

# SLAM-Former: Putting SLAM into One Transformer

Yijun Yuan   Zhuoguang Chen   Kenan Li   Weibang Wang   Hang Zhao  
IIIS, Tsinghua University

<https://tsinghua-mars-lab.github.io/SLAM-Former>

{yuanyj, hangzhao}@mail.tsinghua.edu.cn

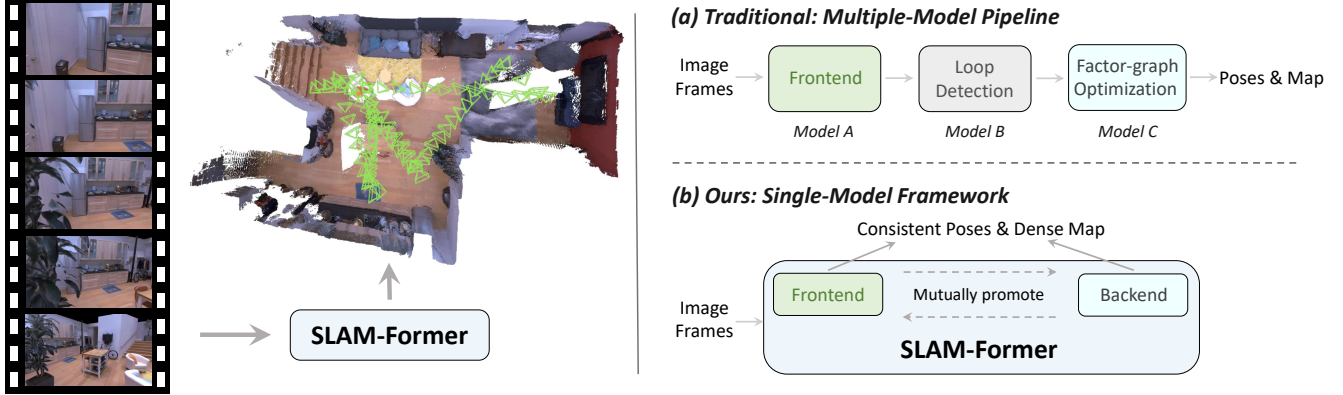


Figure 1. SLAM-Former is a unified Transformer for SLAM. Traditional SLAM employs a multi-model pipeline for frontend and backend tasks. In contrast, SLAM-Former integrates the full SLAM functionality within one transformer, achieving consistent poses and dense maps.

## Abstract

We present *SLAM-Former*, a novel neural approach that integrates full SLAM capabilities into a single transformer. Similar to traditional SLAM systems, *SLAM-Former* comprises both a frontend and a backend that operate in tandem. The frontend processes sequential monocular images in real-time for incremental mapping and tracking, while the backend performs global refinement to ensure a geometrically consistent result. This alternating execution allows the frontend and backend to mutually promote one another, enhancing overall system performance. Comprehensive experimental results demonstrate that *SLAM-Former* achieves superior or highly competitive performance compared to state-of-the-art dense SLAM methods.

## 1. Introduction

In the field of robotic perception, Simultaneous Localization and Mapping (SLAM) plays a great significance. It allows robots to construct a map of an unknown environment

while simultaneously tracking their own location. This capability is essential for robots to navigate autonomously and carry out tasks in a variety of environments. Early SLAM algorithms primarily focused on localization and mapping using sparse points, such as ORB-SLAM [1] and LSD-SLAM [2]. These methods are efficient and robust, but they may not offer detailed information about the surroundings. In contrast, dense mapping techniques aim to create a more detailed and continuous representation of the environment, mainly relying on LiDAR and RGB-D [3, 4].

With the rapid advancement of optical flow and multiview depth estimation techniques, recent research has achieved high-quality dense monocular SLAM using only images as input [5–8]. These approaches leverage the capability of neural networks and computer vision algorithms to estimate depth and motion from a single camera, thereby creating dense maps without the need for additional sensors. Especially noteworthy is the trend of utilizing geometric foundation models such as DUST3R [9] and VGGT [10]. These models reveal the high potential of data-driven 3D structure prediction. Their streaming variants, StreamVGGT [11] and Stream3R [12], by carefully lever-

aging the attention key-value cache (KV cache), enable the model to process incremental visual inputs.

We observe that SLAM methods using geometric foundation models as their reconstruction module, such as MAST3R-SLAM and VGGT-SLAM suffer from global inconsistency, since they rely on the alignment of local submaps. On the other hand, streaming methods such as StreamVGGT and Stream3R process incremental inputs without remapping the past, which can cause a significant mismatch between past and newly incoming data.

In this work, we introduce a visual SLAM framework implemented within a single unified transformer architecture, named SLAM-Former. SLAM-Former consists of a frontend and a backend in the same transformer, working in cooperation, as depicted in Fig. 1 (b). The frontend operates in real-time on the sequential RGB images for keyframe selection and incremental map and pose updates. With the incremental output from the frontend, our backend periodically refines the map and poses globally at a lower frequency.

The frontend and the backend promote each other in this alternating process. After each backend run, the transformer’s KV cache is updated to the frontend for further incremental operation. In return, the frontend provides initial results and the sequential order, assisting the backend in refinement. To empower a single transformer with all SLAM capabilities, we propose three training modes for SLAM-Former.

Compared to the traditional SLAM pipeline which requires an additional loop detection module to close its pose graph, SLAM-Former’s backend accomplishes this with full attention and is equivalent to processing loop detection on a dense factor graph.

