

Geo6D: Geometric Constraints Learning for 6D Pose Estimation

Jianqiu Chen¹, Mingshan Sun², Ye Zheng³, Tianpeng Bao², Zhenyu He^{*1},
Donghai Li², Guoqiang Jin², Zhao Rui², Liwei Wu², Xiaoke Jiang⁴

¹Harbin Institute of Technology, Shenzhen

²SenseTime Research

³JD.com, Inc

⁴International Digital Economy Academy (IDEA)

Abstract

Numerous 6D pose estimation methods have been proposed that employ end-to-end regression to directly estimate the target pose parameters. Since the visible features of objects are implicitly influenced by their poses, the network allows inferring the pose by analyzing the differences in features in the visible region. However, due to the unpredictable and unrestricted range of pose variations, the implicitly learned visible feature-pose constraints are insufficiently covered by the training samples, making the network vulnerable to unseen object poses. To tackle these challenges, we proposed a novel geometric constraints learning approach called Geo6D for direct regression 6D pose estimation methods. It introduces a pose transformation formula expressed in relative offset representation, which is leveraged as geometric constraints to reconstruct the input and output targets of the network. These reconstructed data enable the network to estimate the pose based on explicit geometric constraints and relative offset representation mitigates the issue of the pose distribution gap. Extensive experimental results show that when equipped with Geo6D, the direct 6D methods achieve state-of-the-art performance on multiple datasets and demonstrate significant effectiveness, even with only 10% amount of data.

Introduction

6D pose estimation has drawn widespread attention as the essential prerequisite for emerging applications, such as robotic manipulation, autonomous driving, and augmented reality (Geiger, Lenz, and Urtasun 2012; Xu, Anguelov, and Jain 2018; Chen et al. 2017). In the computer vision community, several approaches have been proposed to estimate the transformation pose from the object frame to the camera frame. These existing methods can be categorized into two distinct groups: indirect and direct approaches. Indirect methods (Wang et al. 2019; He et al. 2020, 2021) usually first predict an intermediate feature and then use post-processing optimization algorithms, such as least-squares fitting and iterative Perspective-n-Point (PnP) algorithms (Su et al. 2022; Haugaard and Buch 2022; Li, Wang, and Ji 2019; Kiru, Patten, and Pix2Pose 2019; Rad and Lepetit 2017), to calculate the target pose based on the transformation or projection equation. The direct methods (Jiang

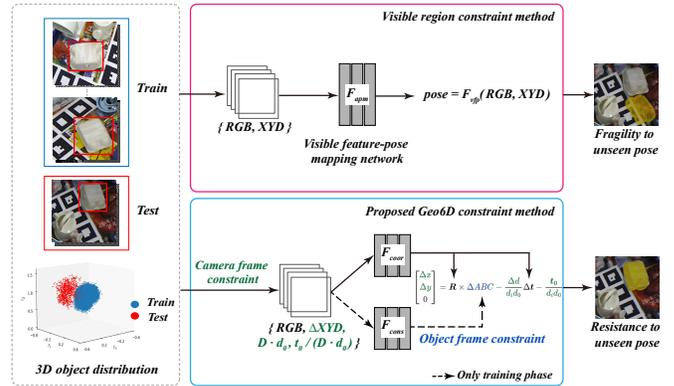


Figure 1: Existing direct pose estimation methods adopt implicit visible feature-pose constraints, that the visible feature of objects depends on its pose, having the fragility to unseen pose. The Geo6D approach introduces novel geometric constraints to rebuild the input and optimization target of the network. It enables the network to regress the pose from explicit geometric constraints and show the resistance to unseen poses.

et al. 2022; Li et al. 2018; Wang et al. 2020, 2021; Sun et al. 2022; Mo et al. 2022) directly predict the final 6D pose parameters (*e.g.*, rotation angles, and translation vectors) by the neural network in an end-to-end manner. However, both methods have their respective weaknesses. Although indirect approaches build the geometric constraints by intermediate geometric features, the detached two-stage pipeline makes the optimization target suboptimal, and the time-consuming iterative fitting in the pose estimation stage is an impediment in reality. On the contrary, the direct methods have the advantage of efficiency and end-to-end optimization. As shown in Fig 1, these leverage the implicit constraint to estimate the pose that the visible region of an object in the input image implicitly depends on the object CAD model, the camera intrinsic, and its pose. Given that the CAD model and camera intrinsic parameters are commonly known, the network is capable of estimating the pose based on the visible region of the object. However, the distribution of poses, specifically the translation component, is unpredictable and unrestricted. The method of mapping visible region appearance features to poses cannot cover ev-

*Corresponding author.

ery possibility. This creates a vulnerability to differences between the distribution of poses in training and testing data, as well as gaps in appearance domains. This limits the robustness and accuracy of the system.

To improve the visible feature-pose constraints, recent methods (Mo et al. 2022; Dong et al. 2019; Zeng et al. 2022) solve the ambiguity of multiple ground-truth poses relating to the same visible feature by modeling symmetric objects. After that, some methods (Sun et al. 2022; Mo et al. 2022) attempt to leverage an instance segmentation to mitigate the impact of visible features differences arising from camera intrinsic factors, the difference in object pose (camera extrinsic) still hampers the capacity of the model to accurately regress the pose from visible features. Besides, several direct methods (Wang et al. 2020, 2021) introduce an auxiliary loss to regress intermediate geometric features, such as 2D-3D correspondences, akin to indirect methods. However, these geometric features are not complete constraints to enable the network to regress the pose parameters based on it.

