



# Convolutional Feature Transfer via Camera-Specific Discriminative Pooling for Person Re-Identification

Tetsu Matsukawa   Einoshin Suzuki

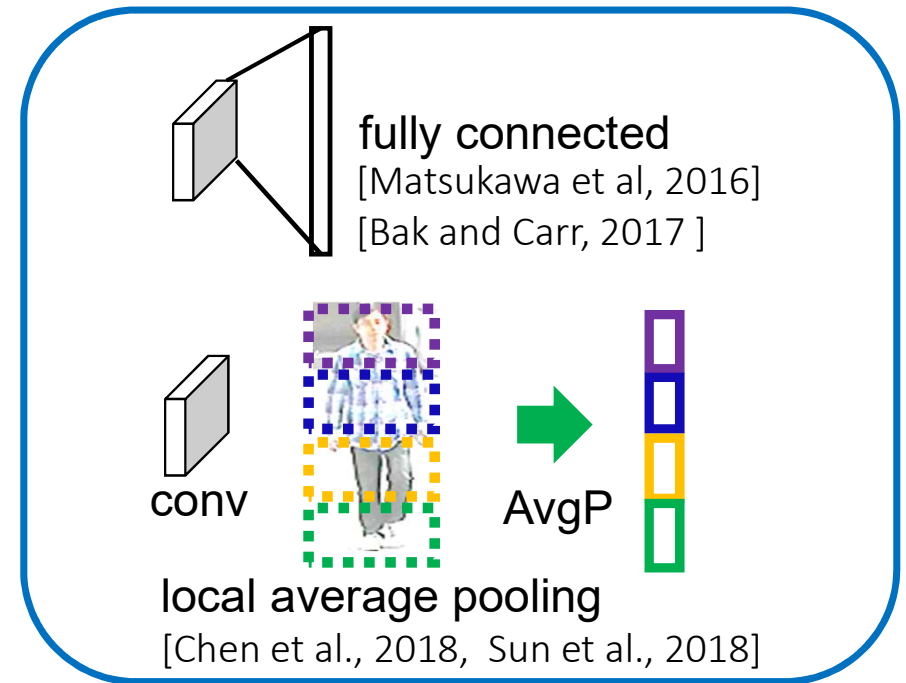
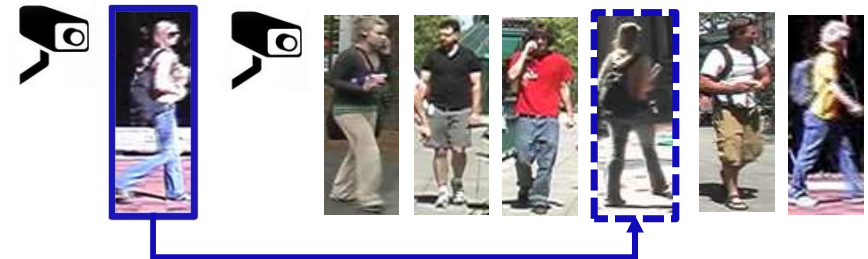
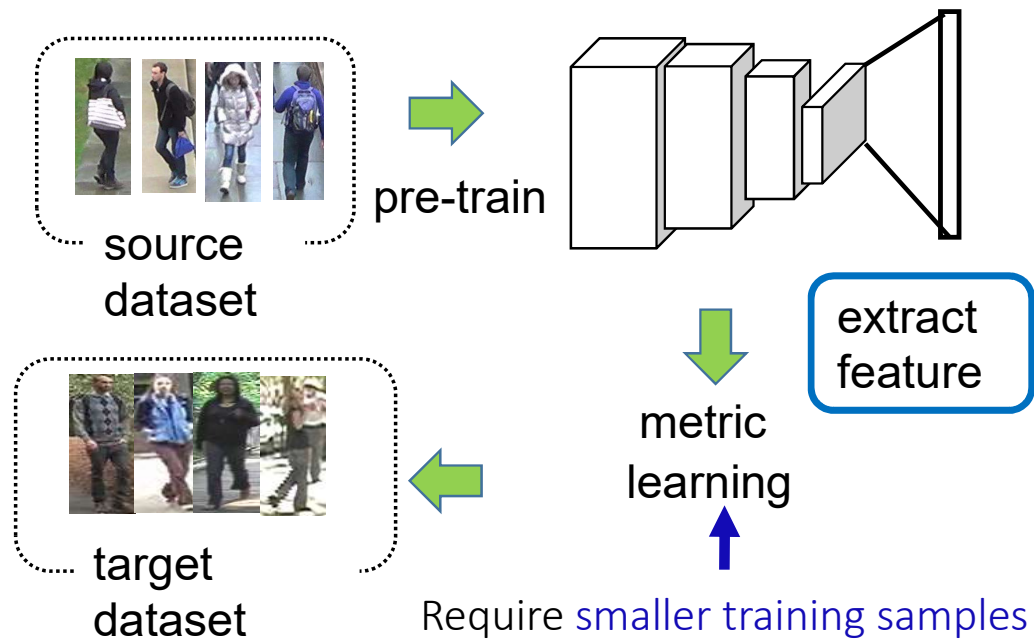
Kyushu University, JAPAN

