



i-VisionGroup@Tsinghua

文献阅读

谭子萌

Pose estimation

- 单人姿态估计：检测关节点，常使用高斯热图回归方法
- 肢体变化、遮挡或背景中相似的物体 导致回归出准确的热图困难
- 主要考虑应用不同关键点之间的相互关系（结构先验）方面
 1. 特征层面融合（显式利用）
 2. GAN结构（隐式利用）



Structured Feature Learning for Pose Estimation

Xiao Chu Wanli Ouyang Hongsheng Li Xiaogang Wang

Department of Electronic Engineering, The Chinese University of Hong Kong

xchu@ee.cuhk.edu.hk

wlouyang@ee.cuhk.edu.hk

hsli@ee.cuhk.edu.hk

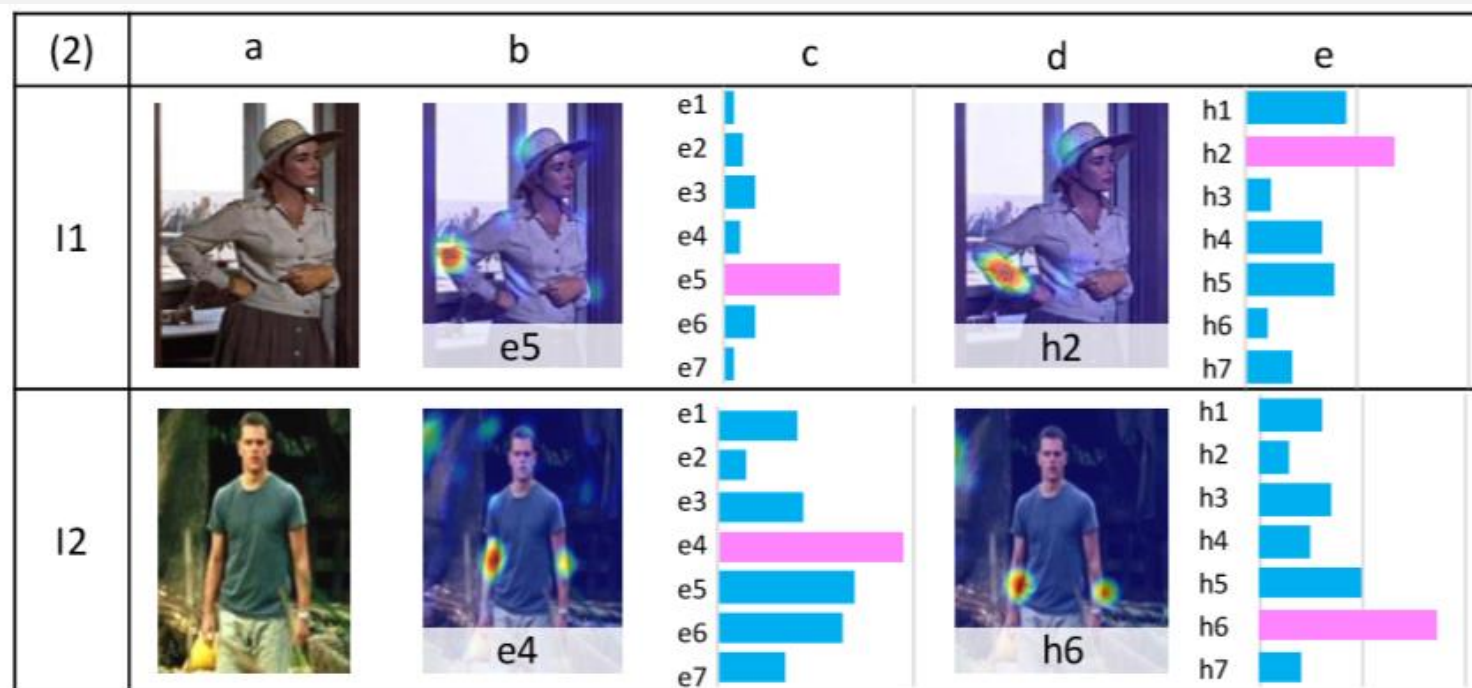
xgwang@ee.cuhk.edu.hk

2016CVPR

从特征层面利用关节点间的结构相关性

（特征层面能够保留除了位置信息以外，其他更为丰富的信息）

不同关节点的不同特征通道上可能具有相关性或反相关性：



Method

Contribution:

