

NeRF/3DGS-based SLAM

Literature Review

Shuqi XIAO

July 1, 2024

- 1 Overview
 - NeRF-based SLAM
 - 3DGS-based SLAM
- 2 NeRF-based SLAM
- 3 MonoGS
 - Methodology

Overview

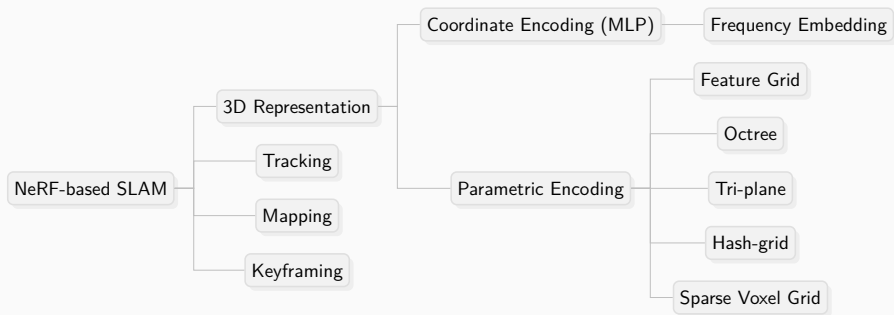
11/2023 - 12/2024

01/2024 - 04/2024

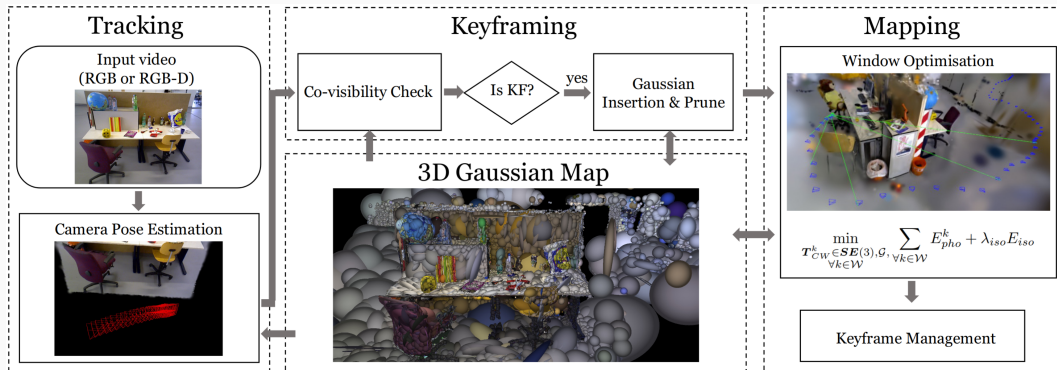
05/2024 - 06/2024

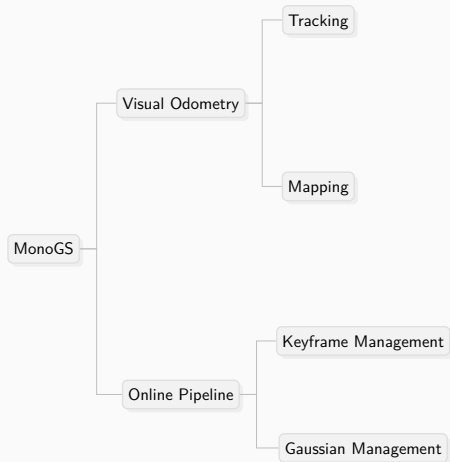


