



Boulder

Computer Vision; Pose Estimation

Maziar Raissi

Assistant Professor

Department of Applied Mathematics

University of Colorado Boulder

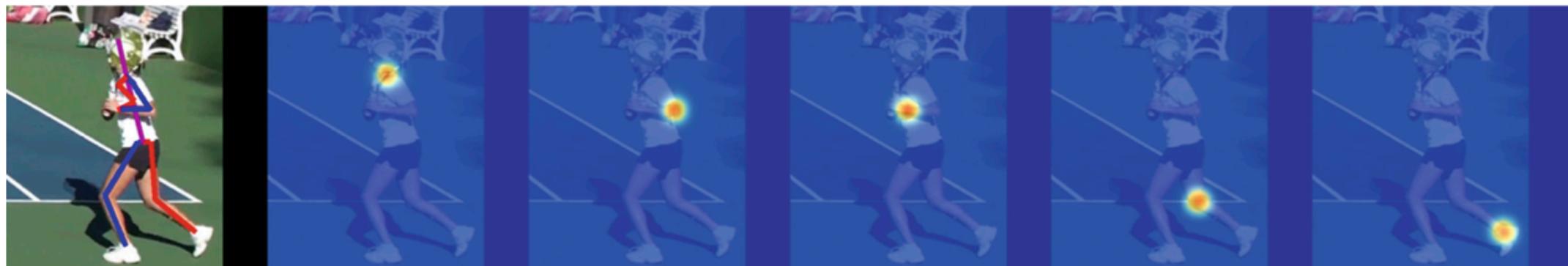
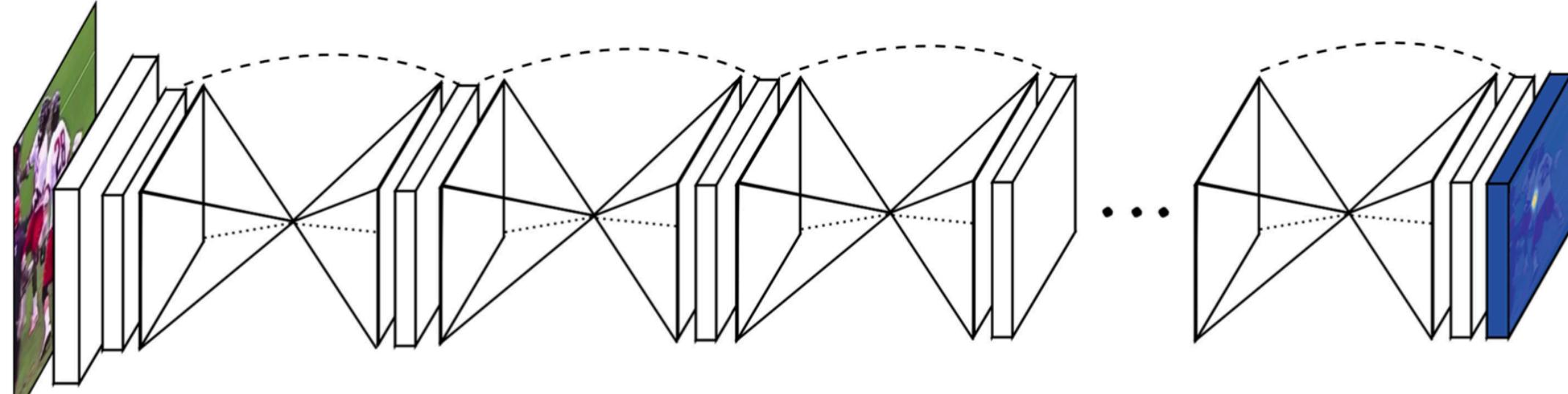
maziar.raissi@colorado.edu