1. 信息传递的方式: geometrical transform kernels
2. 信息传递的路径: Bi-directional tree model
3. 特征层面, 实现端到端的训练, 而非后处理

以VGG为骨架, 共享fc6层 (4096channel)

其后fc7层每个关节点分别得到128个特征图

对第k个关节点 (x, y) 像素位置

$$\mathbf{h}_{fc7}^k(x, y) = f(\mathbf{h}_{fc6}(x, y) \otimes \mathbf{w}_{fc7}^k + \mathbf{b}_{fc6}),$$

Geometrical transform kernels

希望用小臂特征图 h_m 减少手肘 e_n 的错误响应、加强正确响应

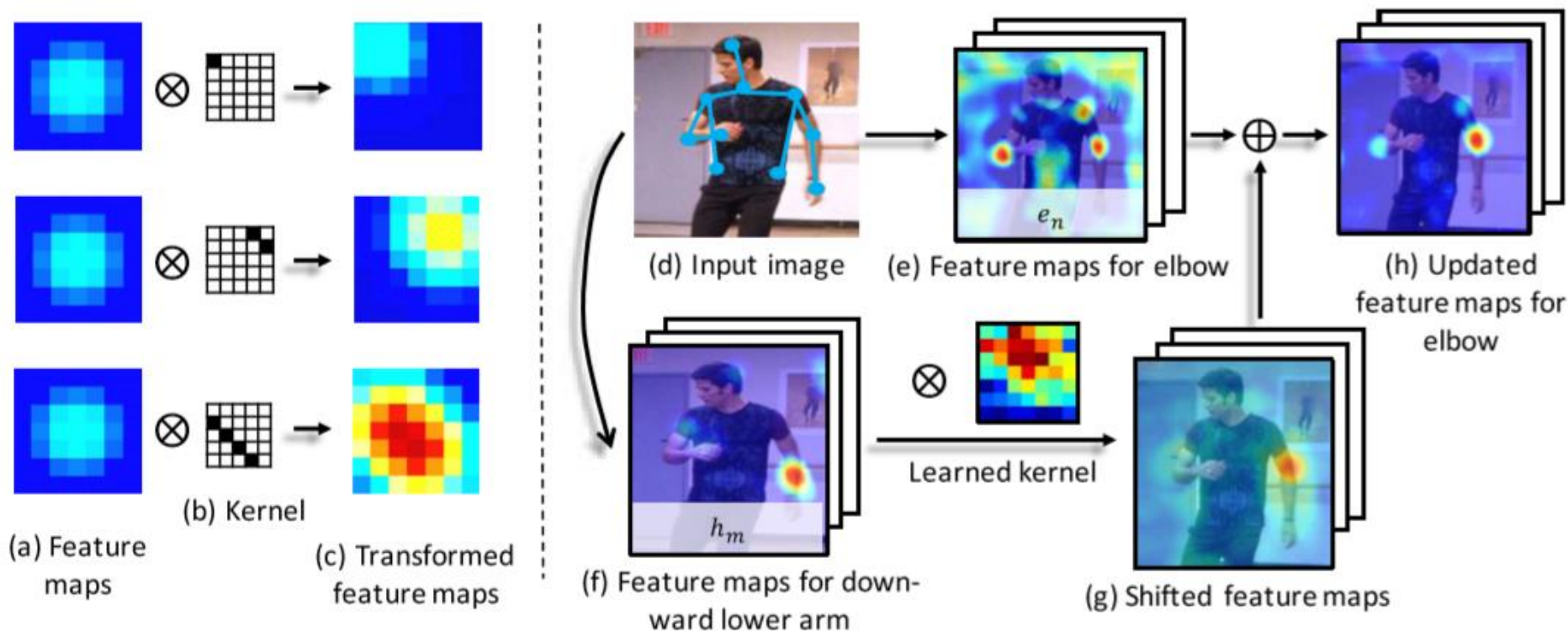
不能单纯地将 e_n 与 h_m 相加：存在空间的不匹配性

→ 不对称的卷积核能够对特征图响应产生几何变换

相邻关节点间的相对空间位置固定，使得几何变换核容易学习到

此外对负相关的情况，使核为负值来抑制错误的响应

连续3个 7×7 的变换核 满足变换尺寸需求



Bi-directional tree model

在那些距离较近、关系比较稳定的关节点传递信息

利用双方向来传递互补的信息：叶子节点 \leftrightarrow 根节点

A_k' 与 B_k'
concat 256

Score map

Refined fcn7
feature map 128

Origin fcn7 128

Shared fcn6 4096

Input image

CNN

1x1 convolution

(1) Part-features

(2) Structured feature learning

(3) Prediction

Joint 3

Joint 4

Joint 6

Joint 5

Score

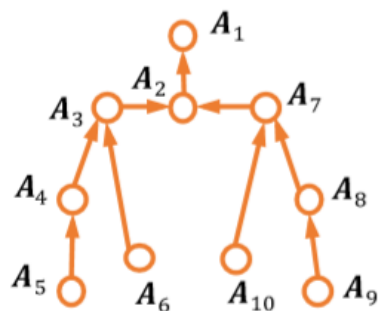
19

56

56

(2,a)

Upward Direction



$$A_6 = f(h_{fcn6} \otimes w^{a6})$$

$$A_6' = A_6$$

$$A_5 = f(h_{fcn6} \otimes w^{a5})$$

$$A_5' = A_5$$

