

Semantic & Neural Rendering & SLAM

Research Notes & Literature Review

Shuqi XIAO

July 3, 2024

1 3DGS-based SLAM

- MonoGS

2 Semantic 3DGS

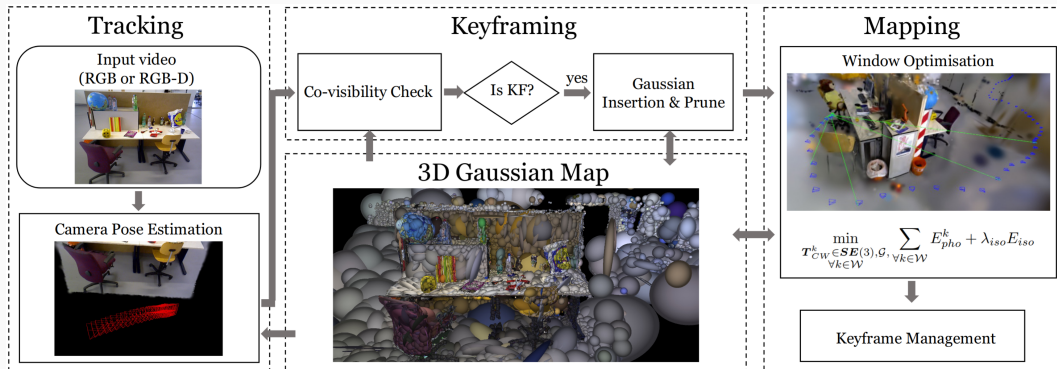
- Feature-3DGS
- LangSplat
- CLIP-GS

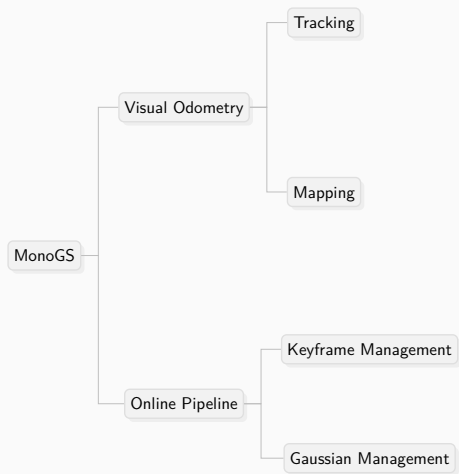
3DGS-based SLAM

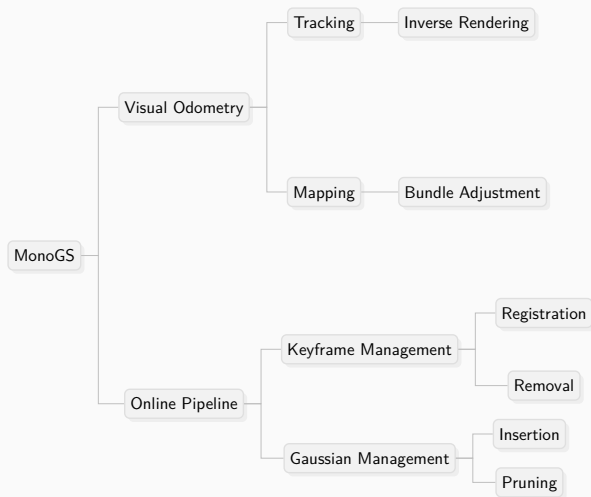


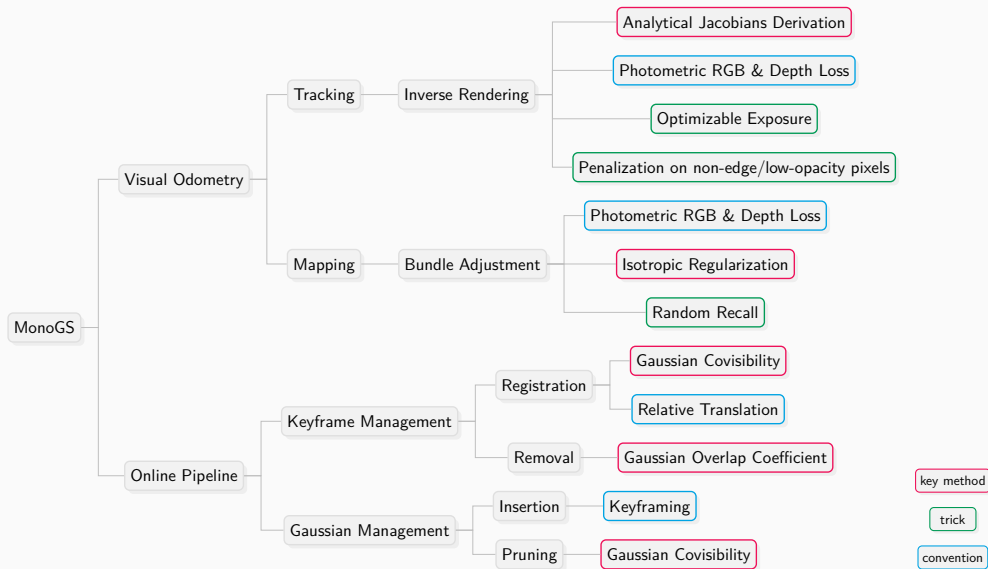
3DGS-based SLAM

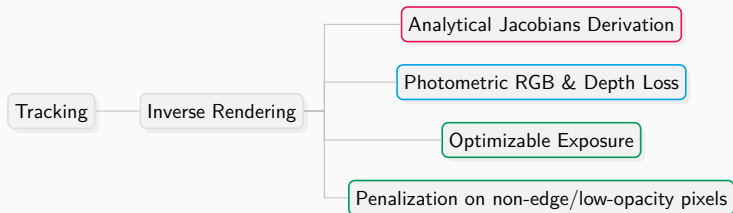
MonoGS





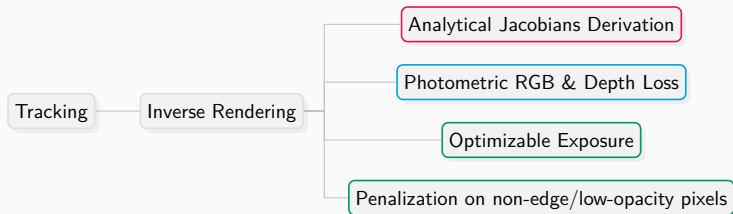






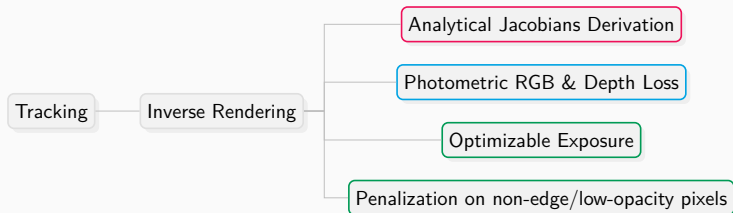
Track camera poses by inverse rendering,

- through the extended differentiable rendering pipeline,
- by a direct optimization against fixed 3D Gaussians,
- with some tricks to be more adaptive to brightness and more robust to noise.



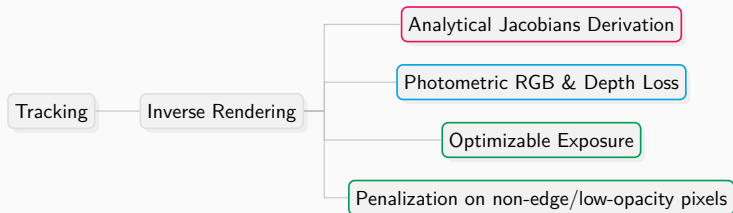
Track camera poses by inverse rendering,

- through the extended differentiable rendering pipeline,
- by a direct optimization against fixed 3D Gaussians,
- with some tricks to be more adaptive to brightness and more robust to noise.



Track camera poses by inverse rendering,

- through the extended differentiable rendering pipeline,
- by a direct optimization against fixed 3D Gaussians,
- with some tricks to be more adaptive to brightness and more robust to noise.



Track camera poses by inverse rendering,

- through the extended differentiable rendering pipeline,
- by a direct optimization against fixed 3D Gaussians,
- with some tricks to be more adaptive to brightness and more robust to noise.

Firstly, let's review the **projection** of 3D Gaussians.

Firstly, let's review the projection of 3D Gaussians.

$$\mathcal{N}(\mu_w, \Sigma_w) \xrightarrow{\pi} \mathcal{N}(\mu_i, \Sigma_i) \quad (1)$$

Firstly, let's review the projection of 3D Gaussians.

$$\mathcal{N}(\mu_w, \Sigma_w) \xrightarrow{\pi} \mathcal{N}(\mu_i, \Sigma_i) \quad (1)$$

is achieved by

$$\mu_i = \pi(\mathbf{T}_{cw} \cdot \mu_w) \quad (2)$$

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^\top \mathbf{J}_\pi^\top \quad (3)$$

Firstly, let's review the projection of 3D Gaussians.

$$\mathcal{N}(\mu_w, \Sigma_w) \xrightarrow{\pi} \mathcal{N}(\mu_i, \Sigma_i) \quad (1)$$

is achieved by

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (2)$$

$\in \mathbb{P}^3$, 3D(world) mean

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (3)$$

Firstly, let's review the projection of 3D Gaussians.

$$\mathcal{N}(\mu_w, \Sigma_w) \xrightarrow{\pi} \mathcal{N}(\mu_i, \Sigma_i) \quad (1)$$

is achieved by

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (2)$$

$\in \mathbb{P}^3$, 3D(world) mean

$\in \text{SE}(3)$, camera pose

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (3)$$

Firstly, let's review the projection of 3D Gaussians.

$$\mathcal{N}(\mu_w, \Sigma_w) \xrightarrow{\pi} \mathcal{N}(\mu_i, \Sigma_i) \quad (1)$$

is achieved by

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (2)$$

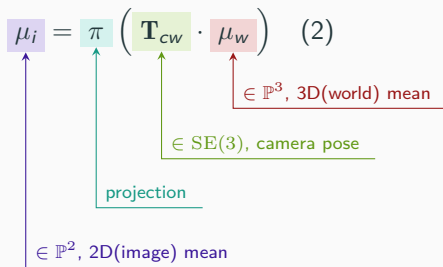
Diagram illustrating the projection of the 3D Gaussian mean μ_w into the camera space to obtain μ_i . The projection is performed by the operator π applied to the transformed mean $\mathbf{T}_{cw} \cdot \mu_w$. The transformation \mathbf{T}_{cw} represents the camera pose, belonging to $\text{SE}(3)$. The world mean μ_w is a point in \mathbb{P}^3 , representing the 3D world mean.

