# GS2Pose: Two-stage 6D Object Pose Estimation Guided by Gaussian Splatting

Jilan Mei, Junbo Li, Cai Meng⋆

Beihang University

*Abstract*— **This paper proposes a new method for accurate and robust 6D pose estimation of novel objects, named GS2Pose. By introducing 3D Gaussian splatting, GS2Pose can utilize the reconstruction results without requiring a high-quality CAD model, which means it only requires segmented RGBD images as input. Specifically, GS2Pose employs a two-stage structure consisting of coarse estimation followed by refined estimation. In the coarse stage, a lightweight U-Net network with a polarization attention mechanism, called Pose-Net, is designed. By using the 3DGS model for supervised training, Pose-Net can generate NOCS images to compute a coarse pose. In the refinement stage, GS2Pose formulates a pose regression algorithm following the idea of reprojection or Bundle Adjustment (BA), referred to as GS-Refiner. By leveraging Lie algebra to extend 3DGS, GS-Refiner obtains a pose-differentiable rendering pipeline that refines the coarse pose by comparing the input images with the rendered images. GS-Refiner also selectively updates parameters in the 3DGS model to achieve environmental adaptation, thereby enhancing the algorithm's robustness and flexibility to illuminative variation, occlusion, and other challenging disruptive factors. GS2Pose was evaluated through experiments conducted on the LineMod dataset, where it was compared with similar algorithms, yielding highly competitive results. The code for GS2Pose will soon be released on GitHub.**

*Index Terms*—**6D pose estimation, 3DGS, light adaptability, novel objects.**

## 1. Introduction

**A**Ccurate 6D object pose estimation is a fundamental problem in the field of computer vision, with broad application prospects in technologies such as robot navigation[1, 2] and virtual reality[3, 4]. However, classical pose estimation algorithms[5–7] lack robustness against environmental interference, such as non-uniform lighting, varying degrees of occlusion, and dynamic blur. Moreover, the lightweight nature of the algorithm is also demanding in the field of embodied intelligence[8–10].

With the widespread application of deep learning methods, the robustness of related algorithms[11–15] against interference has continually improved. Early works[16–19] have achieved high-precision instance-level pose estimation. However, these models can only handle a specific object after the training session and cannot generalize to others. Additionally, they require datasets with precise ground truth poses, which are difficult to obtain in practical applications.

The emergence of novel pose representation methods [20], such as NOCS, has led to breakthroughs in category-level pose estimation methods[21–26], achieving notable intra-class generalization. Trained models can perform high-precision pose estimation on objects with similar geometric and color

features. However, these methods typically require a substantial number of CAD models of the same category during the training phase, results in huge time expenditure. Additionally, since the 6D pose of the target object is bound to the objects coordinate system under the CAD model, which can lead to issues, such as parameter ambiguity in the estimation results during the inference phase.

In recent years, with the development of large models[27–29], some research[30–32] have introduced the concept of pre-training on large datasets into the field of 6D pose estimation. These methods construct large datasets by collecting numerous CAD models of common objects from different categories, enabling effective generalization to unseen objects. They require only the CAD model of the target object during inference, allowing for the artificial setting of strict coordinate relationships without the need for additional training on the target object. However, these models also have drawbacks, such as the inability to generalize to uncommon objects, high consumption of computational resources, and their accuracy being heavily dependent on the quality of CAD modeling.

To address the aforementioned shortcomings of existing algorithms, a novel pose estimation method is proposed that eliminates the need for artificially designed CAD models. This method is designed for application scenarios where high-quality CAD models of the target object are unavailable, and only untextured scanned models or structure-from-motion (SFM) point cloud models can be obtained. To achieve lightweight training, accurate reference relationships, and robustness to interference, GS2Pose consists of a two-stage pose estimation approach comprising coarse estimation followed by pose refinement.

The detailed process of GS2Pose is illustrated in Fig. 1. The 3DGS point cloud model of the object (hereafter referred to as the 3DGS model) is obtained using existing 3DGS reconstruction techniques, with the object coordinate system manually specified. Utilizing insights from the GS-SLAM model[33], the commonly used reprojection-based pose optimization iterative approach from the SLAM domain is introduced [34–36], also known as Bundle Adjustment (BA). By representing object poses using Lie algebra and integrating this representation with the differentiable 3DGS rendering pipeline, an approach is implemented that utilizes reprojection and backpropagation. This enables an iterative optimization algorithm that can regress both the object pose and the camera pose, referred to as GS-Refiner.

Since the iterative optimization algorithm requires a reasonable initial pose as a starting point, it is necessary to design an

algorithm that can provide a rough pose estimate based solely on the segmented object image. Inspired by the NeRF-Pose model[37], a rough pose estimation network named Pose-Unet was developed. RGB images and their corresponding NOCS images are obtained from the camera perspective using 3DGS. These images are subsequently input into a pre-trained coarse pose estimation network (Pose-Unet) for fine-tuning, resulting in a coarse pose estimation for any novel rendering view of the object.

On the other hand, GS-Refiner leverages the parameter interpretability of the 3DGS model to selectively optimize and refine parameters, such as higher-order spherical harmonic color parameters, transparency, and ellipsoid orientation through backpropagation. This allows the surface colors to adaptively adjust to environmental factors encountered during actual capture, such as lighting, occlusion, and motion blur.

The primary contributions of the paper can be summarized as follows:

i) By incorporating 3DGS reconstruction technology, lightweight 6D pose estimation of previously unseen objects is achieved in the absence of CAD models.

ii) By employing Lie algebra to modify the differentiable rendering pipeline of 3DGS, a reprojection iterative algorithm called GS-Refiner has been developed developed, enabling the correction of both object poses and camera poses.

iii) By selectively regressing the parameters of 3DGS, a 6D pose estimation algorithm was developed with robust resistance to complex lighting, motion blur, and occlusions.

iv) Through experiments on datasets such as LineMod, the GS2Pose model demonstrated substantial advantages over comparable algorithms, particularly in terms of accuracy, inference speed, and computational resource efficiency.

## 2. RELATED WORKS

This section provides a brief summary of the development on the 6D pose estimation. We first review the 6D pose prediction about known rigid objects. Then we focus on the progress of 6D pose estimation about novel objects in recent years. We summarize the recent development of Gaussian models finally.

