



## Goal

Realistic 3D hand-object reconstruction from **monocular images**.

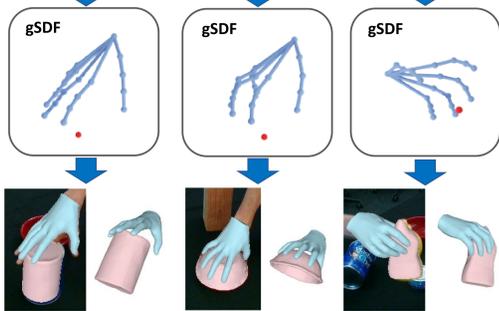


## Contribution

- Leverage **hand pose estimation** as a guidance for SDF-based 3D hand-object reconstruction.
- Use monocular **videos** to alleviate **occlusion and motion blur** issues and improve the performance.

## Motivation

- Deep SDFs can generalize to different shape resolutions but lack explicit modeling of the underlying 3D geometry.
- 3D hand-object reconstruction from a single RGB image is intrinsically hard, especially under occlusion or motion blur.

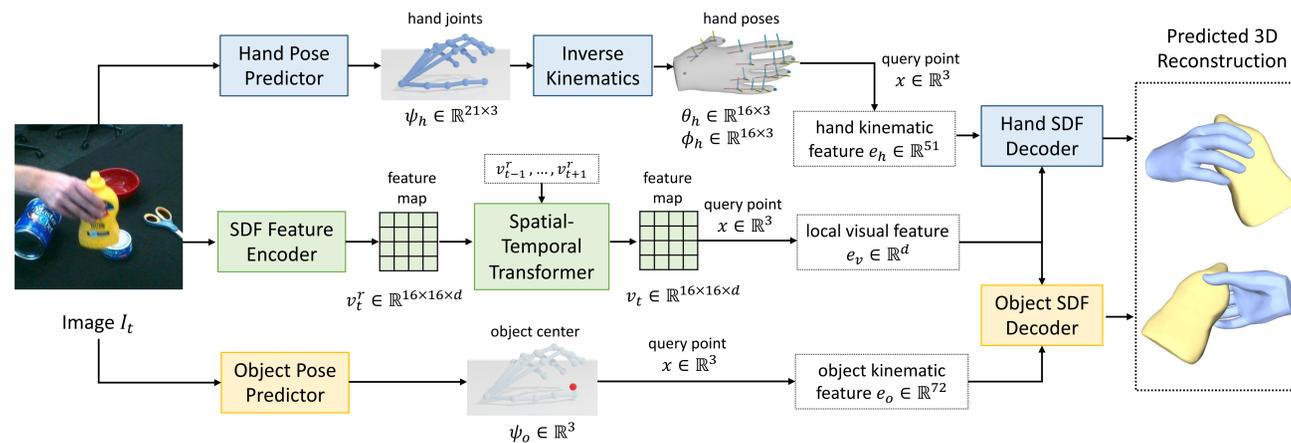


## Related work

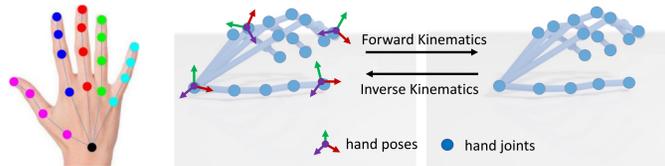
- [1] Y. Hasson, G. Varol, D. Tzionas, I. Kalevtykh, M. Black, I. Laptev, and C. Schmid. Learning joint reconstruction of hands and manipulated objects. *In Proc. CVPR, 2019*.
- [2] K. Karunratanakul, J. Yang, Y. Zhang, M. Black, K. Muandet, and Siyu Tang. Grasping Field: Learning Implicit Representations for Human Grasps. *In Proc. 3DV, 2020*.
- [3] Y. Ye, A. Gupta, and S. Tulsiani. What's in your hands? 3D Reconstruction of Generic Objects in Hands. *In Proc. CVPR, 2022*.
- [4] Z. Chen, Y. Hasson, C. Schmid, I. Laptev. AlignSDF: Pose-Aligned Signed Distance Fields for Hand-Object Reconstruction. *In Proc. ECCV, 2022*.

## Approach

We formulate the joint hand-object 3D reconstruction task as a multi-task learning framework.

