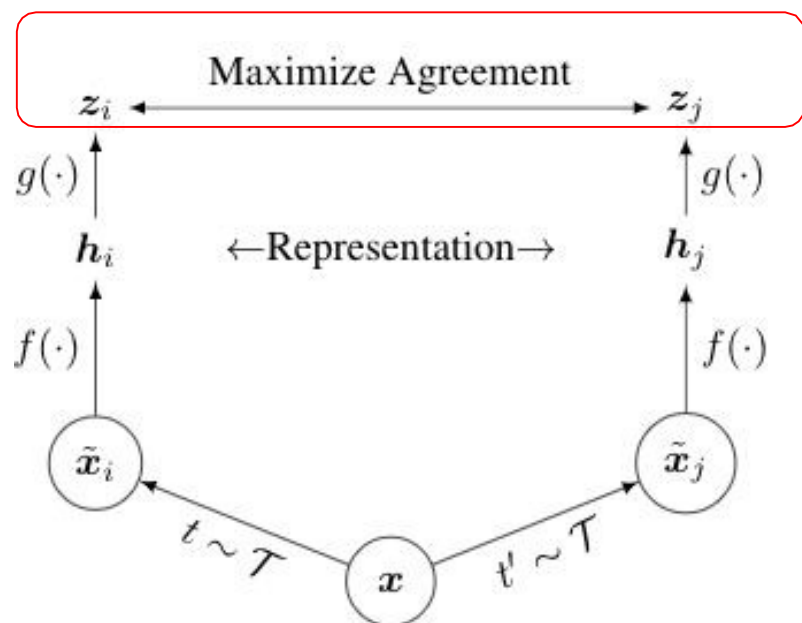
We use a MLP (with non-linearity).



# SimCLR component: contrastive loss

Maximize agreement using a contrastive task:

Given  $\{x_k\}$  where two different examples  $x_i$  and  $x_j$  are a positive pair, identify  $x_j$  in  $\{x_k\}_{k \neq i}$  for  $x_i$ .



Original image    crop 1    crop 2    contrastive image

Loss function:

$$\text{Let } \text{sim}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$$

$$l_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}$$

# Part-based Pseudo Label Refinement for Unsupervised Person Re-identification

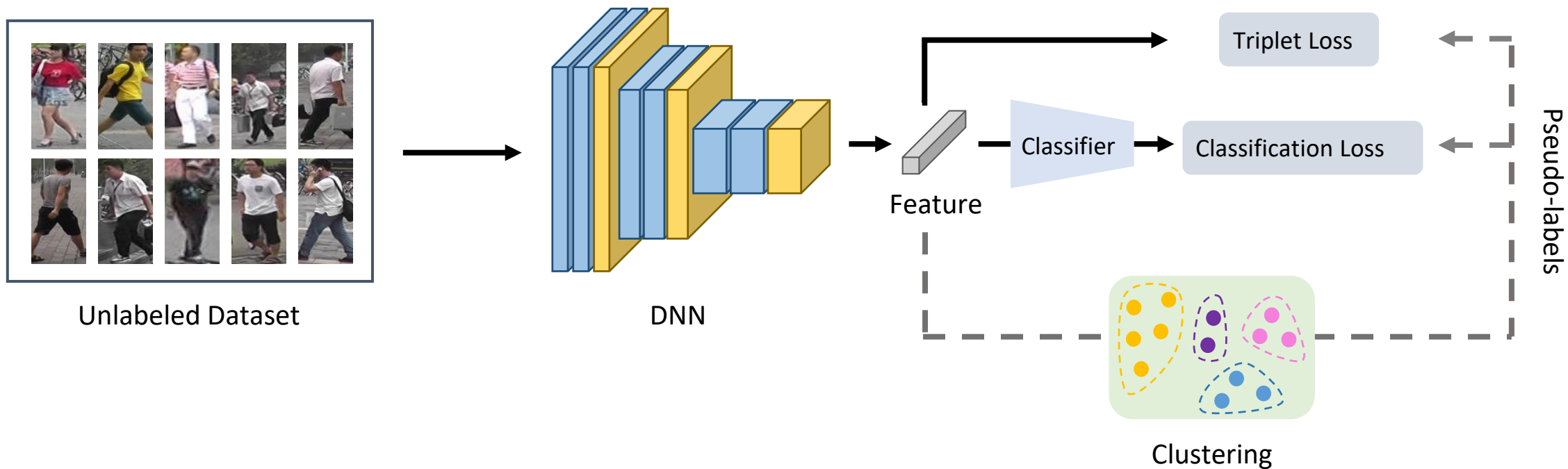
Yoonki Cho, Woo Jae Kim, Seunghoon Hong, Sung-Eui Yoon

