

StaRGS-SLAM: Stable and Robust Gaussian SLAM via Coverage-Aware Fixed-Topology Initialization

Anonymous SPAR-3D submission

Paper ID 6

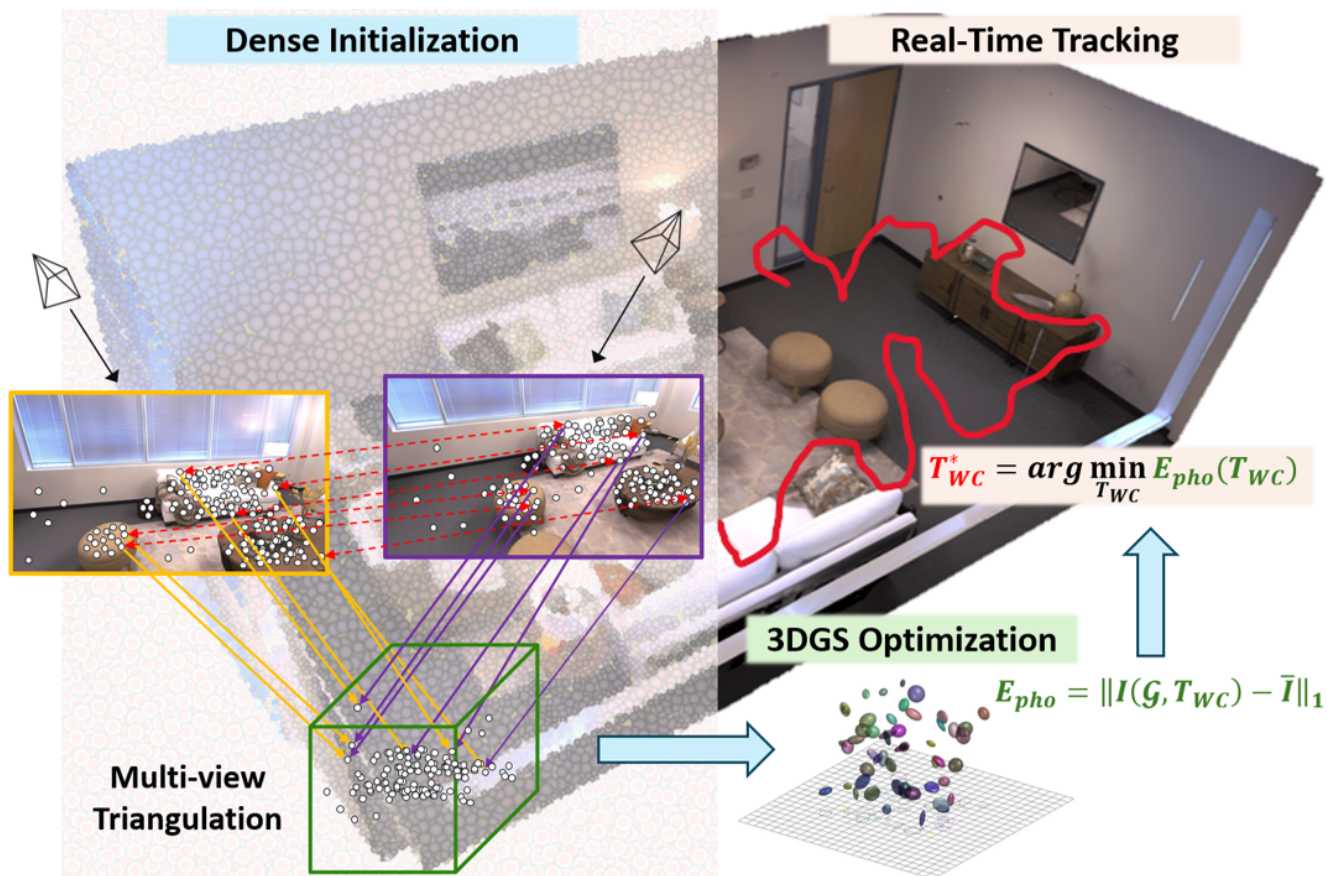


Figure 1. Overview of the proposed StaRGS-SLAM pipeline. The system integrates dense feature matching and multi-view triangulation for one-shot Gaussian initialization, followed by differentiable 3DGS optimization and real-time tracking.

Abstract

001 Reliable 3D scene modeling requires a mapping process
 002 that remains stable when early geometry is incomplete or
 003 visually ambiguous. We present StaRGS-SLAM, a Gaussian
 004 splatting SLAM system that improves robustness and
 005 efficiency through a training-free correspondence-guided
 006 Gaussian initialization scheme. Instead of relying on
 007 residual-driven densification during optimization, StaRGS-
 008 SLAM constructs a geometry-aware prior before iterative
 009 updates. From keyframes, it extracts dense correspondences
 010 from DINOv3 features, suppresses unreliable matches with

confidence-aware inlier filtering, and triangulates the fil-
 tered observations to produce a well-covered Gaussian
 representation in a single initialization stage. This de-
 sign provides stronger geometric support for early map-
 ping, reduces sensitivity to unstable optimization behav-
 ior, and shortens convergence time by about 20%. Exper-
 iments on TUM RGB-D and Replica show that StaRGS-
 SLAM achieves competitive or superior localization and re-
 construction performance compared with recent Gaussian-
 based and point-based SLAM methods, while maintaining
 real-time mapping throughput of up to 925 FPS.

011
012
013
014
015
016
017
018
019
020
021

022 **1. Introduction**

023 Recent advances in 3D Gaussian Splatting (3DGS) have
024 enabled high-quality view synthesis and real-time map-
025 ping. However, most pipelines still rely on residual-driven
026 densification, where Gaussians are iteratively spawned and
027 merged as errors are detected. This causes non-stationary
028 objectives, unstable convergence, and sensitivity to texture-
029 rich or cluttered regions due to delayed coverage and un-
030 even geometry.

031 We take a different approach by initializing from a com-
032 plete and well-distributed Gaussian set rather than growing
033 it incrementally. Using Dense Feature Matching (DFM),
034 we obtain confidence-weighted correspondences within a
035 short keyframe window, triangulate them into structure-
036 aware matches, and instantiate the Gaussian set before opti-
037 mization begins. Subsequent updates refine means, covari-
038 ances, opacities, and colors while keeping topology fixed,
039 resulting in stable and stationary optimization with strong
040 spatial support even in high-frequency regions.

041 Integrated into monocular SLAM, this single-step seed-
042 ing shortens time to usable maps, stabilizes early pose and
043 shape estimation, and removes the need for additional net-
044 works or losses. It directly replaces densification with a
045 brief feature-matching pass at keyframes while leaving the
046 rest of the pipeline unchanged. An overview of the StaRGS-
047 SLAM framework, including the initialization pipeline, is
048 shown in Figure 1. Our contributions are summarized be-
049 low.