To solve these issues, we proposed the geometric constraints (Geo6D) learning approach that introduces a reformulated pose transformation to establish robust constraints on both camera and object frames by a relative offset representation. Specifically, the proposed Geo6D constraints are built upon the pose transformation formula. The rigid object points' 3D coordinates on the different frames can be transformed based on the pose. To address the distribution gap, we introduce a reference point and reformulate the pose transformation formula from the camera frame 3D coordinate representation (the offset for the visible point to the camera) to a relative offset for the visible point to the selected reference point. For making the formula learning-friendly and mathematically correct during network fitting, we separate the variables based on coordinate frames as explicit geometric constraints, demonstrated in Fig 1. For the camera frame variables, we supply and linearize all required variables in the camera frame as input and feed them to the network. For the object frame constraints, we introduce an additional regression network output head to predict the corresponding relative offset value in the object frame.

We encapsulate the Geo6D mechanism as a plugin, which rebuilds the input and output targets of the network and integrates it with two pose estimation networks. Extensive experiments demonstrate the effectiveness of our method, without sacrificing efficiency in both training and inference to enhance accuracy and stability and reduce the required amount of training data. It only requires 10% of training data to reach the comparable performance of full training data. Furthermore, we analyze the impact of the Geo6D mechanism from the perspective of the loss function.

To summarize, our main contributions are:

- Introducing a pose transformation formula in a relative offset presentation to establish explicit geometric constraints for direct methods.
- Proposing the Geo6D mechanism, a plugin module that processes input data and optimization targets to adhere to the geometric constraints, making the network learning-friendly and mathematically correct.

- Extensive experimental results demonstrate that the proposed Geo6D effectively improves the accuracy of existing direct pose estimation methods achieving state-of-the-art overall results and reducing the training data requirement, thus making it more practical for real-world applications.

Related work

Indirect 6D pose estimation

Indirect methods first predict intermediate geometric information and then exploit the projection constraints to estimate the 6D pose by optimization function. Recent methods (Peng et al. 2019; He et al. 2020, 2021) introduce the keypoints mechanism in 6D pose estimation and then estimate the 6D pose by a least-squares fitting algorithm, which takes advantage of the geometric constraints of rigid objects to train the keypoint prediction network. Different from the keypoints-based methods, 2D-3D correspondence-based methods (Su et al. 2022; Hodan, Barath, and Matas 2020; Haugaard and Buch 2022; Li, Wang, and Ji 2019; Kiru, Patten, and Pix2Pose 2019; Rad and Lepetit 2017) first establish the correspondences between 2D coordinates in the image plane and 3D coordinates in the object coordinate system by the neural network and then solve the 6D pose by a PnP or RANSAC algorithm. However, these indirect methods are only optimized in the first stage rather than the final pose regression, which is suboptimal compared with direct methods. Moreover, the optimization is time-consuming and computationally expensive in practical applications.

Direct 6D pose estimation

To estimate 6D pose efficiently, recent approaches (Mo et al. 2022; Jiang et al. 2022; Li et al. 2018; Wang et al. 2020, 2021) directly regress the final 6D pose parameters from the neural network instead of intermediate results. Densefusion (Wang et al. 2019) extracts the visible region features information from RGB-D images by two separate backbones to extract the features from 2D and 3D spaces and fuses them with a dense fusion network. Uni6D (Jiang et al. 2022) simplifies the architecture with a homogeneous single backbone to process RGB-D data, by introducing the extra UV data into input to preserve the projection constraints. Since the corresponding visible features of the object and pose are sensitive to the visual ambiguity of the symmetric object, there are multiple ground-truth poses related to the same visible features that confused the network fitting. ES6D individually models different types of symmetric objects to solve the issue of multiple pose mapping to the same visible features. Besides, the camera intrinsic is another factor for visible features of the object, Uni6Dv2 adopts an instance segmentation method to mitigate the impact of visible features difference from the camera intrinsic difference. However, since the pose parameters are unpredictable and unrestricted and the visible features to pose mapping can not be exhaustive and fragile for the unseen pose of the object in the test scene. To enhance network training, some methods (Wang et al. 2020, 2021) leverage the intermediate geometric features, i.e. 2D-3D correspondences, akin to indirect methods as an

auxiliary task to help network fit the pose transformation. However, the constraints from intermediate geometric features are insufficient to enable the network to regress the pose parameters based on it.

Geometric constraints in 6D pose estimation

The indirect methods usually utilize different geometric constraints to train a neural network to predict intermediate features such as predefined keypoints and 2D-3D correspondence. It has an explicit correlation between the input and output target features, which is easier for network optimization. Unlike the indirect method, the direct method needs to regress the pose parameters by the network, which restricts the network output target. To provide more geometric constraints information, some methods (Wang et al. 2020, 2021) apply an additional network output head and auxiliary loss to regress intermediate geometric features such as 2D-3D correspondences matrix or render alignment. However, the geometric correlations between the network output pose parameters and input data still are not explicitly established. Hence, we propose the Geo6D mechanism to establish an explicit geometric constraint by reformulating the input and output to hold the proposed relative offset representation pose transformation formula and make the pose estimation process easily trainable by the network.

Method

The Geo6D mechanism is proposed to enable the network to regress the pose in an end-to-end manner with explicit geometric constraints. In this section, we first introduce our geometric constraint-based 6D pose estimation method. Then we present how to adopt the proposed constraints in these direct methods, and finally, we analyze the impact of the Geo6D mechanism from the perspective of the loss function that balances of rotation and translation part training losses.

Geometric constraints for direct pose estimation

Current end-to-end direct regression methods (Sun et al. 2022; Mo et al. 2022) usually adopt a two-stage pipeline. In the first stage, an instance segmentation is utilized to crop and mask the candidate objects from the RGB-D image. In the second stage, the 3D coordinates XYD of each object, projected from the depth image and concatenated with the RGB channel are fed into the pose regression network to directly regress the pose parameters, i.e. the rotation matrix $\mathbf{R} \in SO(3)$ and translation vector $\mathbf{t} \in \mathbb{R}^3$.