$$A_4 = f(h_{fcn6} \otimes w^{a4})$$

$$A_4' = f(A_4 + A_5' \otimes w^{a5,a4})$$

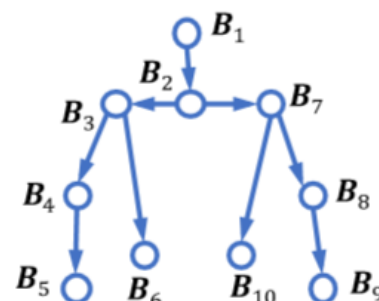
$$A_3 = f(h_{fcn6} \otimes w^{a3})$$

$$A_3' = f(A_3 + A_4' \otimes w^{a4,a3} + A_6' \otimes w^{a6,a3})$$

fcn7中节点3的卷积核

(2,b)

Downward Direction



A6到A3的几何变换核

Method

Score map: 1×1 卷积

$$\mathbf{z}_k = [\mathbf{A}'_k, \mathbf{B}'_k] \otimes \mathbf{w}_{pred}^k.$$

Loss: 分类问题 (18个关节点+1个背景)

当 (x, y) 像素属于第 k 类时, $t_k(x, y)=1$, 否则=0

为解决类别不均衡问题, 引入binary mask来随机采样0.05%的负样本

$$\sum_x \sum_y m(x, y) \sum_k t_k(x, y) \log\left(\frac{e^{z_k(x, y)}}{\sum_{k'} e^{z_{k'}(x, y)}}\right)$$

后处理: $[(dx)^2, (dy)^2]$

$$\begin{aligned} dx &= (x_i - x_j - x_r) \\ dy &= (y_i - y_j - y_r) \end{aligned}$$


Test image



Score map

Results



Adversarial PoseNet:

A Structure-aware Convolutional Network for Human Pose Estimation*

Yu Chen¹ Chunhua Shen² Xiu-Shen Wei³ Lingqiao Liu² Jian Yang¹

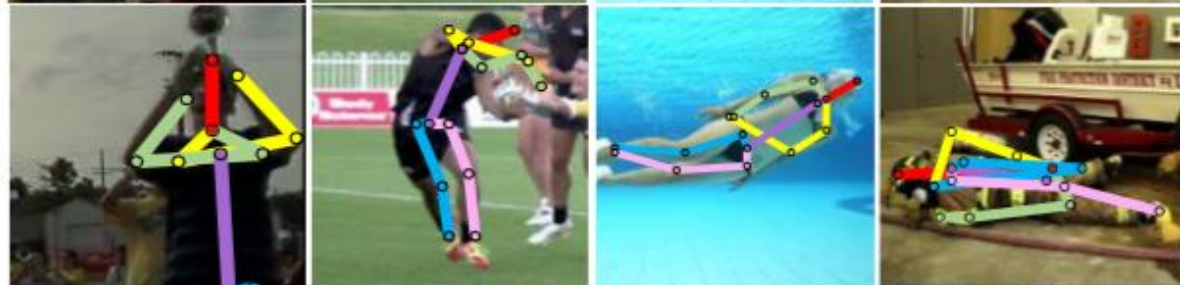
¹Nanjing University of Science and Technology ²University of Adelaide ³Nanjing University 2017 ICCV

