



# Training a Feedback Loop for Hand Pose Estimation

Markus Oberweger, Paul Wohlhart and Vincent Lepetit

Institute for Computer Vision and Graphics Graz University of Technology

### Motivation





- Important for Human Computer Interaction, Augmented Reality
- Hot topic: [Tang et al., ICCV'15], [Tzionas&Gall, ICCV'15], [Li et al., ICCV'15], [Supančič et al., ICCV'15], [Choi et al., ICCV'15], [Rogez et al., ICCV'15], [Xu&Cheng, ICCV'13], [Tompson et al., ToG'14], [Tang et al., ICCV'13], [Sridhar et al., ICCV'13], [Sharp et al., CHI'15], [Qian et al., CVPR'14], [Oikonomidis et al., ICCV'11], [Melax et al., GIC'13], [Kuznetsova et al., ICCV'13], [Keskin et al., ICCV'11], [de La Gorce et al., PAMI'11], [Ballan et al., ECCV'12], Leap Motion
- Our goal: Accurate 3D pose from single depth image 2

### Challenges

- Self-occlusions
- Self-similarity
- Many degrees-of-freedom
- Noisy input



### A Predictor



#### Starting a Feedback Loop: Training a CNN to synthesize depth images from



## Iterating the Feedback Loop



### The Predictor

• Simple and very fast CNN for pose initialization



#### The Synthesizer



#### The Synthesizer



# Directly using the Synthesizer for Optimization?

 $\widehat{\text{pose}} = \arg\min \|\text{depth} - \text{synthesizer}(\text{pose}))\|_2^2$ 

pose

Initialization

**Direct minimization** 

A C

Predicting updates (our method)

### The Updater



## The Updater

Predict updates to increase accuracy of pose



Get closer

 $\|\mathbf{p} + \mathrm{updater}(\mathrm{depth}, \mathrm{synthesizer}(\mathbf{p})) - \mathbf{p}_{\mathrm{GT}}\|_2 < \lambda \|\mathbf{p} - \mathbf{p}_{\mathrm{GT}}\|_2$ 

### The Updater

• Training is actually more complex...

 $\arg\min_{\Omega}\sum_{(\mathcal{D},\mathbf{p})\in\mathcal{T}}\sum_{\mathbf{p}'\in\mathcal{T}_{\mathcal{D}}}\max(0,\|\mathbf{p}'+\mathrm{updater}_{\Omega}(\mathcal{D},\mathrm{synth}(\mathbf{p}'))-\mathbf{p}\|_{2}-\lambda\|\mathbf{p}'-\mathbf{p}\|_{2})$ 







#### Our Results on NYU Dataset





# Conclusion

- Everything is learned: No hand-crafted 3D model and similarity measure needed!
- Trained feedback loop provides updates for pose
- Predicting updates much more robust than minimizing image difference
- Efficient implementation possible (~400fps on GPU)
- General approach, applicable to other pose estimation problems





# Thanks! Questions?

{oberweger,wohlhart,lepetit}@icg.tugraz.at