# Background

Practical issue

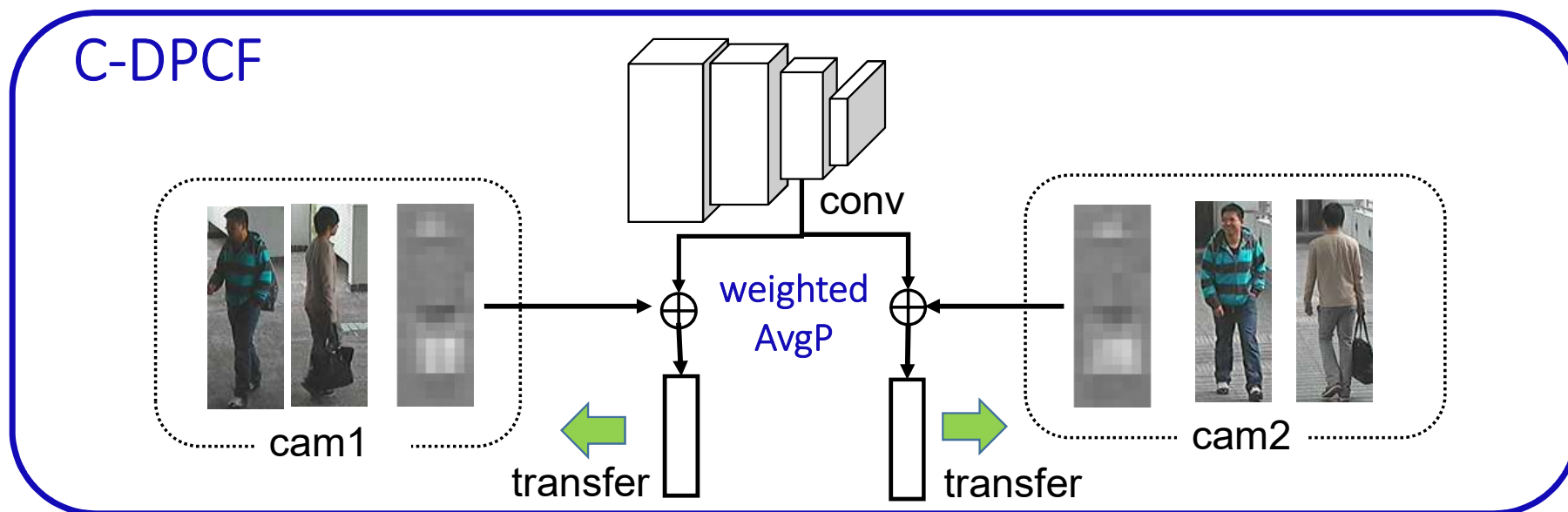
- **Lack of training person IDs**
- Require GPU for fine-tuning