- 050 • **Single-Step Dense Initialization.** A one-shot triangula-
051 tion replaces residual-driven densification within the stan-
052 dard GS-SLAM pipeline, enabling stationary optimization
053 and 20% faster convergence.
- 054 • **Improved Localization Accuracy.** Confidence-weighted
055 correspondences stabilize early pose estimation, reducing
056 drift by over 30%.
- 057 • **Lightweight and Efficient Mapping.** Spatially balanced
058 Gaussians lower computation and memory, achieving 20%
059 higher rendering throughput in real time.
- 060 • **High-Fidelity Reconstruction.** Dense seeding enhances
061 early coverage and geometric consistency, yielding about
062 20% better reconstruction accuracy and completeness.

063 **2. Related Work**

064 **3D Gaussian Splatting and Densification.** 3D Gaus-
065 sian splatting (3DGS) enables real-time view synthesis with
066 anisotropic splats and a visibility-aware renderer, yet most
067 systems rely on residual-driven densification [8]. We in-
068 stead seed a fixed topology from dense multi-view corre-
069 spondences and then refine only splat parameters.

070 **Gaussian Splats for SLAM.** SLAM systems mapping with
071 Gaussians include GS-SLAM, MonoGS, SplaTAM, and
072 Gauss-SLAM [7, 11, 20]. They rely on densification, caus-

ing early non-stationarity. Our training-free dense seed re- 073
moves this stage and drops into MonoGS with minimal 074
modifications. 075

Differentiable Rendering and Real-Time SLAM. Photo- 076
SLAM, GLORIE-SLAM, and Point-SLAM couple differ- 077
entiable rendering with pose optimization for fast updates 078
via analytic/lightweight gradients [4, 12, 23]. We retain this 079
in a Gaussian renderer and use a stationary initialization so 080
early steps do not change topology. 081

Feature Matching and Dense Correspondence. Super- 082
Point/SuperGlue remain strong baselines [1, 13]. Detector- 083
free transformers (LoFTR, LightGlue) extend coverage in 084
low-texture regions, and dense matchers (DKM) provide 085
broad two-view coverage with confidence for geometry [2, 086
9, 17]. We aggregate confidence-weighted dense correspon- 087
dences over a short keyframe window into an explicit Gaus- 088
sian seed. 089

Initialization and Training-Free Priors. SfM and multi- 090
view stereo stabilize early optimization via geometric priors 091
and learned depth [14, 22]. In Gaussian splatting, sched- 092
ules and regularizers typically retain densification. Our 093
correspondence-to-Gaussian initialization follows training- 094
free priors and yields a stationary objective from the start. 095

096 **3. Method**097 **3.1. Gaussian Splatting Representation**

098 We map the scene with a set of anisotropic Gaussians $\mathcal{G} =$ 099
 $\{G_i\}$. Each G_i carries optical properties, a color vector 100
 c_i and an opacity $\alpha_i \in [0, 1]$, and geometric properties, a 101
mean $\mu_i^W \in \mathbb{R}^3$ and a symmetric positive definite covariance 102
 $\Sigma_i^W \in \mathbb{R}^{3 \times 3}$ expressed in world coordinates. For brevity we 103
describe color as a single vector and later allow spherical 104
harmonics for view dependence. 105

Let $T_{WC} = [R \mid t]$ be the world-to-camera pose of the 106
current view and let $\pi(\cdot)$ be the calibrated perspective pro- 107
jection. A 3D Gaussian $N(\mu_i^W, \Sigma_i^W)$ induces a 2D Gaus- 108
sian on the image plane through first-order linearization of 109
the projection around μ_i^W . The projected mean and covari- 110
ance are

$$\mu_i^I = \pi(T_{WC} \mu_i^W), \quad \Sigma_i^I = J_i R \Sigma_i^W R^\top J_i^\top, \quad (1) \quad 111$$

where $J_i = \frac{\partial \pi(RX+t)}{\partial X} \Big|_{X=\mu_i^W}$ is the Jacobian of the projec- 112
tion at μ_i^W . This relates the ellipsoid in 3D to an ellipse on 113
the sensor. 114

Rendering is performed by rasterizing Gaussians instead 115
of ray marching. For pixel p we collect the front-to-back 116
ordered set $N(p)$ of screen-space Gaussians whose 2D 117
footprints overlap p . The pixel color is obtained by 118
 α -compositing 119

$$C_p = \sum_{i \in N(p)} c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (2) \quad 120$$

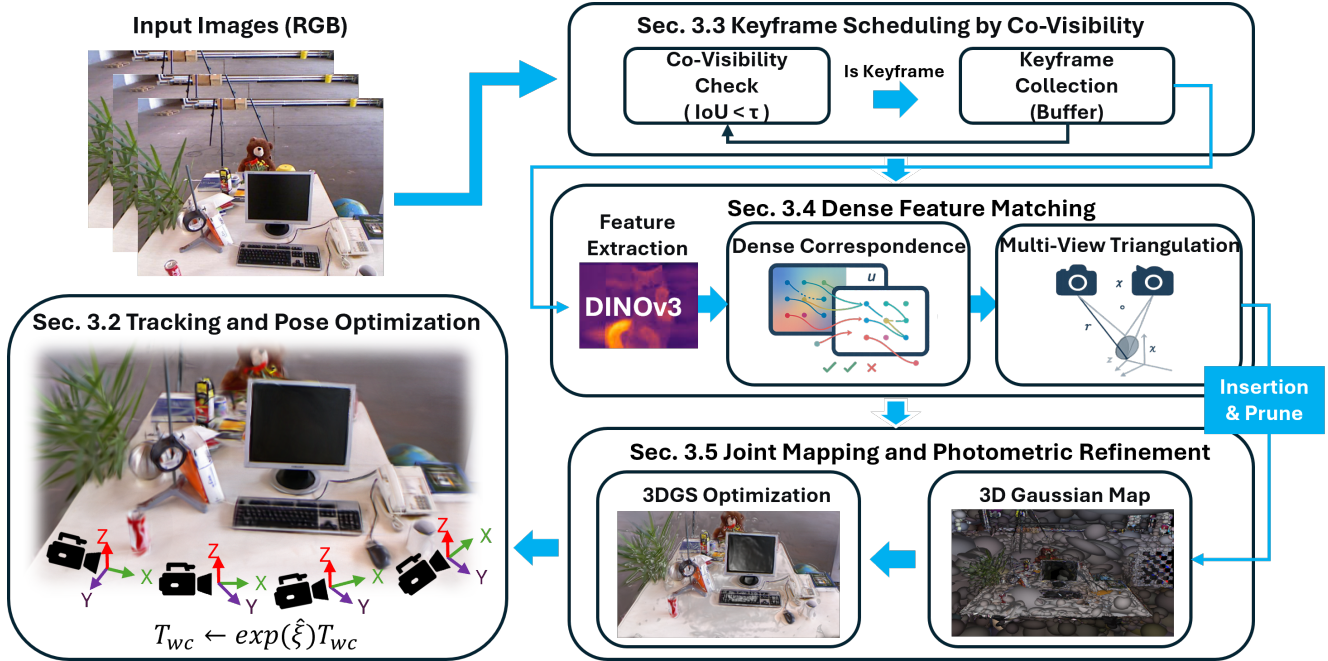


Figure 2. Detailed StaRGS-SLAM pipeline. Each keyframe triggers dense multi-view triangulation that yields a one-shot Gaussian initialization, subsequently refined through joint tracking and mapping within a differentiable 3DGS renderer using analytic SE(3) Jacobians.