SLAM-Former achieves significantly better reconstruction and state-of-the-art tracking performance on widely used Dense Mono SLAM benchmarks, compared to both calibrated and uncalibrated state-of-the-art methods.

## 2. Related Works

### 2.1. Dense RGB SLAM

In recent years, research on dense SLAM using monocular cameras has achieved significant progress [5–8, 13], thanks to the application of deep learning techniques. Due to the absence of depth sensors, dense RGB SLAM requires optimizing the entire sequence of geometry and camera as a whole.

Early works focus on reducing the computational cost of depth estimation. For instance, CodeSLAM [14] and DeepFactors [15] optimize the depth latent as an alternative. Borrowing the strength of MVSNet [16], Tandem relies on an external model, but it breaks the co-optimization structure [17]. Conversely, DROID-SLAM [5] and Scene-

Factory [6] incorporate deep optical flow model into the pipeline and co-optimize both with a speed-dense bundle adjustment. On the other hand, NeRF [18] and Gaussian Splatting [19] based methods have emerged as a trend to reshape Dense SLAM. NeRF-SLAM methods [13, 20] and GS-SLAM methods [21] optimize the scene as a whole for a highly realistic novel view synthesis objective. Nonetheless, those rendering-based SLAMs are usually time-consuming, not reconstruction-affordable, and highly sensitive to blur and noise, which substantially restricts their use in real life.

With the emergence of recent foundational geometry techniques, such as DUST3R [9] and VGGT [10], researchers have found new inspiration. MAST3R-SLAM [7] utilizes the advanced pairwise model MAST3R [22] for high-quality calibration-free matching and geometry construction, demonstrating state-of-the-art performance under a traditional SLAM pipeline. On the other hand, VGGT-SLAM [8] feeds submaps into VGGT and connects them using a novel  $SL(4)$  manifold, modeling the geometry distortion in foundational geometry for the first time. However, these methods rely on pair- or submap-wise optimization of geometry, which often causes conflict structure among frames. MAST3R-SLAM attempts to address this issue through TSDF-fusion [4], which could only fix small mismatches. Meanwhile, VGGT-SLAM cannot deal with this issue because it only has connections of submaps in the front- and end-nodes.

This inspires us to develop a frontend-backend SLAM structure that can properly and neatly address this issue.

### 2.2. Feed-forward 3D Reconstruction

Recently, DUST3R [9] has led a trend to regress the 3D structure with scalable training data directly. However, with processing image pairs, DUST3R requires a global optimization for larger scenes, which lowers the inference efficiency. Several works have been proposed to address this limitation. Fast3R [23], VGGT [10], and Pi3 [24] process multi-view images in a single forward pass, avoiding time-consuming post-processing global optimization. All three are Transformer-based models for multi-view pointmap estimation. Fast3R highlights the ability to efficiently handle thousands of images, while VGGT demonstrates that a simple architecture with 3D multi-task learning and scalable training data can achieve state-of-the-art results. Pi3 further introduces a permutation-equivariant design that removes the dependence on a fixed reference view, enhancing robustness to input ordering and scalability.

Besides feed-forward multi-view methods, recent feed-forward streaming approaches tackle 3D reconstruction online. Spann3R [25] extends Dust3R to streaming by maintaining and interacting with a spatial memory. CUT3R [26] introduces persistent state tokens with transformer-based

