# SparseGrasp: Robotic Grasping via 3D Semantic Gaussian Splatting from Sparse Multi-View RGB Images

Junqiu Yu[1*], Xinlin Ren[1*], Yongchong Gu[1], Haitao Lin[1], Tianyu Wang[1], Yi Zhu[2†],
Hang Xu[2], Yu-Gang Jiang[1], Xiangyang Xue[1], Yanwei Fu[1†]

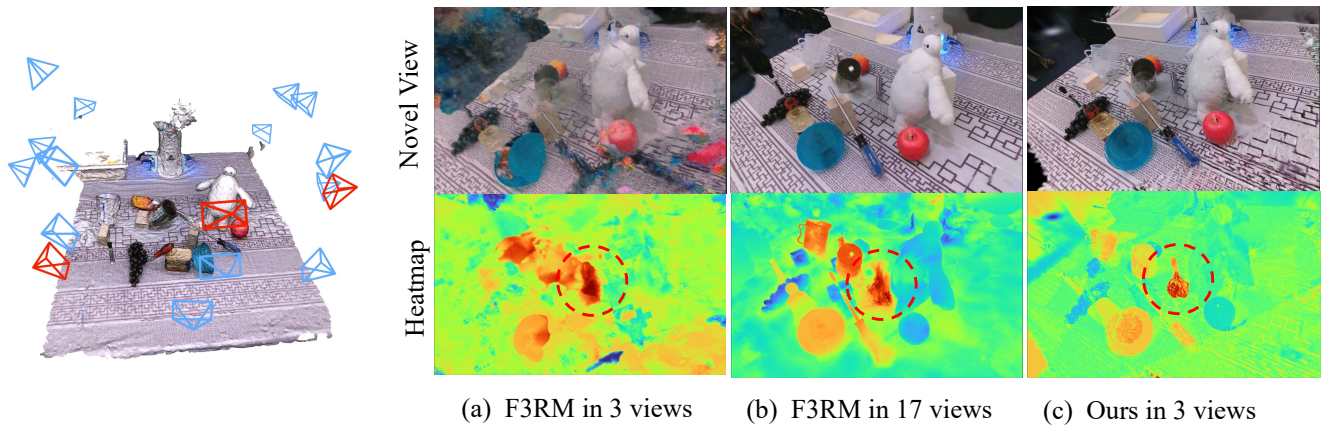(a) F3RM in 3 views     (b) F3RM in 17 views     (c) Ours in 3 views

Fig. 1: We present a comparison between our SparseGrasp and F3RM under both sparse and dense view settings. The top row shows the novel view images, while the bottom row displays the heat map of feature field using the text "whisk" as a query. Remarkably, our method, utilizing only **3** view images, achieves performance comparable to F3RM, which is trained with 17 views.

*Abstract*— Language-guided robotic grasping is a rapidly advancing field where robots are instructed using human language to grasp specific objects. However, existing methods often depend on dense camera views [1], [2] and struggle to quickly update scenes, limiting their effectiveness in changeable environments. In contrast, we propose SparseGrasp, a novel open-vocabulary robotic grasping system that operates efficiently with sparse-view RGB images and handles scene updates fastly. Our system builds upon and significantly enhances existing computer vision modules in robotic learning. Specifically, SparseGrasp utilizes DUSt3R [3] to generate a dense point cloud as the initialization for 3D Gaussian Splatting (3DGS), maintaining high fidelity even under sparse supervision. Importantly, SparseGrasp incorporates semantic awareness from recent vision foundation models. To further improve processing efficiency, we repurpose Principal Component Analysis (PCA) to compress features from 2D models. Additionally, we introduce a novel render-and-compare strategy that ensures rapid scene updates, enabling multi-turn grasping in changeable environments. Experimental results show that SparseGrasp significantly outperforms state-of-the-art methods in terms of both speed and adaptability, providing a robust solution for multi-turn grasping in changeable environment.

## I. INTRODUCTION

Language-guided robotic grasping is an emerging field where robots use human language instructions to grasp specific objects. Imagine a scenario where a robot can swiftly and reliably grasp objects based on verbal commands while simultaneously adapting to changes in the environment. This would enable the robot to seamlessly follow up on new instructions. Achieving this requires the robot to not only locate objects accurately through language but also to understand their geometry, ensuring more precise and effective grasping.

Recent efforts [1], [2], [4], [5] aim to reconstruct scenes while distilling semantic information from 2D foundation models like CLIP [6]. However, these methods heavily rely on dense view reconstruction, which demands significant training time and multi-view capturing via the robot's on-board camera. Some viewpoints are challenging for robotic arms to reach, requiring precise path planning, while performance drops when only limited views are available. For instance, as shown in Fig. 1, F3RM [1] struggles with object geometry using only 3 views, and even with 17 views, semantic distillation remains suboptimal. Additionally, these

†: Corresponding Authour.
∗: Equal Contribution.
[1] Junqiu Yu, Xinlin Ren, Yongchong Gu, Haitao Lin, Tianyu Wang, Yu-Gang Jiang, Xiangyang Xue and Yanwei Fu are with Fudan University. {jqyu20, xlren20,htlin19,ygj ,xyxue,yanweifu}@fudan.edu.cn, {yongchonggu22, tywang22}@m.fudan.edu.cn
[2] Yi Zhu and Hang Xu are with the Department of Noah's Ark Lab, Huawei Technology, Shanghai 200433, China. {zhuyi36, xu.hang}@huawei.com

methods are designed for static environments, requiring dense view capture even for minor changes. This hinders their ability to adapt to changeable environments and perform multi-turn grasping, limiting their real-world applicability. Here, 'changeable environment' refers to scenarios where objects may be moved, while 'multi-turn' indicates the robot executing a sequence of language commands.

