

---

# OnePose++: Keypoint-Free One-Shot Object Pose Estimation without CAD Models

---

Xingyi He<sup>1\*</sup> Jiaming Sun<sup>2\*</sup> Yuang Wang<sup>1</sup> Di Huang<sup>3</sup>  
Hujun Bao<sup>1</sup> Xiaowei Zhou<sup>1†</sup>

<sup>1</sup>Zhejiang University

<sup>2</sup>Image Derivative Inc.

<sup>3</sup>The University of Sydney

## 1 Compare with Related Works

### 1.1 Compare with PatchFlow.

Different from PatchFlow [3] that performs complicated multiview refinement by interpolating previously estimated two-view local patch flows, the proposed refinement strategy is simple yet effective in achieving consistent matches for 3D model refinement. Given a coarse 3D model, we keep the selected reference node of each feature track fixed and search around each query node for the fine-level match by the transformer-based matching module. The advantages are that our graph structure is significantly simpler than PatchFlow, which is efficient for matching, and we do not need to store and interpolate dense flow fields for optimization.

### 1.2 Compare with PixSfM.

The main difference between our refinement module and PixSfM [8] is that we leverage fine-level matching with the transformer to refine the 2D locations of coarse feature tracks and then optimize the 3D model with geometric error, while PixSfM uses pre-stored dense feature maps and feature-metric BA to refine the 3D model and 2D keypoints globally. Our advantage is the more accurate refinement thanks to the transformed features. Moreover, we are more memory efficient by avoiding feature or cost map storage in the BA stage. Please refer to Sec. 5.3 for more qualitative and quantitative comparison results.

## 2 Method Details

### 2.1 Keypoint-Free SfM

**Reference nodes selection strategy.** The proposed keypoint-free SfM establishes 3D structures in a coarse-to-fine manner. It first reconstructs a complete 3D model leveraging the semi-dense matches of the coarse-level LoFTR, then refines the initial 3D model to higher accuracy. We refine the initial point cloud by refining “keypoints” in the coarse feature tracks to sub-pixel accuracy and optimize the point cloud based on the refined feature tracks. The fine-level LoFTR is used for coarse feature track refinement, which refines all points in a feature track with Transformers with reference to a fixed reference point.

We find that the fine-level LoFTR is robust and insensitive to the reference point selection strategy. This is reasonable since the fine-level LoFTR is exactly trained to find a sub-pixel correspondence on a local feature patch for an arbitrarily given feature point. Consequently, we design the reference

---

\*The first two authors contributed equally. The authors from Zhejiang University are affiliated with the State Key Lab of CAD&CG and the ZJU-SenseTime Joint Lab of 3D Vision.

†Corresponding author.

node selection strategy mainly for the ease of implementation and lower memory consumption. To refine all feature tracks with the fine-level LoFTR, we need to extract fine-level CNN feature maps for all images and store them for further use. This would take a ton of storage, varied according to the total number of frames, which can hardly fit into the RAM of consumer-grade GPUs.

To make our reconstruction pipeline broadly usable, we treat an image, instead of a feature track, as the minimum processing unit. More specifically, we recursively select the frame containing the maximum number of ‘‘keypoints’’ (i.e., involved in the most feature tracks), and select all ‘‘keypoints’’ in this frame as the fixed reference nodes of their belonging feature tracks. All feature tracks in the selected frame are then refined, and marked as processed. We repeat this process until all feature tracks are refined. This strategy avoids the need to store dense feature maps of all frames and minimizes the number of frames whose feature maps are repeatedly extracted.