121 where $\alpha_i \in [0, 1]$ denotes the *screen-space* opacity at pixel
 122 p , obtained by modulating the primitive opacity with the
 123 value of the 2D Gaussian density at p using the param-
 124 eters (μ_i^I, Σ_i^I) (the dependence on p is omitted for brevity).
 125 Empty space does not contribute because the renderer iter-
 126 ates over primitives that actually cover the pixel.

127 Equations (1) and (2) are differentiable in color, opac-
 128 ity, mean, covariance, and pose. Gradients flow through
 129 the splat weights and the projection Jacobian, which allows
 130 first-order optimizers to refine both optical and geometric
 131 parameters until the rendered image matches the observa-
 132 tion with high fidelity. In subsequent sections we will ini-
 133 tialize $\{\mu_i^W, \Sigma_i^W, \alpha_i, c_i\}$ densely in a single-step and then
 134 refine them under this differentiable renderer.

135 3.2. Tracking and Camera Pose Optimization

136 3.2.1. Objective and Per-Frame Update

137 For each incoming frame, we estimate the camera pose by
 138 minimizing an image-domain objective under the 3DGS
 139 renderer in Figure 2. Let $I(\mathcal{G}, T_{WC}) = \mathcal{S}(\mathcal{G}, T_{WC})$ be
 140 the rendered image and \bar{I} the observation. The photometric
 141 residual is

$$142 E_{\text{pho}} = \| I(\mathcal{G}, T_{WC}) - \bar{I} \|_1, \quad (3)$$

143 augmented with an affine brightness model to absorb ex-
 144 posure changes, i.e., we jointly estimate gain and bias and
 145 substitute $gI(\cdot) + b$ into Eq. (3). We optimize $E = E_{\text{pho}}$.

146 Pixels with low screen-space opacity or low image gradi-
 147 ent are downweighted to reduce the influence of textureless
 148 regions. In practice we perform tens of gradient steps per
 149 frame to reach a stable update.

150 3.2.2. Alpha Compositing for Color

151 Rendering is carried out by rasterizing Gaussians in screen
 152 space and compositing them front to back. For a pixel p ,
 153 let $N(p)$ be the set of overlapping Gaussians sorted from
 154 near to far. The color follows the standard α -compositing
 155 in Eq. (2), which naturally handles occlusion via transmit-
 156 tance $\prod_{j < i} (1 - \alpha_j)$. No depth map is produced or used in
 157 our tracking objective.

158 3.2.3. Minimal Pose Jacobians on SE(3)

159 We update the world-to-camera pose by a left-multiplicative
 160 twist SE(3),

$$161 T_{WC} \leftarrow \exp(\hat{\xi}) T_{WC}, \quad (4)$$

162 where we differentiate the objective with respect to ξ in
 163 minimal coordinates. Let μ^W be a 3D Gaussian mean in
 164 world coordinates and $\mu^C = R\mu^W + t$ its camera-space
 165 position for $T_{WC} = [R | t]$. The 3D point Jacobian with
 166 respect to the pose twist is the standard 3×6 form

$$167 \frac{\partial \mu^C}{\partial \xi} = [I \quad -[\mu^C]_{\times}], \quad (5)$$

168 where $[\cdot]_{\times}$ is the skew-symmetric matrix. With calibrated
 169 projection π , the image-plane mean $\mu^I = \pi(\mu^C)$ has Jaco-

170 bian

$$171 \quad \frac{\partial \mu^I}{\partial \xi} = J_\pi(\mu^C) [I \quad -[\mu^C]_\times], \quad (6)$$

172 where J_π is the 2×3 projection Jacobian evaluated at μ^C .
173 The screen-space covariance Σ^I from Eq. (1) depends on
174 both the projection Jacobian and the rotation, using the
175 chain rule,

$$176 \quad \frac{\partial \Sigma^I}{\partial \xi} = \frac{\partial \Sigma^I}{\partial J} \frac{\partial J}{\partial \mu^C} \frac{\partial \mu^C}{\partial \xi} + \frac{\partial \Sigma^I}{\partial R} \frac{\partial R}{\partial \xi}, \quad (7)$$

177 with $\partial R/\partial \xi$ obtained from the Lie algebra relation $\delta R \approx$
178 $[\delta\omega]_\times R$ for an infinitesimal rotation $\delta\omega$. These analytic Ja-
179 cobians remove the overhead of generic autodiff and match
180 the degrees of freedom of the pose, which is essential for
181 fast and stable tracking under a tight per-frame budget.

182 3.2.4. Optimization Solver and Weighting Scheme

183 We minimize the photometric objective in Eq. (3) using a
184 first-order optimizer with a cosine learning rate schedule,
185 and apply a robust penalty to per-pixel residuals. The per-
186 pixel weight combines exposure correction, edge aware-
187 ness, and visibility (via screen-space opacity) so that in-
188 formative regions dominate the update. Because the 3DGS
189 renderer and the pose Jacobians are fully analytic, gradients
190 propagate through Eq. (2) and Eq. (6) without resorting to
191 expensive automatic differentiation.

192 3.3. Keyframe Scheduling by Co-Visibility

193 Given the last accepted keyframe I_{k^*} , we measure co-
194 visibility between the current frame I_k and I_{k^*} by the
195 intersection-over-union of visible Gaussians

$$196 \quad \text{IoU}(I_k, I_{k^*}) = \frac{|V(I_k) \cap V(I_{k^*})|}{|V(I_k) \cup V(I_{k^*})|}, \quad (8)$$

197 where $V(I)$ collects Gaussians whose screen-space opacity
198 exceeds a small threshold on a sufficient fraction of pixels.
199 A new keyframe is created when $\text{IoU}(I_k, I_{k^*})$ is less than
200 τ and the inter-view parallax is above a bound. Accepted
201 keyframes are stored in a bounded buffer \mathcal{B} that provides
202 neighbours for multi-view initialization.