KAIST

CVPR 2022

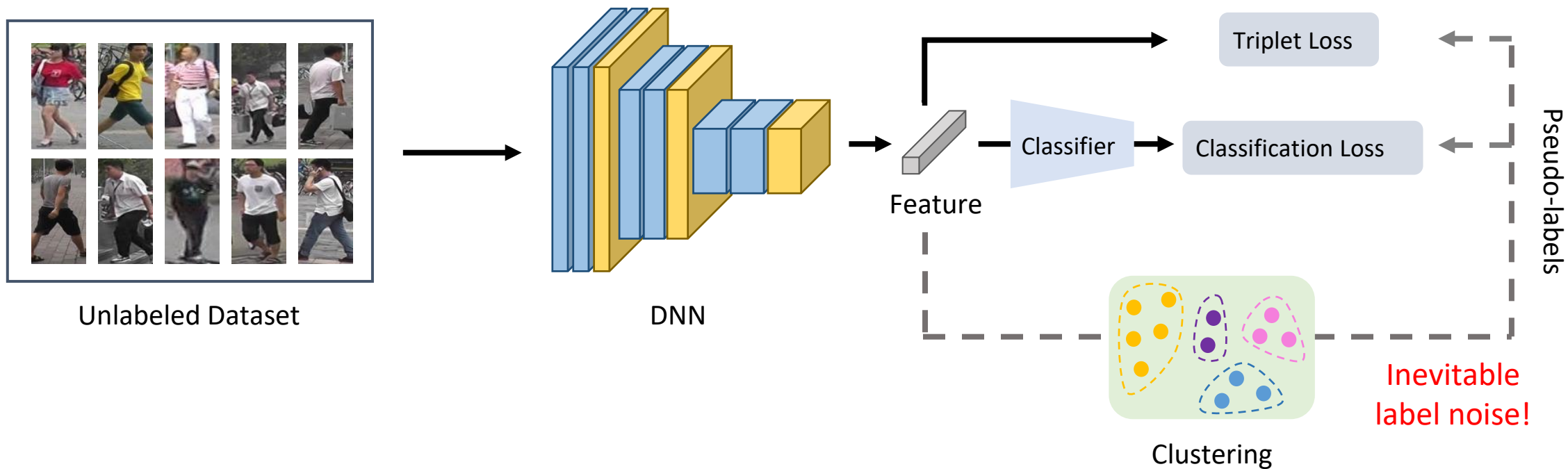
# Unsupervised Person Re-identification

- Learn the discriminative features for person re-ID from unlabeled data



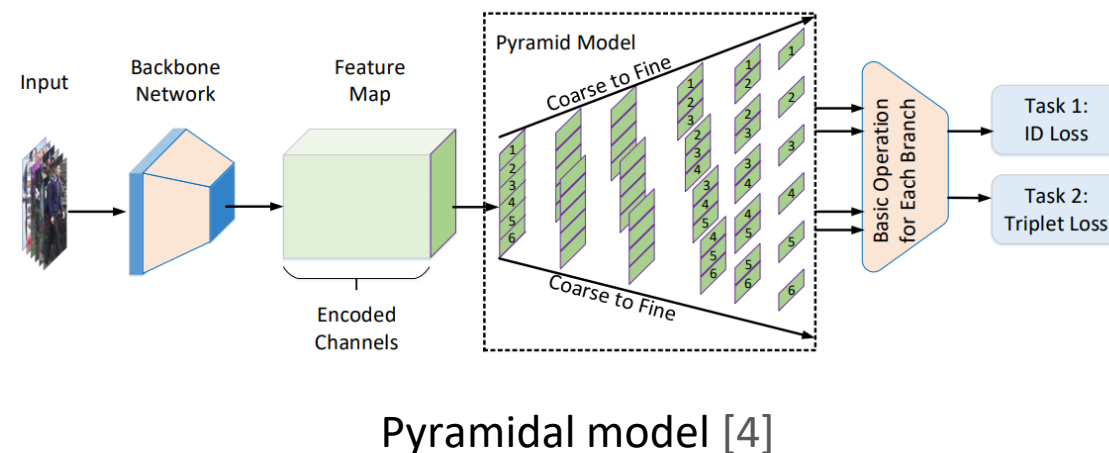
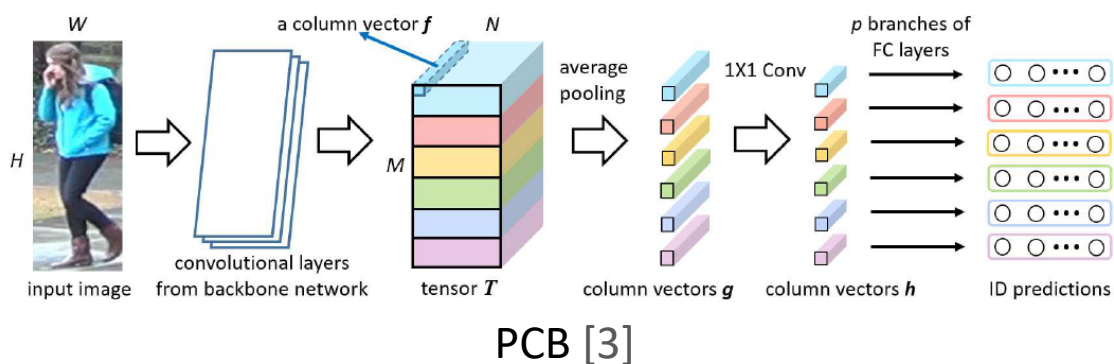
# Unsupervised Person Re-identification