Boulder

Stacked Hourglass Networks for Human Pose Estimation

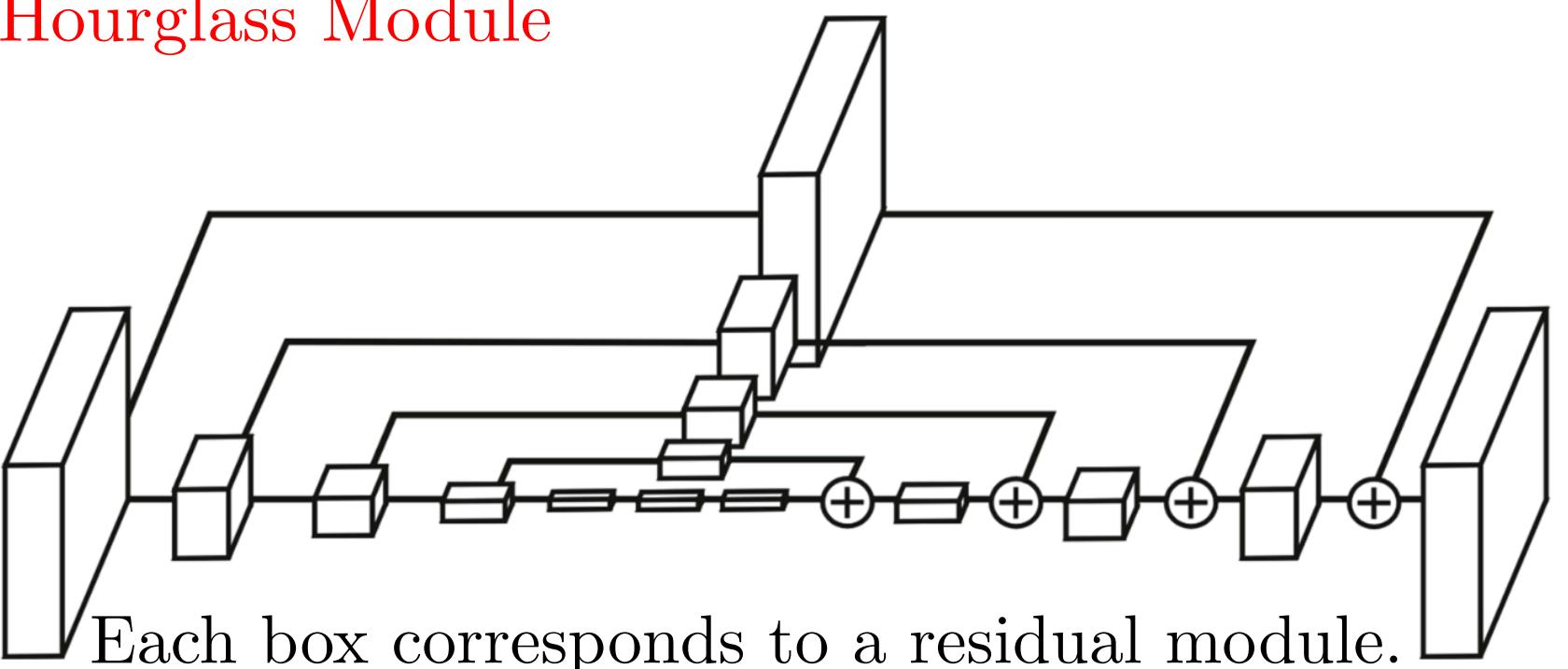
Human-computer interaction and animation