To address these challenges, we introduce SparseGrasp, which enables fast scene reconstruction (around 240 seconds) and fast updation using only **sparse-view RGB images**. Since 3D Gaussian Splatting (3DGS) [7] tends to overfit with sparse views, we integrate DUSt3R [3] to generate a dense point cloud as initialization, offering greater robustness than sparse point clouds, as in Fig. 3. Besides, we further incorporate 2D foundation models like MaskCLIP [8] and Segment Anything (SAM) [9] to extract dense semantic features, which are distilled into 3DGS. Unlike prior methods [2] that require multiple CLIP calls for each object identified by SAM, we propose a more efficient approach: applying patch-level CLIP and SAM once per image to generate comprehensive semantic features. Further, given the high dimensionality of semantic features, directly distilling them into 3DGS is impractical. Therefore, we use PCA to compress the features, reducing dimensions significantly (from 768 to 16). In addition, we improve GraspNet [10] by generating grasps directly from 3DGS, eliminating the voxelization and depth back-projection required by methods like F3RM [1] and LERF-TOGO [2]. Finally, for objects that have changed positions in the scene, we propose a render-and-compare strategy to efficiently update their scene representations, eliminating the need for full scene reconstruction.

Our contributions are as follows: 1) We present the *SparseGrasp* system for rapid scene reconstruction (240 seconds) and fast updates using sparse-view RGB images, overcoming the dense-view dependency in prior methods. 2) We propose *3D Semantic Gaussian Splatting*, incorporating DUSt3R for robust dense point cloud initialization, and enhancing semantic distillation by efficiently applying MaskCLIP and SAM once per image to extract dense semantic features, followed by PCA for feature reduction. 3) We improve GraspNet by generating grasps directly from 3DGS, eliminating voxelization and depth back-projection. 4) Finally, we introduce a *render-and-compare strategy* for efficiently updating scene representations when objects change, avoiding the need for full scene reconstruction.

## II. RELATED WORK

**Rendering Methods for Robotics.** Neural Radiance Fields (NeRF) [11] have recently advanced to enable high-quality scene reconstruction from RGB images, leading to their integration into various robotics applications such as grasping [12]–[15] and navigation [16]–[19]. Despite their exceptional reconstruction quality, NeRF methods typically focus on complete scene reconstruction and are limited by the time-consuming process of capturing multiple images and training models, making them suitable only for static environments. While Evo-NeRF [12] aimed to speed up scene updates

for sequential grasping, it compromised on object geometry accuracy and still required multiple perspectives, reducing its efficiency for language-guided grasping. In contrast, 3DGS [7] offers efficient scene reconstruction through adaptive density control and multi-view images, proving valuable in SLAM [20]–[22]. However, its application to robotic grasping has been limited till now, with notable use only in simulated environments [23]. Our approach innovatively applies 3DGS to *RGB images from sparse camera views*, achieving faster and more efficient scene reconstruction and updates. This enables precise grasping of objects in both static and dynamically changing environments. To the best of our knowledge, SparseGrasp is the first method to leverage only *sparse view RGB images* for scene reconstruction and object grounding in changeable environment.

**Language-guided Grasping.** The integration of Computer Vision (CV) and Natural Language Processing (NLP) enhances robots' language comprehension. Early studies [24]–[28] achieved object grounding with 2D models. More recent approaches [29], [30] combine visual grounding with 6D pose estimation [31] for robotic grasping, though they may struggle with precision in fine-grained object manipulation.

Recent advancements have integrated semantic information into 3D representations. Works like CLIPort [32] align semantic features with point clouds or scene depth, but rely exclusively on depth data, resulting in suboptimal alignment. Other approaches [1], [2], [5] improve semantic information embedding within 3D scene reconstruction. Specifically, F3RM [1] and LERF-TOGO [2] train NeRF using RGB images, while GNFactor [5] leverages RGB-D data to generate voxelized scene representations and aligns semantic features with voxels. However, these methods may struggle when scene changes and require additional time to adapt. Our method uniquely leverages 3DGS to learn semantic information directly from sparse RGB images and employs render-and-compare strategy, enabling robots to perform sequential operations effectively even in changeable environment.

## III. METHOD

Given sparse view images $I$, our goal is to continuously pick and place objects according to the open-vocabulary languages in static and changeable scenes. Formally, this involves generating 6-DoF grasp poses $\mathbf{T} = (R, t)$ where $R, t$ represent the rotation and translation components of the grasp pose, respectively. Our approach interprets language commands to guide robotic actions in diverse and changing environments. In Fig. 2, we start by collecting RGB images from sparse viewpoints and use DUSt3R [3] to initialize 3DGS with dense point clouds. We then integrate semantic features into 3DGS, focusing on the extraction and compression of pixel-level CLIP features, and employ the Render-and-Compare strategy for rapid scene updates. Additionally, we introduce language-guided robotic grasping.

