

# HDNet: Human Depth Estimation for Multi-Person Camera-Space Localization



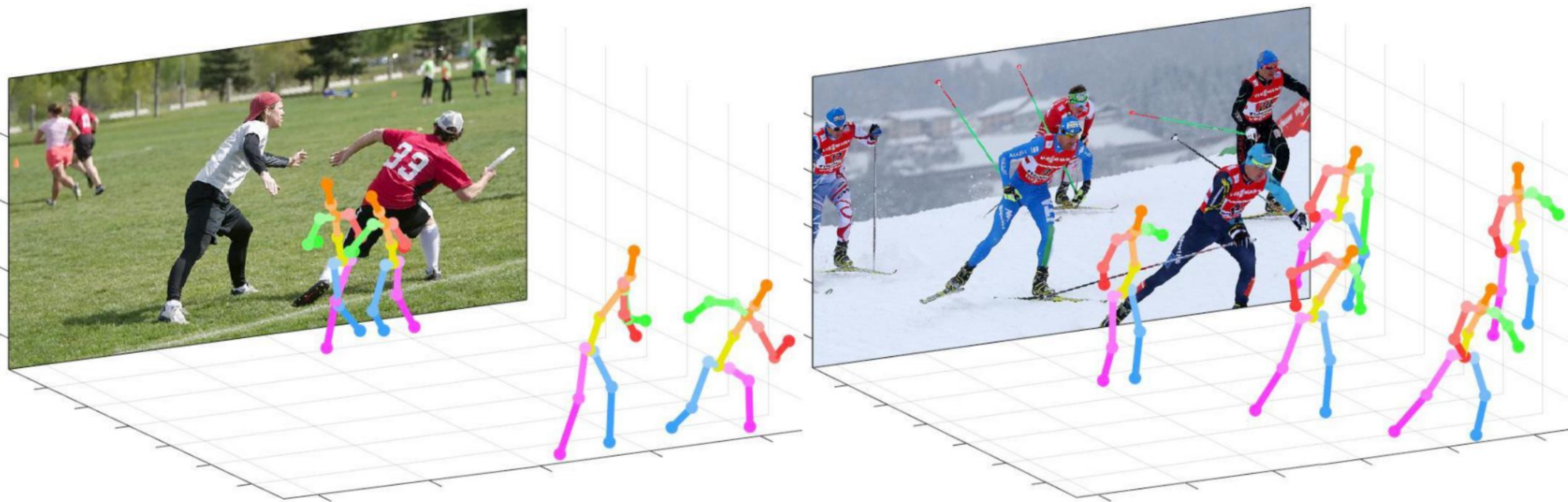
Jiahao Lin



Gim Hee Lee



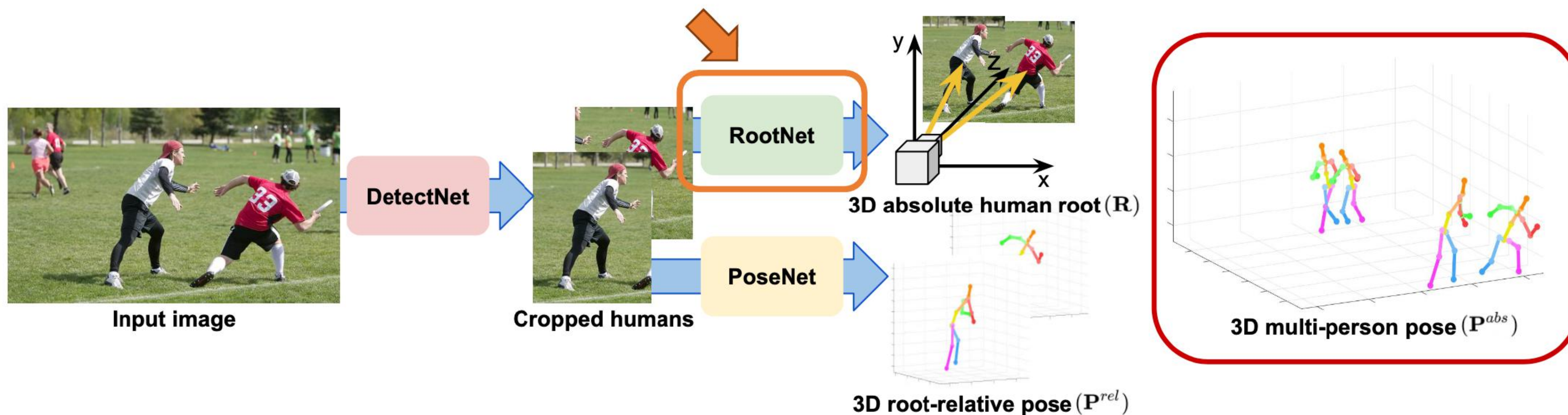
# Multi-Person 3D Pose Estimation





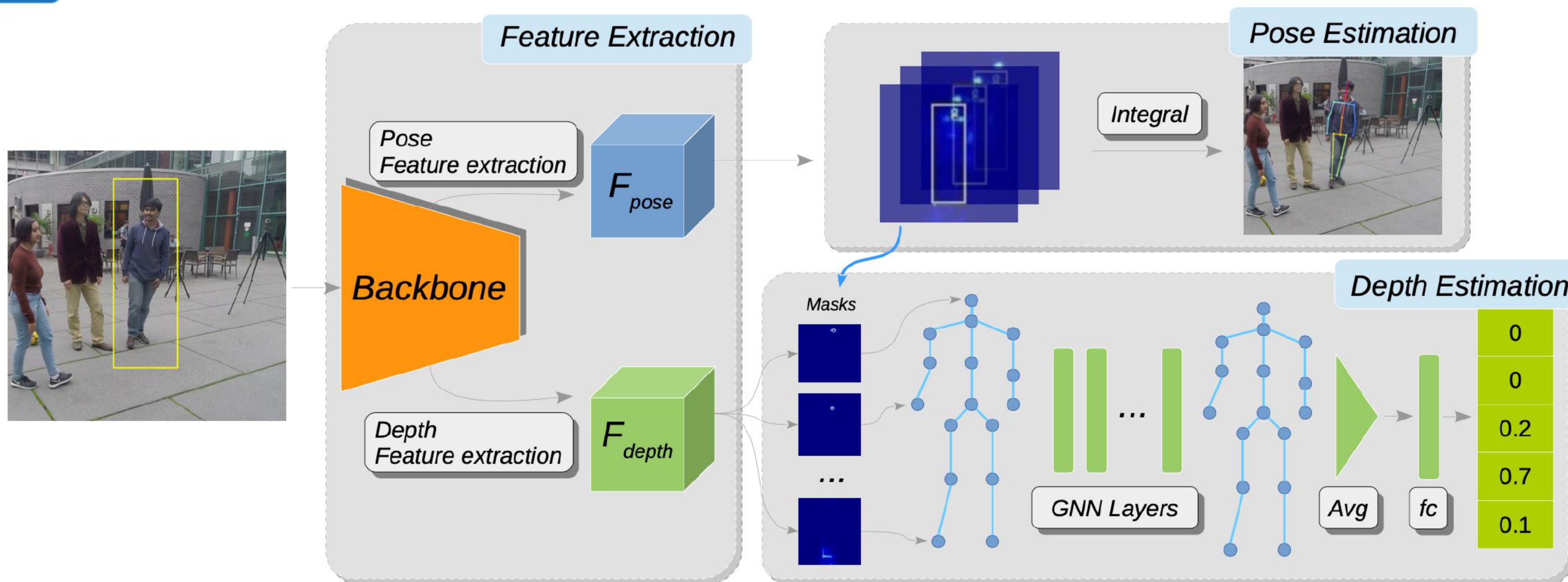