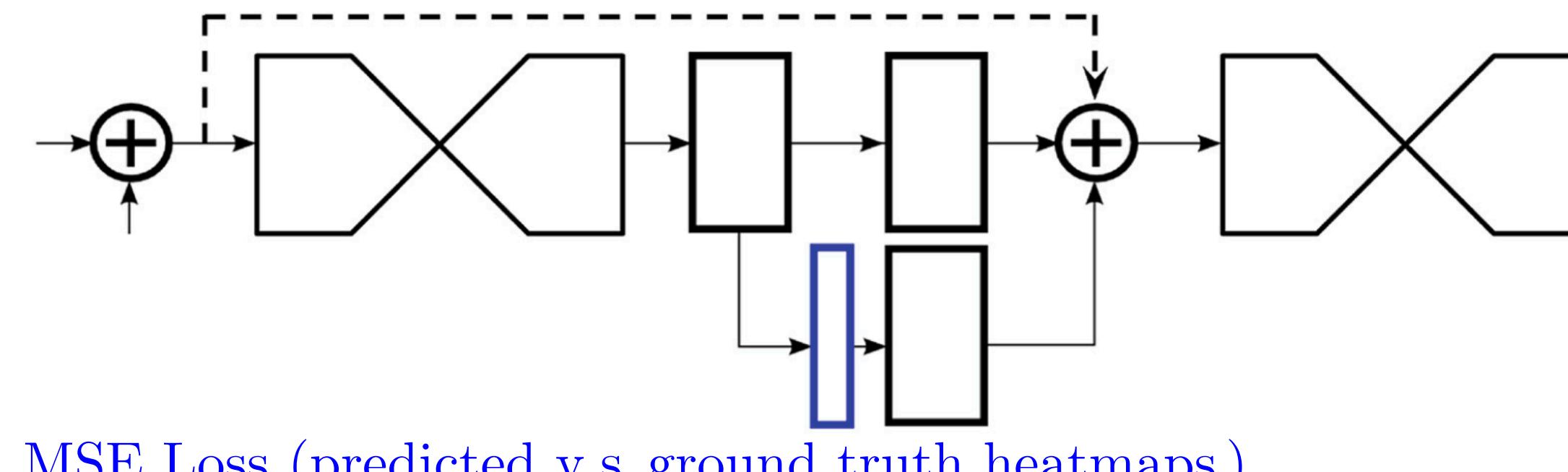
Motivation: Capturing information at every scale!

“Local evidence is essential for identifying features like faces and hands, while a final pose estimate requires a coherent understanding of the full body.”

Hourglass Module



Intermediate Supervision

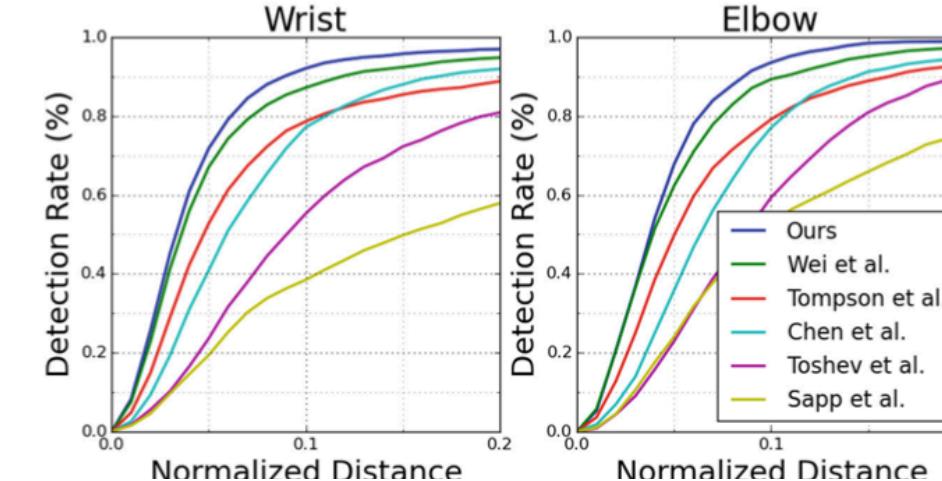


MSE Loss (predicted v.s. ground truth heatmaps)