**Preliminary: 3D Gaussian Splatting.** 3DGS reconstructs scenes from images by approximating objects in 3D space using Gaussians, initialized with point clouds generated by COLMAP [33] from dense view images. Each 3D Gaussian
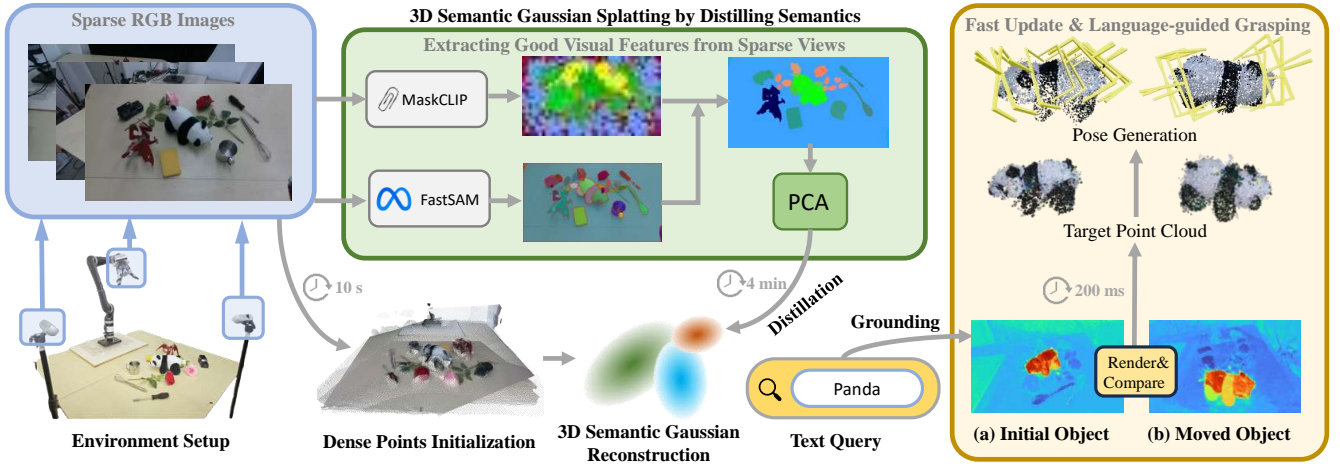
Fig. 2: Our architecture. It starts with collecting sparse view images and generating dense point clouds to initialize 3DGS. Next, we integrate FastSAM and MaskCLIP to generate average features within each mask. Then, PCA is applied to compress the whole average features in a low dimension, then distilled into 3DGS. Given an open-vocabulary language instruction, our system can locate the target object and generate appropriate grasp poses. When scene changes, the Render-and-Compare strategy enables fast scene updates.



Dense Point Initialization | Novel View Image | Sparse Point Initialization | Novel View Image
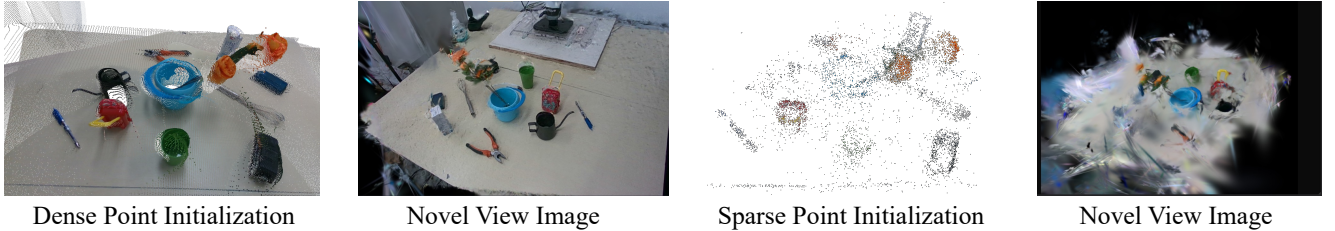
Fig. 3: Comparison of dense point initialization v.s. sparse point initialization in sparse view images. Initializing with sparse points often leads to overfitting with sparse view images.

$g_i$ is characterized by a 3D coordinate $p_i \in \mathbb{R}^3$, a scaling factor $s_i \in \mathbb{R}^3$, a rotation quaternion $q_i \in \mathbb{R}^4$, and an opacity value $\alpha_i \in \mathbb{R}$. Color $c_i \in \mathbb{R}^3$ can be derived from spherical harmonics with given direction. Using tile-based rasterization, the pixel color $C$ is calculated:

$$C = \sum_{i \in \mathcal{N}} T_i c_i \alpha_i \qquad (1)$$

where $T_i = \prod_{j=1}^{i-1}(1 - \alpha_j)$. During the training phase, the loss is computed to optimize the parameters. Leveraging this rasterization technique, 3DGS facilitates real-time rendering.

### A. Extracting Good Visual Features from Sparse Views

**Dense Point Initialization**. We initialize with dense points using DUSt3R [3] to generate dense point clouds in sparse view scene reconstruction. Traditional 3DGS requires dense-view images and uses COLMAP [33] to create sparse point clouds, which can be problematic in sparse-view settings. NeRF-based methods also tend to overfit with sparse-view images. While many methods [34], [35] use RGB-D sensors to supplement depth information, this adds hardware requirements. Our approach overcomes these limitations by enabling fast reconstruction from sparse views and demon-

strates robustness in novel views, as illustrated in Fig. 3.

**Extracting Dense Visual Features.** To extract dense visual features, we address two key challenges: 1) achieving per-pixel feature extraction with clear object boundaries, and 2) minimizing time expenditure.

To overcome CLIP's limitation of extracting features only at the image level, we adapt the approach from F3RM [1] by using the MaskCLIP [8] reparameterization trick for patch-level alignment. However, patch-level CLIP features alone often yield imprecise object boundaries, complicating accurate object grounding.