### A. 6D pose estimation of seen objects

Traditional 6D pose estimation methods[5, 38–40] rely on extracting local invariant features and establishing correspondences by template matching. Researchers have made innovative explorations in the features robustness and the template matching performance in complex occlusion scenarios. However, these traditional methods still struggle to solve challenges related to large variations in lighting and the accurate pose estimation of symmetric objects. As a result, the 6D pose estimation becomes inefficient and unsuitable for widespread development and practical applications.

Conversely, deep learning methods have gained attention in 6D pose estimation due to their powerful ability to automatically learn features from datasets. The PoseCNN model[19] introduced a novel loss function, enabling the network to

better handle symmetric objects, thereby enhancing the robot's ability to interact with the real world. As for the applications without depth information, BB8[41] model proposed a classifier to restrict the range of poses, which can compensates the lack of depth information. Moreover, RADet[42] proposed a rigidity-aware detection method to better address occlusion issues, which created a visibility map using the minimum barrier between each pixel in the detection bounding box and the box boundary.

Recently, Generative Adversarial Networks (GAN) have demonstrated exceptional capabilities in denoising and recovering missing parts of images. UnrealDA[43] proposed a GAN-based network, which transformed real depth maps with background occlusion into synthetic depth maps to improve pose estimation performance. Apart from that, the Pix2Pose[44] model based on GAN network, introduced a transformer loss to guide predictions toward the closest pose, addressing pose estimation for symmetric objects.

### B. 6D pose estimation of unseen objects

To improve the generalization ability and robustness of 6D pose estimation with CAD models, some researchers aim to address pose estimation for novel objects. MegaPose[31] network proposed a 6D pose estimator based on a rendering and comparison strategy, which trains the network on a large synthetic dataset. Moreover, GigaPose[30] network proposed a novel solution by leveraging templates to recover out-of-plane rotations, then utilizing patches correspondences to estimate the four remaining pose parameters. Although above foundation methods have strong generalization capabilities, their robustness remains insufficient for specialized devices in industries and medical, such as surgical instruments and precision constructions. We proposed a course-refine 6D pose estimation network. For each new object, coarse estimation network provides an approximate pose by rapid training, followed by precise correction in the refine estimation network.

### C. 6D pose estimation with 3D reconstruction model

3DGS[45] demonstrates significant advantages in high-quality and real-time rendering. This work represents scene with 3D Gaussian ellipsoids and efficiently renders by rasterizing the Gaussian ellipsoids into images, achieving state-of-the-art (SOTA) level visual quality. At the same time, 3DGS employs an explicit construction method, possessing clear geometric structure and appearance. This technology has already been applied in multiple fields, including autonomous navigation[33, 36, 46, 47], virtual human body reconstruction[48, 49] and 3D generation[50–52].

However, there are few works that apply 3D Gaussian Splatting (3DGS) to the field pose estimation currently. Although GSPose network attempts to apply 3DGS model, it still requires training a DINO network to create a database, which rely on dataset training. So this research is difficult to fine-tune for new objects. Moreover, it does not fully utilize the differentiable advantages of 3D Gaussian.
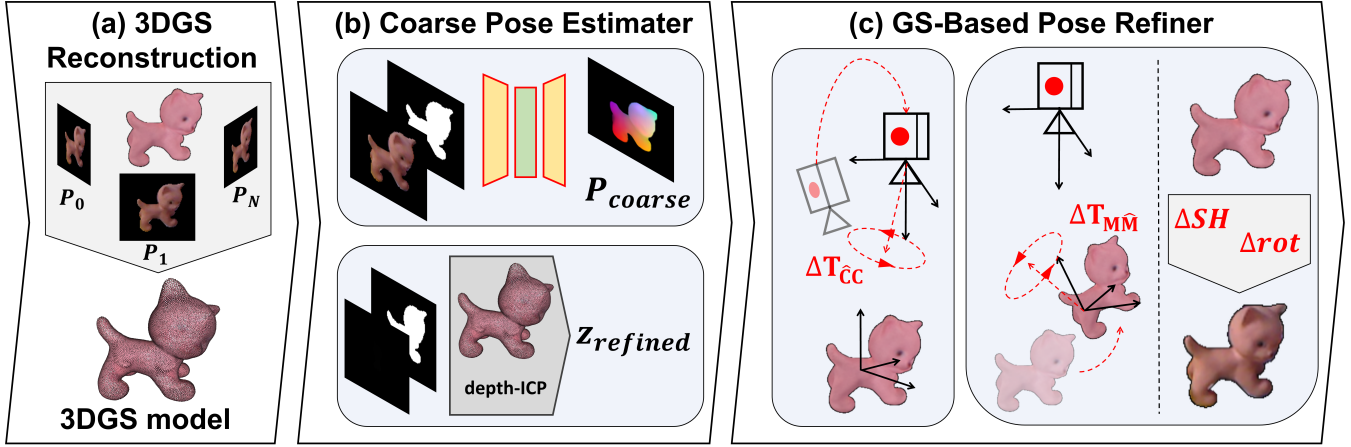
Fig. 1. The structure of the GS2POSE

## 3. METHODOLOGY

### A. Overview

In this chapter, we provide a detailed overview of the framework and principles of pose estimation methods. Our objective is to determine the relative pose of an object with respect to the camera, based on the input RGB-D image $I_{\text{in}}$ and the 3D geometric reconstruction model $G_m$ of the object. This involves computing the transformation matrix $T_{cm}$ from the coordinate system of the reconstructed model $m$ to the camera coordinate system $c$, which is composed of a translation vector $t_{cm}$ and a rotation matrix $R_{cm}$.

$$T_{cm} = \begin{bmatrix} R_{cm} & t_{cm} \\ 0^T & 1 \end{bmatrix}, \quad t_{cm} = \begin{pmatrix} x_{cm} & y_{cm} & z_{cm} \end{pmatrix}^T \quad (1)$$

To achieve the aforementioned objective, we first reconstruct the 3D Gaussian Splatting (3DGS) model of the target object. Subsequently, under the supervision of this 3DGS model, we train a coarse estimation network, referred to as Pose-net, which is capable of generating NOCS images from novel viewpoints and predicting the coarse pose $T_{cm}^{\text{coarse}}$ of the object in RGB images captured from arbitrary angles. Finally, we propose a novel refinement algorithm that utilizes the coarse predicted pose as an initial estimate, following an iterative optimization approach based on 3DGS reprojection. By continuously minimizing the differences between the rendered images and the input images, we refine and optimize the pose to obtain an accurate final output $T_{cm}^{\text{refined}}$.