- Learn the discriminative features for person re-ID from unlabeled data



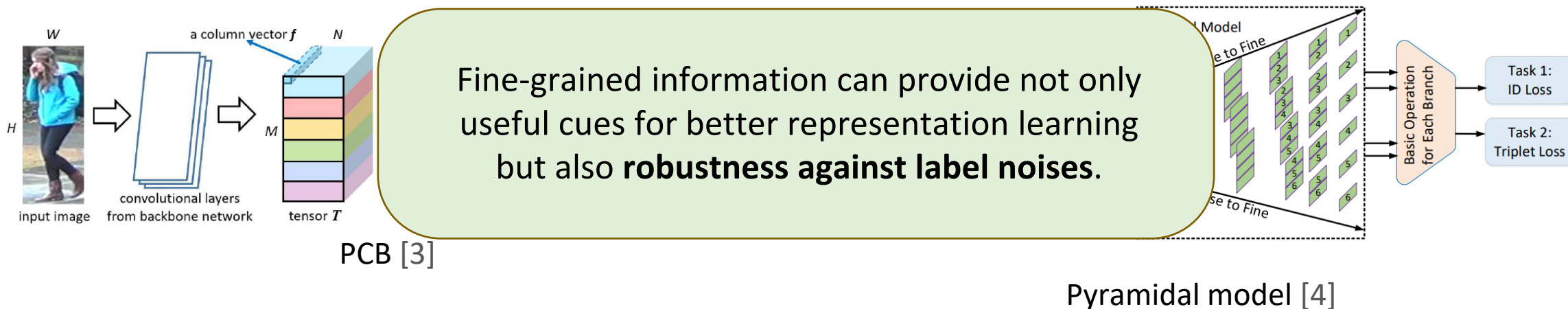
# Motivation and Idea

- Existing works neglect fine-grained information essential to person re-ID



# Motivation and Idea

- Existing works neglect fine-grained information essential to person re-ID



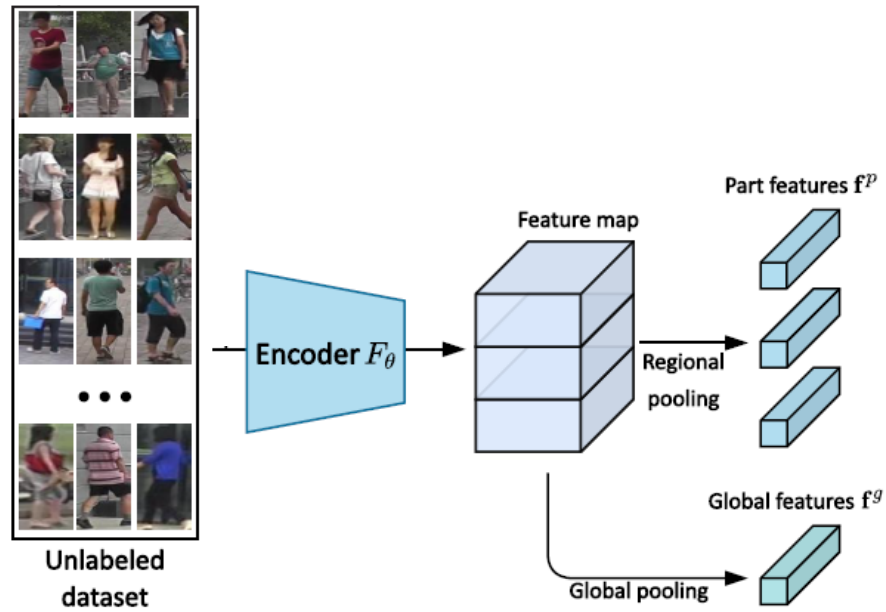
[3] Sun et al. Beyond Part Models: Person Retrieval with Refined Part Pooling (and a strong convolutional baseline). In ECCV 2018.

[4] Zheng et al. Pyramidal Person Re-Identification via Multi-Loss Dynamic Training. In CVPR 2019.