So, instead of using DINO features for regularization [36] or invoking CLIP on numerous masked segments from SAM [9], which can be inefficient, our method extracts patch-level features with MaskCLIP and applies nearest neighbor upsampling to match the input image resolution. We then average the features within each mask generated by FastSAM [37], producing a dense feature matrix $\mathbf{F}^{\text{sem}} \in \mathbb{R}^{N \times C}$, where N is the number of masks.

A key consideration is that the masks produced by Fast-SAM may overlap. To resolve this, our method first sorts the masks by size in ascending order, ensuring that smaller masks are prioritized. We then create a new mask where

the smaller masks take precedence in overlapping regions by selecting the first non-zero mask at each pixel. It ensures that the final mask resolves any overlaps by favoring smaller regions. Through this approach, we successfully extract precise and dense pixel-level semantic features at around 180ms. In addition, since our method requires only sparse views as input, the time spent on feature extraction is negligible.

### B. 3D Semantic Gaussian Splatting by Distilling Semantics

**Compression of Language Features.** Unlike NeRF, 3DGS uses over 100,000 Gaussians, leading to high memory and computation costs with high-dimensional CLIP features. While Feature-3DGS [38] uses a lightweight decoder and LangSplat [39] employs a scene-specific autoencoder, both methods either reduce rendering speed or require extensive training (over 30 minutes). In contrast, our method uses PCA for efficient feature compression. By averaging object features, we apply PCA to reduce the feature set to as few as 16 dimensions, maintaining sufficient accuracy for scene representation, as illustrated in Fig. 9. The compressed feature is denoted as $\mathbf{F}_{pca}^{sem}$.

**Distilling Semantic Features into 3D Gaussians.** We enhance 3DGS by integrating language features into each 3D Gaussian. Specifically, we use differential rasterization to derive dense semantic features $\hat{\mathbf{F}}^{sem}$ for each pixel.

$$\hat{F}_i^{sem} = \sum_{i \in N} f_i \alpha_i T_i \qquad (2)$$

Where, $f_i$ denotes the semantic feature within each 3D Gaussian. Rather than rasterizing the RGB image and feature map separately, we use a joint optimization approach where both are processed through the same tile-based rasterization. Unlike LangSplat [39], which calculates the gradient of $\alpha$ using language features, our method derives this gradient solely from RGB images, given their more reliable supervision. We optimize our 3DGS with the following objective:

$$\mathcal{L} = \mathcal{L}_{rec} + \lambda_2 \cdot \mathcal{L}_{sem}, \qquad (3)$$
$$\mathcal{L}_{rec} = (1 - \lambda_1)\mathcal{L}_1 + \lambda_1 \cdot \mathcal{L}_{D-SSIM}, \qquad (4)$$
$$\mathcal{L}_{sem} = \mathcal{L}_1(\mathbf{F}_{pca}^{sem}, \hat{\mathbf{F}}^{sem}) \qquad (5)$$

where $\mathcal{L}_1$ denotes the L1 loss and $\mathcal{L}_{D-SSIM}$ is the SSIM loss between the rendered image and the ground truth. The reconstruction loss $L_{rec}$ is formulated as in the original 3DGS [7], and $L_{sem}$ represents the L1 loss between the rendered and compressed semantic features. The terms $\lambda_1$ and $\lambda_2$ are set to 0.2 and 1, respectively.

### C. Render and Compare for Fast Scene Updating

In changeable environments where objects may be moved, resulting in unknown translations and rotations, we here propose a method that uses sparse view images, $\{I_{mov}\}$ (three for experiment), to predict the object's translation and rotation. These predictions are utilized as optimization parameters. Specifically, we begin by applying MOG2 algo-

rithm [40][1] to detect moved pixels between the initial and current frames, whose center coordinates in current frames is denoted as $\{d_{gt}\}$. Next, we compute the mean semantic features of these pixels and use their cosine similarity to identify the 3D Gaussians of the moved objects, $\{g_i\}$. During optimization, we adjust $\{g_i\}$ using the predicted translation and rotation, while keeping other Gaussians fixed. Finally, we render two images: $I_{obj}$ with the adjusted Gaussians and $I_{pred}$ with all Gaussians, and define the optimization $\mathcal{L}$ as,

$$\mathcal{L} = \mathcal{L}_1(I_{mov}, I_{pred}) + \lambda_3\mathcal{L}_1(d_{pred}, d_{gt}) \qquad (6)$$

where $d_{pred}$ is the predicted location of the object on $I_{obj}$. We set the $\lambda_3$ to 0.1 to balance the loss of pixels and 2D distance. Since only the translation and rotation parameters of the objects require optimization, the process is efficient and typically completes in around 200 ms.
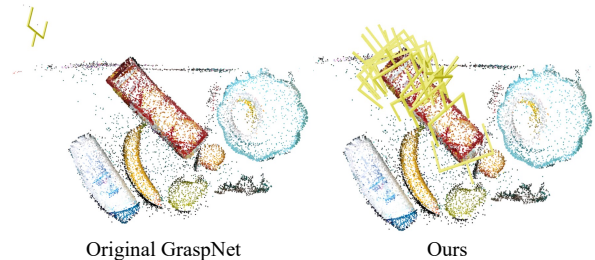
### D. Language-guided Robotic Grasping



Fig. 4: Effectiveness of our grasp model: Unlike original GraspNet, which failed to generate grasp poses using 3DGS's centers, our model successfully generates grasp poses.