单个关节的热图估计可能导致出现生物学上不可能的姿态

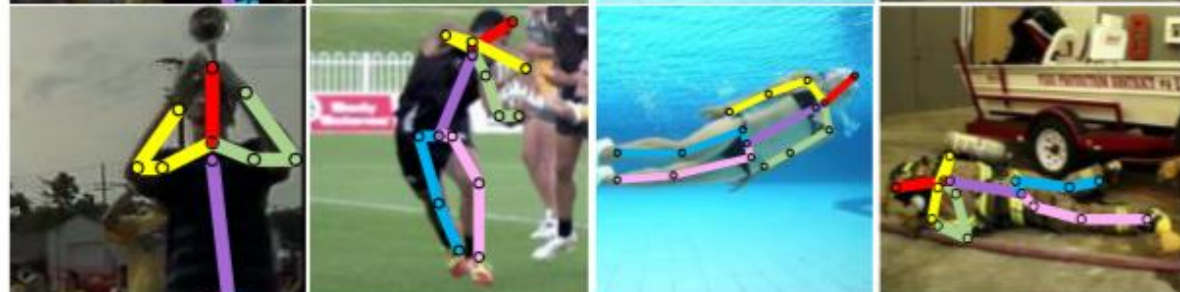
Original Images



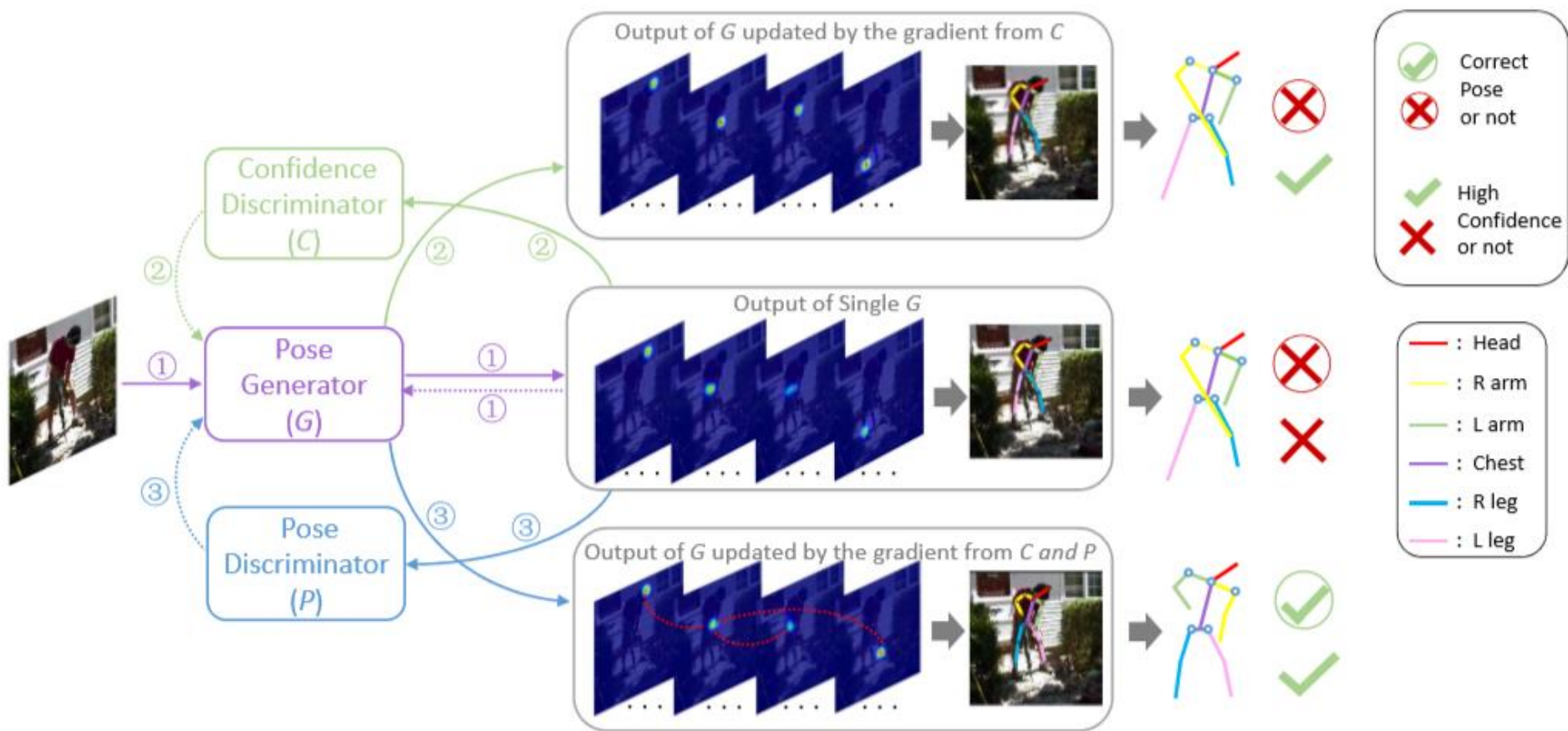
Pose predictions
without Structure



Structure-aware
pose predictions



网络遵循GAN结构，由一个生成器+2个判别器构成



Pose generator G

Multi-task: 输出32channel 16个pose heatmap + 对应的occlusion heatmap

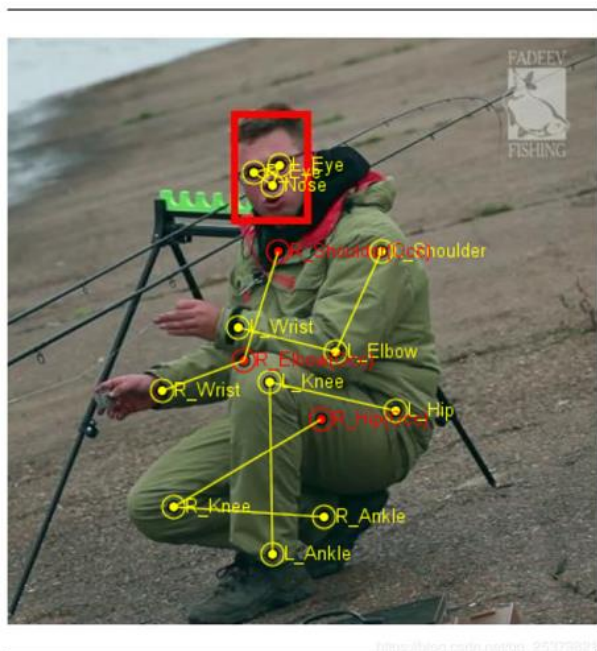