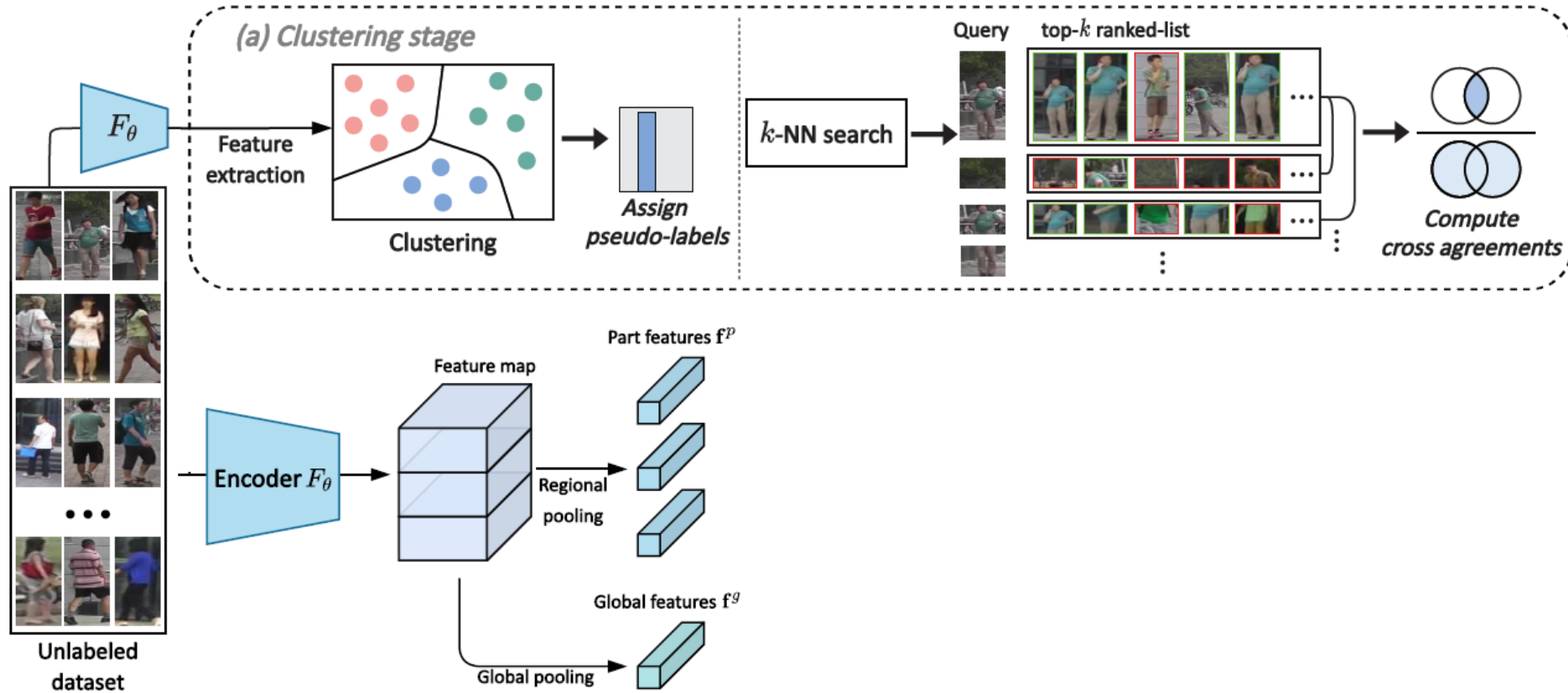
# Overview

- Part-based Pseudo Label Refinement (PPLR) framework



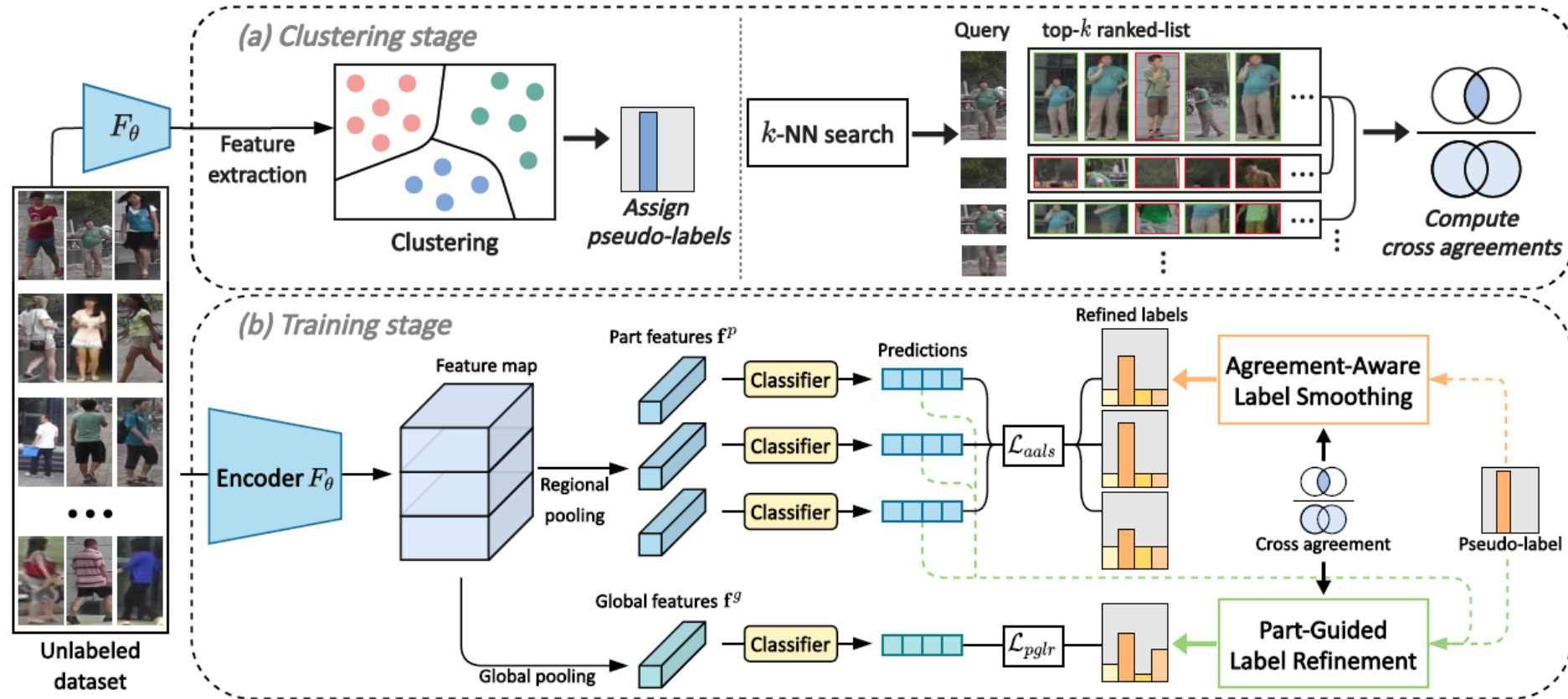
# Overview

- Part-based Pseudo Label Refinement (PPLR) framework



# Overview

- Part-based Pseudo Label Refinement (PPLR) framework



# Cross Agreement Score

---

- Global and part features in the same image can capture very different semantic information



Examples of ID-166 of Market-1501

# Cross Agreement Score

---

- Global and part features in the same image can capture very different semantic information