In language-guided robotic grasping, existing methods [1], [2] often require processes like voxelization or multi-view depth and RGB information for scene reconstruction, which can introduce approximation errors. Unlike these methods, which resample the scene, our approach avoids direct use of sparse 3D Gaussians as point clouds for GraspNet, as illustrated in Fig. 4. Instead, we retrain GraspNet using $p_i, s_i, q_i$ as inputs. For training, we modify the original GraspNet dataset by selecting 100 scenes, segmenting objects and backgrounds, and reconstructing the scenes separately using RGB images of either objects or backgrounds. We then combine these to label the 3D Gaussians from the object segments as 'objectness'. To account for real-world noise, we randomly add some 3D gaussian noise to each scene and use varying densification thresholds to enhance robustness.

## IV. EXPERIMENTS

**Environment and Setup.** In our physical robotic experiment, we use a KINOVA Gen2 robot with a 6-DoF curved wrist and a KG-3 gripper. The robot is equipped with three common cameras with good RGB images. The system runs on a desktop with an NVIDIA GTX A6000 GPU.

**Implementation Details.** To reconstruct the scene at a real-world scale, we use the calibrated camera poses as input

---

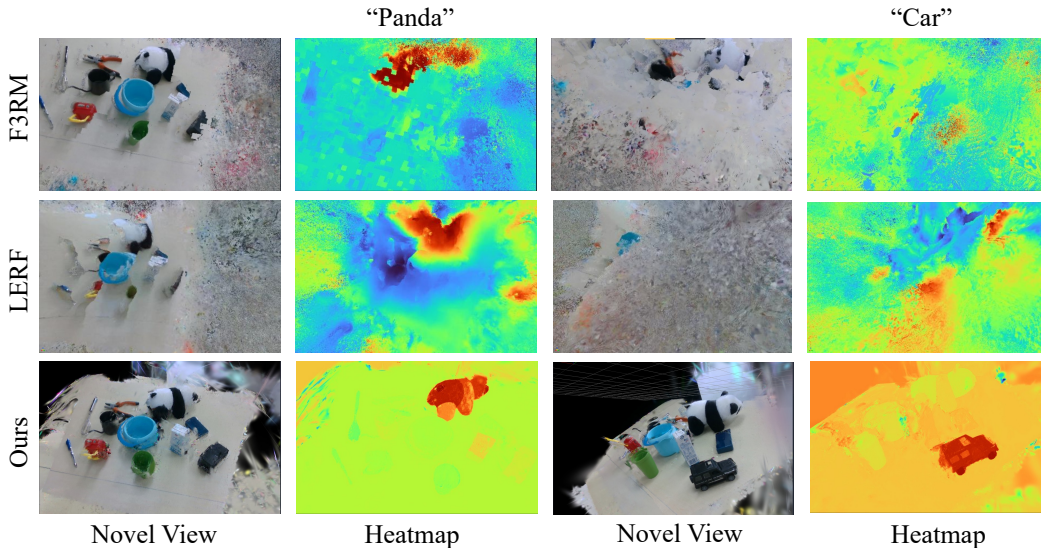[1]We use this algorithm, as it is simple and good enough in our task.

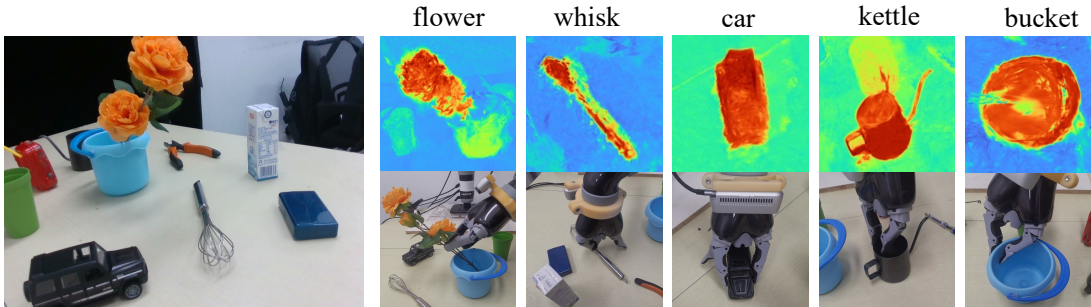Fig. 5: Qualtitative results of reconstruction and semantic distillation results.



Fig. 6: Qualtitative results of language-guided grasping in the static environment. (Top Row) Heatmaps of given text queries. (Bottom Row) Robot executing grasps sequentially without rescanning.

to DUSt3R. The image resolution is also resized to 336 to serve as input for the MaskCLIP ViT-B/16 model, which extracts language features from each image. For 2D mask segmentation, we employ the SAM ViT-H model. Besides, we resample the initial point cloud to approximately 100,000 points and jointly train our 3DGS model while distilling semantic features. The pipeline is trained over 7,000 iterations, taking about 4 minutes. Due to the good initialization, the densification interval is set to 200 iterations.

### A. Results of Reconstruction and Semantic Distillation

We compare our method with other state-of-art methods such as F3RM [1] and LERF-TOGO [2] in terms of scene reconstruction and semantic distillation. Our method, F3RM and LERF are trained with 3 RGB images. Under this condition, both F3RM and LERF-TOGO exhibit overfitting. As illustrated in Fig. 5, these methods fail to maintain the geometric integrity of objects, let alone produce reliable semantic outcomes.

