



ICCV23
PARIS

Zolly: Zoom Focal Length Correctly for Perspective-Distorted Human Mesh Reconstruction

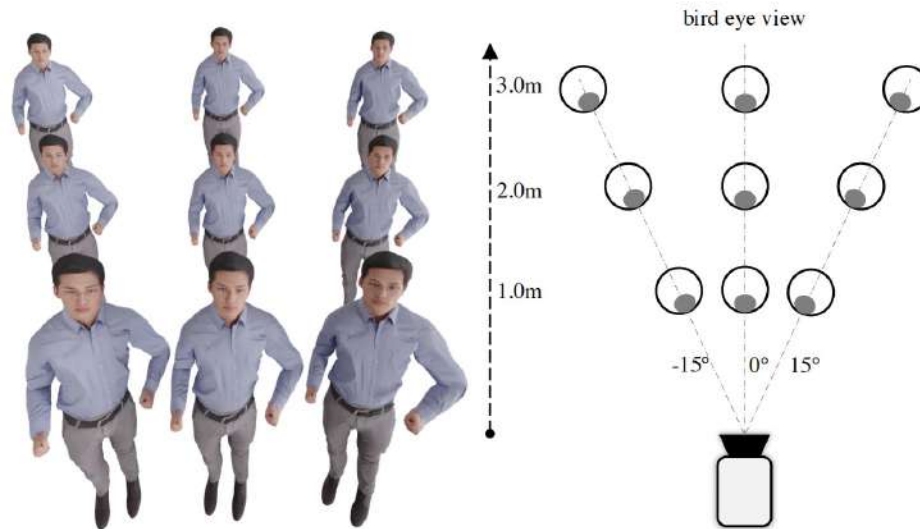
Wenjia Wang^{1,2}, Yongtao Ge³, Haiyi Mei⁴, Zhongang Cai⁴,
Qingping Sun⁴, Yanjun Wang⁴, Chunhua Shen⁵, Lei Yang^{2,4}, Taku Komura¹

¹The University of Hong, ²Shanghai AI Laboratory

³The University of Adelaide, ⁴SenseTime Research, ⁵The Zhejiang University

Motivation

- Existing 3D human mesh reconstruction methods use a constant focal length or estimate one based on the background context.
- Such focal length deviates a lot from distorted images caused by perspective projection. The distortion is directly caused by the distance and the facing angle to the camera center.

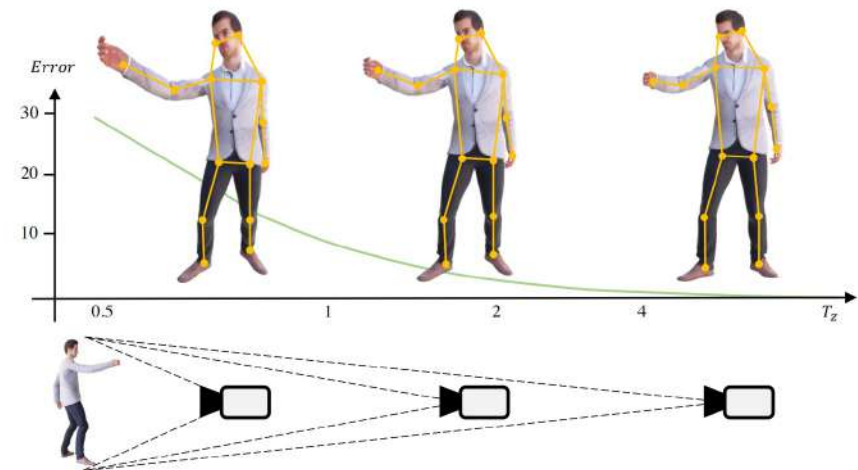


Motivation

- Close-view shots could cause distortion on human bodies, which could be used to calculate the distance. The closer to the camera, the server the distortion. When human body is far enough, the large focal length.
- This is called dolly zoom. We find the inversion of this phenomenon and use it to solve the distance and focal length.