Examples of ID-166 of Market-1501

**Using complementary relationship naïvely can result in unreliable information!**

# Cross Agreement Score

---

- Cross agreement score  $C_i$  between global feature space  $g$  and part feature space  $p_n$  for the image  $x_i$ 
  - Jaccard similarity of nearest neighbors between global and part features

$$C_i(g, p_n) = \frac{|R_i(g, k) \cap R_i(p_n, k)|}{|R_i(g, k) \cup R_i(p_n, k)|} \in [0, 1]$$

$k$ -NN of global feature

$k$ -NN of  $n$ -th part feature

# Cross Agreement Score

---

- Cross agreement score  $C_i$  between global feature space  $g$  and part feature space  $p_n$  for the image  $x_i$ 
  - Jaccard similarity of nearest neighbors between global and part features

$$C_i(g, p_n) = \frac{|R_i(g, k) \cap R_i(p_n, k)|}{|R_i(g, k) \cup R_i(p_n, k)|} \in [0, 1]$$

$k$ -NN of global feature       $k$ -NN of  $n$ -th part feature

- $C_i(g, p_n) \uparrow$  :  $g$  and  $p_n$  are highly correlated around the data  $i$  = reliable
- $C_i(g, p_n) \downarrow$  :  $g$  and  $p_n$  are not correlated around the data  $i$  = unreliable

# Agreement-Aware Label Smoothing (AALS)

---

- Smooths pseudo-labels according to a cross agreement score of each part

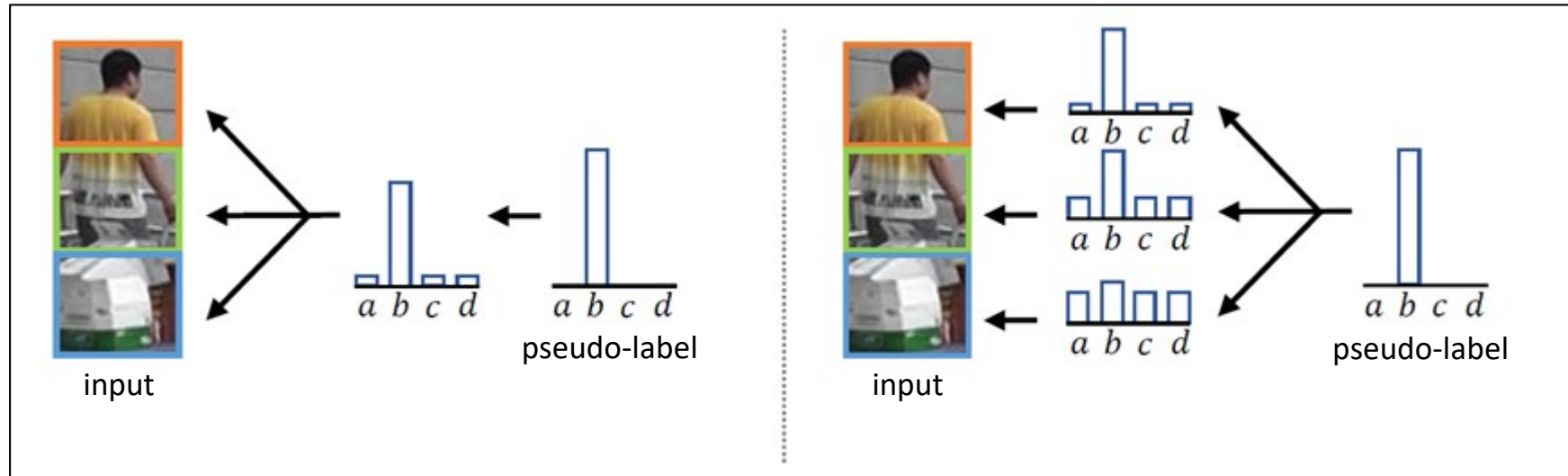
$$\tilde{y}_i^{p_n} = C_i(g, p_n) * \underbrace{y_i}_{\text{pseudo-label}} + (1 - C_i(g, p_n)) * \underbrace{u}_{\text{uniform distribution}}$$

- Cross agreement  $C_i \uparrow$  : prediction should be close to pseudo-label
- Cross agreement  $C_i \downarrow$  : prediction should be close to uniform distribution



# Agreement-Aware Label Smoothing (AALS)

- Calibrates the predictions of part features leading to reliable part feature learning