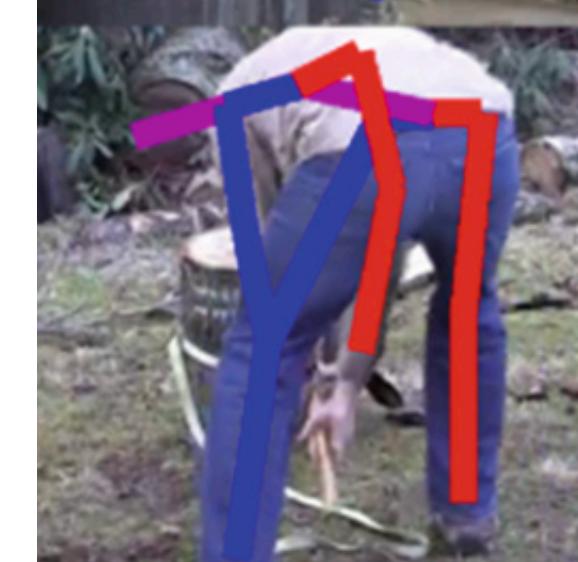
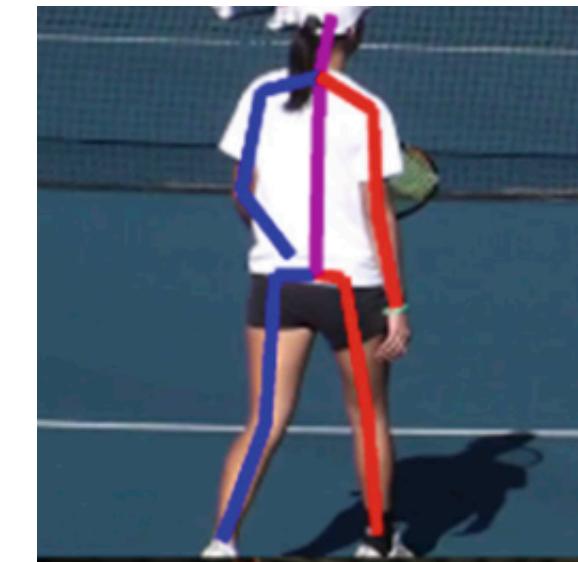
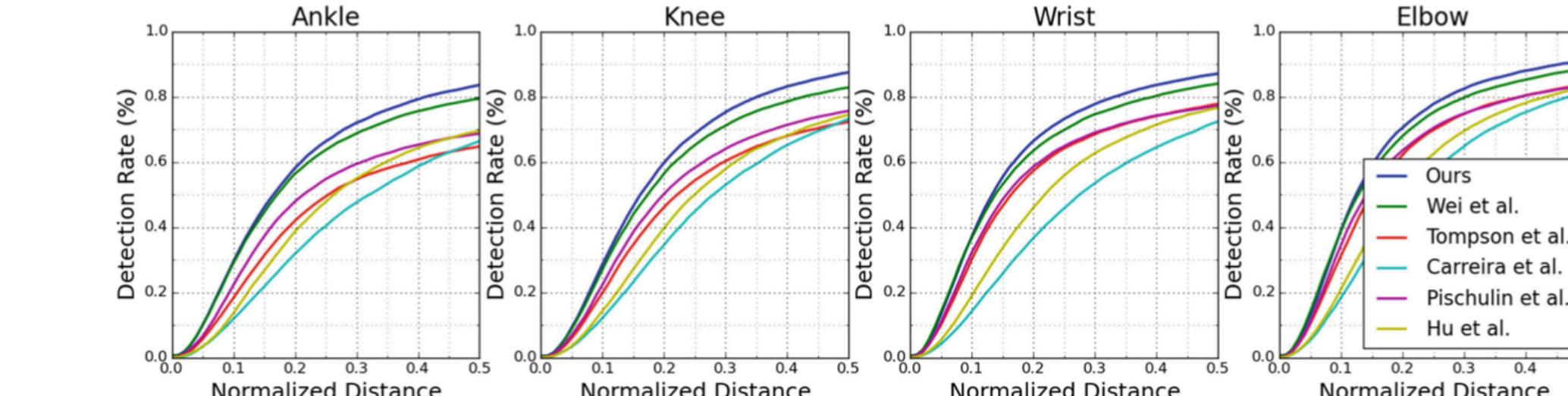
Gaussian

Percentage of Correct Keypoints (PCK): Percentage of detections that fall within a normalized distance of the ground truth.