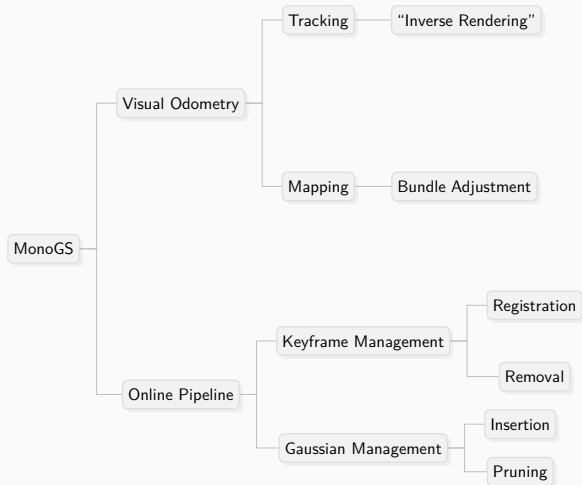
NeRF-based SLAM

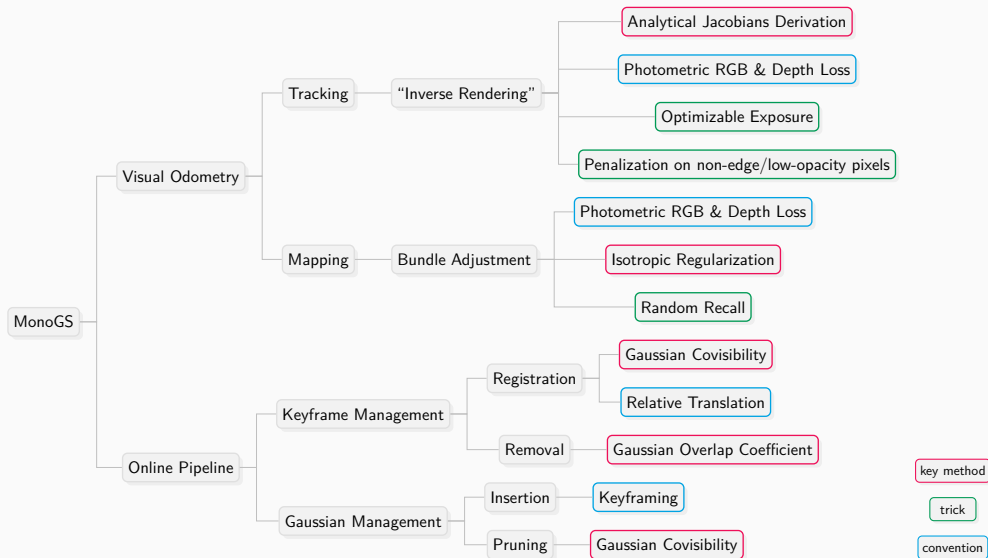


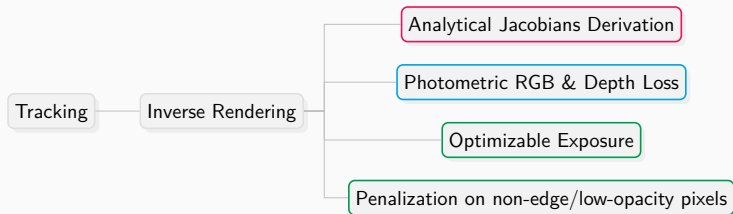
MonoGS









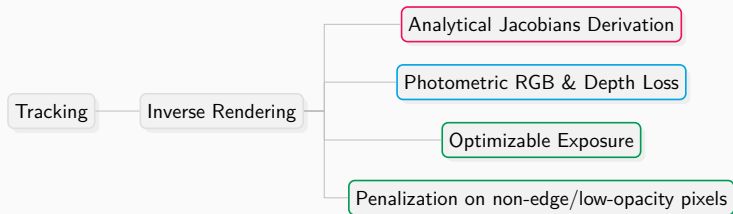


Track camera poses by inverse rendering,

key method

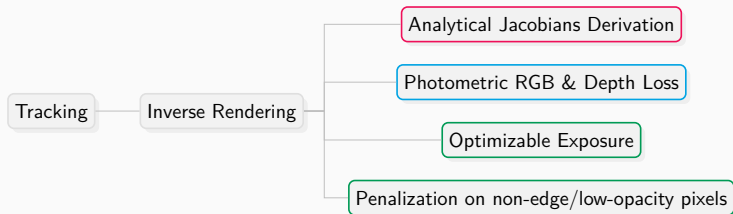
trick

convention



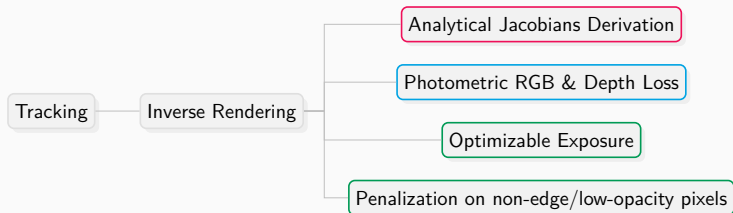
Track camera poses by inverse rendering,

