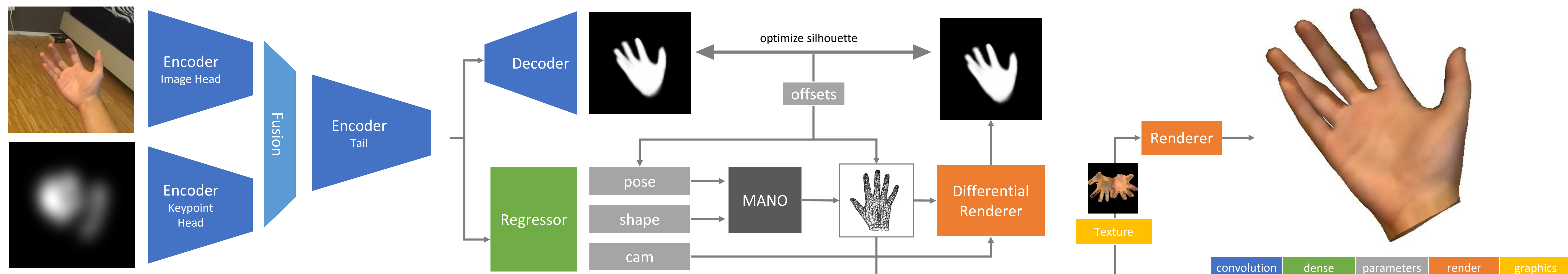


# RealisticHands: A Hybrid Model for 3D Hand Reconstruction

Michael Seeber<sup>1</sup>, Roi Poranne<sup>2</sup>, Marc Pollefeys<sup>1,4</sup>, Martin R. Oswald<sup>1,3</sup>  
<sup>1</sup>ETH Zurich, <sup>2</sup>University of Haifa, <sup>3</sup>University of Amsterdam, <sup>4</sup>Microsoft

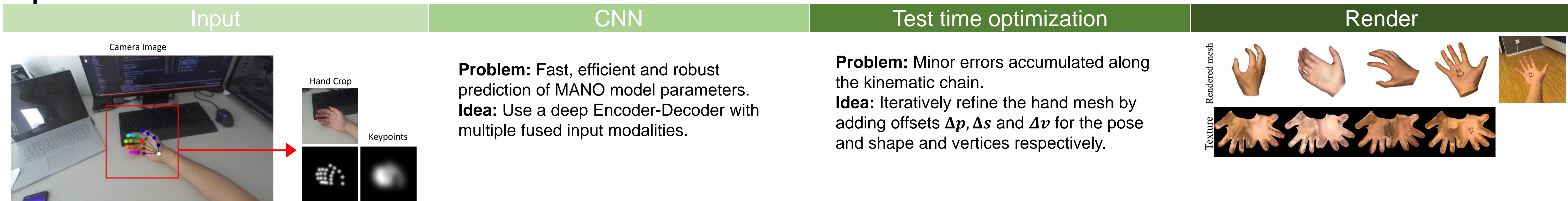
REPLICATING THE USER'S PERSONAL HANDS IN 3D FOR BETTER EMBODIMENT

**Summary:** A novel hybrid approach for 3D hand reconstruction that combines the performance of learning based methods with the accurate image-model alignment of optimization based methods.



A deep encoder-decoder network regresses MANO & camera parameters as well as a segmentation mask. The hand mesh is refined via differential rendering to optimize the silhouette consistency.