FLIC Results

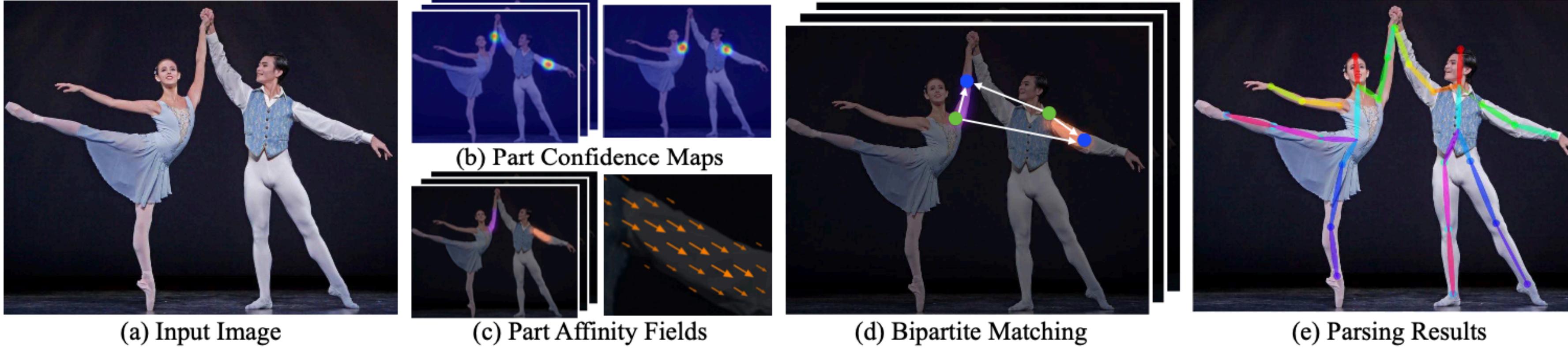


MPII Results



Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields

Human 2D pose estimation → localizing anatomical keypoints or “parts”



(a) Input Image

(c) Part Affinity Fields

(d) Bipartite Matching

(e) Parsing Results

$S \rightarrow$ set of 2D confidence maps of body part locations

$L \rightarrow$ set of 2D vector fields of part affinities
(encode the degree of association between parts)

$S = (S_1, \dots, S_J)$ has J confidence maps (one per part)

$S_j \in \mathbb{R}^{w \times h}$

$L = (L_1, \dots, L_C)$ has C vector fields (one per limb)

$L_c \in \mathbb{R}^{w \times h \times 2} \rightarrow$ each image location encodes a 2D vector

$F \leftarrow$ image \triangleright VGG-19 (first 10 layers)

$$S^1 = \rho^1(F) \quad S^t = \rho^t(F, S^{t-1}, L^{t-1}) \text{ for } t \geq 2$$