### B. 3D Gaussian Splatting

3D Gaussian Spheres (3DGS) is a scene representation method that describes objects in the world coordinate system using Gaussian spheres. All attributes of the 3D Gaussian Spheres are learnable, including the position parameters $\mu$, opacity $a$, the 3D covariance matrix $r$, and the spherical harmonics $sh$. Given any point $\mathbf{x}$ in the world coordinate system, the 3D Gaussian sphere defined at point $\mathbf{x}$ according to the Gaussian distribution is as follows:

$$f(\mathbf{x}; \mu, \Sigma) = \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^{\text{T}}\Sigma^{-1}(\mathbf{x} - \mu)\right) \quad (2)$$

$$\Sigma = RSS^{\text{T}}R^{\text{T}} \quad (3)$$

where $R$ denotes the rotation matrix computed from $r$, and $S$ represents the diagonal matrix derived from $s$. Subsequently, a fast rasterization approach is employed to project the 3D Gaussian points onto a 2D plane for rendering.

### C. Coarse Pose Estimation Network

Inspired by the NeRF-Pose model[37], which is currently the state-of-the-art approach in 6D pose estimation, we have designed a lightweight NOCS image generation network ( Pose-Unet ) to predict the coarse pose of objects. The 3DGS method generates RGB images from the camera viewpoint along with the corresponding NOCS images, which are used as training inputs for Pose-Unet. Through fine-tuning, the model can rapidly generalize to new objects. Subsequently, the test RGB images (segmented using the CNOS model) are input to obtain the corresponding NOCS images, from which a coarse pose is estimated. Since the NOCS image predictions exhibit significant deviations along the z-axis, the improved ICP algorithm is utilized to align the point cloud model in the observed viewpoint (acquired from RGB-D images) with the Gaussian model, correcting the z-axis in the coarse pose.

Pose-Unet utilizes ResNet50 as the encoder. While in the decoder stage, three transposed convolution layers are employed for up sampling. As most encoder-decoder based network models, Pose-Unet incorporates skip connections ( Mobile-ASPP) during the down sampling to minimize information loss. Mobile-ASPP optimizes the ASPP structure, which consists of three parallel atrous convolutions. Specifically, the dilated convolution layers have kernel sizes of $1 \times 1$, $3 \times 3$, and $3 \times 3$, with corresponding dilation rates of 1, 1, and 2, respectively. This module enables the network to fully capture shallow information while reducing computational resource consumption. In the deep feature extraction module, based on the PPM structure, four parallel pooling layers with different kernel sizes are constructed to effectively capture dependencies between pixels.
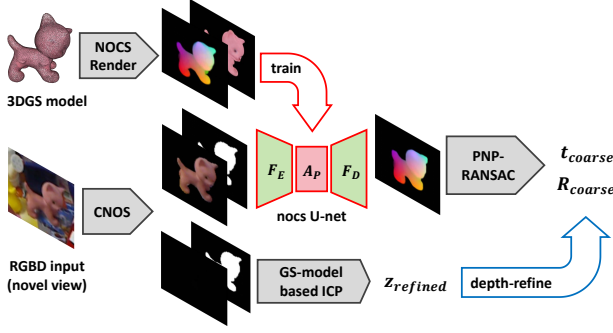
Fig. 2. The structure of the Pose-Unet

A deep estimation method for point cloud registration on object surfaces is proposed. The point cloud $P_1$ is generated by combining the RGB image and depth image from the target viewpoint. By combining the CAD model with camera pose estimation, the point cloud $P_2$ of the object's surface facing the camera is generated. By calculating the average z-values, $Z_1$ and $Z_2$, of the two point clouds $P_1$ and $P_2$ in the camera coordinate system along the z-axis, the estimated $Z_2$ is corrected to $Z_1$. This completes a pose correction along the depth direction.

### D. Refine Pose Estimation Network

*1) Overview:* After obtaining a coarse estimation $T_{cm}^{\text{coarse}}$ with limited accuracy, we designed a multi-stage refinement algorithm, termed GS-refiner, which leverages the 3DGS representation model of the object. This algorithm employs an iterative reprojection method to provide a precise pose estimation of the object.

Inspired by 3D Gaussian Splatting SLAM[33], we represent the pose changes between coordinate systems using Lie algebra. We compute the error through reprojection for backpropagation, aiming to regress the precise pose of the object.

Thanks to the differentiable rendering pipeline of 3DGS, we can differentiate most parameters of the 3DGS, including the rendering pose, by calculating the differences between the reprojection images $I_{\text{pred}}(T_{cm}^{\text{iter}})$ under the coarse estimated pose and the ground truth images $I_{\text{in}}$. Following the approach of 3D Gaussian Splatting, we design the loss function as follows:

$$\text{loss}(I_{\text{in}}, I_{\text{pred}}) = \lambda L_1 + (1 - \lambda) L_{\text{DSSIM}} \tag{4}$$

where $\lambda$ is a hyperparameter, $\mathcal{L}_1$ represents the L1 loss between two images, and $\mathcal{L}_{D-\text{SSIM}}$ represents the D-SSIM loss. Specifically, let the loss at a certain pixel $p(u, v)$ be determined by the value of that pixel:

$$L(u, v) = \text{loss}(p_{\text{pred}}, p_{\text{gt}}) \tag{5}$$

which is influenced by the 2D elliptical projections of multiple 3D Gaussian Splatting ellipsoids projected onto it. This can be expressed using the ray casting formula:

$$p_{\text{pred}} = \sum_{i=1}^{N} c_i \alpha_i \prod_{j=i-1}^{1} (1 - \alpha_j) \tag{6}$$