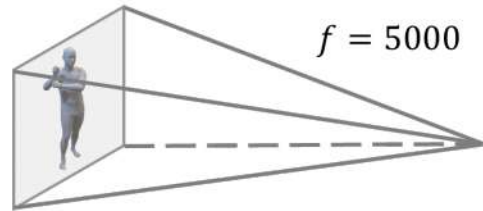
Jaws, a 1975 American thriller film directed by Steven Spielberg



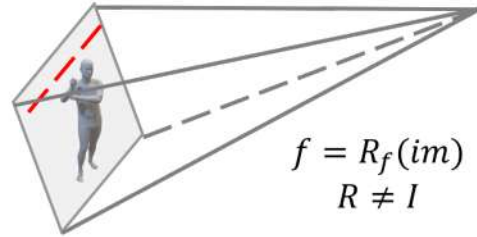
Contribution

- A novel camera system for the perspective-distorted human mesh reconstruction task.
- A new neural model, consists of a translation estimation module, a mesh reconstruction module, and a hybrid re-projection loss.
- A new and the first synthetic dataset, PDHuman, for perspective-distorted 3D human pose estimation.

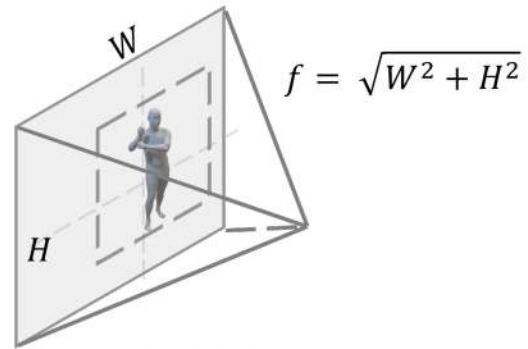