## Pipeline

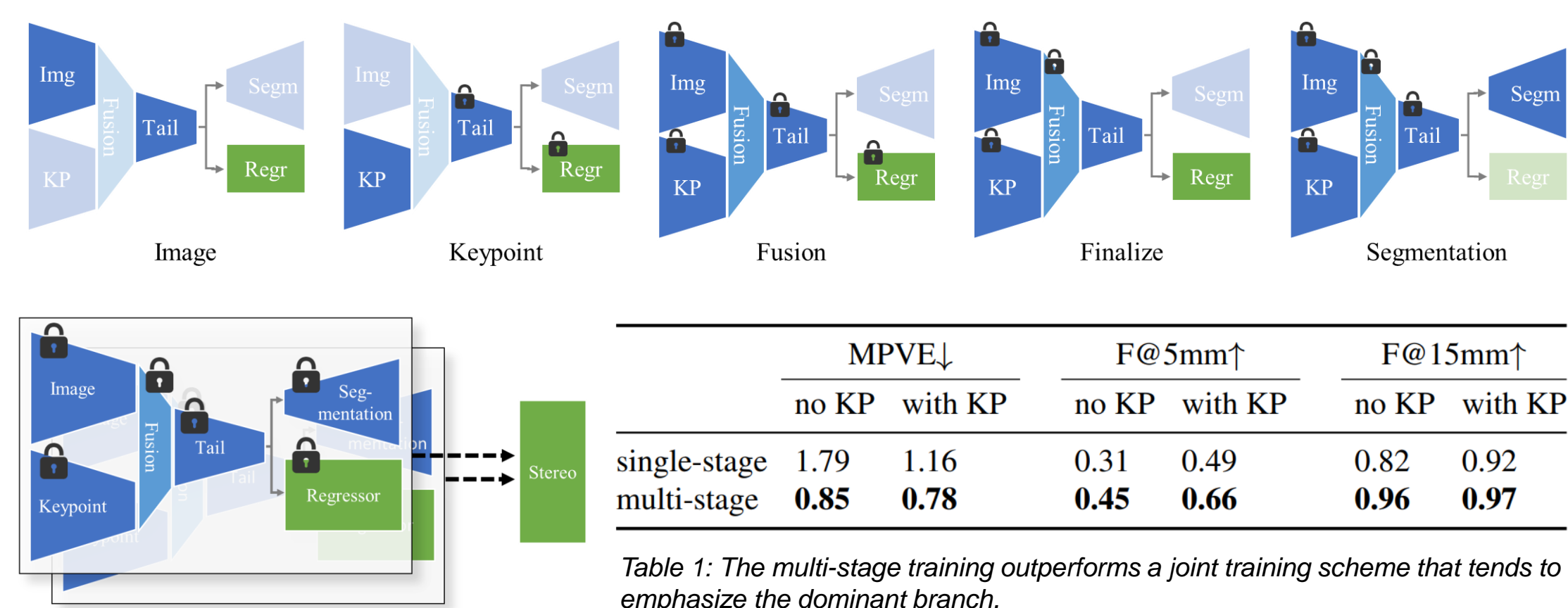


## Motivation & Challenges

- Related Work:**
- 3D hand pose estimation: optimization  $\leftrightarrow$  learning
  - 3D hand reconstruction: model based  $\leftrightarrow$  model free
- Challenges:**
- High DoF:** hands are highly articulated with multiple degrees of freedom
  - Self Similarity:** hard to distinguish between fingers due to similar appearance
  - Occlusion:** especially the egocentric setting comes with lot of self-occlusion
  - Variation of hands:** human hands exhibit a lot of variation in size, shape etc.
  - Fast motion:** Fast moving fingers lead to motion blur & dilution of contours
  - Limited data:** Most datasets on hand shapes include only a few participants.

## Training Framework

In our training scheme selective components are frozen / removed to maximize performance of individual network branches.  
 → Significant performance increase (Table 1)

