

RGS-SLAM: Robust Gaussian Splatting SLAM with One-Shot Dense Initialization

Supplementary Material

001 1. Reproducibility and Code Release

002 All components of the system are implemented in PyTorch
003 with custom CUDA kernels for rasterization and analytic
004 Jacobians. The repository provides the full SLAM pipeline,
005 configuration files, scripts for evaluation on TUM RGB-D
006 and Replica, and instructions for reproducing all quantita-
007 tive tables and qualitative renderings in the main paper. The
008 training and optimization settings match those described in
009 Sections 3 and 4 of the main paper. An anonymized version
010 of the codebase is available at:

011 [https://anonymous.4open.science/r/RGS-
012 SLAM](https://anonymous.4open.science/r/RGS-SLAM)

013 2. Discussion

014 2.1. Limitations

015 RGS-SLAM still exhibits several practical limitations de-
016 spite the gains over residual-driven densification. The cur-
017 rent evaluation focuses on indoor static scenes in Replica
018 and TUM RGB-D, so the robustness of the one-shot Gaus-
019 sian initialization in highly dynamic scenes, or under severe
020 rolling shutter remains unclear. The dense correspondences
021 are derived from a single pretrained DINOv3 backbone to-
022 gether with a confidence aware inlier classifier, which can
023 degrade under drastic appearance changes, strong motion
024 blur, or camera viewpoints far outside the training distri-
025 bution, and in these regimes the triangulated seeds may be-
026 come biased or incomplete. The fixed topology after each
027 keyframe initialization improves stability but can leave per-
028 sistent coverage gaps when large textureless surfaces, fine
029 specular structures, or objects with weak visual support
030 never accumulate enough consistent matches, which con-
031 strains the reconstruction quality. The current implemen-
032 tation also assumes a calibrated pinhole camera and synchro-
033 nized RGB-D streams, without exploiting lidar cues or in-
034 erial measurements that are available on many robotic plat-
035 forms. Finally, the system still requires a GPU with mod-
036 erate memory to sustain high frame rates, and the practi-
037 cal impact of memory budgets, long-term map growth, and
038 large-scale loop closure has not yet been characterized on
039 resource constrained devices.

040 2.2. Societal Impact

041 RGS-SLAM advances camera based dense SLAM with a
042 training-free one-shot Gaussian initialization that stabilizes
043 optimization and improves throughput, which can benefit
044 robotics, extended reality, and digital twin systems through

045 more reliable mapping and safer physical interaction. At
046 the same time, dense reconstruction of indoor environments
047 raises privacy risks whenever RGB or RGB-D streams are
048 captured and stored without informed consent, since metri-
049 cally consistent maps can reveal layouts, personal belong-
050 ings, and usage patterns. The method is training-free and
051 easily integrated into existing SLAM stacks, which lowers
052 the barrier for large-scale deployment and makes responsi-
053 ble use dependent on appropriate safeguards such as trans-
054 parent data handling, limited retention of raw sensor data,
055 access control to stored maps, and a preference for local
056 processing. Our experiments rely on public benchmarks
057 without personally identifiable information and the code re-
058 lease is intended for reproducible research, although future
059 work should pair technical advances in robust mapping with
060 privacy aware data design and interdisciplinary guidelines
061 for ethical deployment.

062 2.3. Training Details

063 **Optimization of Gaussian maps.** The optimization set-
064 tings used by our Gaussian SLAM system are summarized
065 in Table 1. For both initialization and online mapping we
066 use the Adam optimizer with a shared parameterization for
067 geometry and appearance. The base learning rate is set to
068 2.0×10^{-3} for color and opacity parameters and 1.0×10^{-3}
069 for positions and rotations, with $(\beta_1, \beta_2) = (0.9, 0.999)$. A
070 cosine decay schedule is applied within each optimization
071 window so that early iterations focus on rapid geometry re-
072 finement while later iterations stabilize the map. Initializa-
073 tion uses a window of five frames and runs for 1050 itera-
074 tions before the system starts live tracking. During online
075 operation every incoming frame is refined for 30 tracking
076 iterations and each accepted keyframe-triggers 60 mapping
077 iterations over the same local window. The loss combines
078 an L_1 term and an SSIM term weighted by $\lambda_{dssim} = 0.2$.
079 Dynamic density control uses the same thresholds across all
080 experiments and adopts the opacity culling and densifica-
081 tion settings listed in Table 1. Mixed precision training and
082 gradient norm clipping are enabled to keep the per frame
083 compute budget and memory stable.