CNN-feature transfer

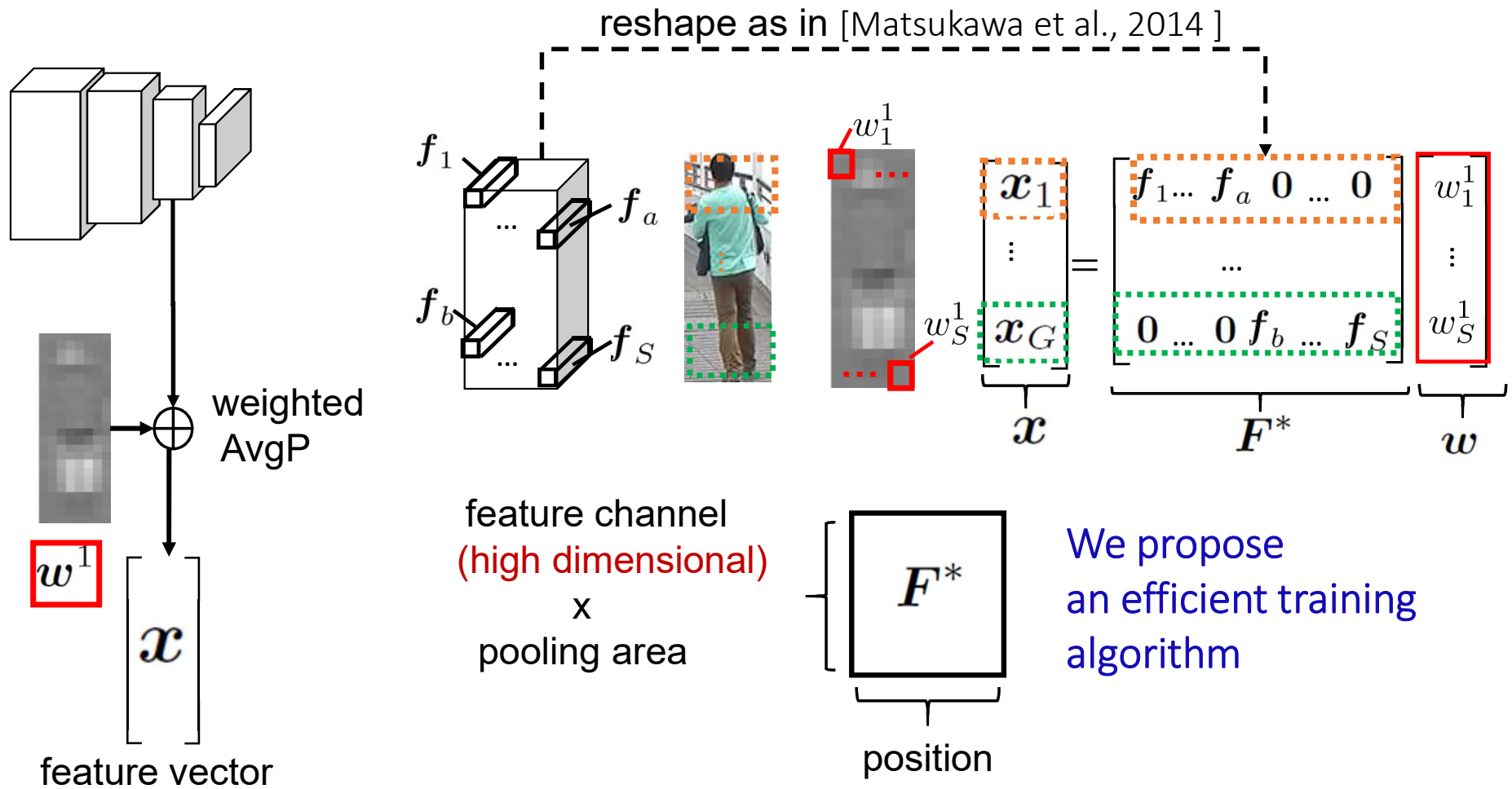


# Approach

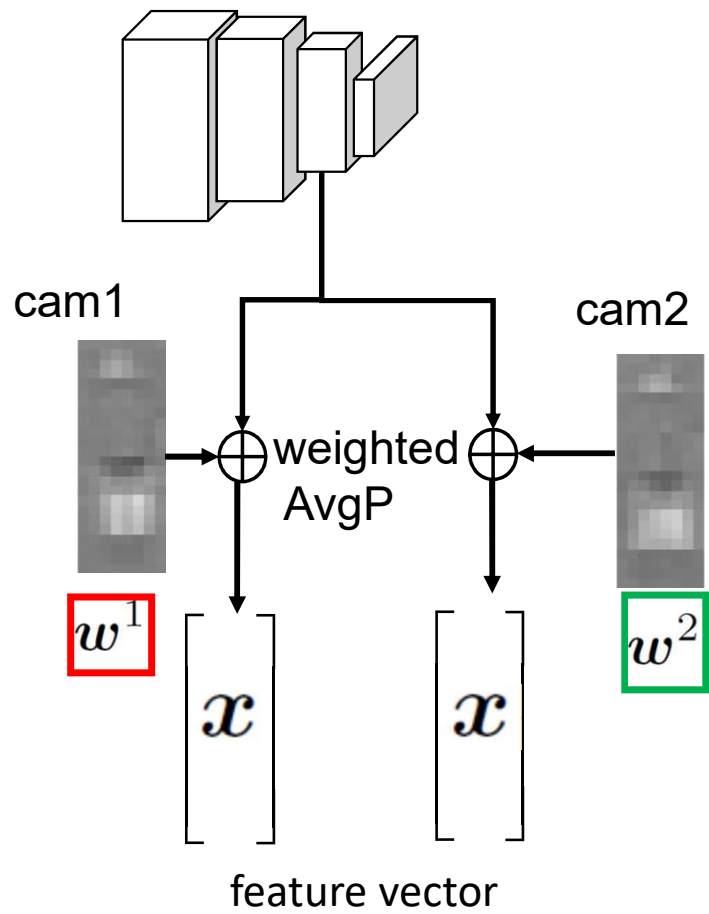
Existing features are less transferable to different camera/datasets due to spatial bias



# Weighted Local AvgP



# Camera-specific pooling



Zero-padding to feature matrix

$$\begin{array}{l}
 \text{cam1} \\
 \left[ \mathcal{X} \right] = \left[ \begin{array}{c|c} F^* & \mathbf{0}_{D \times S} \\ \hline & \end{array} \right] \begin{array}{c} \boxed{w^1} \\ \boxed{w^2} \end{array} \\