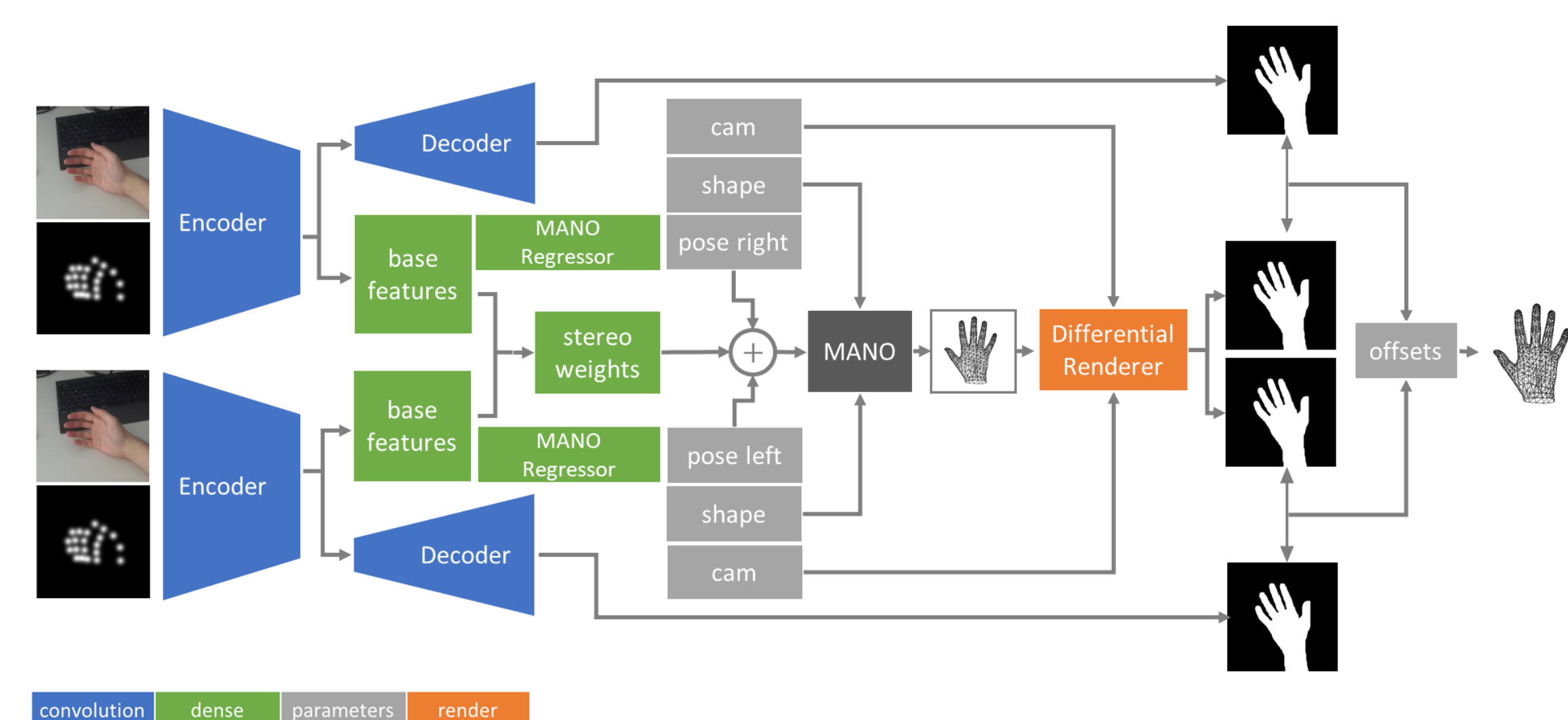


	MPVE↓		F@5mm↑		F@15mm↑	
	no KP	with KP	no KP	with KP	no KP	with KP
single-stage	1.79	1.16	0.31	0.49	0.82	0.92
multi-stage	<b>0.85</b>	<b>0.78</b>	<b>0.45</b>	<b>0.66</b>	<b>0.96</b>	<b>0.97</b>

Table 1: The multi-stage training outperforms a joint training scheme that tends to emphasize the dominant branch.

## Stereo Extension

Estimation a 3D hand pose from a single RGB image is an ill-posed problem due to depth and scale ambiguities. We regress stereo weights  $w$  that are used to fuse the predicted poses from the left and right in a meaningful way.



## Results

### MANO Regression

We outperform baselines, such as the mean shape and inverse kinematic fits, and also rank higher than other model based approaches such as [1, 2, 3, 4]. Further, we perform quantitatively similarly to I2L [5], a model-free approach that occasionally produces collapsed results. The qualitative superiority of our methods over I2L is also well demonstrated in the supplementary video.