# Multi-Person 3D Pose Estimation

## Root joint localization



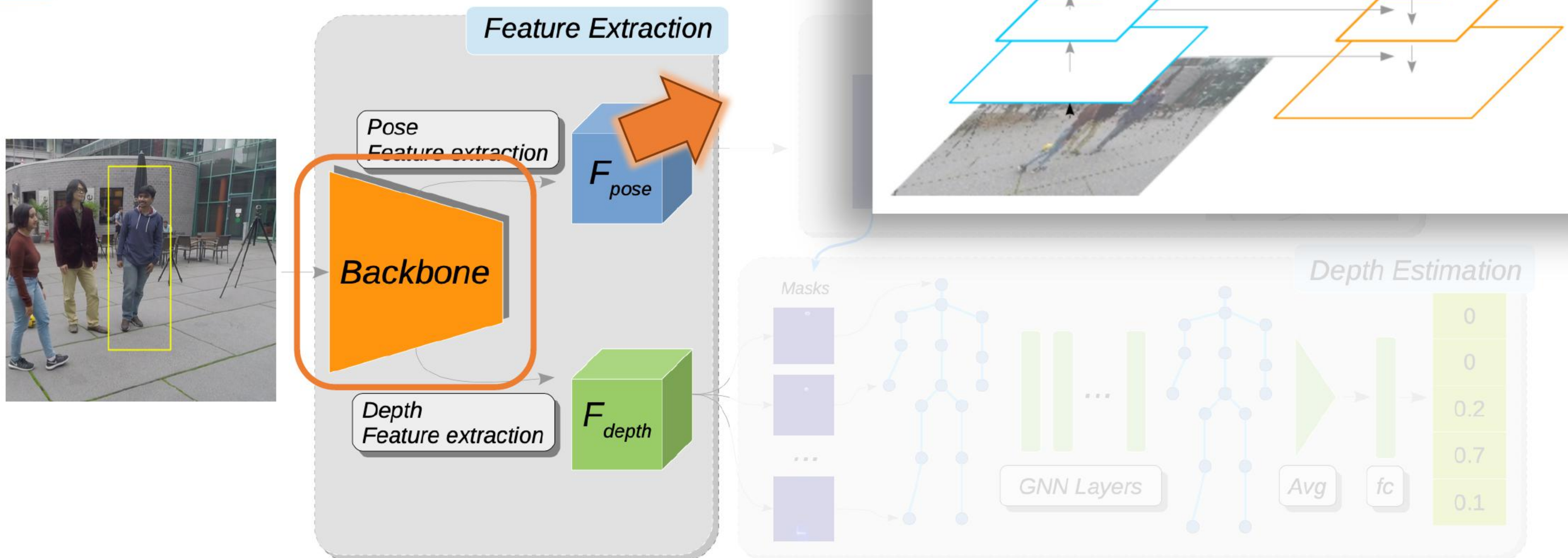


# Overview



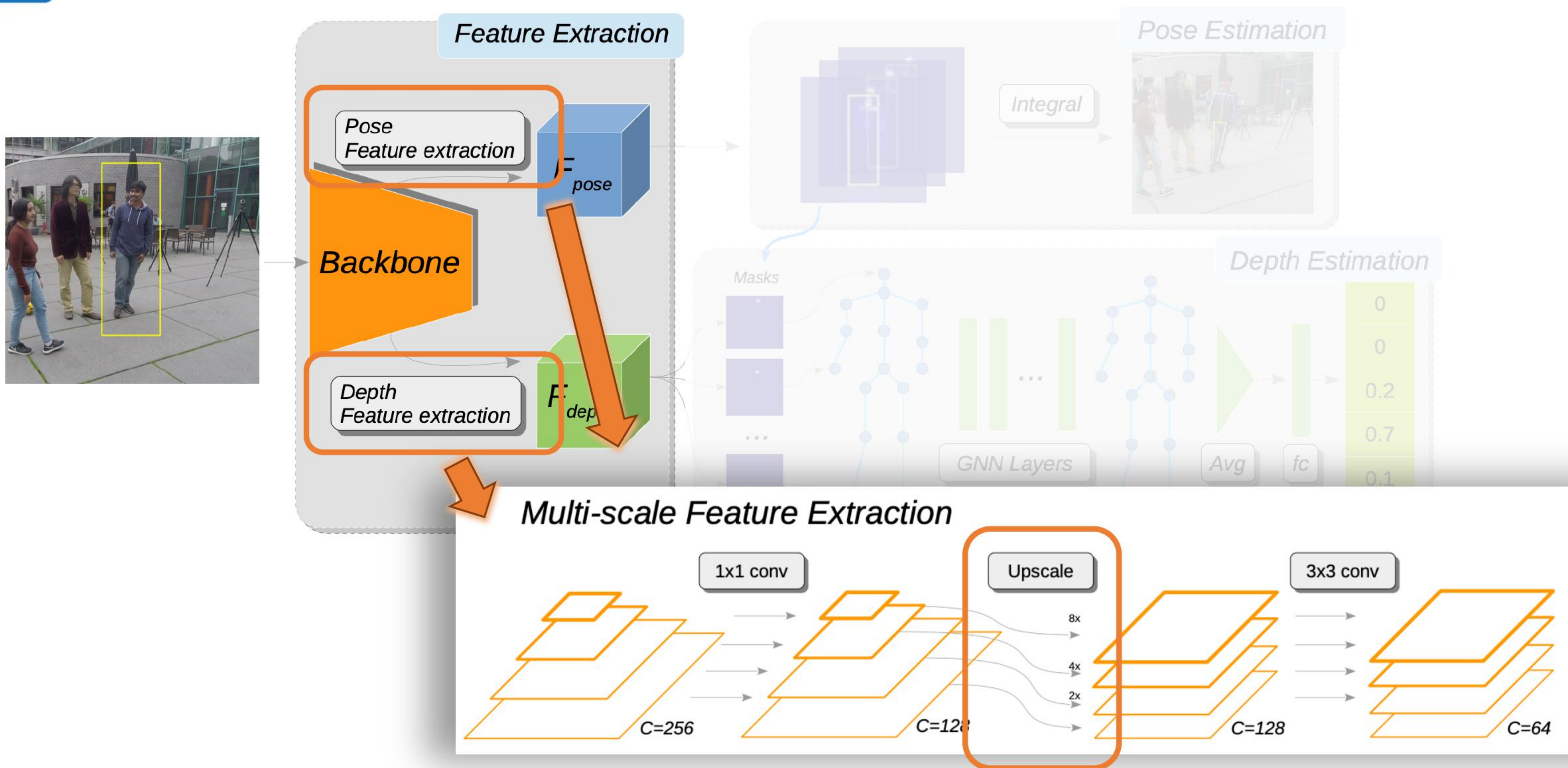


# Feature Extraction





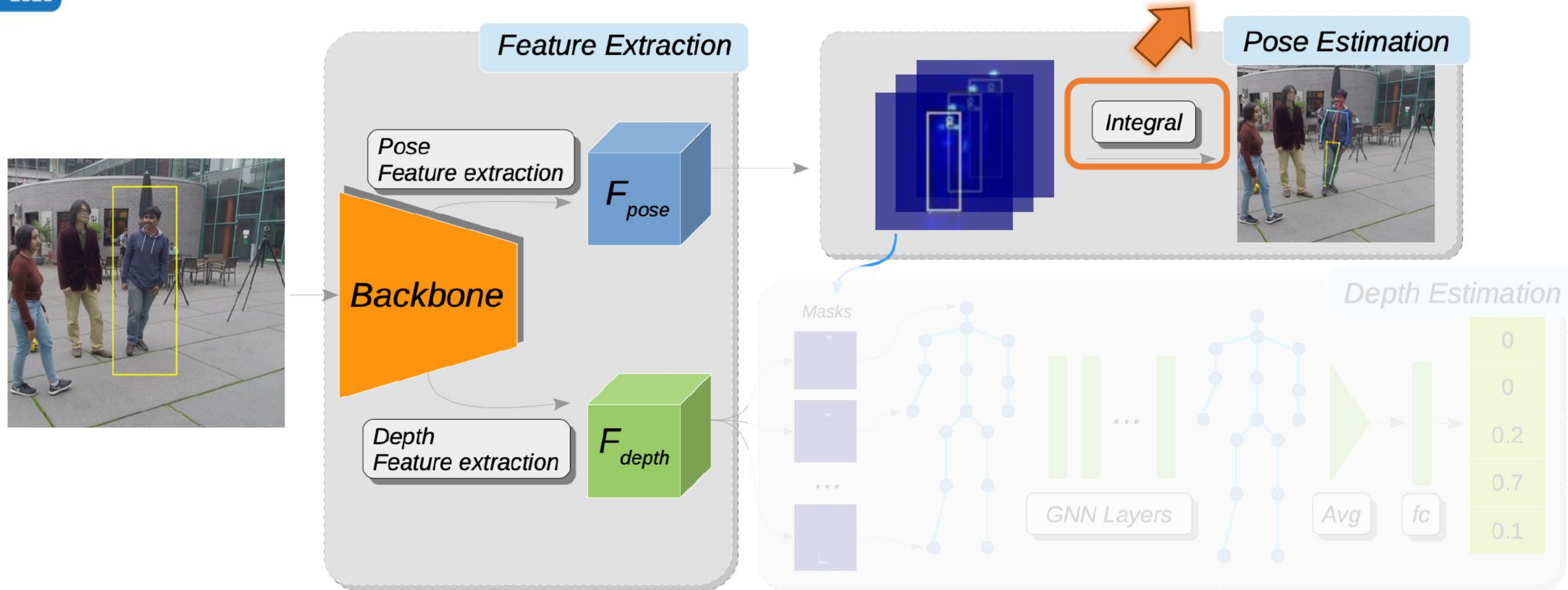
# Feature Extraction





# 2D Pose Estimation

$$\text{Soft-argmax} \quad (\hat{u}_j, \hat{v}_j) = \sum_{(u,v)=(0,0)}^{(W-1,H-1)} \hat{\mathbf{H}}_{u,v}^{(j)} \cdot (u, v),$$



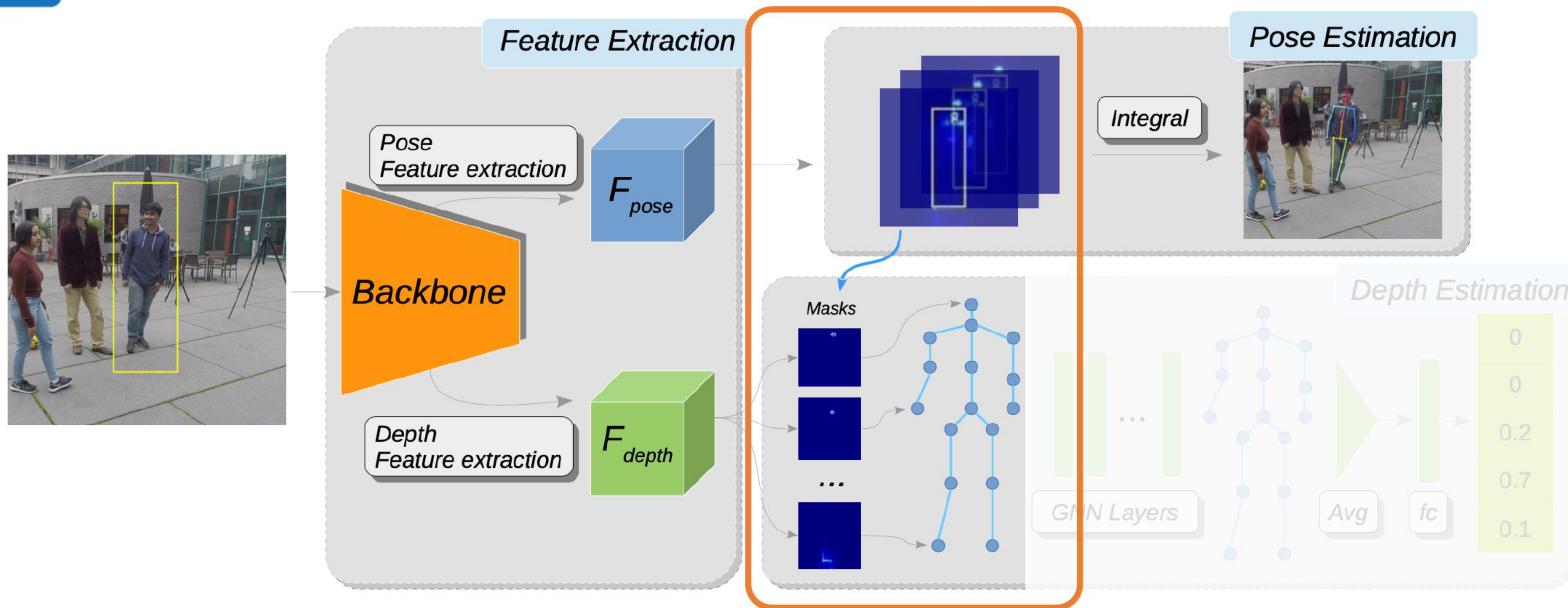
$$\mathcal{L}_{\text{hm}} = \frac{1}{N_J H W} \sum_j^{N_J} \sum_{(u,v)=(0,0)}^{(W-1,H-1)} \left\| \mathbf{H}_{u,v}^{(j)\text{GT}} - \hat{\mathbf{H}}_{u,v}^{(j)} \right\|^2$$