Camera System Design



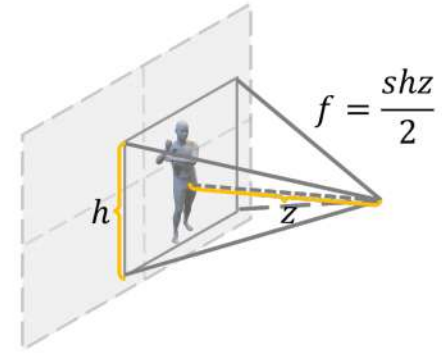
(a). HMR



(b). SPEC



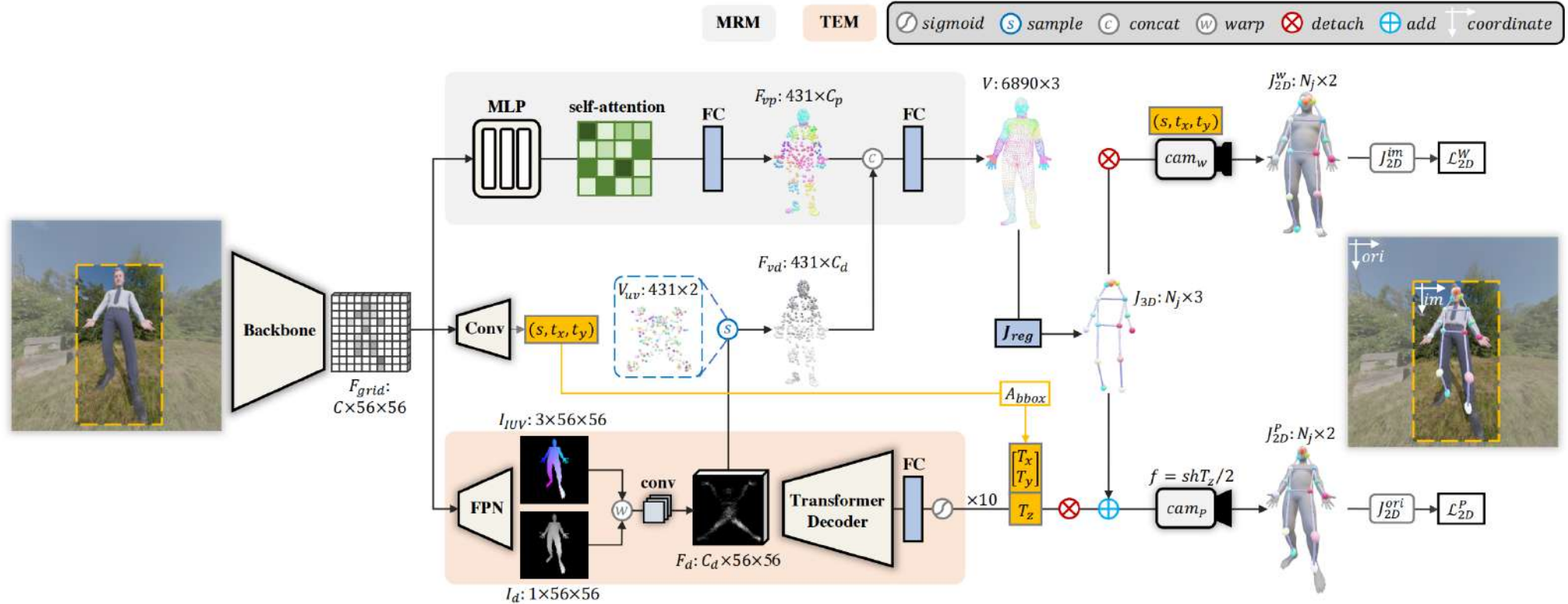
(c). CLIFF



(d). Zolly

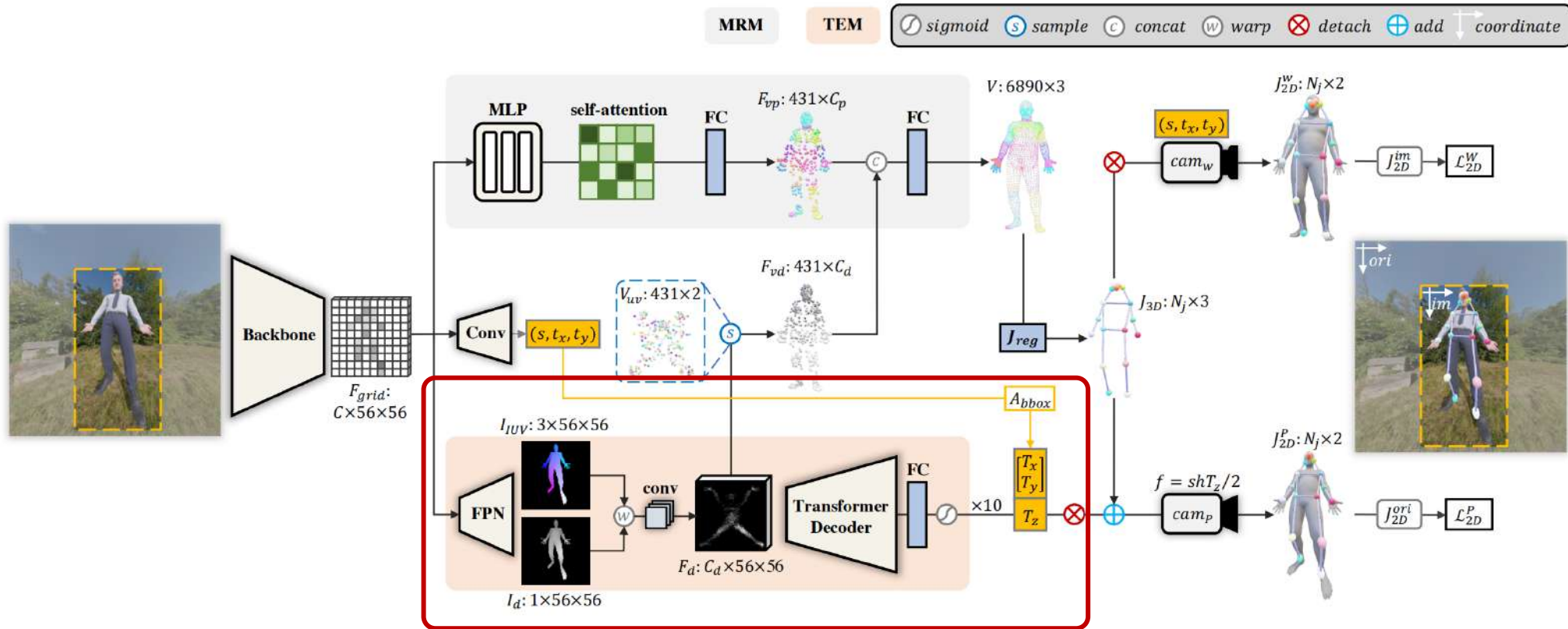
- HMR. $f = 5000$ (*pixels*). Most methods follow this setting.
- SPEC. The f is estimated by a network pre-trained on other datasets.
- CLIFF. Use the diagonal length as f if no ground truth f .
- Zolly. $f = \frac{shT_z}{2}$ (*pixels*). Where T_z is the z-axis distance. ($f = sT_z$ in NDC space.) (Please find detailed information in our Paper)