## 2.2 Object Pose Estimation

**Positional encoding.** We apply positional encoding modules on top of the coarse 3D features  $\tilde{\mathbf{F}}_{3D} \in \mathbb{R}^{N \times D_c}$  and 2D feature maps  $\tilde{\mathbf{F}}_{2D}$  to make them position-dependent, which is proved to boost the matching performance [11, 14]. Because the 2D-3D matching involves two modalities, we use the standard sinusoidal encoding for the 2D feature maps and leverage a learned positional encoding for the 3D features. More specifically, for the 3D features, we embed the normalized 3D coordinates  $\mathbf{x}_{3D} \in \mathbb{R}^{N \times 3}$  into a high-dimensional vector with an MLP:

$$\tilde{\mathbf{F}}'_{3D} = \tilde{\mathbf{F}}_{3D} + \text{MLP}_{\text{pe}}(\mathbf{x}_{3D}). \quad (1)$$

For the 2D feature maps, we use a 2D extension of the standard sinusoidal positional encoding proposed in Transformers following DETR [1]:

$$\mathcal{PE}_{x,y}^i = f(x,y)^i := \begin{cases} \sin(\omega_k \cdot x), & i = 4k \\ \cos(\omega_k \cdot x), & i = 4k + 1 \\ \sin(\omega_k \cdot y), & i = 4k + 2 \\ \cos(\omega_k \cdot y), & i = 4k + 3, \end{cases} \quad (2)$$

where  $\omega_k = \frac{1}{10000^{2k/d}}$ ,  $d$  is the number of feature channels on which positional encoding is applied,  $i$  is the index of the feature channel.

The positional encoding modules enable the later attention modules to jointly reason about visual appearances and positions, benefiting 2D-3D matching. Note that the positional encodings are only applied once before the first attention module.

**Attention module.** Directly using the vanilla Transformer [16] to our model is not applicable because its computation cost grows quadratically with the length of input features. Following [14], we use the Linear Transformer [4] to efficiently transform 2D and 3D features. It reduces the computational complexity of the Transformer [16] from  $O(N^2)$  to  $O(N)$  by substituting the exponential kernel with an alternative kernel function  $\text{sim}(Q, K) = \phi(Q) \cdot \phi(K)^T$ , where  $\phi(\cdot) = \text{elu}(\cdot) + 1$ . Please refer to the original paper [4] for more details.

We denote a set of self- and cross-attention layers as an attention block:

$$\begin{cases} \mathbf{F}'_{2D}{}^{(l+1)} = \text{SelfAtten}(\mathbf{F}_{2D}^{(l)}, \mathbf{F}_{2D}^{(l)}), \\ \mathbf{F}'_{3D}{}^{(l+1)} = \text{SelfAtten}(\mathbf{F}_{3D}^{(l)}, \mathbf{F}_{3D}^{(l)}), \\ \mathbf{F}_{2D}^{(l+1)}, \mathbf{F}_{3D}^{(l+1)} = \text{CrossAtten}(\mathbf{F}'_{2D}{}^{(l+1)}, \mathbf{F}'_{3D}{}^{(l+1)}). \end{cases} \quad (3)$$

The indices of intermediate features are indicated by  $\cdot^{(l)}$ .  $\mathbf{F}'$  represents an intermediate feature processed by a self-attention layer. Our attention module sequentially performs the attention block  $N_c = 3$  times to transform the 3D and 2D features.

**Supervision.** We jointly train the coarse and fine modules in our 2D-3D matching framework with different supervisions. We project the observable 3D points to the 2D frame to build the ground-truth 2D-3D correspondences  $\mathcal{M}_{gt}^f$  for our fine-level matching module. For the coarse matching module, we round the projected 2D points to their nearest grid points to obtain the ground-truth coarse 2D-3D correspondences  $\mathcal{M}_{gt}^c$ . We optimize the coarse module by minimizing the focal loss [7] between the



Figure 1: **CAD models in the proposed OnePose-LowTexture dataset.** We capture CAD models for a subset of eight objects from the OnePose-LowTexture dataset. These CAD models can be used to train *Instance-level* methods such as PVNet [9] and CDPN [6] and enable further comparisons between CAD-model-free methods and CAD-model-based methods.