084 **Dense initialization baseline.** As an additional baseline we
085 adopt a dense correspondence based initialization strategy
086 that shares the same differentiable Gaussian representation
087 as our system. The training configuration is summarized
088 in Table 2. The model is optimized for 30000 iterations
089 with Adam and separate learning rates for spatial and ap-
090pearance parameters. Position updates follow a decayed

Table 1. Optimization configuration for the Gaussian SLAM.

Config	Value
Optimizer	Adam
Adam betas	$\beta_1 = 0.9, \beta_2 = 0.999$
Base LR (color, opacity)	2.0×10^{-3}
Base LR (position, rotation)	1.0×10^{-3}
Learning rate schedule	Cosine decay
Loss	$L_1 + \lambda_{\text{dssim}} L_{\text{SSIM}}$
SSIM weight λ_{dssim}	0.2
Initialization iterations	1050
Tracking iterations per frame	30
Mapping iterations per keyframe	60
Local window size	5 frames
Densification interval	100 iterations
Opacity reset interval	3000 iterations
Opacity culling threshold	0.05
Gradient clipping max norm	1.0
Precision	Mixed FP16 / FP32

schedule from 1.6×10^{-4} to 1.6×10^{-6} , while feature, opacity, scaling, and rotation parameters use constant rates that match the public configuration. The reconstruction loss combines an ℓ_1 term with a structural similarity component with $\lambda_{\text{dssim}} = 0.2$. Densification is enabled from iteration 500 to 15000 with an interval of 100 iterations and a gradient based threshold of 2.0×10^{-4} . All experiments use degree three spherical harmonics, batch size 64, and a fixed background color without random perturbations on a CUDA device.

3. Implementation Details

Datasets. We evaluate on the TUM RGB-D and Replica benchmarks, consistent with Section 4 of the main paper. TUM RGB-D is used in both monocular and RGB-D configurations, following standard practice in visual SLAM. Replica is employed for photometric map evaluation and trajectory accuracy on the eight standard indoor scenes *room0* to *room2* and *office0* to *office4*. These splits match those used in the main tables so that rendering metrics, throughput, and localization error remain directly comparable across methods.

Evaluation Metrics. Camera tracking accuracy is measured using the root mean square error of the Absolute Trajectory Error (ATE) computed on keyframes. Photometric map quality is evaluated with three rendering metrics: PSNR, SSIM, and LPIPS. Geometric reconstruction quality is assessed with accuracy (Acc. [cm]), completeness (Comp. [cm]), and completeness ratio (Comp.Ratio [%]). Accuracy is defined as the mean nearest-neighbour distance from reconstructed points to the ground-truth surface.

Table 2. Training configuration of the dense initialization baseline.

Config	Value
Total iterations	30000
Optimizer	Adam
Spherical harmonics degree	3
Batch size	64
Position LR (init → final)	$1.6 \times 10^{-4} \rightarrow 1.6 \times 10^{-6}$
Feature learning rate	2.5×10^{-3}
Opacity learning rate	2.5×10^{-2}
Scaling learning rate	5.0×10^{-3}
Rotation learning rate	1.0×10^{-3}
SSIM weight λ_{dssim}	0.2
Dense correspondence ratio	0.01
Densification interval	100 iterations
Densification range	iter. 500 to 15000
Densification gradient threshold	2.0×10^{-4}
Opacity reset iteration	30000
Exposure LR (init → final)	$1.0 \times 10^{-2} \rightarrow 1.0 \times 10^{-4}$
Random background	disabled
Logging train camera index	50
Logging test camera index	10
Dataset resolution	original
White background	false
Data device	CUDA

Completeness is the mean nearest-neighbour distance from ground-truth surface samples to the reconstruction. The completeness ratio is the proportion of ground-truth samples within a distance threshold τ from the reconstruction. Unless stated otherwise, we uniformly sample 50K points on each surface, set $\tau = 5$ cm, and average per-scene scores over the benchmark sequences. For photometric rendering metrics, we evaluate every fifth frame and exclude keyframes to avoid bias toward training views. Reconstruction metrics are computed with the same surface-sampling protocol. Each experiment is repeated three times on identical hardware and stopping criteria, and all tables report the mean scores. In every table, the best result is typeset in bold and the second best is underlined.