According to the pose transformation formula

$$\begin{bmatrix} x \\ y \\ d \end{bmatrix} = \mathbf{R} \times \begin{bmatrix} a \\ b \\ c \end{bmatrix} + \mathbf{t}, \quad (1)$$

any visible point on the target object surface with the camera frame coordinates (x_i, y_i, d_i) has its corresponding object frame coordinates (a_i, b_i, c_i) . Similarly, the camera frame coordinate (x_0, y_0, d_0) and the corresponding object frame coordinate (a_0, b_0, c_0) of the centroid of all visible points should satisfy Eq (1). We can obtain the following naive geometric constraints by substituting coordinates of any visible point (x_i, y_i, d_i) and the reference point (x_0, y_0, d_0) into

Eq (1) successively and then subtracting both sides equally:

$$\begin{bmatrix} x_i - x_0 \\ y_i - y_0 \\ d_i - d_0 \end{bmatrix} = \mathbf{R} \times \begin{bmatrix} a_i - a_0 \\ b_i - b_0 \\ c_i - c_0 \end{bmatrix}. \quad (2)$$

However, these constraints only preserve the rotation part of the pose parameters and eliminate the translation part, which provides no benefit to regressing the \mathbf{t} during training. To preserve constraints of the rotation and translation at the same time, we scale all coordinates with respect to their corresponding depth value d_i or d_0 before subtracting and resulting in the following geometric constraints:

$$\begin{bmatrix} \Delta x \\ \Delta y \\ 0 \end{bmatrix} = \mathbf{R} \times \Delta ABC - \frac{\Delta d}{d_i d_0} \Delta \mathbf{t} - \frac{\mathbf{t}_0}{d_i d_0} \quad (3)$$

where

$$\Delta x = \frac{x_i}{d_i} - \frac{x_0}{d_0}, \quad \Delta y = \frac{y_i}{d_i} - \frac{y_0}{d_0}, \quad \Delta d = d_i - d_0,$$

$$\Delta ABC = \begin{bmatrix} \frac{a_i}{d_i} - \frac{a_0}{d_0} \\ \frac{b_i}{d_i} - \frac{b_0}{d_0} \\ \frac{c_i}{d_i} - \frac{c_0}{d_0} \end{bmatrix}, \quad \Delta \mathbf{t} = \mathbf{t} - \mathbf{t}_0, \quad \mathbf{t}_0 = \begin{bmatrix} x_0 \\ y_0 \\ d_0 \end{bmatrix}.$$

As shown in Fig 2, we divide the variables in Eq (3) into two groups according to whether they can be captured during the inference phase. Since Δx , Δy , Δd , $d_i d_0$ and $\frac{\mathbf{t}_0}{d_i d_0}$ can be calculated from the camera frame coordinates, these variables can be fed into the network directly as the **camera frame constraints**. On the contrary, the object frame coordinates a_i, b_i, c_i can not be captured during the inference phase and these variables are regressed by the auxiliary network head as the **object frame constraints**. As illustrated in Fig 2, the object frame constraints and the Geo Head are activated during the training phase but deactivated during inference. The camera frame constraints and object frame constraints together make up the proposed geometric constraints. Because the proposed Geo6D mechanism just modifies the input and output without any assumption of the network architecture, it can be plugged into 6D pose direct regression methods described in the following section.

Learning framework

In this section, we show how to integrate the Geo6D mechanism into current RGB-D direct regression methods. The overall framework is depicted in Fig 2, with our proposed Geo6D mechanism in yellow. Different from the current methods taking the RGB and XYD as the network input, the Geo6D mechanism is based on the geometric constraints to adopt the camera frame constraint variables $\Delta XYD, D \cdot d_0, \frac{\mathbf{t}_0}{D \cdot d_0}$ concatenated with RGB as input, where D and ΔXYD stands for the set of all points' depth value d_i and the offset value $\Delta x, \Delta y, \Delta d$. An additional output head (Geo head) introduces object frame constraints in order to regress the ΔABC . The Geo head is made up of convolutional blocks that are used to regress the visible points' relative offset in the object frame and optimize using an L2 loss. To establish explicit geometric restrictions, the object

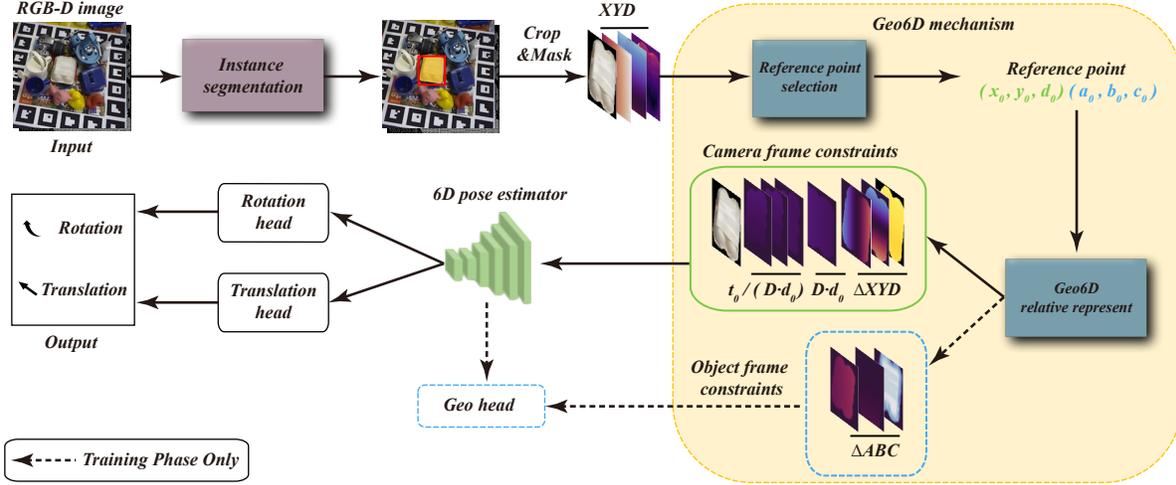


Figure 2: Overview of our proposed Geo6D method. The yellow area depicts the Geo6D mechanism, which reformulates the network input and output target based on Eq (3). The proposed Geo6D constraints are divided into two parts (the camera frame and object frame constraints) according to whether the introduced variables can be captured during the inference phase.

