

# Perspective Crop Based Egocentric Hand Pose Estimation via Fisheye Stereo Vision

Hyejin Hur<sup>1,2</sup>, Seongmin Baek<sup>1</sup>, Younhee Gil<sup>1</sup>, Sangpil Kim<sup>2</sup>

<sup>1</sup>Electronics and Telecommunications Research Institute (ETRI), South Korea

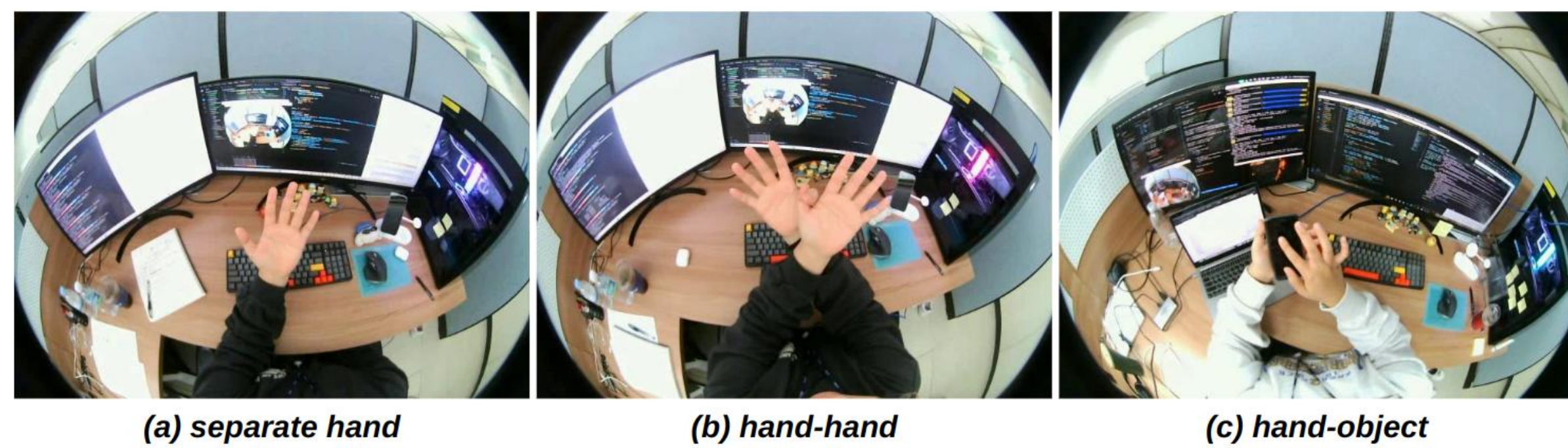
<sup>2</sup>Korea University, South Korea

## INTRODUCTION

With the advancement of Virtual Reality (VR) and Augmented Reality (AR), the importance of hand pose estimation from egocentric view has also increased. To capture diverse hand movements in daily activities, fisheye cameras with a wide Field of View (FoV) are essential. While Fisheye camera offer a broader capture range but they introduce distortion. Thus, this leads to inaccurate hand pose estimation with fisheye cameras. Additionally, most existing hand pose estimation research has focused on third-person view, resulting in limited datasets for egocentric and fisheye-based hand pose estimation.

To address these challenges, we propose two-stage egocentric hand pose estimation method using a fisheye stereo camera. This method generates undistorted hand crop images through perspective cropping [4]. We employed U-Net [3] for HandNet and compared the performance of our model against SimpleBaseline [1] and HRNet [2].

## DATASET



We constructed fisheye camera-based egocentric hand dataset, which we refer to as **FisheyeEgoHAND**. To capture various hand interactions, we designed three categories of scenarios:

- **separate hand** (7 scenarios, 5,142 images)
- **hand-hand** (9 scenarios, 7,008 images)
- **hand-object** (4 scenarios, 4,358 images)

This dataset consists of total 20 scenarios with 16,508 images. For each scenario, videos were recorded at 30 fps, and each frame was extracted as image data.

## RESULTS

Model	Avg EPE ↓	PCK ↑	AUC ↑
SimpleBaseline (r-50)	17.974	0.968	0.890
SimpleBaseline (r-50, w/ PC)	13.178	0.951	0.893
SimpleBaseline (r-101)	18.062	0.958	0.883
SimpleBaseline (r-101, w/ PC)	13.013	0.937	0.894
SimpleBaseline (r-152)	18.296	0.958	0.887
SimpleBaseline (r-152, w/ PC)	13.411	0.927	0.892
HRNet	17.918	0.971	0.890
HRNet (w/ PC)	12.724	0.937	0.896
<b>HandNet (ours)</b>	<b>12.626</b>	<b>0.984</b>	<b>0.917</b>
<b>HandNet (ours, w/ PC)</b>	<b>9.634</b>	<b>0.969</b>	<b>0.914</b>

Table 1: 2D pose estimation on fisheye images and perspective cropped images. "w/PC" refers to "with perspective cropping"

We adopted simplified version of U-Net as HandNet, which enables learning high-resolution representation features essential for hand pose estimation. We trained SimpleBaseline, HRNet, and HandNet for 200 epochs on FisheyeEgoHAND dataset using a single RTX 4090. The batch size was 64 for SimpleBaseline and 128 for HRNet and HandNet, with IoU loss. We used SGD with a 0.1 learning rate and a multi-step scheduler at 40 and 100 epochs ( $\gamma=0.5$ ).