采用stack的方式 均为带有跳接结构的encoder-decoder网络

$$\begin{cases} \{Y_n, Z_n, X\} = \mathcal{G}_n(Y_{n-1}, Z_{n-1}, X) & \text{if } n \geq 2 \\ \{Y_n, Z_n, X\} = \mathcal{G}_n(X) & \text{if } n = 1 \end{cases}$$

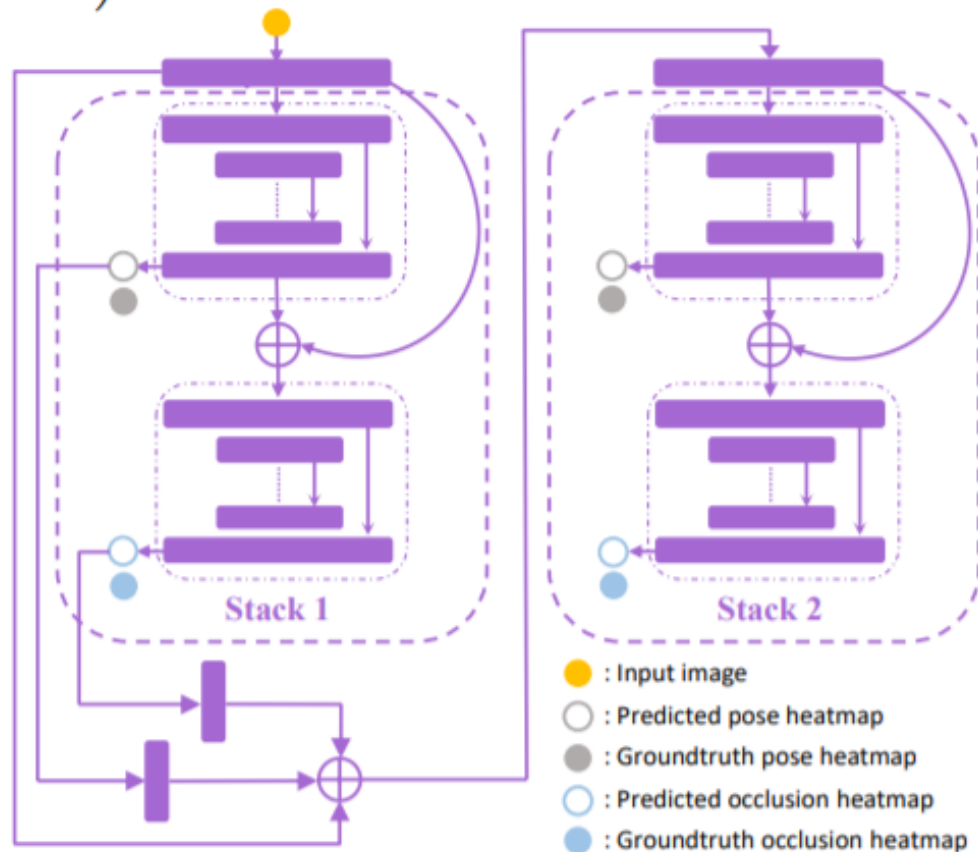
$$\mathcal{L}_G(\Theta) = \frac{1}{2MN} \sum_{n=1}^N \sum_{i=1}^M \left(\|y^i - \hat{y}_n^i\|^2 + \|z^i - \hat{z}_n^i\|^2 \right)$$