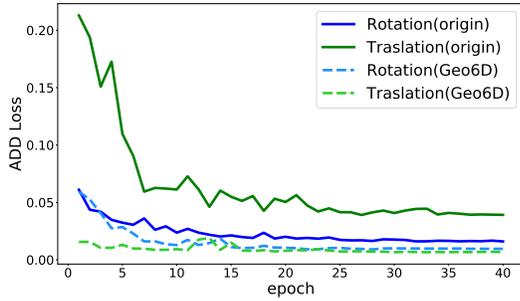


Figure 3: The Geo6D balances ADD loss of the rotation part and translation part on the Occlusion LineMOD dataset.

frame constraints collaborate with the camera frame constraints in the input data. It explicitly presents the necessary variables as input or optimization targets, allowing the network to regress the pose parameters \mathbf{R} and \mathbf{t} from them. Furthermore, the pose estimator output changes the translation vector \mathbf{t} with $\ominus \mathbf{t}$, which the final translation vector \mathbf{t} can be deduced by adding the reference point coordinates \mathbf{t}_0 in the camera frame. The pose estimation network can be summarised as the following function after utilizing the proposed Geo6D mechanism:

$$\mathbf{R}, \ominus \mathbf{t}, \Delta ABC = f(\text{RGB}, \Delta XYD, D \cdot d_0, \frac{\mathbf{t}_0}{D \cdot d_0}). \quad (4)$$

Balanced ADD loss

For most approaches, ADD loss is used to train the network to predict 6D pose. It transforms the points sampled on the object CAD model by the predicted pose and ground truth pose in a heterogeneous manner, and minimizes the distances between matching points in two separate transformation point sets as follows:

$$L_{ADD} = \frac{1}{m} \sum_j \left\| (\hat{\mathbf{R}}p_j + \hat{\mathbf{t}}) - (\bar{\mathbf{R}}p_j + \bar{\mathbf{t}}) \right\|_2^2 \quad (5)$$

where p_j is the j^{th} point from m randomly sampled CAD model's 3D points in the object frame, $(\hat{\mathbf{R}}, \hat{\mathbf{t}})$ is the predicted pose and $(\bar{\mathbf{R}}, \bar{\mathbf{t}})$ is the ground truth pose. Recombining Eq (5), it can be reformulated as:

$$\begin{aligned} L_{ADD} &= \frac{1}{m} \sum_j \left\| (\ominus \mathbf{R}p_j + \ominus \mathbf{t}) \right\|_2^2 \\ &= \frac{1}{m} \sum_j \left(\left\| \ominus \mathbf{R}p_j \right\|_2^2 + 2 \ominus \mathbf{R}p_j \ominus \mathbf{t} + \left\| \ominus \mathbf{t} \right\|_2^2 \right) \\ &= \frac{1}{m} \sum_j \left(\left\| \ominus \mathbf{R}p_j \right\|_2^2 \right) + \frac{2}{m} \ominus \mathbf{R} \sum_j p_j \ominus \mathbf{t} + \left\| \ominus \mathbf{t} \right\|_2^2 \end{aligned} \quad (6)$$

where $\ominus \mathbf{R} = \hat{\mathbf{R}} - \bar{\mathbf{R}}$, $\ominus \mathbf{t} = \hat{\mathbf{t}} - \bar{\mathbf{t}}$. Since p_j is sampled from the object frame, the centroid of sampled points is approx to the origin that $\sum_j p_j \approx 0$. Hence, the ADD loss is approximately equal to

$$L_{ADD} \approx \frac{1}{m} \sum_j \left(\left\| \ominus \mathbf{R}p_j \right\|_2^2 \right) + \left\| \ominus \mathbf{t} \right\|_2^2 \quad (7)$$

which consists of the rotation part $\frac{1}{m} \sum_j \left(\left\| \ominus \mathbf{R}p_j \right\|_2^2 \right)$ and translation part $\left\| \ominus \mathbf{t} \right\|_2^2$.

Because any point on the object is still inside the circumscribed sphere of the object after rotating, the rotation component is constrained by the object's diameter d . The translation portion, on the other hand, is unbounded since $\hat{\mathbf{t}}$ and $\bar{\mathbf{t}}$ are absolute 3D translate vectors, signifying the centroid of the object in the camera frame, which can be any place in 3D space. As a result, when $\left\| \ominus \mathbf{t} \right\|_2^2 \gg d$, the translation part may dominate the optimization while the rotation part is ignored. To accommodate the absolute magnitude of the translation vector, the loss will be downgraded to the l_2 loss. The

ADD loss has an apparent imbalance between the rotation and translation parts, as seen in Fig 3, with the translation part occupying the majority.

After adopting the Geo6D, the regression target translation vector is replaced by an equivalent offset value between the object centroid and the reference point, which range is transformed into the range $(-d/2, d/2)$ comparable to the rotation part. As a result, the rotation and translation parts are balanced in the same magnitude, allowing the rotation part to converge faster and the balanced loss to produce better optimization results. The balanced ADD loss optimizes the regression network to find correspondences between input data, the CAD model rotated, and the offset translate vector. Hence, the Geo6D mechanism corrects the ADD loss degradation and simplifies the pose estimation task.