Methods	MPVE (PA)↓	F@5mm (PA)↑	F@5mm (PA)↑	Model based	GT scale
Mean shape	1.64	0.336	0.837	✓	
Inverse Kinematics [1]	1.37	0.439	0.892	✓	
Hasson <i>et al.</i> [2]	1.33	0.429	0.907	✓	✓
Boukhayma <i>et al.</i> [3]	1.32	0.427	0.894	✓	✓
Mano CNN [1]	1.09	0.516	0.934	✓	✓
Kulon <i>et al.</i> [4]	0.86	0.614	0.966	✗	✗
I2L [5]	0.76	0.681	0.973	✗	✗
Ours (Mono, MP Hands)	0.97	0.575	0.949	✓	✗
Ours (Mono, PoseNet(I2L))	0.78	0.662	0.971	✓	✗

Table 2: Quantitative comparison of our approach with other methods on the task of monocular hand pose and shape estimation.

### Segmentation

We compared the segmentation performance to various baseline such as HSV, CNN [6] and UMA [6].

	Combined		FreiHAND		Synth	
	IoU↑	Accuracy↑	IoU↑	Accuracy↑	IoU↑	Accuracy↑
HSV	0.26	0.81	0.33	0.90	0.18	0.72
CNN [6]	0.85	0.98	0.77	0.97	0.93	0.99
UMA [6]	0.86	0.98	0.78	0.98	0.93	0.99
Ours	<b>0.88</b>	<b>0.99</b>	<b>0.80</b>	<b>0.98</b>	<b>0.95</b>	<b>0.99</b>

Table 3: Quantitative segmentation results on the FreiHAND and Synthetic Stereo dataset.

### Test Time Optimization

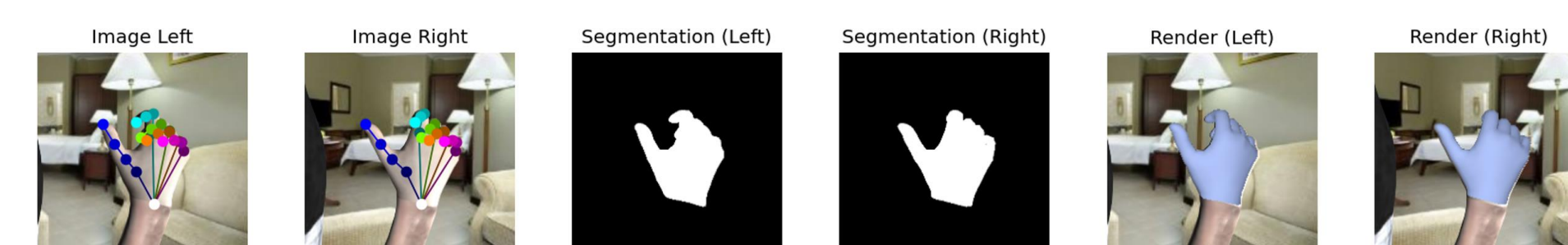
To evaluate the benefit of the optimization stage, we compare the IoU of projected silhouettes with the segmentation mask.

### Stereo Extension

The stereo version greatly outperforms the monocular, showing that multiple views can resolve ambiguities.

### Synthetic Stereo Hands Dataset

Stereo hand datasets from an egocentric viewpoint are lacking.  
 → We generated a large scale synthetic dataset



## References

- [1] Christian Zimmermann, Duygu Ceylan, Jimei Yang, Bryan Russell, Max Argus, and Thomas Brox. FreiHAND: A dataset for markerless capture of hand pose and shape from single rgb images.
- [2] Yana Hasson, Gui Varol, Dimitris Tzionas, Igor Kalevatykh, Michael J. Black, Ivan Laptev, and Cordelia Schmid. Learning joint reconstruction of hands and manipulated objects.
- [3] Adnane Boukhayma, Rodrigo de Bam, and Philip HS Torr. 3d hand shape and pose from images in the wild.
- [4] Dominik Kulon, Riza Alp Guler, Iasonas Kokkinos, Michael M Bronstein, and Stefanos Zafeiriou. Weaklysupervised mesh-convolutional hand reconstruction in the wild.
- [5] Gyeongsik Moon and Kyoung Mu Lee. I2L-meshnet: Image-to-level prediction network for accurate 3d human pose and mesh estimation from a single rgb image.
- [6] Minjie Cai, Feng Lu, and Yoichi Sato. Generalizing hand segmentation in egocentric videos with uncertainty-guided model adaptation.



Figure 1: Qualitative results on the FreiHAND test set in comparison to I2L [5].