203 3.4. Dense Feature Matching

204 We extract dense visual descriptors using DINOv3 [15],
205 which provide semantically consistent features across
206 views. These descriptors are used to establish multi-view
207 dense correspondences, replacing the residual-driven den-
208 sification process in GS-SLAM. A confidence-aware inlier
209 classifier is then applied to filter unreliable matches, ensur-
210 ing stable multi-view geometry. Finally, a one-shot triangulation
211 is performed to initialize a uniformly distributed set
212 of 3D Gaussian seeds.

Dense Correspondence. Let I_r be the current keyframe
and let $\mathcal{N}_r \subset \mathcal{B}$ be K neighbours selected by pose prox-
imity and parallax. A dense matcher outputs, for each pair
(r, n) with $n \in \mathcal{N}_r$, a displacement field $\mathbf{u}_{r \rightarrow n}(p)$ on I_r
and a confidence map $\kappa_{r \rightarrow n}(p) \in [0, 1]$. A correspondence
is represented as the pixel pair

$$(p, p + \mathbf{u}_{r \rightarrow n}(p)), \quad (9)$$

and low-confidence matches are filtered by a symmetric
epipolar test and spatial blue-noise thinning. We aggregate
per-view confidence for each retained reference pixel by

$$\bar{\kappa}(p) = \frac{1}{|\mathcal{N}_r|} \sum_{n \in \mathcal{N}_r} \kappa_{r \rightarrow n}(p). \quad (10)$$

Multi-view Triangulation. For each retained pixel p and
neighbour $n \in \mathcal{N}_r$, we solve a two-view linear triangulation.
Let $P_r, P_n \in \mathbb{R}^{3 \times 4}$ be the camera projection matrices
and $\tilde{x}_r, \tilde{x}_n \in \mathbb{P}^2$ the homogeneous pixel coordinates. We
solve

$$A = \begin{bmatrix} \tilde{x}_r^x P_r^{3\top} - P_r^{1\top} \\ \tilde{x}_r^y P_r^{3\top} - P_r^{2\top} \\ \tilde{x}_n^x P_n^{3\top} - P_n^{1\top} \\ \tilde{x}_n^y P_n^{3\top} - P_n^{2\top} \end{bmatrix}, \quad \tilde{X} = \arg \min_{\|\tilde{X}\|=1} \|A\tilde{X}\|_2, \quad (11)$$

then $\hat{X} = \tilde{X}_{1:3}/\tilde{X}_4$ (obtained as the right singular vector of
 A associated with the smallest singular value). Among all
neighbours we keep the hypothesis with the lowest repro-
jection error, breaking ties by larger baseline angle. Can-
didates with small parallax or large reprojection error are
rejected.

Gaussian Parameter Initialization. Each valid triangulation
spawns one Gaussian G_i with world mean

$$\mu_i^W = \hat{X}(p). \quad (12)$$

Construct a local orthonormal frame $U_i = [\mathbf{t}_1, \mathbf{t}_2, \mathbf{v}]$, where
 \mathbf{v} is the surface normal estimated by a plane fit over neigh-
bouring triangulated points, and $\mathbf{t}_1, \mathbf{t}_2$ span the tangent
plane. Initialize the covariance as an anisotropic ellipsoid
aligned with this frame

$$\Sigma_i^W = U_i \text{diag}(s_\perp^2, s_\perp^2, s_\parallel^2) U_i^\top, \quad (12)$$

where s_\perp is obtained by back-projecting a one-pixel im-
age uncertainty via the projection Jacobian at the reference
view, and s_\parallel is set larger to encode depth uncertainty that
increases when the baseline angle is small or the triangulation
residual is high. The color c_i is the median of bilinearly
sampled RGB values across supporting views after applying
the exposure parameters estimated by tracking. The initial
opacity α_i is a monotone mapping of the aggregated cor-
respondence confidence $\bar{\kappa}(p)$ so that unreliable candidates
remain visually weak at insertion. Finally, Gaussians are
subsampling to maintain uniform spatial coverage before be-
ing inserted into the map.

257 3.5. Joint Mapping and Photometric Refinement

258 At each accepted keyframe, we first perform the single-step
259 dense initialization (Sec. 3.4) to generate Gaussian seeds
260 from multi-view correspondences. The surviving seeds are
261 immediately inserted into \mathcal{G} without any densification stage.
262 Each newly inserted Gaussian participates in mapping right
263 away.

264 **Insertion, Lightweight Merging, and Pruning.** Each
265 Gaussian G_i tracks its observation count m_i , cumulative
266 visibility v_i , and an exponential moving average of its
267 screen-space footprint. To keep memory bounded and re-
268 move unstable outliers, we apply a lightweight periodic
269 cleanup and prune splats that violate

$$270 \quad \begin{aligned} m_i &< m_{\min}, & v_i &< v_{\min}, \\ \text{tr}(\Sigma_i^W) &> \sigma_{\max}^2, & \alpha_i &< \alpha_{\min}. \end{aligned}$$

271 Neighbouring Gaussians with highly overlapping foot-
272 prints and similar colors are merged, retaining a visibility-
273 weighted mean of color and covariance to avoid over-
274 population.

275 **Sliding-Window Photometric Refinement.** Let \mathcal{W} denote
276 a window around the latest keyframe. We jointly refine
277 $\{T_{WC}\}_{I \in \mathcal{W}}$ and $\{c_i, \alpha_i, \mu_i^W, \Sigma_i^W\}$ by minimizing

$$278 \quad \mathcal{L} = \sum_{I \in \mathcal{W}} \lambda_{\text{pho}} \|g_I \mathcal{S}(\mathcal{G}, T_{WC}) + b_I - \bar{I}\|_1 + \mathcal{R}, \quad (13)$$

279 where g_I, b_I compensate exposure changes. The regularizer

$$280 \quad \begin{aligned} \mathcal{R} = & \lambda_{\text{iso}} \sum_i \left\| \Sigma_i^W - \frac{\text{tr}(\Sigma_i^W)}{3} I_3 \right\|_F \\ & + \lambda_{\alpha} \sum_i \psi(\alpha_i) + \lambda_{\mu} \sum_i \left\| \mu_i^W - \bar{\mu}_i^W \right\|_2^2 \end{aligned} \quad (14)$$

281 discourages needle-shaped ellipsoids, avoids degenerate
282 transmittance, and stabilizes early iterations via an EMA
283 anchor $\bar{\mu}_i^W$. We alternate pose-only updates and full map
284 updates with robust per-pixel weights. Gradients propagate
285 through the analytic α -compositing in Eq. (2) and the pose
286 Jacobians in Eq. (6).