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (3)$$

Firstly, let's review the projection of 3D Gaussians.

$$\mathcal{N}(\mu_w, \Sigma_w) \xrightarrow{\pi} \mathcal{N}(\mu_i, \Sigma_i) \quad (1)$$

is achieved by

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (2)$$


$\mu_w \in \mathbb{P}^3, 3D(\text{world}) \text{ mean}$

$\mathbf{T}_{cw} \in SE(3), \text{ camera pose}$

π projection

$\mu_i \in \mathbb{P}^2, 2D(\text{image}) \text{ mean}$

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (3)$$

Firstly, let's review the projection of 3D Gaussians.

$$\mathcal{N}(\mu_w, \Sigma_w) \xrightarrow{\pi} \mathcal{N}(\mu_i, \Sigma_i) \quad (1)$$

is achieved by

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (2)$$

$\mu_i \in \mathbb{P}^2, 2D(\text{image}) \text{ mean}$

π projection

$\mathbf{T}_{cw} \in \text{SE}(3), \text{ camera pose}$

$\mu_w \in \mathbb{P}^3, 3D(\text{world}) \text{ mean}$

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (3)$$

$\Sigma_w \in \mathbb{R}^{3 \times 3}, 3D(\text{world}) \text{ covariance}$

Firstly, let's review the projection of 3D Gaussians.

$$\mathcal{N}(\mu_w, \Sigma_w) \xrightarrow{\pi} \mathcal{N}(\mu_i, \Sigma_i) \quad (1)$$

is achieved by

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (2)$$

$\mu_i \in \mathbb{P}^2, 2D(\text{image}) \text{ mean}$
 π projection
 $\mathbf{T}_{cw} \in \text{SE}(3), \text{ camera pose}$
 $\mu_w \in \mathbb{P}^3, 3D(\text{world}) \text{ mean}$

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (3)$$

$\Sigma_i \in \mathbb{R}^{2 \times 2}, 2D(\text{image}) \text{ covariance}$
 \mathbf{J}_π Jacobian of π
 $\mathbf{R}_{cw} \in \text{SO}(3), \text{ rotation component of } \mathbf{T}_{cw}$
 $\Sigma_w \in \mathbb{R}^{3 \times 3}, 3D(\text{world}) \text{ covariance}$

Firstly, let's review the projection of 3D Gaussians.

$$\mathcal{N}(\mu_w, \Sigma_w) \xrightarrow{\pi} \mathcal{N}(\mu_i, \Sigma_i) \quad (1)$$

is achieved by

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (2)$$

$\mu_i \in \mathbb{P}^2, 2D(\text{image}) \text{ mean}$
 π projection
 $\mathbf{T}_{cw} \in \text{SE}(3), \text{ camera pose}$
 $\mu_w \in \mathbb{P}^3, 3D(\text{world}) \text{ mean}$

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (3)$$

$\Sigma_i \in \mathbb{R}^{2 \times 3}, \text{ Jacobian of the linear approximation of } \pi$
 $\mathbf{J}_\pi \in \mathbb{R}^{2 \times 3}, \text{ Jacobian of the linear approximation of } \pi$
 $\mathbf{R}_{cw} \in \text{SO}(3), \text{ rotation component of } \mathbf{T}_{cw}$
 $\Sigma_w \in \mathbb{R}^{3 \times 3}, 3D(\text{world}) \text{ covariance}$

Firstly, let's review the projection of 3D Gaussians.

$$\mathcal{N}(\mu_w, \Sigma_w) \xrightarrow{\pi} \mathcal{N}(\mu_i, \Sigma_i) \quad (1)$$

is achieved by

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (2)$$

Diagram illustrating the projection of the 3D world mean μ_w to the 2D image mean μ_i . The components are:

- $\mu_i \in \mathbb{P}^2$, 2D(image) mean (purple box)
- π (cyan box) is the projection function.
- $\mathbf{T}_{cw} \in \text{SE}(3)$, camera pose (green box).
- $\mu_w \in \mathbb{P}^3$, 3D(world) mean (red box).

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (3)$$

Diagram illustrating the derivation of the 2D image covariance Σ_i . The components are:

- $\Sigma_i \in \mathbb{R}^{2 \times 2}$, 2D(image) covariance (purple box)
- $\mathbf{J}_\pi \in \mathbb{R}^{2 \times 3}$, Jacobian of the linear approximation of π (cyan box).
- $\mathbf{R}_{cw} \in \text{SO}(3)$, rotation component of \mathbf{T}_{cw} (green box).
- $\Sigma_w \in \mathbb{R}^{3 \times 3}$, 3D(world) covariance (red box).

The chain rule,

$$\frac{\partial \mu_i}{\partial \mathbf{T}_{cw}} = \frac{\partial \mu_i}{\partial \mu_c} \frac{\partial \mu_c}{\partial \mathbf{T}_{cw}} \quad (4)$$

$$\frac{\partial \Sigma_i}{\partial \mathbf{T}_{cw}} = \frac{\partial \Sigma_i}{\partial \mathbf{J}_\pi} \frac{\partial \mathbf{J}_\pi}{\partial \mu_c} \frac{\partial \mu_c}{\partial \mathbf{T}_{cw}} + \frac{\partial \Sigma_i}{\partial \mathbf{R}_{cw}} \frac{\partial \mathbf{R}_{cw}}{\partial \mathbf{T}_{cw}} \quad (5)$$

The Lie Algebra,

$$\frac{\partial \mu_c}{\partial \mathbf{T}_{cw}} = [\mathbf{I} \quad -\mu_c^\times] \quad (6)$$

$$\frac{\partial \mathbf{R}_{cw}}{\partial \mathbf{T}_{cw}} = \begin{bmatrix} \mathbf{0} & -\mathbf{R}_{cw}^\times(:, 1) \\ \mathbf{0} & -\mathbf{R}_{cw}^\times(:, 2) \\ \mathbf{0} & -\mathbf{R}_{cw}^\times(:, 3) \end{bmatrix} \quad (7)$$

The chain rule,

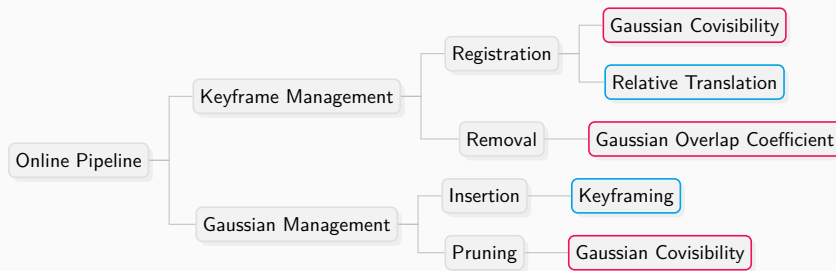
$$\frac{\partial \mu_i}{\partial \mathbf{T}_{cw}} = \frac{\partial \mu_i}{\partial \mu_c} \frac{\partial \mu_c}{\partial \mathbf{T}_{cw}} \quad (4)$$

$$\frac{\partial \Sigma_i}{\partial \mathbf{T}_{cw}} = \frac{\partial \Sigma_i}{\partial \mathbf{J}_\pi} \frac{\partial \mathbf{J}_\pi}{\partial \mu_c} \frac{\partial \mu_c}{\partial \mathbf{T}_{cw}} + \frac{\partial \Sigma_i}{\partial \mathbf{R}_{cw}} \frac{\partial \mathbf{R}_{cw}}{\partial \mathbf{T}_{cw}} \quad (5)$$

The Lie Algebra,

$$\frac{\partial \mu_c}{\partial \mathbf{T}_{cw}} = [\mathbf{I} \quad -\mu_c^\times] \quad (6)$$

$$\frac{\partial \mathbf{R}_{cw}}{\partial \mathbf{T}_{cw}} = \begin{bmatrix} \mathbf{0} & -\mathbf{R}_{cw}^\times(:, 1) \\ \mathbf{0} & -\mathbf{R}_{cw}^\times(:, 2) \\ \mathbf{0} & -\mathbf{R}_{cw}^\times(:, 3) \end{bmatrix} \quad (7)$$



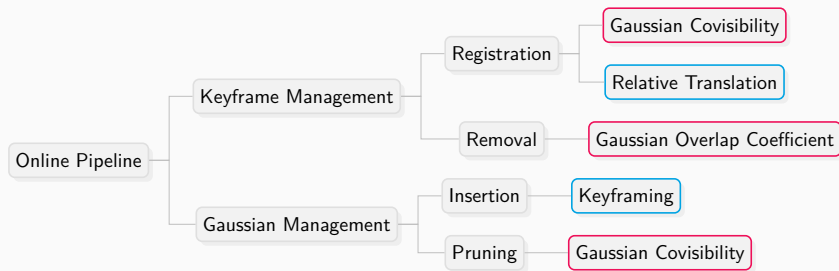
Keyframe Management:

- Classic strategies, e.g. covisibility & overlap, from DSO [5].
- Off-the-shelf occlusion-aware Gaussian visibility is leveraged to construct metrics.

key method trick convention

(arXiv, 2016) DSO: Direct Sparse Odometry

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM



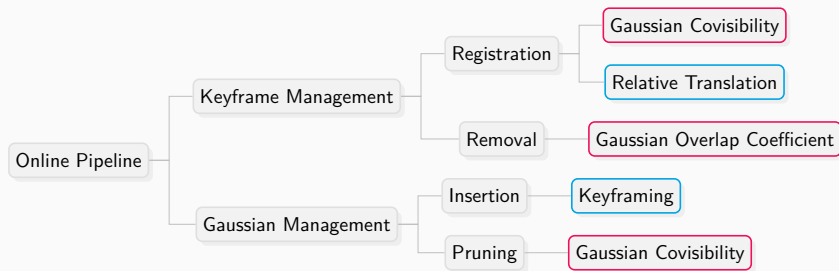
Keyframe Management:

- **Classic** strategies, e.g. covisibility & overlap, from DSO [5].
- Off-the-shelf occlusion-aware Gaussian visibility is leveraged to construct metrics.

key method trick convention

(arXiv, 2016) DSO: Direct Sparse Odometry

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM



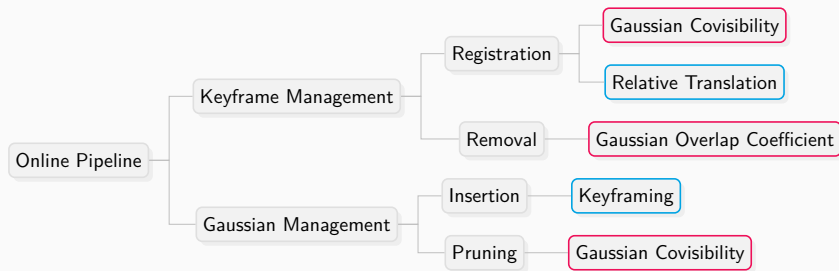
Keyframe Management:

- Classic strategies, e.g. covisibility & overlap, from DSO [5].
- **Off-the-shelf** occlusion-aware Gaussian visibility is leveraged to construct metrics.

key method trick convention

(arXiv, 2016) DSO: Direct Sparse Odometry

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM



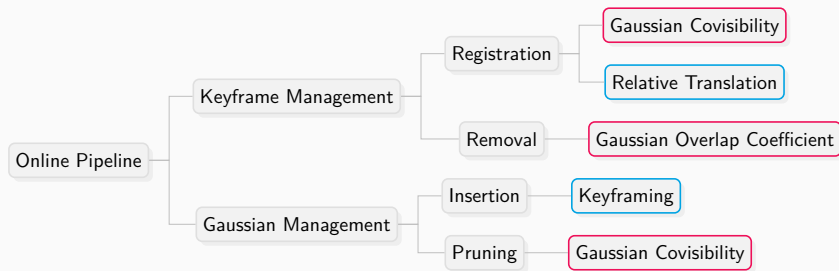
Gaussian Management:

- Insertion: triggered by keyframing, followed by Gaussian initialization.
- Pruning: to remove unstable/incorrect Gaussians by covisibility in a monocular setting.

key method trick convention

(arXiv, 2016) DSO: Direct Sparse Odometry

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM



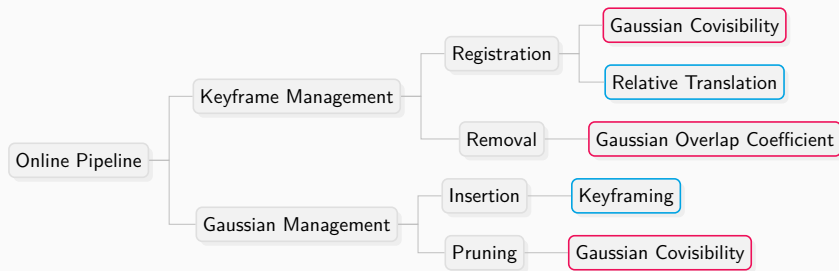
Gaussian Management:

- **Insertion:** triggered by **keyframing**, followed by **Gaussian initialization**.
- **Pruning:** to remove unstable/incorrect Gaussians by covisibility in a monocular setting.

key method trick convention

(arXiv, 2016) DSO: Direct Sparse Odometry

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM



Gaussian Management:

- **Insertion:** triggered by keyframing, followed by Gaussian initialization.
- **Pruning:** to remove unstable/incorrect Gaussians by covisibility in a monocular setting.

key method trick convention

(arXiv, 2016) DSO: Direct Sparse Odometry

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

1 What is keyframing or keyframe management?

- A strategy of selecting and utilizing a crucial subset of frames.

2 Why do we need keyframing?

- Infeasible to optimize jointly on all frames online.

3 How should we select keyframes?

- non-redundant and observing the same area.
- spanning a wide baseline for better multi-view constraints.