3.1. Compare Model Settings

We compare RGS-SLAM with a set of representative dense SLAM pipelines that cover implicit volumes, voxel grids, point clouds, and Gaussian splats under the same scene types and benchmarks. NICE-SLAM, Co-SLAM, VoxFusion, DI-Fusion, and SNI-SLAM operate on RGB-D input and maintain volumetric implicit or voxel based representations that are optimized per scene. iMAP and Point-SLAM map monocular or RGB-D streams to neural fields or point clouds with scene specific training and joint pose refinement. SplaTAM, Gauss-SLAM, MonoGS, and RK-

146 SLAM adopt 3D Gaussian splatting and couple a Gaussian
 147 renderer with pose and appearance optimization, while
 148 Photo-SLAM and GLORIE-SLAM combine differentiable
 149 rendering with explicit point or mesh structures.

150 3.2. Computing Resource Configuration

151 All experiments are run on a workstation equipped with
 152 two NVIDIA L40 GPUs and an Intel Xeon Platinum
 153 8362 CPU at 2.80 GHz. Time-critical components, including
 154 3D Gaussian rasterization and gradient computation,
 155 are implemented in CUDA, while the remaining SLAM
 156 pipeline is implemented in PyTorch. Mixed precision is en-
 157 abled for rendering and backpropagation whenever this im-
 158 proves throughput without degrading stability. The tracking
 159 loop operates in real-time, and mapping is executed asyn-
 160 chronously within a bounded local window so that latency
 161 remains stable as the map grows.

162 4. Methodology Details

163 4.1. Gaussian Splatting Representation

164 Each Gaussian G_i is represented by a compact tuple con-
 165 taining the world space mean $\mu_i^W \in \mathbb{R}^3$, an opacity param-
 166 eter $\alpha_i \in [0, 1]$, a set of view dependent color coefficients,
 167 and a covariance parameterization. In the underlying model
 168 the covariance appears as a full matrix Σ_i^W , while in the im-
 169 plementation it is encoded as a rotation matrix $R_i \in \text{SO}(3)$
 170 and three axis aligned scales $s_i \in \mathbb{R}_+^3$ such that

$$171 \Sigma_i^W = R_i \text{diag}(s_i^2) R_i^\top. \quad (1)$$

172 The rotation is stored as a unit quaternion and the scales are
 173 optimized in logarithmic space, which guarantees positive
 174 definiteness under gradient updates.

175 Colors are represented by second order spherical har-
 176 monics in camera space. For a viewing direction $\mathbf{v} \in \mathbb{S}^2$,
 177 the color of G_i is

$$178 C_i(\mathbf{v}) = \sum_{\ell=0}^2 \sum_{m=-\ell}^{\ell} \mathbf{c}_{i,\ell m} Y_\ell^m(\mathbf{v}), \quad (2)$$

179 where $\mathbf{c}_{i,\ell m} \in \mathbb{R}^3$ are learned RGB coefficients and Y_ℓ^m are
 180 real spherical harmonics.

181 Projection to the image plane uses the calibrated camera
 182 intrinsics together with the Gaussian projection model, and
 183 screen space compositing applies standard alpha compositing.
 184 All attributes are packed into contiguous GPU buffers
 185 and updated in place. This layout lets the renderer handle
 186 several million Gaussians without fragmentation and keeps
 187 memory access patterns coherent during both forward and
 188 backward passes.

189 4.2. Tracking and Camera Pose Optimization

190 Given the current map \mathcal{G} and a new RGB or RGB-D frame,
 191 pose tracking alternates between rendering and gradient

192 based refinement of the camera pose T_{WC} . We render
 193 a synthesized image $\hat{I}(x; \mathcal{G}, T_{WC})$ at the native resolution
 194 and apply an affine brightness model with gain g_I and bias
 195 b_I for each frame,

$$196 \tilde{I}_I(x; \mathcal{G}, T_{WC}) = g_I \hat{I}(x; \mathcal{G}, T_{WC}) + b_I. \quad (3)$$

197 The parameters (g_I, b_I) are estimated by a small least
 198 squares problem

$$199 (g_I, b_I) = \arg \min_{g, b} \sum_{x \in \Omega_I} \left(I(x) - g \hat{I}(x; \mathcal{G}, T_{WC}) - b \right)^2, \quad (4)$$

200 and the resulting solution is substituted into the photometric
 201 objective so that pose updates remain invariant to slow
 202 exposure drift.