 \underbrace{\hspace{10em}}_F \quad \underbrace{\hspace{2em}}_w \\
 \text{cam2} \\
 \left[ \mathcal{X} \right] = \left[ \begin{array}{c|c} \mathbf{0}_{D \times S} & F^* \\ \hline & \end{array} \right] \begin{array}{c} \boxed{w^1} \\ \boxed{w^2} \end{array}
 \end{array}$$

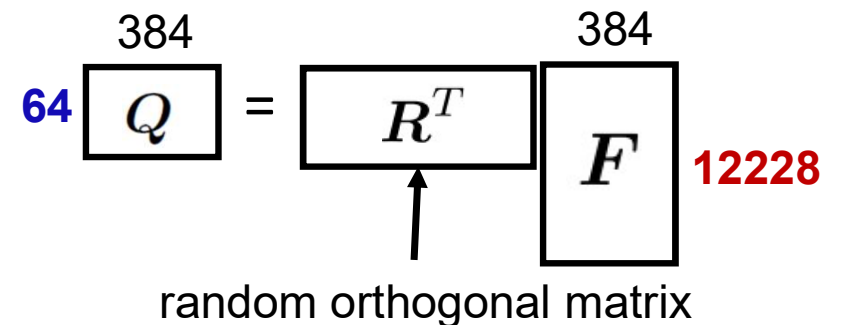
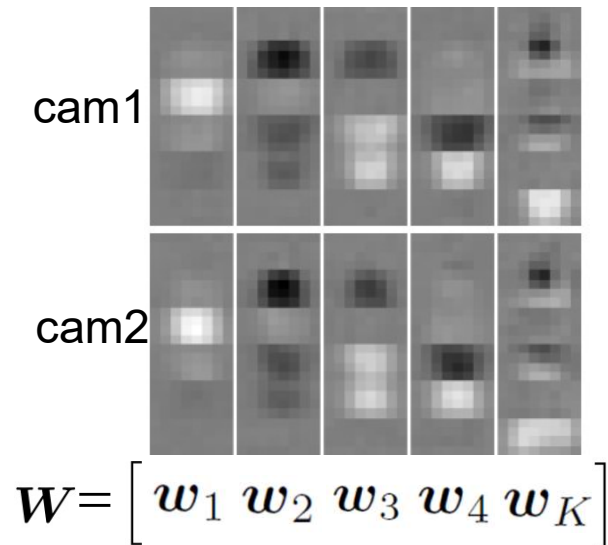
# Weight map learning: Problem formulation