1 What is keyframing or keyframe management?

- A strategy of **selecting** and **utilizing** a **crucial subset** of frames.

2 Why do we need keyframing?

- Infeasible to optimize jointly on all frames online.

3 How should we select keyframes?

- non-redundant and observing the same area.
- spanning a wide baseline for better multi-view constraints.

1 What is keyframing or keyframe management?

- A strategy of selecting and utilizing a crucial subset of frames.

2 Why do we need keyframing?

- Infeasible to optimize jointly on all frames online.

3 How should we select keyframes?

- non-redundant and observing the same area.
- spanning a wide baseline for better multi-view constraints.

1 What is keyframing or keyframe management?

- A strategy of selecting and utilizing a crucial subset of frames.

2 Why do we need keyframing?

- **Infeasible** to optimize jointly on all frames online.

3 How should we select keyframes?

- non-redundant and observing the same area.
- spanning a wide baseline for better multi-view constraints.

1 What is keyframing or keyframe management?

- A strategy of selecting and utilizing a crucial subset of frames.

2 Why do we need keyframing?

- Infeasible to optimize jointly on all frames online.

(a **trade-off** between efficiency and accuracy/robustness/...)

3 How should we select keyframes?

- non-redundant and observing the same area.
- spanning a wide baseline for better multi-view constraints.

1 What is keyframing or keyframe management?

- A strategy of selecting and utilizing a crucial subset of frames.

2 Why do we need keyframing?

- Infeasible to optimize jointly on all frames online.

(a trade-off between efficiency and accuracy/robustness/...)

3 How should we select keyframes?

- non-redundant and observing the same area.
- spanning a wide baseline for better multi-view constraints.

1 What is keyframing or keyframe management?

- A strategy of selecting and utilizing a crucial subset of frames.

2 Why do we need keyframing?

- Infeasible to optimize jointly on all frames online.

(a trade-off between efficiency and accuracy/robustness/...)

3 How should we select keyframes?

- **non-redundant** and observing the **same area**.
- spanning a wide baseline for better multi-view constraints.

- 1 What is keyframing or keyframe management?
 - A strategy of selecting and utilizing a crucial subset of frames.
- 2 Why do we need keyframing?
 - Infeasible to optimize jointly on all frames online.
(a trade-off between efficiency and accuracy/robustness/...)
- 3 How should we select keyframes?
 - non-redundant and observing the same area.
 - spanning a **wide baseline** for better multi-view constraints.

If **any** of the following conditions **is true**...

Small Gaussian Covisibility

Condition i, Keyframe Registration

Gaussian covisibility between the current frame and the previous keyframe drops below a threshold.

$$\frac{|\mathbf{v}(\mathcal{G}, \mathcal{F}_i) \cap \mathbf{v}(\mathcal{G}, \mathcal{F}_j)|}{|\mathbf{v}(\mathcal{G}, \mathcal{F}_i) \cup \mathbf{v}(\mathcal{G}, \mathcal{F}_j)|} < \tau_1 \quad (8)$$

In practice, $\tau_1 = 0.95$.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

If any of the following conditions is true...

Small Gaussian Covisibility

Condition i, Keyframe Registration

Gaussian covisibility between the current frame and the previous keyframe drops below a threshold.

$$\frac{|\mathbf{v}(\mathcal{G}, \mathcal{F}_i) \cap \mathbf{v}(\mathcal{G}, \mathcal{F}_j)|}{|\mathbf{v}(\mathcal{G}, \mathcal{F}_i) \cup \mathbf{v}(\mathcal{G}, \mathcal{F}_j)|} < \tau_1 \quad (8)$$

In practice, $\tau_1 = 0.95$.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

If any of the following conditions is true...

Small Gaussian Covisibility

Condition i, Keyframe Registration

Gaussian covisibility between the current frame and the previous keyframe drops below a threshold.

$$\frac{\left| \mathbf{v}(\mathcal{G}, \mathcal{F}_i) \cap \underbrace{\mathbf{v}(\mathcal{G}, \mathcal{F}_j)}_{\substack{\subset \mathcal{G}, \text{ visible Gaussians from frame } j}} \right|}{\left| \mathbf{v}(\mathcal{G}, \mathcal{F}_i) \cup \mathbf{v}(\mathcal{G}, \mathcal{F}_j) \right|} < \tau_1 \quad (8)$$

In practice, $\tau_1 = 0.95$.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

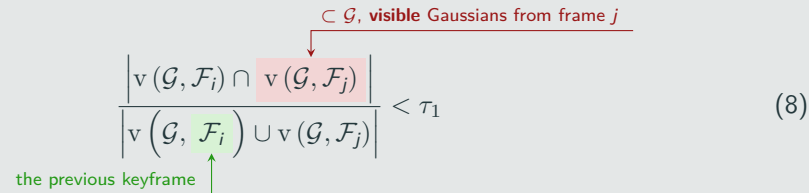
If any of the following conditions is true...

Small Gaussian Covisibility

Condition i, Keyframe Registration

Gaussian covisibility between the current frame and the previous keyframe drops below a threshold.

$$\frac{\left| \mathbf{v}(\mathcal{G}, \mathcal{F}_i) \cap \mathbf{v}(\mathcal{G}, \mathcal{F}_j) \right|}{\left| \mathbf{v}(\mathcal{G}, \mathcal{F}_i) \cup \mathbf{v}(\mathcal{G}, \mathcal{F}_j) \right|} < \tau_1 \quad (8)$$



In practice, $\tau_1 = 0.95$.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

If any of the following conditions is true...

Small Gaussian Covisibility

Condition i, Keyframe Registration

Gaussian covisibility between the current frame and the previous keyframe drops below a threshold.

$$\frac{\left| v(\mathcal{G}, \mathcal{F}_i) \cap v(\mathcal{G}, \mathcal{F}_j) \right|}{\left| v(\mathcal{G}, \mathcal{F}_i) \cup v(\mathcal{G}, \mathcal{F}_j) \right|} < \tau_1 \quad (8)$$

Diagram illustrating the condition for keyframe registration based on Gaussian covisibility. The equation shows the ratio of the number of visible Gaussians from the intersection of the current frame \mathcal{F}_i and the previous keyframe \mathcal{F}_j to the total number of visible Gaussians from both frames. The threshold τ_1 is set to 0.95 in practice.

Annotations:

- $\subset \mathcal{G}$, visible Gaussians from frame j (points to the intersection term in the numerator)
- the previous keyframe (points to \mathcal{F}_i in the denominator)
- the current frame (points to \mathcal{F}_j in the denominator)

In practice, $\tau_1 = 0.95$.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

Large Relative Translation

Condition ii, Keyframe Registration

Translation from the previous keyframe w.r.t. to the median depth reaches a threshold.

$$\frac{\|\mathbf{t}_{\mathcal{F}_i\mathcal{F}_j}\|_2}{\bar{D}_{\mathcal{F}_i\mathcal{F}_j}} > \tau_2, \quad \bar{D}_{\mathcal{F}_i\mathcal{F}_j} = \frac{1}{2HW} \sum_{\{\mathcal{F}_i, \mathcal{F}_j\}} \sum_{h=0}^H \sum_{w=0}^W d(h, w) \quad (9)$$

In practice, $\tau_2 = 0.04$. Additionally, evaluate the Gaussian covisibility only if the relative translation is not too small (> 0.02) for efficiency.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

Large Relative Translation

Condition ii, Keyframe Registration

Translation from the previous keyframe w.r.t. to the median depth reaches a threshold.

$$\frac{\|\mathbf{t}_{\mathcal{F}_i\mathcal{F}_j}\|_2}{\bar{D}_{\mathcal{F}_i\mathcal{F}_j}} > \tau_2, \quad \bar{D}_{\mathcal{F}_i\mathcal{F}_j} = \frac{1}{2HW} \sum_{\{\mathcal{F}_i, \mathcal{F}_j\}} \sum_{h=0}^H \sum_{w=0}^W d(h, w) \quad (9)$$

$\in \mathbb{R}^3$, translation from \mathcal{F}_i to \mathcal{F}_j

In practice, $\tau_2 = 0.04$. Additionally, evaluate the Gaussian covisibility only if the relative translation is not too small (> 0.02) for efficiency.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

Large Relative Translation

Condition ii, Keyframe Registration

Translation from the previous keyframe w.r.t. to the median depth reaches a threshold.

$$\frac{\left\| \mathbf{t}_{\mathcal{F}_i \mathcal{F}_j} \right\|_2}{\bar{D}_{\mathcal{F}_i \mathcal{F}_j}} > \tau_2, \quad \bar{D}_{\mathcal{F}_i \mathcal{F}_j} = \frac{1}{2HW} \sum_{\{\mathcal{F}_i, \mathcal{F}_j\}} \sum_{h=0}^H \sum_{w=0}^W d(h, w) \quad (9)$$

$\in \mathbb{R}^3$, translation from \mathcal{F}_i to \mathcal{F}_j

$\in \mathbb{R}$, the median depth

In practice, $\tau_2 = 0.04$. Additionally, evaluate the Gaussian covisibility only if the relative translation is not too small (> 0.02) for efficiency.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

Large Relative Translation

Condition ii, Keyframe Registration

Translation from the previous keyframe w.r.t. to the median depth reaches a threshold.

$$\frac{\left\| \mathbf{t}_{\mathcal{F}_i \mathcal{F}_j} \right\|_2}{\bar{D}_{\mathcal{F}_i \mathcal{F}_j}} > \tau_2, \quad \bar{D}_{\mathcal{F}_i \mathcal{F}_j} = \frac{1}{2HW} \sum_{\{\mathcal{F}_i, \mathcal{F}_j\}} \sum_{h=0}^H \sum_{w=0}^W d(h, w) \quad (9)$$

$\in \mathbb{R}^3$, translation from \mathcal{F}_i to \mathcal{F}_j
 $\in \mathbb{R}$, the median depth
 $d(h, w)$ depth of pixel (h, w)

In practice, $\tau_2 = 0.04$. Additionally, evaluate the Gaussian covisibility only if the relative translation is not too small (> 0.02) for efficiency.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

Large Relative Translation

Condition ii, Keyframe Registration

Translation from the previous keyframe w.r.t. to the median depth reaches a threshold.

$$\frac{\|\mathbf{t}_{\mathcal{F}_i\mathcal{F}_j}\|_2}{\bar{D}_{\mathcal{F}_i\mathcal{F}_j}} > \tau_2, \quad \bar{D}_{\mathcal{F}_i\mathcal{F}_j} = \frac{1}{2HW} \sum_{\{\mathcal{F}_i, \mathcal{F}_j\}} \sum_{h=0}^H \sum_{w=0}^W d(h, w) \quad (9)$$

$\in \mathbb{R}^3$, translation from \mathcal{F}_i to \mathcal{F}_j
 $\in \mathbb{R}$, the median depth
 image height
 depth of pixel (h, w)

In practice, $\tau_2 = 0.04$. Additionally, evaluate the Gaussian covisibility only if the relative translation is not too small (> 0.02) for efficiency.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

Large Relative Translation

Condition ii, Keyframe Registration

Translation from the previous keyframe w.r.t. to the median depth reaches a threshold.

$$\frac{\|\mathbf{t}_{\mathcal{F}_i\mathcal{F}_j}\|_2}{\bar{D}_{\mathcal{F}_i\mathcal{F}_j}} > \tau_2, \quad \bar{D}_{\mathcal{F}_i\mathcal{F}_j} = \frac{1}{2 H W} \sum_{\{\mathcal{F}_i, \mathcal{F}_j\}} \sum_{h=0}^H \sum_{w=0}^W d(h, w) \quad (9)$$

$\in \mathbb{R}^3$, translation from \mathcal{F}_i to \mathcal{F}_j
 $\in \mathbb{R}$, the median depth
 image height
 image width
 depth of pixel (h, w)

In practice, $\tau_2 = 0.04$. Additionally, evaluate the Gaussian covisibility only if the relative translation is not too small (> 0.02) for efficiency.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

If **any** of the following conditions **is true**...

Beyond Window Capacity

Condition i, Keyframe Removal

Remove one of previous keyframes

If any of the following conditions is true...

