

Representation Learning in Person Re-ID

Depu Meng

2019.01.19

Content

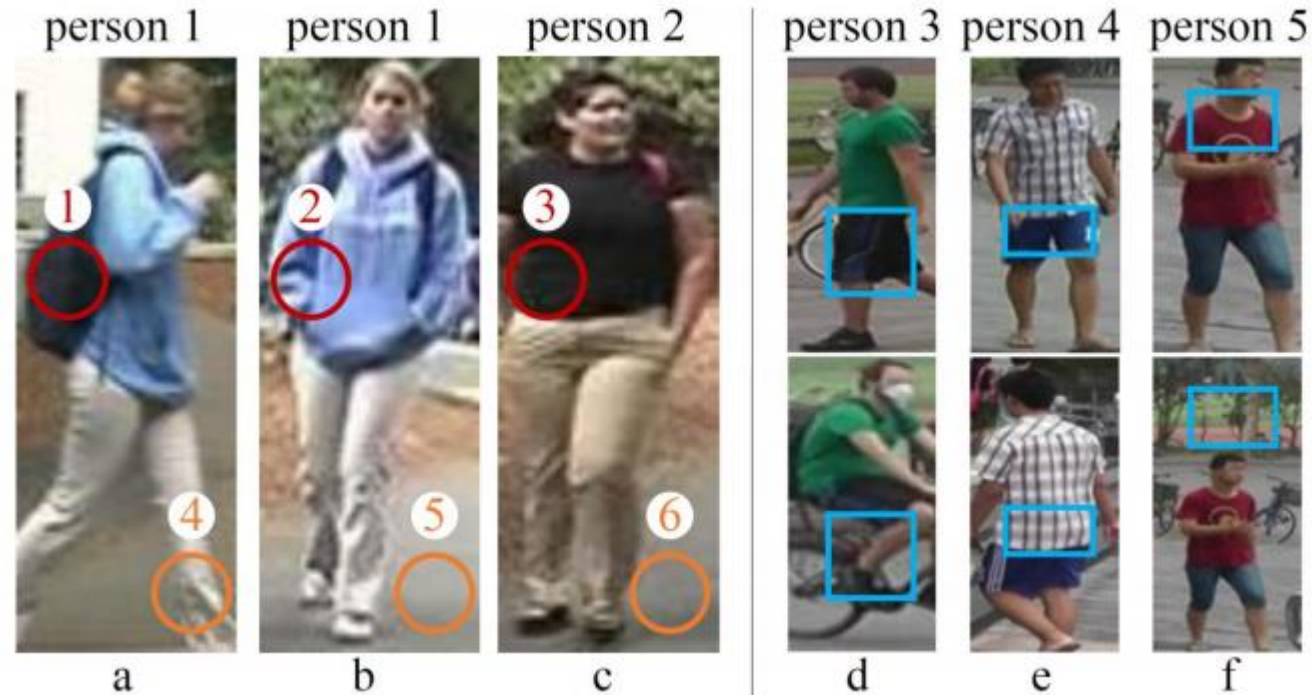
- Part-aligned Representation Learning
- Image generation in Re-ID
- GAN as supervisor

Content

- Part-aligned Representation Learning
- Image generation in Re-ID
- GAN as supervisor

Part-aligned representation

- Person Re-ID is based on comparison of body parts
- Body parts are misaligned in the most cases



Related Works

- An Improved Deep Learning Architecture for Person Re-Identification, CVPR 2015
- Deeply-Learned Part-Aligned Representations for Person Re-Identification, ICCV 2017
- Attention-Aware Compositional Network for Person Re-Identification, CVPR 2018
- Part-Aligned Bilinear Representations for Person Re-Identification, ECCV 2018

Related Works

- An Improved Deep Learning Architecture for Person Re-Identification, w/o pose CVPR 2015
- Deeply-Learned Part-Aligned Representations for Person Re-Identification, ICCV 2017
- Attention-Aware Compositional Network for Person Re-Identification, CVPR 2018
- Part-Aligned Bilinear Representations for Person Re-Identification, ECCV 2018

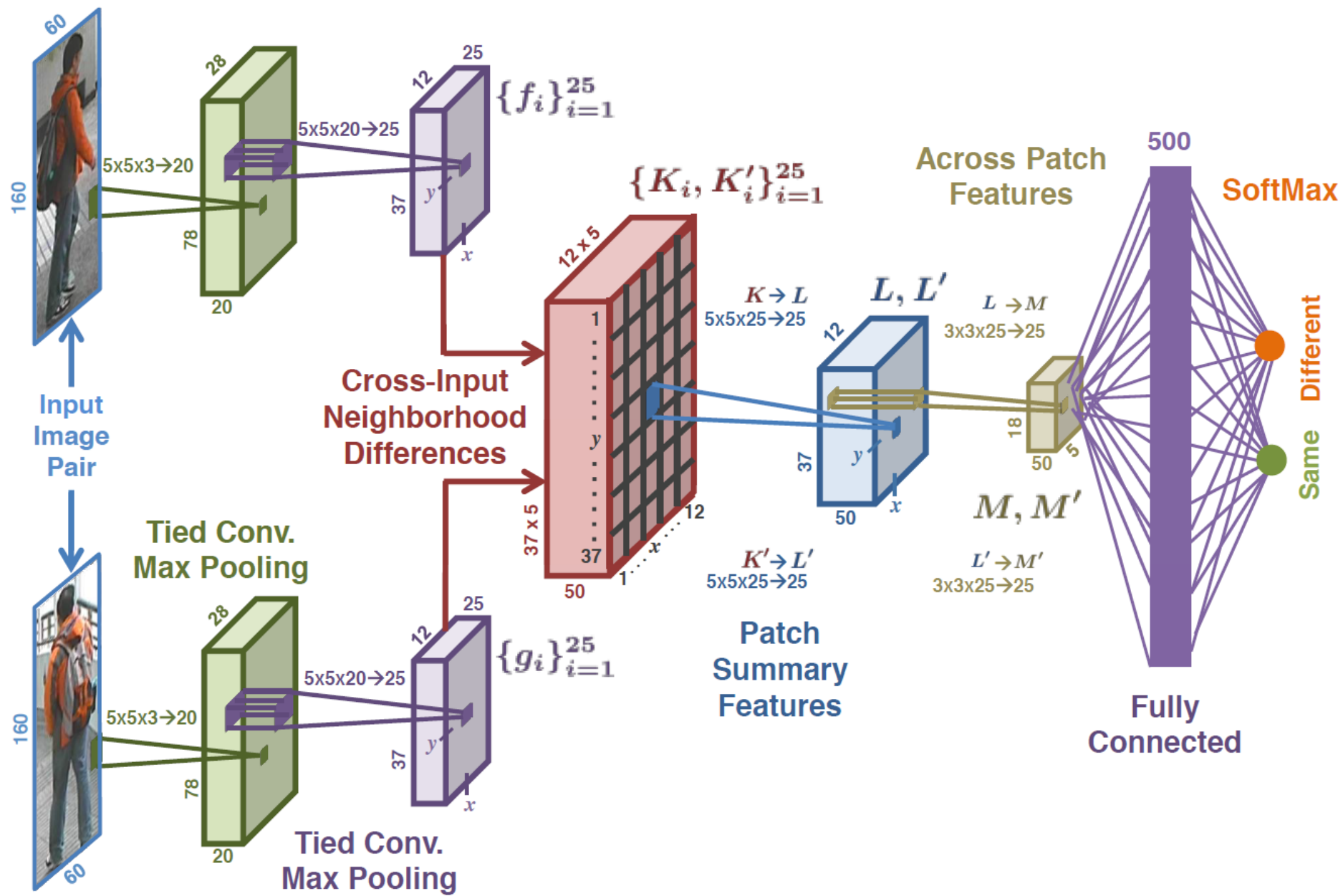
Related Works

- An Improved Deep Learning Architecture for Person Re-Identification, w/o pose CVPR 2015
- Deeply-Learned Part-Aligned Representations for Person Re-Identification, ICCV 2017
- Attention-Aware Compositional Network for Person Re-Identification, w/ pose CVPR 2018
- Part-Aligned Bilinear Representations for Person Re-Identification, ECCV 2018

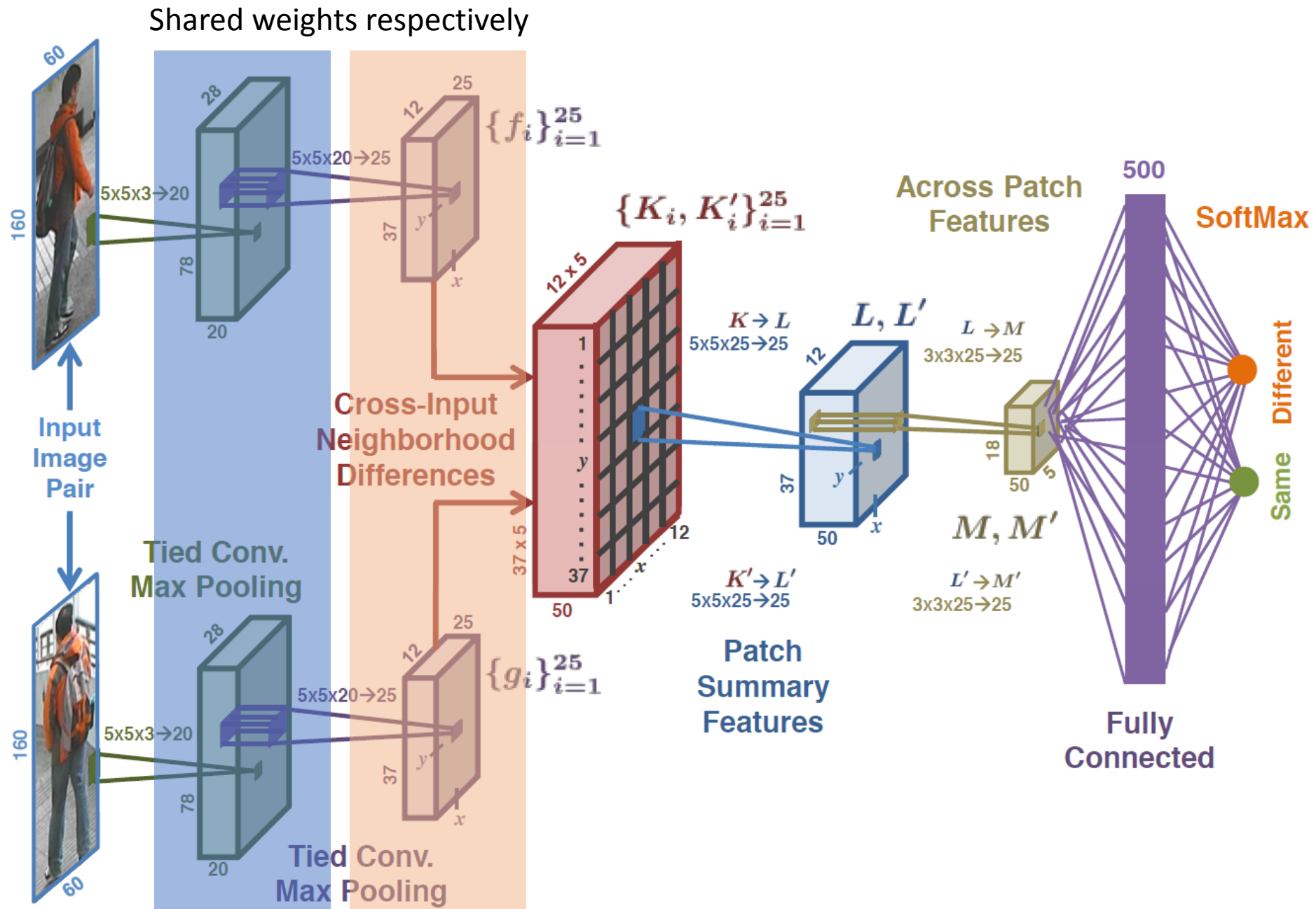
Related Works

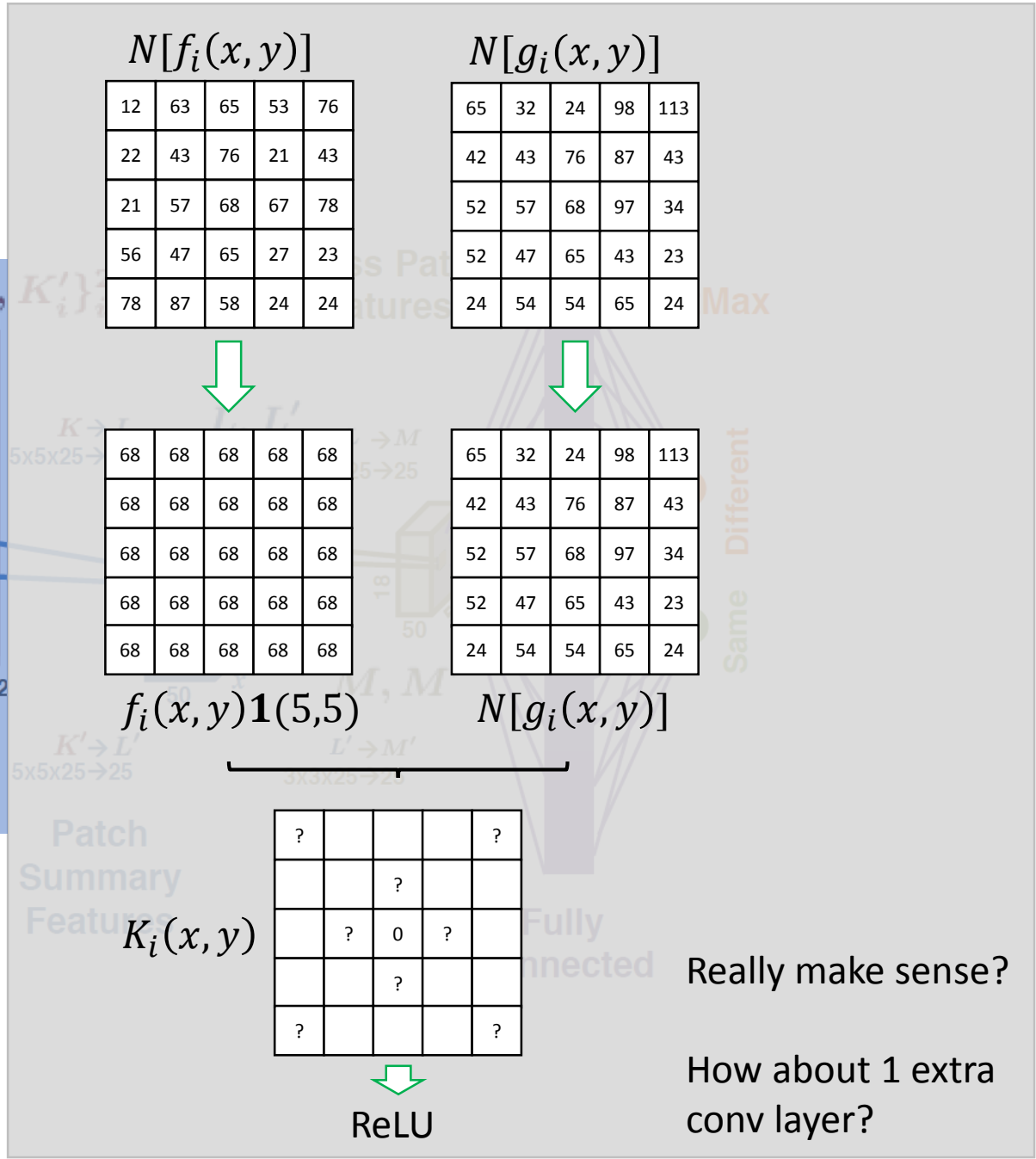
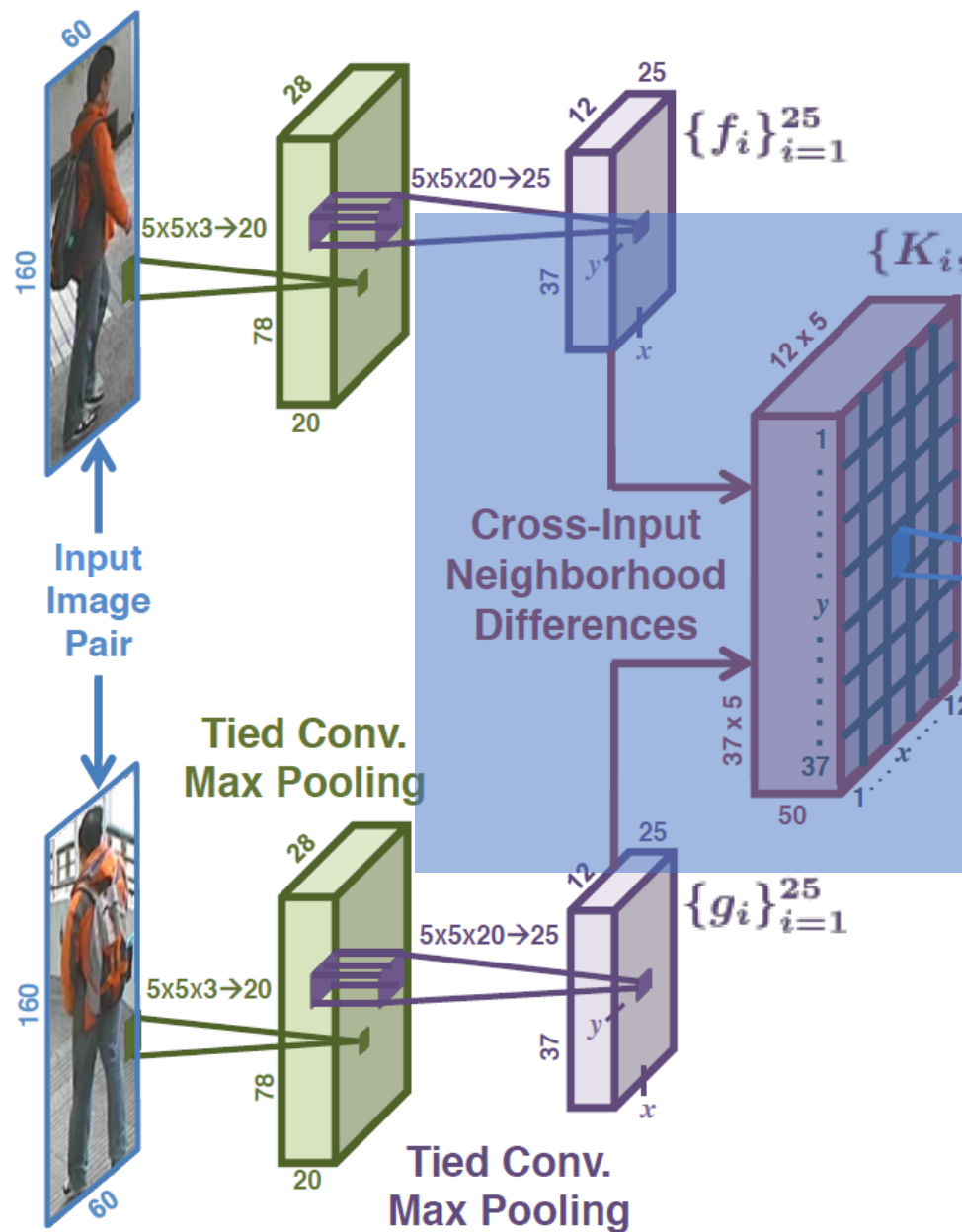
- An Improved Deep Learning Architecture for Person Re-Identification, ^{w/o pose} CVPR 2015
- Deeply-Learned Part-Aligned Representations for Person Re-Identification, ICCV 2017
- Attention-Aware Compositional Network for Person Re-Identification, ^{w/ pose} CVPR 2018
- Part-Aligned Bilinear Representations for Person Re-Identification, ECCV 2018