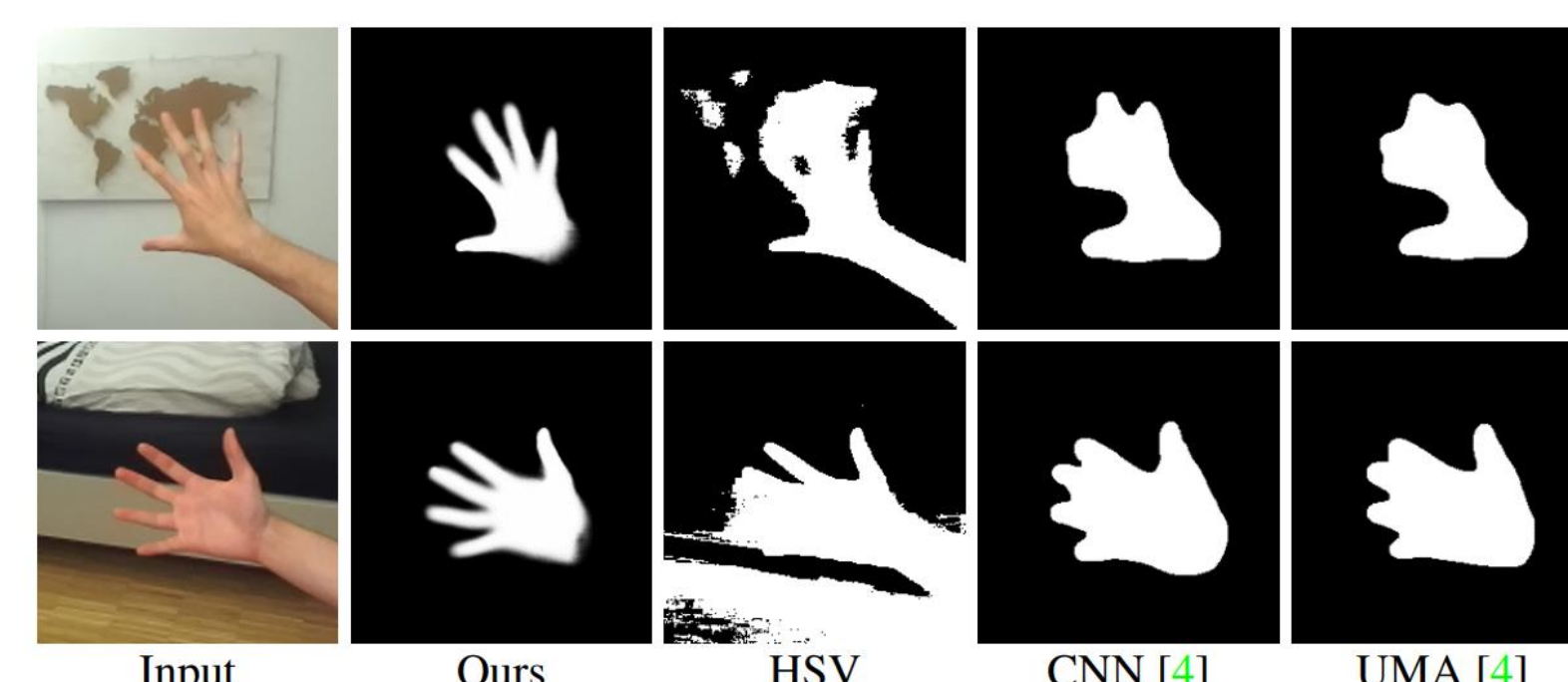


Figure 2: Qualitative segmentation results showing our method performs best.

	IoU↑	
	MP Hands	PoseNet
no refinement	0.677	0.722
3 iterations	0.706	0.758
10 iterations	0.735	0.793
15 iterations	<b>0.748</b>	<b>0.806</b>
I2L [27]		0.737

Figure 3: The test time optimization improves the image-model alignment.

Version	MPVE (PA)↓	F@5mm (PA)↑	F@15mm (PA)↑
Monocular	1.65	0.352	0.860
Stereo	<b>0.94</b>	<b>0.578</b>	<b>0.967</b>

Figure 4: Quantitative comparison results of monocular architecture and stereo architecture.