

Semantic & Neural Rendering & SLAM

Research Notes & Literature Review

Shuqi XIAO

July 9, 2024

1 Semantic 3DGS

- Overview
- Feature-3DGS
- LangSplat
- CLIP-GS

2 3DGS SLAM

- Overview

Semantic 3DGS

Semantic 3DGS

Overview

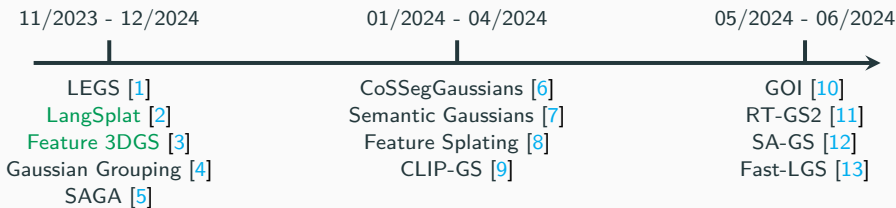


Figure 1: A timeline of Semantic 3DGS papers

Consensus

- What do we care about?
 - Accuracy; Consistency¹; Efficiency; Interactivity².
- How can we achieve it?

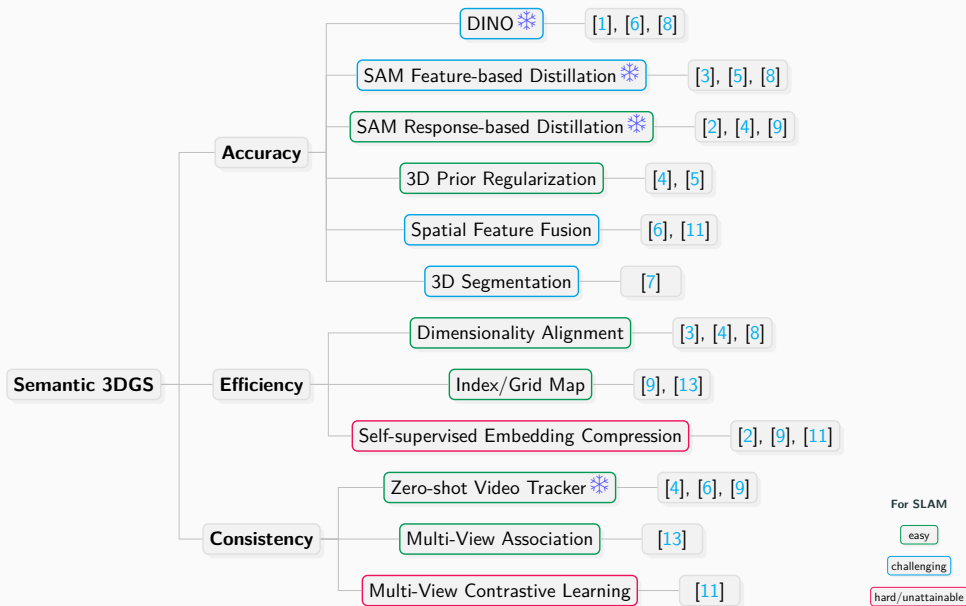
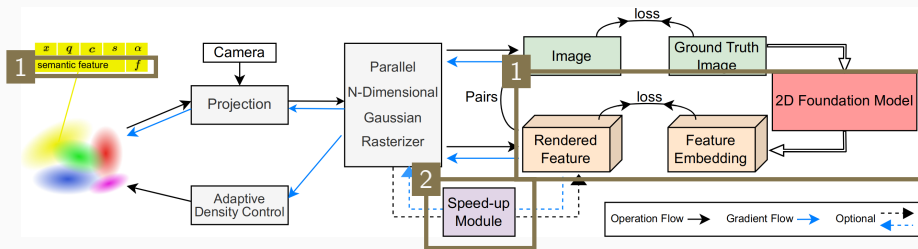


Figure 2: A taxonomy of Semantic 3DGS

Semantic 3DGS

Feature-3DGS



To render semantic embeddings, i.e.