287 3.6. System Schedule and Computational Profile

288 Each incoming frame is tracked for K_t gradient steps using
289 the photometric objective in Eq. (3). When the co-visibility
290 test Eq. (8) accepts a keyframe, we select K neighbours
291 from \mathcal{B} and execute the single-step dense initialization of
292 Sec. 3.4 (dense correspondence, weighted multi-view trian-
293 gulation, and parameter initialization) in one pass. The re-
294 sulting Gaussians are immediately inserted into \mathcal{G} , followed
295 by K_m mapping iterations over the current window \mathcal{W} op-
296 timizing Eq. (13), and a lightweight cleanup as described in
297 Sec. 3.5. Replacing iterative densification with this single-
298 step initialization reduces the drift of newly added paramet-
299 ers and lowers the number of mapping iterations required to

reach the same photometric fidelity, improving wall-clock
throughput without changing the objective or renderer.

4. Experiments

4.1. Experimental Setup

300 **Datasets.** We evaluate on TUM RGB-D and Replica. TUM
301 RGB-D is evaluated in both monocular and RGB-D set-
302 tings. Replica is employed for photometric map evalua-
303 tion on *room0-2* and *office0-4*, matching the splits used in
304 our tables to ensure comparability of rendering metrics and
305 throughput.

306 **Implementation.** Gaussian rasterization and gradients are
307 implemented in CUDA, and the remaining pipeline is in
308 PyTorch. Mixed precision is enabled where beneficial.
309 Tracking runs in real time, while mapping executes asyn-
310 chronously within a bounded local window. Non-standard
311 hyperparameters (learning rate schedule, window sizes,
312 keyframe and culling thresholds) are provided in the sup-
313 plementary.

314 **Evaluation Metrics.** Tracking accuracy uses RMSE of Ab-
315 solute Trajectory Error (ATE) on keyframes. Photometric
316 quality adopts PSNR [3], SSIM [19], and LPIPS [24]. Re-
317 construction quality reports *Acc.* [cm]↓, *Comp.* [cm]↓, and
318 *Comp.Ratio* (%)↑. Unless specified, we uniformly sample
319 50K surface points, set $\tau = 5$ cm, and average per scene.
320 Photometric metrics are computed on every fifth frame ex-
321 cluding keyframes. Reconstruction metrics use the same
322 sampling protocol. Each experiment is repeated three times,
323 and the mean results are reported.

324 **Baseline Methods.** We compare with iMAP [16],
325 NICE-SLAM [26], Vox-Fusion [21], ESLAM [6],
326 Point-SLAM [12], Co-SLAM [18], SplatTAM [7],
327 Gauss-SLAM [20], and MonoGS [11]. We also include
328 SNI-SLAM [25] for reconstruction and Photo-SLAM [4],
329 GLORIE-SLAM [23], RK-SLAM [10] for rendering.
330 RGB-D-only methods run in RGB-D, with monocular
331 results reported only when supported. Hyperparameters
332 follow official documented defaults on the same splits.

4.2. Training and Convergence Analysis

333 The system optimizes scene specific Gaussian parameters
334 and camera poses using the same optimizer, schedule, and
335 window size as MonoGS, with densification removed. Each
336 frame is tracked for K_t steps. Each accepted keyframe trig-
337 gers the single-step dense initialization followed by K_m
338 mapping steps with exposure compensation. Wall clock
339 time is measured on identical hardware and stopping cri-
340 teria. Results are summarized in Table 1. On TUM RGB-D
341 the average training time decreases from 14.8 to 12 minutes
342 while maintaining localization and rendering quality.

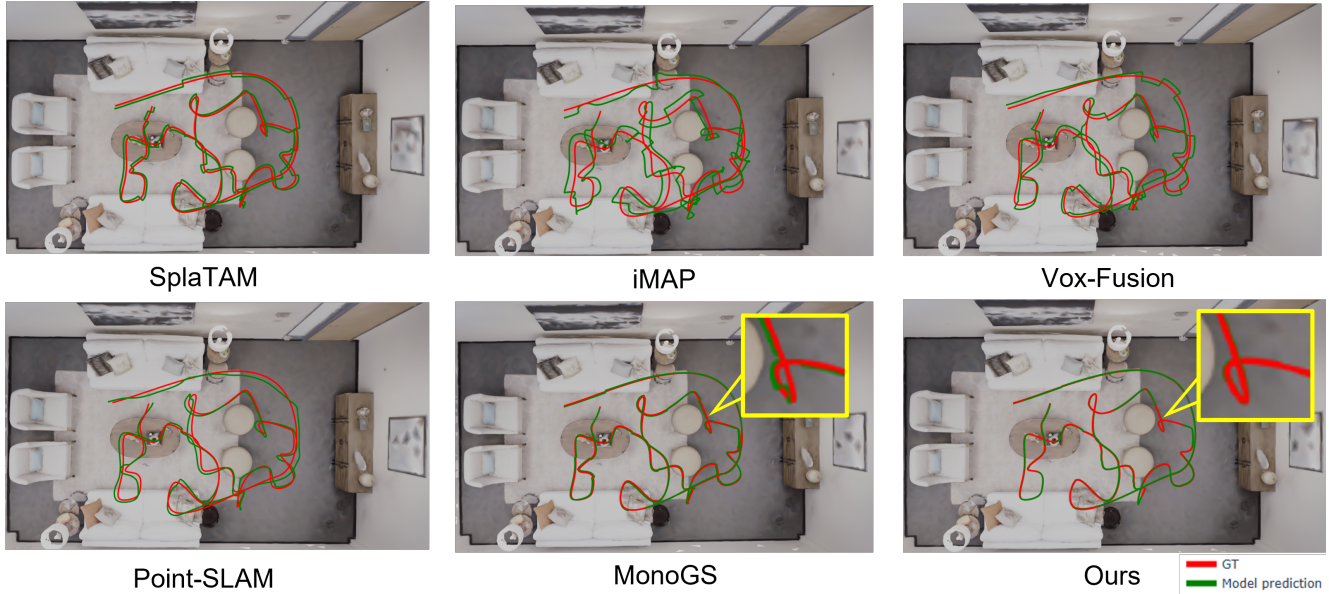


Figure 3. Trajectory comparison on a living-room scene. The red line indicates the ground-truth path and the green line shows the estimated trajectory. Our method aligns more closely with the ground-truth and exhibits fewer large drifts than previous systems.

Table 1. Optimization time (min) on TUM RGB-D sequences *fr1/desk*, *fr2/xyz*, and *fr3/office*.