N:stack

M:data number



MPII数据集

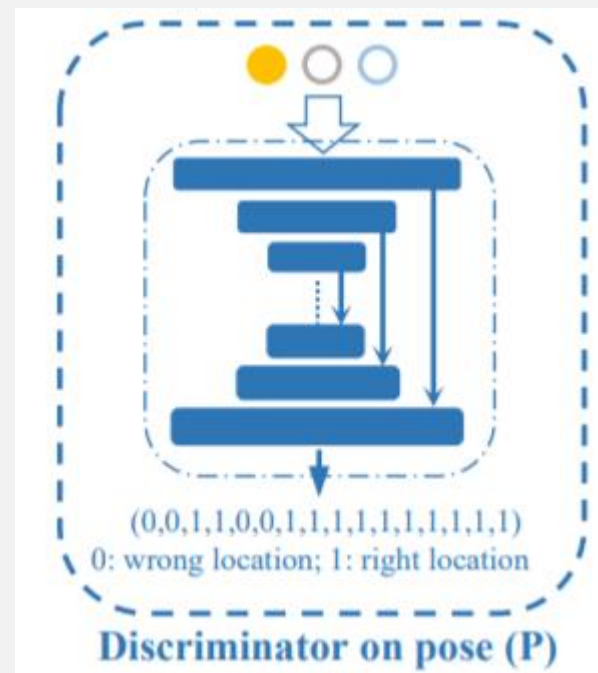


Pose Discriminator P

区分那些不满足人体关节点约束的pose

为保证对输入图片姿态合理，需要同时送入输入图片

$$\mathcal{L}_P(G, P) = \mathbb{E}[\log P(\mathbf{y}, \mathbf{z}, \mathbf{x})] + \mathbb{E}[\log(1 - |P(G(\mathbf{x}), \mathbf{x}) - \mathbf{p}_{\text{fake}}|)].$$



传统的gt=0/1会导致训练困难 → 对16个关节点分别分析

在传统GAN中 $p_{\text{fake}}=0$ → 可以认为那些和真实姿态接近的预测结果为真样本

(当一个关节点偏离真值很远时，会导致整个姿态出错)

$$\mathbf{p}_{\text{fake}}^i = \begin{cases} 1 & \text{if } d_i < \delta \\ 0 & \text{if } d_i \geq \delta \end{cases}$$

δ 为normalized distance

Confidence Discriminator C

区分输出热图的low-confidence与high-confidence (Gaussian centered)