Beyond Window Capacity

Condition i, Keyframe Removal

Remove one of previous keyframes

If any of the following conditions is true...

Beyond Window Capacity

Condition i, Keyframe Removal

Remove **one** of previous keyframes

If any of the following conditions is true...

Beyond Window Capacity

Condition i, Keyframe Removal

Remove one of previous keyframes that **minimize** the impact on the **overall baseline length**.

If any of the following conditions is true...

Beyond Window Capacity

Condition i, Keyframe Removal

Remove one of previous keyframes that minimize the impact on the overall baseline length.

$$\mathcal{F}^* = \arg \max_{\mathcal{F} \in \mathcal{W}} l(\mathcal{W} \setminus \{\mathcal{F}\}) \quad (10)$$

If any of the following conditions is true...

Beyond Window Capacity

Condition i, Keyframe Removal

Remove one of previous keyframes that minimize the impact on the overall baseline length.

$$\mathcal{F}^* = \arg \max_{\mathcal{F} \in \mathcal{W}} l(\mathcal{W} \setminus \{\mathcal{F}\}), \quad l(\mathcal{W}) = \sum_{i=1}^{|\mathcal{W}|} \sum_{j=1}^i \|\mathbf{t}_{\mathcal{F}_i \mathcal{F}_j}\| \quad (10)$$

If any of the following conditions is true...

Beyond Window Capacity

Condition i, Keyframe Removal

Remove one of previous keyframes that minimize the impact on the overall baseline length.

$$\mathcal{F}^* = \arg \max_{\mathcal{F} \in \mathcal{W}} l(\mathcal{W} \setminus \{\mathcal{F}\}), \quad l(\mathcal{W}) = \sum_{i=1}^{|\mathcal{W}|} \sum_{j=1}^i \|\mathbf{t}_{\mathcal{F}_i \mathcal{F}_j}\| \quad (10)$$

Remark: for the best multi-view constraints.

Low Gaussian Overlap Coefficient

Condition ii, Keyframe Removal

Remove multiple previous keyframes if the “Gaussian overlap coefficient” drops below a threshold.

Szymkiewicz–Simpson coefficient

In practice, $\tau_4 = 0.4$.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

Low Gaussian Overlap Coefficient

Condition ii, Keyframe Removal

Remove **multiple** previous keyframes if the “Gaussian overlap coefficient” drops below a threshold.

Szymkiewicz–Simpson coefficient

In practice, $\tau_4 = 0.4$.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

Low Gaussian Overlap Coefficient

Condition ii, Keyframe Removal

Remove multiple previous keyframes if the “Gaussian overlap coefficient” drops **below** a threshold.

Szymkiewicz–Simpson coefficient

In practice, $\tau_4 = 0.4$.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

Low Gaussian Overlap Coefficient

Condition ii, Keyframe Removal

Remove multiple previous keyframes if the “Gaussian overlap coefficient” drops below a threshold.

$$\frac{|\mathbf{v}(\mathcal{G}, \mathcal{F}_i) \cap \mathbf{v}(\mathcal{G}, \mathcal{F}_j)|}{\min(|\mathbf{v}(\mathcal{G}, \mathcal{F}_i)|, |\mathbf{v}(\mathcal{G}, \mathcal{F}_j)|)} < \tau_4 \quad (11)$$

Szymkiewicz-Simpson coefficient

In practice, $\tau_4 = 0.4$.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

Low Gaussian Overlap Coefficient

Condition ii, Keyframe Removal

Remove multiple previous keyframes if the “Gaussian overlap coefficient” drops below a threshold.

$$\frac{|\mathbf{v}(\mathcal{G}, \mathcal{F}_i) \cap \mathbf{v}(\mathcal{G}, \mathcal{F}_j)|}{\min(|\mathbf{v}(\mathcal{G}, \mathcal{F}_i)|, |\mathbf{v}(\mathcal{G}, \mathcal{F}_j)|)} < \tau_4 \quad (11)$$

Remark: not observing the same area.

Szymkiewicz–Simpson coefficient

In practice, $\tau_4 = 0.4$.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

■ Why do we need “Gaussian insertion”?

■ SLAM is for robotic exploration.

■ When do we need “Gaussian insertion”?

Keyframing

Condition i, Gaussian Insertion

Insertion is triggered for every new keyframe.

- Why do we need “Gaussian insertion”?

- SLAM is for robotic exploration.

- When do we need “Gaussian insertion”?

Keyframing

Condition i, Gaussian Insertion

Insertion is triggered for every new keyframe.

- Why do we need “Gaussian insertion”?
 - SLAM is for robotic exploration.
- When do we need “Gaussian insertion”?

Keyframing

Insertion is triggered for every new keyframe.

Condition i, Gaussian Insertion

- Why do we need “Gaussian insertion”?
 - SLAM is for robotic exploration.
- When do we need “Gaussian insertion”?

Keyframing

Condition i, Gaussian Insertion

Insertion is triggered for every new keyframe.

■ How do we insert Gaussians?

- Gaussian insertion is Gaussian initialization.

If Depth Available

Gaussian Initialization

Back-project in a per-pixel, per-Gaussian approach.

If Depth Unavailable

Gaussian Initialization

Leverage the rendered depth map.

- for pixels with depth: use the rendered depth and assign a “low” covariance.
- for pixels w/o depth: use the median of rendered depth and assign a “high” covariance.

In practice, “low”: 0.2σ ; “high”: 0.5σ , where σ is the standard deviation of the rendered depth map.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

- How do we insert Gaussians?
 - Gaussian insertion is Gaussian **initialization**.

If Depth Available

Gaussian Initialization

Back-project in a per-pixel, per-Gaussian approach.

If Depth Unavailable

Gaussian Initialization

Leverage the rendered depth map.

- for pixels with depth: use the rendered depth and assign a “low” covariance.
- for pixels w/o depth: use the median of rendered depth and assign a “high” covariance.

In practice, “low”: 0.2σ ; “high”: 0.5σ , where σ is the standard deviation of the rendered depth map.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

- How do we insert Gaussians?
 - Gaussian insertion is Gaussian initialization.

If Depth Available

Gaussian Initialization

Back-project in a per-pixel, per-Gaussian approach.

If Depth Unavailable

Gaussian Initialization

Leverage the rendered depth map.

- for pixels with depth: use the rendered depth and assign a “low” covariance.
- for pixels w/o depth: use the median of rendered depth and assign a “high” covariance.

In practice, “low”: 0.2σ ; “high”: 0.5σ , where σ is the standard deviation of the rendered depth map.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

- How do we insert Gaussians?
 - Gaussian insertion is Gaussian initialization.

If Depth Available

Gaussian Initialization

Back-project in a per-pixel, per-Gaussian approach.

If Depth Unavailable

Gaussian Initialization

Leverage the rendered depth map.

- for pixels with depth: use the rendered depth and assign a “low” covariance.
- for pixels w/o depth: use the median of rendered depth and assign a “high” covariance.

In practice, “low”: 0.2σ ; “high”: 0.5σ , where σ is the standard deviation of the rendered depth map.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

- How do we insert Gaussians?
 - Gaussian insertion is Gaussian initialization.

If Depth Available

Gaussian Initialization

Back-project in a per-pixel, per-Gaussian approach.

If Depth Unavailable

Gaussian Initialization

Leverage the rendered depth map.

- for pixels with depth: use the rendered depth and assign a “low” covariance.
- for pixels w/o depth: use the median of rendered depth and assign a “high” covariance.

In practice, “low”: 0.2σ ; “high”: 0.5σ , where σ is the standard deviation of the rendered depth map.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

- How do we insert Gaussians?
 - Gaussian insertion is Gaussian initialization.

If Depth Available

Gaussian Initialization

Back-project in a per-pixel, per-Gaussian approach.

If Depth Unavailable

Gaussian Initialization

Leverage the rendered depth map.

- for pixels with depth: use the rendered depth and assign a “low” covariance.
- for pixels w/o depth: use the median of rendered depth and assign a “high” covariance.

In practice, “low”: 0.2σ ; “high”: 0.5σ , where σ is the standard deviation of the rendered depth map.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

■ Why do we need “Gaussian Pruning”?

- if depth unavailable, too many incorrect newly inserted Gaussians.

Low Gaussian Opacity

Condition i, Gaussian Pruning

Low opacity Gaussians are pruned.

$$\{\mathcal{G}_i \in \mathcal{G} \mid \alpha(\mathcal{G}_i) < \tau_\alpha\} \quad (12)$$

Low Gaussian Covisibility

Condition ii, Gaussian Pruning

For “just” inserted Gaussians but unobserved by “some other” keyframes, are pruned out.

If no pruning, although the majority of incorrect Gaussians vanish quickly in following optimization, there are some survivals.

In practice, $\tau_\alpha = 0.7$.

In practice, the pruned Gaussians are inserted in the last 3 keyframes and unobserved by any other 3 keyframes in the sliding window.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

- Why do we need “Gaussian Pruning”?
 - if depth **unavailable**, too many **incorrect** newly inserted Gaussians.

Low Gaussian Opacity

Condition i, Gaussian Pruning

Low opacity Gaussians are pruned.

$$\{\mathcal{G}_i \in \mathcal{G} \mid \alpha(\mathcal{G}_i) < \tau_\alpha\} \quad (12)$$

Low Gaussian Covisibility

Condition ii, Gaussian Pruning

For “just” inserted Gaussians but unobserved by “some other” keyframes, are pruned out.

If no pruning, although the majority of incorrect Gaussians vanish quickly in following optimization, there are some survivals.

In practice, $\tau_\alpha = 0.7$.

In practice, the pruned Gaussians are inserted in the last 3 keyframes and unobserved by any other 3 keyframes in the sliding window.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

- Why do we need “Gaussian Pruning”?
 - if depth unavailable, too many incorrect newly inserted Gaussians.

Low Gaussian Opacity

Condition i, Gaussian Pruning

Low opacity Gaussians are pruned.

$$\{\mathcal{G}_i \in \mathcal{G} \mid \alpha(\mathcal{G}_i) < \tau_\alpha\} \quad (12)$$

Low Gaussian Covisibility

Condition ii, Gaussian Pruning

For “just” inserted Gaussians but unobserved by “some other” keyframes, are pruned out.

If no pruning, although the majority of incorrect Gaussians vanish quickly in following optimization, there are some survivals.

In practice, $\tau_\alpha = 0.7$.

In practice, the pruned Gaussians are inserted in the last 3 keyframes and unobserved by any other 3 keyframes in the sliding window.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

- Why do we need “Gaussian Pruning”?
 - if depth unavailable, too many incorrect newly inserted Gaussians.

Low Gaussian Opacity

Condition i, Gaussian Pruning

Low opacity Gaussians are pruned.

$$\{\mathcal{G}_i \in \mathcal{G} \mid \alpha(\mathcal{G}_i) < \tau_\alpha\} \quad (12)$$

Low Gaussian Covisibility

Condition ii, Gaussian Pruning

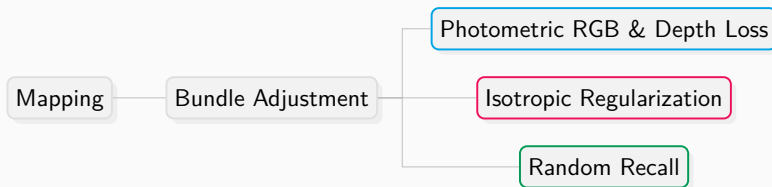
For “just” inserted Gaussians but unobserved by “some other” keyframes, are pruned out.

If no pruning, although the majority of incorrect Gaussians vanish quickly in following optimization, there are some survivals.

In practice, $\tau_\alpha = 0.7$.

In practice, the pruned Gaussians are inserted in the last 3 keyframes and unobserved by any other 3 keyframes in the sliding window.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM



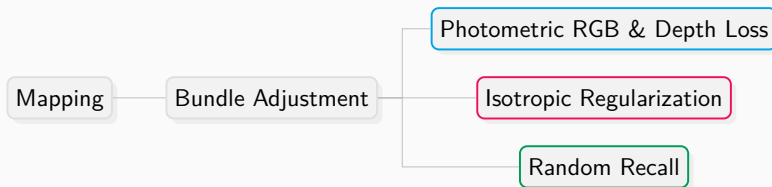
■ Why do we need mapping in 3DGS SLAM?