The RGB color vector $c_i$ can be obtained from the spherical harmonic parameters of the 3DGS ellipsoid and the relative pose $T_{cm}$ with respect to the camera. The transparency $\alpha_i$ is determined by the distance between the current pixel and the center point $p_i$ of the 2D ellipse projection, as well as the Gaussian covariance parameters $\Sigma_i$ of the projected ellipse. Furthermore, the center point $p_i$ is determined by the camera's relative pose $T_{cm}$, namely:

$$p_i = \pi(T_{cm}, p_m) \tag{7}$$

Based on the $T_{cm}$ and the camera intrinsic parameter matrix $K$, a Jacobian matrix $J$ can be generated for the purpose of flattening the ellipsoid parameters into the plane. $\Sigma_i$ can be determined using the Jacobian matrix $J$ and the rotational part of the relative pose $R_{cm}$:

$$\Sigma_i = J\Sigma_c J^{\text{T}} \tag{8}$$

$$\Sigma_c = R_{cm}\Sigma_m R_{cm}^{\text{T}} \tag{9}$$

according to the chain rule of differentiation:

$$\frac{\partial p_i}{\partial T_{cm}} = \frac{\partial p_i}{\partial p_c} \frac{\partial p_c}{\partial T_{cm}} \tag{10}$$

$$\frac{\partial \Sigma_i}{\partial T_{cm}} = \frac{\partial \Sigma_i}{\partial J} \frac{\partial J}{\partial p_c} \frac{\partial p_c}{\partial T_{cm}} + \frac{\partial \Sigma_i}{\partial R_{cm}} \frac{\partial R_{cm}}{\partial T_{cm}} \tag{11}$$

$$\frac{\partial c_i}{\partial T_{cm}} = \frac{\partial c_i}{\partial t_{cm}} \frac{\partial t_{cm}}{\partial T_{cm}} \tag{12}$$

Due to the discontinuity of the matrix form of $T_{cm}$ in $\mathbb{R}^{4\times 4}$, $\frac{\partial p_c}{\partial T_{cm}}$ and $\frac{\partial R_{cm}}{\partial T_{cm}}$ cannot be directly differentiated. Therefore, we need to convert $T_{cm}$ into the Lie algebra form before performing the differentiation.

Let the homogeneous coordinates of any point in the point cloud model in the object coordinate system be denoted as $p_m = [x_m, y_m, z_m, 1]^{\text{T}}$ and the homogeneous coordinates in the camera coordinate system be denoted as $p_c = [x_c, y_c, z_c, 1]^{\text{T}}$. When the non-homogeneous form of the coordinates is used, it will be indicated by a subscript, such as $p_c^{:3}$. According to the definition of the transformation matrix, we have:

$$p_c = T_{cm}p_m \tag{13}$$

It is important to note that, based on the knowledge of Lie algebra, altering the camera pose in the object coordinate system (perturbing the shooting perspective) yields fundamentally different effects compared to changing the object pose in the camera coordinate system when refining relative poses $T_{cm}$. Subsequent formula derivations and experiments demonstrate that these two approaches efficiently correct the translational and angular relationships of the object relative to the camera.

Consequently, we have separated the Refiner into two components: perspective pose correction (Camera refiner) and object pose correction (Object refiner). Below, we will introduce the principles of these two components, provide a brief derivation of the relevant formulas, and finally explain how we integrate these two components to form the GS Refiner.

*2) Camera Refiner:* In the Camera Refiner, the object being updated is the camera coordinate system $c$. During each iteration, a new camera coordinate system $c'$ is obtained, and the coordinates of any object point $p_m$ in the new camera coordinate system are given by:

$$p_{c'} = T_{c'c}T_{cm}p_m = T_{c'c}p_c \tag{14}$$

Here, $T_{c'c} \in SE(3)$ can be viewed as a left perturbation applied to $T_{cm}$. Let the Lie algebra corresponding to $T_{c'c}$ be denoted as:

$$\tau_c = \begin{bmatrix} \rho_c & \varphi_c \end{bmatrix}^T \in \mathfrak{se}(3) \quad p_{c'} = \exp(\tau_c)p_c \tag{15}$$

Since the rendered coordinate system at this point is the transformed camera coordinate system, that is:

$$\frac{\partial p_i}{\partial T_{cm}} = \frac{\partial p_i}{\partial p_c} \cdot \frac{\partial p_c}{\partial T_{cm}} \tag{16}$$

Let $p_{c'}$ take the derivative of $\tau_c$, that is:

$$\frac{\partial p_{c'}}{\partial \tau_c} = \begin{bmatrix} I, & -p_{c'}^3 \end{bmatrix} \tag{17}$$

On the other hand, the updated rotation matrix part is:

$$R_{c'm} = R_{c'c}R_{cm} \tag{18}$$

Where $R_{c'c}$ corresponds to the Lie algebra $\phi_c$, so we can obtain the derivative of the matrix $R_{cm}$ with respect to $\phi_c$:

$$\frac{\partial R_{c'm}}{\partial \phi_c} = \begin{bmatrix} -R_{cm}^{:1}, & -R_{cm}^{:2}, & -R_{cm}^{:3} \end{bmatrix}^T \tag{19}$$

Finally, from $t_{c'} = (\phi_c + I)R_{cm} + \rho_c$, we can conclude that:

$$\frac{\partial t_{c'm}}{\partial \rho_c} = I \tag{20}$$

Through the above derivation, we have obtained $\frac{\partial p_c}{\partial T_{cm}}$, $\frac{\partial R_{cm}}{\partial T_{cm}}$, $\frac{\partial t_{cm}}{\partial T_{cm}}$, $\frac{\partial p_i}{\partial T_{cm}}$, $\frac{\partial \Sigma_i}{\partial T_{cm}}$, $\frac{\partial c_i}{\partial T_{cm}}$, which allows us to derive completing the construction of the back propagation chain and enabling the gradient descent update for the pose $T_{cm}$.