### B. Results of Language-guided Grasping.

We first present results for language-guided grasping in static environments, followed by validation of our render-

TABLE I: Quantitative results of success rate by using estimated grasping poses

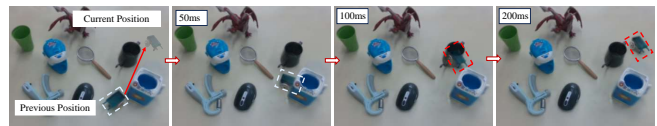| | Flower | Whisk | Dragon | Car | Kettle's Lip | Panda | ScrewDriver | Total |
|---|---|---|---|---|---|---|---|---|
| F3RM(3 Views) | 1/10 | 0/10 | 1/10 | 1/10 | 0/10 | 0/10 | 1/10 | 4/70 |
| LERF-TOGO(3 Views) | 0/10 | 1/10 | 0/10 | 1/10 | 1/10 | 1/10 | 1/10 | 5/70 |
| F3RM(17 Views) | 6/10 | 1/10 | 5/10 | 7/10 | 7/10 | 5/10 | 7/10 | 38/70 |
| LERF-TOGO(17 Views) | 6/10 | 3/10 | 4/10 | 5/10 | 5/10 | 3/10 | 7/10 | 33/70 |
| Ours(3 Views) | **8/10** | **8/10** | **9/10** | **8/10** | 7/10 | **8/10** | 7/10 | **55/70** |



Fig. 7: Optimization process of our Render&Compare module. We can generate the accurate position of the cart fastly.

and-compare method for scene updating.

**Results of Grasping Accuracy in Static Environment.** As shown in Tab. I, both F3RM and LERF-TOGO demonstrate poor grasping performance in sparse views, and even in dense views, they struggle to successfully grasp the whisk. Their failure in sparse views arises from their inability to accurately reconstruct the scene. Additionally, the difficulty

TABLE II: Result of grasping accuracy after movements.

| Methods | Time | Scene1 | Scene2 | Scene3 | Scene4 | Scene5 |
|---|---|---|---|---|---|---|
| F3RM [1] | 10min | 4/5 | **4/5** | **5/5** | **4/5** | 3/5 |
| LERF-TOGO [2] | 34min | 3/5 | 4/5 | 4/5 | 3/5 | 3/5 |
| | 50 ms | 0/5 | 0/5 | 0/5 | 1/5 | 0/5 |
| Ours | 100 ms | 3/5 | 2/5 | 3/5 | 2/5 | 3/5 |
| | 200 ms | **5/5** | **4/5** | **5/5** | **4/5** | **4/5** |

in grasping the whisk can be attributed to the implicit representation used by NeRF, which requires voxelization to generate the 6D pose. This voxelization process brings a loss of precision. In contrast, our method achieves high-precision grasping across all objects, even in the sparse view setting. This is due to its superior reconstruction accuracy and more effective semantic distillation. Additionally, our language-guided grasping results are presented in Fig. 6.

**Results of Grasping Accuracy in Changeable Scene.** In this scenario, objects could be randomly shifted by people. Rather than necessitating dense view images for scene updates like F3RM and LERF-TOGO, our approach employs Render-and-Compare for quicker updates. As illustrate in Tab. II, F3RM and LERF-TOGO are unable to update specific areas of the scene, requiring a full scene reconstruction instead. Note that the times listed for F3RM and LERF-TOGO in Tab.II exclude the additional 2 minutes required for image capturing, a step our method does not require. In contrast, our render-and-compare strategy can quickly update the scene, taking approximately 200ms.

**Result of Semantic Distillation Effect.** We use 2D IoU as a metric to evaluate the effectiveness of our semantic distillation in both training and novel views. Additionally, we compare our results with other methods using either sparse or dense views, as shown in Tab. III. F3RM and LERF-TOGO fail to preserve object geometries, leading to significantly lower IoU scores in novel views. In contrast, our approach excels at maintaining geometric integrity across novel views, demonstrating its ability to infuse semantic information into 3D objects.

TABLE III: Quantitative results of 2D IOU.

| Methods | F3RM (3 view) | LERF-TOGO (3 view) | F3RM (17 view) | LERF-TOGO (17 view) | Ours (3 view) |
|---|---|---|---|---|---|
| Training Views | 0.75 | 0.69 | 0.81 | 0.74 | **0.83** |
| Novel Views | 0.13 | 0.08 | **0.75** | 0.71 | 0.71 |

**Results of Time Consumption.** In Tab. IV, show the results of all methods in both static scenes and changeable environments. In contrast, our SparseGrasp not only significantly outperforms F3RM and LERF-TOGO, but also is much faster than LangSplat, which requires extensive preprocessing time. Furthermore, when the scene changes, these methods need to recapture images and reconstruct the entire scene, which is time-consuming. While our can quickly update the changed objects, offering a more efficient solution.

### C. More Analysis

**Results on the Different Number of Views.** Fig. 8 shows a qualitative comparison of RGB images and heatmaps of the 'metal mug' from the novel view in F3RM dataset. As the

TABLE IV: Comparison of Time Consumption.

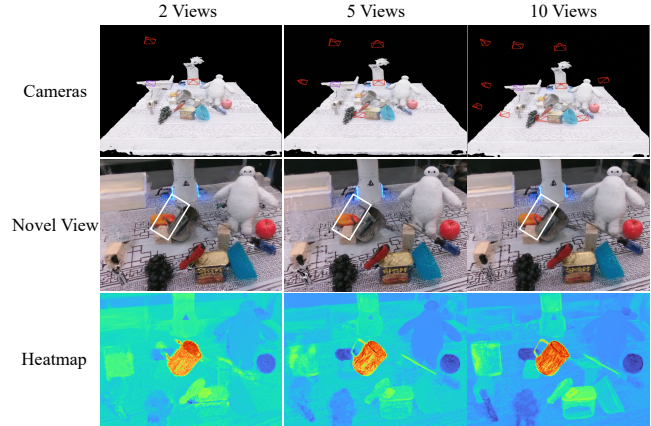| Methods | F3RM | LERF-TOGO | LangSplat | Ours |
|---|---|---|---|---|
| Static Scene | 10min | 34min | 3600min+ | **4min** |
| Scene Updation | - | - | - | 200ms |



Fig. 8: Rendered images and heatmaps for the query 'metal hug' across different camera views. Red and purple indicate training and testing views, respectively (top row).

number of training images increases, the boundary around the slim part of the metal mug becomes noticeably clearer. Additionally, Even with only two views, our method is still able to accurately distinguish the object.