Proposed Method



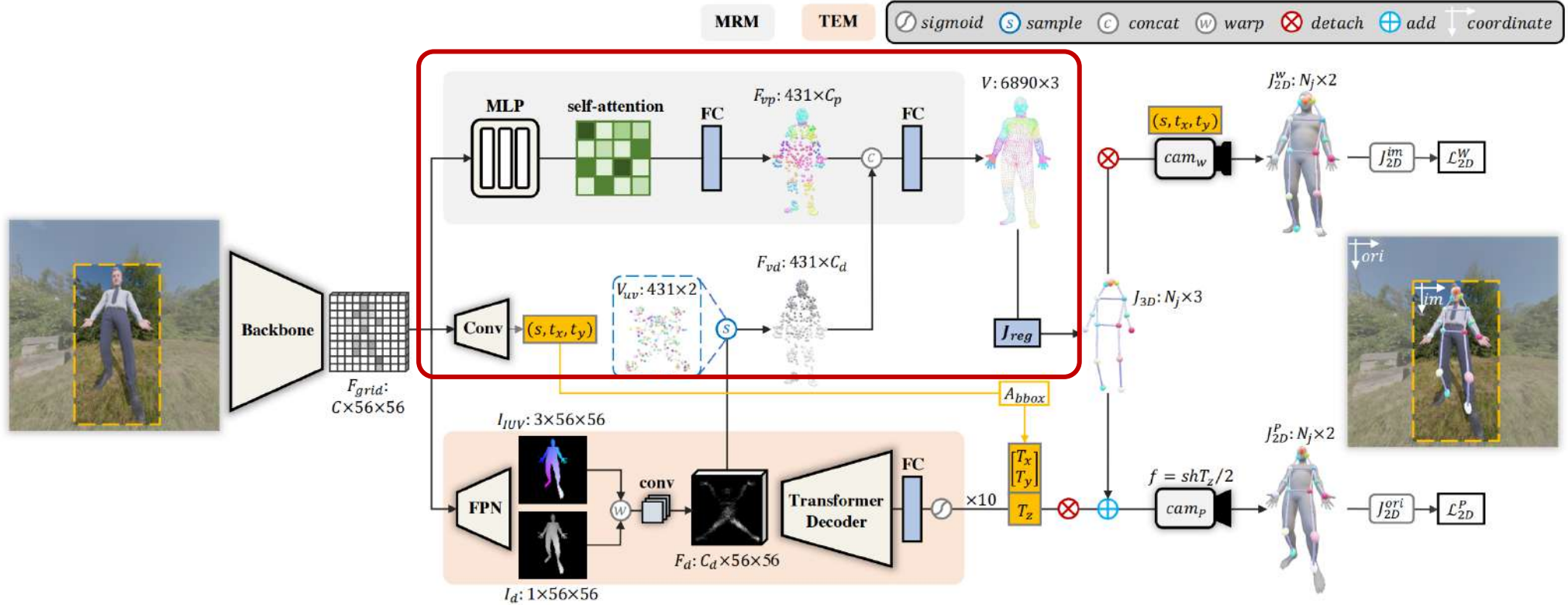
- According to our setting $f = sT_z$, we need to estimated the distance T_z and orthorgraphics scale s . Our pipeline consists of 3 main parts: Translation Estimation Module, Mesh Reconstruction Module, and a hybrid reprojection supervision loss.