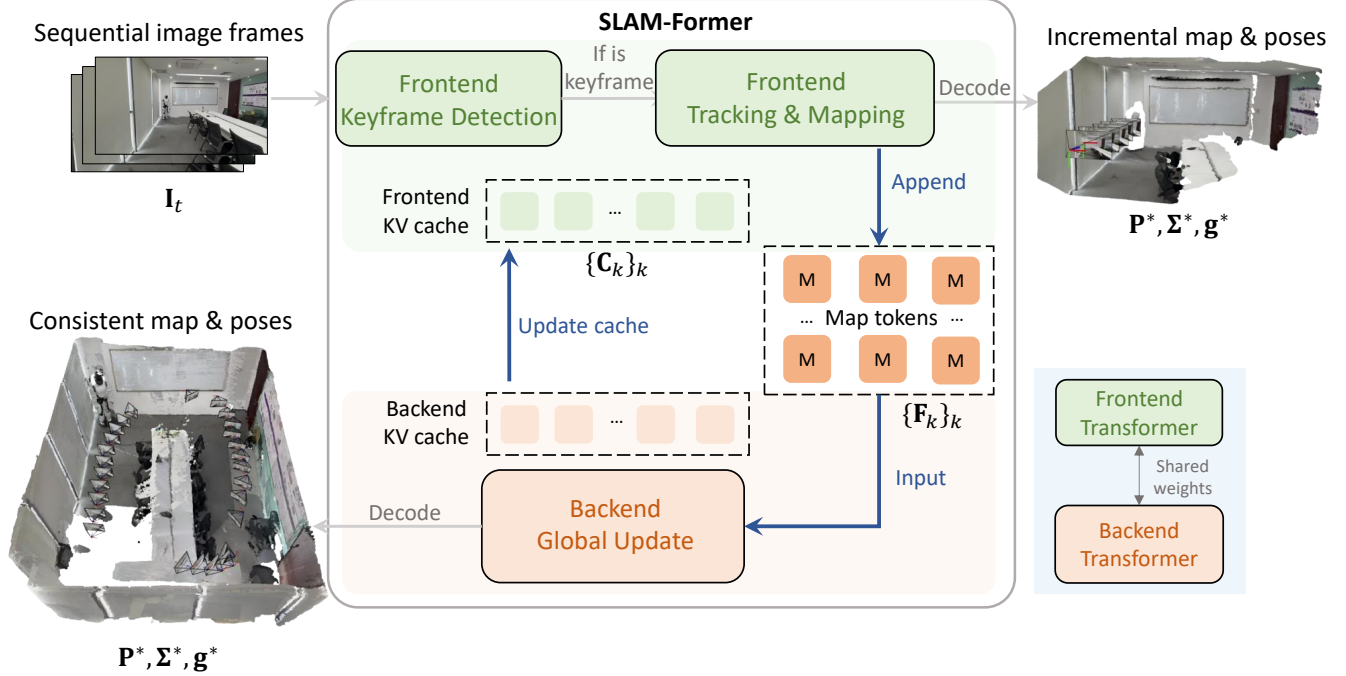


Figure 2. The working pipeline of SLAM-Former. The frontend detects keyframes and performs incremental pose and map updates, while the backend performs global pose and map updates. The shared map token memory and the KV cache update mechanism ensure that the frontend and the backend promote each other, and this process is marked by blue arrows  $\rightarrow$ .

recurrent updates for online reconstruction. LONG3R [27] employs a 3D spatio-temporal memory and a coarse-to-fine pipeline to handle long sequence streaming reconstruction. StreamVGGT [11] and SStream3R [12] further incorporate causal attention, drawing inspiration from modern language models to enable real-time streaming reconstruction.

However, existing streaming methods focus solely on incremental updates without revisiting past estimates, leading to drift and limited global consistency. To address this, we propose SLAM-Former, a unified neural SLAM pipeline that integrates both frontend and backend for efficient incremental updates and periodic global refinement.

### 3. SLAM-Former

This section introduces our proposed *SLAM-Former*. We first describe the underlying transformer architecture in Sec. 3.1, then detail its roles in the SLAM frontend (Sec. 3.2) and backend (Sec. 3.3). We further present a joint training strategy that unifies these tasks within a single model in Sec. 3.4, followed by the inference pipeline in Sec. 3.5.

#### 3.1. Transformer Architecture

SLAM-Former is built upon a single transformer model, where a **transformer backbone**  $f$  aggregates both intra-frame and inter-frame information, and task-specific **heads**

$h$  decode scene geometry and camera poses. For clarity, we assume that image features are pre-encoded, and the input to  $f$  consists of a set of image patch tokens augmented with register tokens. Following Pi3-like designs, we employ shared register tokens across all frames, thereby eliminating the need to designate a reference frame. The backbone contains  $L$  layers, each equipped with intra-frame and inter-frame attention to jointly capture local image context and temporal correspondences.

SLAM-Former integrates both *frontend* for incremental frame processing, and *backend* for global map and pose refinement within a shared transformer backbone (see Fig. 2).

#### 3.2. Frontend

We illustrate the frontend process in Fig. 2. When a new frame arrives, the frontend first decides whether it should be a new keyframe. If so, the system proceeds with tracking and mapping.

Formally, given an image sequence  $\{\mathbf{I}_t \in \mathbb{R}^{3 \times H \times W}\}_{t \in \{1, 2, \dots\}}$ , the **frontend**  $f_{\text{fn}}$  maps each frame into a set of *map tokens*:

$$\mathbf{F}_t = f_{\text{fn}}(\mathbf{I}_t)_{\{\mathbf{C}_k\}_{k \in \mathcal{S}}} \quad (1)$$

where  $\{\mathbf{C}_k\}_{k \in \mathcal{S}}$  denotes the **KV cache** of previous keyframes, storing the key (K) and value (V) tensors prior to the inter-frame attention layer. Here,  $\mathcal{S}$  is the set of