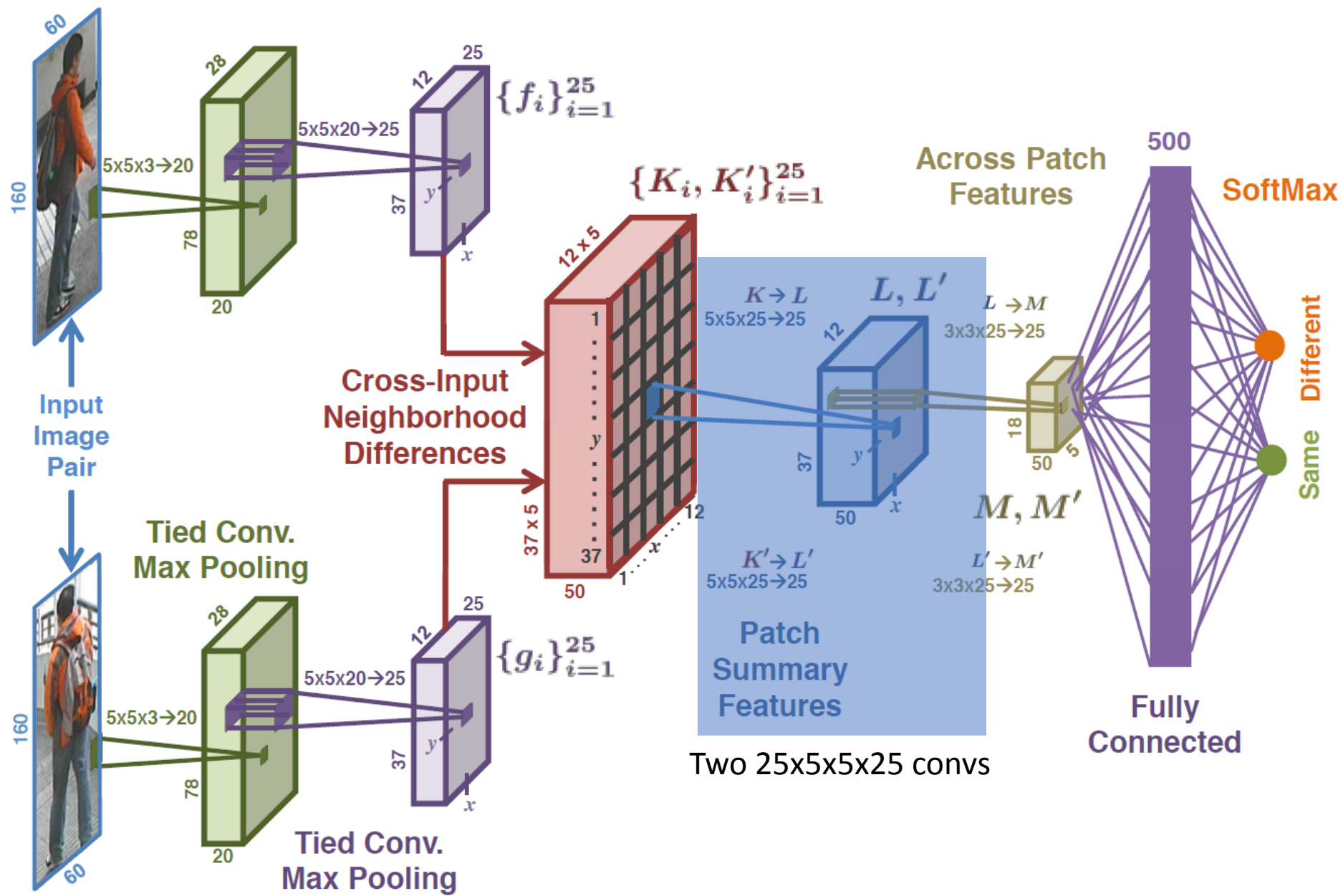
Architecture

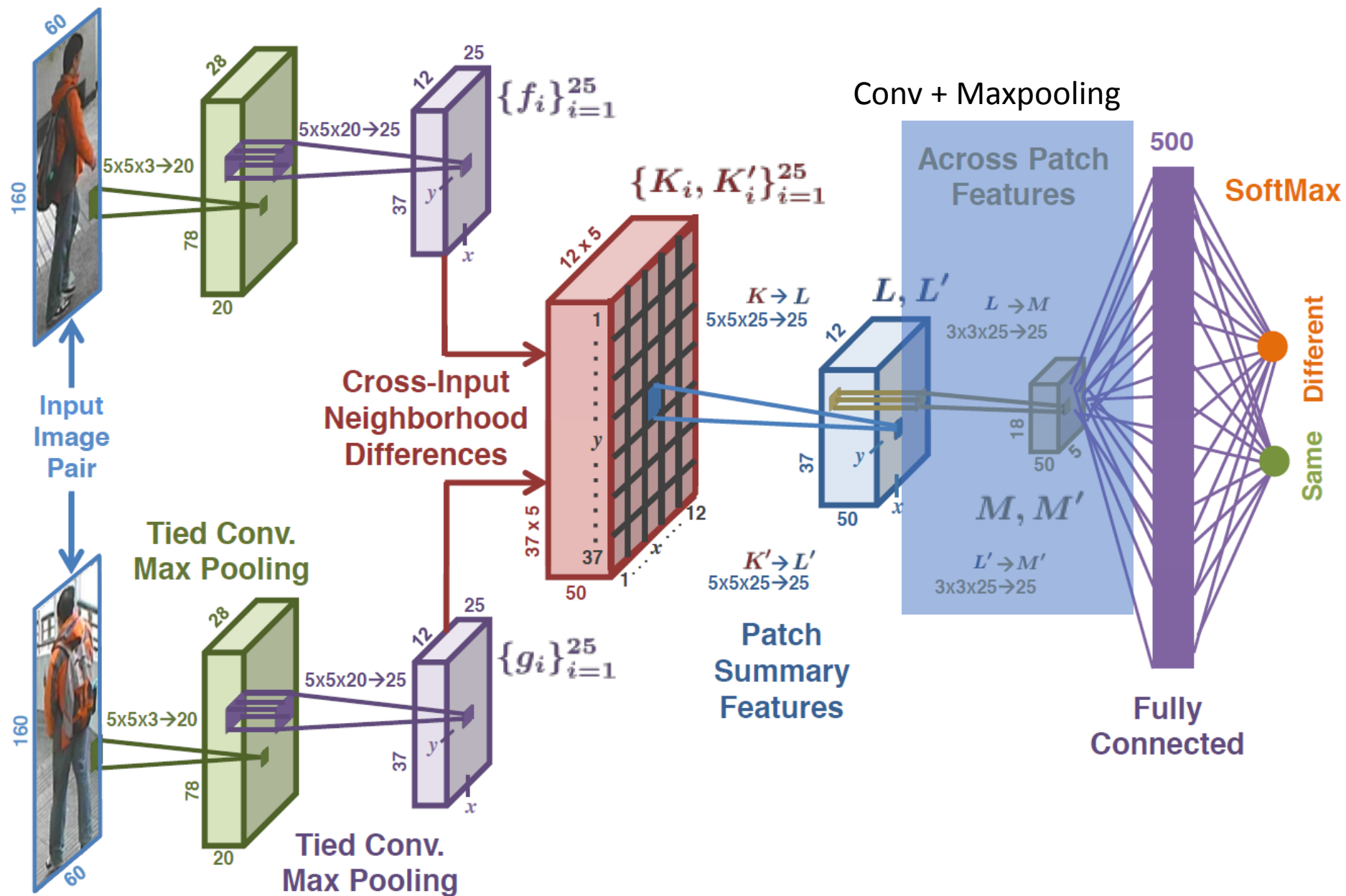


Ejaz Ahmed, Michael Jones, Tim K. Marks: An Improved Deep Learning Architecture for Person Re-Identification, CVPR 2015



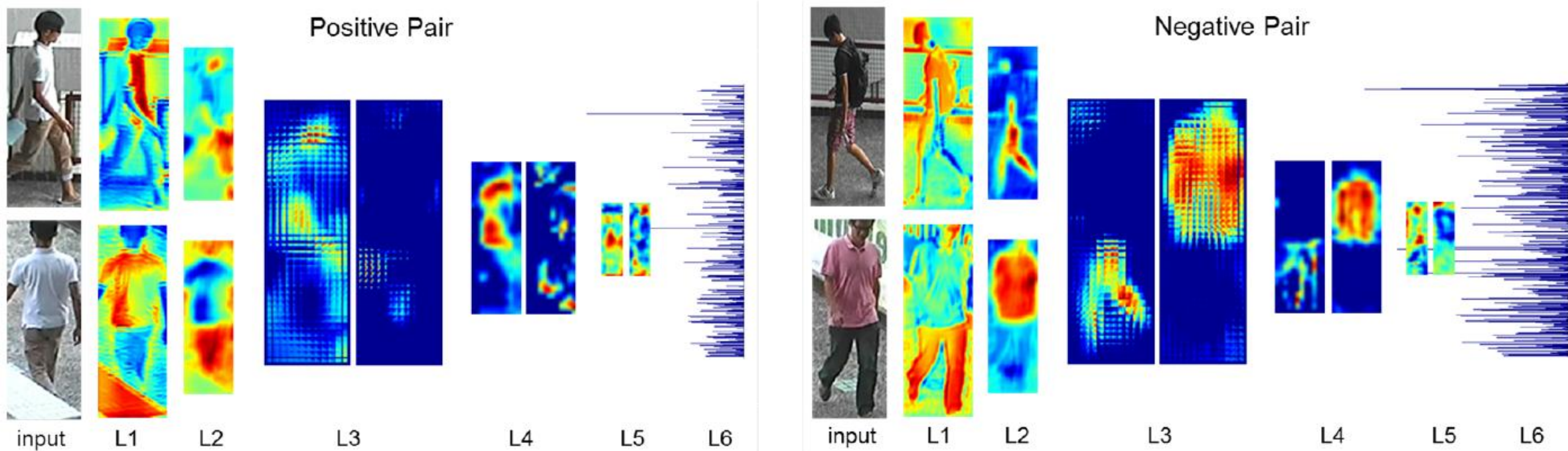






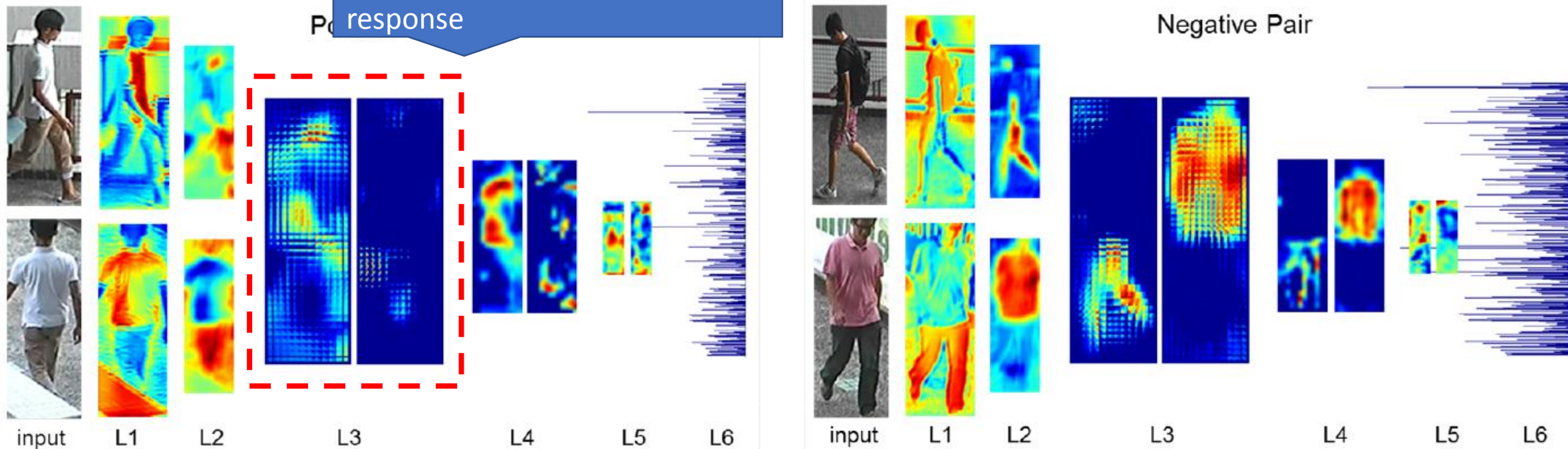
Ejaz Ahmed, Michael Jones, Tim K. Marks: An Improved Deep Learning Architecture for Person Re-Identification, CVPR 2015

Response Maps



Response Maps

Q: Why these two maps look complementary?
A: ReLU removes the negative response



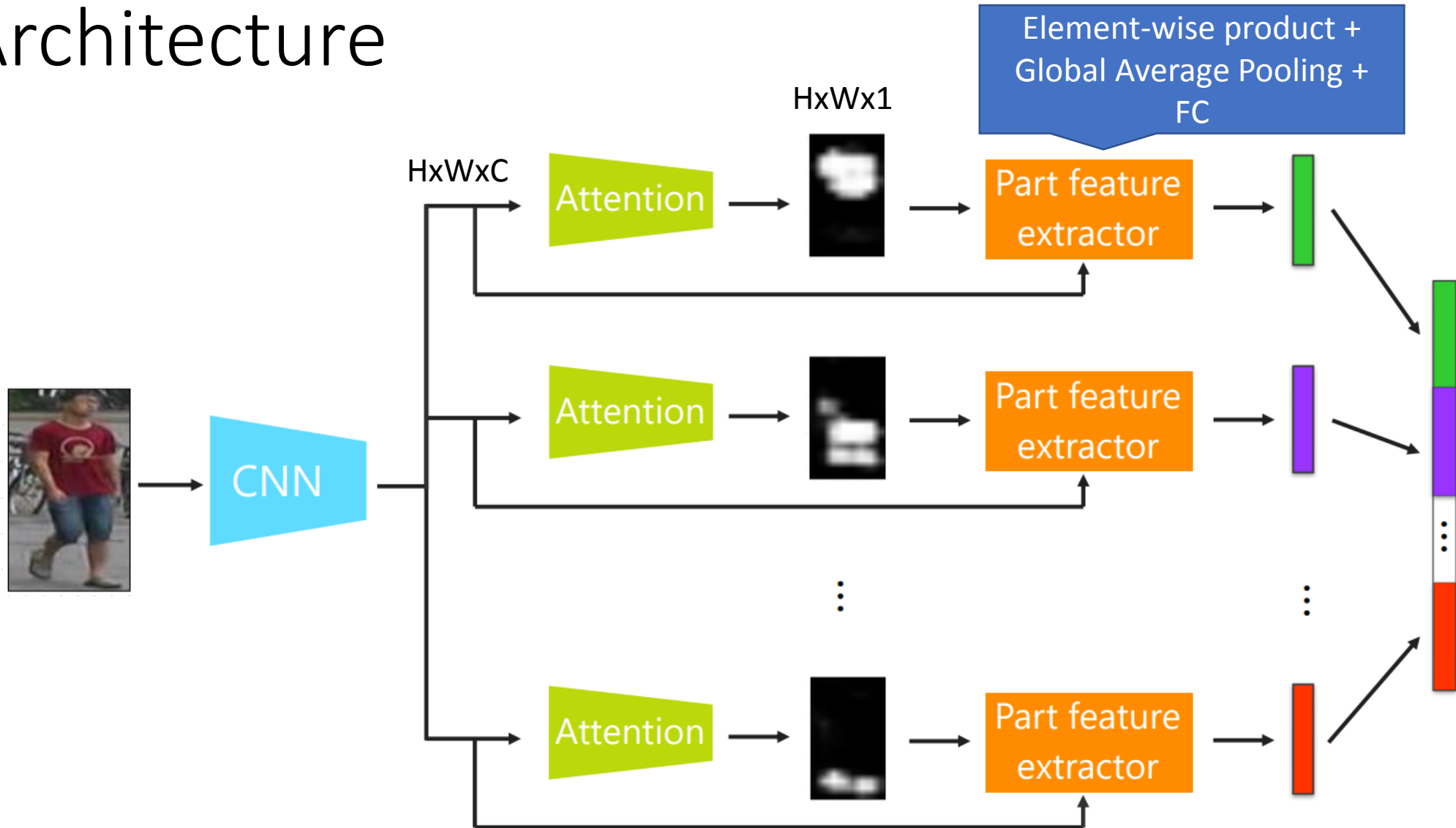
Overview

- CNN based model
- Considered misalignment of different parts
- A naïve solution

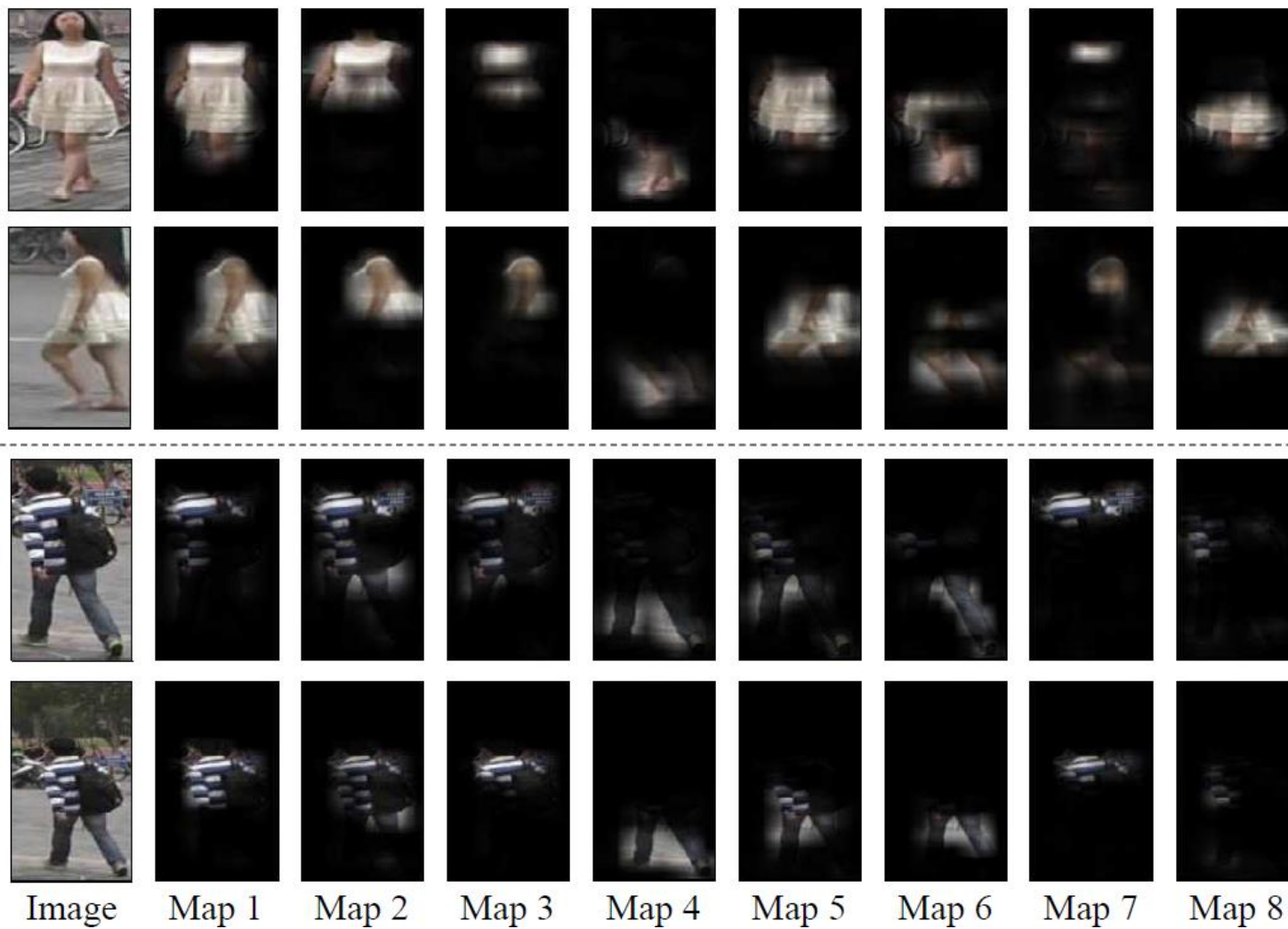
Related Works

- An Improved Deep Learning Architecture for Person Re-Identification, w/o pose
CVPR 2015
- Deeply-Learned Part-Aligned Representations for Person Re-Identification, ICCV 2017
- Attention-Aware Compositional Network for Person Re-Identification, w/ pose
CVPR 2018
- Part-Aligned Bilinear Representations for Person Re-Identification, ECCV 2018

Architecture



Response Maps



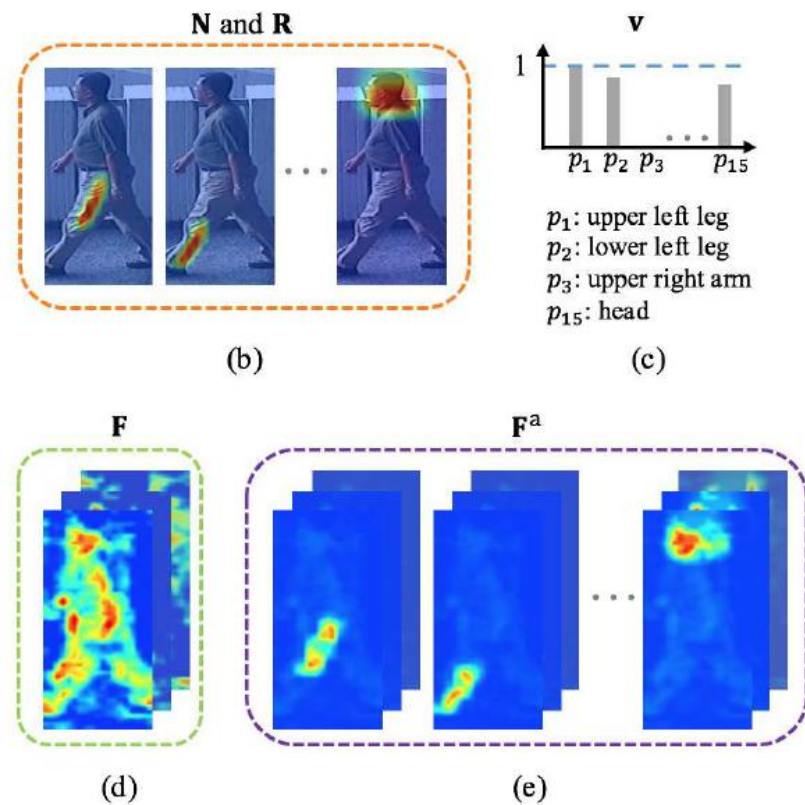
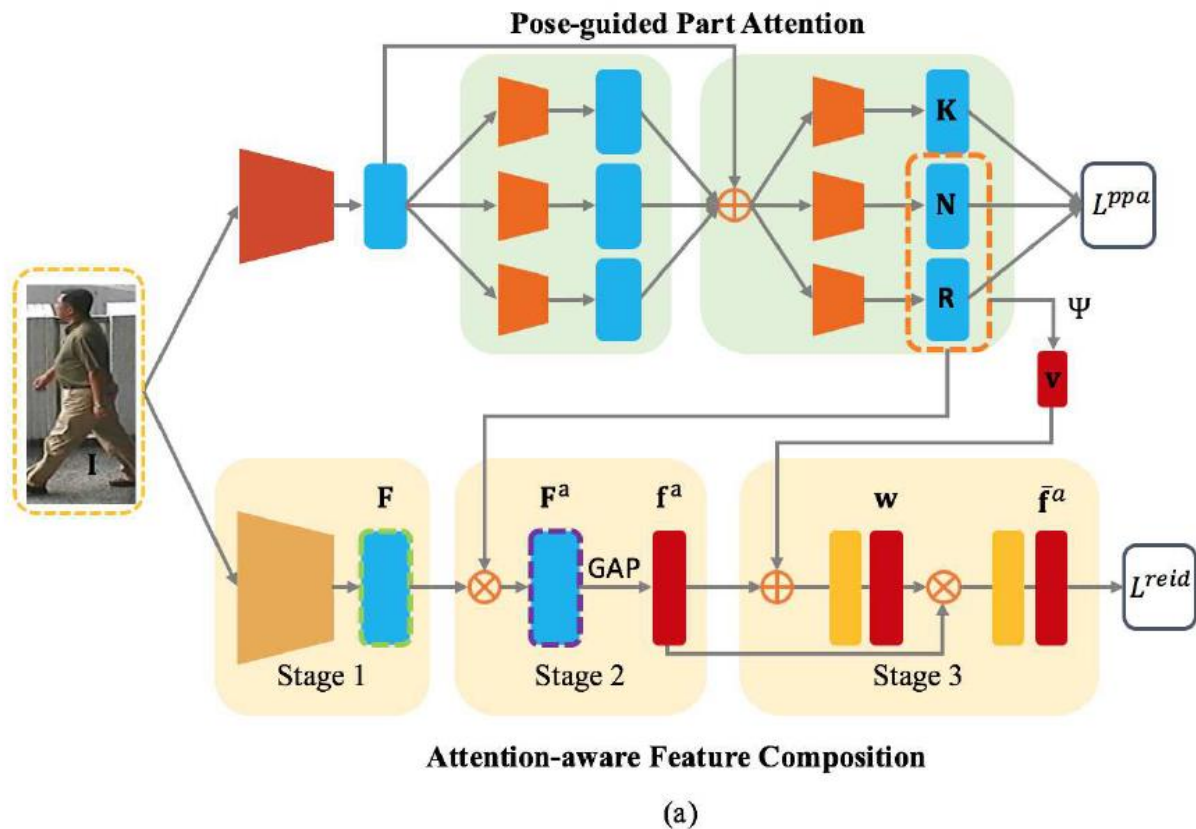
Overview