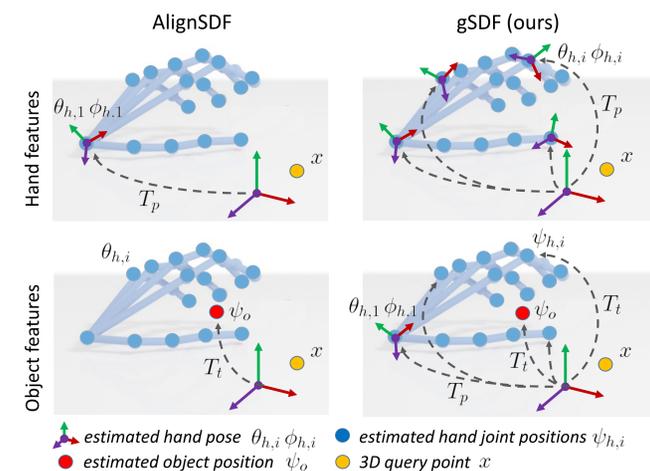


### Pose Estimation



- (a) Hand Joints
- As shown in (a), since deep CNNs is good at detecting interested points, we first use neural networks to predict 3D hand joint locations from single-view images.
- (b) Hand Kinematics
- As shown in (b), we use inverse kinematics to recover the pose transformations for each hand bone.

### Kinematic Features

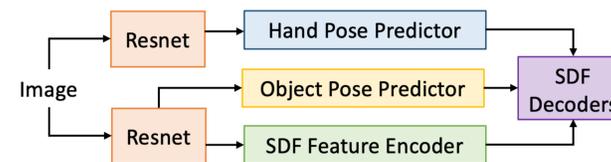


- For hand kinematic features, compared with a recent work [4], we use full kinematic chains of pose transformations.
- For object kinematic features, instead of only considering the object translation as in [4], we additionally consider relative positions between the query point  $x$  and each hand joint.

### Visual Feature

- We use the spatial and temporary transformer to aggregate features from multiple frames.
- Then, we project the query point  $x$  onto the plane of the feature map and obtain the refined local feature for the shape reconstruction.

### Network Architecture



- Our model consists of hand pose predictor, object pose predictor, SDF feature encoder and SDF decoders.
- We use two backbones to handle the task of 3D shape reconstructions and the task of pose predictions separately.
- We observe that our model can achieve the best performance when the object pose predictor and SDF feature encoder shares the same backbone.

## Results

We validate the method by conducting experiments on ObMan and DexYCB benchmarks. We employ metrics including Chamfer Distance (CD) and F-score (FS) to evaluate the quality of results.

### Quantitative Comparison on ObMan

Methods	CD <sub>h</sub> ↓	FS <sub>h</sub> @1 ↑	FS <sub>h</sub> @5 ↑	CD <sub>o</sub> ↓	FS <sub>o</sub> @5 ↑	FS <sub>o</sub> @10 ↑	E <sub>h</sub> ↓	E <sub>o</sub> ↓
Hasson <i>et al.</i> [1]	0.415	0.138	0.751	3.60	0.359	0.590	1.13	-
Karunratanakul <i>et al.</i> [2]	0.261	-	-	6.80	-	-	-	-
Ye <i>et al.</i> [3]	-	-	-	-	0.420	0.630	-	-
Chen <i>et al.</i> [4]	0.136	0.302	0.913	3.38	0.404	0.636	1.27	<b>3.29</b>
gSDF (Ours)	<b>0.112</b>	<b>0.332</b>	<b>0.935</b>	<b>3.14</b>	<b>0.438</b>	<b>0.660</b>	<b>0.93</b>	3.43

### Quantitative Comparison on DexYCB

Methods	CD <sub>h</sub> ↓	FS <sub>h</sub> @1 ↑	FS <sub>h</sub> @5 ↑	CD <sub>o</sub> ↓	FS <sub>o</sub> @5 ↑	FS <sub>o</sub> @10 ↑	E <sub>h</sub> ↓	E <sub>o</sub> ↓
Hasson <i>et al.</i> [1]	0.537	0.115	0.647	1.94	0.383	0.642	1.67	-
Karunratanakul <i>et al.</i> [2]	0.364	0.154	0.764	2.06	0.392	0.660	-	-
Chen <i>et al.</i> [4]	0.358	0.162	0.767	1.83	0.410	0.679	1.58	<b>1.78</b>
Chen <i>et al.</i> [4] <sup>††</sup>	0.344	0.167	0.776	1.81	0.413	0.687	1.57	1.93
gSDF (Ours)	<b>0.302</b>	<b>0.177</b>	<b>0.801</b>	<b>1.55</b>	<b>0.437</b>	<b>0.709</b>	<b>1.44</b>	1.96

### Qualitative Results