203 The tracking loss can be written in the form

$$204 \mathcal{L}_{\text{track}}(T_{WC}) = \sum_{x \in \Omega_I} w_I(x) \|I(x) - \tilde{I}_I(x; \mathcal{G}, T_{WC})\|_1, \quad (5)$$

205 where $w_I(x)$ denotes a per pixel weight. In practice this
 206 weight is factored into opacity and gradient terms,

$$207 w_I(x) = w_\alpha(x) w_\nabla(x), \quad (6)$$

208 with

$$209 w_\alpha(x) = \text{clip}\left(\frac{\hat{\alpha}(x)}{\tau_\alpha}, 0, 1\right), \quad (7)$$

$$210 w_\nabla(x) = \text{clip}\left(\frac{\|\nabla I(x)\|_2}{\tau_\nabla}, 0, 1\right), \quad (8)$$

211 where $\hat{\alpha}(x)$ is the accumulated opacity at pixel x , τ_α and τ_∇
 212 are fixed thresholds, and $\text{clip}(z, a, b) = \min(\max(z, a), b)$.
 213 Pixels with low opacity or weak gradients therefore have
 214 reduced influence in the optimization.

215 The pose is updated in the minimal twist coordinates
 216 $\xi \in \mathbb{R}^6$ using a standard left multiplicative update rule on
 217 $\text{SE}(3)$. The Jacobians of the camera projection and the im-
 218 age formation model are implemented analytically and are
 219 reused across all Gaussians that share the same pose, which
 220 reduces both computation and memory traffic. The deriva-
 221 tive of the projected covariance with respect to pose is eval-
 222 uated by an explicit chain rule involving the projection Ja-
 223 cobian and the local rotation. This avoids generic automatic
 224 differentiation through the entire renderer.

225 In practice, between thirty and sixty gradient steps per
 226 frame are performed using the Adam optimizer with a co-
 227 sine learning rate schedule centered around 5×10^{-3} . This
 228 configuration yields stable tracking even when large parts
 229 of the image are textureless or underexposed.

230 4.3. Keyframe Scheduling by Co-Visibility

231 Keyframe scheduling uses a co-visibility based policy. For
 232 each incoming frame we maintain a binary visibility mask

that records, for every pixel, whether its accumulated screen space opacity exceeds a small threshold. Let V_a and V_b denote the sets of visible pixels in images I_a and I_b . The co-visibility score is defined as

$$\text{IoU}(I_a, I_b) = \frac{|V_a \cap V_b|}{|V_a \cup V_b|}. \quad (9)$$

The intersection and union are counted entirely on the GPU.

A frame I_k is promoted to a keyframe when the co-visibility $\text{IoU}(I_k, I_{k^*})$ with the last keyframe I_{k^*} falls below a user defined threshold τ and when the relative translation and rotation exceed small geometric bounds. These additional bounds avoid accepting nearly redundant viewpoints that would increase memory and computation without improving triangulation baselines.

Accepted keyframes are stored in a bounded buffer \mathcal{B} together with their poses. Between eight and twelve recent keyframes are kept, which is sufficient to form well-conditioned local triangulation baselines while keeping all multi view operations inexpensive. The same buffer also provides the neighbour set used in the mapping stage.

4.4. Dense Feature Matching and Triangulation

Dense matching and multi view initialization rely on DINoV3 features computed at a fixed feature resolution obtained by downsampling the RGB images. The descriptors are ℓ_2 normalized per pixel. For each reference keyframe I_r a neighbour set $\mathcal{N}_r \subset \mathcal{B}$ is selected based on pose proximity and parallax, using the current camera estimates.

Let $f_r(p)$ and $f_n(q)$ denote DINoV3 descriptors at pixels p and q in images I_r and I_n . For each neighbour $n \in \mathcal{N}_r$ and displacement $\mathbf{u}_{r \rightarrow n}(p)$, the raw descriptor similarity is

$$s_{r \rightarrow n}(p) = \langle f_r(p), f_n(p + \mathbf{u}_{r \rightarrow n}(p)) \rangle. \quad (10)$$

This similarity is combined with forward backward consistency and epipolar agreement into a scalar score

$$\tilde{\kappa}_{r \rightarrow n}(p) = w_{\text{sim}} s_{r \rightarrow n}(p) + w_{\text{fb}} \rho_{\text{fb}}(p) + w_{\text{epi}} \rho_{\text{epi}}(p), \quad (11)$$

where ρ_{fb} and ρ_{epi} quantify consistency of the forward backward displacement and the epipolar distance, and w_{sim} , w_{fb} , and w_{epi} are fixed weights. The confidence $\kappa_{r \rightarrow n}(p)$ is obtained by a piecewise linear mapping to $[0, 1]$,

$$\kappa_{r \rightarrow n}(p) = \text{clip}\left(\frac{\tilde{\kappa}_{r \rightarrow n}(p) - \gamma_0}{\gamma_1 - \gamma_0}, 0, 1\right), \quad (12)$$

with fixed thresholds γ_0 and γ_1 . This design keeps the inlier classifier training-free and dataset agnostic.