- Attention-based method
- Lack of guidance in splitting parts

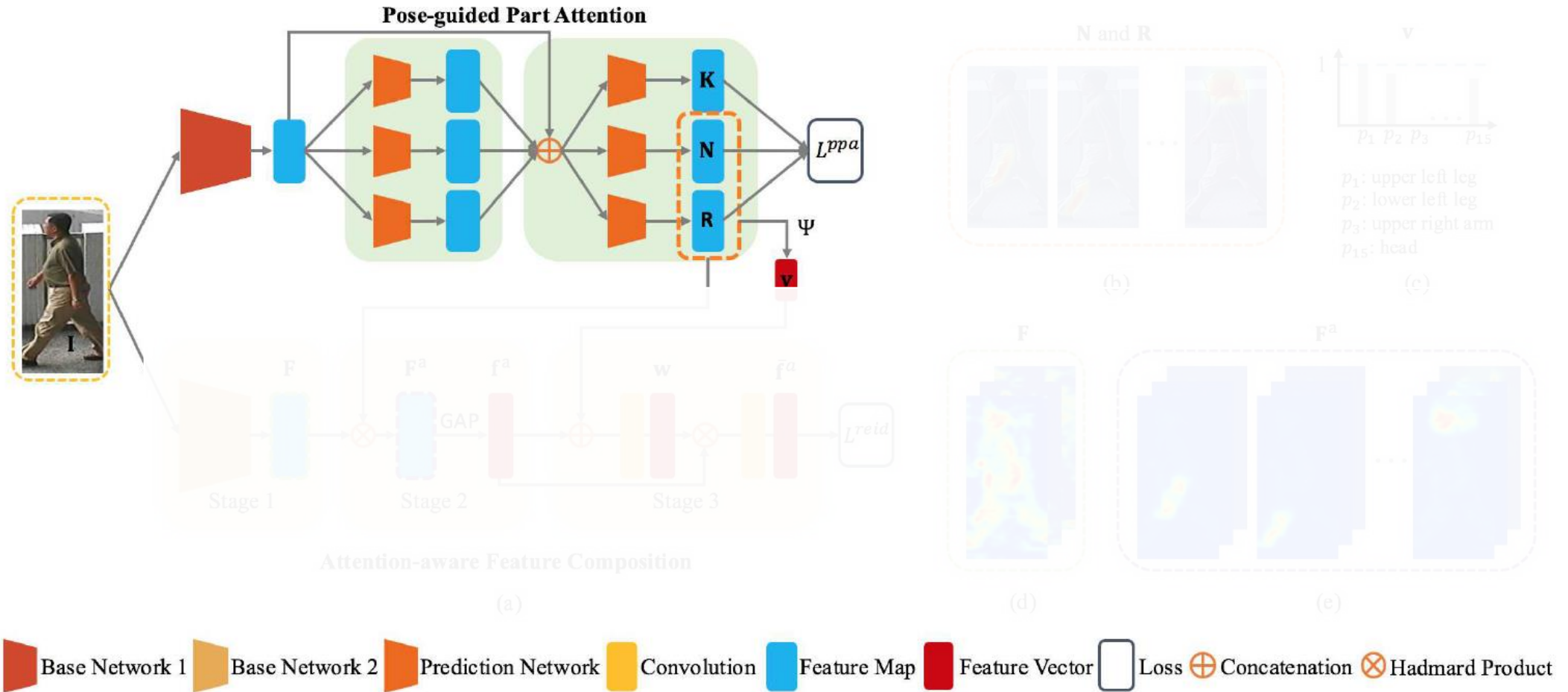
Related Works

- An Improved Deep Learning Architecture for Person Re-Identification, **w/o pose**, CVPR 2015
- Deeply-Learned Part-Aligned Representations for Person Re-Identification, ICCV 2017
- **Attention-Aware Compositional Network for Person Re-Identification, w/ pose**, CVPR 2018
- Part-Aligned Bilinear Representations for Person Re-Identification, ECCV 2018

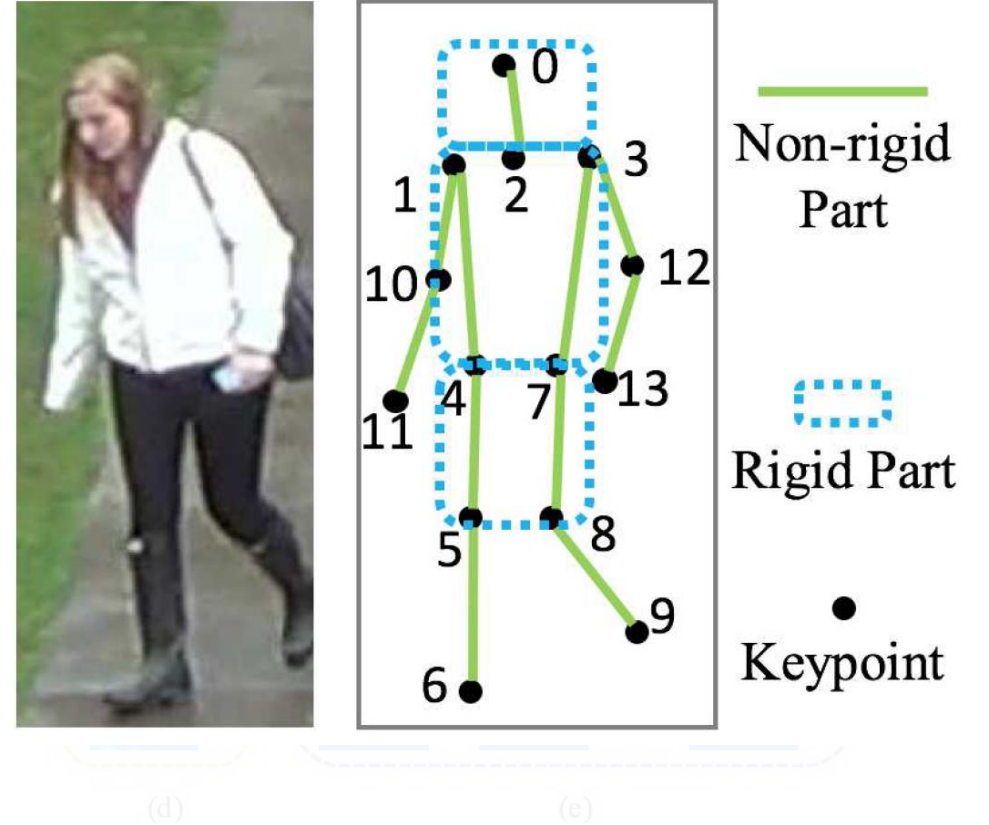
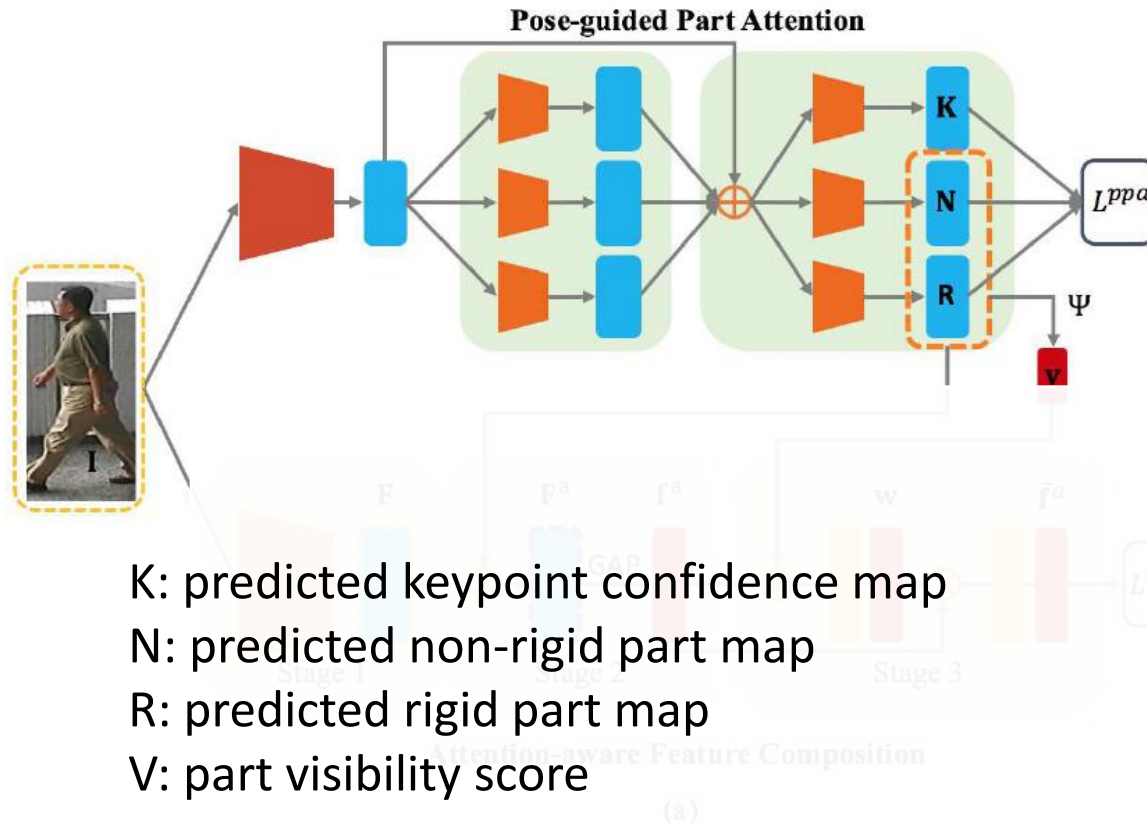
Architecture



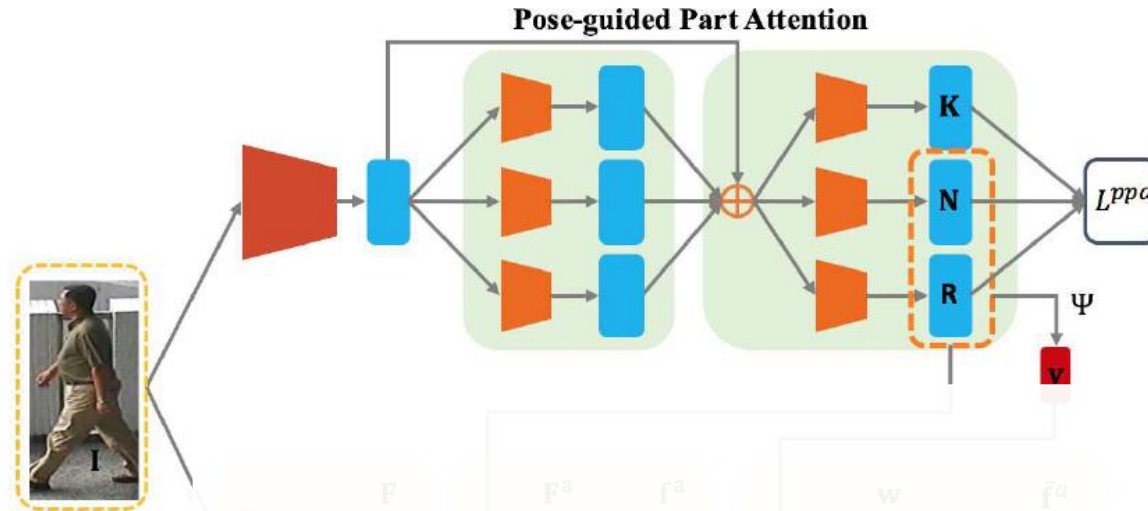
Architecture



Architecture

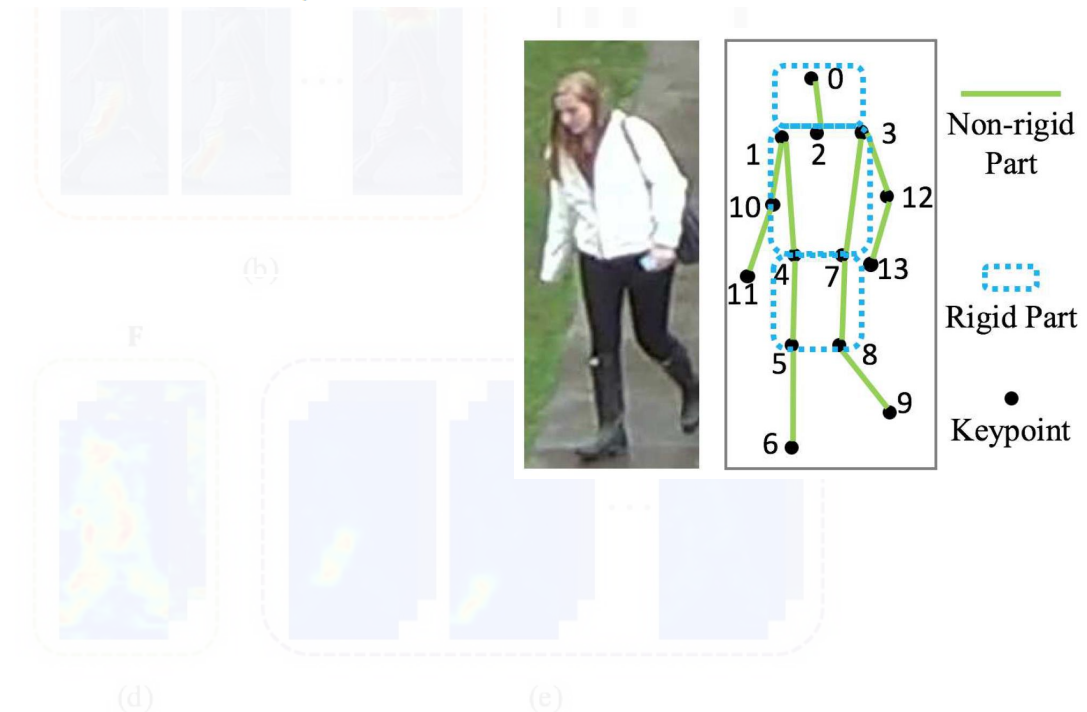


Architecture

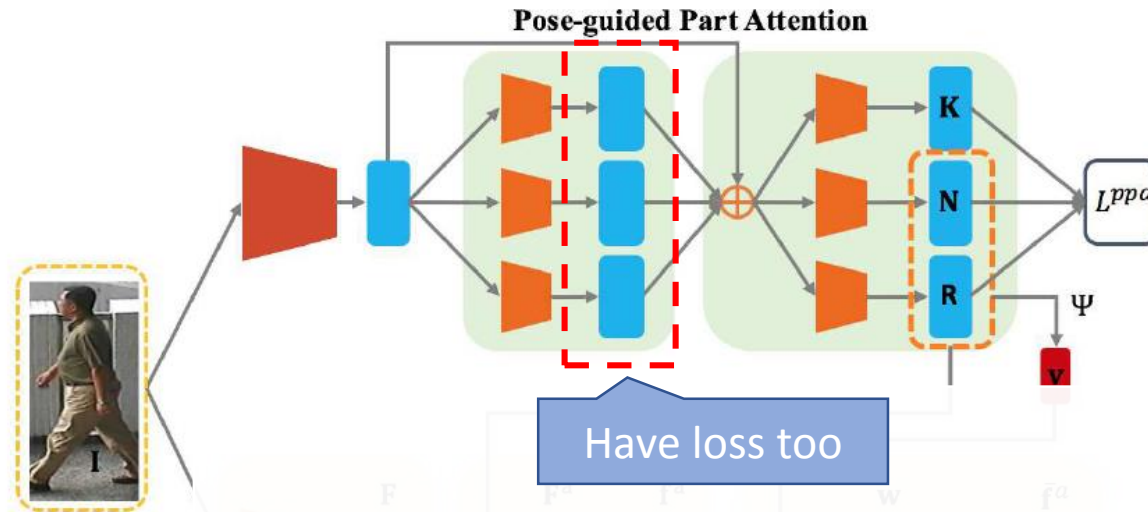


K: predicted keypoint confidence map
 N: predicted non-rigid part map
 R: predicted rigid part map
 V: part visibility score

$$L^{ppa}(\rho, \phi, \psi) = \sum_{t=1,2} L^k(\mathbf{K}^t) + \mu_1 L^n(\mathbf{N}^t) + \mu_2 L^r(\mathbf{R}^t).$$

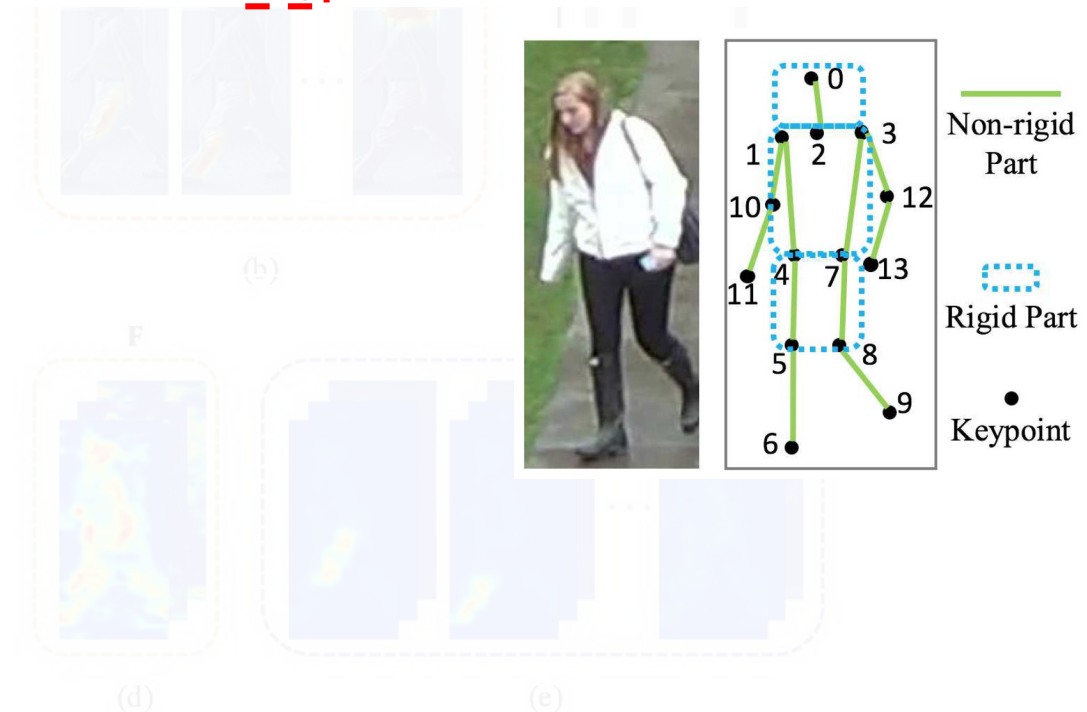


Architecture

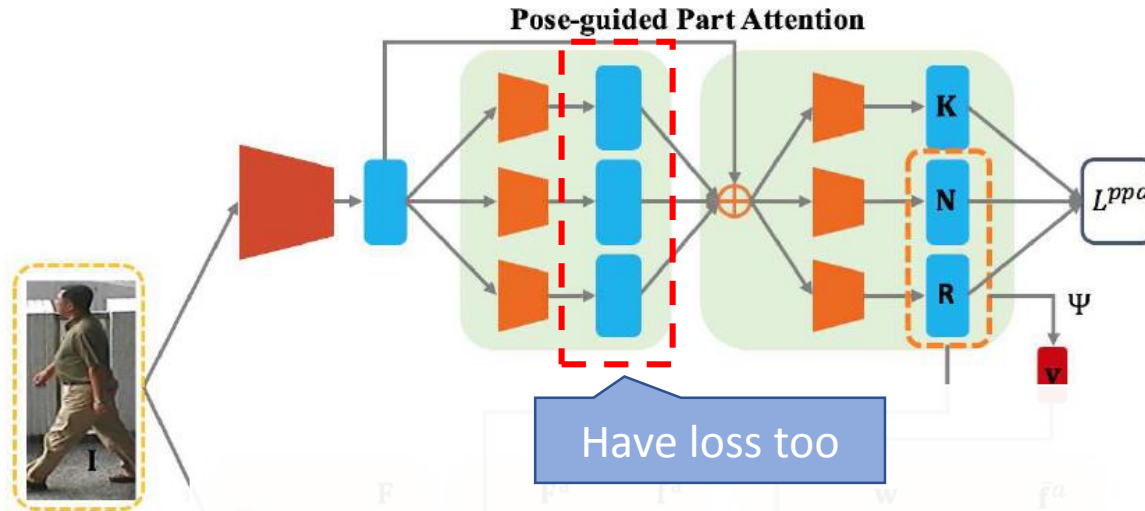


K: predicted keypoint confidence map
 N: predicted non-rigid part map
 R: predicted rigid part map
 V: part visibility score

$$L^{ppa}(\rho, \phi, \psi) = \sum_{t=1,2} L^k(\mathbf{K}^t) + \mu_1 L^n(\mathbf{N}^t) + \mu_2 L^r(\mathbf{R}^t).$$



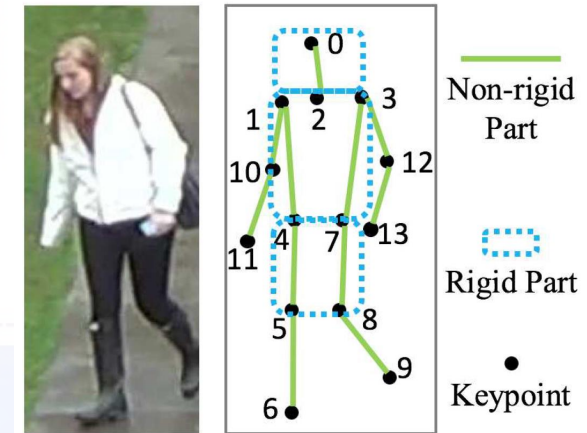
Architecture



K: predicted keypoint confidence map
 N: predicted non-rigid part map
 R: predicted rigid part map
 V: part visibility score

$$L^{ppa}(\rho, \phi, \psi) = \sum_{t=1,2} L^k(\mathbf{K}^t) + \mu_1 L^n(\mathbf{N}^t) + \mu_2 L^r(\mathbf{R}^t)$$