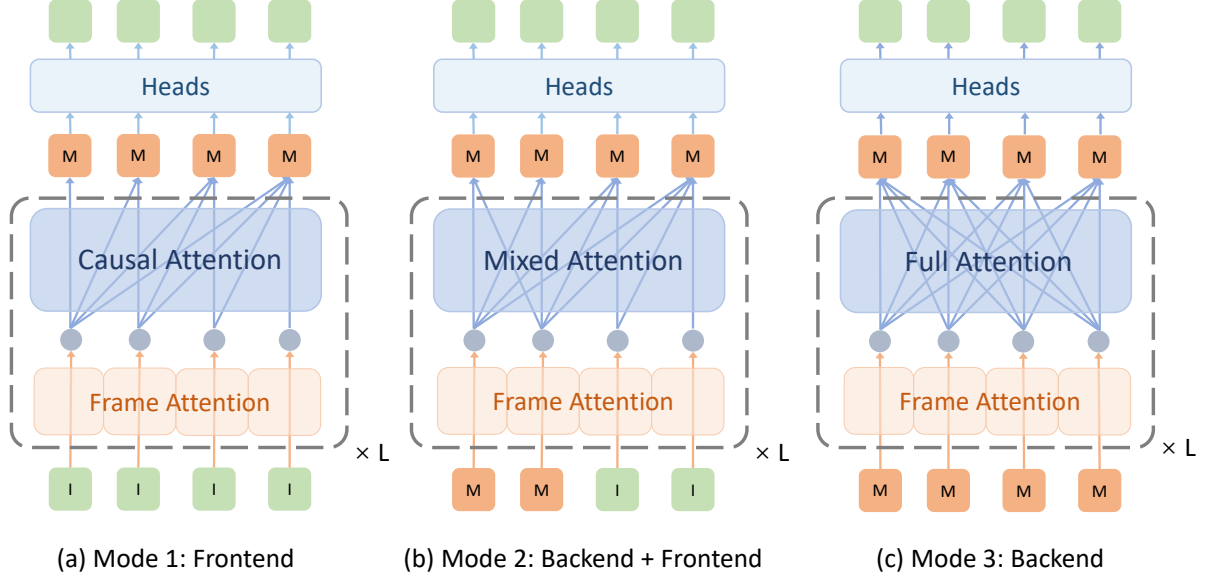


Figure 3. Three training modes of SLAM-Former.  $\text{I}$  and  $\text{M}$  represent the patch-wise image tokens and map tokens of a frame.  $\bullet$  and  $\blacksquare$  represent the layer intermediate and final output. In each mode, either  $\text{I}$  or  $\text{M}$  tokens, or both, are fed into the transformer backbone  $f$ , which contains  $L$  layers of frame attention and various inter-frame attentions. Finally, pose and pointmap are regressed by the heads  $h$ .

keyframe indices,  $S = |S|$  its cardinality, and  $\mathbf{F}_t$  the resulting map tokens of frame  $t$ , which serve as an implicit neural representation of the scene. The new KV cache generated in this process,  $\mathbf{C}_t = \text{Cache}(f(\mathbf{F}_t))$ , will also be extended to  $\{\mathbf{C}_k\}_{k \in S}$  for later use.

**Keyframe Detection.** After generating map tokens  $\mathbf{F}_t$ , the frontend estimates the camera pose with the pose head  $h_{\text{pose}}$ :

$$\mathbf{g}_t = h_{\text{pose}}(\mathbf{F}_t). \quad (2)$$

A frame is marked as a new keyframe if its relative pose to the latest keyframe  $k_{\text{prev}}$ ,  $\mathbf{g}_{k_{\text{prev}}, t} = \mathbf{g}_{k_S}^{-1} \mathbf{g}_t$ , exceeds a translation threshold  $\tau$ .

In practice, for keyframe detection we do not rely on the KV cache; instead, we directly apply  $f_{\text{fn}}(\mathbf{I}_{k_{\text{prev}}}, \mathbf{I}_t)$  to the frame pair  $(k_{\text{prev}}, t)$ , which improves efficiency and avoids the need to specify a reference frame (see Sec. 3.4).

**Frontend Tracking and Mapping.** If a new keyframe is confirmed,  $\mathbf{F}_t$  is recomputed using the full KV cache from Eq. (1), and the token map  $(\mathbf{M}, \mathbf{S})$  is updated:

$$\mathbf{M} \leftarrow \mathbf{M} \cup \{\mathbf{F}_t\}, \quad \mathbf{S} \leftarrow \mathbf{S} \cup \{t\}, \quad S \leftarrow S + 1. \quad (3)$$

The frontend only relies on past frames, making it causal and suitable for online tracking. However, such causality inevitably leads to error accumulation and local inconsistencies. To mitigate this, we introduce a **backend** module to perform the global refinement.

### 3.3. Backend

The backend is responsible for refining the map tokens to enforce global consistency. As illustrated in Fig. 2, traditional SLAM pipelines typically rely on loop closure detection and graph optimization for this purpose. In contrast, our approach employs a transformer-based **backend**  $f_{\text{bn}}$  that directly refines all map tokens in a single pass:

$$\bar{\mathbf{M}} = f_{\text{bn}}(\mathbf{M}). \quad (4)$$

The effectiveness of this design lies in the full attention mechanism within  $f_{\text{bn}}$ , which establishes dense connections across all map tokens. This global receptive field enables the backend to correct accumulated drift and enforce structural coherence across the reconstructed scene.

**Cache Sharing.** To inherit the benefits of backend refinement, the frontend reuses shared KV cache  $\mathbf{C}_{\mathbf{M}}$  from the backend:

$$\{\mathbf{C}_k\}_{k \in S} \leftarrow \mathbf{C}_{\mathbf{M}}. \quad (5)$$

In this way, subsequent frames are tracked and mapped with respect to the refined global structure, reducing the risk of error accumulation in long sequences.

### 3.4. Training Strategy

The training strategy is designed to enable a single transformer to handle both frontend and backend SLAM functionalities, as illustrated in Fig. 3. We jointly train SLAM-

Former across three modes, each corresponding to a different input–output relationship, within a single iteration.

**Training Frontend.** The frontend is trained with a causal attention mask (Fig. 3 (a) Mode 1). At inference time, it reuses KV caches from previous frames, which enables efficient, end-to-end learning in a single pass:  $\mathbf{F} = f(\mathbf{I})_{KV}$ .

However, pure causal attention inherently designates the first frame as the reference frame. When operating on two or more frames jointly, no single frame defines the coordinates, removing the dependence on a fixed reference view and avoiding the need for frames to be close to it. We therefore apply *full attention* to the first two frames, while still leveraging causal attention for all subsequent frames. With this design, during inference, keyframe detection is performed by passing the last keyframe and the incoming frame together through  $f_{\text{in}}$ . For tracking and mapping, the first two keyframes are jointly processed to determine the global coordinate.

**Training Frontend with Backend Cooperation.** To bridge frontend and backend operations (Mode 2), we train  $f$  with *mixed attention* to simultaneously process backend and cache-sharing functionality. Specifically, the backend refines map tokens with full attention,  $\bar{\mathbf{M}} = f_{\text{bn}}(\mathbf{M})$ , while the frontend processes new images in the same forward pass as the backend, using causal attention that is equivalent to conditioning on the backend-refined KV cache:  $\mathbf{F} = f_{\text{in}}(\mathbf{I})_{C_M}$ .