Proposed Method



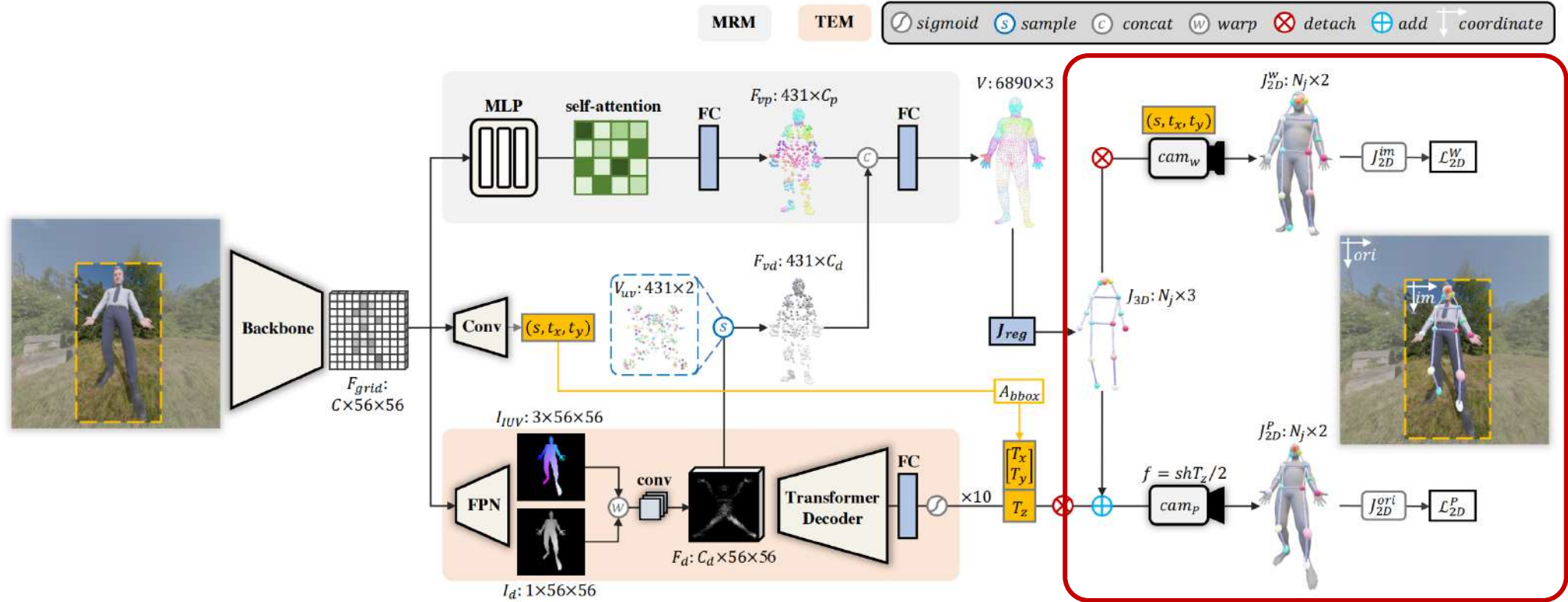
- **Translation Estimation Module.** We use a Transformer to regress the z-axis distance from the warped IUV image and use $10 \times \text{sigmoid}()$ to restrict T_z less than 10m.