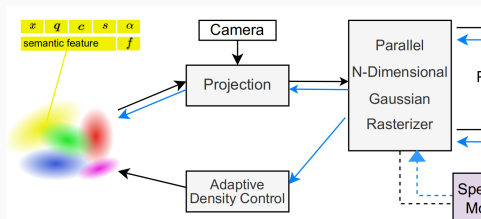
$$\mathbb{R}^D \times \{0, 1, \dots, N\} \times \{\mathcal{F}\} \mapsto \{0, 1, \dots, H\} \times \{0, 1, \dots, W\} \times \mathbb{R}^D \quad (1)$$

Annotations for Equation (1):

- D : dimension of semantic embedding
- N : number of 3D Gaussians
- \mathcal{F} : camera frames
- H : image height
- W : image width

5 things,

- 1 representation
- 2 projection
- 3 blending
- 4 rasterization
- 5 inverse rendering



1. Representation: 3D Gaussian augmented with a latent embedding.

$$\mathcal{G}_i = \{ \mathbf{x}, \mathbf{q}, \mathbf{s}, \alpha, \mathbf{c}, \mathbf{f} \} \quad (2)$$

Diagram illustrating the representation of a 3D Gaussian augmented with a latent embedding. The components are defined as follows:

- \mathcal{G} : Optimizable attributes of a 3D Gaussian (indicated by a red arrow).
- i : Index of 3D Gaussian, $i \in \mathbb{N}$ (indicated by a brown arrow).
- \mathbf{x} : Position, $\mathbf{x} \in \mathbb{R}^3$ (indicated by a green arrow).
- \mathbf{q} : Rotation, $\mathbf{q} \in \text{SO}(3)$ (indicated by a green arrow).
- \mathbf{s} : Scale, $\mathbf{s} \in \mathbb{R}^3$ (indicated by a green arrow).
- α : Opacity, $\alpha \in [0, 1]$ (indicated by a blue arrow).
- \mathbf{c} : Color, $\mathbf{c} \in \mathbb{R}^{3n}$ (indicated by a blue arrow).
- \mathbf{f} : Semantic feature, $\mathbf{f} \in \mathbb{R}^3$ (indicated by a purple arrow).

n : the maximal order of spherical harmonics to represent a color channel. In practice, $n = 4$.

(CVPR Highlight, 2024) Feature 3DGS: Supercharging 3D Gaussian Splatting to Enable Distilled Feature Fields

2. Projection: from 3D ellipsoids to 2D ellipses.

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (3)$$

Diagram illustrating the projection of a 3D world mean μ_w (red box) into a 2D image mean μ_i (purple box) using the camera pose \mathbf{T}_{cw} (green box). The projection function π (cyan box) maps the 3D point to 2D. The camera pose \mathbf{T}_{cw} is an element of $\text{SE}(3)$. The world mean μ_w is an element of \mathbb{P}^3 , 3D(world) mean. The image mean μ_i is an element of \mathbb{P}^2 , 2D(image) mean. The projection function π is labeled "projection".

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (4)$$

Diagram illustrating the projection of a 3D world covariance Σ_w (red box) into a 2D image covariance Σ_i (purple box) using the camera pose \mathbf{T}_{cw} (green box). The Jacobian of the linear approximation of π , \mathbf{J}_π (cyan box), is used to map the 3D covariance to 2D. The rotation component of \mathbf{T}_{cw} , \mathbf{R}_{cw} (green box), is an element of $\text{SO}(3)$. The world covariance Σ_w is an element of $\mathbb{R}^{3 \times 3}$, 3D(world) covariance. The image covariance Σ_i is an element of $\mathbb{R}^{2 \times 2}$, 2D(image) covariance. The Jacobian \mathbf{J}_π is an element of $\mathbb{R}^{2 \times 3}$, Jacobian of the linear approximation of π .

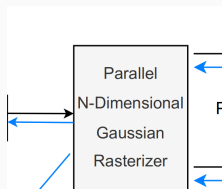
3. Blending: α -blending of semantic embeddings.

$$\mathbf{f}(h, w) = \sum_{i=1}^N T_i \alpha_i \mathbf{f}_i(h, w), \quad T_i = \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (5)$$

Diagram illustrating the blending process:

- N : number of the sorted & visible subset of 3D Gaussians
- $\mathbf{f}(h, w)$: semantic feature on pixel (h, w)
- T_i : background opacity for i -th Gaussian
- α_i : opacity of i -th Gaussian
- $\mathbf{f}_i(h, w)$: semantic feature of i -th Gaussian on pixel (h, w)

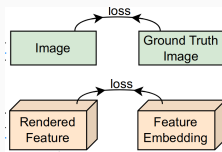
4. **Rasterization**: tiled and implemented in CUDA.