**Training Backend.** The backend (Mode 3) refines map tokens originating from different runs or KV cache states. Here, full attention is always applied, allowing the model to resolve drift and enforce global consistency:  $\bar{\mathbf{M}} = f_{\text{bn}}(\mathbf{M})$ .

**Joint Training.** In all modes, the resulting tokens serve as implicit representations of geometry and camera poses. Task-specific heads predict the pointmap  $\mathbf{P}^*$ , confidence  $\Sigma^*$ , and camera pose  $\mathbf{g}^*$ :

$$\mathbf{P}^*, \Sigma^*, \mathbf{g}^* = h(\mathbf{F}). \quad (6)$$

Unlike VGGT, which predicts global geometry, SLAM-Former produces local pointmaps for each frame to avoid the need to define a specific world coordinate.

The overall loss combines depth, pointmap, and camera supervision:

$$L = L_{\text{depth}} + L_{\text{pmap}} + \lambda L_{\text{cam}}. \quad (7)$$

For depth loss, the predicted depth  $\mathbf{D}^* = \mathbf{P}_z^*$  is supervised against ground-truth depth  $\mathbf{D}$ , weighted by confidence  $\Sigma^*$ :  $L_{\text{depth}} = \sum_t \left( \|\Sigma_t^* \odot (s^* \mathbf{D}_t^* - \mathbf{D}_t)\| + \|\Sigma_t^* \odot (\nabla s^* \mathbf{D}_t^* - \nabla \mathbf{D}_t)\| - \alpha \log \Sigma_t^* \right)$ , where  $\odot$  is element-wise multiplication,  $\nabla$  the spatial gradient, and  $s^*$  a scale factor estimated following Pi3:  $s^* = \text{argmin}_s \sum_t \|s(\mathbf{P}_t^* - \mathbf{P}_t)/\mathbf{D}_t\|_1$ .

For pointmap loss, similar to depth loss but defined on transformed local pointmaps aligned to the first frame:  $\mathbf{P}_{t,1}^* = \mathbf{g}_1^{*-1} \mathbf{g}_t^* \mathbf{P}_t^*$ , the loss is designed as  $L_{\text{pmap}} = \sum_t \left( \|\Sigma_t^* \odot (s^* \mathbf{P}_{t,1}^* - \mathbf{P}_t)\| + \|\Sigma_t^* \odot (\nabla s^* \mathbf{P}_{t,1}^* - \nabla \mathbf{P}_t)\| - \alpha \log \Sigma_t^* \right)$ .

For camera loss, relative pose consistency is supervised using a scaled Huber loss:  $L_{\text{cam}} = \sum_{i,j} \|s_i^* \otimes (\mathbf{g}_i^{*-1} \mathbf{g}_j^*) - (\mathbf{g}_i^{*-1} \mathbf{g}_j^*)\|_\epsilon$ , where  $\otimes$  scales translations and  $\|\cdot\|_\epsilon$  denotes the Huber norm.

All three modes are executed sequentially in a single iteration with shared weights. The final training objective is:

$$L_{\text{all}} = L_1 + L_2 + \beta L_3.$$

### 3.5. Execution Pipeline

The execution pipeline integrates frontend and backend to perform online SLAM inference.

**Frontend.** Each incoming frame is first passed to the keyframe detector. If identified as a keyframe, it undergoes further processing. The first two keyframes are jointly processed for initialization, producing map tokens  $\{\mathbf{F}_1, \mathbf{F}_2\}$  and KV caches  $\{\mathbf{C}_1, \mathbf{C}_2\}$ , which are stored. For the  $k$ -th keyframe ( $k > 2$ ), the frontend leverages the cached tokens  $\{\mathbf{C}_i\}_{i < k}$  to generate  $\mathbf{F}_k$  and its cache  $\mathbf{C}_k$ , which are appended to the storage.

**Backend.** After every  $T$  keyframes, the backend is triggered. The accumulated map tokens are refined, and the resulting KV caches are used to update the first  $T$  frontend caches  $\{\mathbf{C}_i\}_{i \leq k}$ .

## 4. Experiments

We evaluate SLAM-Former on multiple tasks, including camera tracking (Sec. 4.2) and dense 3D reconstruction (Sec. 4.3). Subsequently, we analyze the impact of our frontend-backend design (Sec. 4.4) and evaluate the time efficiency (Sec. 4.5).

### 4.1. Experimental Setup

**Implementation Details.** SLAM-Former has 36 transformer layers of both frame- and global-attention in total. We initialize SLAM-Former with Pi3 pre-trained weights and train 10 epochs with a batch size of 32 (excluding the frozen image encoder and camera head). For training, we utilize the AdamW optimizer with a learning rate of  $1e-5$ , and a cosine learning rate scheduler. In the loss function, the hyperparameters are set as  $\lambda = 100$  and  $\beta = 10$ . Regarding datasets, SLAM-Former is trained on ARKitScenes [35], ScanNet [36], ScanNet++ [37], HyperSim [38], Blended-MVS [39], MegaDepth [40], and MVS-Synth [41]. During each iteration, all three modes of a single SLAM-Former are trained. The entire training process takes 11 hours on