Proposed Method



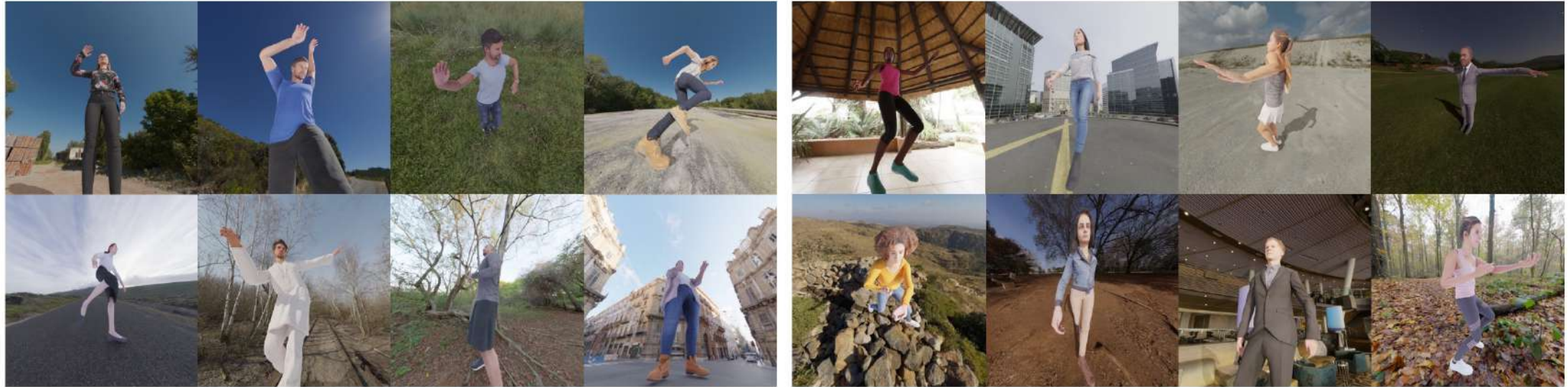
- **Mesh Estimation Module.** We adopt an MLP structure to predict the coordinates of a coarse mesh of the body, then up-sample the mesh using two fully connected layers.

Proposed Method



- Loss Functions. The total loss function is the summation of mesh loss, translation loss, and hybrid re-projection loss.

PDHuman Dataset



We propose the first dataset which aims for perspective-distorted 3D human pose estimation.

- Annotations: Camera intrinsic matrix, 2D/3D keypoints, SMPL parameters β , θ , and translation.
- Camera: Use the dolly-zoom effect to generate random camera intrinsic matrices
- Rendering: Use human models from RenderPeople and body pose sequences from Mixamo, with HDRi images as backgrounds. Use Blender to render the RGB images.
- Amount: 126,198 images in training and 27,448 images in testing split.

Experiments

- We achieve SOTA results on 3DPW test set and comparable results on Human3.6M validation set. Images in 3DPW are actually captured from a close distance, which is suitable for our model. Human3.6M shows that our translation estimation could also generalize to images from a large distance.

| Method | Backbone | 3DPW | | | Human3.6M | |
|--------|----------|-------------|-------------|-------------|-------------|-------------|
| | | PA-JPE | MPJPE | PVE | PA-JPE | MPJPE |
| HMR | Res50 | 72.6 | 116.5 | - | 56.8 | 88.0 |
| SPEC | Res50 | 52.7 | 96.4 | - | | |
| CLIFF | HR48 | 43.0 | 69.0 | 81.2 | 32.7 | 47.1 |
| Zolly | HR48 | 39.8 | 65.0 | 76.3 | 32.3 | 49.4 |

Experiments

- We achieve SOTA results on PDHuman test set and SPEC-MTP. We chose the protocol with most serve distortion from each dataset here. mIoU here is the mIoU between rendered mask and GT masks, which could be used to describe the overlay performance vividly.

| Method | Backbone | PDHuman (p5) | | | SPEC-MTP (p3) | | |
|--------|----------|--------------|-------------|-------------|---------------|--------------|-------------|
| | | PA-JPE | PVE | mIoU | PA-JPE | PVE | mIoU |
| HMR | Res50 | 62.5 | 106.7 | 21.7 | 73.9 | 145.6 | 16.0 |
| SPEC | Res50 | 65.8 | 109.6 | 19.6 | 76.0 | 144.6 | 18.8 |
| CLIFF | HR48 | 66.2 | 115.2 | 24.8 | 74.3 | 132.4 | 23.7 |
| Zolly | HR48 | 49.9 | 82.0 | 26.5 | 67.4 | 126.7 | 30.4 |