*3) Object Refiner:* In the second stage, the object of the update changes from the pose of the camera relative to the object to the pose of the object relative to the camera. Let the updated object coordinate system be $m'$. Since each object point $p_m$ on the object is rigidly attached to the object coordinate system, its coordinate values in the object coordinate system will not change, that is:

$$p_m = p_{m'} \quad p_c = T_{cm}T_{mm'}p_m \tag{21}$$

$T_{mm'} \in SE(3)$ can be viewed as a right perturbation applied to $T_{cm}$. Let the Lie algebra corresponding to $T_{mm'}$ be denoted as:

$$\tau_m = \begin{bmatrix} \rho_m & \varphi_m \end{bmatrix}^T \in se(3), \quad p_c = T_{cm}\exp(\tau_m)p_m \tag{22}$$

By drawing an analogy to the derivation in Camera Refiner, we can obtain:

$$\frac{\partial p_c}{\partial \tau_m} = \begin{bmatrix} T_{cm} \bullet I, & -T_{cm} \bullet p_m^{:3} \end{bmatrix}, \tag{23}$$

$$\frac{\partial R_{cm}}{\partial \phi_m} = \begin{bmatrix} -R_{cm_{1,:}}, & -R_{cm_{2,:}}, & -R_{cm_{3,:}} \end{bmatrix} \tag{24}$$

$$\frac{\partial t_{c'm}}{\partial \rho_c} = R_{cm} \tag{25}$$

Through the derivation above, the backpropagation for Object Refiner has also been implemented. In the experiments, we will implement the aforementioned backpropagation chain using CUDA programming, enabling efficient computation for pose correction through the application of gradient descent algorithms.

*4) Environment adoption:* Since the 3DGS model is a type of self-emissive model, the lighting and shading characteristics of the model are not derived from its relative pose to the light source. Instead, they are obtained through the superposition of the RGB colors of each Gaussian sphere. This is a distinctive feature of the RayCast rendering algorithm. To address issues such as reflections and shadows under varying lighting conditions, we leverage the learnable nature of the 3DGS color parameters and the anisotropic properties of color parameters expressed by spherical harmonics. This allows the model to adapt to changes in lighting while adjusting its pose, thereby enhancing the accuracy of the model and preventing angle miscorrections due to lighting or shadow issues.

In this step, we set the 16 spherical harmonic parameters of the Gaussian model as learnable parameters, along with the rotational pose parameter, rot, of the Gaussian spheres. We have observed that during the model's color learning process, there is a tendency for the back side of the Gaussian spheres to be assigned a black color. Allowing the Gaussian spheres to rotate freely can accelerate the learning efficiency and prevent overfitting of the colors, as well as mitigate the issue of vanishing gradients.

Additionally, we lock other parameters, such as the scale parameter of the Gaussian spheres, their position parameters (xyz) relative to the object coordinate system, and their transparency. This is to prevent the model from compromising its original structure during iterations in an attempt to forcefully fit the target image. Such compromises could negatively impact the accuracy of angle estimation.

By carefully managing these parameters, we ensure that the model retains its integrity while effectively adapting to various lighting conditions. This approach not only enhances the model's performance but also maintains the precision required for accurate angle calculations, ultimately leading to improved results in rendering.

## 4. EXPERIMENTAL RESULTS AND ANALYSES

In order to evaluate the effectiveness of the proposed model, this section conducts a comparative analysis of its performance against a range of state-of-the-art deep-learning 6D pose estimation models, including Pix2Pose [44], SSD-6D [53], Lienet [54], Cai [55], DPOD [56], PVNet [57], CDPN [13]

### A. Experimental Dataset and Settings

Experiments were conducted on two publicly accessible datasets for 6D pose estimation: Linemod (LM) [7] and

**Linemod (LM)** [7]: The LM dataset consists of 15 registered video sequences, each containing over 1100 frames. The object scales range from 100 mm to 300 mm. There are significant variations in illumination intensity of the images captured under the same model, along with minimal occlusion phenomena. We referenced the majority of 6D pose estimation methods [31, 37] and selected 13 categories to evaluate the performance of the model, including ape, bvise, cam, can, cat, driller, duck, eggbox, glue, holep, iron, lamp and phone.
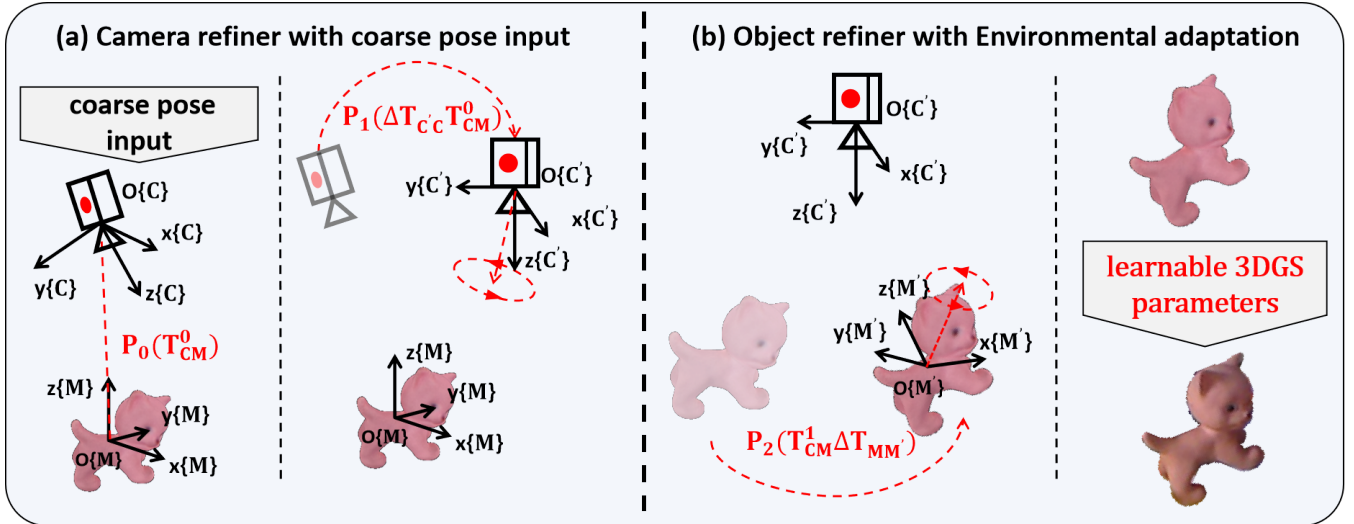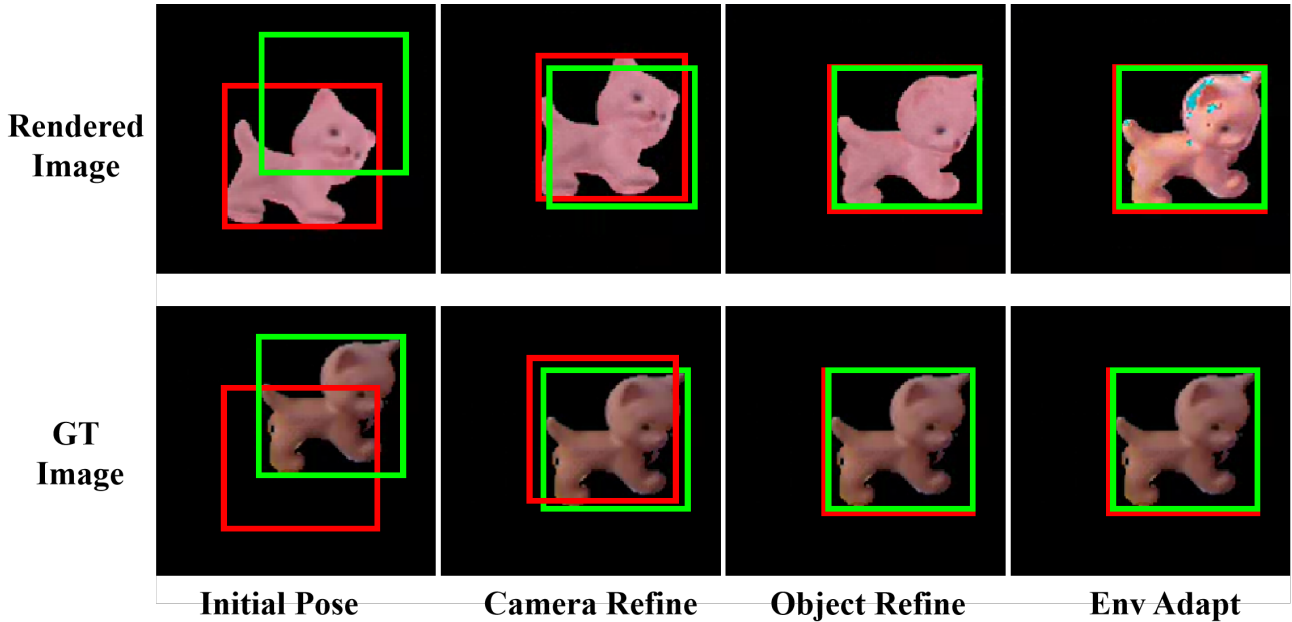
Fig. 3. The structure of the GS-Refiner