$$\mathcal{L}_{\text{pose}} = \frac{1}{N_J} \sum_j^{N_J} \left( \left| u_j^{\text{GT}} - \hat{u}_j \right| + \left| v_j^{\text{GT}} - \hat{v}_j \right| \right)$$



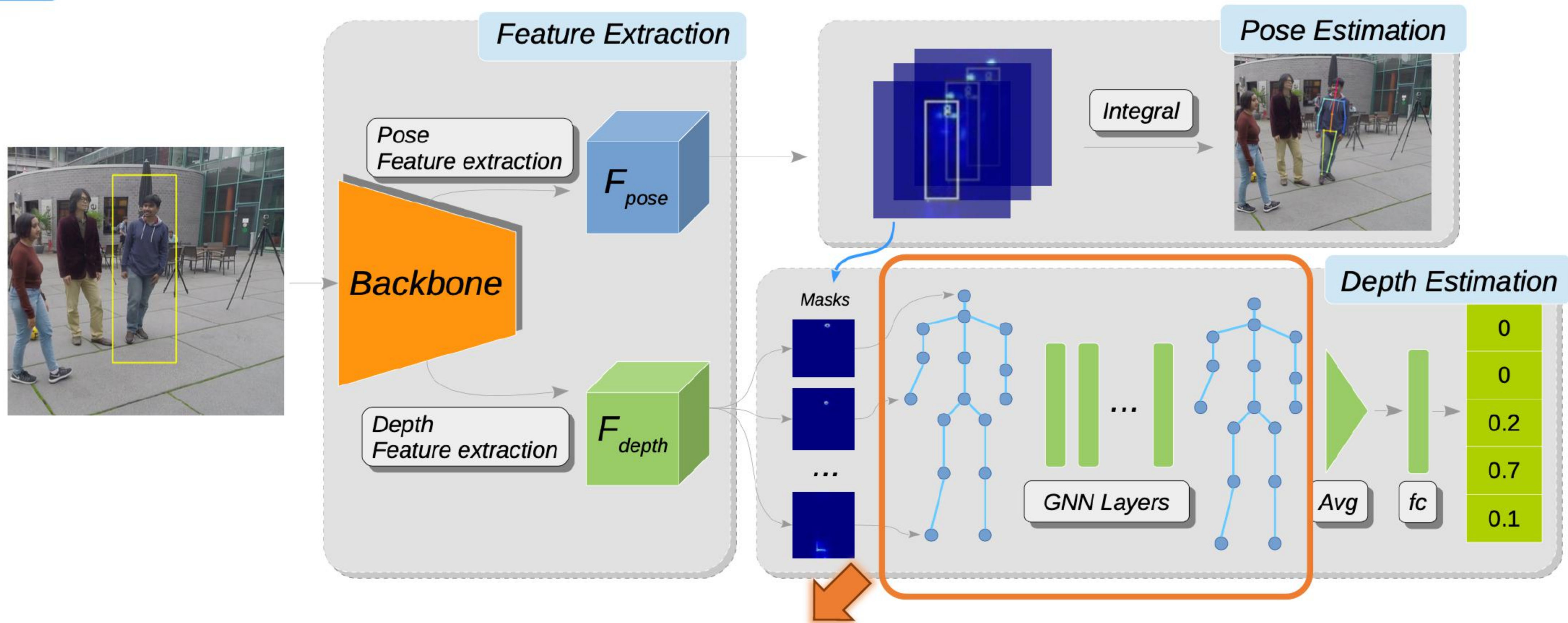
# Depth Estimation

$$\mathbf{d}^{(j)} = \sum_{(u,v)=(0,0)}^{(W-1,H-1)} \hat{\mathbf{H}}_{u,v}^{(j)} \cdot \mathbf{F}_{\text{depth}_{u,v}}$$





# Depth Estimation



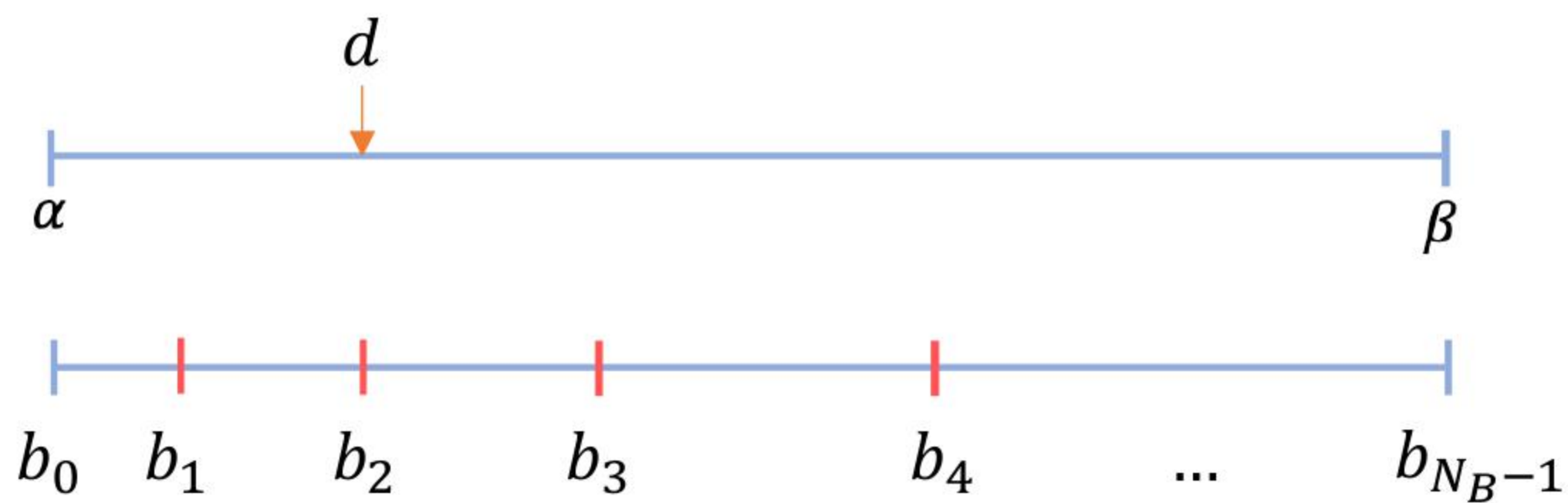
$$X_{out}^{(i)} = \sigma \left( \tilde{a}_{ii} \underline{f_{self}}(X_{in}^{(i)}; \Theta_{self}) + \sum_{j \neq i} \tilde{a}_{ij} \underline{f_{inter}}(X_{in}^{(j)}; \Theta_{inter}) \right)$$