	Method	Sequence									Avg
		360	desk	desk2	floor	plant	room	rpy	teddy	xyz	
Calib.	ORB-SLAM3 [1]	×	0.017	0.210	×	0.034	×	×	×	<b>0.009</b>	N/A
	DeepV2D [28]	0.243	0.166	0.379	1.653	0.203	0.246	0.105	0.316	0.064	0.375
	DeepFactors [15]	0.159	0.170	0.253	0.169	0.305	0.364	0.043	0.601	0.035	0.233
	DPV-SLAM [29]	0.112	0.018	0.029	0.057	0.021	0.330	0.030	0.084	0.010	0.076
	DPV-SLAM++ [29]	0.132	0.018	0.029	0.050	0.022	0.096	0.032	0.098	0.010	0.054
	GO-SLAM [30]	0.089	<b>0.016</b>	0.028	0.025	0.026	<b>0.052</b>	<b>0.019</b>	0.048	0.010	0.035
	DROID-SLAM [5]	0.111	0.018	0.042	<b>0.021</b>	<b>0.016</b>	<b>0.049</b>	0.026	0.048	0.012	0.038
	MASt3R-SLAM [7]	<b>0.049</b>	<b>0.016</b>	<b>0.024</b>	0.025	0.020	0.061	0.027	<b>0.041</b>	<b>0.009</b>	<b>0.030</b>
Uncalib.	DROID-SLAM* [5]	0.202	0.032	0.091	0.064	0.045	0.918	0.056	0.045	0.012	0.158
	MASt3R-SLAM* [7]	0.070	0.035	0.055	<b>0.056</b>	0.035	0.118	0.041	0.114	0.020	0.060
	VGGT-SLAM [8]	0.071	0.025	0.040	0.141	0.023	0.102	0.030	0.034	0.014	0.053
	CUT3R+ [26]	0.102	0.054	0.118	0.211	0.083	0.264	0.044	0.120	0.020	0.113
	StreamVGGT+ [11]	0.088	0.063	0.105	0.604	0.070	0.633	0.025	0.081	0.015	0.187
	VGGT-SLAM+ [10]	0.053	0.024	0.040	0.335	0.022	0.166	0.040	0.037	0.041	0.084
	Ours	<b>0.067</b>	<b>0.018</b>	<b>0.026</b>	0.079	<b>0.021</b>	<b>0.082</b>	<b>0.017</b>	<b>0.030</b>	<b>0.011</b>	<b>0.039</b>

Table 1. Root mean square error (RMSE) of absolute trajectory error (ATE) on TUM RGB-D [31] (unit: m). The \* symbol indicates that the baseline is evaluated in the uncalibrated mode from the VGGT-SLAM paper [8], the + symbol indicates that the baseline is tested on our machine.

	Method	Sequence							Avg
		chess	fire	heads	office	pumpkin	kitchen	stairs	
Calib.	NICER-SLAM [13]	<b>0.033</b>	0.069	0.042	0.108	0.200	<b>0.039</b>	0.108	0.086
	DROID-SLAM [5]	0.036	0.027	0.025	<b>0.066</b>	0.127	0.040	0.026	0.049
	MASt3R-SLAM [7]	0.053	<b>0.025</b>	<b>0.015</b>	0.097	<b>0.088</b>	0.041	<b>0.011</b>	<b>0.047</b>
Uncalib.	DROID-SLAM* [5]	0.047	0.038	0.034	0.136	0.166	0.080	0.044	0.078
	MASt3R-SLAM* [7]	0.063	0.046	0.029	0.103	0.114	0.074	0.032	0.066
	VGGT-SLAM [10]	0.036	0.028	<b>0.018</b>	0.103	0.133	0.058	0.093	0.067
	CUT3R+ [26]	0.046	0.043	0.055	0.120	0.096	0.061	0.086	0.073
	StreamVGGT+ [11]	0.048	0.036	0.030	0.117	0.094	0.063	0.179	0.081
	VGGT-SLAM+ [8]	<b>0.037</b>	<b>0.027</b>	<b>0.018</b>	0.101	0.138	0.057	0.095	0.068
	Ours	0.039	0.033	<b>0.018</b>	<b>0.070</b>	<b>0.065</b>	<b>0.035</b>	<b>0.033</b>	<b>0.042</b>

Table 2. Root mean square error (RMSE) of absolute trajectory error (ATE) on 7-Scenes [32] (unit: m). The \* symbol indicates that the baseline is evaluated in the uncalibrated mode from the VGGT-SLAM paper [8], the + symbol indicates that the baseline is tested on our machine.

Method		Sequence								Avg
		Rm0	Rm1	Rm2	Of0	Of1	Of2	Of3	Of4	
Calib.	NICER-SLAM [13]	0.013	0.016	0.011	0.021	0.032	0.021	0.014	0.020	0.019
	DROID-SLAM [5]	<b>0.003</b>	<b>0.001</b>	<b>0.003</b>	<b>0.003</b>	<b>0.004</b>	<b>0.003</b>	<b>0.005</b>	<b>0.004</b>	<b>0.003</b>
Uncalib.	SLAM3R [33]	0.046	0.059	0.057	0.112	0.063	0.062	0.050	0.081	0.066
	CUT3R+ [26]	0.145	0.243	0.127	0.159	0.230	0.162	0.088	0.204	0.170
	StreamVGGT+ [11]	0.113	0.163	0.077	0.076	0.070	0.180	0.153	0.168	0.125
	VGGT-SLAM+ [8]	0.030	0.167	0.086	0.042	0.064	0.095	0.039	0.043	0.071
	Ours	<b>0.030</b>	<b>0.026</b>	<b>0.027</b>	<b>0.028</b>	<b>0.029</b>	<b>0.038</b>	<b>0.028</b>	<b>0.031</b>	<b>0.030</b>

Table 3. Root mean square error (RMSE) of absolute trajectory error (ATE) on Replica [34] (unit: m). The + symbol indicates that the baseline is tested on our machine.