- Given a training data  $\{\mathbf{F}_i, p_i, c_i\}_{i=1}^N$ 
  - Person ID
  - Cam ID
- Optimize sum of distances of  $K$ -weight map pairs

$$\delta_{\mathbf{W}}^2(i, j) = \sum_{k=1}^K \delta_{\mathbf{w}_k}^2(i, j)$$

➤ Random projection distance

$$\begin{aligned} \delta_{\mathbf{w}_k}^2(i, j) &= \|\mathbf{R}^T \mathbf{x}_{k,i} - \mathbf{R}^T \mathbf{x}_{k,j}\|_2^2 \\ &= \|\mathbf{R}^T \mathbf{F}_i \mathbf{w}_k - \mathbf{R}^T \mathbf{F}_j \mathbf{w}_k\|_2^2 \\ &= \|\mathbf{Q}_i \mathbf{w}_k - \mathbf{Q}_j \mathbf{w}_k\|_2^2 \end{aligned}$$



# Weight map learning: Optimization

- Maximum margin with orthogonal constraint

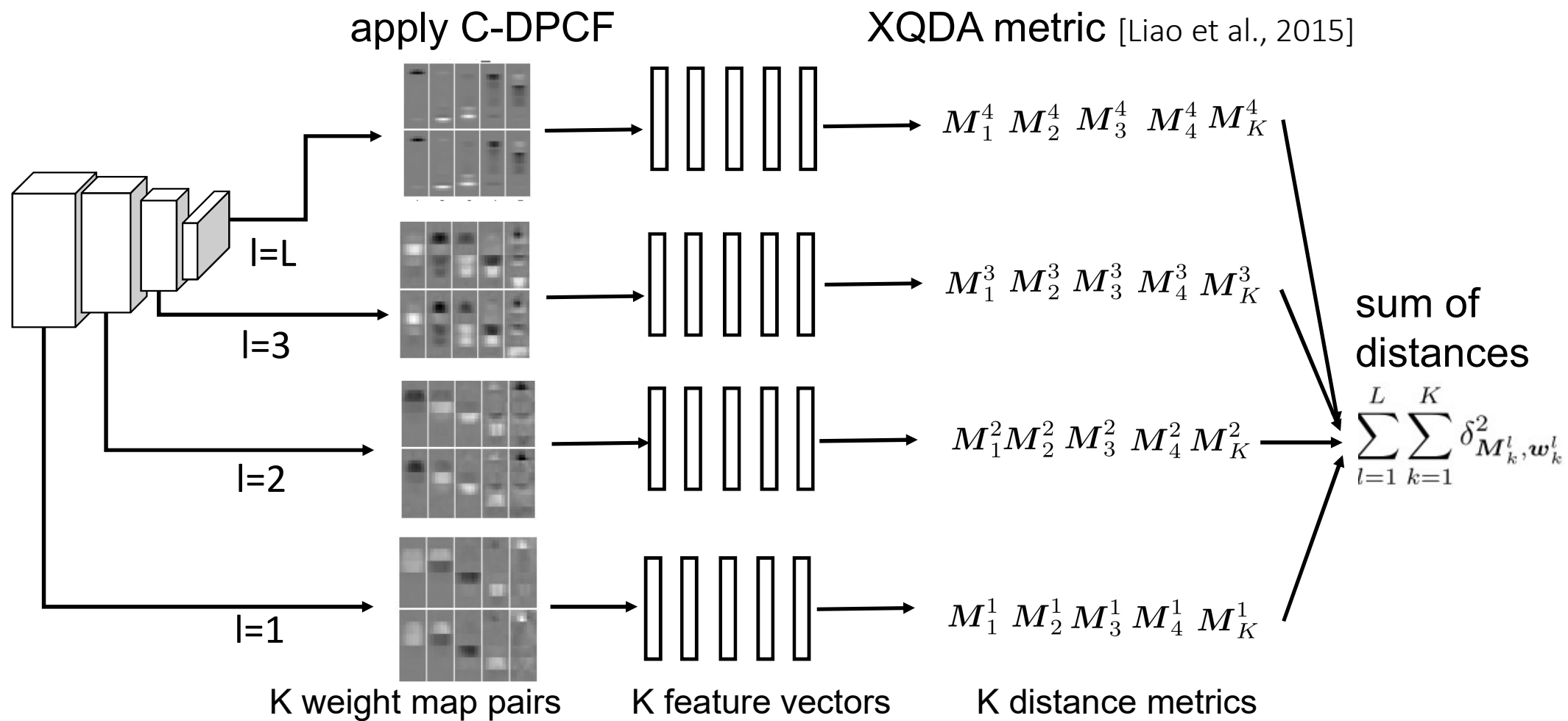
$$\max_{\mathbf{W}} \quad J(\mathbf{W}) = \underbrace{\text{Tr} \left[ \mathbf{W}^T \boldsymbol{\Sigma}_{\mathcal{D}} \mathbf{W} \right]}_{\text{avg. distance of different persons}} - \underbrace{\text{Tr} \left[ \mathbf{W}^T \boldsymbol{\Sigma}_{\mathcal{S}} \mathbf{W} \right]}_{\text{avg. distance of same person}}$$