$$A \in \{0, 1\}^{N_j \times N_j}$$

$$\rightarrow \tilde{A}$$

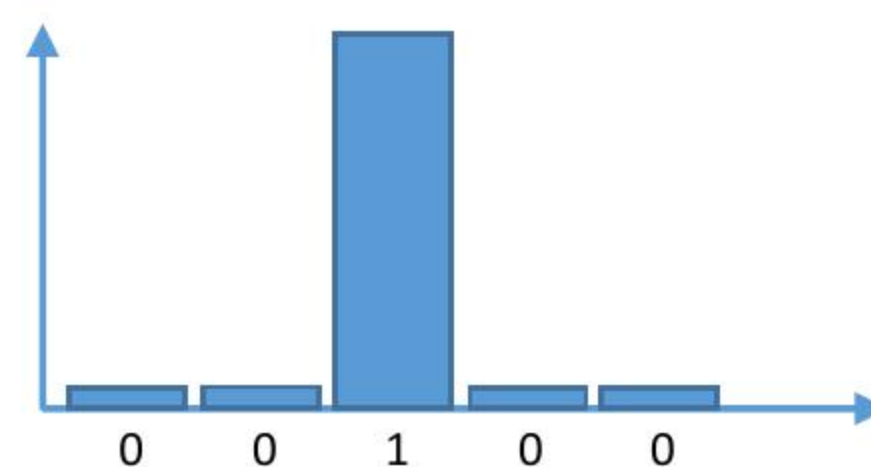


# Depth Estimation

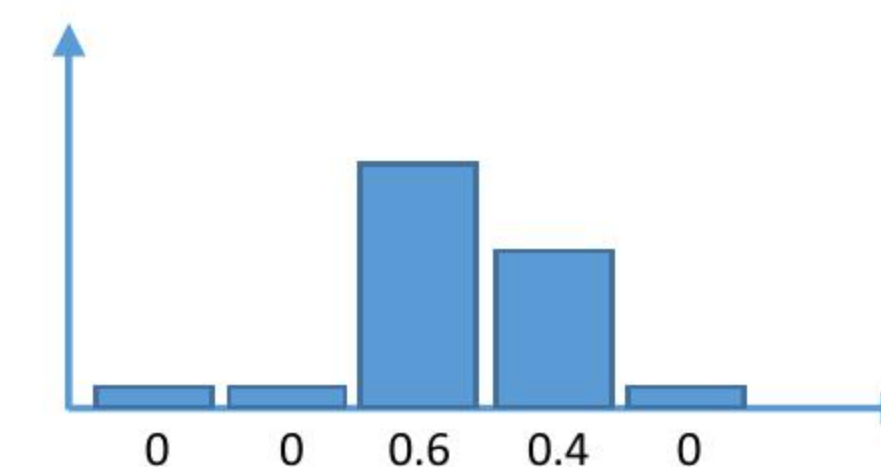


$$b(d) = \frac{\log d - \log \alpha}{\log \beta - \log \alpha} \cdot (N_{\mathbf{B}} - 1)$$

$N_{\mathbf{B}} = 5$  and  $b = 2.4$ .



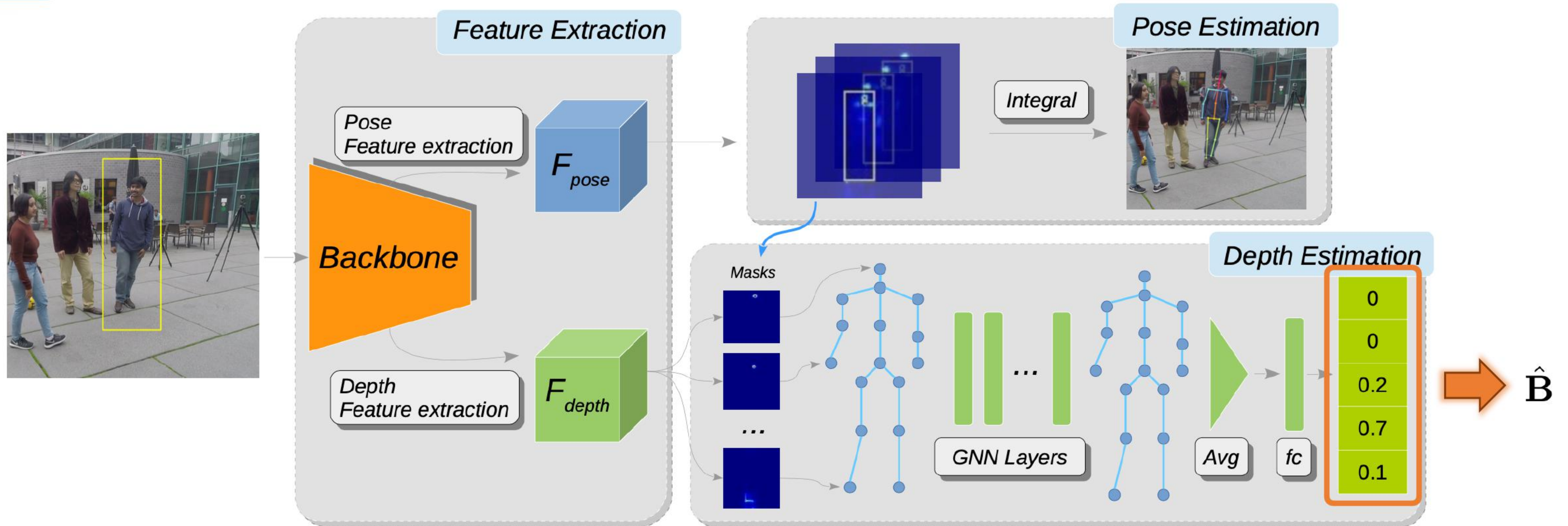
One-hot



Bi-linear



# Depth Estimation



$$d = \left[ \frac{\hat{b}}{N_B - 1} \cdot (\log \beta - \log \alpha) + \log \alpha \right] \cdot \sqrt{f_x \cdot f_y}, \text{ where } \hat{b} = \sum_{i=0}^{N_B-1} \hat{\mathbf{B}}_i \cdot i$$

$$\mathcal{L}_{\text{bins}} = - \sum_{i=0}^{N_B-1} \mathbf{B}_i^{\text{GT}} \cdot \log \hat{\mathbf{B}}_i, \text{ and } \mathcal{L}_{\text{idx}} = \left| b^{\text{GT}} - \hat{b} \right|$$