Quality Results



(9516, 3.1°, 51.06m) (276, 85.7°, 1.66m) (724, 39.0°, 4.12m) (134, 124.8°, 0.82m) (173, 112.0°, 1.02m)



(12544, 2.5°, 42.07m) (1686, 18.7°, 5.33m) (789, 41.7°, 2.61m) (269, 91.8°, 0.88m) (none, none, none)



(14977, 2.5°, 41.59m) (1385, 27.1°, 3.95m) (945, 39.0°, 2.84m) (669, 53.0°, 1.86m) (none, none, none)



(42692, 2.6°, 55.37m) (1476, 66.1°, 1.87m) (2202, 47.1°, 3.00m) (802, 100.3°, 0.97m) (685, 109.0°, 0.87m)

HMR

SPEC

CLIFF

Zolly

Input



(11428, 2.6°, 33.33m) (521, 52.4°, 1.43m) (724, 39.0°, 2.08m) (330, 75.6°, 0.88m) (330, 75.6°, 0.83m)



(11160, 2.6°, 36.22m) (1151, 18.4°, 3.74m) (707, 46.0°, 2.75m) (415, 48.5°, 1.33m) (none, none, none)



(15089, 2.6°, 43.24m) (922, 23.2°, 2.75m) (833, 39.6°, 2.82m) (352, 56.2°, 1.09m) (none, none, none)



(53359, 2.3°, 43.46m) (273, 32.7°, 2.80m) (2202, 47.1°, 1.91m) (1373, 70.0°, 1.14m) (990, 88.3°, 0.99m)