- Divide the screen space into tiles (CUDA thread blocks).
- Group the Gaussians by view frustum and tile index.
- Sort the Gaussians by front-to-back depth order.
- Blend each pixel within a tile in parallel (CUDA threads).

Figure 4: A brief summary of the tiled rasterization in 3DGS

5. **Inverse rendering**: guided by image-wise photometric loss.



$$\mathcal{L} = \mathcal{L}_{appearance} + \gamma \mathcal{L}_{semantics} \quad (6)$$

captured RGB image (GT) rendered RGB image

$$\mathcal{L}_{appearance} = (1 - \lambda) \mathcal{L}_1 \left(\text{C}, \hat{\text{C}} \right) + \lambda \mathcal{L}_{D-SSIM} \left(\text{C}, \hat{\text{C}} \right) \quad (7)$$

$$\mathcal{L}_{semantics} = \mathcal{L}_1 \left(\text{F}, \hat{\text{F}} \right) = \sum_{h=1}^H \sum_{w=1}^W \| \mathbf{f}(h, w) - \hat{\mathbf{f}}(h, w) \|_1 \quad (8)$$

inferred semantic feature map rendered semantic feature map

Motivation

Too **inefficient** to embed naively,

- 1 **High dimension:** latent features in large foundation models.
- 2 **Large quantities:** millions of Gaussians in a scene.

Solution

- 1 **Compactness:** to embed Gaussians with more compact vectors, $\dim = D' < D$.
- 2 **Alignment:** to align the dimensionalities using a lightweight decoder.

$D = 512$ in CLIP; $D = 256$ in SAM.

In practice, $D' = 128$.

Lightweight decoder: In practice, a 1×1 convolutional layer or a fully-connected layer.

(CVPR Highlight, 2024) [Feature 3DGS: Supercharging 3D Gaussian Splatting to Enable Distilled Feature Fields](#)

Limitations

1 Inefficiency

- “Speed-up module” is not enough,
dim = 128 embedding for millions of Gaussians.

2 3D Inconsistency & Inaccuracy

- 2D foundation models are still 2D.

Semantic 3DGS

LangSplat

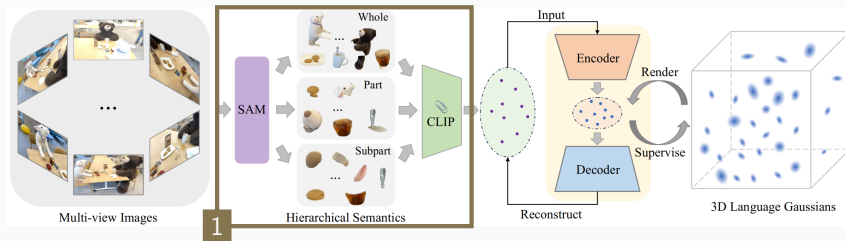


Figure 5: Overview of LangSplat

- 1 **Accuracy:** SAM outputs to enhance CLIP features.
 - **CLIP:** image-aligned training leads to “point-ambiguity”.
 - **SAM:** pixel-aligned & object-centered & multi-granularity.

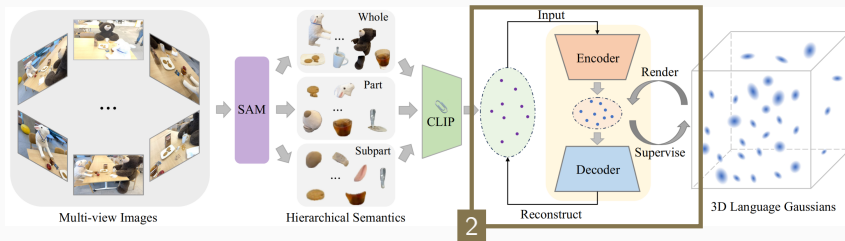


Figure 5: Overview of LangSplat

2 Efficiency: an **auto-encoder** to compress latent features.

- More complexity and better compression, compared with “speed-up module” in Feature 3DGS [3].

(CVPR Highlight, 2024) [LangSplat: 3D Language Gaussian Splatting](#)

In practice, the latent dimension is 3 in the auto-encoder.