**Results of Different Number of PCA's Components.** To evaluate the compression ability of PCA, we explore varying numbers of PCA components. As shown in Fig. 9, we visualize the heatmaps of the metal mug with different numbers of PCA components. We observe that using as few as 16 components already achieves effective segmentation of the metal mug.
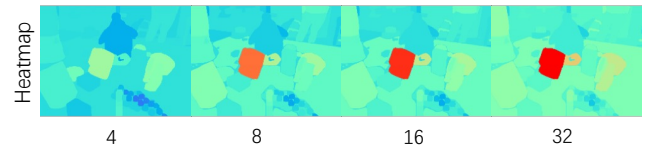


Fig. 9: Qualitative results of heatmaps when applying PCA with different numbers of components .

**Discussion and Future Work.** *1) Can SparseGrasp handle adding or replacing objects in the scene?* Yes, but requires a 'refresh'. Our system primarily focuses on multi-turn grasping in changeable environments. While it needs re-initialization when adding or removing objects, it still outperforms *F3RM* and *LERF-TOGO* in speed, as shown in Table IV. In future work, we aim to explore 3D inpainting methods [41] to improve processing speed for newly added or removed objects. *2) How about directly using optical flow or tracking?* Possible, but depth estimation and semantic integration are still needed. Some works [42], [43] use tracking or optical flow for scene updation, but these methods rely on depth data, which is difficult to estimate from sparse RGB views and increases processing time. Moreover, they

lack the semantic understanding necessary for text-based queries like ours. Our method enables fast scene updates using only sparse RGB images, without requiring depth estimation, while maintaining semantic capabilities.

## V. CONCLUSIONS

We propose **SparseGrasp**, a system for rapid scene reconstruction, understanding and updating using sparse-view RGB images. By integrating *DUSt3R* for dense point cloud initialization and *MaskCLIP* with *SAM* for efficient semantic extraction, we overcome the limitations of dense-view dependencies. Dimensionality reduction via *PCA* boosts efficiency, while our render-and-compare strategy enables fast scene updates without full reconstructions. Finally, our 3DGS-based grasping method streamlines grasp generation, avoiding voxelization issues. **SparseGrasp** improves scene understanding, speed and robustness, paving the way for future advancements.

## REFERENCES

[1] W. Shen, G. Yang, A. Yu, J. Wong, L. P. Kaelbling, and P. Isola, "Distilled feature fields enable few-shot language-guided manipulation," *arXiv preprint arXiv:2308.07931*, 2023.

[2] A. Rashid, S. Sharma, C. M. Kim, J. Kerr, L. Y. Chen, A. Kanazawa, and K. Goldberg, "Language embedded radiance fields for zero-shot task-oriented grasping," in *7th Annual Conference on Robot Learning*, 2023. [Online]. Available: https://openreview.net/forum?id=k-Fg8JDQmc

[3] S. Wang, V. Leroy, Y. Cabon, B. Chidlovskii, and J. Revaud, "Dust3r: Geometric 3d vision made easy," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 20 697–20 709.

[4] Y. Li and D. Pathak, "Object-aware gaussian splatting for robotic manipulation," in *ICRA 2024 Workshop on 3D Visual Representations for Robot Manipulation*, 2024.

[5] Y. Ze, G. Yan, Y.-H. Wu, A. Macaluso, Y. Ge, J. Ye, N. Hansen, L. E. Li, and X. Wang, "Gnfactor: Multi-task real robot learning with generalizable neural feature fields," in *Conference on Robot Learning*. PMLR, 2023, pp. 284–301.

[6] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.

[7] 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.

[8] C. Zhou, C. C. Loy, and B. Dai, "Extract free dense labels from clip," in *European Conference on Computer Vision*. Springer, 2022, pp. 696–712.

[9] 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.

[10] H.-S. Fang, C. Wang, M. Gou, and C. Lu, "Graspnet-1billion: A large-scale benchmark for general object grasping," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 444–11 453.

[11] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[12] J. Kerr, L. Fu, H. Huang, Y. Avigal, M. Tancik, J. Ichnowski, A. Kanazawa, and K. Goldberg, "Evo-nerf: Evolving nerf for sequential robot grasping of transparent objects," in *Conference on Robot Learning*. PMLR, 2023, pp. 353–367.

[13] J. Ichnowski, Y. Avigal, J. Kerr, and K. Goldberg, "Dex-nerf: Using a neural radiance field to grasp transparent objects," *arXiv preprint arXiv:2110.14217*, 2021.

[14] Q. Dai, Y. Zhu, Y. Geng, C. Ruan, J. Zhang, and H. Wang, "Graspnerf: Multiview-based 6-dof grasp detection for transparent and specular objects using generalizable nerf," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 1757–1763.

[15] L. Yen-Chen, P. Florence, J. T. Barron, T.-Y. Lin, A. Rodriguez, and P. Isola, "Nerf-supervision: Learning dense object descriptors from neural radiance fields," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 6496–6503.