Before triangulation, correspondences are thinned in image space with a blue noise pattern in order to avoid redundant seeds in locally homogeneous regions. The aggregated confidence $\bar{\kappa}(p)$ is cached for each surviving refer-

ence pixel and encodes agreement across multiple neighbours. Linear triangulation uses a homogeneous linear system solved with a double precision singular value decomposition. Candidates with very small parallax or a reprojection error above a conservative threshold are discarded, and among hypotheses obtained from different neighbours the one with the smallest reprojection error and a sufficiently large baseline angle is kept.

4.5. Gaussian Initialization and Joint Mapping

Each valid triangulated point spawns a Gaussian G_i whose world mean is set to the reconstructed point $\mu_i^W = \hat{X}(p)$. A local orthonormal frame $U_i = [\mathbf{t}_1, \mathbf{t}_2, \mathbf{v}]$ is constructed by fitting a plane to triangulated neighbours inside a fixed radius around $\hat{X}(p)$, where \mathbf{v} is the normal and $\mathbf{t}_1, \mathbf{t}_2$ span the tangent plane in the neighbourhood. The covariance is initialized as an anisotropic ellipsoid aligned with this frame,

$$\Sigma_i^W = U_i \text{ diag}(s_\perp^2, s_\perp^2, s_\parallel^2) U_i^\top, \quad (13)$$

where the tangential scale s_\perp is obtained by back projecting a one pixel footprint through the projection Jacobian at the reference view and the axial scale s_\parallel is a calibrated function of the baseline angle and triangulation residual, which increases depth uncertainty in poorly conditioned configurations.

The initial color c_i is the median of bilinear RGB samples across all supporting views after applying the exposure parameters estimated during tracking. The initial opacity α_i is obtained as a monotone mapping of $\bar{\kappa}(p)$,

$$\alpha_i = \alpha_{\min} + (\alpha_{\max} - \alpha_{\min}) \bar{\kappa}(p), \quad (14)$$

with fixed bounds α_{\min} and α_{\max} in $(0, 1)$. Unreliable seeds with low aggregated confidence therefore enter the map as visually weak Gaussians and are easy to prune. Before insertion, Poisson disk subsampling in world space is applied to promote uniform coverage and to avoid excessive density in locally flat regions.

After insertion, the mapping module maintains an observation count m_i , a cumulative visibility score v_i , and an exponential moving average of the world position $\bar{\mu}_i^W$ for each Gaussian. The moving average is updated after each mapping step according to

$$\bar{\mu}_i^{W(t+1)} = (1 - \eta) \bar{\mu}_i^{W(t)} + \eta \mu_i^{W(t+1)}, \quad (15)$$

with a fixed smoothing factor $\eta \in (0, 1)$.

The mapping loss is evaluated over a sliding window \mathcal{W} around the latest keyframe with per frame brightness parameters g_I and b_I and per pixel weights $w_I(x)$ as in Eq. (6). A regularizer discourages highly elongated covariances, avoids degenerate transmittance, and anchors early updates to $\bar{\mu}_i^W$. The coefficients of this regularizer are kept fixed across all sequences and are not adapted per scene,

325 which keeps the behaviour of the optimizer comparable on
 326 TUM and Replica.

327 Pose only updates and full map updates are alternated,
 328 and robust per pixel weights are used so that outliers in the
 329 photometric residual have limited influence. A lightweight
 330 merging operation averages color and covariance for pairs
 331 of Gaussians that have almost identical screen space foot-
 332 prints and similar appearance. Pruning removes Gaussians
 333 that violate bounds on m_i , v_i , $\text{tr}(\Sigma_i^W)$, or α_i . These
 334 maintenance steps keep the representation compact and prevent
 335 numerical instabilities caused by extremely large or ex-
 336 tremely transparent splats.