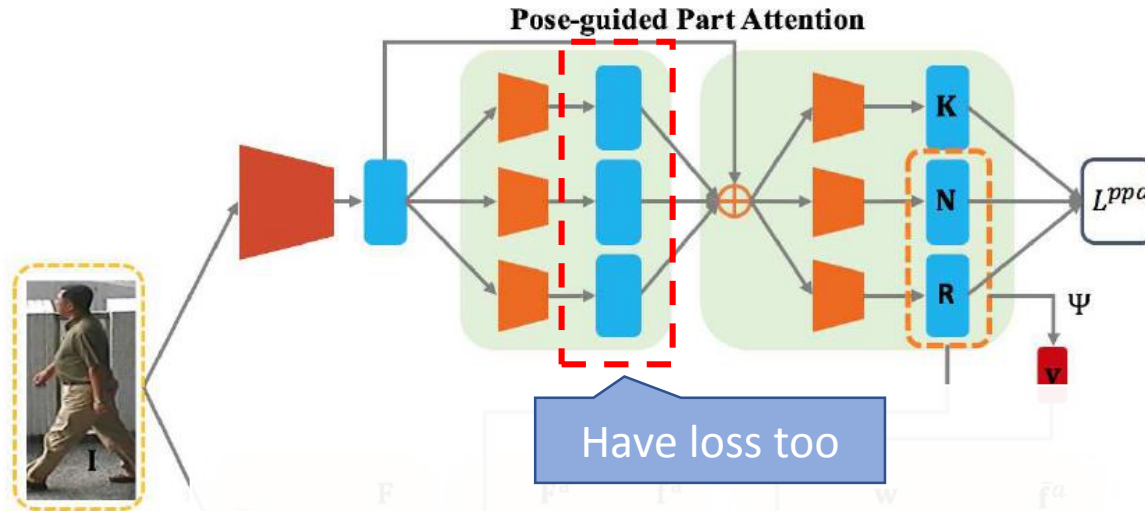
$$L^k(\mathbf{K}) = \frac{1}{C^k} \sum_{i=1}^{C^k} \|\mathbf{K}_i^* - \mathbf{K}_i\|^2$$



K_i^* is generated by applying Gaussian kernel at the true location of keypoint i



Architecture

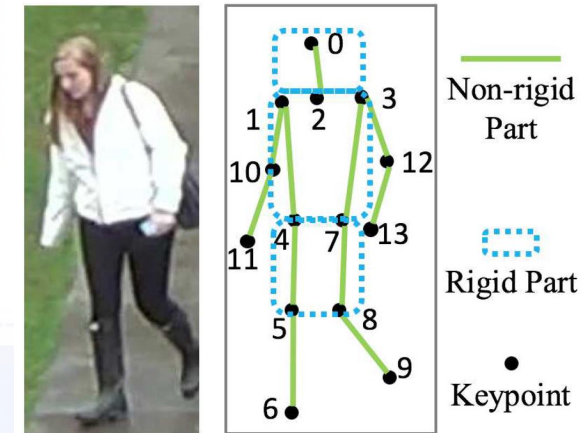


K: predicted keypoint confidence map
 N: predicted non-rigid part map
 R: predicted rigid part map
 V: part visibility score

$$L^{ppa}(\rho, \phi, \psi) = \sum_{t=1,2} L^k(\mathbf{K}^t) + \mu_1 L^n(\mathbf{N}^t) + \mu_2 L^r(\mathbf{R}^t)$$

$$L^k(\mathbf{K}) = \frac{1}{C^k} \sum_{i=1}^{C^k} \|\mathbf{K}_i^* - \mathbf{K}_i\|^2$$

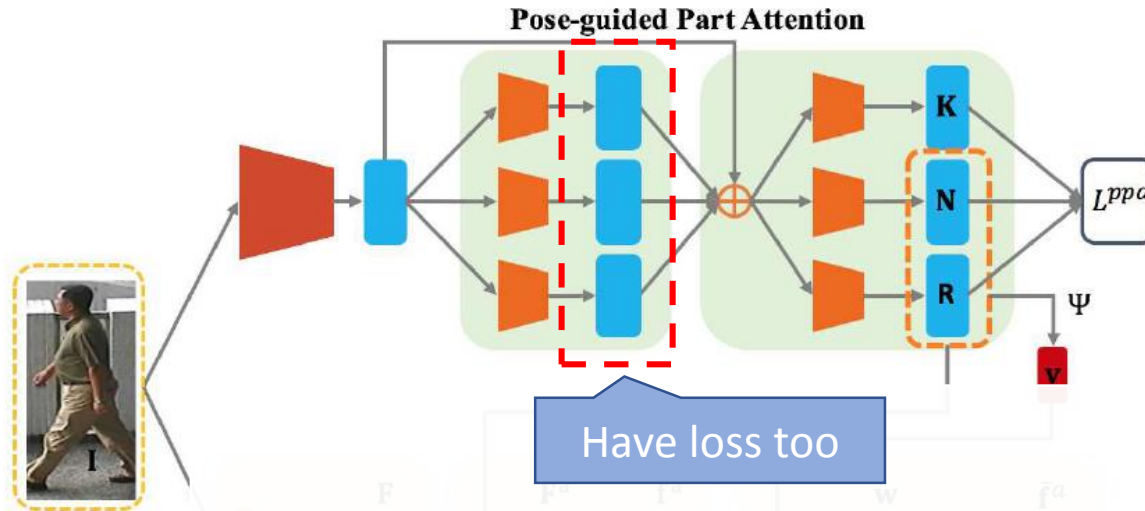
$$L^n(\mathbf{N}) = \frac{1}{C^n} \sum_{p=1}^{C^n} \|\mathbf{N}_p^* - \mathbf{N}_p\|^2$$



N_p^* is the p -th non-grid part, defined as a rectangle area connecting two keypoints with bandwidth σ



Architecture



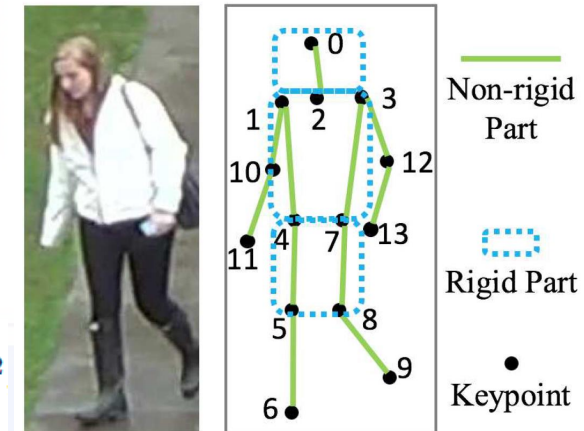
K: predicted keypoint confidence map
 N: predicted non-rigid part map
 R: predicted rigid part map
 V: part visibility score

$$L^{ppa}(\rho, \phi, \psi) = \sum_{t=1,2} L^k(\mathbf{K}^t) + \mu_1 L^n(\mathbf{N}^t) + \mu_2 L^r(\mathbf{R}^t)$$

$$L^k(\mathbf{K}) = \frac{1}{C^k} \sum_{i=1}^{C^k} \|\mathbf{K}_i^* - \mathbf{K}_i\|^2$$

$$L^n(\mathbf{N}) = \frac{1}{C^n} \sum_{p=1}^{C^n} \|\mathbf{N}_p^* - \mathbf{N}_p\|^2$$

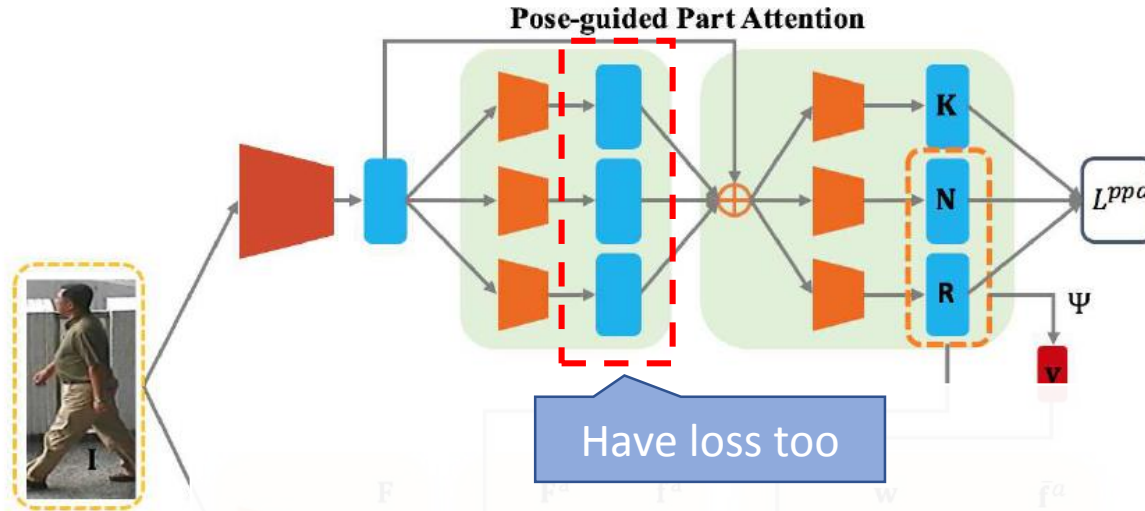
$$L^r(\mathbf{R}, \mathbf{N}) = \frac{1}{C^r} \sum_{p=1}^{C^r} \|\mathbf{R}_p^* - \hat{\mathbf{R}}_p\|^2$$



R_p^* is defined by a rectangle tightly contains specified keypoints



Architecture



K: predicted keypoint confidence map
 N: predicted non-rigid part map
 R: predicted rigid part map
 V: part visibility score

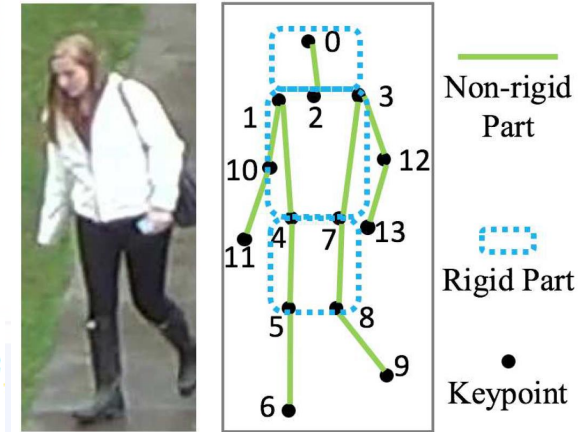
$$L^{ppa}(\rho, \phi, \psi) = \sum_{t=1,2} L^k(\mathbf{K}^t) + \mu_1 L^n(\mathbf{N}^t) + \mu_2 L^r(\mathbf{R}^t)$$

$$L^k(\mathbf{K}) = \frac{1}{C^k} \sum_{i=1}^{C^k} \|\mathbf{K}_i^* - \mathbf{K}_i\|^2$$

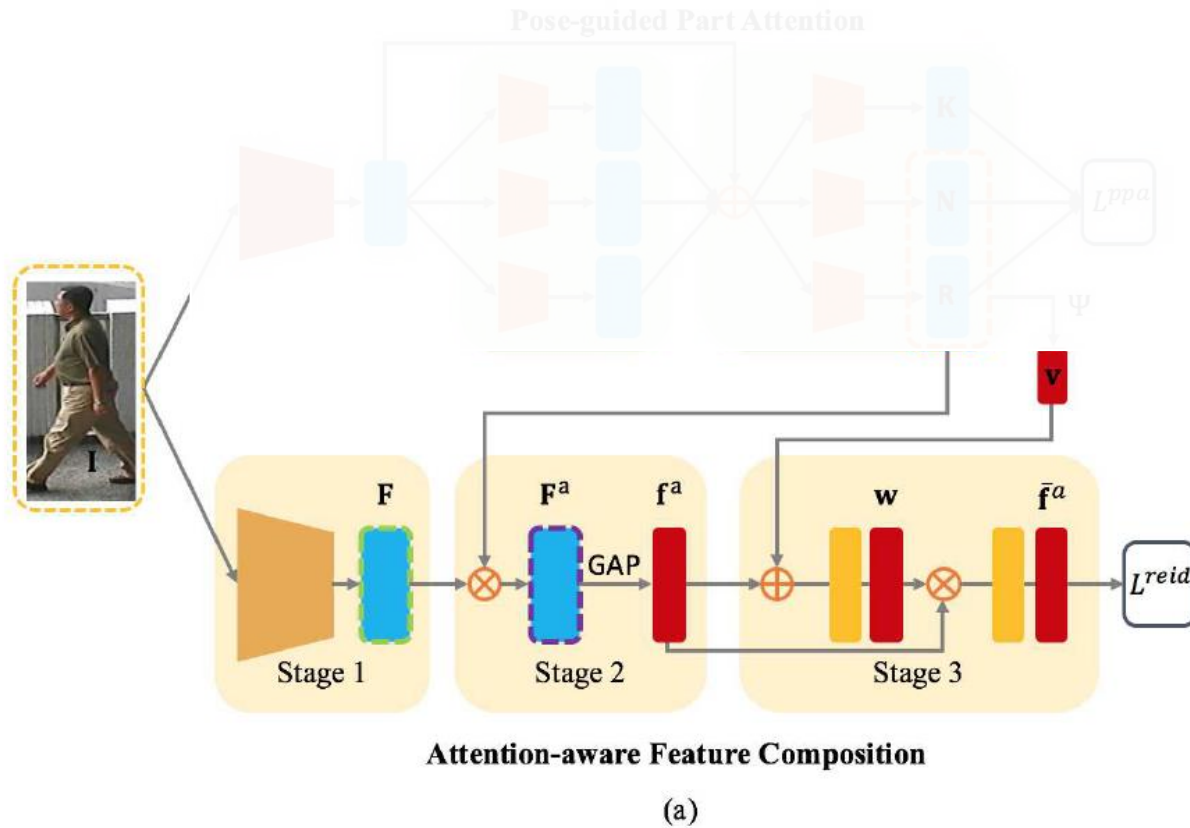
$$L^n(\mathbf{N}) = \frac{1}{C^n} \sum_{p=1}^{C^n} \|\mathbf{N}_p^* - \mathbf{N}_p\|^2$$

$$L^r(\mathbf{R}, \mathbf{N}) = \frac{1}{C^r} \sum_{p=1}^{C^r} \|\mathbf{R}_p^* - \hat{\mathbf{R}}_p\|^2$$

$$v_p = \sum_{x,y} |\mathbf{R}_p(x,y)|, \text{ or } v_p = \sum_{x,y} |\mathbf{N}_p(x,y)|$$

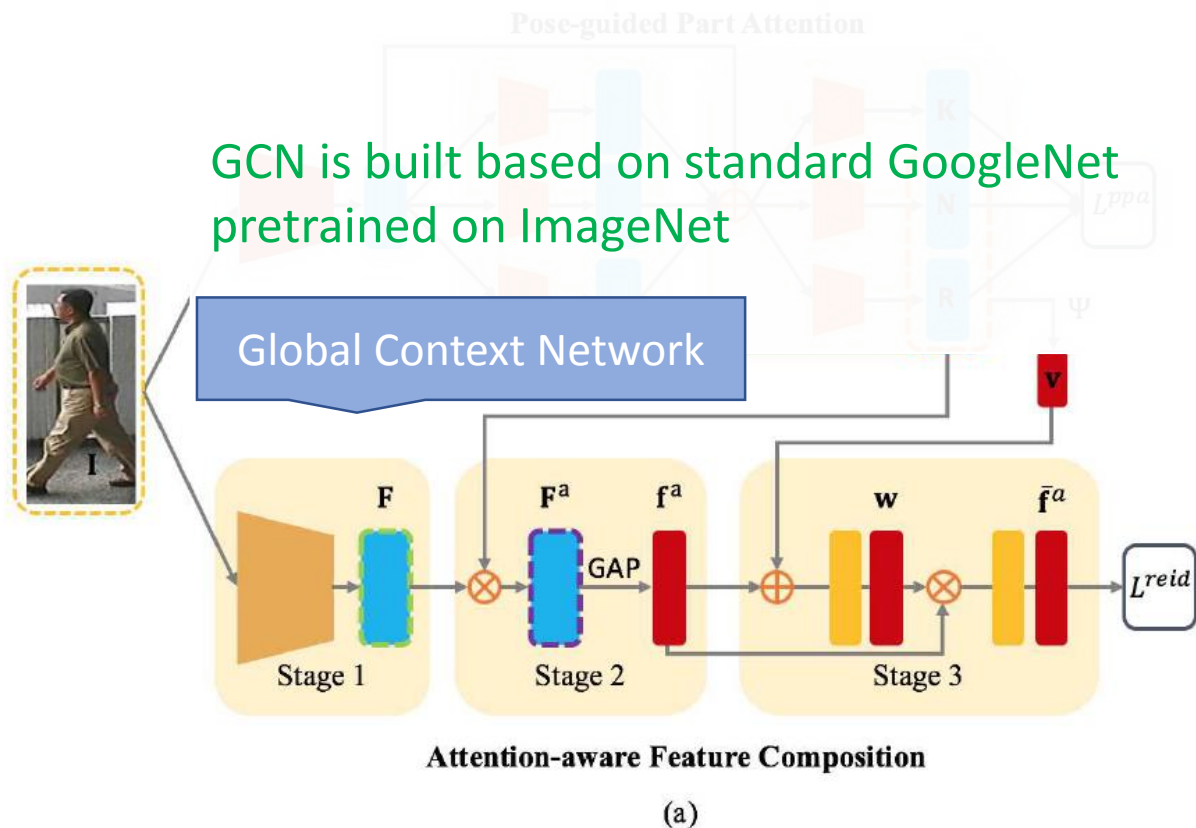


Architecture



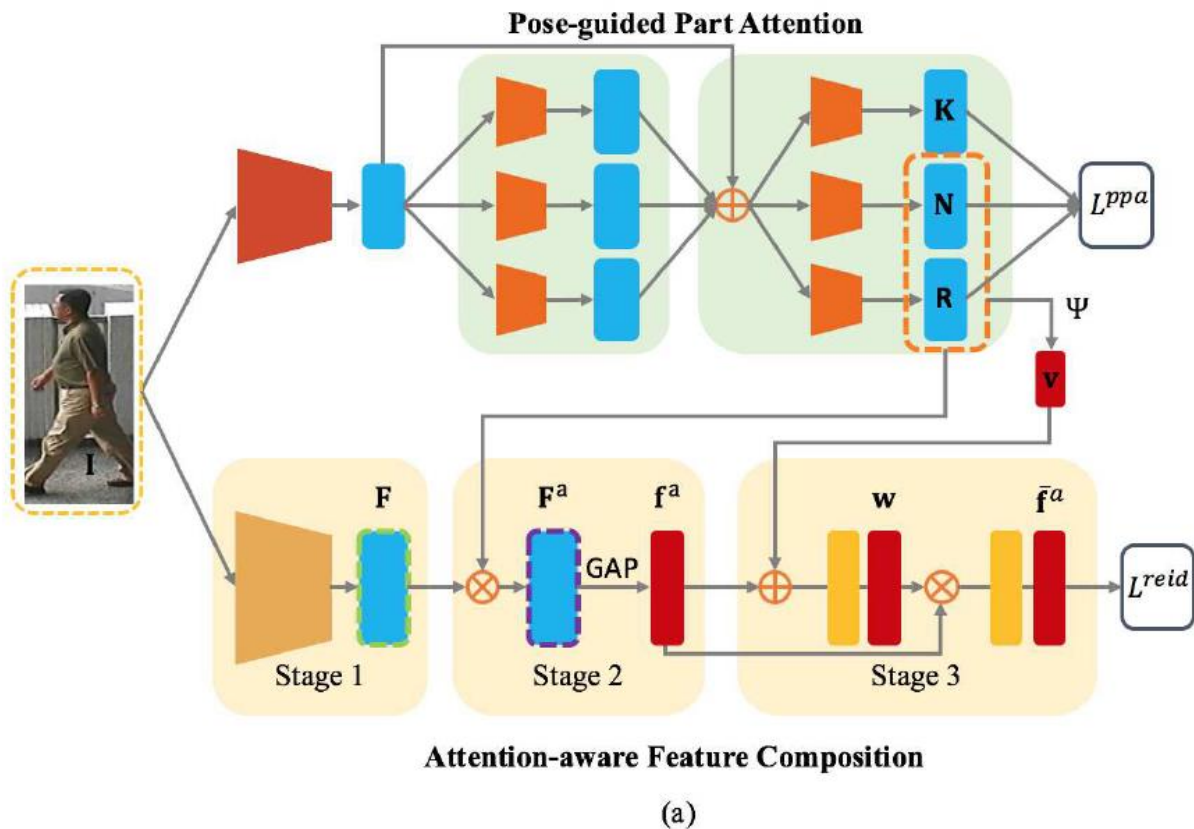
Architecture

Step 1: Train PPA independently
 Step 2: Train GCN independently



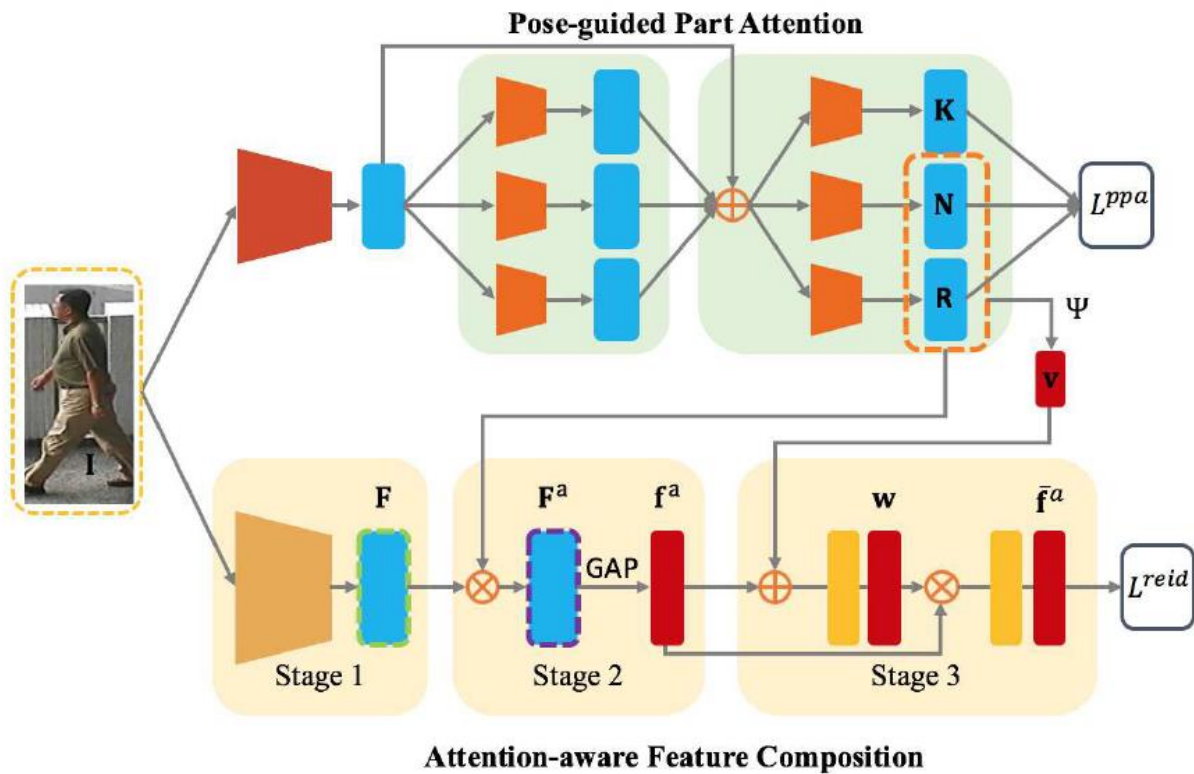
Architecture

- Step 1: Train PPA independently
- Step 2: Train GCN independently
- Step 3: Attention-Aware Feature Alignment