Method	<i>fr1/desk</i>	<i>fr2/xyz</i>	<i>fr3/office</i>	Avg.
MonoGS [11]	6.4	20.6	17.5	14.8
Ours (RGB)	4.9	16.1	15.0	12.0

348

4.3. Localization Accuracy

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

We evaluate trajectory accuracy on Replica and TUM (Tables 2, 3). On Replica, the average ATE is **0.61 cm** across *r0-r2*, *o0-o4*, outperforming iMAP (2.58), NICE-SLAM (1.07), Vox-Fusion (3.09), and ESLAM (0.90 cm), while remaining competitive with Point-SLAM and MonoGS. On TUM RGB-D, the average ATE is **1.02 cm** on *fr1/desk*, *fr2/xyz*, and *fr3/office*, achieving better results than MonoGS and surpassing the same baselines. Figure 3 shows a trajectory comparison on a living-room scene from the Replica dataset, where the red line indicates the ground-truth (GT) path and the green line represents the estimated trajectory. Among existing methods, SplaTAM, iMAP, Vox-Fusion, and Point-SLAM exhibit large localization drift, with evident deviations from the GT path. MonoGS and our method perform significantly better. In this visualization, the red line is drawn above the green line, so greater overlap, where the red line covers the green one, intuitively reflects smaller localization error. Our method achieves a higher overlap ratio, indicating closer adherence to the GT trajectory and better pose consistency.

The zoom-in comparison further shows smoother alignment and reduced drift compared with MonoGS, demonstrating stronger robustness in long-term tracking and loop-closure maintenance.

369

370

371

372

Table 2. Camera tracking results on the Replica dataset under the RGB-D setting. Reported values denote RMSE of ATE across *room0-2* and *office0-4*.

Method	room0	room1	room2	office0	office1	office2	office3	office4	Avg.
iMAP [16]	3.12	2.54	2.31	1.69	1.03	3.99	4.05	1.93	2.58
NICE-SLAM [26]	0.97	1.31	1.07	0.88	1.00	1.06	1.10	1.13	1.07
Vox-Fusion [21]	1.37	4.70	1.47	8.48	2.04	2.58	1.11	2.94	3.09
ESLAM [6]	0.71	0.70	0.52	0.57	0.55	0.58	0.72	0.63	0.63
Point-SLAM [12]	0.61	0.41	0.37	0.38	0.48	0.54	0.69	0.72	0.53
MonoGS [11]	0.62	0.62	0.77	0.44	0.52	0.23	0.62	2.25	0.76
Ours (RGB)	0.45	0.51	0.53	0.52	0.78	1.03	0.45	0.63	0.61

Table 3. Camera tracking results on the TUM RGB-D dataset. Values denote RMSE of ATE over *fr1/desk*, *fr2/xyz*, and *fr3/office*.

Method	<i>fr1/desk</i>	<i>fr2/xyz</i>	<i>fr3/office</i>	Avg.
iMAP [16]	4.90	2.00	5.80	4.23
NICE-SLAM [26]	4.26	6.19	3.87	4.77
DI-Fusion [5]	4.40	2.00	5.80	4.07
Vox-Fusion [21]	3.52	1.49	26.01	10.34
ESLAM [6]	2.47	1.11	2.42	2.00
Co-SLAM [18]	2.40	1.70	2.40	2.17
Point-SLAM [12]	4.34	1.31	3.48	3.04
MonoGS [11]	1.50	1.44	1.49	1.47
Ours (RGB)	1.02	0.98	1.05	1.02



Figure 4. Rendering results on the TUM dataset. The proposed keyframe-triggered single-step initialization produces sharper edges, fewer transparency artifacts, and more consistent colors than residual-driven densification.

373 4.4. Rendering Quality and Throughput

374 We report fidelity and throughput in Table 4 and 5.
 375 On Replica, our initializer averages **925 FPS**, exceed-
 376 ing MonoGS (769 FPS) while maintaining competitive
 377 PSNR [3], SSIM [19], and LPIPS [24] across *room0–2* and
 378 *office0–4*. On TUM, the system runs in real time (2.5–3.2
 379 FPS) with PSNR/SSIM comparable to SplaTAM, Photo-
 380 SLAM, and GLORIE-SLAM, and low LPIPS. The through-
 381 put gain arises from the keyframe-triggered single-step den-

Table 4. Rendering quality results on the Replica dataset across room0–2 and office0–4.

Method (FPS)	Metric	room0	room1	room2	office0	office1	office2	office3	office4	Avg.
NICE-SLAM [26] (6.54)	PSNR[dB]↑	22.12	22.47	24.52	29.07	30.34	19.66	22.23	24.94	24.42
	SSIM↑	0.689	0.757	0.814	0.874	0.868	0.797	0.801	0.856	0.809
	LPIPS↓	0.330	0.271	0.208	0.229	0.181	0.235	0.209	0.198	0.233
Vox-Fusion [21] (2.17)	PSNR[dB]↑	22.39	22.36	23.92	27.79	29.83	20.33	23.47	25.21	24.41
	SSIM↑	0.683	0.751	0.798	0.857	0.876	0.794	0.803	0.847	0.801
	LPIPS↓	0.303	0.269	0.234	0.241	0.184	0.243	0.213	0.199	0.236
Point-SLAM [12] (1.33)	PSNR[dB]↑	32.40	34.08	35.50	38.26	39.16	33.98	33.48	33.49	35.17
	SSIM↑	0.974	0.977	0.979	0.982	0.986	0.962	0.960	0.979	0.975
	LPIPS↓	0.113	0.116	0.111	0.100	0.118	0.156	0.132	0.142	0.124
Co-SLAM [18]	PSNR[dB]↑	28.88	28.51	29.37	35.44	34.63	26.56	28.79	32.16	28.42
	SSIM↑	0.892	0.843	0.851	0.854	0.826	0.814	0.866	0.856	0.837
	LPIPS↓	0.213	0.205	0.215	0.177	0.161	0.172	0.163	0.176	0.185
SplaTAM [7]	PSNR[dB]↑	32.49	33.72	34.65	38.29	39.04	31.91	30.05	31.83	30.98
	SSIM↑	0.975	0.970	0.980	0.982	0.982	0.965	0.952	0.949	0.953
	LPIPS↓	0.072	0.096	0.078	0.086	0.093	0.100	0.110	0.150	0.179
Gauss-SLAM [20]	PSNR[dB]↑	29.57	31.61	33.46	38.39	39.62	32.91	33.62	34.26	30.90
	SSIM↑	0.944	0.952	0.973	0.985	0.991	0.974	0.982	0.979	0.972
	LPIPS↓	0.197	0.184	0.148	0.099	0.097	0.158	0.123	0.138	0.229
MonoGS [11] (769)	PSNR[dB]↑	34.83	36.43	37.49	39.95	42.09	36.24	36.70	36.07	37.50
	SSIM↑	0.954	0.959	0.965	0.971	0.974	0.964	0.963	0.957	0.960
	LPIPS↓	0.068	0.076	0.075	0.072	0.055	0.078	0.065	0.099	0.070
Ours (RGB) (925)	PSNR[dB]↑	35.95	33.55	32.45	34.45	35.45	34.87	34.02	35.85	34.57
	SSIM↑	0.852	0.945	0.985	0.952	0.925	0.952	0.855	0.961	0.928
	LPIPS↓	0.085	0.092	0.112	0.088	0.096	0.078	0.101	0.096	0.093

se initialization, which fixes Gaussian topology upfront and
 removes residual-driven densification, reducing per-frame
 cost. Qualitative results in Figure 4 show sharper edges,
 fewer transparency artifacts, and more consistent colors
 than residual-driven baselines.