HMR

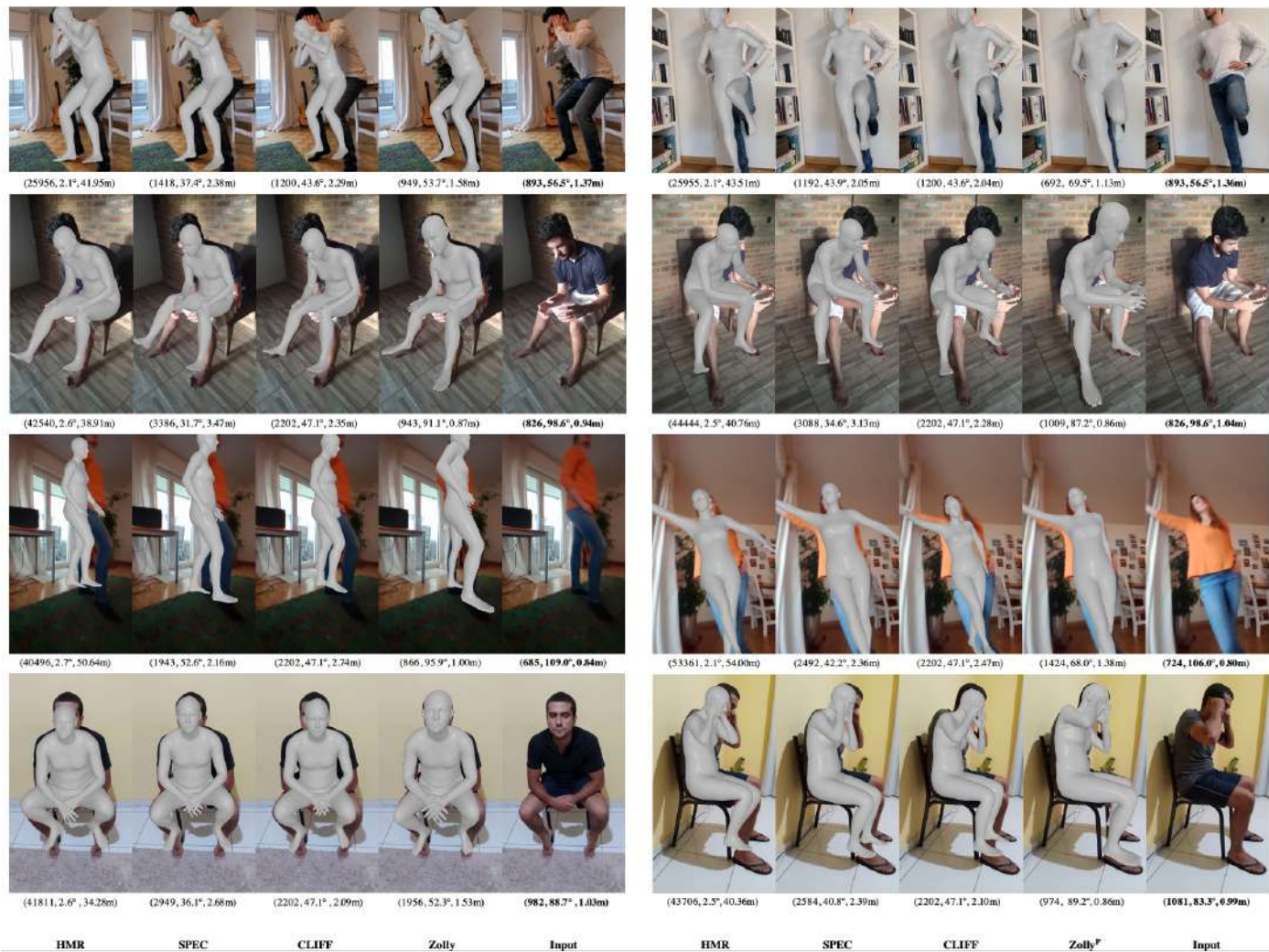
SPEC

CLIFF

Zolly^p

Input

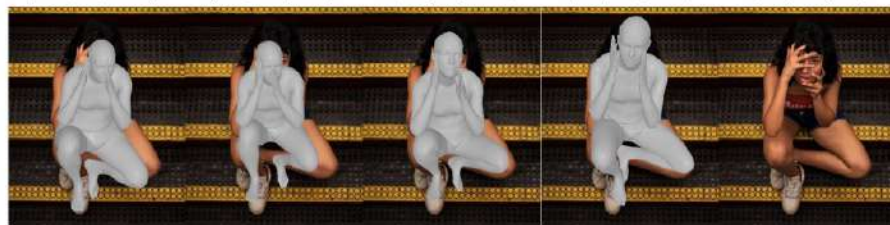
Quality Results on SPEC-MTP



Quality Results on Web Images



(14375, 2.4°, 53.41m) (2174, 15.8°, 8.64m) (882, 37.7°, 3.57m) (494, 62.8°, 1.84m) (none, none, none)



(22321, 2.6°, 29.02m) (2678, 21.1°, 3.90m) (1283, 42.6°, 1.88m) (1538, 36.0°, 1.87m) (none, none, none)



(23816, 2.6°, 55.18m) (833, 65.3°, 2.16m) (1333, 43.6°, 3.38m) (785, 68.4°, 1.55m) (none, none, none)



(64441, 2.6°, 43.65m) (2378, 62.5°, 1.72m) (1283, 96.7°, 1.88m) (2360, 62.9°, 1.60m) (none, none, none)

HMR

SPEC

CLIFF

Zolly

Input



(14486, 2.5°, 50.20m) (1183, 29.9°, 3.63m) (906, 38.5°, 2.94m) (525, 62.1°, 1.60m) (none, none, none)



(14352, 2.6°, 45.45m) (1118, 32.5°, 3.40m) (823, 43.2°, 2.63m) (473, 69.1°, 1.51m) (none, none, none)



(4620, 2.6°, 48.72m) (736, 16.2°, 7.87m) (267, 42.9°, 2.74m) (131, 77.4°, 1.30m) (none, none, none)



(22321, 2.6°, 56.99m) (689, 72.7°, 1.89m) (1280, 43.2°, 3.57m) (399, 103.6°, 1.00m) (none, none, none)

HMR

SPEC

CLIFF

Zolly^v

Input

Quality Results on Standard Benchmark