Soft-argmax



# Evaluation



Human3.6M

- Mean Root Position Error (MRPE)



MuCo-3DHP



MuPoTS-3D

- Average Precision (AP) and Recall (AR)
- $3DPCK_{rel}$  and  $3DPCK_{abs}$



# Results

**Table 1.** MRPE results comparison with state-of-the-arts on the Human3.6M dataset.  $MRPE_x$ ,  $MRPE_y$ , and  $MRPE_z$  are the average errors in  $x$ ,  $y$ , and  $z$  axes, respectively.

| Method                   | MRPE        | $MRPE_x$    | $MRPE_y$    | $MRPE_z$    |
|--------------------------|-------------|-------------|-------------|-------------|
| Baseline                 | 267.8       | 27.5        | 28.3        | 261.9       |
| Baseline w/o limb joints | 226.2       | 24.5        | 24.9        | 220.2       |
| Baseline with RANSAC     | 213.1       | 24.3        | 24.3        | 207.1       |
| RootNet [21]             | 120.0       | 23.3        | 23.0        | 108.1       |
| <b>Ours</b>              | <b>77.6</b> | <b>15.6</b> | <b>13.6</b> | <b>69.9</b> |

**Table 2.** Root joint localization accuracy comparison in average precision and recall with state-of-the-arts on MuPoTS-3D dataset.

| Method       | $AP_{25}^{root}$ | $AP_{20}^{root}$ | $AP_{15}^{root}$ | $AP_{10}^{root}$ | $AR_{25}^{root}$ | $AR_{20}^{root}$ | $AR_{15}^{root}$ | $AR_{10}^{root}$ |
|--------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|
| RootNet [21] | 31.0             | 21.5             | 10.2             | 2.3              | 55.2             | 45.3             | 31.4             | 15.2             |
| <b>Ours</b>  | <b>39.4</b>      | <b>28.0</b>      | <b>14.6</b>      | <b>4.1</b>       | <b>59.8</b>      | <b>50.0</b>      | <b>35.9</b>      | <b>19.1</b>      |

**Table 4.** Joint-wise  $3DPCK_{abs}$  comparison with state-of-the-arts on MuPoTS-3D dataset. Accuracy is measured on matched ground-truths.

| Method       | Head        | Neck        | Shoulder    | Elbow       | Wrist       | Hip         | Knee        | Ankle       | Avg         |
|--------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| RootNet [21] | 37.6        | 35.6        | 34.0        | 34.1        | 30.7        | 30.6        | 31.3        | 25.3        | 31.8        |
| <b>Ours</b>  | <b>38.3</b> | <b>37.8</b> | <b>36.2</b> | <b>37.4</b> | <b>34.0</b> | <b>34.9</b> | <b>36.4</b> | <b>29.2</b> | <b>35.2</b> |

**Table 5.** Sequence-wise  $3DPCK_{rel}$  comparison with state-of-the-arts on MuPoTS-3D dataset. Accuracy is measured on matched ground-truths.

| Method                   | S1          | S2          | S3          | S4          | S5          | S6          | S7          | S8          | S9          | S10         | - |
|--------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|---|
| Rogez <i>et al.</i> [26] | 69.1        | 67.3        | 54.6        | 61.7        | 74.5        | 25.2        | 48.4        | 63.3        | 69.0        | 78.1        | - |
| Mehta <i>et al.</i> [20] | 81.0        | 65.3        | 64.6        | 63.9        | 75.0        | 30.3        | 65.1        | 61.1        | 64.1        | 83.9        | - |
| Rogez <i>et al.</i> [27] | 88.0        | 73.3        | 67.9        | 74.6        | 81.8        | 50.1        | 60.6        | 60.8        | 78.2        | 89.5        | - |
| RootNet [21]             | <b>94.4</b> | 78.6        | 79.0        | 82.1        | 86.6        | 72.8        | <b>81.9</b> | 75.8        | <b>90.2</b> | 90.4        | - |
| <b>Ours</b>              | <b>94.4</b> | <b>79.6</b> | <b>79.2</b> | <b>82.4</b> | <b>86.7</b> | <b>73.0</b> | 81.6        | <b>76.3</b> | 90.1        | <b>90.5</b> | - |