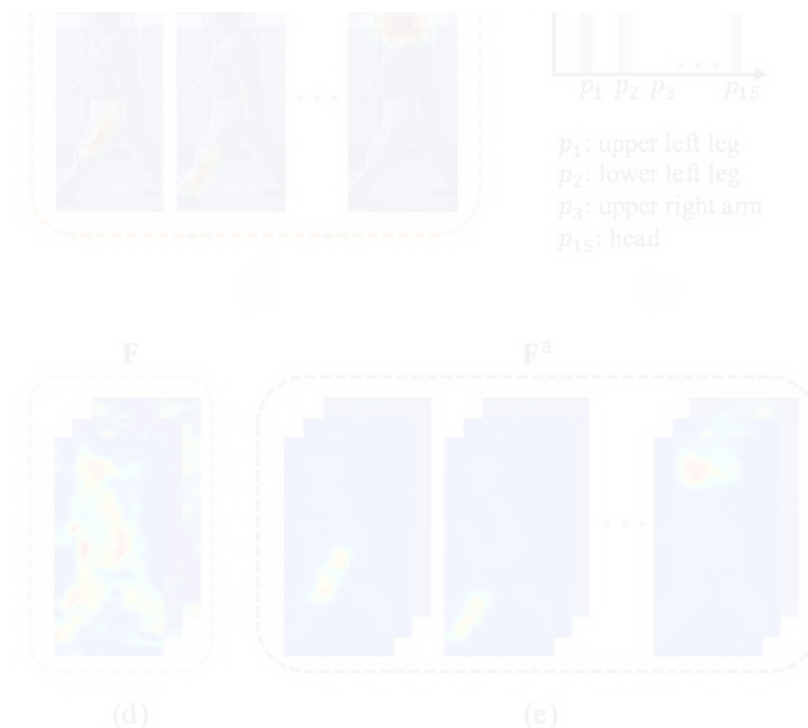


Architecture

- Step 1: Train PPA independently
- Step 2: Train GCN independently
- Step 3: Attention-Aware Feature Alignment

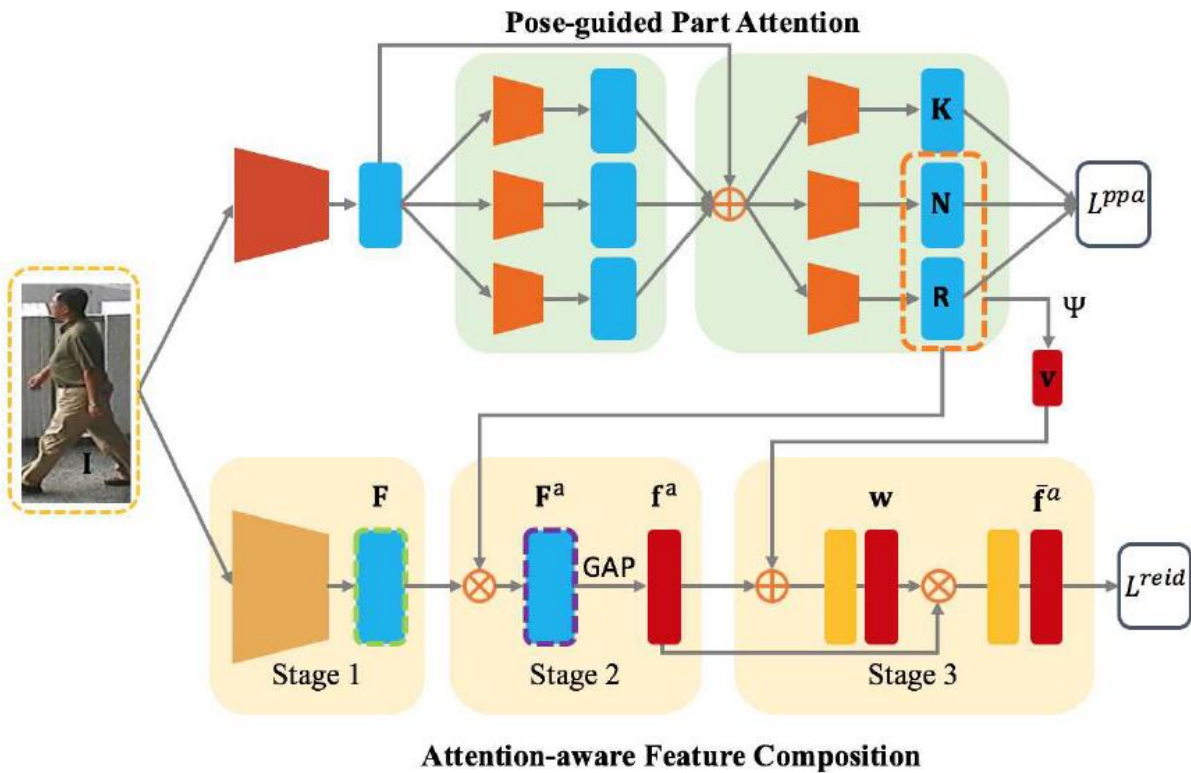


$M_p \in \{N_p, R_p\}$ is the attention map for body parts,



Architecture

- Step 1: Train PPA independently
- Step 2: Train GCN independently
- Step 3: Attention-Aware Feature Alignment

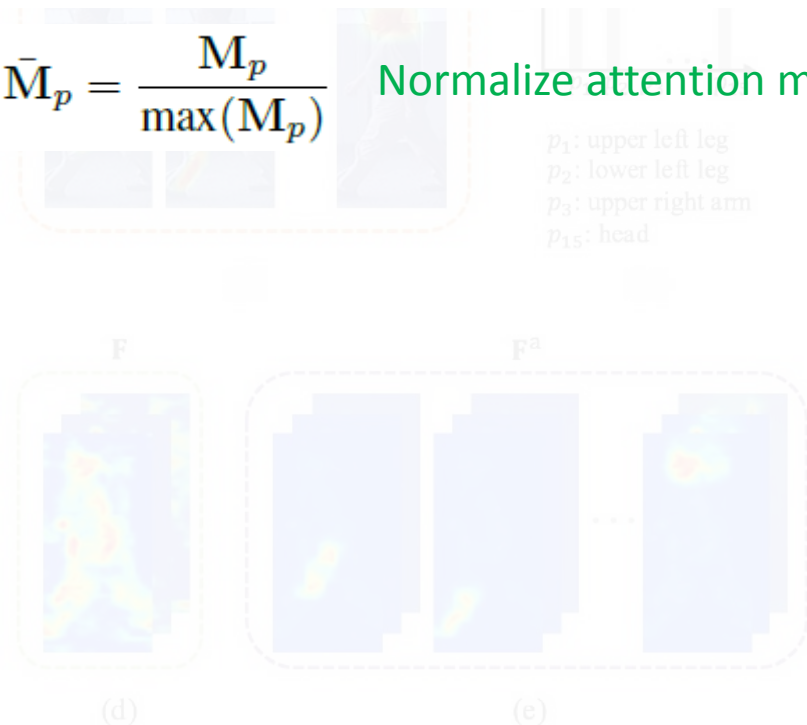


(a)

$M_p \in \{N_p, R_p\}$ is the attention map for body parts,

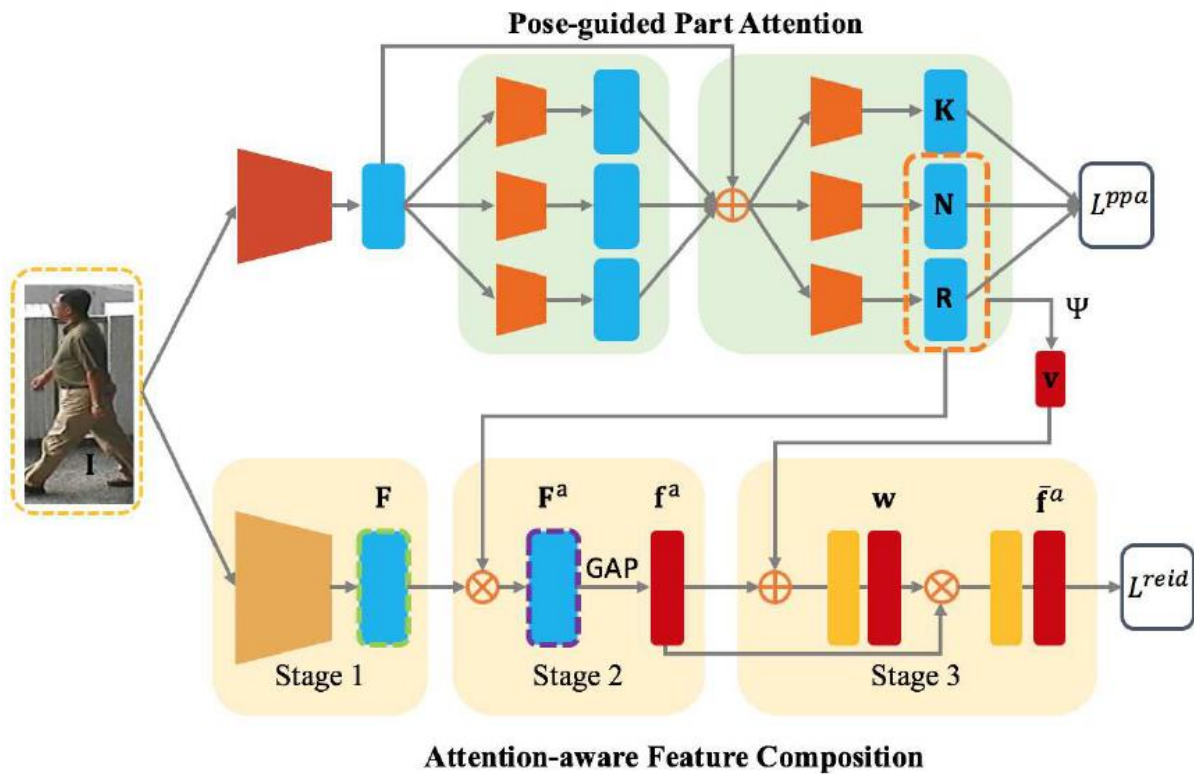
$$\bar{M}_p = \frac{M_p}{\max(M_p)}$$

Normalize attention map



Architecture

- Step 1: Train PPA independently
- Step 2: Train GCN independently
- Step 3: Attention-Aware Feature Alignment

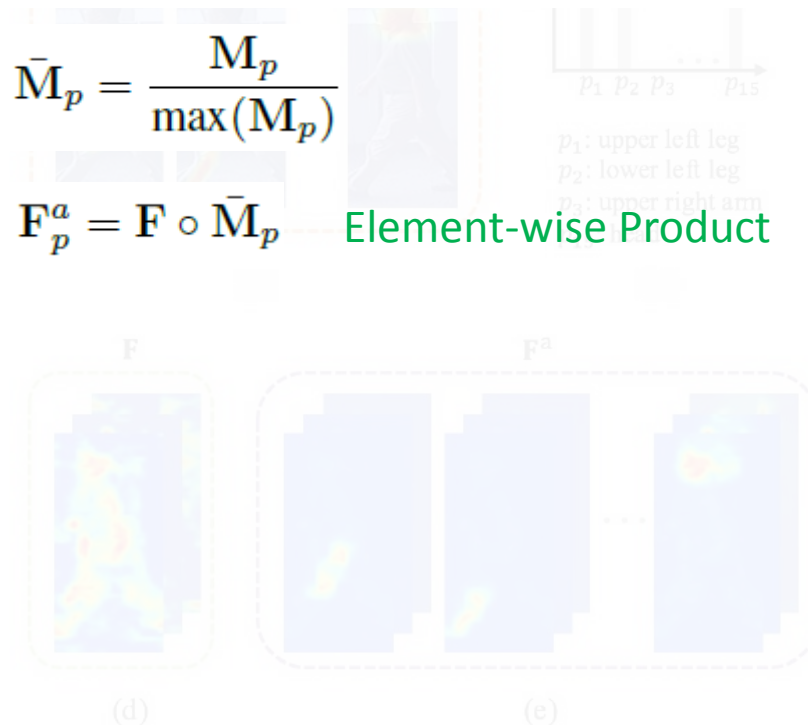


(a)

$M_p \in \{N_p, R_p\}$ is the attention map for body parts,

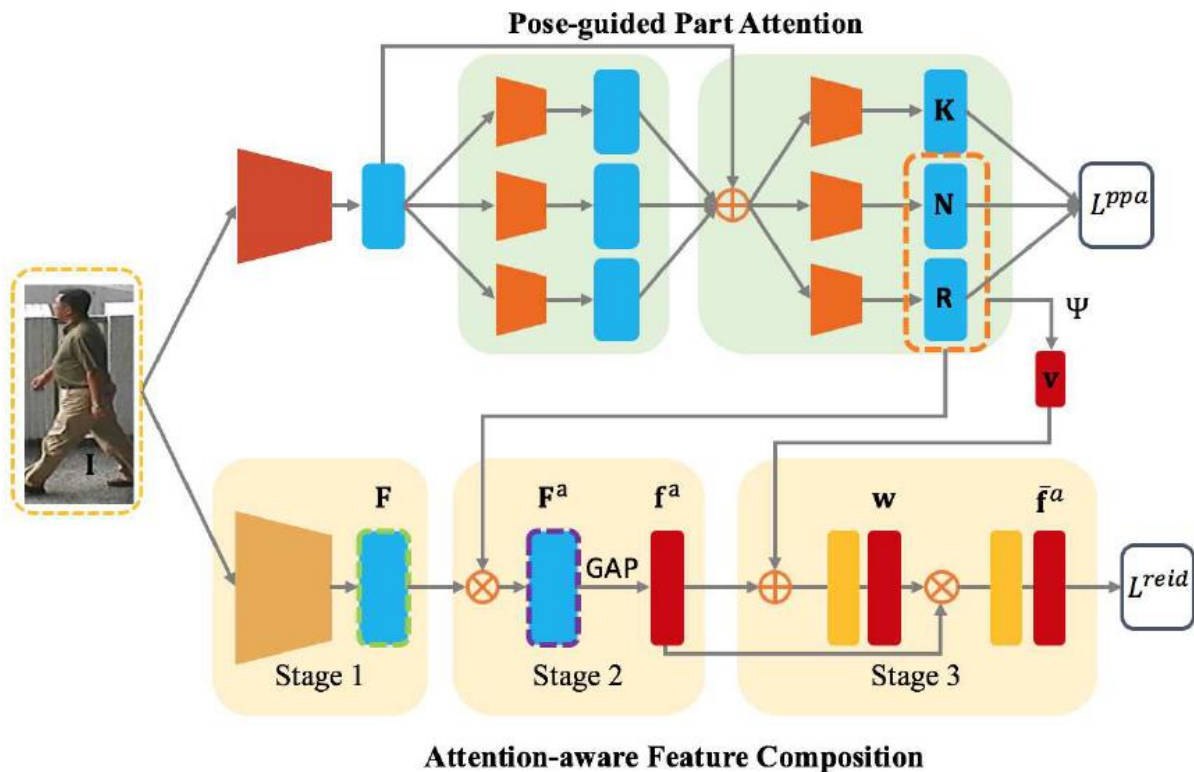
$$\bar{M}_p = \frac{M_p}{\max(M_p)}$$

$F_p^a = F \circ \bar{M}_p$ — Element-wise Product



Architecture

- Step 1: Train PPA independently
- Step 2: Train GCN independently
- Step 3: Attention-Aware Feature Alignment



(a)

$M_p \in \{N_p, R_p\}$ is the attention map for body parts,

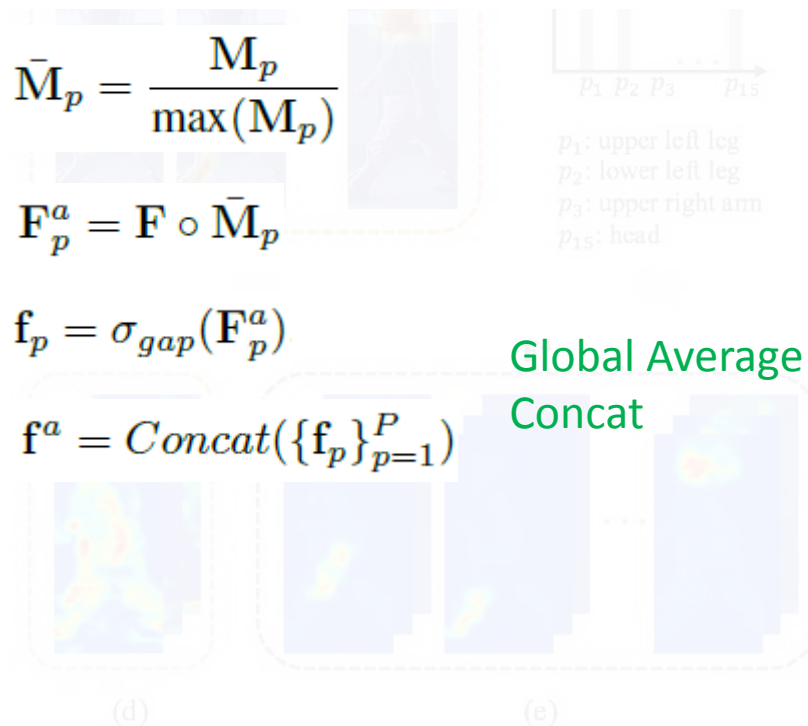
$$\bar{M}_p = \frac{M_p}{\max(M_p)}$$

$$F_p^a = F \circ \bar{M}_p$$

$$f_p = \sigma_{gap}(F_p^a)$$

$$f^a = \text{Concat}(\{f_p\}_{p=1}^P)$$

Global Average Pooling +
Concat



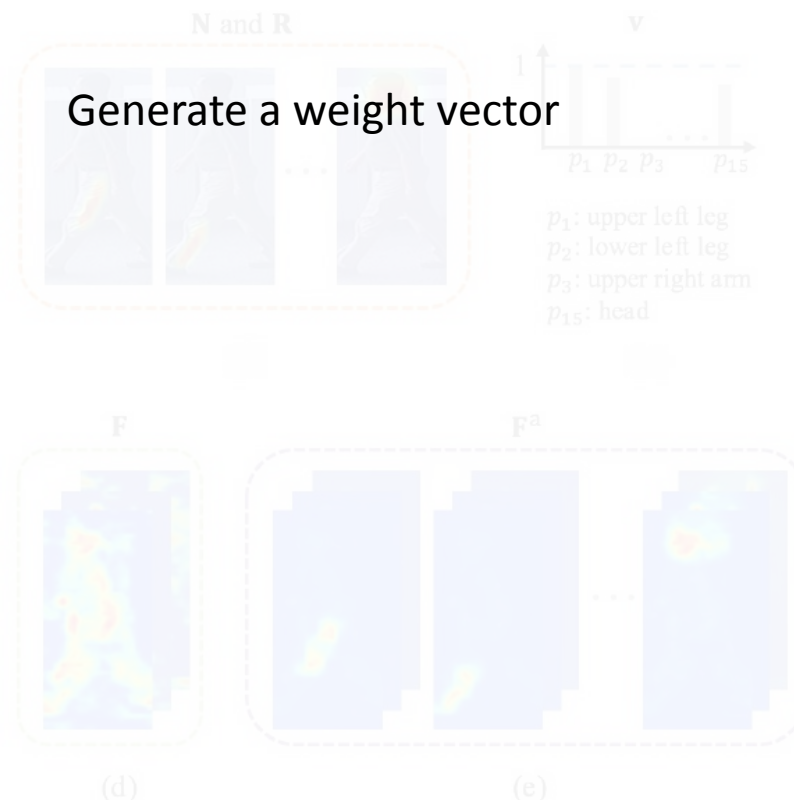
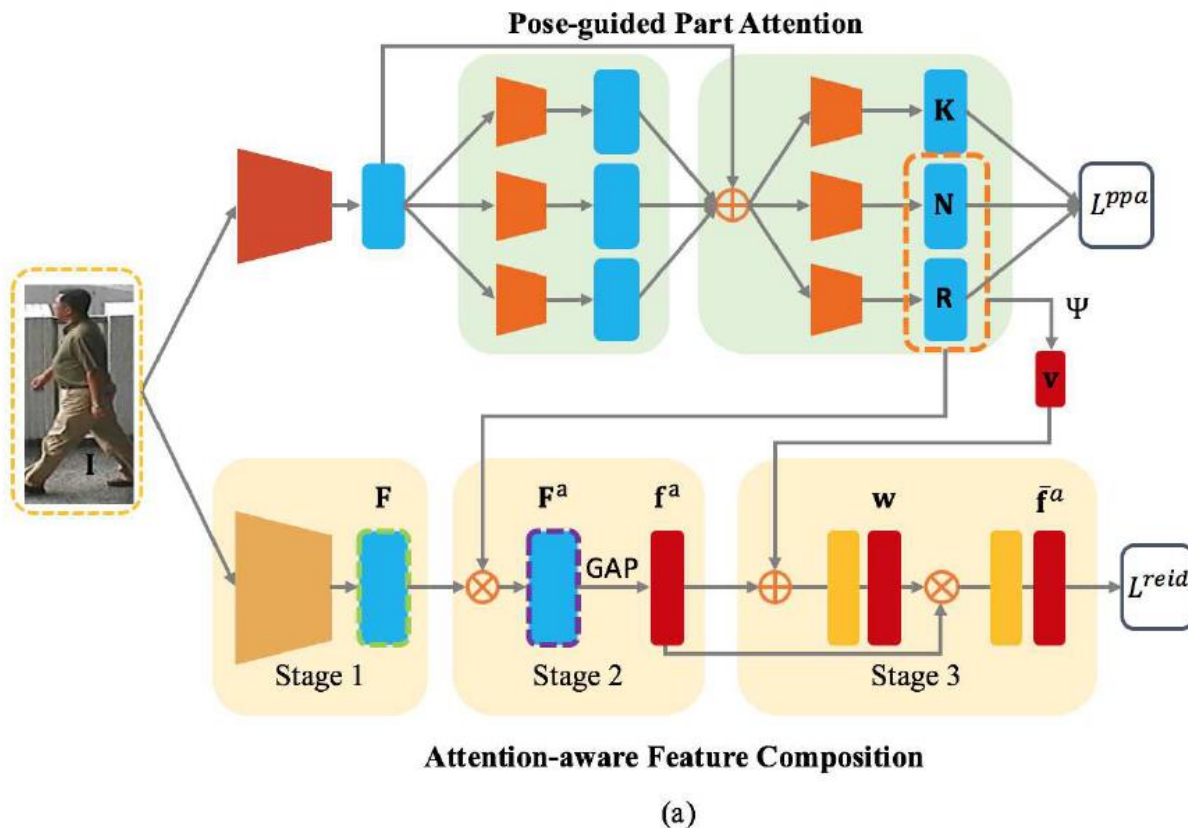
(d)