337 **4.6. System Schedule and Computational Profile**

338 The runtime schedule separates tracking and mapping.
 339 Each incoming frame is tracked for K_t gradient steps using
 340 the photometric loss described above. This stage updates
 341 only the current pose and does not modify the map. If the
 342 co-visibility test does not accept the frame as a keyframe,
 343 the system immediately proceeds to the next image.

344 Let ρ denote the empirical fraction of frames that are se-
 345 lected as keyframes. The average number of gradient based
 346 optimization steps per frame is then approximately

$$347 N_{\text{steps}}^{\text{frame}} \approx K_t + \rho K_m, \quad (16)$$

348 where K_m is the number of joint refinement iterations exe-
 349 cuted after each keyframe.

350 When a keyframe is accepted, a set of neighbours from
 351 \mathcal{B} is selected and dense matching, confidence aggregation,
 352 triangulation, and Gaussian initialization are executed in a
 353 single GPU pass. The resulting Gaussians are inserted into
 354 \mathcal{G} and immediately participate in mapping. The system then
 355 runs K_m iterations of joint refinement over the window \mathcal{W}
 356 using the mapping loss and regularizer defined above, fol-
 357 lowed by a maintenance pass that performs merging and
 358 pruning.

359 This schedule replaces iterative residual-driven densi-
 360 fication with a one-shot dense seed and thereby reduces
 361 the early drift of newly added parameters. The number
 362 of mapping iterations required to reach a given photomet-
 363 ric fidelity decreases, which improves wall clock through-
 364 put while keeping the renderer and optimization objectives
 365 identical to those used in the main method description.

366 **5. Additional Experiments**

367 **5.1. Additional Rendering on Replica**

368 On the Replica dataset, extended qualitative comparisons in
 369 Figure 1 show that our keyframe-triggered single-step ini-
 370 tialization yields sharper object boundaries and more sta-
 371 ble shading than the residual-driven densification baseline
 372 across living room and office scenes. The reconstructed

373 views exhibit fewer transparency artifacts around thin struc-
 374 tures such as chair legs and table edges, and color transi-
 375 tions remain consistent across viewpoints, which confirms
 376 that the proposed initialization creates a well-conditioned
 377 Gaussian map for subsequent optimization. Challenging
 378 regions for Gaussian splatting, including large textureless
 379 walls and slanted ceilings, also show reduced blotchy arti-
 380 facts because the one-shot dense seed avoids early gaps in
 381 coverage. These qualitative trends agree with the quanti-
 382 tative gains in reconstruction metrics reported in the main
 383 paper and indicate that the initializer improves both con-
 384 vergence speed and the final visual fidelity of the radiance
 385 field.

386 **5.2. Camera Tracking on Replica Offices**

387 To assess tracking robustness, we visualize top view camera
 388 trajectories for two Replica office scenes in Figure 2 and 3.
 389 The proposed system maintains tight alignment with ground
 390 truth over long paths that include turns, loops, and revisits,
 391 and the strong overlap between the predicted and reference
 392 trajectories indicates that the jointly optimized poses and
 393 Gaussians provide accurate geometric constraints for down-
 394 stream mapping and loop closure. In the office0 sequence
 395 the path combines slow pans and rapid rotations around the
 396 desk area, yet our trajectory returns to previously visited
 397 regions without noticeable misalignment, while competitor
 398 methods accumulate drift near corners and walls. In the of-
 399 fice2 sequence the camera passes through narrow corridors
 400 before entering a wider workspace, and the estimated path
 401 from our method preserves the global layout without evi-
 402 dent shearing or scale distortion. These qualitative patterns
 403 are consistent with the Absolute Trajectory Error reported
 404 in the main paper and support the claim that dense Gaus-
 405 sian initialization yields a stable optimization landscape for
 406 pose refinement.

407 **5.3. Cluttered Desk Reconstruction**

408 In a cluttered desk sequence with strong self occlusion and
 409 fine scale objects such as cables, pencils, and plush toys,
 410 shown in Figure 4, we evaluate an additional reconstruc-
 411 tion. The left image shows the input RGB view and the
 412 right image illustrates the corresponding Gaussian map ren-
 413 dered from a novel viewpoint. The reconstruction preserves
 414 thin structures and surface boundaries while avoiding the
 415 truncation and over smoothing artifacts observed in resi-
 416 dual driven pipelines, which demonstrates that the proposed
 417 initialization and refinement strategy scales to scenes with
 418 complex object layouts and high frequency details. Fine el-
 419 ements such as tripod legs, monitor edges, and scattered sta-
 420 tionery remain clearly separated from the background even
 421 when objects move partially in and out of view, so the map
 422 integrates evidence from multiple viewpoints without dupli-
 423 cated Gaussians.