## CONCLUSION AND FUTURE WORK

We propose an effective method for egocentric hand pose estimation using a fisheye camera. We created FisheyeEgoHAND dataset and trained stage-1 HandNet on it for perspective cropping. The undistorted hand crop images were then used to estimate the final 2D hand pose, achieving better performance than one-stage approach. To enable diverse applications in VR and AR, 3D hand pose is more useful but our approach ultimately estimates 2D hand pose. In future work, we plan to explore a distortion-robust 3D hand pose estimation approach.

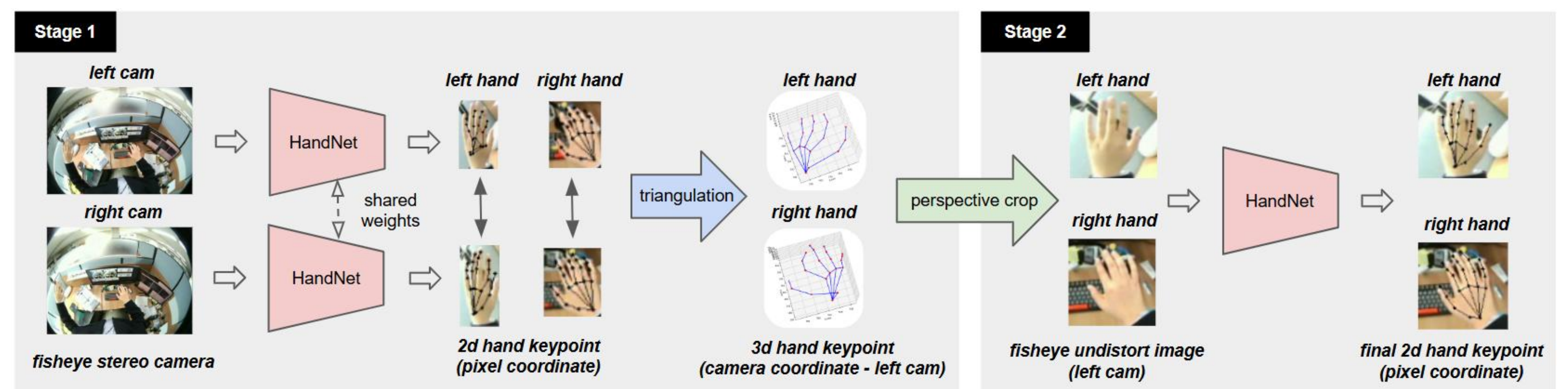
## AFFILIATIONS



## Acknowledge

This work was supported by Electronics and Telecommunications Research Institute (ETRI) grant funded by the Korean government [25ZC1210, Research on hyper-realistic interaction technology for five senses and emotional experience, 100%]. This work was supported by Artificial intelligence industrial convergence cluster development project funded by the Ministry of Science and ICT(MSIR, Korea) & Gwangju Metropolitan City.

## METHODOLOGY



### • 2D hand pose estimation from two fisheye camera views

Using images from fisheye stereo camera as input, stage-1 HandNet predicts the 2D hand keypoints of left/right hands in each view.

### • 3D keypoints via triangulation

The predicted 2D hand keypoints are used for triangulation with camera parameters. 3D hand keypoints are computed in the camera coordinate system with the left (or right) camera as the origin.

### • Perspective cropping

Virtual cameras are generated from 3D keypoints and undistorted hand crop images are obtained through a warping technique.

### • 2D hand pose estimation from undistorted hand crop images

Stage-2 HandNet predicts 2D hand keypoints from undistorted hand crop images obtained in the previous step.

Perspective cropping is a technique that generated a virtual camera for the Region of Interest (RoI) and crops the image accordingly. A virtual camera can be created when 3D points within the RoI are available. We set hand keypoints as the required points for perspective cropping and aimed to obtain 3D points from a stereo camera. To efficiently acquire 3D hand keypoints, we used triangulation for fast and accurate computation. For accurate depth values of these points, precise 2D hand pose estimation in stereo views is crucial for triangulation. Since we use a fisheye stereo camera, training 2D hand pose estimator on distorted images is essential.



Since our proposed method employs perspective cropping, we generated a new ground truth derived from FisheyeEgoHAND. Specifically, the original 2D hand keypoints from FisheyeEgoHAND were projected onto virtual cameras obtained from perspective cropping. This dataset was then used to train stage-2 HandNet in the same way as stage-1.

The results in Table 1. show that the two-stage 2D hand pose estimation with perspective cropping achieved an EPE of 9.634, outperforming the one-stage estimation using the original fisheye camera image as input, which had an EPE of 12.626. In Figure (h), perspective cropping result from stage-1 HandNet trained on FisheyeEgoHAND demonstrating improved distortion correction over Figure (g). Also, Figure (h) shows that the 2D hand pose estimation results improve with perspective cropping.



## REFERENCES

- [1] Xiao B., Wu H., Wei Y.: Simple Baselines for Human Pose Estimation and Tracking. In ECCV (2018).
- [2] Sun K., Xiao B., Liu D., Wang J.: Deep High-Resolution Representation Learning for Human Pose Estimation. In CVPR (2019).
- [3] Ronneberger O., Fischer P., Brox T.: U-Net: Convolutional Networks for Biomedical Image Segmentation. In MICCAI (2015).
- [4] Han S., Wu P., Zhang Y., Liu B., Zhand L., Wang A., Si W., Zhang P., Cai Y., Hodan T., Cabezas R., Tran L., Akbay M., Yu T., Keskin C., Wang R.: UmeTrack: Unified multi-view end-to-end hand tracking for VR. In SIGGRAPH (2022).