Vanilla label smoothing

AALS

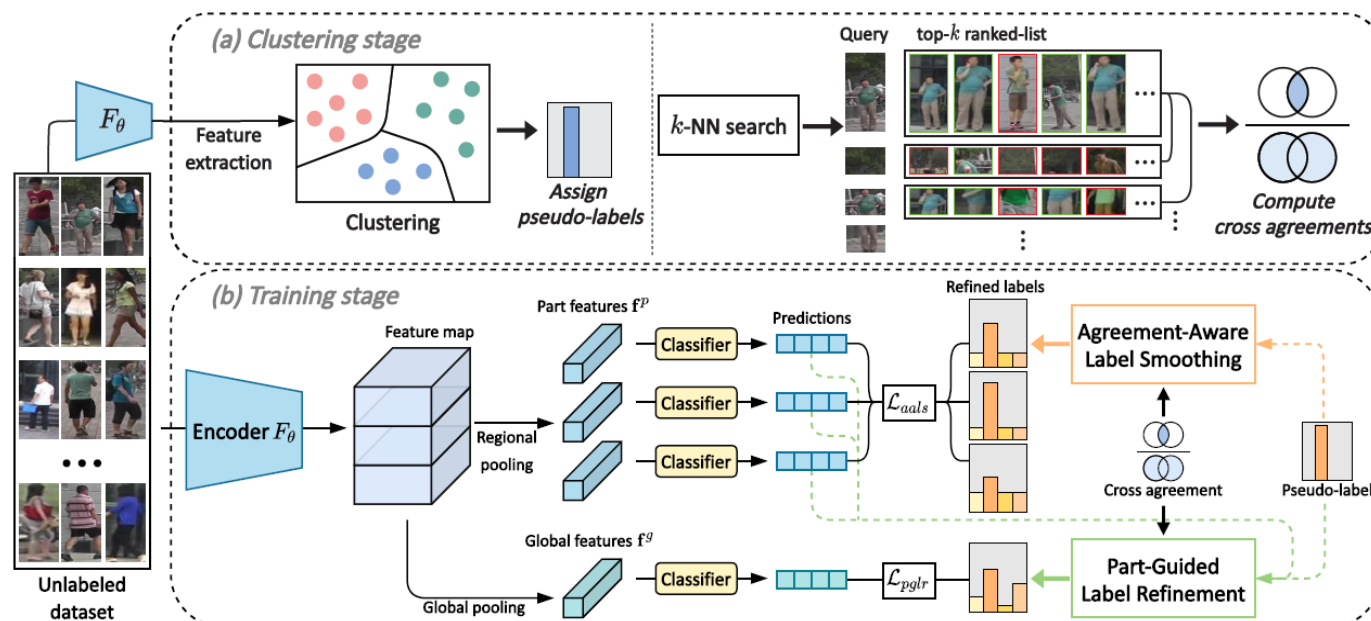
# Part-Guided Label Refinement (PGLR)

- Refines pseudo-labels by aggregating predictions of part features with different weights depending on each cross agreement score

$$\tilde{y}_i^g = \beta y_i + (1 - \beta) \sum_{n=1}^{N_p} \underbrace{w_i^{p_n}}_{\text{ensemble weight}} \underbrace{q_i^{p_n}}_{\text{prediction of } n\text{-th part feature}}, \text{ where } w_i^{p_n} = \frac{\exp(C_i(g, p_n))}{\sum_k \exp(C_i(g, p_k))}$$

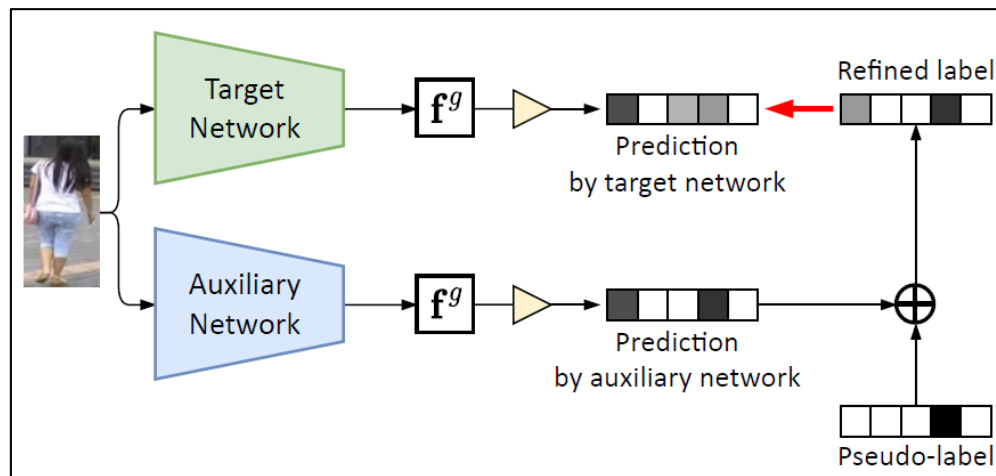
ensemble weight      prediction of  $n$ -th part feature

- $\beta$ : weighting parameter

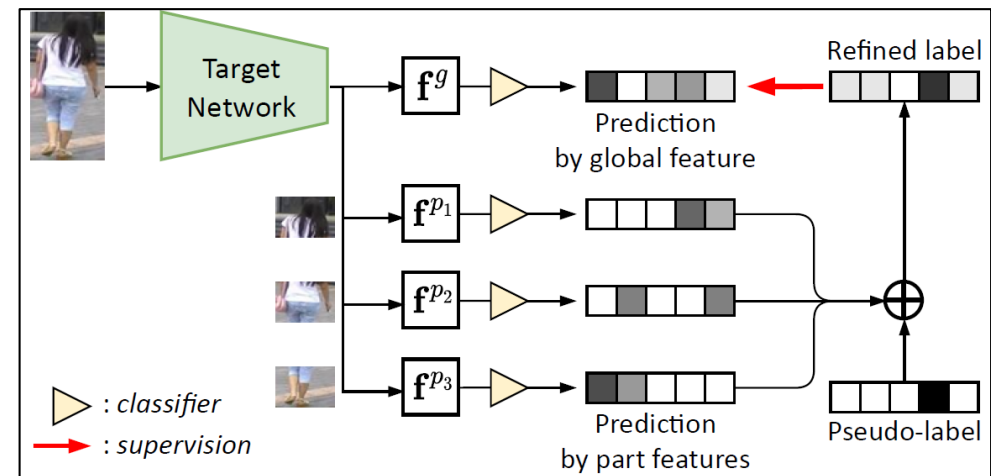


# Part-Guided Label Refinement (PGLR)

- Global features learn from the ensembled part predictions with rich fine-grained information without additional teacher networks