- through the extended differentiable rendering pipeline,



Track camera poses by inverse rendering,

- through the extended differentiable rendering pipeline,
- by a direct optimization against fixed 3D Gaussians,



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 to the 2D image plane μ_i . The projection function π (teal box) takes the camera pose \mathbf{T}_{cw} (green box, $\in \text{SE}(3)$) and the 3D world mean μ_w (red box, $\in \mathbb{P}^3$) as inputs. The result is the 2D image mean μ_i .

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

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 a 3D Gaussian mean μ_w (red box) to a 2D image mean μ_i (purple box) using the camera pose \mathbf{T}_{cw} (green box). The projection is denoted by π . The 3D mean μ_w is in \mathbb{P}^3 , the camera pose \mathbf{T}_{cw} is in $\text{SE}(3)$, and the 2D mean μ_i is in \mathbb{P}^2 .

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

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

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

Σ_i
 \mathbf{J}_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)
 $\mathbf{T}_{cw} \in \text{SE}(3), \text{ camera pose}$
 $\mu_w \in \mathbb{P}^3, 3D(\text{world}) \text{ mean}$

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

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

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

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

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

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

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

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 chain rule,

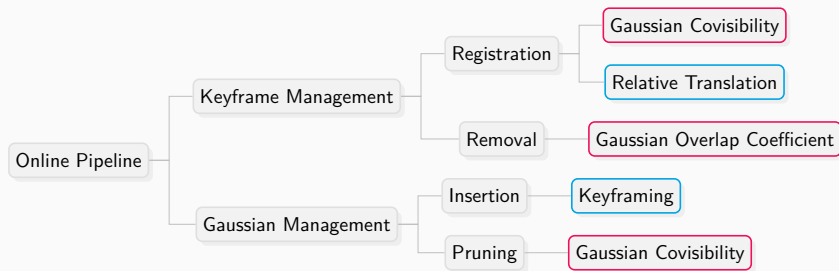
$$\frac{\partial \mu_i}{\partial \mathbf{T}_{cw}} = \frac{\partial \mu_i}{\partial \mu_c} \frac{\partial \mu_c}{\partial \mathbf{T}_{cw}} \quad (4)$$

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

The Lie Algebra,

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

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

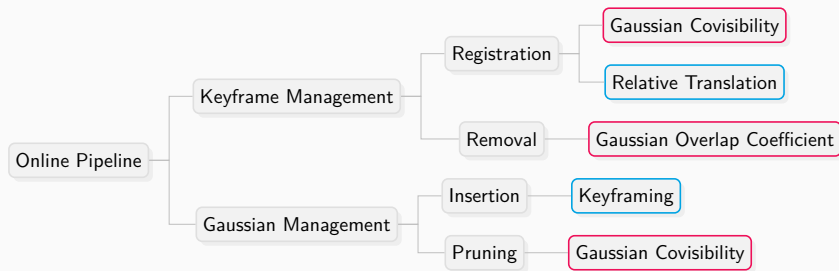


Keyframe Management:

key method trick convention

(arXiv, 2016) DSO: Direct Sparse Odometry

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM



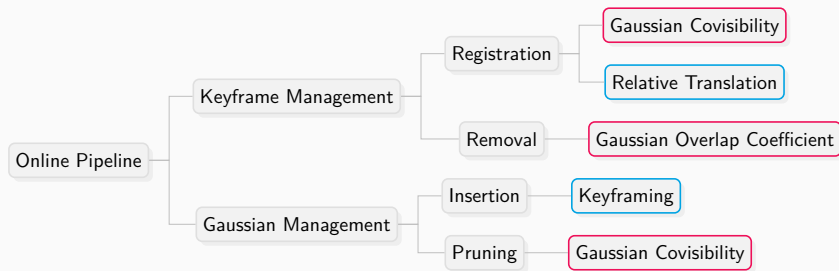
Keyframe Management:

- Classic strategies, e.g. covisibility & overlap, from DSO [5].

key method trick convention

(arXiv, 2016) DSO: Direct Sparse Odometry

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM



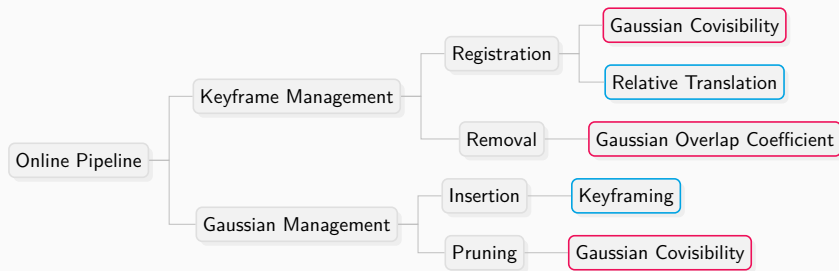
Keyframe Management:

- 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

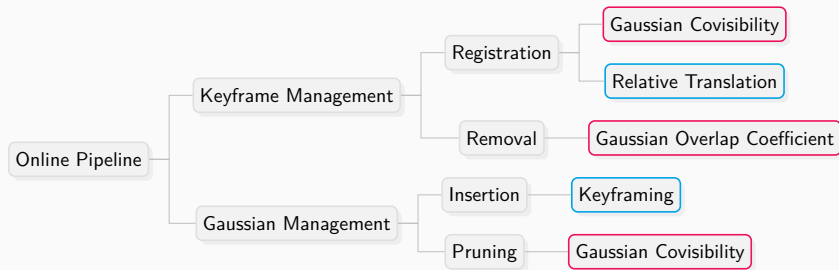


Gaussian Management:

key method trick convention

(arXiv, 2016) DSO: Direct Sparse Odometry

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM



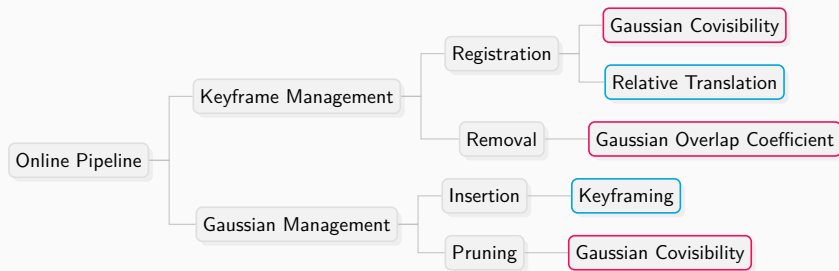
Gaussian Management:

- **Insertion:** triggered by **keyframing**, followed by **Gaussian initialization**.

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?

- 1 What is keyframing or keyframe management?
 - A strategy of selecting and utilizing a crucial subset of frames.

- 1 What is keyframing or keyframe management?
 - A strategy of selecting and utilizing a crucial subset of frames.
- 2 Why do we need keyframing?

- 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.

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/...)

- 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?

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**.

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**...

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.

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

$\subset \mathcal{G}$, **visible** Gaussians from frame j

In practice, $\tau_1 = 0.95$.

(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM

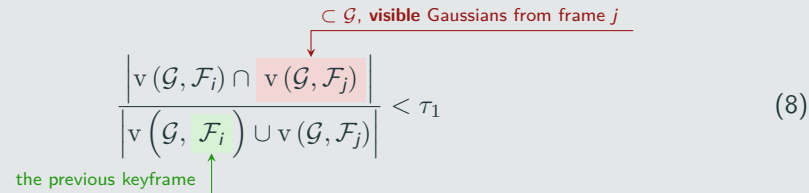
If any of the following conditions is true...