Table 5. Rendering quality results on the TUM RGB-D dataset.

Method (FPS)	Metric	fr1/desk	fr2/xyz	fr3/office	Avg.
Point-SLAM [12]	PSNR[dB]↑	13.79	17.62	18.29	16.57
	SSIM↑	0.625	0.710	0.749	0.695
	LPIPS↓	0.545	0.584	0.452	0.527
Photo-SLAM [4]	PSNR[dB]↑	20.97	21.07	19.59	20.54
	SSIM↑	0.740	0.730	0.690	0.720
	LPIPS↓	0.230	0.170	0.240	0.213
MonoGS [11]	PSNR[dB]↑	19.67	16.17	20.63	18.82
	SSIM↑	0.730	0.720	0.770	0.740
	LPIPS↓	0.330	0.310	0.340	0.327
GLORIE-SLAM [23]	PSNR[dB]↑	20.26	25.62	21.21	22.36
	SSIM↑	0.790	0.720	0.720	0.743
	LPIPS↓	0.310	0.090	0.320	0.240
SplaTAM [7]	PSNR[dB]↑	21.49	25.06	21.17	22.57
	SSIM↑	0.839	0.950	0.861	0.883
	LPIPS↓	0.255	0.099	0.221	0.192
RK-SLAM [10]	PSNR[dB]↑	22.31	22.47	20.67	21.82
	SSIM↑	0.741	0.729	0.710	0.727
	LPIPS↓	0.254	0.220	0.251	0.242
Ours (RGB)	PSNR[dB]↑	23.11	24.85	23.59	23.85
	SSIM↑	0.853	0.896	0.801	0.850
	LPIPS↓	0.232	0.198	0.219	0.216

387 **4.5. Reconstruction Fidelity**

388 Geometric fidelity is evaluated using accuracy, complete-
 389 ness, and completeness ratio in Table 6, computed on
 390 aligned point clouds under standard thresholds. Our method
 391 attains **1.537 cm** accuracy and **1.477 cm** completeness with
 392 a **97.843%** completeness ratio. Relative to SNI-SLAM, ac-
 393 curacy improves by 20.9% and completeness by 13.2%, to-
 394 gether with a 1.22-point gain in completeness ratio. The
 395 margins over ESLAM and Vox-Fusion are larger, including
 396 a 42% reduction in completeness error against Vox-Fusion.
 397 The improvements are consistent across scenes with thin
 398 structures and clutter, where coverage gaps and over-
 399 regularization commonly inflate geometric error. Qualita-
 400 tive inspection shows reduced truncation at object bound-
 401 aries, cleaner reconstruction of high-frequency edges, and
 402 better recovery of small appendages. We attribute these
 403 outcomes to anisotropic Gaussian primitives with visibility-
 404 aware α -compositing, which sharpen depth gradients and
 405 limit color bleeding, and to a bounded, keyframe-related
 406 optimization that preserves spatial coverage without topol-
 407 ogy changes. By keeping the Gaussian set fixed after dense
 408 seeding, the optimization remains stationary and avoids
 409 late-map artifacts, which stabilizes surface inference and
 410 suppresses oversmoothing during refinement.

Table 6. Reconstruction results on the Replica dataset. Lower is better for Acc./Comp., higher for Comp.Ratio.

Methods	Reconstruction		
	Acc. [cm] ↓	Comp. [cm] ↓	Comp.Ratio (%) ↑
iMAP [16]	3.624	4.934	80.515
NICE-SLAM [26]	2.373	2.645	91.137
Vox-Fusion [21]	1.882	2.563	90.936
Co-SLAM [18]	2.104	2.082	93.435
ESLAM [6]	2.082	1.754	96.427
SNI-SLAM [25]	1.942	1.702	96.624
Ours	1.537	1.477	97.843

411 **4.6. Ablation Study**

412 **Effect of Dense Initialization.** Consistent rendering
 413 gains on TUM RGB, with higher PSNR/SSIM and lower
 414 LPIPS/RMSE across all sequences. On *fr1.desk*, *fr2.xyz*,
 415 and *fr3.office*, PSNR improves to 23.11, 24.85, and 23.59
 416 with SSIM gains and LPIPS/RMSE drops (Table 7). Dis-
 417 tributed Gaussian seeds, whose multi-view triangulation
 418 stabilizes mapping under larger motion and maintains cov-
 419 erage in low-parallax segments. This yields faster conver-
 420 gence and fewer artifacts on thin structures and cluttered
 421 regions. Without dense initialization, residual driven den-
 422 sification converges slowly, exhibits early spatial inconsis-
 423 tency, and tends to over-smooth before adequate coverage
 424 is established.

Table 7. Impact of DFM on the TUM RGB-D dataset.

	Method	PSNR ↑	SSIM ↑	LPIPS ↓	RMSE ↓
fr1_desk	w/o DFM	19.67	0.73	0.33	1.5
	Ours	23.11	0.853	0.232	1.02
fr2_xyz	w/o DFM	16.17	0.72	0.31	1.44
	Ours	24.85	0.896	0.198	0.98
fr3_office	w/o DFM	20.63	0.77	0.34	1.49
	Ours	23.59	0.801	0.219	1.05

425 **Effect of Gaussian Count per Keyframe on Tracking.**

426 We vary the number of newly triangulated 3D points per
 427 keyframe, with each verified point instantiated as a Gaus-
 428 sian primitive, so the abscissa in Figure 5 corresponds to the
 429 count of Gaussians. Increasing the budget from 200 to 1000
 430 points reduces the tracking error sharply, reaching about
 431 **0.7cm** at 1000. Beyond this regime the curve plateaus and
 432 improvements are marginal, approaching roughly 0.6 cm at
 433 2000. We therefore adopt 1000 points per keyframe as the
 434 default trade-off between accuracy, memory, and runtime.