TABLE I
COMPARISON WITH OTHER METHODS ON THE LINEMOD TEST SET. (ADD-0.1D)

| Method | Publish | ape | bvise | cam | can | cat | driller | duck | eggbox* | glue* | holep | iron | lamp | phone | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SSD-6D[53] | [ICCV 2017] | 65.0 | 80.0 | 78.0 | 86.0 | 70.0 | 73.0 | 66.0 | 100.0 | 100.0 | 49.0 | 78.0 | 73.0 | 79.0 | 76.7 |
| Pix2Pose[44] | [ICCV 2019] | 58.1 | 91.0 | 60.9 | 84.4 | 65.0 | 76.3 | 43.8 | 96.8 | 79.4 | 74.8 | 83.4 | 82.0 | 45.0 | 72.4 |
| DPOD[56] | [ICCV 2019] | 53.3 | 95.2 | 90.0 | 94.1 | 60.4 | 97.4 | 66.0 | 99.6 | 93.8 | 64.9 | 99.8 | 88.1 | 71.4 | 82.6 |
| CDPN[13] | [ICCV 2019] | 64.4 | 97.8 | 91.7 | 95.9 | 83.8 | 96.2 | 66.8 | 99.7 | 99.6 | 85.8 | 97.9 | 97.9 | 90.8 | 89.9 |
| Cai[55] | [CVPR 2020] | 52.9 | 96.5 | 87.8 | 86.8 | 67.3 | 88.7 | 54.7 | 94.7 | 91.9 | 75.4 | 94.5 | 96.6 | 89.2 | 82.9 |
| Lienet[54] | [TCDS 2022] | 38.8 | 71.2 | 52.5 | 86.1 | 66.2 | 82.3 | 32.5 | 79.4 | 63.7 | 56.4 | 65.1 | 89.4 | 65.0 | 65.2 |
| PVNet[57] | [CVPR 2022] | 43.6 | 99.9 | 86.9 | 95.5 | 79.3 | 96.4 | 52.6 | 99.2 | 95.7 | 81.9 | 98.9 | 99.3 | 92.4 | 86.3 |
| OnePose[58] | [CVPR 2022] | 11.8 | 92.6 | 88.1 | 77.2 | 47.9 | 74.5 | 34.2 | 71.3 | 37.5 | 54.9 | 89.2 | 87.6 | 60.6 | 63.6 |
| OnePose++[59] | [NeurIPS 2022] | 31.2 | 97.3 | 88.0 | 89.8 | 70.4 | 92.5 | 42.3 | 99.7 | 48.0 | 69.7 | 97.4 | 97.8 | 76.0 | 76.9 |
| NeRF-Pose[37] | [ICCVW 2023] | 69.4 | 99.4 | 98.3 | 97.8 | 77.8 | 99.6 | 69.7 | 99.9 | 98.9 | 89.4 | 99.8 | 99.8 | 94.8 | 91.8 |
| GS-Pose[2] | [Arxiv 2024] | 71.0 | 99.8 | 98.2 | 97.7 | 86.7 | 96.2 | 77.2 | 99.6 | 98.4 | 87.4 | 99.2 | 98.9 | 85.0 | 92.0 |
| GS2pose | | **95.4(24.4↑)** | 99.2 | 97.4 | 95.1 | **93.2(6.5↑)** | 95.3 | **83.8(6.6↑)** | 100.0 | 98.4 | **93.5(4.1↑)** | 94.7 | 93.9 | **95.0(0.2↑)** | **95.0(3.0↑)** |



## 5. CONCLUSION

In conclusion, this paper presents GS2Pose, a novel method for accurate and robust 6D pose estimation of novel objects that effectively addresses the limitations of traditional approaches reliant on high-quality CAD models. By leveraging 3D Gaussian splatting and segmented RGBD images, GS2Pose demonstrates a significant advancement in the efficiency and accessibility of pose estimation. The two-stage architecture, comprising the coarse estimation via the Pose-Net and the refined estimation through the GS-Refiner, showcases a well-

integrated approach that enhances the precision of pose estimation. The experimental results on the LineMod dataset confirm the effectiveness of GS2Pose, positioning it as a competitive alternative to existing algorithms in the field.