Experiments

Benchmark datasets

We conduct our experiments on three benchmark datasets.

LineMOD (Hinterstoisser et al. 2011) contains 13 sequences of 13 low-textured objects. Since there is only about 1.2k real training data annotated with 6D pose, we follow previous work (Peng et al. 2019; Xiang et al. 2018; He et al. 2020; Jiang et al. 2022) to add synthetic images for training with 99.71% synthetic ratio.

Occlusion LineMOD (Brachmann et al. 2014) consists of 1214 testing images that are selected from the LineMOD dataset in the occlusion scene. Since there is no extra real training data, the training dataset is the same as LineMOD. The domain gap between training and testing datasets, as well as the heavily occluded objects, make it more challenging for 6D pose estimation.

YCB-Video (Calli et al. 2015) is a large and challenging dataset that contains 21 objects with 92 RGB-D sequences. It provides a large amount of real training data holding the same pose distribution between training and testing datasets.

Evaluation metrics

We adopt the commonly used average distance metrics ADD, ADD-S, and ADD(S) to evaluate different methods. ADD evaluates the mean of pairwise distance between two object point clouds which are transformed according to the ground truth pose $[\bar{\mathbf{R}}, \bar{\mathbf{t}}]$ and the predicted pose $[\hat{\mathbf{R}}, \hat{\mathbf{t}}]$ respectively:

$$\text{ADD} = \frac{1}{m} \sum_{p \in \mathcal{O}} \|(\hat{\mathbf{R}}p + \hat{\mathbf{t}}) - (\bar{\mathbf{R}}p + \bar{\mathbf{t}})\| \quad (8)$$

where \mathcal{O} denotes the 3D model of a object, p denotes any point on the model and m denotes total point number in the model. To alleviate the ambiguous matching of the symmetric objects, ADD-S is adopted to estimate the closest point distance between two point clouds:

$$\text{ADD-S} = \frac{1}{m} \sum_{p_1 \in \mathcal{O}} \min_{p_2 \in \mathcal{O}} \|(\hat{\mathbf{R}}p_1 + \hat{\mathbf{t}}) - (\bar{\mathbf{R}}p_2 + \bar{\mathbf{t}})\|. \quad (9)$$

For convenience, we introduce ADD(S) metric:

$$\text{ADD(S)} = \begin{cases} \text{ADD} & \mathcal{O} \text{ is asymmetric} \\ \text{ADD-S} & \mathcal{O} \text{ is symmetric} \end{cases} \quad (10)$$

	Geo6D	Occlusion LineMOD	LineMOD	YCB-Video
Uni6Dv2	w/o	40.6	97.2	91.5
Uni6Dv2	w	79.8(+39.2)	99.6(+2.4)	91.6(+0.1)

Table 1: Evaluation results of Uni6Dv2 with Geo6D on Occlusion LineMOD, LineMOD and YCB-Video datasets.

For LineMOD and Occlusion LineMOD datasets, we choose the threshold with 0.1d (10% of the diameter of the object) to calculate the accuracy of ADD(S), following (Peng et al. 2019; He et al. 2021). For the YCB-Video dataset, following (Xiang et al. 2018; Wang et al. 2019; He et al. 2020, 2021; Mo et al. 2022), we calculate the AUC (area under the accuracy-threshold curve) of ADD(S) with a maximum threshold of 0.1 meters.

Apply Geo6D to different methods

To verify the extensibility of the proposed Geo6D, we apply it to ES6D (Mo et al. 2022) and Uni6Dv2 (Sun et al. 2022). These two methods have different frameworks, output formats, and loss functions.

Apply Geo6D mechanism to Uni6Dv2. The input of Uni6Dv2 consists of an image patch (RGB, X, Y, D, NRM) which is cropped by the prediction of the segmentation network in the first stage. The sparse regression network based on the image patch then regresses the rotation and 3D location of each object in the camera frame and trains it using ADD Loss. To apply the Geo6D mechanism to Uni6Dv2, we replace the visible points’ absolute positional encoding part (X, Y, D) with relative values $(\Delta X, \Delta Y, \Delta D)$ and follow Eq (3) supplements the additional component concatenating with input data among channels. For the output part, we use an additional Geo head with Conv blocks to regress visible points’ ΔABC and adjust the regression target of the original translation head reformulated as the offset of the reference point to the centroid Δt . The comparison results are shown in Tab 1.

Apply Geo6D mechanism to ES6D. The original ES6D uses a dense regression network to estimate the offsets of visible points to an object’s centroid and a confidence score by feeding input cropped RGB images with vanilla normalized (x, y, d) . To reach the final pose, it substitutes a loss for the ADD loss and chooses the point with the highest confidence score. We alter the input to be Eq (3) and add an additional Geo head with the same structure as the translation head to regress ΔABC in order to apply the Geo6D method to ES6D. Because ES6D does not give the Occlusion LineMOD result, we implement it on the YCB-Video dataset, and all results are based on Ground Truth segmentation masks due to its sensitivity to mask results. We provide the results about the different amounts of training data in Tab 2.

Comparison with SOTA methods

We provide comprehensive and detailed comparison results on Occlusion LineMOD, LineMOD, and YCB-Video datasets. For the brevity of the paper, category-level exper-

Setting	Geo6D	10% R+10%S	10% R+100%S	100% R+100%S
ES6D	w/o	83.6	90.1	93.2
ES6D	w	84.7(+1.1)	92.3(+2.2)	93.6(+0.4)
Uni6Dv2	w/o	79.7	88.0	91.5
Uni6Dv2	w	86.5(+6.8)	89.3(+1.3)	91.6(+0.1)

Table 2: Evaluation results of ES6D and Uni6Dv2 with Geo6D on YCB-Video dataset with different amounts of training data. R means real data and S means synthetic data.