32 A100 GPUs. Evaluation is conducted on a single RTX 4090 GPU during testing.

**Baselines.** The baselines utilized in our experiments are categorized as Calibrated and Uncalibrated:

- Calibrated: ORB-SLAM3 [42], DeepV2D [28], DeepFactors [15], DPV-SLAM [29], DPV-SLAM++ [29], GO-SLAM [30], DROID-SLAM [5], MASt3R-SLAM [7] and NICER-SLAM [13].
- Uncalibrated: DROID-SLAM and MASt3R-SLAM, VGGT-SLAM, SLAM3R [33], as well as our method, SLAM-Former. Additionally, we also test related methods CUT3R [26] and StreamVGGT [11] using our keyframes.

## 4.2. 3D Tracking Evaluation

We first evaluate the tracking performance of SLAM-Former on the TUM RGB-D, 7-Scenes, and Replica. We compute the Root Mean Square Error (RMSE) of Absolute Trajectory Error (ATE) for various methods under both calibrated and uncalibrated settings.

**TUM RGB-D Track.** In the TUM test, evaluations are conducted on a widely used subset of scenes. The results are summarized in Tab. 1. As shown, our model consistently outperforms most baselines in the uncalibrated setting. The superior performance in these more complex sequences, such as the room and floor that involve significant camera rotation and potential for loop closure, suggests that our backend’s global refinement is particularly effective at mitigating accumulated drift. More importantly, it significantly reduced the error relative to calibrated baselines, achieving a highly competitive level.

**7-scenes Track.** Following a similar protocol as with the TUM RGB-D tracks, we evaluate the 7-scene as shown in Tab. 2. Our method outperforms most baselines under both the uncalibrated and calibrated settings. In more complex scenes, such as the office, pumpkin and kitchen, our model achieves a notably higher performance gap compared to the rest. On average, our method outperforms all baselines.



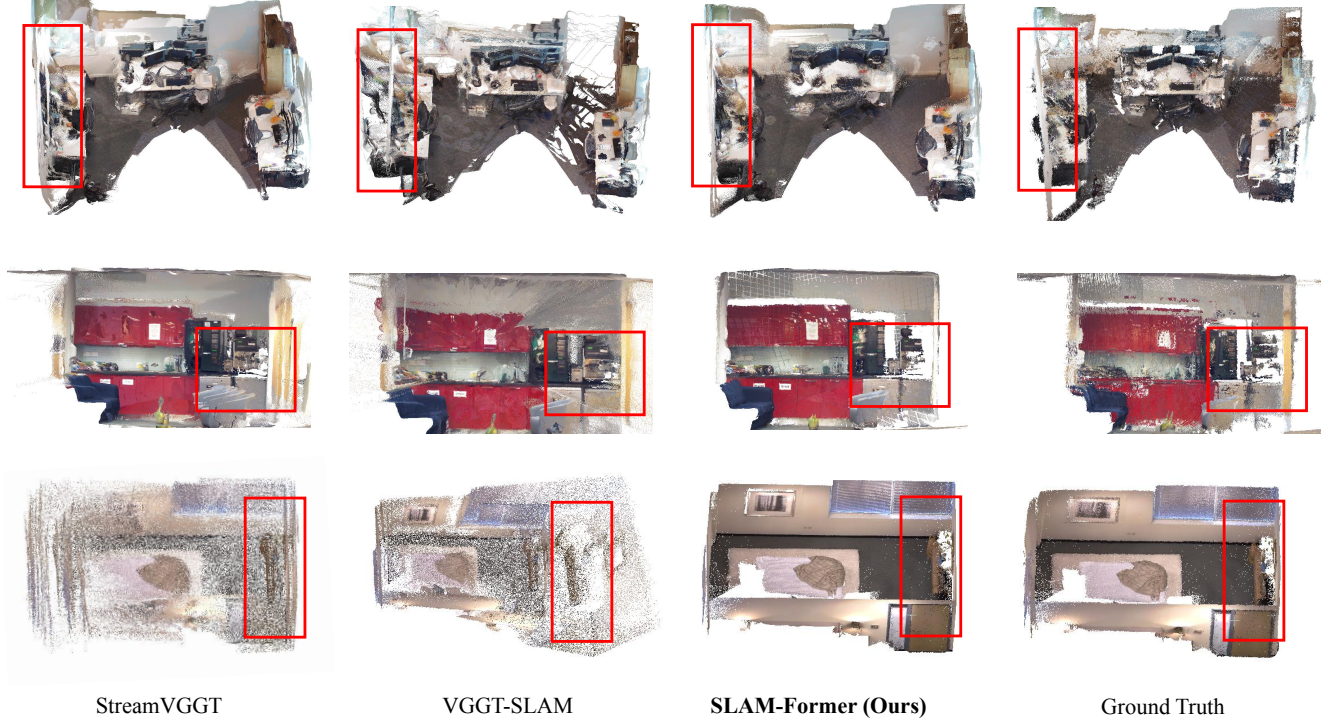


Figure 4. Qualitative reconstruction comparison. Note the significant structural errors, such as misalignments, from baseline methods (red boxes), which are corrected by SLAM-Former’s globally consistent refinement.