## REFERENCES

[1] X. Deng, Y. Xiang, A. Mousavian, C. Eppner, T. Bretl, and D. Fox, "Self-supervised 6d object pose estimation for robot manipulation," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 3665–3671.

[2] D. Cai, J. Heikkilä, and E. Rahtu, "Gs-pose: Cascaded framework for generalizable segmentation-based 6d object pose estimation," *arXiv preprint arXiv:2403.10683*, 2024.

[3] E. Marchand, H. Uchiyama, and F. Spindler, "Pose estimation for augmented reality: a hands-on survey," *IEEE transactions on visualization and computer graphics*, vol. 22, no. 12, pp. 2633–2651, 2015.

[4] Y. Su, J. Rambach, N. Minaskan, P. Lesur, A. Pagani, and D. Stricker, "Deep multi-state object pose estimation for augmented reality assembly," in *2019 IEEE International Symposium on Mixed and Augmented Reality Adjunct (ISMAR-Adjunct)*. IEEE, 2019, pp. 222–227.

[5] D. G. Lowe, "Object recognition from local scale-invariant features," in *Proceedings of the seventh IEEE international conference on computer vision*, vol. 2. Ieee, 1999, pp. 1150–1157.

[6] V. Lepetit, P. Fua *et al.*, "Monocular model-based 3d tracking of rigid objects: A survey," *Foundations and Trends® in Computer Graphics and Vision*, vol. 1, no. 1, pp. 1–89, 2005.

[7] S. Hinterstoisser, V. Lepetit, S. Ilic, S. Holzer, G. Bradski, K. Konolige, and N. Navab, "Model based training, detection and pose estimation of texture-less 3d objects in heavily cluttered scenes," in *Computer Vision–ACCV 2012: 11th Asian Conference on Computer Vision, Daejeon, Korea, November 5-9, 2012, Revised Selected Papers, Part I 11*. Springer, 2013, pp. 548–562.

[8] W. Huang, C. Wang, R. Zhang, Y. Li, J. Wu, and L. Fei-Fei, "Voxposer: Composable 3d value maps for robotic manipulation with language models," *arXiv preprint arXiv:2307.05973*, 2023.

[9] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu *et al.*, "Palm-e: An embodied multimodal language model," *arXiv preprint arXiv:2303.03378*, 2023.

[10] Y. Long, X. Li, W. Cai, and H. Dong, "Discuss before moving: Visual language navigation via multi-expert discussions," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 17 380–17 387.

[11] J. Lin, L. Liu, D. Lu, and K. Jia, "Sam-6d: Segment anything model meets zero-shot 6d object pose estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 27 906–27 916.

[12] Y. Li, G. Wang, X. Ji, Y. Xiang, and D. Fox, "Deepim: Deep iterative matching for 6d pose estimation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 683–698.

[13] Z. Li, G. Wang, and X. Ji, "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*, 2019, pp. 7678–7687.

[14] Y. He, W. Sun, H. Huang, J. Liu, H. Fan, and J. Sun, "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*, 2020, pp. 11 632–11 641.

[15] Y. He, H. Huang, H. Fan, Q. Chen, and J. Sun, "Ffb6d: A full flow bidirectional fusion network for 6d pose estimation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 3003–3013.

[16] Y. Su, M. Saleh, T. Fetzer, J. Rambach, N. Navab, B. Busam, D. Stricker, and F. Tombari, "Zebrapose: Coarse to fine surface encoding for 6dof object pose estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 6738–6748.

[17] C. Wang, D. Xu, Y. Zhu, R. Martín-Martín, C. Lu, L. Fei-Fei, and S. Savarese, "Densefusion: 6d object pose estimation by iterative dense fusion," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 3343–3352.

[18] G. Wang, F. Manhardt, F. Tombari, and X. Ji, "Gdr-net: Geometry-guided direct regression network for monocular 6d object pose estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 16 611–16 621.

[19] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, "Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes," *arXiv preprint arXiv:1711.00199*, 2017.

[20] H. Wang, S. Sridhar, J. Huang, J. Valentin, S. Song, and L. J. Guibas, "Normalized object coordinate space for category-level 6d object pose and size estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2642–2651.

[21] S. Song and J. Xiao, "Deep sliding shapes for amodal 3d object detection in rgb-d images," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 808–816.

[22] C. R. Qi, W. Liu, C. Wu, H. Su, and L. J. Guibas, "Frustum pointnets for 3d object detection from rgb-d data," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 918–927.

[23] X. Chen, K. Kundu, Z. Zhang, H. Ma, S. Fidler, and R. Urtasun, "Monocular 3d object detection for autonomous driving," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2147–2156.

[24] A. Mousavian, D. Anguelov, J. Flynn, and J. Kosecka, "3d bounding box estimation using deep learning and geometry," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2017, pp. 7074–7082.

[25] Y. Xiang, W. Choi, Y. Lin, and S. Savarese, "Data-driven 3d voxel patterns for object category recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1903–1911.

[26] M. Tian, M. H. Ang, and G. H. Lee, "Shape prior deformation for categorical 6d object pose and size estimation," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXI 16*. Springer, 2020, pp. 530–546.

[27] C. Leiter, R. Zhang, Y. Chen, J. Belouadi, D. Larionov, V. Fresen, and S. Eger, "Chatgpt: A meta-analysis after 2.5 months," *Machine Learning with Applications*, vol. 16, p. 100541, 2024.

[28] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, "Segment anything," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 4015–4026.

[29] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong *et al.*, "A survey of large language models," *arXiv preprint arXiv:2303.18223*, 2023.

[30] V. N. Nguyen, T. Groueix, M. Salzmann, and V. Lepetit, "Gigapose: Fast and robust novel object pose estimation via one correspondence," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 9903–9913.

[31] Y. Labbé, L. Manuelli, A. Mousavian, S. Tyree, S. Birchfield, J. Tremblay, J. Carpentier, M. Aubry, D. Fox, and J. Sivic, "Megapose: 6d pose estimation of novel objects via render & compare," *arXiv preprint arXiv:2212.06870*, 2022.

[32] B. Wen, W. Yang, J. Kautz, and S. Birchfield, "Foundationpose: Unified 6d pose estimation and tracking of novel objects," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 17 868–17 879.