$$s.t. \quad \mathbf{W}^T \mathbf{W} = \mathbf{I}_K$$

$$\text{where} \quad \boldsymbol{\Sigma}_{\mathcal{D}} = \frac{1}{N_{\mathcal{D}}} \sum_{(i,j) \in \mathcal{D}} (\mathbf{Q}_i - \mathbf{Q}_j)^T (\mathbf{Q}_i - \mathbf{Q}_j)$$
$$\boldsymbol{\Sigma}_{\mathcal{S}} = \frac{1}{N_{\mathcal{S}}} \sum_{(i,j) \in \mathcal{S}} (\mathbf{Q}_i - \mathbf{Q}_j)^T (\mathbf{Q}_i - \mathbf{Q}_j)$$

The solution is given by **eigen decomposition** of  $\boldsymbol{\Sigma}_{\mathcal{D}} - \boldsymbol{\Sigma}_{\mathcal{S}}$

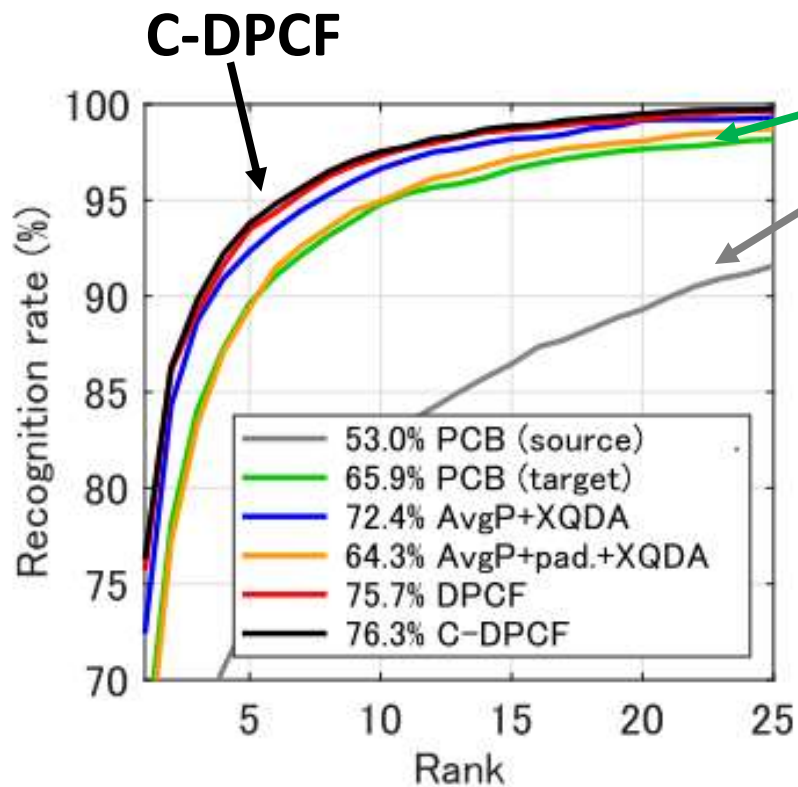
# Distance for re-id





# Comparison

Baseline: PCB [Sun et al., 2018] trained on (source/target) dataset



Source dataset: MSMT17

Target dataset: VIPeR

also compared on other 3 target datasets

Training time

C-DPCF	PCB(target)
42.5 sec (CPU)	312 sec (GPU)