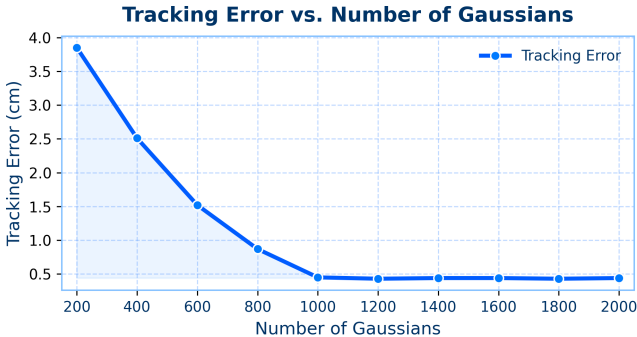


Figure 5. Tracking error versus Gaussian count. Error decreases rapidly with denser seeding and plateaus near 1000 Gaussians, indicating diminishing returns beyond this density.

435 **5. Conclusion**

436 We presented StaRGS-SLAM, a Gaussian splatting SLAM
 437 framework that replaces residual-driven densification with a
 438 keyframe-based correspondence-guided initialization strat-
 439 egy. Through dense multi-view matching and triangula-
 440 tion, the method forms a geometry-aware Gaussian prior
 441 that strengthens early mapping and improves the reliability
 442 of subsequent optimization under analytic SE(3) Jacobians.
 443 Experiments on Replica and TUM RGB-D show that this
 444 design reduces computational overhead while preserving
 445 strong localization accuracy and rendering quality. These
 446 results suggest that a more structured initialization stage can
 447 make 3D scene modeling more stable and dependable in vi-
 448 sually challenging conditions, while remaining compatible
 449 with existing Gaussian splatting SLAM pipelines.

450

References

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

[1] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2018. 2

[2] Johan Edstedt, Ioannis Athanasiadis, Mårten Wadenbäck, and Michael Felsberg. DKM: Dense kernelized feature matching for geometry estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2

[3] Alain Hore and Djemel Ziou. Image quality metrics: PSNR vs. SSIM. In *Proceedings of the IEEE International Conference on Image Processing (ICIP)*, 2010. 5, 7

[4] Huajian Huang, Longwei Li, Hui Cheng, and Sai-Kit Yeung. Photo-SLAM: Real-time simultaneous localization and photorealistic mapping for monocular, stereo, and RGB-D cameras. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2, 5, 7

[5] Jiahui Huang, Shi-Sheng Huang, Haoxuan Song, and Shi-Min Hu. DI-Fusion: Online implicit 3D reconstruction with deep priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 6

[6] Mohammad Mahdi Johari, Camilla Carta, and François Fleuret. ESLAM: Efficient dense SLAM system based on hybrid representation of signed distance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 5, 6, 8

[7] Nikhil Keetha, Jay Karhade, Krishna Murthy Jatavallabhula, Gengshan Yang, Sebastian Scherer, Deva Ramanan, and Jonathon Luiten. Splatam: Splat, track & map 3D Gaussians for dense RGB-D SLAM. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2, 5, 7

[8] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics*, 2023. 2

[9] Philipp Lindenberger, Paul-Edouard Sarlin, and Marc Pollefeys. Lightglue: Local feature matching at light speed. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 2

[10] Xiasheng Ma, Ci Song, Yimin Ji, and Shanlin Zhong. Related keyframe optimization Gaussian-simultaneous localization and mapping: A 3D Gaussian Splatting-based simultaneous localization and mapping with related keyframe optimization. *Applied Sciences*, 2025. 5, 7

[11] Hidenobu Matsuki, Riku Murai, Paul Kelly, and Andrew J. Davison. Gaussian Splatting SLAM. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2, 5, 6, 7

[12] Erik Sandström, Yue Li, Luc Van Gool, and Martin R. Oswald. Point-SLAM: Dense neural point cloud-based SLAM. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 2, 5, 6, 7

[13] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[14] Johannes L. Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2

[15] Oriane Siméoni, Huy V. Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, Francisco Massa, Daniel Haziza, Luca Wehrstedt, Jianyuan Wang, Timothée Darcet, Théo Moutakanni, Leonel Sentana, Claire Roberts, Andrea Vedaldi, Jamie Tolan, John Brandt, Camille Couprie, Julien Mairal, Hervé Jégou, Patrick Labatut, and Piotr Bojanowski. DINOv3. *arXiv preprint arXiv:2508.10104*, 2025. 4

[16] Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew J. Davison. imap: Implicit mapping and positioning in real-time. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 5, 6, 8

[17] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. LoFTR: Detector-free local feature matching with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[18] Hengyi Wang, Jingwen Wang, and Lourdes Agapito. Co-SLAM: Joint coordinate and sparse parametric encodings for neural real-time SLAM. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 5, 6, 7, 8

[19] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004. 5, 7

[20] Chao Yan, Zirui Wang, Zhiqiang Li, Wei Gao, Hao Wang, Guofeng Zhang, Hujun Bao, and Xiaowei Zhou. GS-SLAM: Dense visual SLAM with 3D Gaussian Splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2, 5, 7

[21] Xingrui Yang, Hai Li, Hongjia Zhai, Yuhang Ming, Yuqian Liu, and Guofeng Zhang. Vox-fusion: Dense tracking and mapping with voxel-based neural implicit representation. In *Proceedings of the IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, pages 499–507, 2022. 5, 6, 7, 8

[22] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. MVSNet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. 2

[23] Ganlin Zhang, Erik Sandström, Youmin Zhang, Manthan Patel, Luc Van Gool, and Martin R. Oswald. GLORIE-SLAM: Globally optimized RGB-only implicit encoding point cloud SLAM. *arXiv preprint arXiv:2403.19549*, 2024. 2, 5, 7

[24] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of

- 563 deep features as a perceptual metric. In *Proceedings of*
564 *the IEEE/CVF Conference on Computer Vision and Pattern*
565 *Recognition (CVPR)*, 2018. 5, 7
- [25] Siting Zhu, Guangming Wang, Hermann Blum, Jiuming Liu,
566 Liang Song, Marc Pollefeys, and Hesheng Wang. SNI-
567 SLAM: Semantic neural implicit SLAM. In *Proceedings of*
568 *the IEEE/CVF Conference on Computer Vision and Pattern*
569 *Recognition (CVPR)*, pages 21167–21177, 2024. 5, 8
- [26] Zihan Zhu, Songyou Peng, Viktor Larsson, Weiwei Xu, Hu-
571 jun Bao, Zhaopeng Cui, Martin R. Oswald, and Marc Polle-
572 feys. NICE-SLAM: Neural implicit scalable encoding for
573 SLAM. In *Proceedings of the IEEE/CVF Conference on*
574 *Computer Vision and Pattern Recognition (CVPR)*, 2022. 5,
575 6, 7, 8