(e)



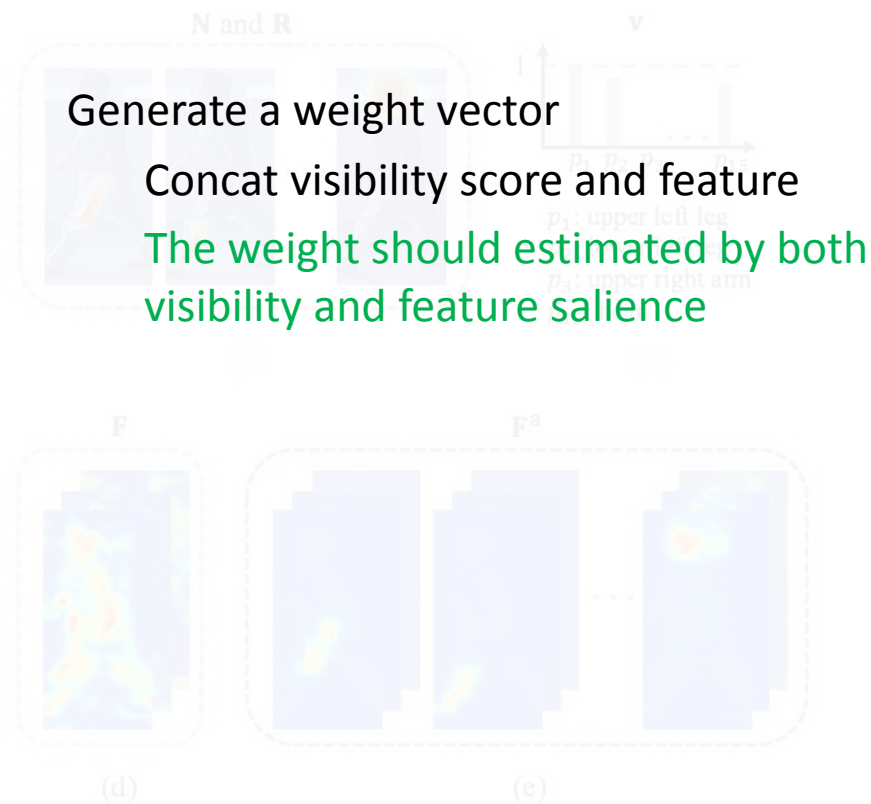
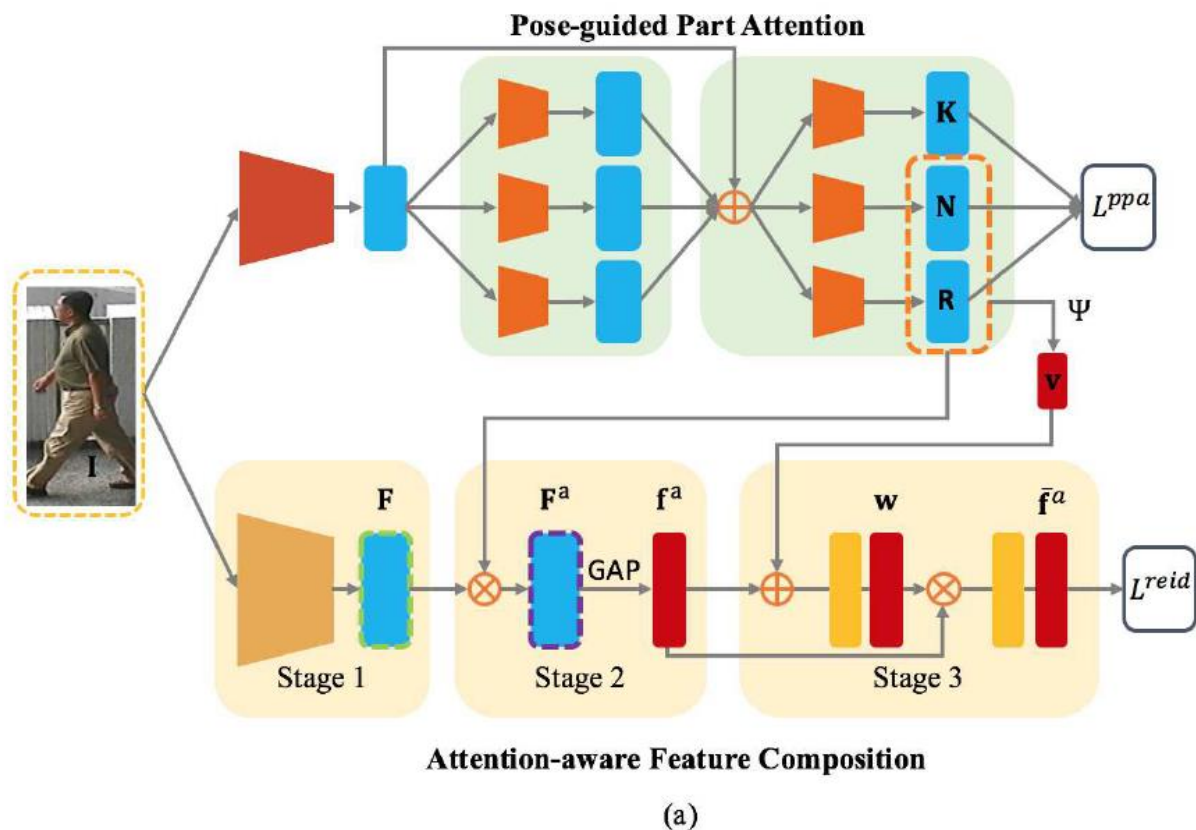
Architecture

- Step 1: Train PPA independently
- Step 2: Train GCN independently
- Step 3: Attention-Aware Feature Alignment
- Step 4: Weighted Feature Composition



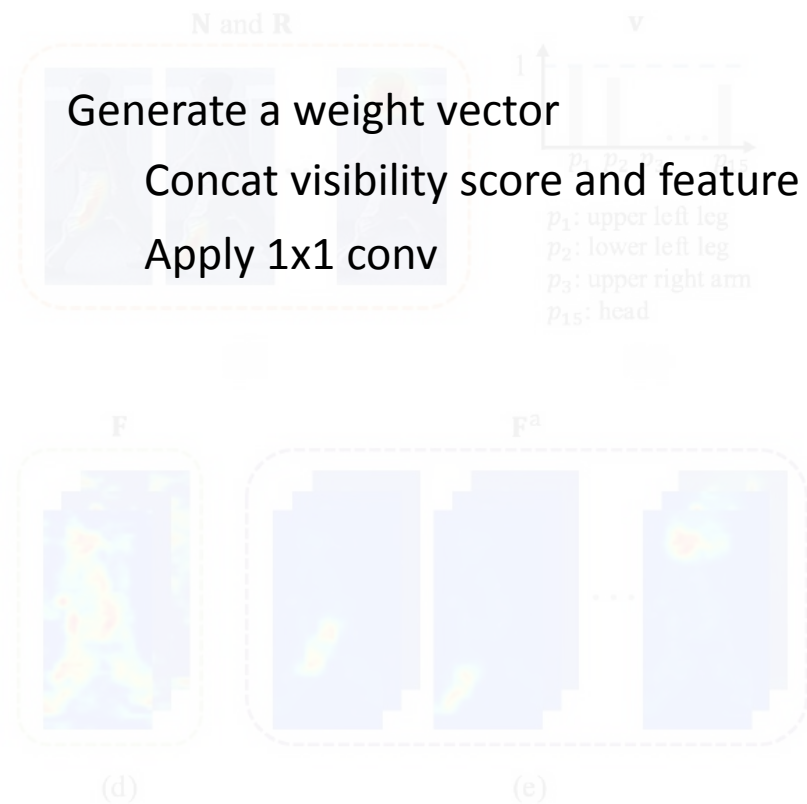
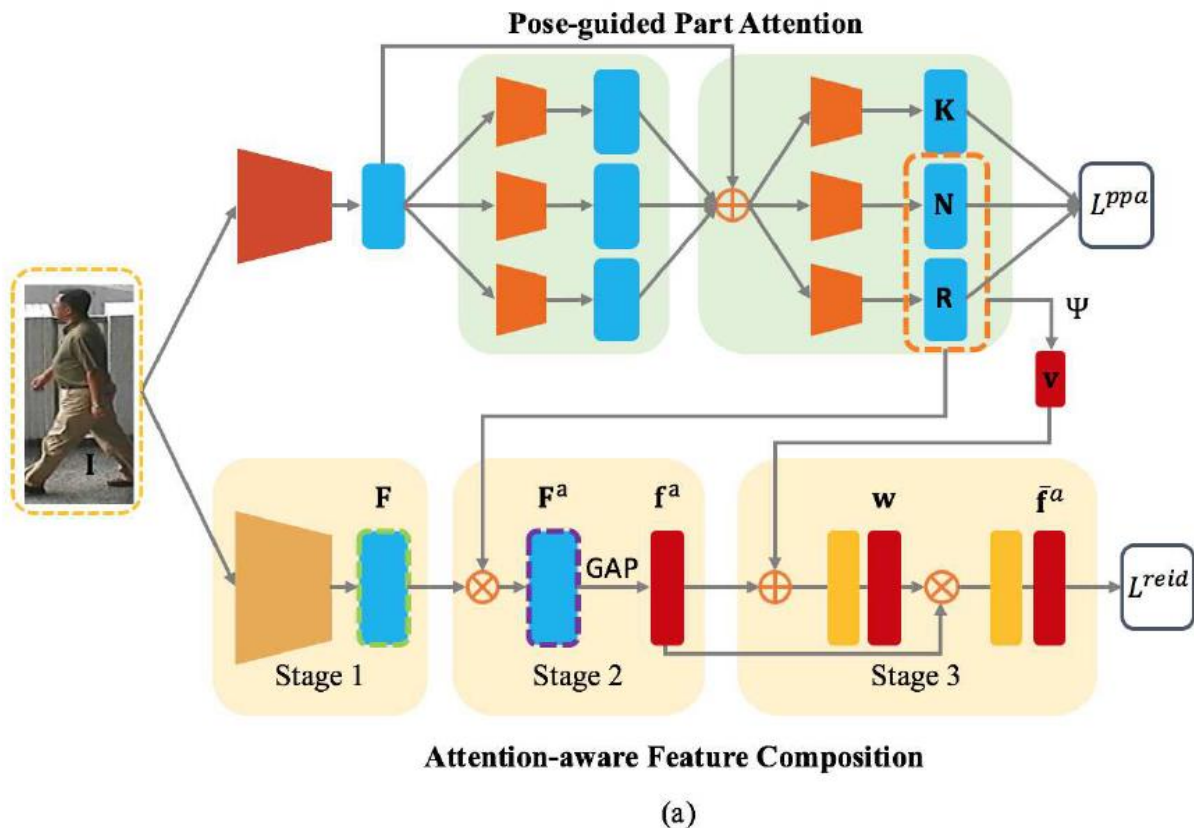
Architecture

- Step 1: Train PPA independently
- Step 2: Train GCN independently
- Step 3: Attention-Aware Feature Alignment
- Step 4: Weighted Feature Composition



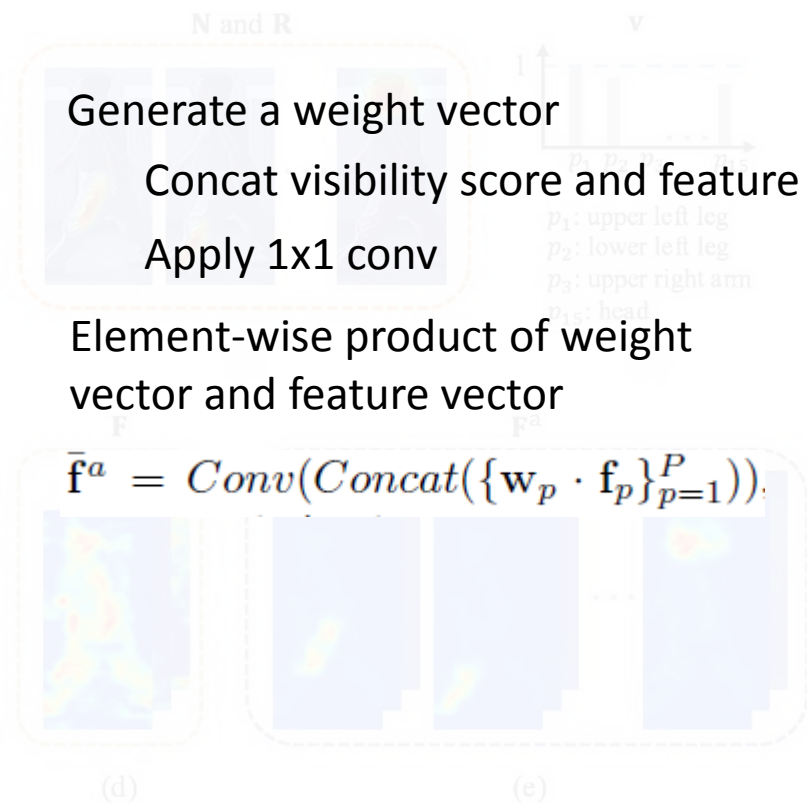
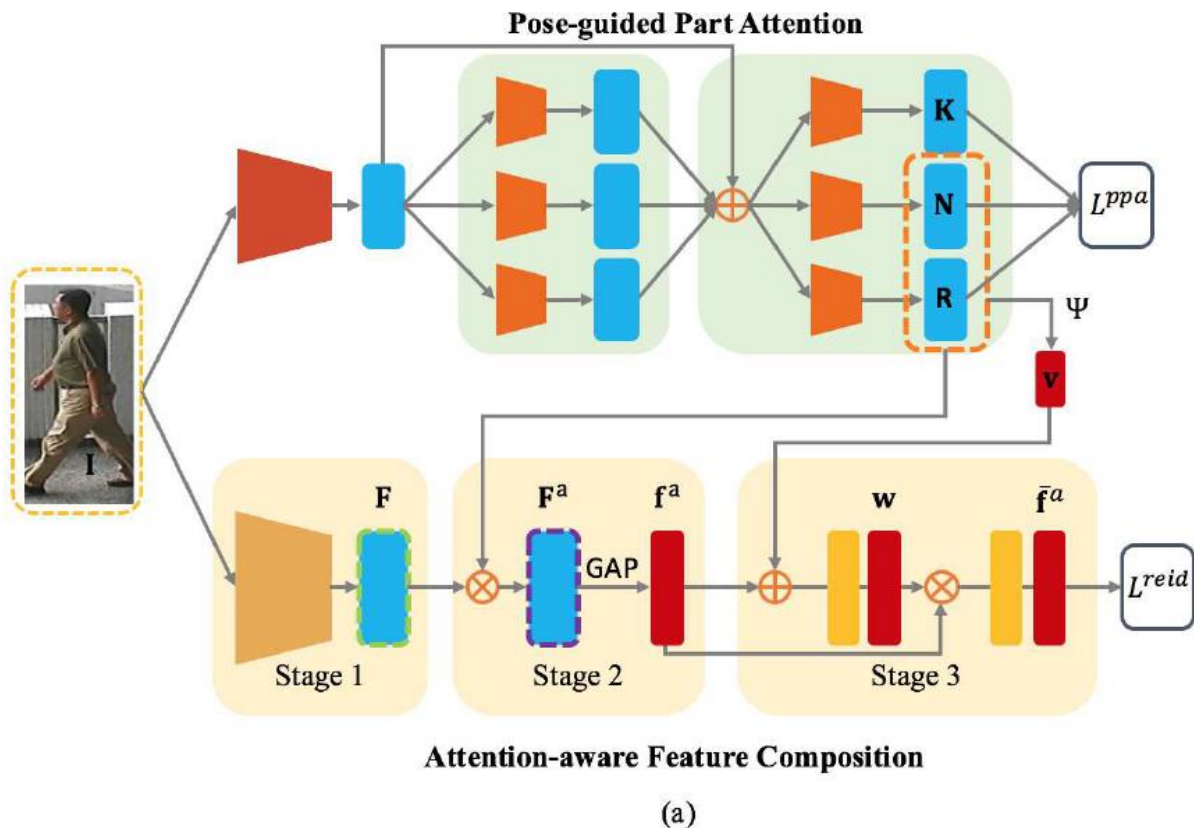
Architecture

- Step 1: Train PPA independently
- Step 2: Train GCN independently
- Step 3: Attention-Aware Feature Alignment
- Step 4: Weighted Feature Composition

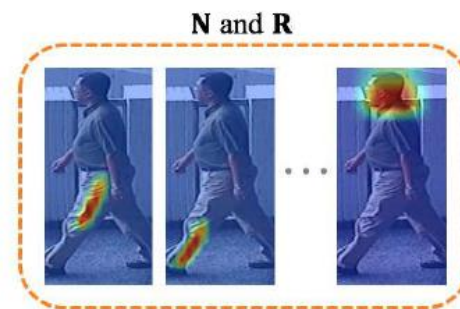
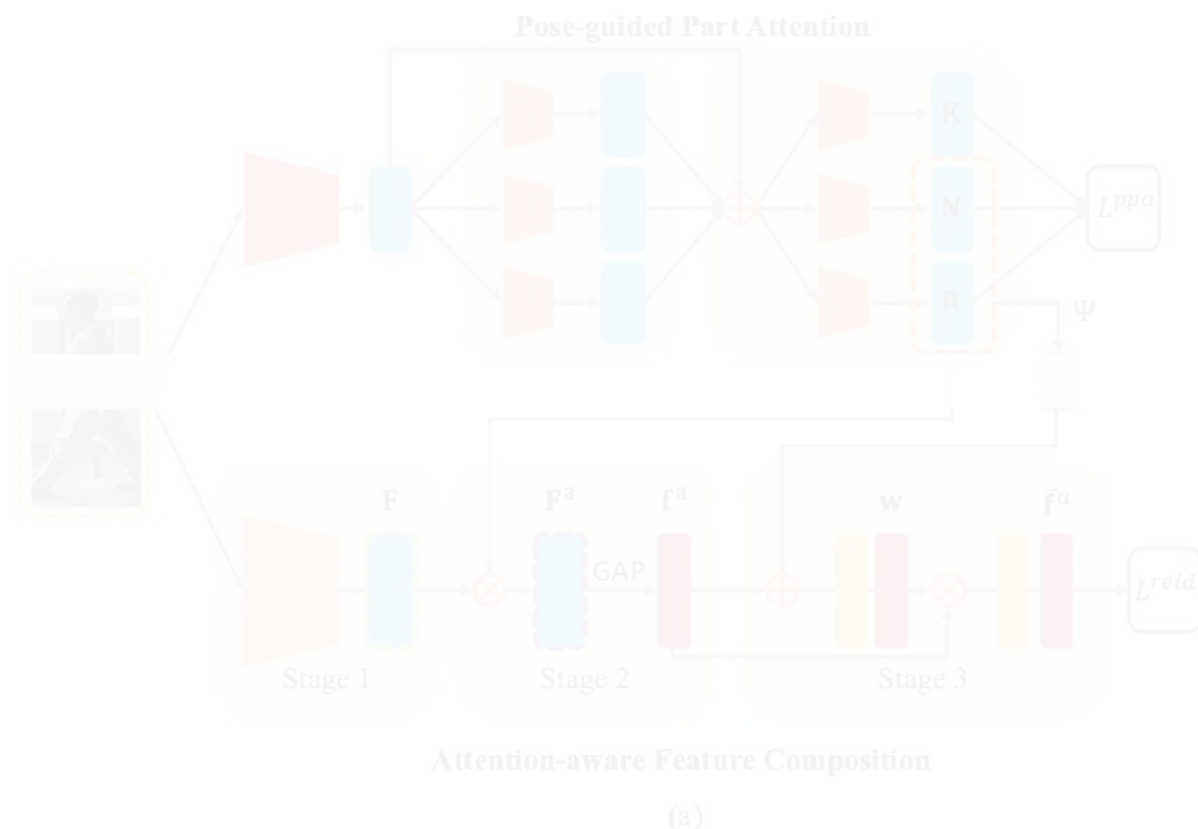


Architecture

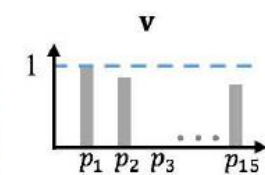
- Step 1: Train PPA independently
- Step 2: Train GCN independently
- Step 3: Attention-Aware Feature Alignment
- Step 4: Weighted Feature Composition



Response Maps

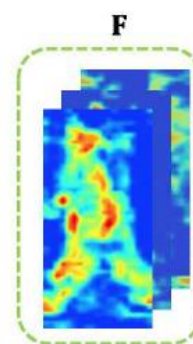


(b)

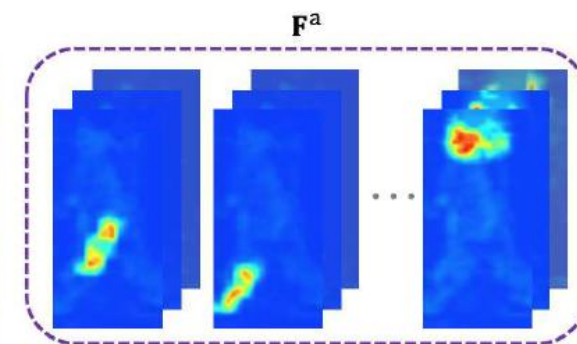


p_1 : upper left leg
 p_2 : lower left leg
 p_3 : upper right arm
 p_{15} : head

(c)



(d)



(e)



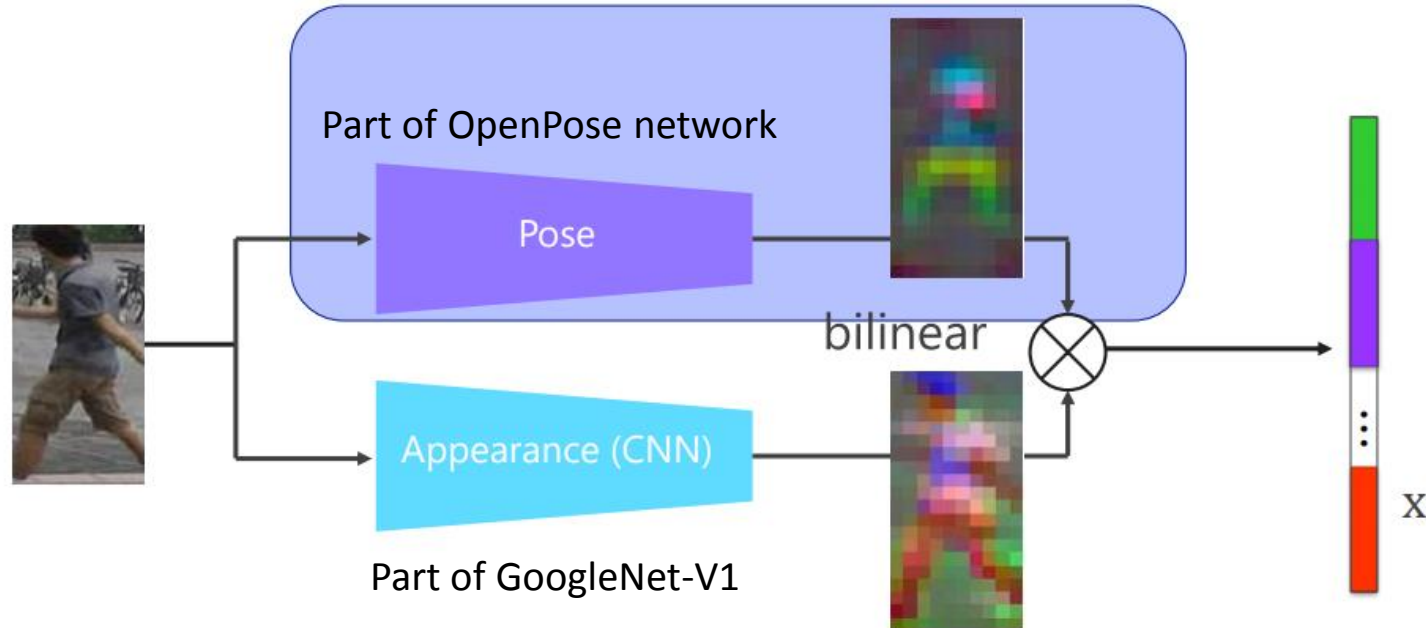
Overview

- Good idea, this can avoid the noise in the RoI methods
- Too complex...