Previous



PGLR

# Experimental Results

---

- Ablation Study
  - Effectiveness of AALS and PGLR

Method		Market-1501		MSMT17	
AALS	PGLR	mAP	Rank-1	mAP	Rank-1
-	-	73.5	88.5	25.1	51.2

# Experimental Results

---

- Comparison with State-of-the-Arts

Method	Market-1501		MSMT17	
	mAP	Rank-1	mAP	Rank-1
SpCL (NeurIPS 20)	73.1	88.1	19.1	42.3
GCL (CVPR 21)	66.8	87.3	21.3	45.7
IICS (CVPR 21)	72.9	89.5	26.9	56.4
RLCC (CVPR 21)	77.7	90.8	27.9	31.4
ICE (ICCV 21)	79.5	92.0	29.8	59.0
<b>PPLR</b>	<b>81.5</b>	<b>92.8</b>	<b>31.4</b>	<b>61.1</b>

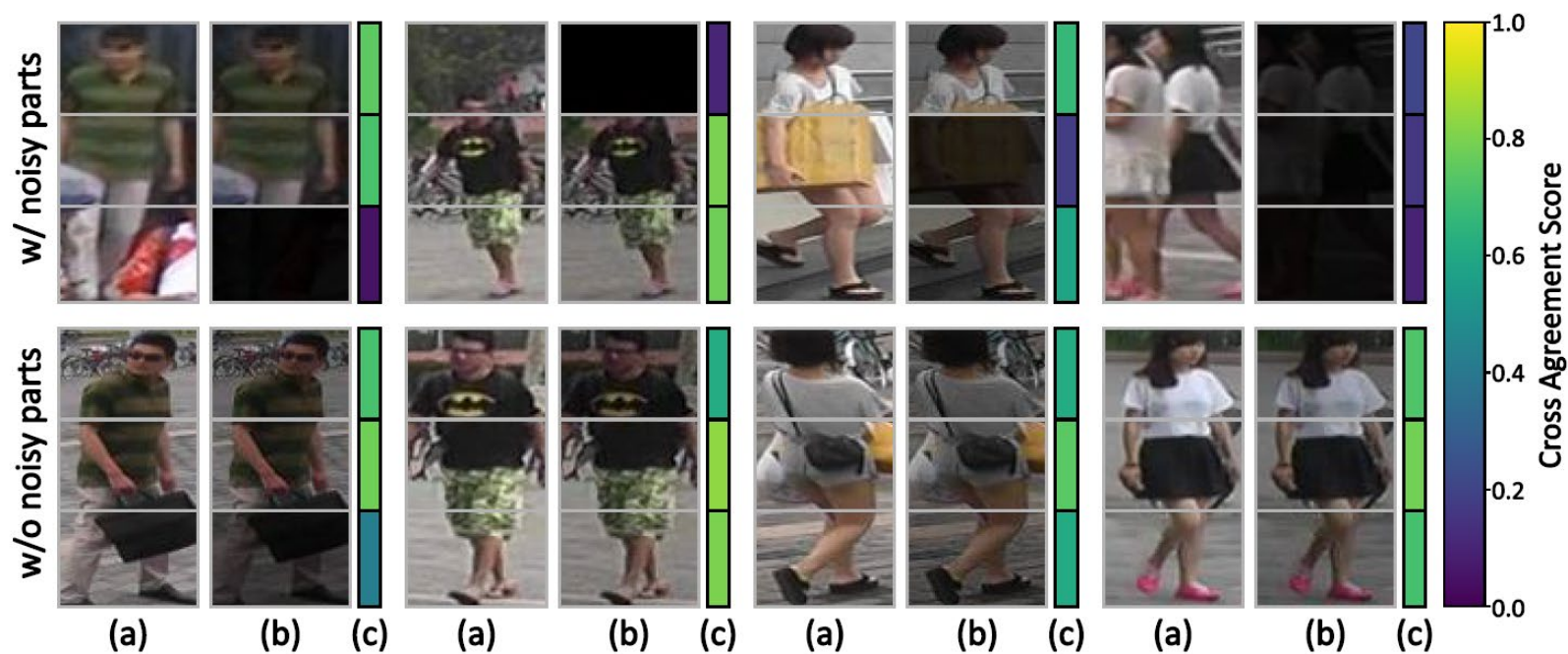
# Experimental Results

- Analysis of Cross Agreement Score

(a) Original image