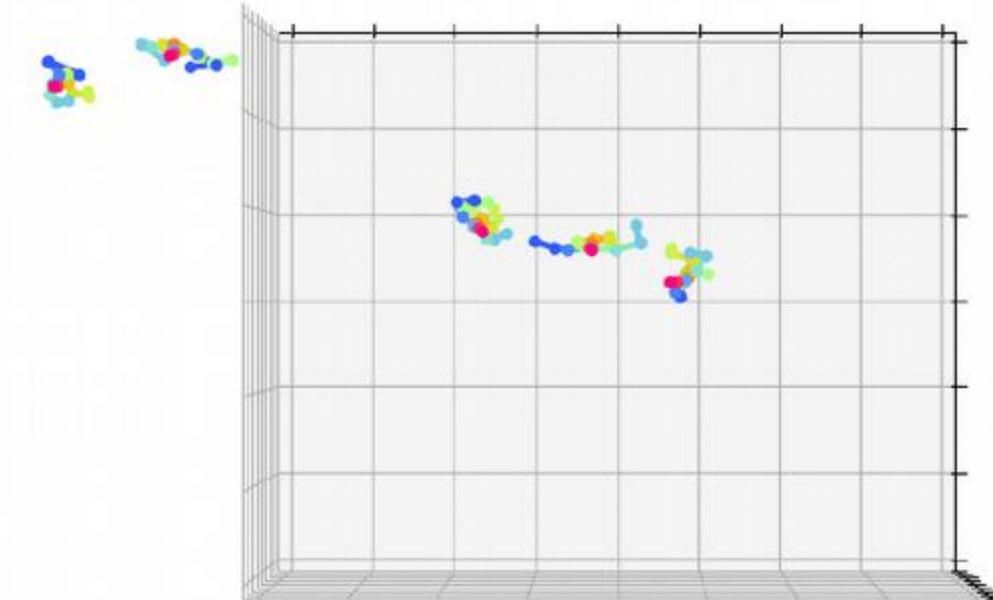
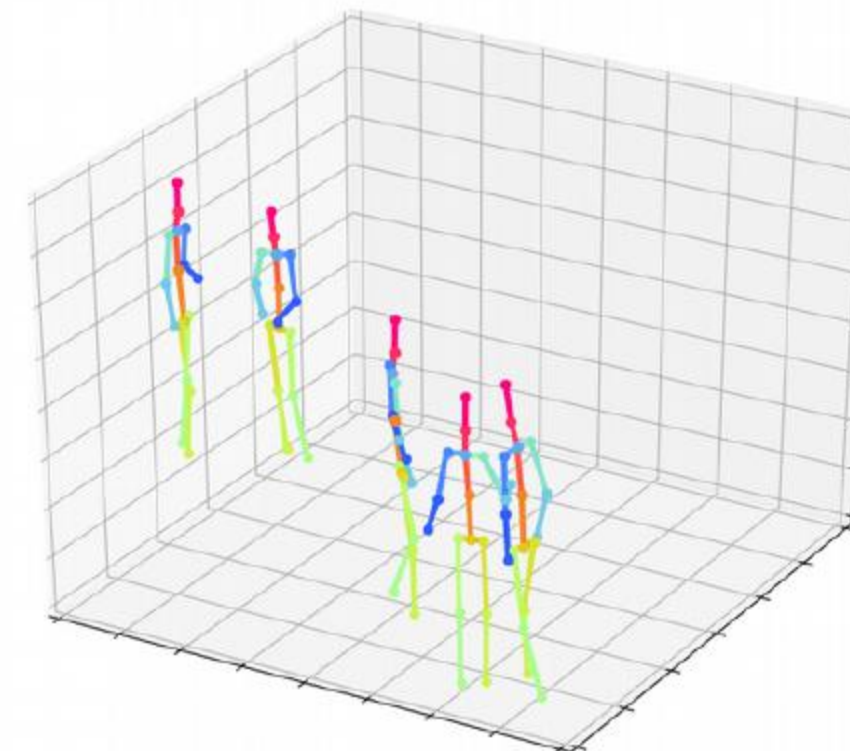
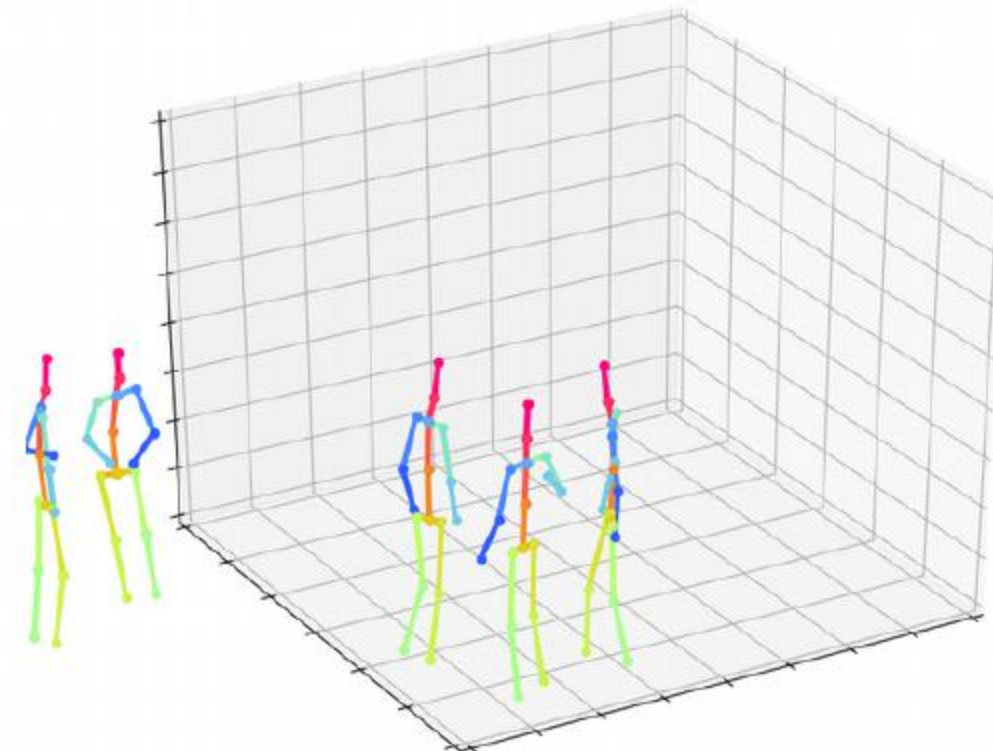
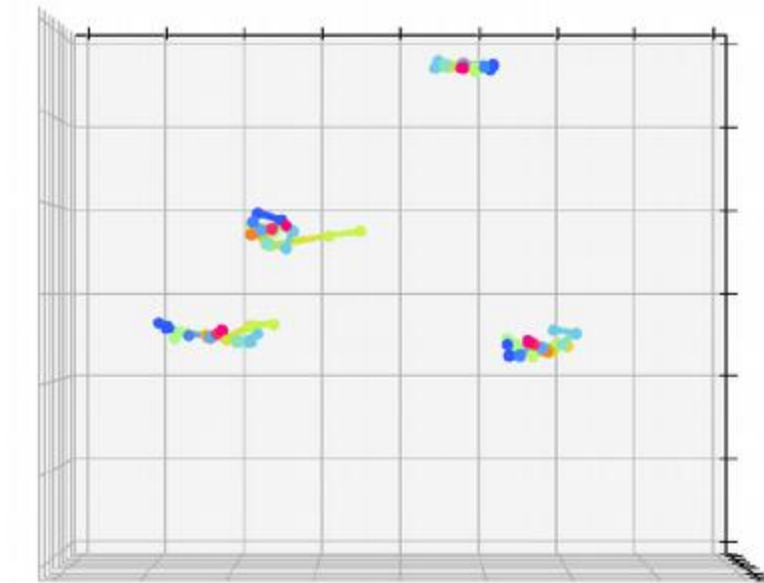
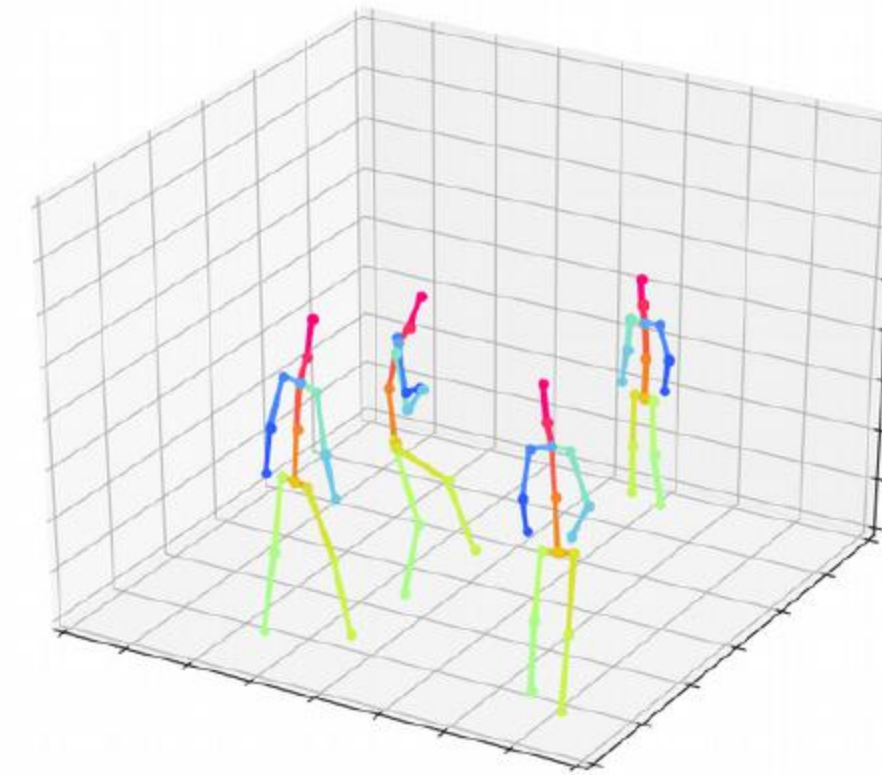
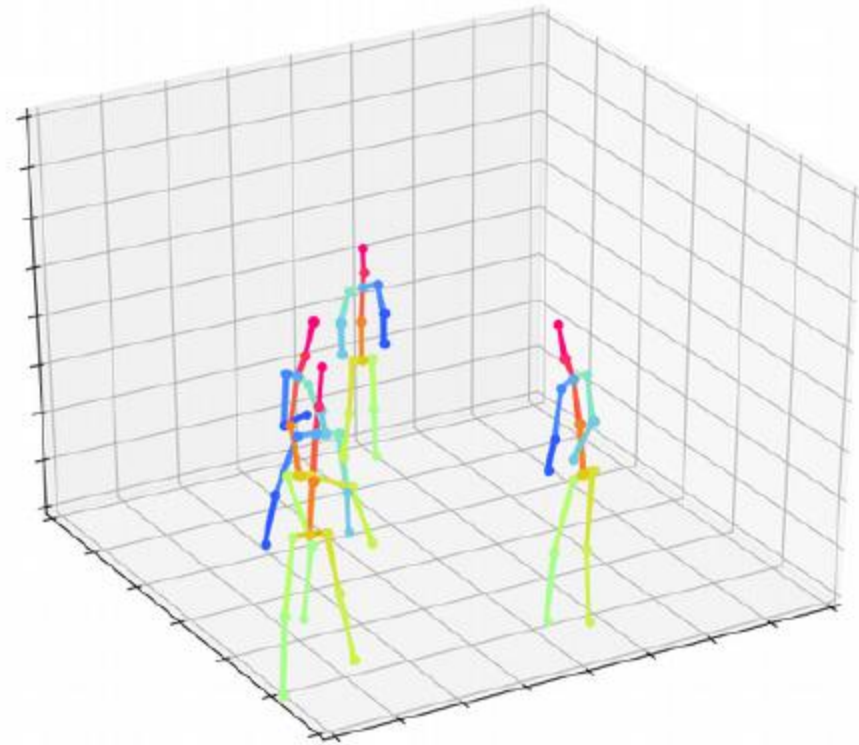
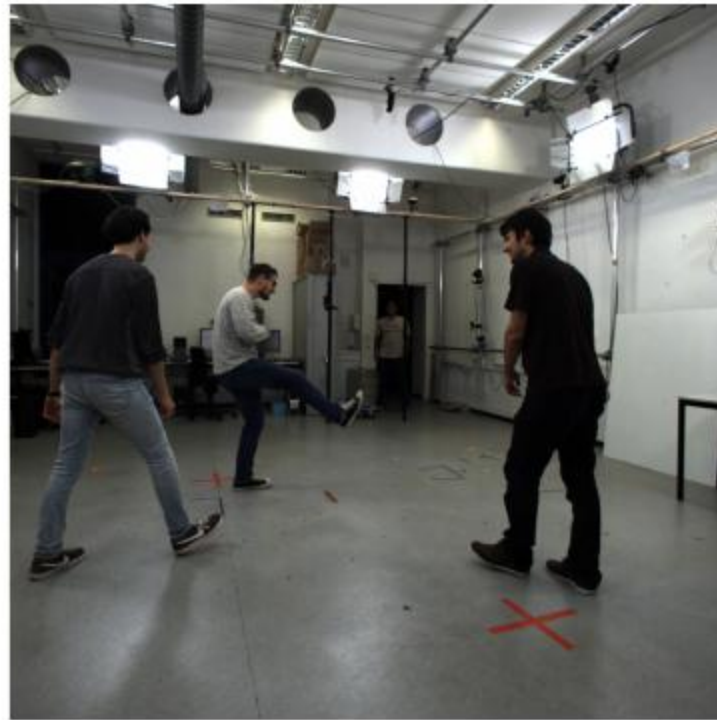
  