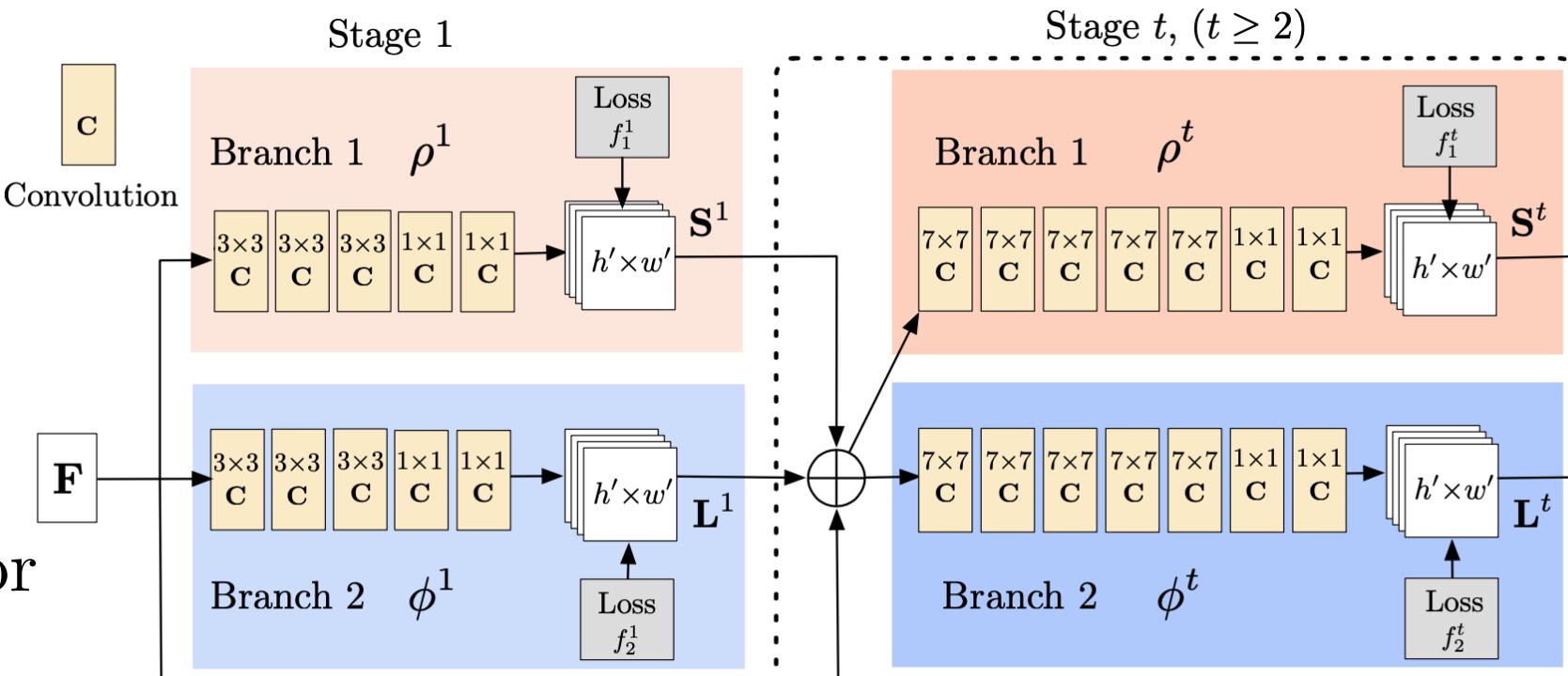
$$L^1 = \phi^1(F) \quad L^t = \phi^t(F, S^{t-1}, L^{t-1}) \text{ for } t \geq 2$$

$W(p) = 0$ when the annotation is missing at an image location p

$$f_S^t = \sum_{j=1}^J \sum_p W(p) \|S_j^t(p) - S_j^*(p)\|_2^2$$

$$f = \sum_{t=1}^T (f_S^t + f_L^t)$$

$$\underbrace{\text{loss}}_{f_L^t} = \sum_{c=1}^C \sum_p W(p) \|L_c^t(p) - L_c^*(p)\|_2^2$$