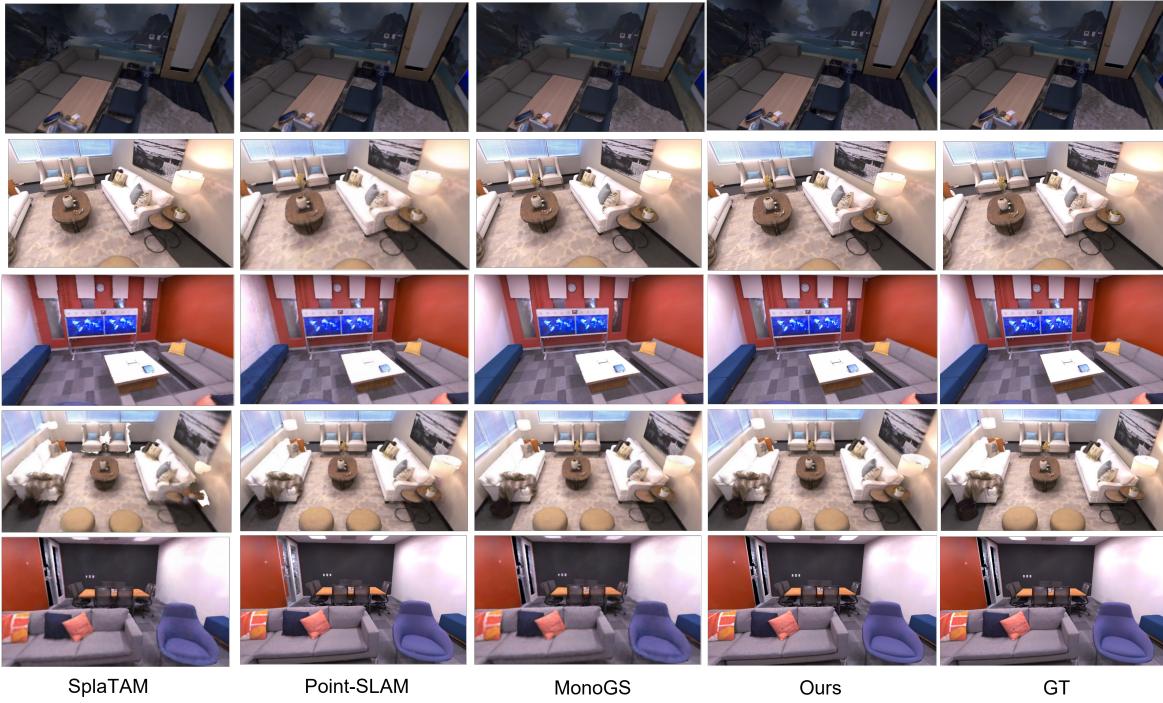


Figure 1. Rendering results on the Replica dataset. The proposed keyframe-triggered single-step initialization produces sharper edges, fewer transparency artifacts, and more consistent colors than residual-driven densification.