predicted matching probability matrix  $\mathcal{P}^c$  and the ground truth  $\mathcal{P}_{gt}^c$  constructed with  $\mathcal{M}_{gt}^c$  similar to [11, 14]:

$$\mathcal{L}_c = \frac{1}{|\mathcal{P}_{gt}^c|} \sum_{j,q} \text{FL}(\mathcal{P}^c(j, q)), \quad (4)$$

$$\text{FL}(\mathcal{P}^c(j, q)) = \begin{cases} -\alpha(1 - \mathcal{P}^c(j, q))^\gamma \log(\mathcal{P}^c(j, q)), & \text{if } \mathcal{P}_{gt}^c(j, q) = 1 \\ -(1 - \alpha)\mathcal{P}^c(j, q)^\gamma \log(1 - \mathcal{P}^c(j, q)), & \text{if } \mathcal{P}_{gt}^c(j, q) \neq 1 \end{cases}$$

For the fine module, we use a  $\ell_2$  loss to minimize the distances between the predicted 2D coordinates  $\hat{\mathbf{u}}^q$  and the ground truth  $\hat{\mathbf{u}}_{gt}^q$ . Following [17, 14], we make our loss uncertainty-weighted with a variance term  $\sigma^2(q)$ :

$$\mathcal{L}_f = \frac{1}{|\mathcal{M}_f|} \sum_{q \in \mathcal{M}_f} \frac{1}{\sigma^2(q)} \|\hat{\mathbf{u}}^q - \hat{\mathbf{u}}_{gt}^q\|_2. \quad (5)$$

Notably, we detach  $\sigma^2(q)$  during training to prevent the network from decreasing the loss by increasing the variance. The total loss is the weighted sum of the coarse and fine losses  $\mathcal{L} = \omega_c \mathcal{L}_c + \omega_f \mathcal{L}_f$ . In the experiment,  $\alpha$  is 0.5,  $\gamma$  is 2.0,  $\omega_c$  is 1.0 and  $\omega_f$  is 1.0.

### 3 OnePose-LowTexture Dataset

In this section, we provide more details of the proposed OnePose-LowTexture evaluation dataset. This test dataset is used for the evaluation, and there are no training objects. It contains 80 sequences of 40 household low-textured objects. For each object, two video sequences with object poses and annotated object 3D bounding boxes are provided. Video sequences of each object are captured with different backgrounds, simulating the real-world using scenario. Each video is recorded at 30 fps for about 30 seconds in  $1920 \times 1440$  resolution. The total number of images in the reference sequences is 35521, and the total number of images in the query sequences is 32477.

The data capture and annotation pipeline follow the setup of OnePose [15]. The camera poses provided by ARKit can be transformed into the object-centric coordinate system induced from the user-annotated object 3D bounding boxes. Following [15], we align multiple captured sequences of an object with the annotated object 3D bounding boxes. Then, we perform a bundle adjustment with COLMAP to reduce the pose drift of ARKit and inconsistency between 3D bounding box annotations in multiple sequences. This offline optimization process leads to more consistent 3D bounding box annotations across sequences and more accurate object poses.

To compare with instance-level methods such as PVNet [9] and CDPN [6], we additionally capture high-quality 3D CAD models for a selected subset of ten objects from the OnePose-LowTexture dataset. We use the SHINING<sup>(R)</sup> scanner for the CAD model capturing. Fig. 1 illustrates all captured CAD models.

## 4 Experiment Details

### 4.1 Training Details

Our model is trained on the OnePose [15] training set, which contains 49 objects. We first reconstruct the semi-dense object point cloud with our keypoint-free SfM for each object using all training

Table 1: Use feature matching as the 2D object detector.

2D Detector Type	Method	OnePose dataset		
		1cm-1deg	3cm-3deg	5cm-5deg
GT box	Ours	<b>51.1</b>	<b>80.8</b>	<b>87.7</b>
	OnePose	49.7	77.5	84.1
Feature matching detector	Ours	<b>49.6</b>	<b>80.0</b>	<b>87.1</b>
	OnePose	47.3	77.0	83.9

Table 2: Compare with HLoc using different feature matching methods.