Semantic 3DGS

CLIP-GS

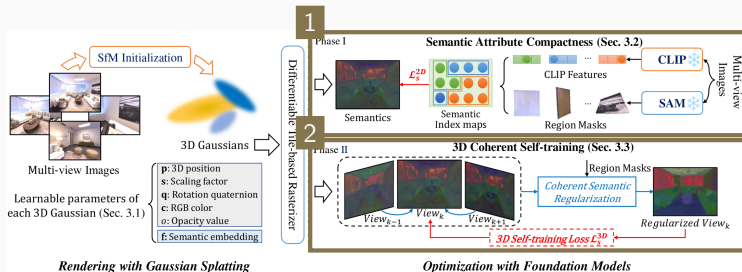


Figure 6: Overview of CLIP-GS

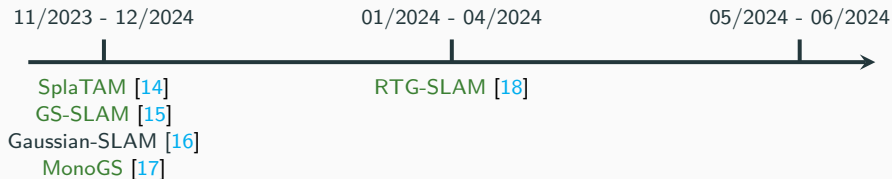
Key Insights

- Efficiency:** unify semantic features within an object by leveraging SAM.
- Consistency:** supervise consecutive frames by video segmentation.

3DGS SLAM

3DGS SLAM

Overview



References

- [1] J.-C. Shi, M. Wang, H.-B. Duan, and S.-H. Guan, *Language embedded 3d gaussians for open-vocabulary scene understanding*, Nov. 30, 2023. arXiv: [2311.18482\[cs\]](https://arxiv.org/abs/2311.18482). [Online]. Available: <http://arxiv.org/abs/2311.18482> (visited on 06/08/2024) (cit. on pp. v, vi).
- [2] M. Qin, W. Li, J. Zhou, H. Wang, and H. Pfister, *LangSplat: 3d language gaussian splatting*, Dec. 26, 2023. arXiv: [2312.16084\[cs\]](https://arxiv.org/abs/2312.16084). [Online]. Available: <http://arxiv.org/abs/2312.16084> (visited on 02/23/2024) (cit. on pp. v, vi).
- [3] S. Zhou, H. Chang, S. Jiang, et al., *Feature 3dgs: Supercharging 3d gaussian splatting to enable distilled feature fields*, Apr. 8, 2024. arXiv: [2312.03203\[cs\]](https://arxiv.org/abs/2312.03203). [Online]. Available: <http://arxiv.org/abs/2312.03203> (visited on 05/22/2024) (cit. on pp. v, vi, xix).
- [4] M. Ye, M. Danelljan, F. Yu, and L. Ke, *Gaussian grouping: Segment and edit anything in 3d scenes*, Dec. 1, 2023. arXiv: [2312.00732\[cs\]](https://arxiv.org/abs/2312.00732). [Online]. Available: <http://arxiv.org/abs/2312.00732> (visited on 01/02/2024) (cit. on pp. v, vi).
- [5] J. Cen, J. Fang, C. Yang, et al., *Segment any 3d gaussians*, Dec. 1, 2023. arXiv: [2312.00860\[cs\]](https://arxiv.org/abs/2312.00860). [Online]. Available: <http://arxiv.org/abs/2312.00860> (visited on 03/12/2024) (cit. on pp. v, vi).
- [6] B. Dou, T. Zhang, Y. Ma, Z. Wang, and Z. Yuan, *CoSSegGaussians: Compact and swift scene segmenting 3d gaussians with dual feature fusion*, Jan. 30, 2024. arXiv: [2401.05925\[cs\]](https://arxiv.org/abs/2401.05925). [Online]. Available: <http://arxiv.org/abs/2401.05925> (visited on 06/08/2024) (cit. on pp. v, vi).
- [7] J. Guo, X. Ma, Y. Fan, H. Liu, and Q. Li, *Semantic gaussians: Open-vocabulary scene understanding with 3d gaussian splatting*, Mar. 22, 2024. arXiv: [2403.15624\[cs\]](https://arxiv.org/abs/2403.15624). [Online]. Available: <http://arxiv.org/abs/2403.15624> (visited on 05/20/2024) (cit. on pp. v, vi).