■ Local: Optimize newly inserted 3D Gaussians

■ Global: Reconstruct globally 3D scene structure

key method trick convention

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

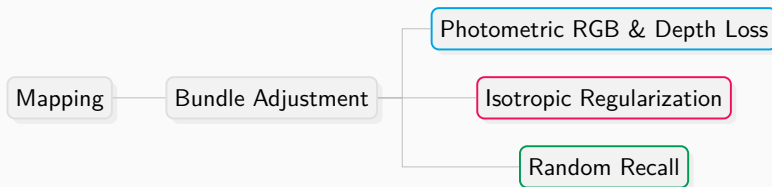


■ Why do we need mapping in 3DGS SLAM?

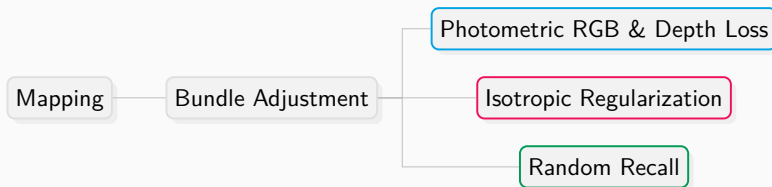
- Local: Optimize newly inserted 3D Gaussians.
- Global: Reconstruct a globally 3D-coherent structure.

key method trick convention

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM



- Why do we need mapping in **3DGS** SLAM?
 - **Local**: Optimize newly inserted 3D Gaussians.
 - Global: Reconstruct a globally 3D-coherent structure.



- Why do we need mapping in **3DGS** SLAM?
 - Local: Optimize newly inserted 3D Gaussians.
 - **Global**: Reconstruct a globally 3D-coherent structure.

Bundle Adjustment

$$\operatorname{argmin}_{\mathcal{G}, \{\mathbf{T}_{cw}(\mathcal{F}_k) | \mathcal{F}_k \in \mathcal{W}\}} \sum_{\mathcal{F}_k}^{\mathcal{W}} \mathcal{L}_{pho}(\mathcal{F}_k) \quad (13)$$

Random Recall

Besides \mathcal{W} , “some” randomly selected historical keyframes are also leveraged in BA to avoid forgetting the global map.

Bundle Adjustment

$$\underset{\mathcal{G}, \{\mathbf{T}_{cw}(\mathcal{F}_k) | \mathcal{F}_k \in \mathcal{W}\}}{\operatorname{argmin}} \sum_{\mathcal{F}_k}^{\mathcal{W}} \mathcal{L}_{pho}(\mathcal{F}_k) \quad (13)$$

3D Gaussians

Random Recall

Besides \mathcal{W} , “some” randomly selected historical keyframes are also leveraged in BA to avoid forgetting the global map.

Bundle Adjustment

$$\underset{\mathcal{G}, \{\mathbf{T}_{cw}(\mathcal{F}_k) | \mathcal{F}_k \in \mathcal{W}\}}{\operatorname{argmin}} \sum_{\mathcal{F}_k \in \mathcal{W}} \mathcal{L}_{pho}(\mathcal{F}_k) \quad (13)$$

3D Gaussians

camera poses of keyframes in the sliding window

Random Recall

Besides \mathcal{W} , “some” randomly selected historical keyframes are also leveraged in BA to avoid forgetting the global map.

Bundle Adjustment

$$\underset{\mathcal{G}, \{\mathbf{T}_{cw}(\mathcal{F}_k) | \mathcal{F}_k \in \mathcal{W}\}}{\operatorname{argmin}} \sum_{\mathcal{F}_k \in \mathcal{W}} \mathcal{L}_{pho}(\mathcal{F}_k) \quad (13)$$

3D Gaussians

camera poses of keyframes in the sliding window

Random Recall

Besides \mathcal{W} , “some” randomly selected historical keyframes are also leveraged in BA to avoid forgetting the global map.

- **Why** do we need “isotropic regularization”?
 - Observation: isotropic Gaussians behave better than anisotropic.
 - Analysis: no constraints on the elongation along the viewing ray direction, **even with depth**.

Isotropic Regularization

$$\mathcal{L}_{iso} = \sum_{i=1}^{|\mathcal{G}|} \|\mathbf{s}(\mathcal{G}_i) - \bar{\mathbf{s}}(\mathcal{G}_i)\|_1, \quad (14)$$

- Why do we need “isotropic regularization”?
 - **Observation:** isotropic Gaussians behave better than anisotropic.
 - Analysis: no constraints on the elongation along the viewing ray direction, **even with depth**.

Isotropic Regularization

$$\mathcal{L}_{iso} = \sum_{i=1}^{|\mathcal{G}|} \|\mathbf{s}(\mathcal{G}_i) - \bar{\mathbf{s}}(\mathcal{G}_i)\|_1, \quad (14)$$

- Why do we need “isotropic regularization”?
 - Observation: isotropic Gaussians behave better than anisotropic.
 - **Analysis:** no constraints on the elongation along the viewing ray direction, **even with depth.**

Isotropic Regularization

$$\mathcal{L}_{iso} = \sum_{i=1}^{|\mathcal{G}|} \|\mathbf{s}(\mathcal{G}_i) - \bar{\mathbf{s}}(\mathcal{G}_i)\|_1, \quad (14)$$

- Why do we need “isotropic regularization”?
 - Observation: isotropic Gaussians behave better than anisotropic.
 - Analysis: no constraints on the elongation along the viewing ray direction, **even with depth**.

Isotropic Regularization

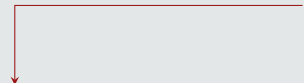
$$\mathcal{L}_{iso} = \sum_{i=1}^{|\mathcal{G}|} \|\mathbf{s}(\mathcal{G}_i) - \bar{\mathbf{s}}(\mathcal{G}_i)\|_1, \quad (14)$$

- Why do we need “isotropic regularization”?
 - Observation: isotropic Gaussians behave better than anisotropic.
 - Analysis: no constraints on the elongation along the viewing ray direction, **even with depth**.

Isotropic Regularization

$$\mathcal{L}_{iso} = \sum_{i=1}^{|\mathcal{G}|} \|\mathbf{s}(\mathcal{G}_i) - \bar{\mathbf{s}}(\mathcal{G}_i)\|_1, \quad (14)$$

$\in \mathbb{N}$, total number of Gaussians



- Why do we need “isotropic regularization”?
 - Observation: isotropic Gaussians behave better than anisotropic.
 - Analysis: no constraints on the elongation along the viewing ray direction, **even with depth**.

Isotropic Regularization

$$\mathcal{L}_{iso} = \sum_{i=1}^{|G|} \| \mathbf{s}(\mathcal{G}_i) - \bar{\mathbf{s}}(\mathcal{G}_i) \|_1, \quad (14)$$

$\in \mathbb{N}$, total number of Gaussians

$\in \mathbb{R}^3$, scale of i -th Gaussian

- Why do we need “isotropic regularization”?
- Observation: isotropic Gaussians behave better than anisotropic.
- Analysis: no constraints on the elongation along the viewing ray direction, **even with depth**.

Isotropic Regularization

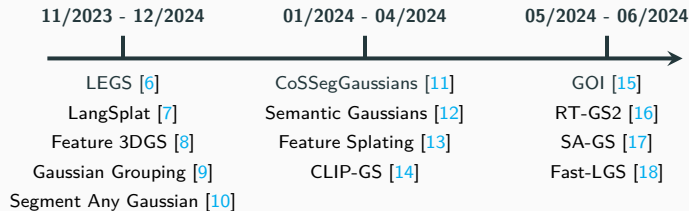
$$\mathcal{L}_{iso} = \sum_{i=1}^{|\mathcal{G}|} \| \mathbf{s}(\mathcal{G}_i) - \bar{\mathbf{s}}(\mathcal{G}_i) \|_1, \quad \bar{\mathbf{s}}(\mathcal{G}_i) = \begin{bmatrix} (s(\mathcal{G}_i)^x + s(\mathcal{G}_i)^y + s(\mathcal{G}_i)^z) / 3 \\ (s(\mathcal{G}_i)^x + s(\mathcal{G}_i)^y + s(\mathcal{G}_i)^z) / 3 \\ (s(\mathcal{G}_i)^x + s(\mathcal{G}_i)^y + s(\mathcal{G}_i)^z) / 3 \end{bmatrix} \quad (14)$$

$|\mathcal{G}| \in \mathbb{N}$, total number of Gaussians
 $\mathbf{s}(\mathcal{G}_i) \in \mathbb{R}^3$, scale of i -th Gaussian
 $\bar{\mathbf{s}}(\mathcal{G}_i) \in \mathbb{R}^3$, averaged scale of i -th Gaussian

Overall Optimization for Mapping

$$\operatorname{argmin}_{\mathcal{G}, \{\mathbf{T}_{cw}(\mathcal{F}_k) | \mathcal{F}_k \in \mathcal{W}^+\}} \sum_{\mathcal{F}_k}^{\mathcal{W}^+} \mathcal{L}_{pho}(\mathcal{F}_k) + \lambda_{iso} \mathcal{L}_{iso} \quad (15)$$

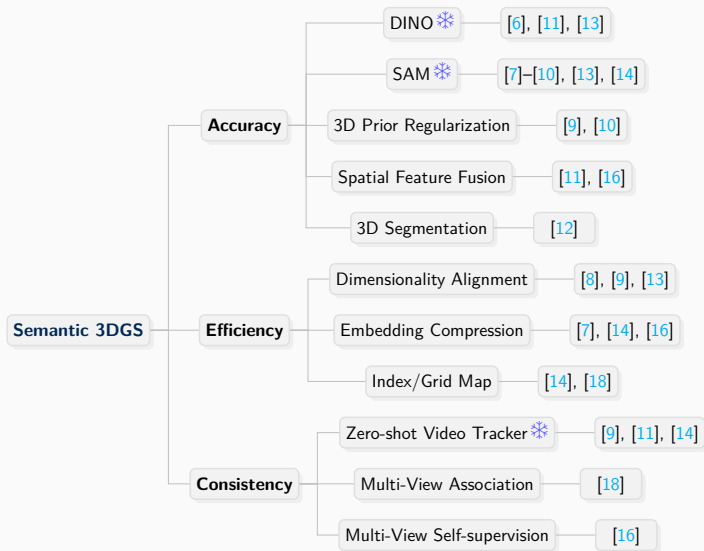
Semantic 3DGS

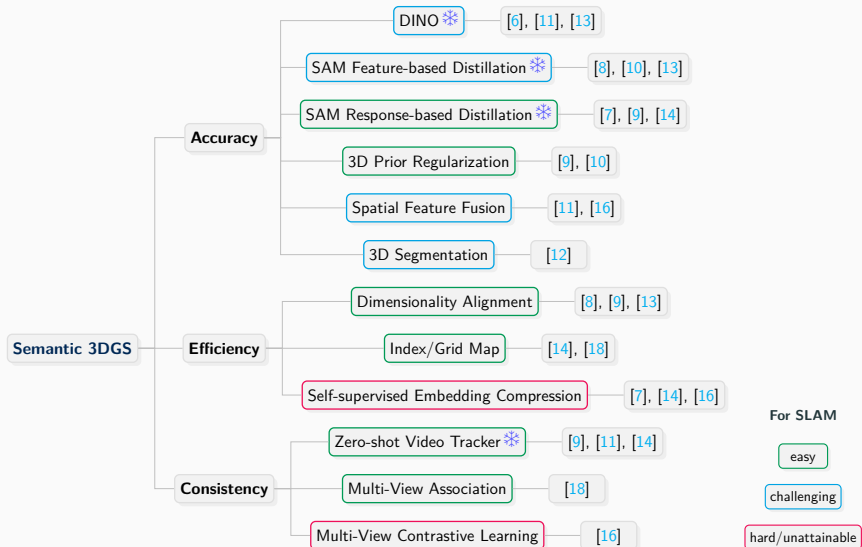


Consensus

Lift 2D foundation models to 3D scene-specific Gaussians under 2D supervision.