Related Works

- An Improved Deep Learning Architecture for Person Re-Identification, w/o pose
CVPR 2015
- Deeply-Learned Part-Aligned Representations for Person Re-Identification, ICCV 2017
- Attention-Aware Compositional Network for Person Re-Identification, w/ pose
CVPR 2018
- Part-Aligned Bilinear Representations for Person Re-Identification,
ECCV 2018

Architecture



Each subvector in \mathbf{x} corresponds to a key point:

$$\text{vec}(\mathbf{a} \otimes \mathbf{p}) = [(p_1 \mathbf{a})^\top \ (p_2 \mathbf{a})^\top \ \dots \ (p_{c_P} \mathbf{a})^\top]^\top$$

Pose estimator *pretrained on COCO*, and re-trained *only with the re-id loss*

Response Maps



Left: Image, middle: appearance map, right: part map

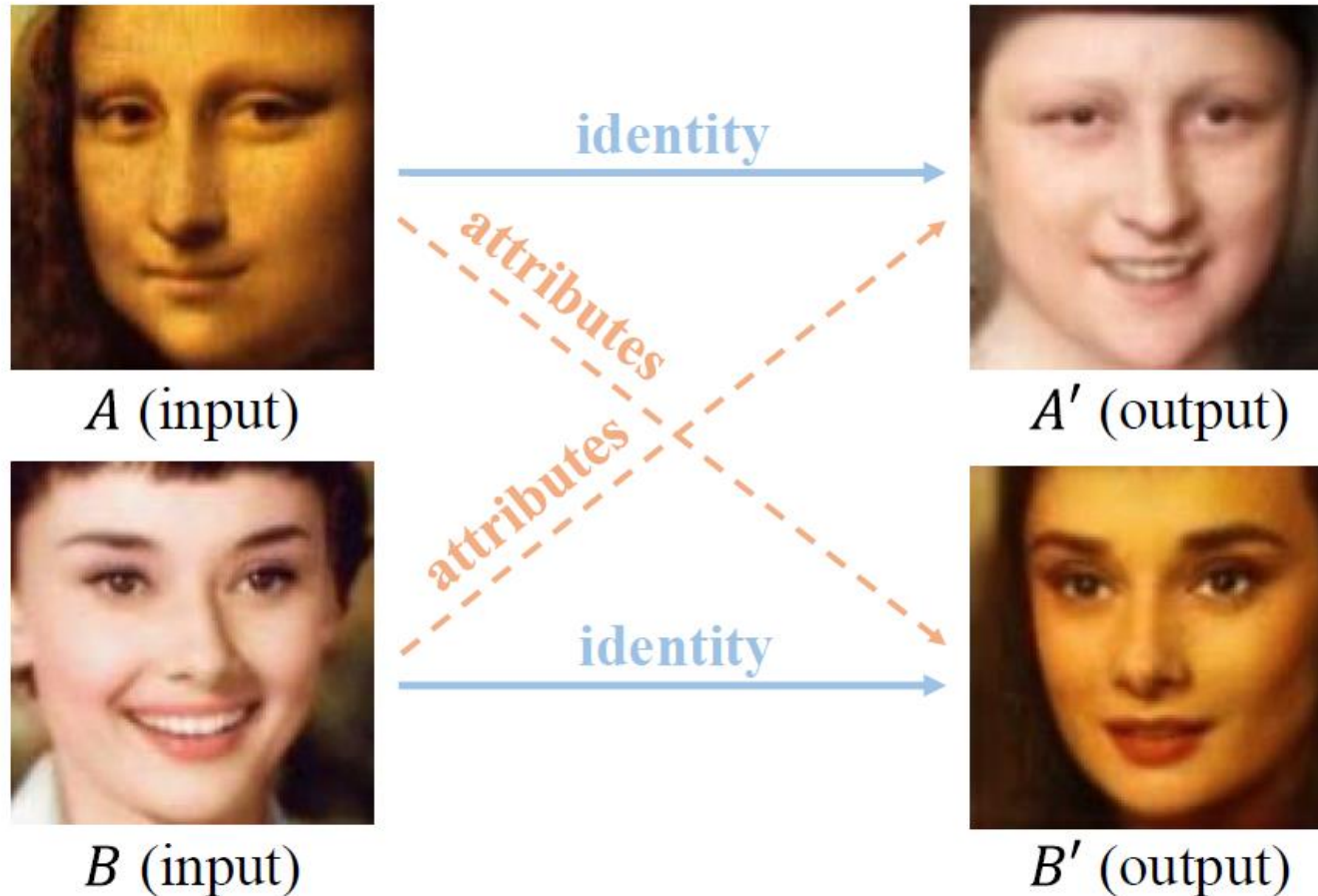
Overview

- Simple, effective
- Not take view changes into consideration

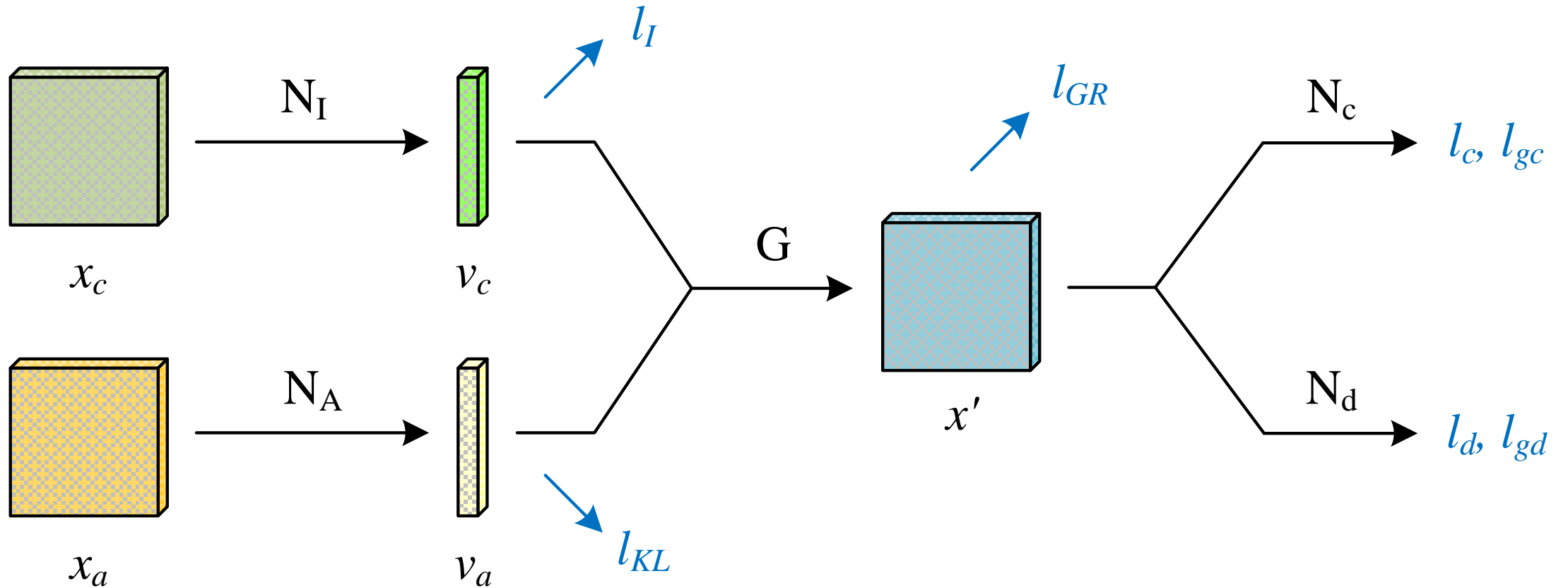
GANs methods

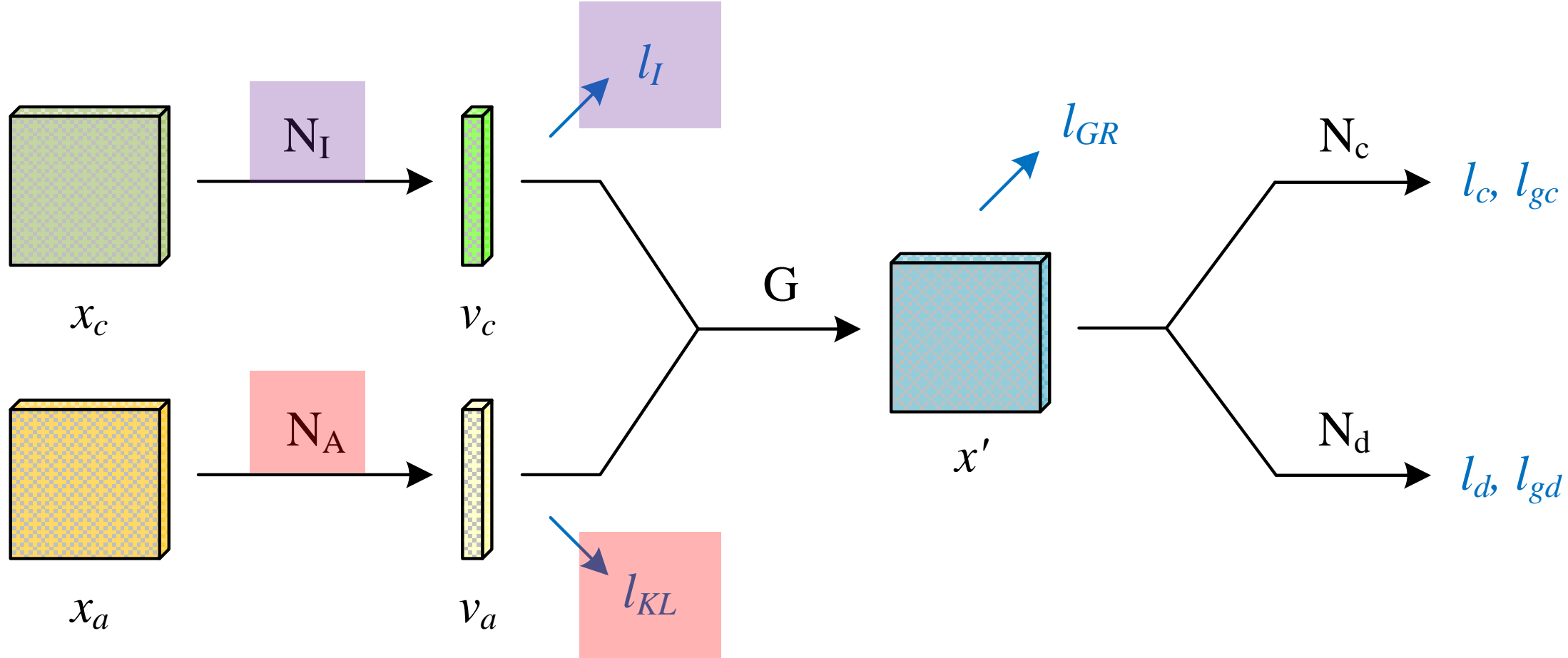
- IP-GAN for face synthesis, CVPR 2018
- PN-GAN for Re-ID, ECCV 2018
- FD-GAN for Re-ID, NIPS 2018

Identity Preserving GAN



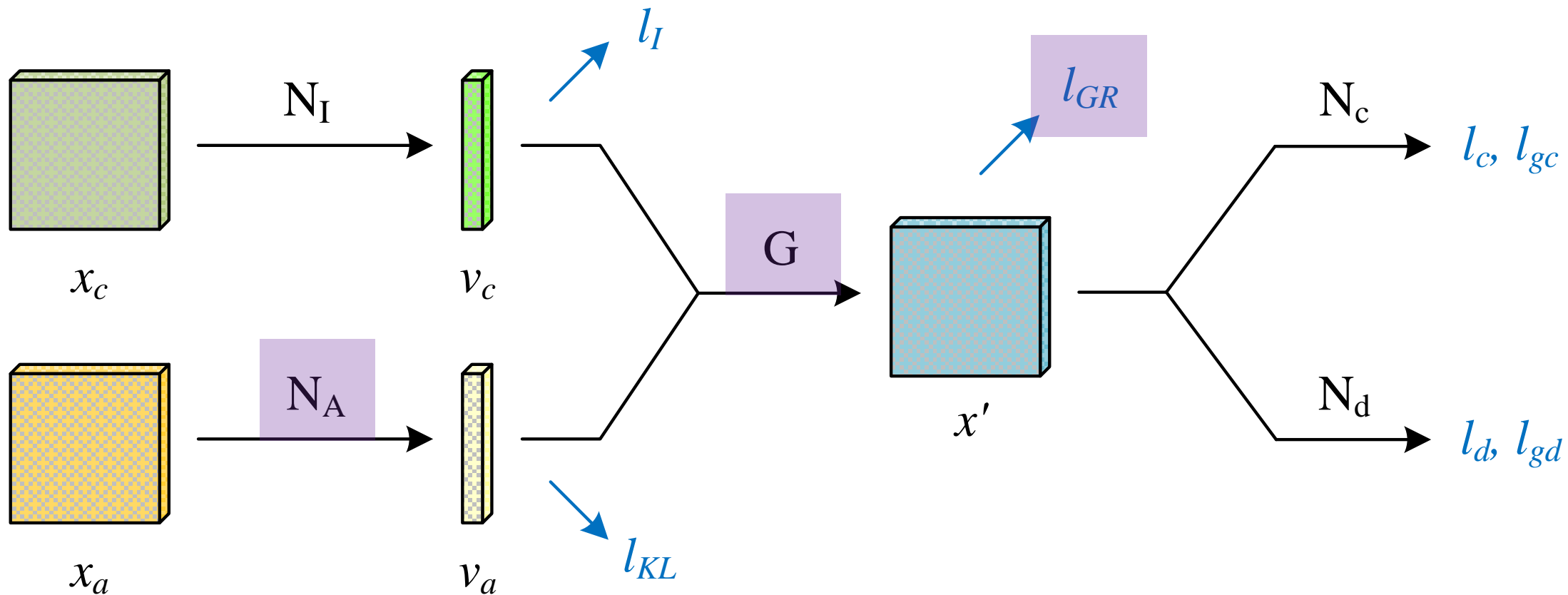
Identity Preserving GAN - Architecture



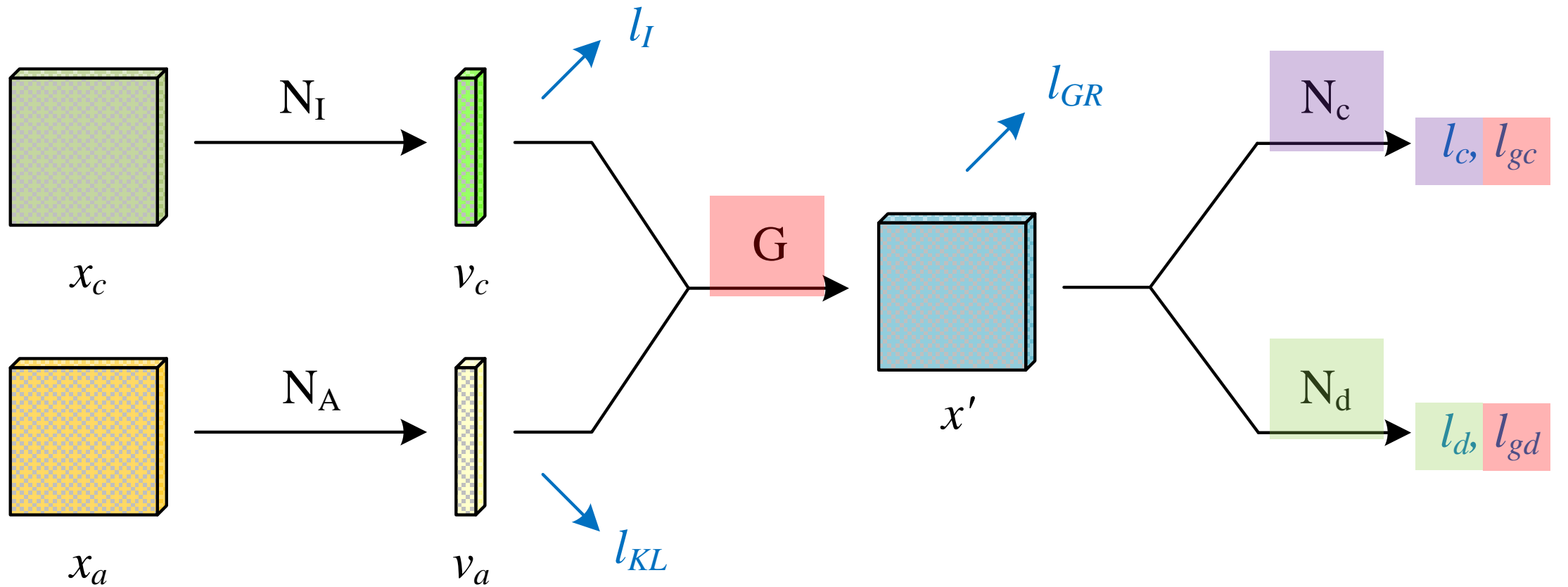


$$l_I = -E[\log P(c|x_c)]$$

$$l_{KL} = \frac{1}{2} (\mu^T \mu + \sum_{j=1}^J (e^{\epsilon_j} - \epsilon_j - 1))$$



$$l_{GR} = \begin{cases} \frac{1}{2} \| (x_a) - (x') \|^2, & x_a = x_c \\ \frac{0.1}{2} \| (x_a) - (x') \|^2, & otherwise \end{cases}$$



$$l_c = -E[\log P(c|x_c)]$$

$$l_{gc} = \frac{1}{2} \| f_c(x_c) - f_c(x') \|^2$$

$$l_d = -E[\log D(x_a)] - E[\log(1 - D(x'))]$$

$$l_{gd} = \frac{1}{2} \| f_d(x_a) - f_d(x') \|^2$$

Identity Preserving GAN - Overview

- Objective: Generate faces with diff. id & attribute
- Weakness: Only can generate faces with known id
- Strength: g.t. for x' is not needed
- Training set: (x_c, x_a) only require x_c id is known
- Identity & attribute feature decoupled

Content

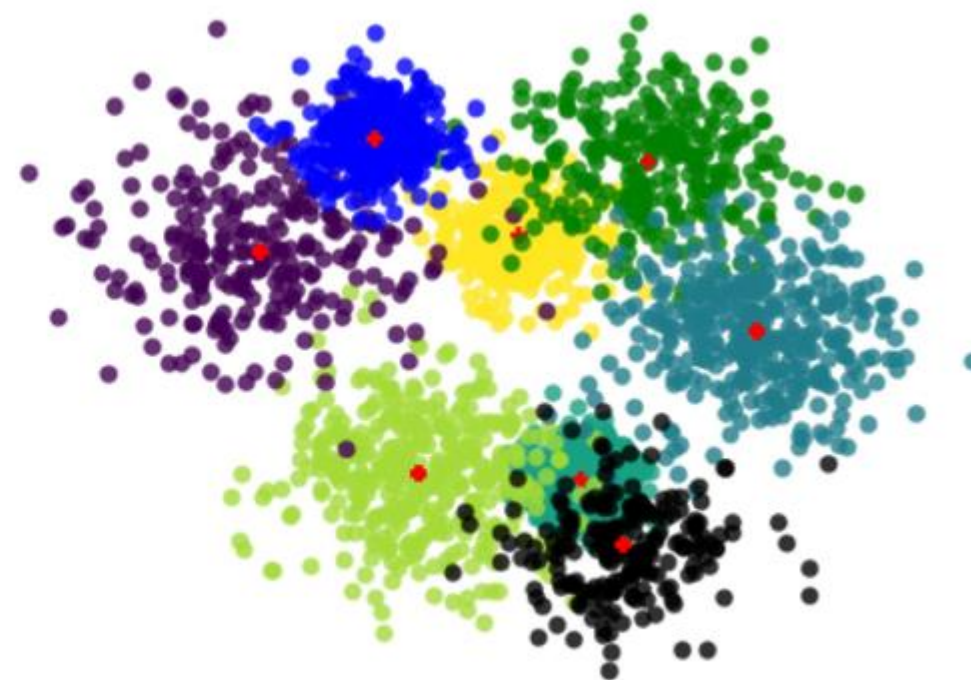
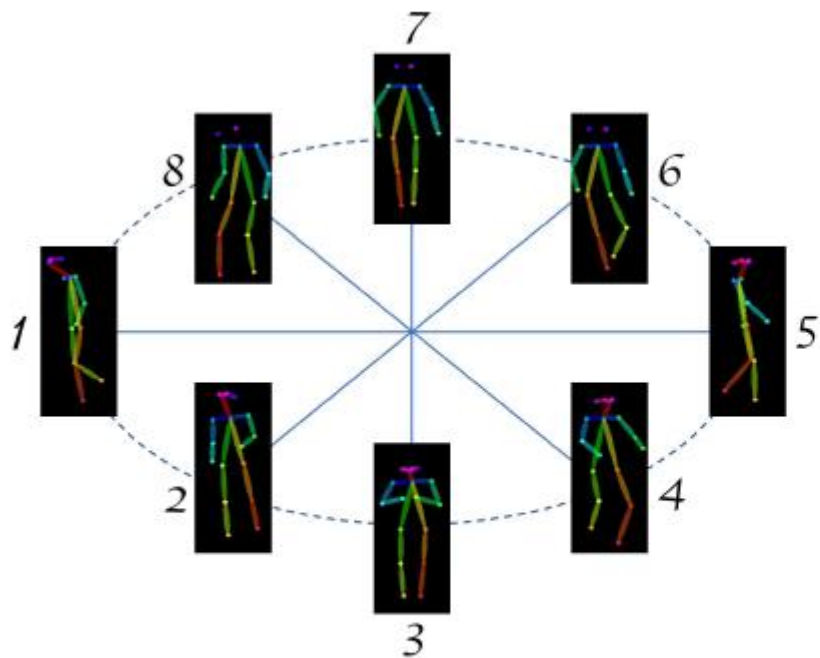
- Part-aligned Representation Learning
- **Image generation in Re-ID**
- GAN as supervisor