[16] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison, "imap: Implicit mapping and positioning in real-time," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6229–6238.

[17] Z. Zhu, S. Peng, V. Larsson, W. Xu, H. Bao, Z. Cui, M. R. Oswald, and M. Pollefeys, "Nice-slam: Neural implicit scalable encoding for slam," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 786–12 796.

[18] A. Rosinol, J. J. Leonard, and L. Carlone, "Nerf-slam: Real-time dense monocular slam with neural radiance fields," in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2023, pp. 3437–3444.

[19] M. Adamkiewicz, T. Chen, A. Caccavale, R. Gardner, P. Culbertson, J. Bohg, and M. Schwager, "Vision-only robot navigation in a neural radiance world," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 4606–4613, 2022.

[20] 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," *arXiv preprint arXiv:2312.02126*, 2023.

[21] 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.

[22] C. Yan, D. Qu, D. Wang, D. Xu, Z. Wang, B. Zhao, and X. Li, "Gs-slam: Dense visual slam with 3d gaussian splatting," *arXiv preprint arXiv:2311.11700*, 2023.

[23] G. Lu, S. Zhang, Z. Wang, C. Liu, J. Lu, and Y. Tang, "Manigaussian: Dynamic gaussian splatting for multi-task robotic manipulation," *arXiv preprint arXiv:2403.08321*, 2024.

[24] S. Guadarrama, E. Rodner, K. Saenko, N. Zhang, R. Farrell, J. Donahue, and T. Darrell, "Open-vocabulary object retrieval." in *Robotics: science and systems*, vol. 2, no. 5, 2014, p. 6.

[25] J. Hatori, Y. Kikuchi, S. Kobayashi, K. Takahashi, Y. Tsuboi, Y. Unno, W. Ko, and J. Tan, "Interactively picking real-world objects with unconstrained spoken language instructions," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 3774–3781.

[26] M. Shridhar and D. Hsu, "Interactive visual grounding of referring expressions for human-robot interaction," *arXiv preprint arXiv:1806.03831*, 2018.

[27] T. Nguyen, N. Gopalan, R. Patel, M. Corsaro, E. Pavlick, and S. Tellex, "Robot object retrieval with contextual natural language queries," *arXiv preprint arXiv:2006.13253*, 2020.

[28] Y. Chen, R. Xu, Y. Lin, and P. A. Vela, "A joint network for grasp detection conditioned on natural language commands," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 4576–4582.

[29] C. Cheang, H. Lin, Y. Fu, and X. Xue, "Learning 6-dof object poses to grasp category-level objects by language instructions," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 8476–8482.

[30] Q. Sun, H. Lin, Y. Fu, Y. Fu, and X. Xue, "Language guided robotic grasping with fine-grained instructions," in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2023, pp. 1319–1326.

[31] H. Lin, Z. Liu, C. Cheang, Y. Fu, G. Guo, and X. Xue, "Sar-net: Shape alignment and recovery network for category-level 6d object pose and size estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 6707–6717.

[32] M. Shridhar, L. Manuelli, and D. Fox, "Cliport: What and where pathways for robotic manipulation," in *Conference on Robot Learning*. PMLR, 2022, pp. 894–906.

[33] J. L. Schonberger and J.-M. Frahm, "Structure-from-motion revisited," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4104–4113.

[34] M. Zollhöfer, P. Stotko, A. Görlitz, C. Theobalt, M. Nießner, R. Klein, and A. Kolb, "State of the art on 3d reconstruction with rgb-d cameras," in *Computer graphics forum*, vol. 37, no. 2. Wiley Online Library, 2018, pp. 625–652.

[35] K. Wang, G. Zhang, and H. Bao, "Robust 3d reconstruction with an rgb-d camera," *IEEE Transactions on Image Processing*, vol. 23, no. 11, pp. 4893–4906, 2014.

[36] J. Kerr, C. M. Kim, K. Goldberg, A. Kanazawa, and M. Tancik, "Lerf: Language embedded radiance fields," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 19 729–19 739.

[37] X. Zhao, W. Ding, Y. An, Y. Du, T. Yu, M. Li, M. Tang, and J. Wang, "Fast segment anything," *arXiv preprint arXiv:2306.12156*, 2023.

[38] S. Zhou, H. Chang, S. Jiang, Z. Fan, Z. Zhu, D. Xu, P. Chari, S. You, Z. Wang, and A. Kadambi, "Feature 3dgs: Supercharging 3d gaussian splatting to enable distilled feature fields," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.

[39] M. Qin, W. Li, J. Zhou, H. Wang, and H. Pfister, "Langsplat: 3d language gaussian splatting," *arXiv preprint arXiv:2312.16084*, 2023.

[40] Z. Zivkovic, "Improved adaptive gaussian mixture model for background subtraction," in *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, vol. 2. IEEE, 2004, pp. 28–31.

[41] M. Ye, M. Danelljan, F. Yu, and L. Ke, "Gaussian grouping: Segment and edit anything in 3d scenes," *arXiv preprint arXiv:2312.00732*, 2023.

[42] H. Bharadhwaj, R. Mottaghi, A. Gupta, and S. Tulsiani, "Track2act: Predicting point tracks from internet videos enables diverse zero-shot robot manipulation," *arXiv preprint arXiv:2405.01527*, 2024.

[43] Y. Li and D. Pathak, "Object-aware gaussian splatting for robotic manipulation," in *ICRA 2024 Workshop on 3D Visual Representations for Robot Manipulation*, 2024. [Online]. Available: https://openreview.net/forum?id=gdRI43hDgo