# SOTA

S: Supervised  
U: Unsupervised  
DG: Domain Generalization

Rank-1 rates

		Type	VIPeR	GRID	PRID	CUHK01	
		CMDL[PAMI18]	S	66.4	30.9	52.0	78.2
		HGD[PAMI20]	S	52.8	28.2	-	-
CNN features	{	Synthesis[ECCV18]	U	43.0	-	-	54.9
		One-shot [CVPR17]	U+S	34.3	-	-	45.6
		CRAFT [PAMI18]	S	50.3	22.4	-	-
		<b>C-DPCF [ours]</b>	S	<b>76.3</b>	<b>34.8</b>	<b>79.4</b>	<b>89.1</b>
Mobilenet -V2	{	DIMN [CVPR19]	DG	51.2	29.3	-	-
		DN [BMVC19]	DG	58.8	39.7	73.6	-
		DN + <b>ours</b>	DG+S	<b>73.9</b>	<b>42.3</b>	<b>84.1</b>	-

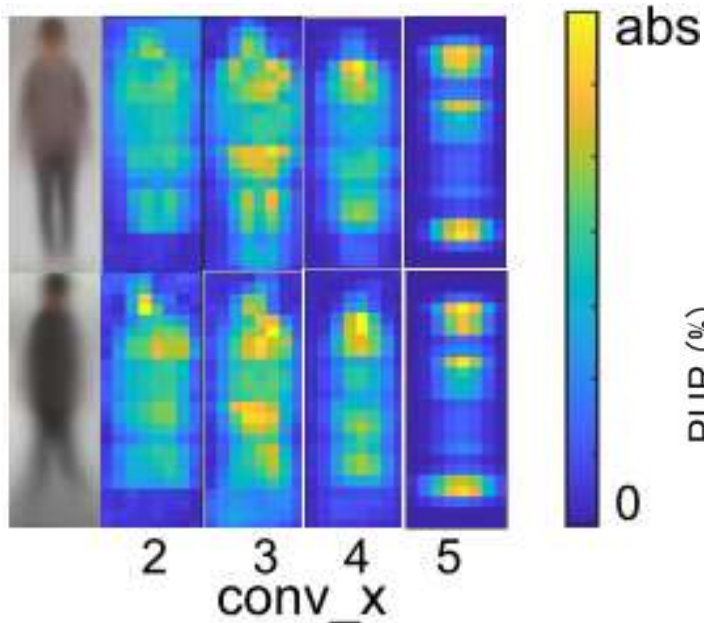
# Analysis

## Random projection

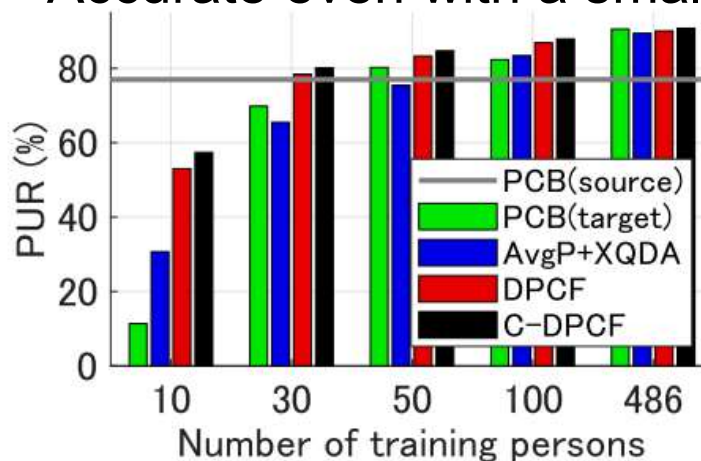
	PUR	Training time
w/o	90.5%	1684.4 sec
w/	91.3%	49.0 sec

34.4x faster

## Weight maps



## Accurate even with a small training data



- C-DPCF improves PCB(source) with 30 persons
- Camera-specific weight maps always outperforms common weight maps