Pose Normalized GAN - Motivation

- Identity-Sensitive View-Insensitive (ISVI) features are needed
- → Part models
 - Lack of scalability
 - Lack of generalizability
 - ISVI features and IIVS features are not independent
 - $p(x, y) \Rightarrow p(x)$? Not easy
 - But we can have... $p(x, y) \Rightarrow p(x, y_1), p(x, y_2), p(x, y_3) \dots$

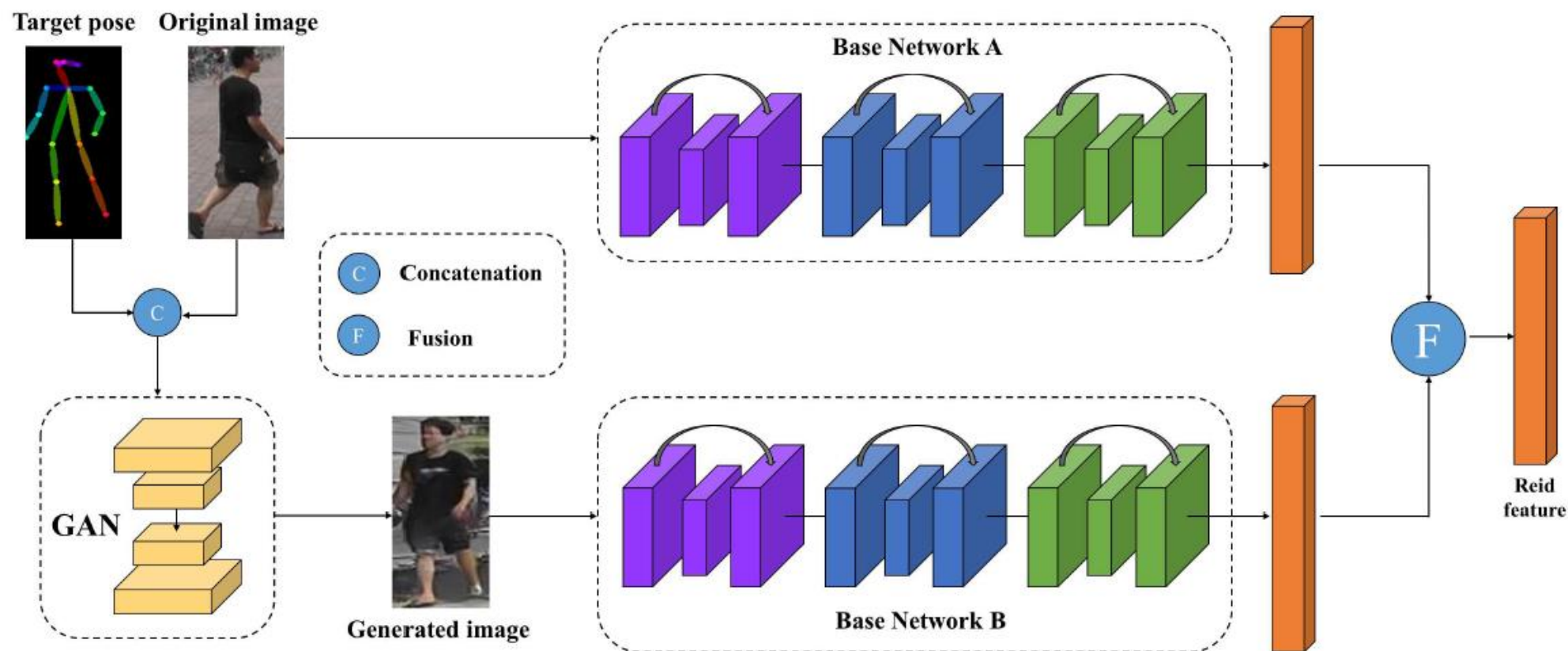
Pose Normalized GAN - Model

- At which poses we want to generate?

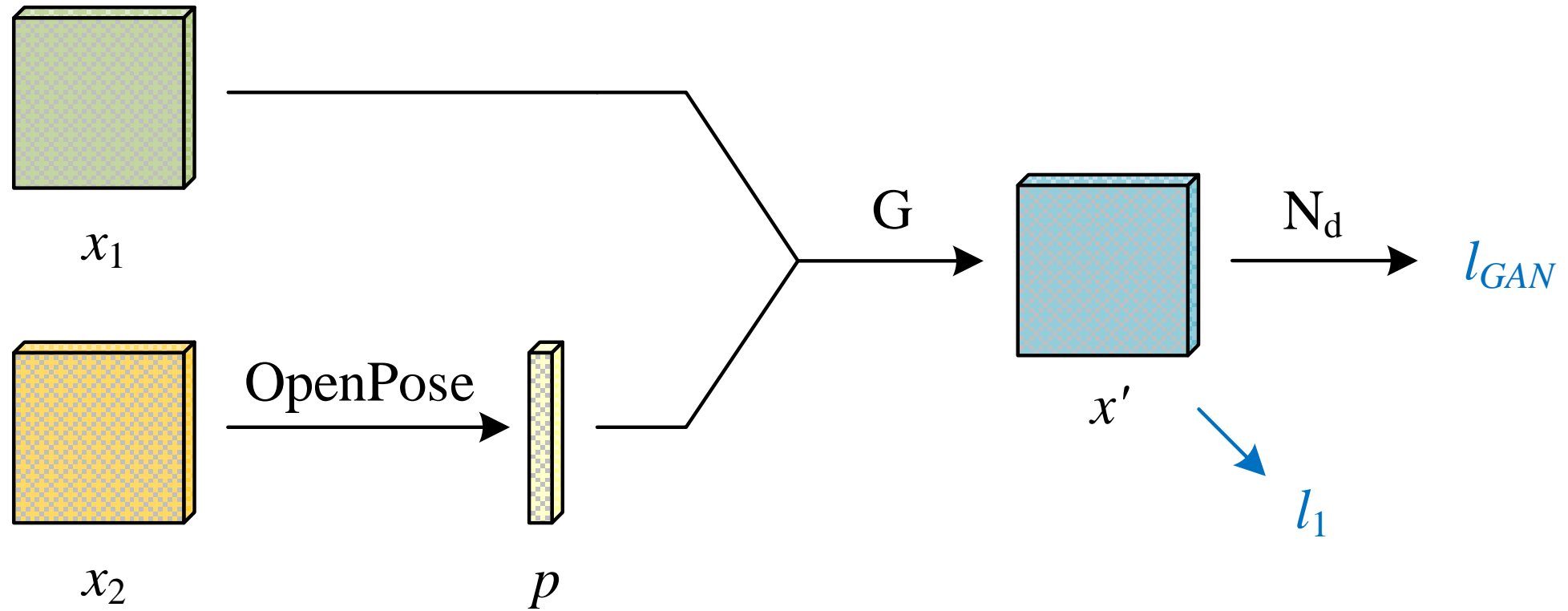


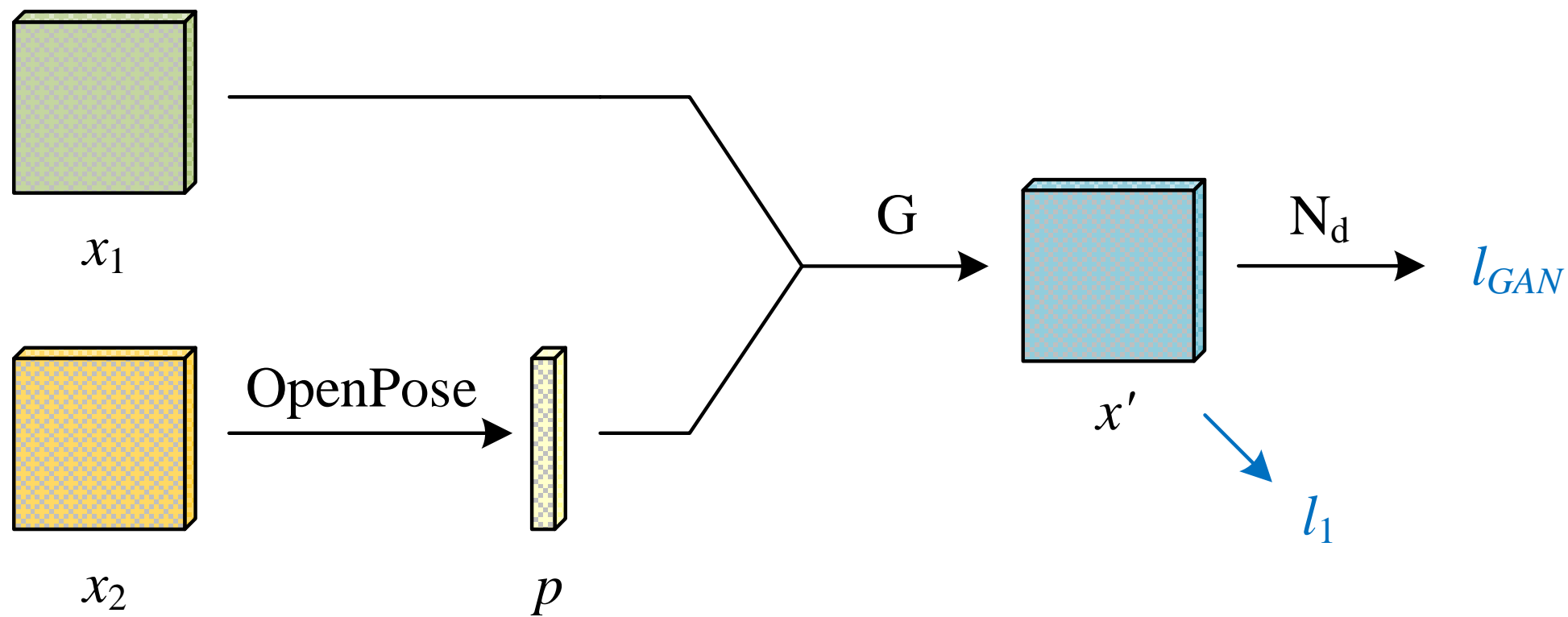
(a) Eight canonical poses on Market-1501 (b) t-SNE visualization of different poses.

Pose Normalized GAN - Architecture



Pose Normalized GAN - Architecture





$$l_1 = \|x_2 - x'\|_1$$

$$l_{GAN} = -E[\log D(x_2)] - E[\log(1 - D(x'))]$$

Pose Normalized GAN - Training steps

- Train PN-GAN
- Do Re-ID using pretrained PN-GAN

Pose Normalized GAN - Experimental Results

- Market-1501
 - Baseline (same structure without PN-GAN):
 - Top-1: 87.26
 - mAP: 69.32
 - PN-GAN
 - Top-1: 89.43
 - mAP: 72.58
- Performance boost is limited (+2 top-1, +3 mAP)
- Some part-models can achieve +10 mAP with similar setting

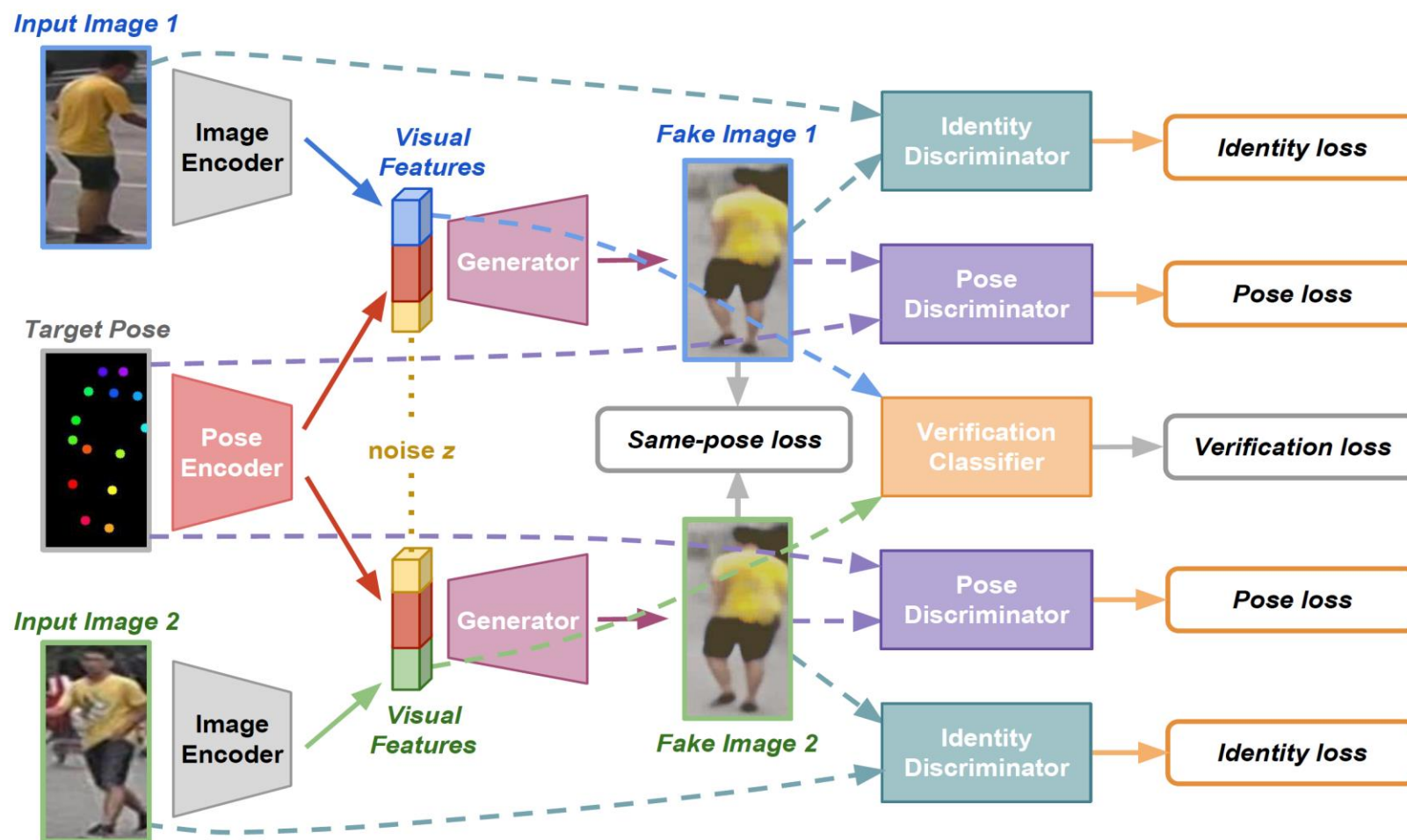
Pose Normalized GAN - Overview

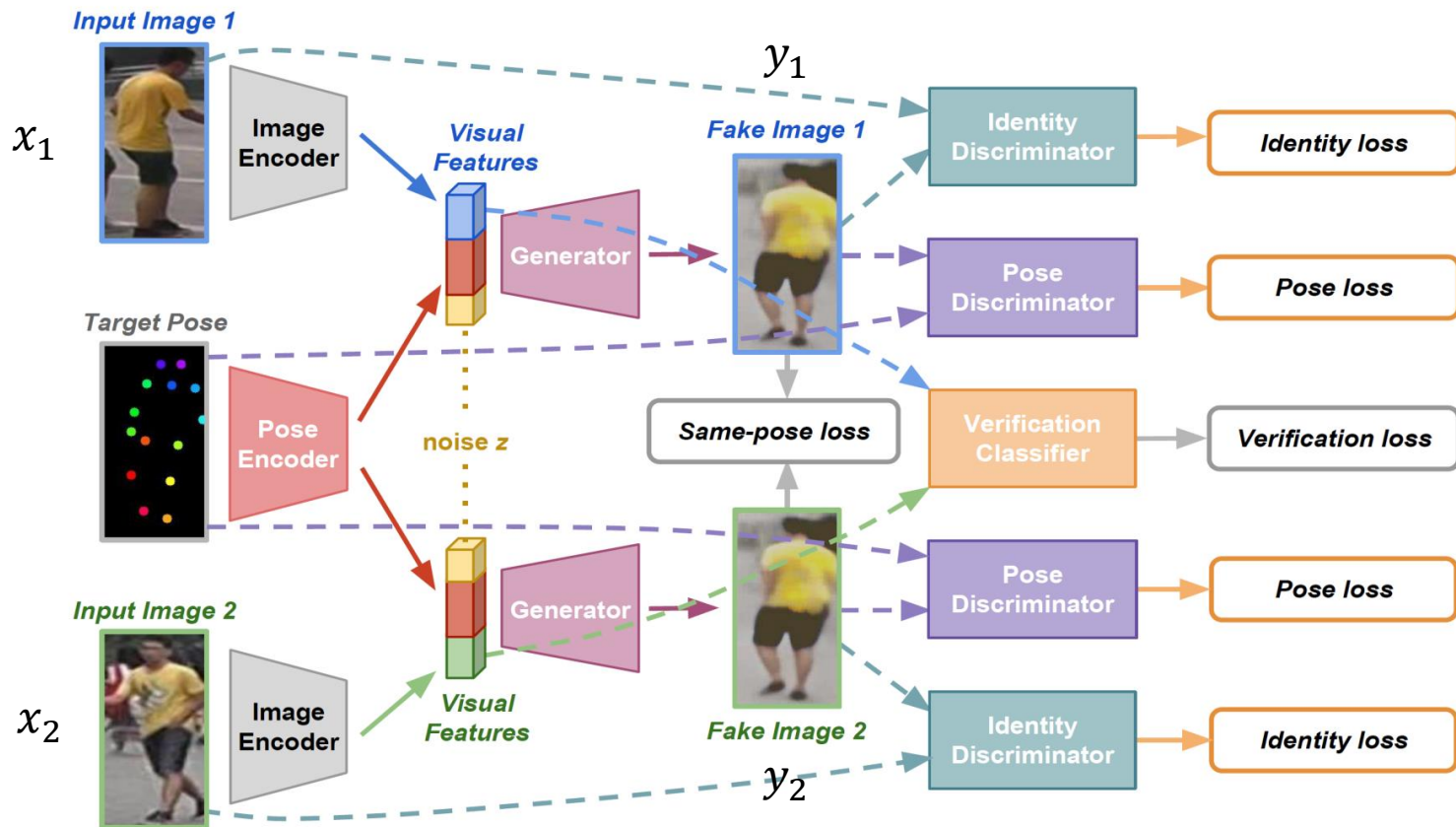
- Objective: Generate new images for better Re-ID
- Weakness: Small available data for training
 - x_1, x_2 required to be same person (g.t. for x' is needed)
- Implicitly decouple id & pose feature

Content

- Part-aligned Representation Learning
- Image generation in Re-ID
- **GAN as supervisor**

Feature Distilling GAN





$\mathcal{L}_v = -C \log d(x_1, x_2) - (1 - C)(1 - \log d(x_1, x_2))$, C is the ground-truth label

$$\mathcal{L}_{id} = \max_{D_{id}} \sum_{k=1}^2 \left(\mathbb{E}_{y'_k \in \mathcal{Y}} [\log D_{id}(x_k, y'_k)] + \mathbb{E}_{y_k \in \mathcal{Z}} [\log(1 - D_{id}(x_k, y_k))] \right)$$

$$\mathcal{L}_{pd} = \max_{D_{pd}} \sum_{k=1}^2 \left(\mathbb{E}_{y'_k \in \mathcal{Y}} [\log D_{pd}([p, y'_k])] + \mathbb{E}_{y_k \in \mathcal{Z}} [\log(1 - D_{pd}([p, y_k]))] \right)$$

$$\mathcal{L}_r = \sum_{k=1}^2 \frac{1}{mn} \|y_k - y'_k\|_1$$

$$\mathcal{L}_{sp} = \frac{1}{mn} \|y_1 - y_2\|_1$$

Feature Distilling GAN - Training steps

- Train feature encoder using Re-ID loss
- Train FD-GAN
- Global finetuning

Feature Distilling GAN - Experimental Results

- Market-1501
 - Baseline (same structure without FD-GAN):
 - Top-1: 88.2
 - mAP: 72.5
 - FD-GAN
 - Top-1: 90.5
 - mAP: 77.7
- Performance boost is limited (+2 top-1, +5 mAP)

Feature Distilling GAN - Overview

- Objective: Learn a better feature encoder
- Weakness: Small available data for training
 - x_1, x_2 required to be same person (g.t. for x' is needed)

Overview

- Rethinking previous models
 - IP-GAN (face)
 - Objective: generate better face with specific id & attribute
 - Pro: g.t. for generated image is not required
 - Con: Can only generate known identities
 - PN-GAN (Re-ID)
 - Objective: generate better image with specific id & pose
 - Pro: can generate unknown id images
 - Con: g.t. is required during training, pose is not precious
 - FD-GAN (Re-ID)
 - Objective: GAN for better feature encoder learning
 - Pro: can generate unknown id images
 - Con: g.t. is required during training

Overview

- Rethinking previous models
 - IP-GAN (face)
 - Objective: generate better face with specific id & attribute
 - Pro: g.t. for generated image is not required
 - Con: Can only generate known identities
 - PN-GAN (Re-ID)
 - Objective: generate better image with specific id & pose
 - Pro: can generate unknown id images
 - Con: g.t. is required during training, pose is not precious
 - FD-GAN (Re-ID)
 - Objective: GAN for better feature encoder learning
 - Pro: can generate unknown id images
 - Con: g.t. is required during training

Thanks

Q&A