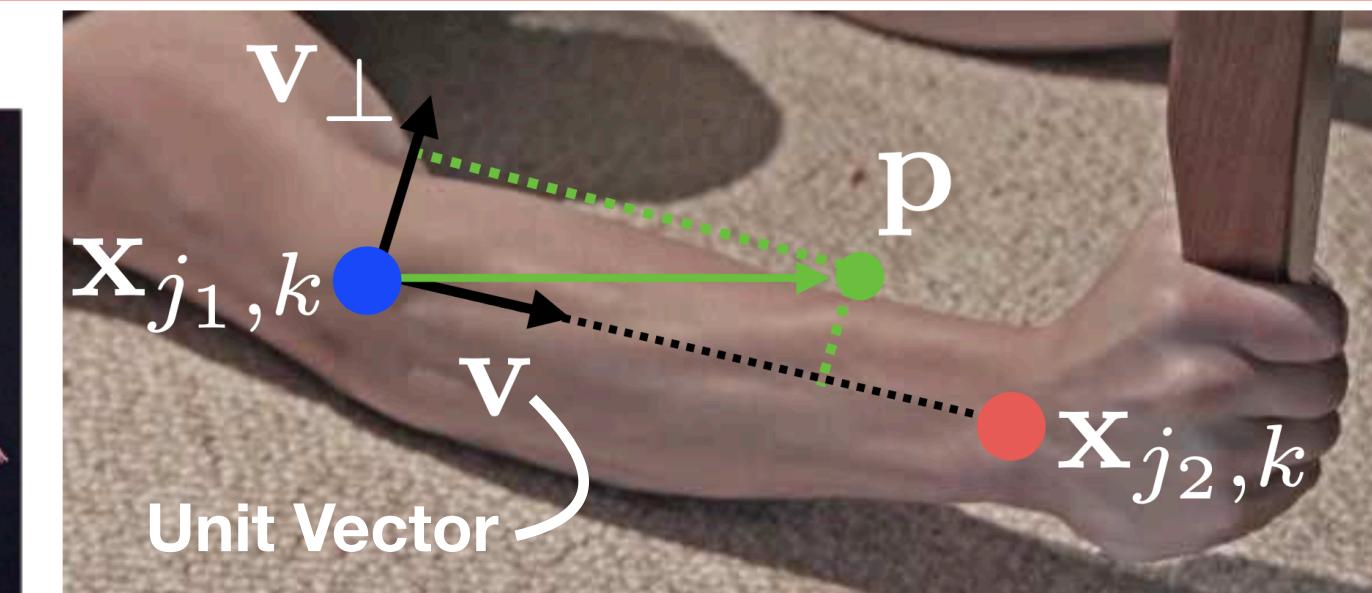


$$S_j^*(p) = \max_k S_{j,k}^*(p) \rightarrow \text{groundtruth confidence map}$$

$S_{j,k}^*(p) \rightarrow$ confidence map for person k

$$S_{j,k}^*(p) = \exp\left(-\frac{\|p - x_{j,k}\|_2^2}{\sigma^2}\right)$$

$x_{j,k} \in \mathbb{R}^2 \rightarrow$ groundtruth position of body part j for person k in the image



$$L_c^*(p) = \frac{1}{n_c(p)} \sum_k L_{c,k}^*(p)$$

groundtruth part affinity field

$n_c(p) \rightarrow$ number of non-zero vectors
at point p across all k people

$$L_{c,k}^*(p) = \begin{cases} v & \text{if } p \text{ on limb } c, k \\ 0 & \text{otherwise} \end{cases}$$

$$0 \leq \mathbf{v} \cdot (\mathbf{p} - \mathbf{x}_{j1,k}) \leq l_{c,k} \text{ and } |\mathbf{v}_\perp \cdot (\mathbf{p} - \mathbf{x}_{j1,k})| \leq \sigma_l$$

Testing $l_{c,k} = \|\mathbf{x}_{j2,k} - \mathbf{x}_{j1,k}\|_2$

$d_{j1}, d_{j2} \rightarrow$ two candidate part locations

$$p(u) = (1 - u)d_{j1} + ud_{j2}$$

$$E = \int_{u=0}^{u=1} L_c(p(u)) \cdot \frac{d_{j2} - d_{j1}}{\|d_{j2} - d_{j1}\|_2} du$$

association confidence





Boulder

Questions?