- [8] R.-Z. Qiu, G. Yang, W. Zeng, and X. Wang, *Feature splatting: Language-driven physics-based scene synthesis and editing*, Apr. 1, 2024. arXiv: [2404.01223\[cs\]](https://arxiv.org/abs/2404.01223). [Online]. Available: <http://arxiv.org/abs/2404.01223> (visited on 06/08/2024) (cit. on pp. v, vi).
- [9] G. Liao, J. Li, Z. Bao, et al., *CLIP-GS: CLIP-informed gaussian splatting for real-time and view-consistent 3d semantic understanding*, Apr. 22, 2024. arXiv: [2404.14249\[cs\]](https://arxiv.org/abs/2404.14249). [Online]. Available: <http://arxiv.org/abs/2404.14249> (visited on 05/20/2024) (cit. on pp. v, vi).
- [10] Y. Qu, S. Dai, X. Li, et al., *GOI: Find 3d gaussians of interest with an optimizable open-vocabulary semantic-space hyperplane*, May 27, 2024. arXiv: [2405.17596\[cs\]](https://arxiv.org/abs/2405.17596). [Online]. Available: <http://arxiv.org/abs/2405.17596> (visited on 06/08/2024) (cit. on p. v).
- [11] M.-B. Jurca, R. Royen, I. Giosan, and A. Munteanu, *RT-GS2: Real-time generalizable semantic segmentation for 3d gaussian representations of radiance fields*, May 28, 2024. arXiv: [2405.18033\[cs\]](https://arxiv.org/abs/2405.18033). [Online]. Available: <http://arxiv.org/abs/2405.18033> (visited on 06/08/2024) (cit. on pp. v, vi).
- [12] B. Xiong, X. Ye, T. H. E. Tse, K. Han, S. Cui, and Z. Li, *SA-GS: Semantic-aware gaussian splatting for large scene reconstruction with geometry constrain*, May 28, 2024. arXiv: [2405.16923\[cs\]](https://arxiv.org/abs/2405.16923). [Online]. Available: <http://arxiv.org/abs/2405.16923> (visited on 06/08/2024) (cit. on p. v).
- [13] Y. Ji, H. Zhu, J. Tang, et al., *FastLGS: Speeding up language embedded gaussians with feature grid mapping*, Jun. 3, 2024. arXiv: [2406.01916\[cs\]](https://arxiv.org/abs/2406.01916). [Online]. Available: <http://arxiv.org/abs/2406.01916> (visited on 06/08/2024) (cit. on pp. v, vi).

- [14] N. Keetha, J. Karhade, K. M. Jatavallabhula, *et al.*, *SplaTAM: Splat, track & map 3d gaussians for dense RGB-d SLAM*, Apr. 16, 2024. arXiv: [2312.02126\[cs\]](https://arxiv.org/abs/2312.02126). [Online]. Available: <http://arxiv.org/abs/2312.02126> (visited on 05/20/2024) (cit. on p. xxiv).
- [15] C. Yan, D. Qu, D. Wang, *et al.*, *GS-SLAM: Dense visual SLAM with 3d gaussian splatting*, Nov. 21, 2023. arXiv: [2311.11700\[cs\]](https://arxiv.org/abs/2311.11700). [Online]. Available: <http://arxiv.org/abs/2311.11700> (visited on 12/26/2023) (cit. on p. xxiv).
- [16] V. Yugay, Y. Li, T. Gevers, and M. R. Oswald, *Gaussian-SLAM: Photo-realistic dense SLAM with gaussian splatting*, Mar. 22, 2024. arXiv: [2312.10070\[cs\]](https://arxiv.org/abs/2312.10070). [Online]. Available: <http://arxiv.org/abs/2312.10070> (visited on 03/27/2024) (cit. on p. xxiv).
- [17] H. Matsuki, R. Murai, P. H. J. Kelly, and A. J. Davison, *Gaussian splatting SLAM*, Apr. 14, 2024. arXiv: [2312.06741\[cs\]](https://arxiv.org/abs/2312.06741). [Online]. Available: <http://arxiv.org/abs/2312.06741> (visited on 05/20/2024) (cit. on p. xxiv).
- [18] Z. Peng, T. Shao, Y. Liu, *et al.*, *RTG-SLAM: Real-time 3d reconstruction at scale using gaussian splatting*, May 8, 2024. DOI: [10.1145/3658233](https://doi.org/10.1145/3658233). arXiv: [2404.19706\[cs\]](https://arxiv.org/abs/2404.19706). [Online]. Available: <http://arxiv.org/abs/2404.19706> (visited on 05/18/2024) (cit. on p. xxiv).