	PoseCNN	PVN3D	FFB6D	Uni6Dv2	ES6D	Uni6Dv2+Geo6D	ES6D+Geo6D
Avg	59.9	91.8	92.7	91.5	93.2	91.6	93.6

Table 3: Evaluation results on the YCB-Video dataset.

imental results on YCB-Video are provided in Supplementary Materials.

Evaluation on Occlusion LineMOD dataset. The quantitative results on the Occlusion LineMOD dataset are presented in Tab 4. Compared with the baseline Uni6Dv2, employing the Geo6D mechanism achieves a significant improvement (+39.2%) on ADD(S) metric and outperforms state-of-the-art method FFB6D (He et al. 2021)(indirect method) by 13.6%. Noteworthy, adopting the Geo6D mechanism makes the network more reliable when there are significant translation distribution gaps.

Evaluation on LineMOD dataset. The quantitative results on the LineMOD dataset are presented in Tab 5. Compared with our baseline Uni6Dv2, introducing the Geo6D mechanism brings a 2.5% performance gain on ADD(S), demonstrating its effectiveness. Compared to FFB6D (He et al. 2021), the performance gap is minimal (99.6% vs. 99.7%), while our Geo6D-based direct method has a more straightforward design and faster speed.

Evaluation on YCB-Video dataset. The quantitative results on the YCB-Video dataset are presented in Tab 3. We provide the implementations on the baseline Uni6Dv2 and ES6D, and introducing the Geo6D mechanism brings a 0.1% and 0.4% performance gain on the AUC of ADD(S) for full training data. Due to the large amount of data in the YCBV dataset, the network can be able to fit some variables, resulting in a relatively small improvement with the full dataset. As the result in Tab 2, the fewer training data, our method can achieve greater improvement.

In summary, our Geo6D mechanism makes direct 6D pose estimation methods outperform indirect methods on Occlusion LineMOD and YCB-Video datasets, and achieve comparable accuracy on the LineMOD dataset.

Ablation study

In this section, we will analyze the effectiveness of the Geo6D mechanism and compare different reference point generation strategies on Uni6Dv2. All experiments are conducted on the Occlusion LineMOD dataset with 10% training data.

Comparison to 3D coordinate normalization. We introduce the 3D normalization method only converting the coordinate into the offset value to reference points. As shown in Tab 6 of part (1), the normalized offset value input holds

class	PoseCNN	PVN3D	FFB6D	Uni6D	Uni6Dv2	Ours
ape	9.6	33.9	47.2	33.0	44.3	64.6
can	45.2	88.6	85.2	51.1	53.3	91.5
cat	0.9	39.1	45.7	4.6	16.7	63.2
driller	41.4	78.4	81.4	58.4	63.0	82.3
duck	19.6	41.9	53.9	34.8	51.6	63.9
eggbox	22.0	80.9	70.2	1.7	4.6	95.4
glue	38.5	68.1	60.1	30.2	40.3	95.0
holepuncher	22.1	74.7	85.9	32.1	50.9	82.6
Avg	24.9	63.2	66.2	30.7	40.6	79.8

Table 4: Evaluation results on the Occlusion LineMOD dataset. Symmetric objects are denoted in bold. "Ours" is Uni6Dv2+Geo6D.

	PoseCNN	PVN3D	FFB6D	Uni6D	Uni6Dv2	Ours
ape	77.0	97.3	98.4	93.7	95.7	98.3
benchvise	97.5	99.7	100.0	99.8	99.9	100.0
camera	93.5	99.6	99.9	96.0	95.8	99.6
can	96.5	99.5	99.8	99.0	96.0	99.9
cat	82.1	99.8	99.9	98.1	99.2	100.0
driller	95.0	99.8	100.0	99.1	99.2	99.8
duck	77.7	97.7	98.4	90.0	92.1	97.4
eggbox	97.1	99.8	100.0	100.0	100.0	100.0
glue	99.4	100.0	100.0	99.2	99.6	100.0
holepuncher	52.8	99.9	99.8	90.2	92.0	99.7
iron	98.3	99.7	99.9	99.5	98.0	100.0
lamp	97.5	99.8	99.9	99.4	98.5	99.9
phone	87.7	99.5	99.7	97.4	97.7	99.8
Avg	88.6	99.4	99.7	97.0	97.2	99.6

Table 5: Results on the LineMOD dataset. Symmetric objects are denoted in bold. "Ours" is Uni6Dv2+Geo6D.

a stable and compact data distribution for the input camera data and achieves 70.7%. However, it only mentions the camera frame representation breaking the transformation constraints. As shown in Eq (2), the translation part is eliminated which is out of optimization. Since the Geo6D endues the data representation with both camera and object constraints, the performance improves to 74.2%.

Effect of different components of Geo6D mechanism
To analyze the effect of each input and output components adjustment, an ablation study is conducted on different components in the proposed geometric constraints equation Eq (3). As reported in Tab 6, the ΔXYD significantly improves by solving the distribution gap issue. Based on the offset representation, introducing the object frame constraints by Geo head shows an advanced performance when choosing the optimization target following the proposed relative offset ΔABC as shown in Eq (3). Taking a channel as an example, the relative offset format $\frac{a_i}{d_i} - \frac{a_0}{d_0}$ follows the equation and achieves the better performance gain. Compared with only applying the Geo head, additional input of the components $D \cdot d_0$ and $\frac{t_0}{D \cdot d_0}$ as the camera frame constraints can indeed advance the performance. Furthermore, as the result in the (3) row, the proposed relative offset representation ($\frac{x_i}{d_i} - \frac{x_0}{d_0}$) outperforms the vanilla normalization ($x_i - x_0$), demonstrating the advantage of introducing the translate vector into optimization. The last row is performed without the RGB input to show that the proposed Geo6D mechanism can perform well without visible feature-pose constraints.