Method	OnePose-LowTexture dataset		
	1cm-1deg	3cm-3deg	5cm-5deg
Ours	<b>16.8</b>	<b>57.7</b>	<b>72.1</b>
HLoc ( <i>LoFTR</i> )	15.4	43.7	53.4
HLoc ( <i>SPP + SPG + Patch2Pix</i> )	10.1	37.2	47.6
HLoc ( <i>SPP + SPG</i> )	13.8	36.1	42.2
HLoc ( <i>DualRC-Net</i> )	11.3	37.0	47.8
HLoc ( <i>NC-Net + Patch2Pix</i> )	2.42	19.0	30.4

sequences with 5x downsampled video frames. Then we leverage the 2D-3D correspondences computed from the annotated poses and the reconstructed 3D model to train our sparse-to-dense 2D-3D matching module. Note that we compute the rough 2D object bounding boxes from the annotated 3D bounding boxes in the dataset and use the cropped images for reconstruction and training, following OnePose [15].

## 4.2 Metrics

We use the commonly used *cm-degree* pose error, the *ADD(S)* and the *Proj2D* metrics to evaluate the estimated object poses. We follow PixSfM [8] to evaluate the reconstructed object point cloud.

**cm-degree metric.** For a predicted pose, the rotation error and translation error are computed separately. A predicted pose is considered correct if both its rotation error and translation error are less than a threshold.

**Proj2D metric.** The *Proj2D* metric computes the mean distance between the projection of 3D model points with given predicted and ground truth object poses. The estimated pose is considered correct if the mean projection distance is less than 5 pixels.

**ADD metric.** We first transform the 3D model points with the ground truth and the predicted poses. Then we compute the *ADD* metric using the mean distance between two transformed point sets. The pose is regarded as correct if the mean distance is less than 10% of the object diameter. Note that for symmetric objects, we use the *ADD(S)* [18] metric for evaluation.

**Point cloud accuracy.** We evaluate the point cloud accuracy in the ablation studies, following the metric in [8, 13]. The accuracy is defined as the percentage of reconstructed points which are within a distance threshold (e.g., *3mm*) with reference to the ground truth point cloud. We use vertices of the scanned object meshes as the ground truth point clouds.

## 4.3 Runtime Analyses of Keypoint-Free SfM

We evaluate the runtime of each part in the proposed keypoint-free SfM. The experiment is conducted on a server with two Intel<sup>(R)</sup> Xeon Gold 6146 CPU and an NVIDIA-V100-32GB GPU. We illustrate the runtime analyses with only one object instance below. The overall runtime varies according to several factors, such as image resolutions, the number of images used for reconstruction, and the number of successfully built coarse matches. For a video sequence with 193 images in  $512 \times 512$  resolution, it takes 135s to perform sequential coarse matching on 1436 image pairs and 40s to load all coarse matches and perform the triangulation [12]. Then, we perform fine matching between 1122 image pairs from the scene graph of coarse reconstructions to refine the feature tracks, which takes 171s. Finally, we optimize the object point cloud, which only consumes 1.03s.

Table 3: Compare with PixSfM.

	Point cloud accuracy			Pose error on OnePose dataset			Feature storage cost
	1mm	3mm	5mm	1cm-1deg	3cm-3deg	5cm-5deg	
LoFTR coarse + Our refinement	<b>30.9</b>	<b>75.8</b>	<b>87.7</b>	<b>51.1</b>	<b>80.8</b>	<b>87.7</b>	-
LoFTR coarse + PixSfM	29.5	73.0	85.9	48.2	79.6	86.7	7.35GB
LoFTR coarse + PixSfM (cost map)	27.9	69.7	82.9	47.9	79.2	86.7	0.17GB