(b) Soft masked image by the cross agreement score

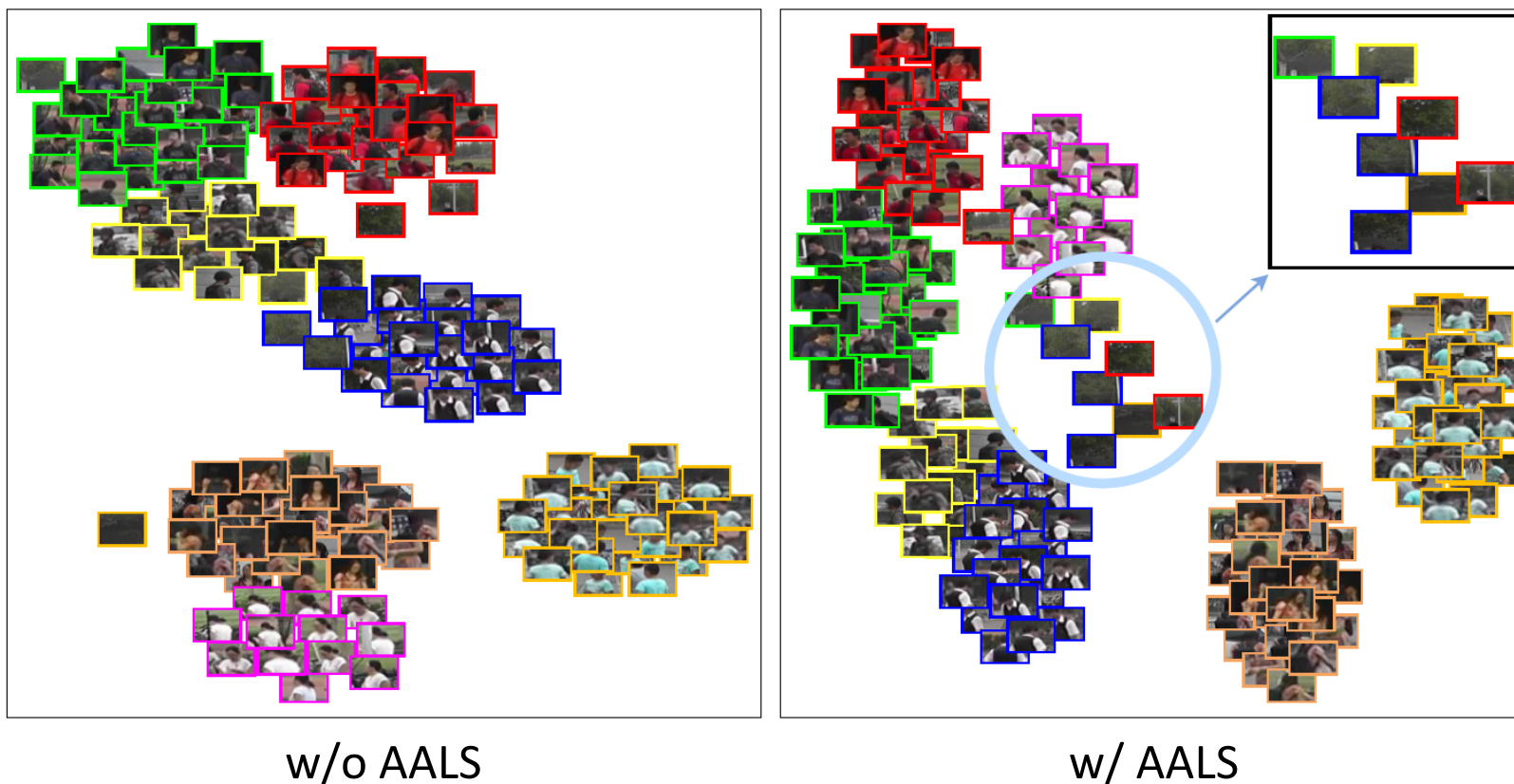
(c) Color jet bar of the cross agreement score



# Experimental Results

- Effect of Agreement-Aware Label Smoothing
  - t-SNE visualization of the topmost part feature space

Meaningless parts are not overfitted to meaningful clusters



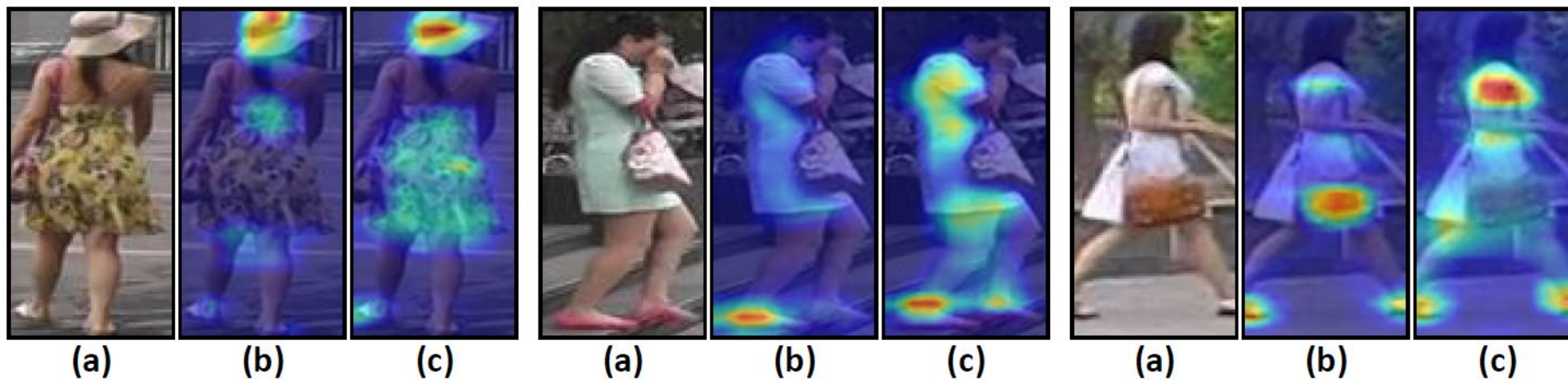
# Experimental Results

- Analysis of Part-Guided Label Refinement

(a) Original image

(b) Without PGLR

(c) With PGLR





# Summary

---

- Person Re-identification
- Unsupervised Approaches
- Part-based Pseudo Label Refinement for Unsupervised Person Re-ID (CVPR 2022)

Project Page: <https://sgvr.kaist.ac.kr/~yoonki/PPLR/>