Comparison to speed and accuracy. When adopting the proposed Geo6D mechanism to Uni6Dv2, as shown in

Effect of different components in Geo6D mechanism		ADD(S)
components	setting	
(1): XYD value	absolute value	8.5
	offset value	70.7
(2): (1) + Geo head	absolute value: a_i	70.8
	offset value: $a_i - a_0$	71.1
	offset value: $\frac{a_i}{d_i} - \frac{a_0}{d_0}$	72.8
(3): (2) + $D \cdot d_0$ input	vanilla normalized XYD with $D \cdot d_0$	73.1
	proposed normalized XYD with $D \cdot d_0$	73.7
(4): (3) + $\frac{t_0}{D \cdot d_0}$ input	$\frac{t_0}{D \cdot d_0}$ input	74.2
(5): (4) - RGB input	only the camera frame constraints	73.7

Table 6: Effect of different components in Geo6D mechanism. The best setting in the previous row serves as the benchmark for the next row.

Reference point generation strategies		ADD(S)
x_0 and y_0	d_0	
Projection values	nearest depth point	66.7
Projection values	mean depth point	72.7
mean of visible points	mean of visible points	74.2

Table 7: Comparison of generation strategies for the reference point (x_0, y_0, d_0) . ‘‘Projection values’’ denotes calculating from the pin-hole projection equation.

Fig 4, it merely does not introduce extra time consumption and helps Uni6Dv2 outperform the current SOTA RGB-D method FFB6D (He et al. 2021) 13.6% while 5.6 times faster than it. The improvement demonstrates that the proposed geometric constraints can facilitate direct network higher performance without scarifying the efficiency.

Comparison to different reference point generation strategies. To build a stable reference point, we attempt three different reference point generation strategies, as shown in Tab 7. The first strategy selects the 2D coordinate UV value (u_0, v_0) of the ROI region center, the nearest point’s depth as reference point depth d_0 , and then calculates x_0 and y_0 based on pin-hole projection equation. The ROI center as the reference point in the 2D image plane maintains a stable UV encoding and attempts to use the center point’s depth to project the reference point in the 3D space. However, since the center point in ROI is likely in the outside or margin of the object due to the occlusion or the irregular, the depth value may be lost or a noise point from the sensor. To alleviate this problem, the second strategy adopts the mean depth value of the object region as the reference depth, which improves the stability of the Geo6D mechanism, achieving a 6% improvement from 66.7% to 72.7%. Besides the reference point selection based on the UV, the third strategy adopts the mean XYD of visible points as the reference point. Compared with the selection from UV, it is prone to consider the 3D shape and shows higher performance. Hence, whatever the reference point generation considers the 2D shape or 3D shape, the Geo6D mechanism can demonstrate its effectiveness.

Comparison to different segmentation results Most of these methods including our two baseline Uni6Dv2 and ES6D adopt the same segmentor MaskRCNN. All experiments hold the same upstream segmentation result with their

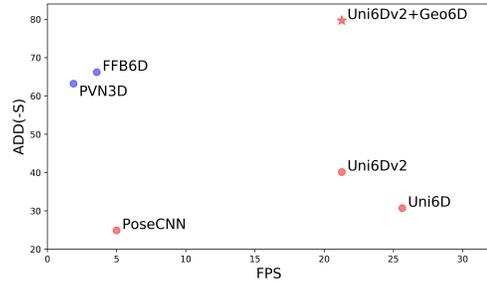


Figure 4: Comparison to speed and accuracy on the Occlusion LineMOD with full amount training data: Uni6Dv2+Geo6D achieves state-of-the-art accuracy in both direct regression methods (red dots) and indirect regression methods (blue dots) methods.

baselines. To evaluate the effect of Geo6D in different segmentation results, we conduct an experiment on the LM-O dataset to adopt the Ground Truth (GT) mask instead of the predicted mask, where the GT mask shows a slight 0.3% improvement from 79.8% to 80.1%. It shows the proposed Geo6D based on the predicted segmentation result reaching nearly the upper bound performance and demonstrates the robustness of the segmentation results.

Comparison to scaling different parts of ADD Loss.

The imbalance of the translation and rotation part in the ADD loss is due to the issue of unlimited output space of the translation part. Scaling the two parts into a comparable level leads to the challenging translation part without enough optimization. We conduct an experiment based on the analysis in Fig 3, increasing rotation part weight 4 times, where the performance decreased from 40.6% to 15.4%. It supports that only scaling the two-part loss can not achieve the same effectiveness as our balanced ADD loss.

Conclusion

In this paper, we propose a novel geometric constraints learning approach for 6D pose estimation that achieves state-of-the-art performance on multiple benchmarks. First, a non-ill-conditioned 6D pose transformation is derived. Geo6D rebuilds the input data based on the derived transformation and employs a Geo head to strengthen the point-to-point relationship as constraints between the camera frame and the object frame in a leaning-friendly manner. Extensive experiments show that Geo6D greatly simplifies the task of 6D pose estimation and can be plugged into various direct methods such as Uni6Dv2 and ES6D. Experiments also indicate that Geo6D achieves higher gains with less training data. We believe the Geo6D mechanism has the potential to inspire other 3D tasks, such as category-level pose estimation and 3D object detection using RGB-D or Lidar data.