1. 2D foundation models: CLIP, SAM, DINO, etc.
2. Interactivity: manipulation, edit, localization, query, simulation, etc.





Semantic 3DGS

Feature-3DGS

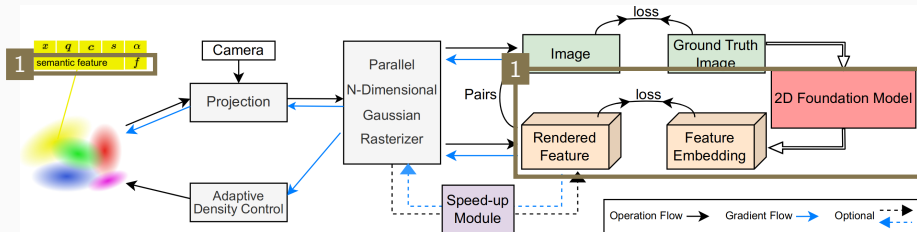


Figure 1: Overview of Feature 3DGS

1 Semantic Rendering Pipeline

Differentiable rendering of Gaussian-wise latent semantic features.

2 Speed-up module

Dimensionality Alignment.

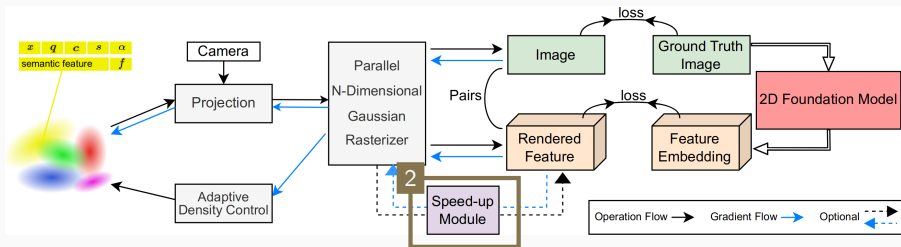


Figure 1: Overview of Feature 3DGS

1 Semantic Rendering Pipeline

Differentiable rendering of Gaussian-wise latent semantic features.

2 Speed-up module

Dimensionality Alignment.

To render semantic embeddings, i.e.

$$\mathbb{R}^{N \times D} \mapsto \mathbb{R}^{H \times W \times D} \quad (16)$$

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

To render semantic embeddings, i.e.

$$\mathbb{R}^{\underbrace{N}_{\text{number of 3D Gaussians}}} \times D \mapsto \mathbb{R}^{H \times W \times D} \quad (16)$$

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

To render semantic embeddings, i.e.

$$\mathbb{R}^{N \times D} \mapsto \mathbb{R}^{H \times W \times D} \quad (16)$$

number of 3D Gaussians

dimension of semantic embedding

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

To render semantic embeddings, i.e.

$$\mathbb{R}^{N \times D} \mapsto \mathbb{R}^{H \times W \times D} \quad (16)$$

number of 3D Gaussians N
dimension of semantic embedding D
image height H

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

To render semantic embeddings, i.e.

$$\mathbb{R}^{N \times D} \mapsto \mathbb{R}^{H \times W \times D} \quad (16)$$

Diagram illustrating the mapping from 3D Gaussian representation to semantic embedding representation:

- N : number of 3D Gaussians (indicated by a red arrow)
- D : dimension of semantic embedding (indicated by a green arrow)
- H : image height (indicated by a teal arrow)
- W : image width (indicated by a purple arrow)

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

To render semantic embeddings, i.e.

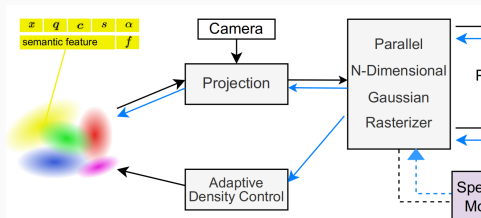
$$\mathbb{R}^{N \times D} \mapsto \mathbb{R}^{H \times W \times D}$$

N : number of 3D Gaussians
 D : dimension of semantic embedding
 H : image height
 W : image width

(16)

5 things,

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



To render semantic embeddings, i.e.

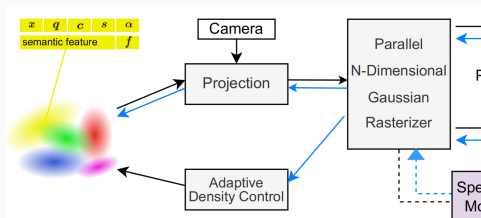
$$\mathbb{R}^{N \times D} \mapsto \mathbb{R}^{H \times W \times D}$$

N : number of 3D Gaussians
 D : dimension of semantic embedding
 H : image height
 W : image width

(16)

5 things,

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



To render semantic embeddings, i.e.

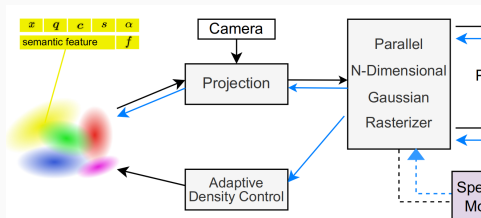
$$\mathbb{R}^{N \times D} \mapsto \mathbb{R}^{H \times W \times D}$$

N : number of 3D Gaussians
 D : dimension of semantic embedding
 H : image height
 W : image width

(16)

5 things,

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



To render semantic embeddings, i.e.

$$\mathbb{R}^{N \times D} \mapsto \mathbb{R}^{H \times W \times D}$$

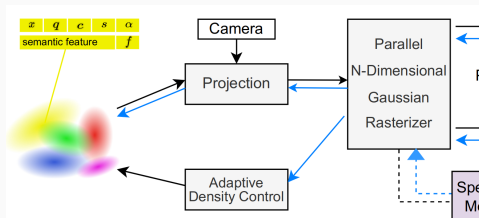
number of 3D Gaussians N → dimension of semantic embedding D

image height H → image width W

(16)

5 things,

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



To render semantic embeddings, i.e.

$$\mathbb{R}^{N \times D} \mapsto \mathbb{R}^{H \times W \times D}$$

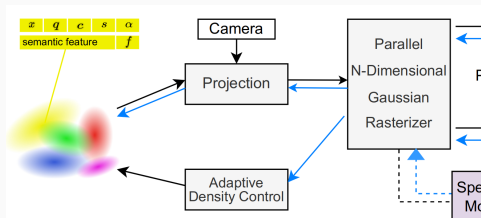
number of 3D Gaussians N → dimension of semantic embedding D

image height H → image width W

(16)

5 things,

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



To render semantic embeddings, i.e.

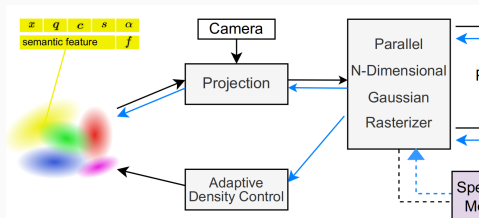
$$\mathbb{R}^{N \times D} \mapsto \mathbb{R}^{H \times W \times D}$$

N : number of 3D Gaussians
 D : dimension of semantic embedding
 H : image height
 W : image width

(16)

5 things,

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



Representation: 3D Gaussian augmented with a latent embedding.

$$\mathcal{G}_i = \{\mathbf{x}, \mathbf{q}, \mathbf{s}, \alpha, \mathbf{c}, \mathbf{f}\} \quad (17)$$

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](#)

Representation: 3D Gaussian augmented with a latent embedding.

$$\mathcal{G}_i = \{\mathbf{x}, \mathbf{q}, \mathbf{s}, \alpha, \mathbf{c}, \mathbf{f}\} \quad (17)$$

optimizable attributes of a 3D Gaussian

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

Representation: 3D Gaussian augmented with a latent embedding.

$$\mathcal{G}_i = \{\mathbf{x}, \mathbf{q}, \mathbf{s}, \alpha, \mathbf{c}, \mathbf{f}\} \quad (17)$$

optimizable attributes of a 3D Gaussian

$\in \mathbb{N}$, index of 3D Gaussian

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

Representation: 3D Gaussian augmented with a latent embedding.

$$\mathcal{G}_i = \{\mathbf{x}, \mathbf{q}, \mathbf{s}, \alpha, \mathbf{c}, \mathbf{f}\} \quad (17)$$

optimizable attributes of a 3D Gaussian

$\in \mathbb{N}$, index of 3D Gaussian

$\in \mathbb{R}^3$, position

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

Representation: 3D Gaussian augmented with a latent embedding.

$$\mathcal{G}_i = \{ \mathbf{x}, \mathbf{q}, \mathbf{s}, \alpha, \mathbf{c}, \mathbf{f} \} \quad (17)$$

Diagram illustrating the components of the 3D Gaussian representation \mathcal{G}_i :

- \mathcal{G} (pink box) and i (yellow box) are labeled "optimizable attributes of a 3D Gaussian".
- i is labeled " $\in \mathbb{N}$, index of 3D Gaussian".
- \mathbf{x} (light green box) is labeled " $\in \mathbb{R}^3$, position".
- \mathbf{q} (light green box) is labeled " $\in \text{SO}(3)$, rotation".

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

Representation: 3D Gaussian augmented with a latent embedding.

$$\mathcal{G}_i = \{ \mathbf{x}, \mathbf{q}, \mathbf{s}, \alpha, \mathbf{c}, \mathbf{f} \} \quad (17)$$

optimizable attributes of a 3D Gaussian

$\in \mathbb{N}$, index of 3D Gaussian

$\in \mathbb{R}^3$, position

$\in \text{SO}(3)$, rotation

$\in \mathbb{R}^3$, scale

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

Representation: 3D Gaussian augmented with a latent embedding.

$$\mathcal{G}_i = \{ \mathbf{x}, \mathbf{q}, \mathbf{s}, \alpha, \mathbf{c}, \mathbf{f} \} \quad (17)$$

Diagram illustrating the representation of a 3D Gaussian augmented with a latent embedding. The representation is defined as $\mathcal{G}_i = \{ \mathbf{x}, \mathbf{q}, \mathbf{s}, \alpha, \mathbf{c}, \mathbf{f} \}$, where i is the index of the 3D Gaussian.

The attributes are defined as follows:

- $\mathbf{x} \in \mathbb{R}^3$, position
- $\mathbf{q} \in \text{SO}(3)$, rotation
- $\mathbf{s} \in \mathbb{R}^3$, scale
- $\alpha \in [0, 1]$, opacity
- \mathbf{c}, \mathbf{f} are optimizable attributes of a 3D Gaussian
- $i \in \mathbb{N}$, index of 3D Gaussian

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

Representation: 3D Gaussian augmented with a latent embedding.

$$\mathcal{G}_i = \{\mathbf{x}, \mathbf{q}, \mathbf{s}, \alpha, \mathbf{c}, \mathbf{f}\} \quad (17)$$

Diagram illustrating the components of a 3D Gaussian representation \mathcal{G}_i :

- $\mathbf{x} \in \mathbb{R}^3$, position
- $\mathbf{q} \in \text{SO}(3)$, rotation
- $\mathbf{s} \in \mathbb{R}^3$, scale
- $\alpha \in [0, 1]$, opacity
- $\mathbf{c} \in \mathbb{R}^{3n}$, color
- \mathbf{f} , latent embedding

The index i is defined as $i \in \mathbb{N}$, index of 3D Gaussian.

The set of attributes $\{\mathbf{x}, \mathbf{q}, \mathbf{s}, \alpha, \mathbf{c}, \mathbf{f}\}$ are the optimizable attributes of a 3D Gaussian.

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

Representation: 3D Gaussian augmented with a latent embedding.

$$\mathcal{G}_i = \{ \mathbf{x}, \mathbf{q}, \mathbf{s}, \alpha, \mathbf{c}, \mathbf{f} \} \quad (17)$$

optimizable attributes of a 3D Gaussian

$\in \mathbb{N}$, index of 3D Gaussian

$\in \mathbb{R}^3$, position

$\in \text{SO}(3)$, rotation

$\in \mathbb{R}^3$, scale

$\in [0, 1]$, opacity

$\in \mathbb{R}^{3n}$, color

$\in \mathbb{R}^3$, semantic feature

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

Projection: from 3D ellipsoids to 2D ellipses.