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:

- $v(\mathcal{G}, \mathcal{F}_i)$ is labeled "the previous keyframe" (green arrow).
- $v(\mathcal{G}, \mathcal{F}_j)$ is labeled "the current frame" (blue arrow).
- The intersection $v(\mathcal{G}, \mathcal{F}_i) \cap v(\mathcal{G}, \mathcal{F}_j)$ is labeled " $\subset \mathcal{G}$, visible Gaussians from frame j " (red arrow).

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.

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

$\in \mathbb{R}^3$, 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**...

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”?

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

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

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**.

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.

- 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.

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.

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 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.

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.

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

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

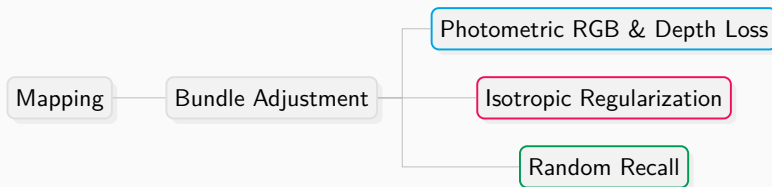
For “just” inserted Gaussians but unobserved by “some other” keyframes, are pruned out.

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

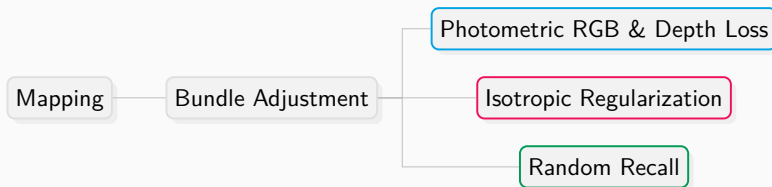


key method

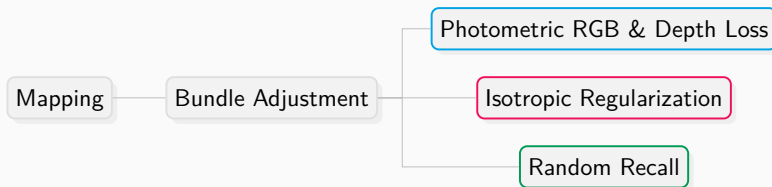
trick

convention

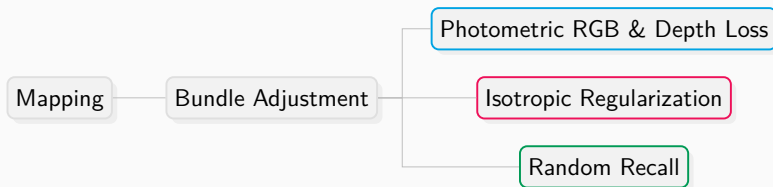
(CVPR Highlight, 2024) MonoGS: Gaussian Splatting SLAM



- Why do we need mapping in **3DGS** SLAM?



- Why do we need mapping in **3DGS** SLAM?
 - **Local**: Optimize newly inserted 3D Gaussians.



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

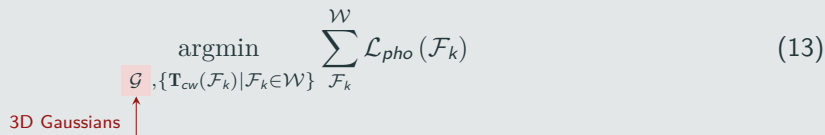
Bundle Adjustment

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

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



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

camera poses of keyframes in the sliding window

- Why do we need “isotropic regularization”?

- Why do we need “isotropic regularization”?
 - **Observation:** isotropic Gaussians behave better than anisotropic.

- 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**.

- 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**.

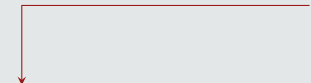
Isotropic Regularization

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

- 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

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


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

- 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

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

$\in \mathbb{N}$, 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

The Overall Optimization for Mapping

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

Appendix

- [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. xxix–xxxi).