In terms of limitations, Geo6D’s current reference point generation strategy requires that the mean point be close to the centroid of the majority of objects. Geo6D struggles to generate a stable reference point on the CAD model when the mean point of irregular objects is far from the centroid. When dealing with irregular objects, more effective generation strategies must be investigated.

References

- Brachmann, E.; Krull, A.; Michel, F.; Gumhold, S.; Shotton, J.; and Rother, C. 2014. Learning 6d object pose estimation using 3d object coordinates. In *European conference on computer vision*, 536–551.
- Calli, B.; Singh, A.; Walsman, A.; Srinivasa, S.; Abbeel, P.; and Dollar, A. M. 2015. The ycb object and model set: Towards common benchmarks for manipulation research. In *2015 international conference on advanced robotics (ICAR)*, 510–517.
- Chen, X.; Ma, H.; Wan, J.; Li, B.; and Xia, T. 2017. Multi-view 3d object detection network for autonomous driving. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 1907–1915.
- Dong, Z.; Liu, S.; Zhou, T.; Cheng, H.; Zeng, L.; Yu, X.; and Liu, H. 2019. PPR-Net: point-wise pose regression network for instance segmentation and 6d pose estimation in bin-picking scenarios. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 1773–1780. IEEE.
- Geiger, A.; Lenz, P.; and Urtasun, R. 2012. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*, 3354–3361.
- Haugaard, R. L.; and Buch, A. G. 2022. SurfEmb: Dense and Continuous Correspondence Distributions for Object Pose Estimation with Learnt Surface Embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 6749–6758.
- He, Y.; Huang, H.; Fan, H.; Chen, Q.; and Sun, J. 2021. Ffb6d: A full flow bidirectional fusion network for 6d pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3003–3013.
- He, Y.; Sun, W.; Huang, H.; Liu, J.; Fan, H.; and Sun, J. 2020. Pvn3d: A deep point-wise 3d keypoints voting network for 6dof pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 11632–11641.
- Hinterstoisser, S.; Holzer, S.; Cagniart, C.; Ilic, S.; Konolige, K.; Navab, N.; and Lepetit, V. 2011. Multimodal templates for real-time detection of texture-less objects in heavily cluttered scenes. In *2011 international conference on computer vision*, 858–865.
- Hodan, T.; Barath, D.; and Matas, J. 2020. Epos: Estimating 6d pose of objects with symmetries. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 11703–11712.
- Jiang, X.; Li, D.; Chen, H.; Zheng, Y.; Zhao, R.; and Wu, L. 2022. Uni6D: A Unified CNN Framework without Projection Breakdown for 6D Pose Estimation. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Kiru, P.; Patten, T.; and Pix2Pose, M. 2019. Pixel-Wise Coordinate Regression of Objects for 6D Pose Estimation. *Proceedings of the ICCV, Seoul, Korea*, 27–28.
- Li, Y.; Wang, G.; Ji, X.; Xiang, Y.; and Fox, D. 2018. Deepim: Deep iterative matching for 6d pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 683–698.
- Li, Z.; Wang, G.; and Ji, X. 2019. Cdpn: Coordinates-based disentangled pose network for real-time rgb-based 6-dof object pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 7678–7687.
- Mo, N.; Gan, W.; Yokoya, N.; and Chen, S. 2022. ES6D: A Computation Efficient and Symmetry-Aware 6D Pose Regression Framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 6718–6727.
- Peng, S.; Liu, Y.; Huang, Q.; Zhou, X.; and Bao, H. 2019. Pvnnet: Pixel-wise voting network for 6dof pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4561–4570.
- Rad, M.; and Lepetit, V. 2017. Bb8: A scalable, accurate, robust to partial occlusion method for predicting the 3d poses of challenging objects without using depth. In *Proceedings of the IEEE international conference on computer vision*, 3828–3836.
- Su, Y.; Saleh, M.; Fetzer, T.; Rambach, J.; Navab, N.; Busam, B.; Stricker, D.; and Tombari, F. 2022. ZebraPose: Coarse to Fine Surface Encoding for 6DoF Object Pose Estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 6738–6748.
- Sun, M.; Zheng, Y.; Bao, T.; Chen, J.; Jin, G.; Zhao, R.; Wu, L.; and Jiang, X. 2022. Uni6Dv2: Noise Elimination for 6D Pose Estimation. *arXiv preprint arXiv:2208.06416*.
- Wang, C.; Xu, D.; Zhu, Y.; Martín-Martín, R.; Lu, C.; Fei-Fei, L.; and Savarese, S. 2019. Densefusion: 6d object pose estimation by iterative dense fusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 3343–3352.
- Wang, G.; Manhardt, F.; Shao, J.; Ji, X.; Navab, N.; and Tombari, F. 2020. Self6d: Self-supervised monocular 6d object pose estimation. In *European Conference on Computer Vision*, 108–125.
- Wang, G.; Manhardt, F.; Tombari, F.; and Ji, X. 2021. Gdrnet: Geometry-guided direct regression network for monocular 6d object pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16611–16621.
- Xiang, Y.; Schmidt, T.; Narayanan, V.; and Fox, D. 2018. PoseCNN: A Convolutional Neural Network for 6D Object Pose Estimation in Cluttered Scenes. In *Robotics: Science and Systems (RSS)*.
- Xu, D.; Anguelov, D.; and Jain, A. 2018. Pointfusion: Deep sensor fusion for 3d bounding box estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 244–253.
- Zeng, L.; Lv, W. J.; Dong, Z. K.; and Liu, Y. J. 2022. PPR-Net++: Accurate 6-D Pose Estimation in Stacked Scenarios. *IEEE Transactions on Automation Science and Engineering*, 19(4): 3139–3151.