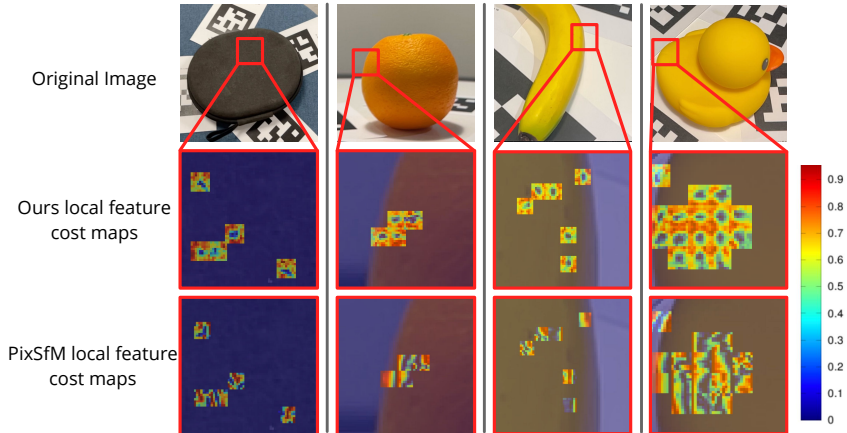


Figure 2: Visualization of local feature patch cost maps of our refinement and PixSfM [8]. Note that small block (■) means a cost map around a coarse match.

## 5 More Experiments

### 5.1 Using feature matching as the object detector

The need for an off-the-shelf 2D object detector can be eliminated by leveraging 2D-2D feature matching. This issue has been explored in OnePose [15]. Following OnePose, We first perform multiple 2D-2D feature matching between reference-query image pairs and then select the image pair with the most inliers to estimate 2D affine transformation. The region of interest (RoI) in the query image is then detected by transforming the corner of RoI in the reference image with the estimated transformation. To validate the effectiveness of this method, we present the evaluation results on the OnePose dataset with the feature-matching-based 2D detector below. For a fair comparison with OnePose, we use the same 2D bounding boxes, which are detected by SuperPoint [2] extractor and SuperGlue [11] matcher. The results in Tab. 1 demonstrate that the performance of the proposed method does not degrade significantly compared with using the ground truth bounding box as the 2D detector.

### 5.2 Compare with HLoc using different feature matching methods

In this part, we conduct experiments to compare our method with HLoc [10] combined with different feature matching and SfM methods, including LoFTR [14], DualRC-Net [5], Patch2Pix [19], on the OnePose dataset. We use the *cm-degree* pose accuracy metric with different thresholds for evaluation. Note that for the comparison with LoFTR, we use the strategy proposed by the original paper that round matches to grid level for SfM. As for Patch2Pix, we follow their points quantization strategy to yield repeatable matches for SfM. Please refer to their papers for more details. The results shown in Tab. 2 demonstrate the superiority of our method.

### 5.3 Compare with PixSfM

We compare the refinement module in our SfM framework with PixSfM [8] by the reconstruction and pose estimation accuracy. Results are shown in Tab. 3. Point cloud accuracy is evaluated on the objects with ground-truth meshes in the OnePose-LowTexture dataset, and pose estimation accuracy is evaluated on the OnePose dataset. The results demonstrate that the 3D models reconstructed by

Table 4: More ablation results.

	OnePose dataset			Time
	1cm-1deg	3cm-3deg	5cm-5deg	
<b>Full</b> ( $N_c = 3, N_f = 1$ , use all 3D points)	<b>51.1</b>	<b>80.8</b>	<b>87.7</b>	88.2ms
w/o Position Encoding	49.0	79.7	86.5	87.6ms
Large model ( $N_c = 6, N_f = 2$ )	48.1	78.5	85.5	133ms
Sample 7000 3D points	48.0	79.5	86.7	57.4ms
Sample 3000 3D points	47.1	78.6	86.1	42.7ms

our refinement achieve higher accuracy, and our refinement also brings improvement for the object pose estimation. As visualized in Fig. 2, we attribute the improvement to the more discriminative features (cost maps) obtained by our transformer-based fine-level matching module. Moreover, our refinement is more memory efficient since we do not need to keep the dense feature map or cost map in memory for the bundle adjustment.