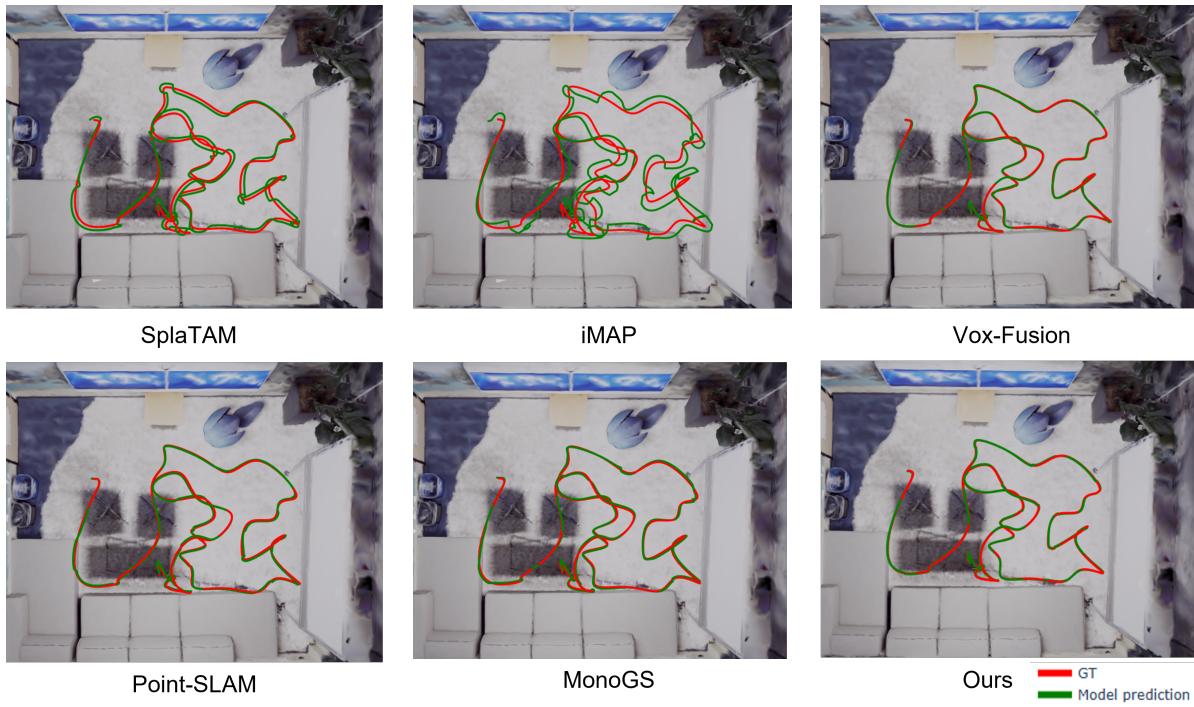


Figure 2. Tracking on Replica office0. Top-view trajectories, with ground truth in red and model predictions in green.

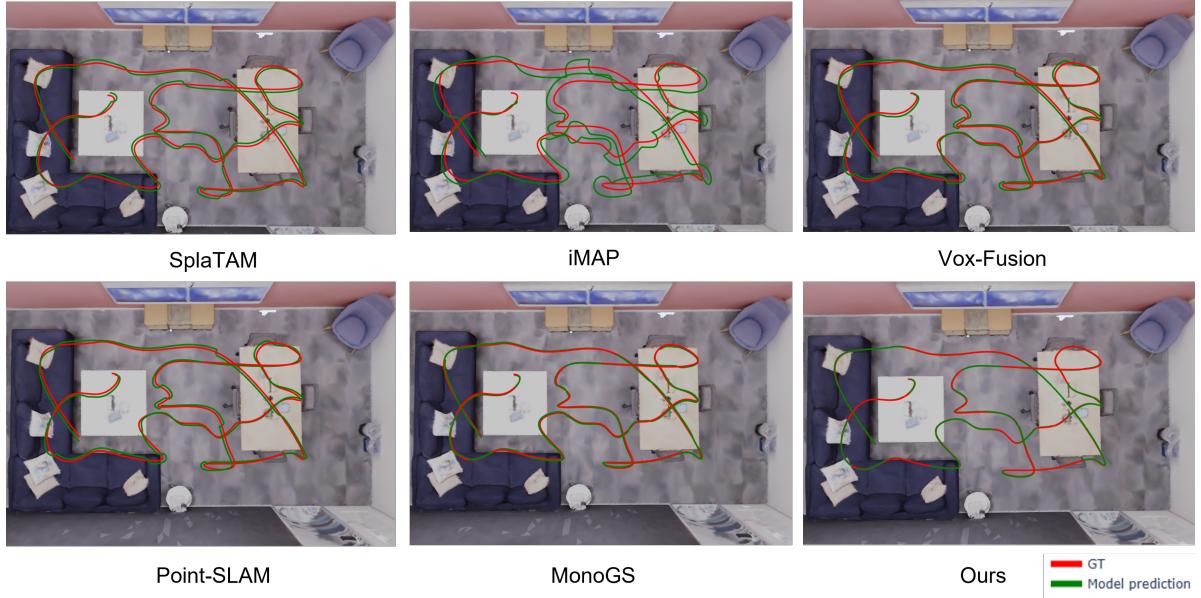


Figure 3. Tracking on Replica office2. Top-view trajectories, with ground truth in red and model predictions in green.



Figure 4. Qualitative reconstruction on a cluttered desk scene, where the left input RGB frame and the right Gaussian map view demonstrate dense coverage with preserved fine structures and reduced truncation artifacts.