即网络在所预测位置是否confident

输入为pose + occlusion heatmap

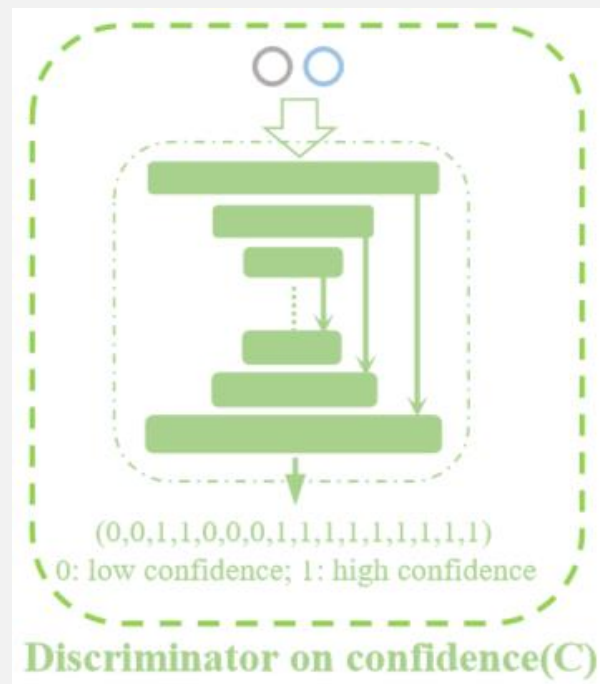
$$\mathcal{L}_C(G, C) = \mathbb{E}[\log C(\mathbf{y}, \mathbf{z})] + \mathbb{E}[\log(1 - |C(G(\mathbf{x})) - \mathbf{c}_{\text{fake}}|)].$$

传统的gt=0/1会导致训练困难 → 对16个关节点分别分析

在传统GAN中cfake=0

→ 可以认为当prediction与gt heatmap较像时为真样本

$$\mathbf{c}_{\text{fake}}^i = \begin{cases} 1 & \text{if } \|\mathbf{y}_i - \hat{\mathbf{y}}_i\| < \varepsilon \\ 0 & \text{if } \|\mathbf{y}_i - \hat{\mathbf{y}}_i\| \geq \varepsilon \end{cases}$$



生成器G的训练

$$\mathcal{L}_G(\Theta) = \frac{1}{2MN} \sum_{n=1}^N \sum_{i=1}^M \left(\|\mathbf{y}^i - \hat{\mathbf{y}}_n^i\|^2 + \|\mathbf{z}^i - \hat{\mathbf{z}}_n^i\|^2 \right)$$

$$\mathcal{L}_P(G, P) = \mathbb{E}[\log P(\mathbf{y}, \mathbf{z}, \mathbf{x})] + \mathbb{E}[\log(1 - |P(G(\mathbf{x}), \mathbf{x}) - \mathbf{p}_{\text{fake}}|)] .$$

$$\mathbf{p}_{\text{fake}}^i = \begin{cases} 1 & \text{if } d_i < \delta \\ 0 & \text{if } d_i \geq \delta \end{cases}$$

$$\mathcal{L}_C(G, C) = \mathbb{E}[\log C(\mathbf{y}, \mathbf{z})] + \mathbb{E}[\log(1 - |C(G(\mathbf{x})) - \mathbf{c}_{\text{fake}}|)] .$$

$$\mathbf{c}_{\text{fake}}^i = \begin{cases} 1 & \text{if } \|\mathbf{y}_i - \hat{\mathbf{y}}_i\| < \varepsilon \\ 0 & \text{if } \|\mathbf{y}_i - \hat{\mathbf{y}}_i\| \geq \varepsilon \end{cases}$$

总Loss:

$$\arg \min_G \max_{P, C} \mathcal{L}_G(\Theta) + \alpha \mathcal{L}_C(G, C) + \beta \mathcal{L}_P(G, P) .$$

当cfake=creal时 $\alpha=0$

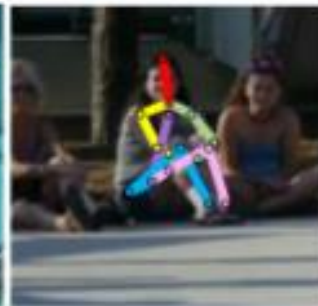
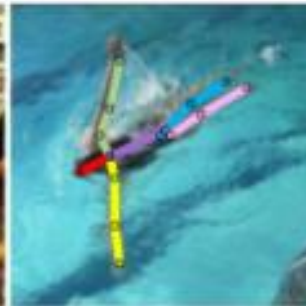
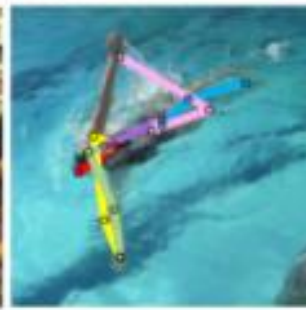
当pfake=freal时 $\beta=0$

排除可能训练得不好的判别器的影响

Results



(b)



(c)



i-VisionGroup@Tsinghua

谢谢大家!