#### 5.4 More Ablation Studies

We further conduct additional ablation studies on variants of the 2D-3D matching network architectures and different numbers of 3D points used for pose estimation.

**Ablation study on 2D-3D matching network architectures.** The results of a large model with more attention layers and a model without positional encoding are shown in Tab. 4. Increasing the number of attention layers by twice barely changes the results.

**Different numbers of 3D points.** We evaluate the effect of using different numbers of 3D points for object pose estimation. The results are illustrated in Tab 4. Our full model uses all reconstructed 3D object points for pose estimation, obtaining the highest accuracy. Decreasing the number of points with subsampling leads to minorly degraded pose estimation accuracy and faster inference speed.

## References

- [1] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. *ArXiv*, abs/2005.12872, 2020. 2
- [2] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. SuperPoint: Self-supervised interest point detection and description. *CVPRW*, 2017. 5
- [3] Mihai Dusmanu, Johannes L. Schönberger, and Marc Pollefeys. Multi-View Optimization of Local Feature Geometry. In *ECCV*, 2020. 1
- [4] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *ICML*, 2020. 2
- [5] Xinghui Li, Kai Han, Shuda Li, and Victor Prisacariu. Dual-resolution correspondence networks. In *NeurIPS*, 2020. 5
- [6] Zhigang Li, Gu Wang, and Xiangyang Ji. Cdpn: Coordinates-based disentangled pose network for real-time rgb-based 6-dof object pose estimation. *ICCV*, 2019. 3
- [7] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 2
- [8] Philipp Lindenberger, Paul-Edouard Sarlin, Viktor Larsson, and Marc Pollefeys. Pixel-perfect structure-from-motion with featuremetric refinement. *ICCV*, 2021. 1, 4, 5
- [9] Sida Peng, Xiaowei Zhou, Yuan Liu, Haotong Lin, Qixing Huang, and Hujun Bao. PVNet: pixel-wise voting network for 6dof object pose estimation. *T-PAMI*, 2020. 3
- [10] Paul-Edouard Sarlin, Cesar Cadena, Roland Siegwart, and Marcin Dymczyk. From coarse to fine: Robust hierarchical localization at large scale. In *CVPR*, 2019. 5
- [11] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. SuperGlue: Learning feature matching with graph neural networks. In *ICCV*, 2020. 2, 3, 5
- [12] Johannes L. Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 4
- [13] Thomas Schöps, Johannes L. Schönberger, S. Galliani, Torsten Sattler, Konrad Schindler, Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution images and multi-camera videos. *CVPR*, 2017. 4
- [14] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. LoFTR: Detector-free local feature matching with transformers. *CVPR*, 2021. 2, 3, 5
- [15] Jiaming Sun, Zihao Wang, Siyu Zhang, Xingyi He, Hongcheng Zhao, Guofeng Zhang, and Xiaowei Zhou. OnePose: One-shot object pose estimation without CAD models. *CVPR*, 2022. 3, 4, 5
- [16] Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017. 2
- [17] Qianqian Wang, Xiaowei Zhou, Bharath Hariharan, and Noah Snavely. Learning feature descriptors using camera pose supervision. In *ECCV*, 2020. 3
- [18] Yu Xiang, Tanner Schmidt, Venkatraman Narayanan, and Dieter Fox. PoseCNN: A convolutional neural network for 6d object pose estimation in cluttered scenes. *RSS*, 2018. 4
- [19] Qunjie Zhou, Torsten Sattler, and Laura Leal-Taixé. Patch2pix: Epipolar-guided pixel-level correspondences. *CVPR*, 2021. 5