[33] C. Yan, D. Qu, D. Xu, B. Zhao, Z. Wang, D. Wang, and X. Li, "Gs-slam: Dense visual slam with 3d gaussian splatting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 19 595–19 604.

[34] J. Engel, T. Schöps, and D. Cremers, "Lsd-slam: Large-scale direct monocular slam," in *European conference on computer vision*. Springer, 2014, pp. 834–849.

[35] R. Mur-Artal and J. D. Tardós, "Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras," *IEEE transactions on robotics*, vol. 33, no. 5, pp. 1255–1262, 2017.

[36] N. Keetha, J. Karhade, K. M. Jatavallabhula, G. Yang, S. Scherer, D. Ramanan, and J. Luiten, "Splatam: Splat track & map 3d gaussians for dense rgb-d slam," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 21 357–21 366.

[37] F. Li, S. R. Vutukur, H. Yu, I. Shugurov, B. Busam, S. Yang, and S. Ilic, "Nerf-pose: A first-reconstruct-then-regress approach for weakly-supervised 6d object pose estimation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 2123–2133.

[38] H. Bay, T. Tuytelaars, and L. Van Gool, "Surf: Speeded up robust features," in *Computer Vision–ECCV 2006: 9th European Conference on Computer Vision, Graz, Austria, May 7-13, 2006. Proceedings, Part I 9*. Springer, 2006, pp. 404–417.

[39] A. Collet and S. S. Srinivasa, "Efficient multi-view object recognition and full pose estimation," in *2010 IEEE International Conference on Robotics and Automation*. IEEE, 2010, pp. 2050–2055.

[40] S. Hinterstoisser, S. Holzer, C. Cagniart, S. Ilic, K. Konolige, N. Navab, and V. Lepetit, "Multimodal templates for real-time detection of texture-less objects in heavily cluttered scenes," in *2011 international conference on computer vision*. IEEE, 2011, pp. 858–865.

[41] M. Rad and V. Lepetit, "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*, 2017, pp. 3828–3836.

[42] Y. Hai, R. Song, J. Li, M. Salzmann, and Y. Hu, "Rigidity-aware detection for 6d object pose estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 8927–8936.

[43] S. Zakharov, B. Planche, Z. Wu, A. Hutter, H. Kosch, and S. Ilic, "Keep it unreal: Bridging the realism gap for 2.5 d recognition with geometry priors only," in *2018 International Conference on 3D Vision (3DV)*. IEEE, 2018, pp. 1–11.

[44] K. Park, T. Patten, and M. Vincze, "Pix2pose: Pixel-wise coordinate regression of objects for 6d pose estimation," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 7668–7677.

[45] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering." *ACM Trans. Graph.*, vol. 42, no. 4, pp. 139–1, 2023.

[46] H. Matsuki, R. Murai, P. H. Kelly, and A. J. Davison, "Gaussian splatting slam," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 18 039–18 048.

[47] V. Yugay, Y. Li, T. Gevers, and M. R. Oswald, "Gaussian-slam: Photo-realistic dense slam with gaussian splatting," *arXiv preprint arXiv:2312.10070*, 2023.

[48] Y. Chen, L. Wang, Q. Li, H. Xiao, S. Zhang, H. Yao, and Y. Liu, "Monogaussianavatar: Monocular gaussian point-based head avatar," in *ACM SIGGRAPH 2024 Conference Papers*, 2024, pp. 1–9.

[49] J. Wang, J.-C. Xie, X. Li, F. Xu, C.-M. Pun, and H. Gao, "Gaussianhead: Impressive head avatars with learnable gaussian diffusion," *arXiv preprint arXiv:2312.01632*,

2023.

[50] Z. Chen, F. Wang, Y. Wang, and H. Liu, "Text-to-3d us-
ing gaussian splatting," in *Proceedings of the IEEE/CVF
Conference on Computer Vision and Pattern Recognition*,
2024, pp. 21 401–21 412.

[51] J. Tang, J. Ren, H. Zhou, Z. Liu, and G. Zeng, "Dream-
gaussian: Generative gaussian splatting for efficient
3d content creation," *arXiv preprint arXiv:2309.16653*,
2023.

[52] T. Yi, J. Fang, G. Wu, L. Xie, X. Zhang, W. Liu, Q. Tian,
and X. Wang, "Gaussiandreamer: Fast generation from
text to 3d gaussian splatting with point cloud priors,"
*arXiv preprint arXiv:2310.08529*, 2023.

[53] W. Kehl, F. Manhardt, F. Tombari, S. Ilic, and N. Navab,
"Ssd-6d: Making rgb-based 3d detection and 6d pose
estimation great again," in *Proceedings of the IEEE
international conference on computer vision*, 2017, pp.
1521–1529.

[54] M. Karnati, A. Seal, A. Yazidi, and O. Krejcar, "Lienet:
A deep convolution neural network framework for de-
tecting deception," *IEEE transactions on cognitive and
developmental systems*, vol. 14, no. 3, pp. 971–984, 2021.

[55] M. Cai and I. Reid, "Reconstruct locally, localize glob-
ally: A model free method for object pose estimation," in
*Proceedings of the IEEE/CVF Conference on Computer
Vision and Pattern Recognition*, 2020, pp. 3153–3163.

[56] S. Zakharov, I. Shugurov, and S. Ilic, "Dpod: 6d pose
object detector and refiner," in *Proceedings of the
IEEE/CVF international conference on computer vision*,
2019, pp. 1941–1950.

[57] S. Peng, Y. Liu, Q. Huang, X. Zhou, and H. Bao, "Pvnet:
Pixel-wise voting network for 6dof pose estimation," in
*Proceedings of the IEEE/CVF conference on computer
vision and pattern recognition*, 2019, pp. 4561–4570.

[58] J. Sun, Z. Wang, S. Zhang, X. He, H. Zhao, G. Zhang,
and X. Zhou, "Onepose: One-shot object pose estimation
without cad models," in *Proceedings of the IEEE/CVF
Conference on Computer Vision and Pattern Recognition*,
2022, pp. 6825–6834.

[59] X. He, J. Sun, Y. Wang, D. Huang, H. Bao, and
X. Zhou, "Onepose++: Keypoint-free one-shot object
pose estimation without cad models," *Advances in Neural
Information Processing Systems*, vol. 35, pp. 35 103–
35 115, 2022.