Projection: from 3D ellipsoids to 2D ellipses.

$$\mu_i = \pi(\mathbf{T}_{cw} \cdot \mu_w) \quad (18)$$

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (19)$$

Projection: from 3D ellipsoids to 2D ellipses.

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (18)$$

$\in \mathbb{P}^3$, 3D(world) mean

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (19)$$

Projection: from 3D ellipsoids to 2D ellipses.

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (18)$$

$\in \mathbb{P}^3$, 3D(world) mean

$\in \text{SE}(3)$, camera pose

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (19)$$

Projection: from 3D ellipsoids to 2D ellipses.

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (18)$$

$\in \mathbb{P}^3$, 3D(world) mean

$\in \text{SE}(3)$, camera pose

projection

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (19)$$

Projection: from 3D ellipsoids to 2D ellipses.

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (18)$$

Diagram illustrating the projection process:

- μ_w (red box) is the 3D(world) mean, $\in \mathbb{P}^3$.
- \mathbf{T}_{cw} (green box) is the camera pose, $\in \text{SE}(3)$.
- π (cyan box) is the projection operation.
- μ_i (purple box) is the 2D(image) mean, $\in \mathbb{P}^2$.

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (19)$$

Projection: from 3D ellipsoids to 2D ellipses.

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (18)$$

$\mu_i \in \mathbb{P}^2$, 2D(image) mean
 π projection
 $\mathbf{T}_{cw} \in \text{SE}(3)$, camera pose
 $\mu_w \in \mathbb{P}^3$, 3D(world) mean

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (19)$$

$\Sigma_w \in \mathbb{R}^{3 \times 3}$, 3D(world) covariance

Projection: from 3D ellipsoids to 2D ellipses.

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (18)$$

Diagram illustrating the projection of 3D ellipsoids to 2D ellipses:

- μ_i (purple box) is the 2D mean, $\in \mathbb{P}^2$, 2D(image) mean.
- π (cyan box) is the projection function.
- \mathbf{T}_{cw} (green box) is the camera pose, $\in \text{SE}(3)$, camera pose.
- μ_w (red box) is the 3D mean, $\in \mathbb{P}^3$, 3D(world) mean.

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (19)$$

Diagram illustrating the projection of 3D ellipsoids to 2D ellipses:

- Σ_i (purple box) is the 2D covariance.
- \mathbf{J}_π (cyan box) is the Jacobian of the projection function.
- \mathbf{R}_{cw} (green box) is the rotation component of the camera pose, $\in \text{SO}(3)$, rotation component of \mathbf{T}_{cw} .
- Σ_w (red box) is the 3D covariance, $\in \mathbb{R}^{3 \times 3}$, 3D(world) covariance.

Projection: from 3D ellipsoids to 2D ellipses.

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (18)$$

Diagram illustrating the projection of a 3D mean vector μ_w to a 2D mean vector μ_i . The 3D mean vector μ_w (red box) is transformed by the camera pose \mathbf{T}_{cw} (green box) and then projected (cyan box) to the 2D mean vector μ_i (purple box). The camera pose \mathbf{T}_{cw} is in $\text{SE}(3)$, the 3D mean vector μ_w is in \mathbb{P}^3 , and the 2D mean vector μ_i is in \mathbb{P}^2 .

$\mu_i \in \mathbb{P}^2$, 2D(image) mean

projection

$\mathbf{T}_{cw} \in \text{SE}(3)$, camera pose

$\mu_w \in \mathbb{P}^3$, 3D(world) mean

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (19)$$

Diagram illustrating the projection of a 3D covariance matrix Σ_w to a 2D covariance matrix Σ_i . The 3D covariance matrix Σ_w (red box) is transformed by the rotation component \mathbf{R}_{cw} (green box) and then projected (cyan box) to the 2D covariance matrix Σ_i (purple box). The rotation component \mathbf{R}_{cw} is in $\text{SO}(3)$, the 3D covariance matrix Σ_w is in $\mathbb{R}^{3 \times 3}$, and the 2D covariance matrix Σ_i is in $\mathbb{R}^{2 \times 3}$.

$\Sigma_i \in \mathbb{R}^{2 \times 3}$, Jacobian of the linear approximation of π

$\mathbf{R}_{cw} \in \text{SO}(3)$, rotation component of \mathbf{T}_{cw}

$\Sigma_w \in \mathbb{R}^{3 \times 3}$, 3D(world) covariance

Projection: from 3D ellipsoids to 2D ellipses.

$$\mu_i = \pi \left(\mathbf{T}_{cw} \cdot \mu_w \right) \quad (18)$$

Diagram illustrating the projection of a 3D mean vector μ_w to a 2D mean vector μ_i . The 3D mean μ_w (red box) is transformed by the camera pose \mathbf{T}_{cw} (green box) and then projected (π , cyan box) to the 2D mean μ_i (purple box). The camera pose \mathbf{T}_{cw} is in $\text{SE}(3)$, the 3D mean μ_w is in \mathbb{P}^3 , and the 2D mean μ_i is in \mathbb{P}^2 .

$$\Sigma_i = \mathbf{J}_\pi \mathbf{R}_{cw} \Sigma_w \mathbf{R}_{cw}^T \mathbf{J}_\pi^T \quad (19)$$

Diagram illustrating the projection of a 3D covariance matrix Σ_w to a 2D covariance matrix Σ_i . The 3D covariance Σ_w (red box) is transformed by the rotation component \mathbf{R}_{cw} (green box) and the Jacobian of the linear approximation of the projection \mathbf{J}_π (cyan box) to the 2D covariance Σ_i (purple box). The rotation component \mathbf{R}_{cw} is in $\text{SO}(3)$, the 3D covariance Σ_w is in $\mathbb{R}^{3 \times 3}$, and the 2D covariance Σ_i is in $\mathbb{R}^{2 \times 2}$.

Blending: α -blending of semantic embeddings.

$$\mathbf{f}(h, w) = \sum_{i=1}^N T_i \alpha_i \mathbf{f}_i(h, w), \quad T_i = \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (20)$$

Blending: α -blending of semantic embeddings.

$$\underbrace{\mathbf{f}(h, w)}_{\text{semantic feature on pixel } (h, w)} = \sum_{i=1}^N T_i \alpha_i \mathbf{f}_i(h, w), \quad T_i = \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (20)$$

Blending: α -blending of semantic embeddings.

$$\mathbf{f}(h, w) = \sum_{i=1}^N T_i \alpha_i \mathbf{f}_i(h, w), \quad T_i = \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (20)$$

semantic feature on pixel (h, w)

semantic feature of i -th Gaussian on pixel (h, w)

Blending: α -blending of semantic embeddings.

$$\mathbf{f}(h, w) = \sum_{i=1}^N T_i \alpha_i \mathbf{f}_i(h, w), \quad T_i = \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (20)$$

Diagram illustrating the semantic feature blending process:

- The final semantic feature on pixel (h, w) is $\mathbf{f}(h, w)$ (highlighted in a pink box).
- This is calculated as a weighted sum of individual Gaussian semantic features $\mathbf{f}_i(h, w)$ (highlighted in a yellow box).
- The weight for each Gaussian is the product of its opacity α_i (highlighted in a green box) and the transmission T_i .
- The transmission T_i is the product of $(1 - \alpha_j)$ for all previous Gaussians $j < i$.

Annotations with arrows:

- Red arrow: semantic feature on pixel (h, w) points to $\mathbf{f}(h, w)$.
- Green arrow: opacity of i -th Gaussian points to α_i .
- Yellow arrow: semantic feature of i -th Gaussian on pixel (h, w) points to $\mathbf{f}_i(h, w)$.

Blending: α -blending of semantic embeddings.

$$\mathbf{f}(h, w) = \sum_{i=1}^N T_i \alpha_i \mathbf{f}_i(h, w), \quad T_i = \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (20)$$

The diagram illustrates the semantic blending equation (20) with color-coded components and labels:

- Red box:** $\mathbf{f}(h, w)$ is labeled "semantic feature on pixel (h, w) ".
- Blue box:** T_i is labeled "background opacity for i -th Gaussian".
- Green box:** α_i is labeled "opacity of i -th Gaussian".
- Yellow box:** $\mathbf{f}_i(h, w)$ is labeled "semantic feature of i -th Gaussian on pixel (h, w) ".

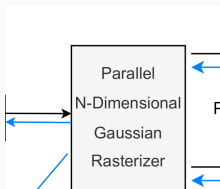
Blending: α -blending of semantic embeddings.

$$\mathbf{f}(h, w) = \sum_{i=1}^N T_i \alpha_i \mathbf{f}_i(h, w), \quad T_i = \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (20)$$

Diagram illustrating the blending process for semantic embeddings:

- 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)

Rasterization: tiled and implemented in CUDA.

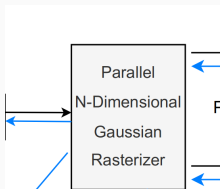


- 1 Divide the screen space into tiles (CUDA thread blocks).
- 2 Group the Gaussians by view frustum and tile index.
- 3 Sort the Gaussians by front-to-back depth order.
- 4 Blend each pixel within a tile in parallel (CUDA threads).

In practice, 16×16 blocks.

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

Rasterization: tiled and implemented in CUDA.

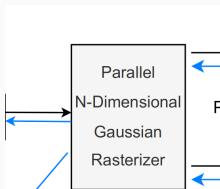


- 1 **Divide** the screen space into tiles (CUDA thread blocks).
- 2 Group the Gaussians by view frustum and tile index.
- 3 Sort the Gaussians by front-to-back depth order.
- 4 Blend each pixel within a tile in parallel (CUDA threads).

In practice, 16×16 blocks.

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

Rasterization: tiled and implemented in CUDA.

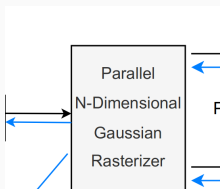


- 1 Divide the screen space into tiles (CUDA thread blocks).
- 2 **Group** the Gaussians by view frustum and tile index.
- 3 Sort the Gaussians by front-to-back depth order.
- 4 Blend each pixel within a tile in parallel (CUDA threads).

In practice, 16×16 blocks.

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

Rasterization: tiled and implemented in CUDA.

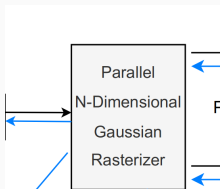


- 1 Divide the screen space into tiles (CUDA thread blocks).
- 2 Group the Gaussians by view frustum and tile index.
- 3 **Sort** the Gaussians by front-to-back depth order.
- 4 Blend each pixel within a tile in parallel (CUDA threads).

In practice, 16×16 blocks.

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

Rasterization: tiled and implemented in CUDA.

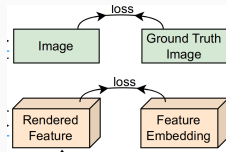


- 1 Divide the screen space into tiles (CUDA thread blocks).
- 2 Group the Gaussians by view frustum and tile index.
- 3 Sort the Gaussians by front-to-back depth order.
- 4 **Blend** each pixel within a tile in parallel (CUDA threads).

In practice, 16×16 blocks.

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

Inverse rendering: guided by image-wise photometric loss,

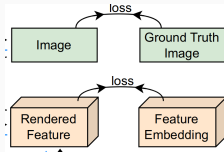


$$\mathcal{L} = \mathcal{L}_{appearance} + \gamma \mathcal{L}_{semantics} \quad (21)$$

In practice, $\gamma = 1$, $\lambda = 0.2$.

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

Inverse rendering: guided by image-wise photometric loss,

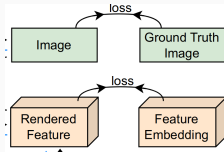


$$\mathcal{L} = \mathcal{L}_{appearance} + \gamma \mathcal{L}_{semantics} \quad (21)$$

In practice, $\gamma = 1$, $\lambda = 0.2$.