| Method                   | S11         | S12         | S13         | S14         | S15         | S16         | S17         | S18         | S19         | S20         | Avg         |
|--------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| Rogez <i>et al.</i> [26] | 53.8        | 52.2        | 60.5        | 60.9        | 59.1        | 70.5        | 76.0        | 70.0        | 77.1        | 81.4        | 62.4        |
| Mehta <i>et al.</i> [20] | 72.4        | 69.9        | 71.0        | 72.9        | 71.3        | 83.6        | 79.6        | 73.5        | 78.9        | <b>90.9</b> | 70.8        |
| Rogez <i>et al.</i> [27] | 70.8        | 74.4        | 72.8        | 64.5        | 74.2        | 84.9        | 85.2        | 78.4        | 75.8        | 74.4        | 74.0        |
| RootNet [21]             | <b>79.4</b> | <b>79.9</b> | 75.3        | 81.0        | <b>81.1</b> | 90.7        | <b>89.6</b> | 83.1        | 81.7        | 77.3        | 82.5        |
| <b>Ours</b>              | 77.9        | 79.2        | <b>78.3</b> | <b>85.5</b> | <b>81.1</b> | <b>91.0</b> | 88.5        | <b>85.1</b> | <b>83.4</b> | 90.5        | <b>83.7</b> |

**Table 6.** Ablation studies on components of the framework. Depth error  $MRPE_z$  (mm) on Human3.6M dataset and  $AP_{25}^{root}$  (%) on MuPoTS-3D dataset are measured.

| Method                     | $MRPE_z(\downarrow)$ | $AP_{25}^{root}(\uparrow)$ |
|----------------------------|----------------------|----------------------------|
| RootNet [21]               | 108.1                | 31.0                       |
| Ours direct regression     | 94.5                 | 27.3                       |
| Ours shared feature branch | 72.0                 | 31.9                       |
| Ours w/o GNN               | 72.9                 | 32.7                       |
| Ours w/o HM pooling        | 71.8                 | 26.0                       |
| <b>Ours (full)</b>         | <b>69.9</b>          | <b>39.4</b>                |





See our project page for more information:  
<https://github.com/jiahaoLjh/HumanDepth>