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

Inverse rendering: guided by image-wise photometric loss,



$$\mathcal{L} = \mathcal{L}_{appearance} + \gamma \mathcal{L}_{semantics} \quad (21)$$

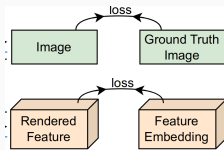
$$\mathcal{L}_{appearance} = (1 - \lambda) \mathcal{L}_1(\mathbf{C}, \hat{\mathbf{C}}) + \lambda \mathcal{L}_{D-SSIM}(\mathbf{C}, \hat{\mathbf{C}}) \quad (22)$$

$$(23)$$

In practice, $\gamma = 1$, $\lambda = 0.2$.

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

Inverse rendering: guided by image-wise photometric loss,



$$\mathcal{L} = \mathcal{L}_{appearance} + \gamma \mathcal{L}_{semantics} \quad (21)$$

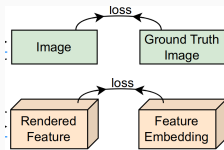
$$\mathcal{L}_{appearance} = (1 - \lambda) \mathcal{L}_1 \left(\overset{\text{captured RGB image (GT)}}{\mathbf{C}}, \overset{\text{rendered RGB image}}{\hat{\mathbf{C}}} \right) + \lambda \mathcal{L}_{D-SSIM}(\mathbf{C}, \hat{\mathbf{C}}) \quad (22)$$

$$(23)$$

In practice, $\gamma = 1$, $\lambda = 0.2$.

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

Inverse rendering: guided by image-wise photometric loss,



$$\mathcal{L} = \mathcal{L}_{appearance} + \gamma \mathcal{L}_{semantics} \quad (21)$$

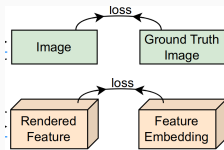
$$\mathcal{L}_{appearance} = (1 - \lambda) \mathcal{L}_1 \left(\overset{\text{captured RGB image (GT)}}{\mathbf{C}}, \overset{\text{rendered RGB image}}{\hat{\mathbf{C}}} \right) + \lambda \mathcal{L}_{D-SSIM}(\mathbf{C}, \hat{\mathbf{C}}) \quad (22)$$

$$\mathcal{L}_{semantics} = \mathcal{L}_1(\mathbf{F}, \hat{\mathbf{F}}) = \sum_{h=1}^H \sum_{w=1}^W \|\mathbf{f}(h, w) - \hat{\mathbf{f}}(h, w)\|_1 \quad (23)$$

In practice, $\gamma = 1$, $\lambda = 0.2$.

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

Inverse rendering: guided by image-wise photometric loss,



$$\mathcal{L} = \mathcal{L}_{appearance} + \gamma \mathcal{L}_{semantics} \quad (21)$$

$$\mathcal{L}_{appearance} = (1 - \lambda) \mathcal{L}_1 \left(\overset{\text{captured RGB image (GT)}}{\mathbf{C}}, \overset{\text{rendered RGB image}}{\hat{\mathbf{C}}} \right) + \lambda \mathcal{L}_{D-SSIM} \left(\mathbf{C}, \hat{\mathbf{C}} \right) \quad (22)$$

$$\mathcal{L}_{semantics} = \mathcal{L}_1 \left(\mathbf{F}, \hat{\mathbf{F}} \right) = \sum_{h=1}^H \sum_{w=1}^W \|\mathbf{f}(h, w) - \hat{\mathbf{f}}(h, w)\|_1 \quad (23)$$

inferred semantic feature map rendered semantic feature map

In practice, $\gamma = 1$, $\lambda = 0.2$.

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

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

- Compactness: to embed Gaussians with more compact vectors, $\dim = D' < D$.
- 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](#)

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

- Compactness: to embed Gaussians with more compact vectors, $\dim = D' < D$.
- 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](#)

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

- Compactness: to embed Gaussians with more compact vectors, $\dim = D' < D$.
- 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](#)

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](#)

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](#)

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](#)

1 Inefficiency

- “Speed-up module” is not enough,

2 3D Inconsistency & Inaccuracy

- 2D foundation models are still 2D.

1 Inefficiency

- “Speed-up module” is not enough,

2 3D Inconsistency & Inaccuracy

- 2D foundation models are still 2D.

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.

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.

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

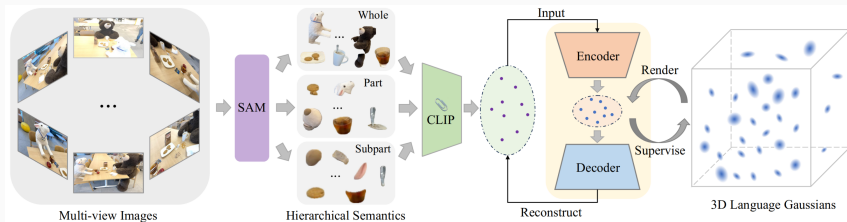


Figure 2: Overview of LangSplat

1 Accuracy: SAM outputs to enhance CLIP features.

- CLIP: image-aligned training leads to “point-ambiguity”.
- SAM: pixel-aligned & object-centered & multi-granularity.

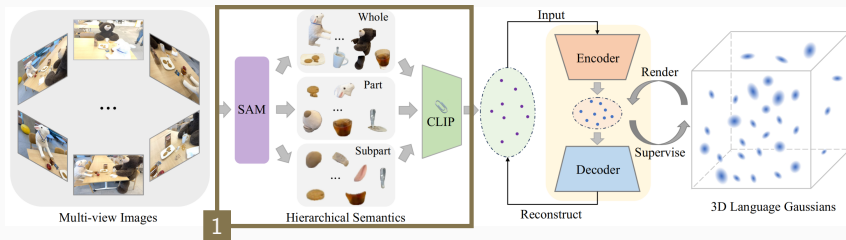
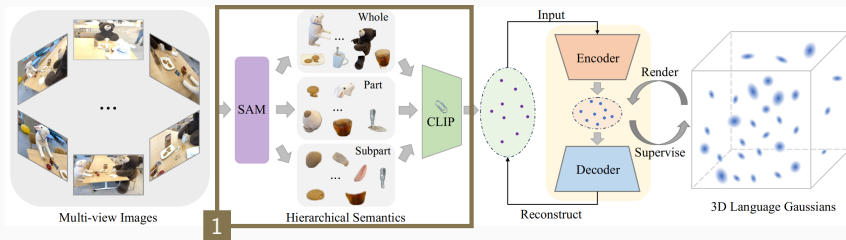


Figure 2: Overview of LangSplat

1 Accuracy: SAM outputs to enhance CLIP features.

- CLIP: image-aligned training leads to “point-ambiguity”.
- SAM: pixel-aligned & object-centered & multi-granularity.



1 Accuracy: SAM outputs to enhance CLIP features.

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

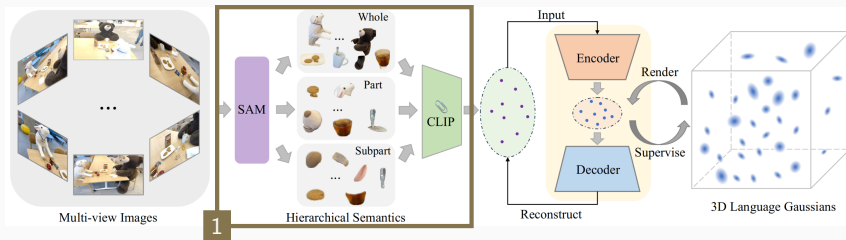


Figure 2: Overview of LangSplat

1 Accuracy: SAM outputs to enhance CLIP features.

- CLIP: image-aligned training leads to “point-ambiguity”.
- SAM: pixel-aligned & object-centered & multi-granularity.

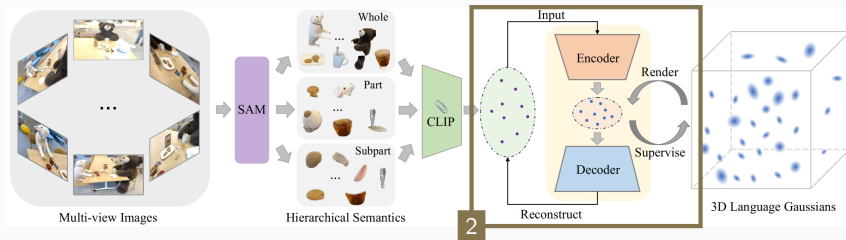


Figure 2: Overview of LangSplat

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

- More complexity and better compression, compared with “speed-up module” in Feature 3DGS [8].

In practice, $\text{dim} = 3$.

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

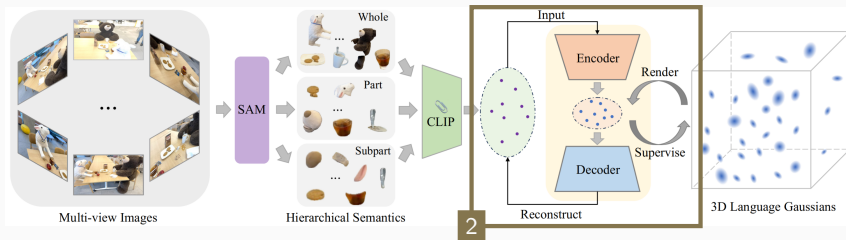


Figure 2: Overview of LangSplat

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

- More complexity and better compression, compared with “speed-up module” in Feature 3DGS [8].

Semantic 3DGS

CLIP-GS

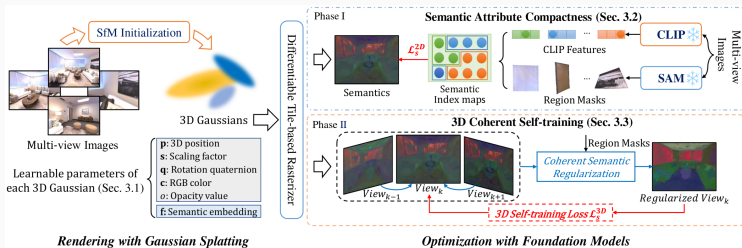


Figure 3: Overview of CLIP-GS

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

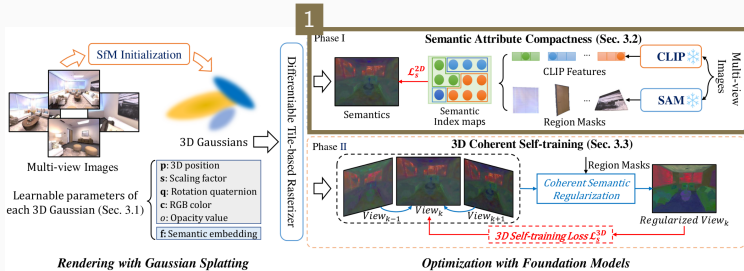


Figure 3: Overview of CLIP-GS

1 **Efficiency:** unify semantic features within an object by leveraging SAM.

2 **Consistency:** supervise consecutive frames by video segmentation.

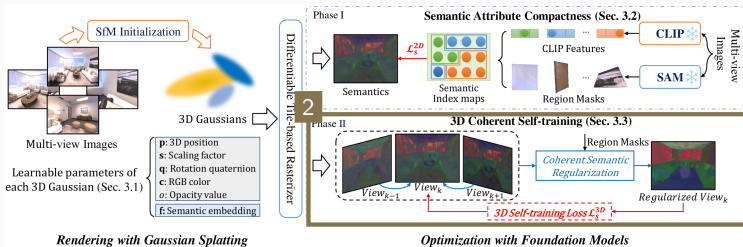


Figure 3: Overview of CLIP-GS

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

References

- [1] 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. iv).
- [2] 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. iv).
- [3] 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. iv).
- [4] 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. iv).
- [5] J. Engel, V. Koltun, and D. Cremers, “Direct sparse odometry,” in *arXiv:1607.02565*, Jul. 2016 (cit. on pp. xxvii–xxix).
- [6] 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. xcv–xcvii).
- [7] 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. xcv–xcvii).
- [8] 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. xcv–xcvii, clxiv, clxv).
- [9] 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. xcv–xcvii).

- [10] 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. xcv–xcvii).
- [11] 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. xcv–xcvii).
- [12] 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. xcv–xcvii).
- [13] 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. xcv–xcvii).
- [14] 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. xcv–xcvii).
- [15] 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. xcv).
- [16] 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. xcv–xcvii).
- [17] 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. xcv).

- [18] 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. xcv–xcvii).