Method	Room 0 Acc. / Comp.	Room 1 Acc. / Comp.	Room 2 Acc. / Comp.	Office 0 Acc. / Comp.	Office 1 Acc. / Comp.	Office 2 Acc. / Comp.	Office 3 Acc. / Comp.	Office 4 Acc. / Comp.	Average Acc. / Comp.
DUSi3R [9]	3.47 / 2.50	2.53 / 1.86	2.95 / 1.76	4.92 / 3.51	3.09 / 2.21	4.01 / 3.10	3.27 / 2.25	3.66 / 2.61	3.49 / 2.48
MASi3R [22]	4.01 / 4.10	3.61 / 3.25	3.13 / 2.15	2.57 / 1.63	12.85 / 8.13	3.13 / 1.99	4.67 / 3.15	3.69 / 2.47	4.71 / 3.36
NICER-SLAM* [13]	2.53 / 3.04	3.93 / 4.10	3.40 / 3.42	5.49 / 6.09	3.45 / 4.42	4.02 / 4.29	3.34 / 4.03	3.03 / 3.87	3.65 / 4.16
DROID-SLAM* [5]	12.18 / 8.96	8.35 / 6.07	3.26 / 16.01	3.01 / 16.19	2.39 / 16.20	5.66 / 15.56	4.49 / 9.73	4.65 / 9.63	5.50 / 12.29
DIM-SLAM* [43]	14.19 / 6.24	9.36 / 6.45	8.41 / 12.17	10.16 / 5.95	7.86 / 8.33	16.50 / 8.28	13.01 / 6.77	13.08 / 8.62	11.60 / 7.85
GO-SLAM <sup>-</sup> [30]	-	-	-	-	-	-	-	-	3.81 / 4.79
Spann3R <sup>-</sup> [25]	9.75 / 12.94	15.51 / 12.94	7.28 / 8.50	5.46 / 18.75	5.24 / 16.64	9.33 / 11.80	16.00 / 9.03	13.97 / 16.02	10.32 / 13.33
SLAM3R <sup>-</sup> [33]	3.19 / 2.40	3.12 / 2.34	2.72 / 2.00	4.28 / 2.60	3.17 / 2.34	3.84 / 2.78	3.90 / 3.16	4.32 / 3.36	3.57 / 2.62
CUT3R <sup>+</sup> [26]	6.09 / 3.09	9.89 / 4.55	5.77 / 2.66	5.23 / 2.46	11.60 / 6.94	8.00 / 3.16	6.12 / 3.09	7.46 / 3.05	7.52 / 3.62
StreamVGGT <sup>+</sup> [11]	9.01 / 4.37	12.22 / 4.66	5.61 / 2.61	6.64 / 3.45	5.14 / 2.55	12.84 / 7.23	12.09 / 6.59	15.48 / 6.35	9.88 / 4.73
VGGT-SLAM <sup>+</sup> [8]	3.23 / 2.66	18.92 / 15.41	15.93 / 12.37	4.01 / 2.49	3.01 / 2.32	6.93 / 5.31	2.97 / 2.53	5.19 / 3.83	7.52 / 5.86
<b>SLAM-Former (Ours)</b>	<b>2.40 / 1.86</b>	<b>1.79 / 1.35</b>	<b>1.83 / 1.39</b>	<b>1.84 / 1.37</b>	<b>1.70 / 1.28</b>	<b>2.44 / 1.67</b>	<b>2.36 / 1.89</b>	<b>2.38 / 1.70</b>	<b>2.09 / 1.56</b>

Table 4. Reconstruction results on the Replica [34] dataset. \* denotes the results reported in NICER-SLAM [13]. <sup>-</sup> shows the results from SLAM3R [33]. <sup>+</sup> represented the results from our run.

Method	7-Scenes			
	ATE ↓	Acc. ↓	Comple. ↓	Chamfer ↓
Calib.				
DROID-SLAM [5]	0.049	0.141	0.048	0.094
MASi3R-SLAM [7]	0.047	0.089	0.085	0.087
Spann3R @20 [25]	N/A	0.069	0.047	0.058
Spann3R @2 [25]	N/A	0.124	0.043	0.084
Uncalib.				
MASi3R-SLAM* [7]	0.066	0.068	0.045	0.056
VGGT-SLAM [8]	0.067	0.052	0.058	0.055
CUT3R <sup>+</sup> [26]	0.073	0.032	0.047	0.040
StreamVGGT <sup>+</sup> [11]	0.081	0.058	0.057	0.057
VGGT-SLAM <sup>+</sup> [8]	0.068	0.054	0.060	0.057
<b>Ours</b>	<b>0.042</b>	<b>0.017</b>	<b>0.037</b>	<b>0.027</b>

Table 5. Reconstruction evaluation on 7-Scenes [32] (unit: m). @n indicates a keyframe every n images.

**Replica Track.** The previous tracking experiments were conducted with real captured data, while the replica dataset

is synthetic. Our approach shows substantial improvements under the uncalibrated setting, achieving an approximate 50% reduction in ATE compared to SLAM3R and outperforming all baselines, as shown in Tab. 3. However, our method performs on par with NICER-SLAM but still lags behind the traditional SLAM method, DROID-SLAM. This is because synthetic data lacks noise and blur, making the matchings sufficiently accurate for pose solving in bundle adjustment. In contrast, in the previous real-life data tests, DROID-SLAM performed on a similar level to our method.

### 4.3. Reconstruction Evaluation

We evaluate the reconstruction performance of our SLAM-Former on the 7-Scenes dataset following the protocol of VGGT-SLAM and on the Replica dataset following the pro-

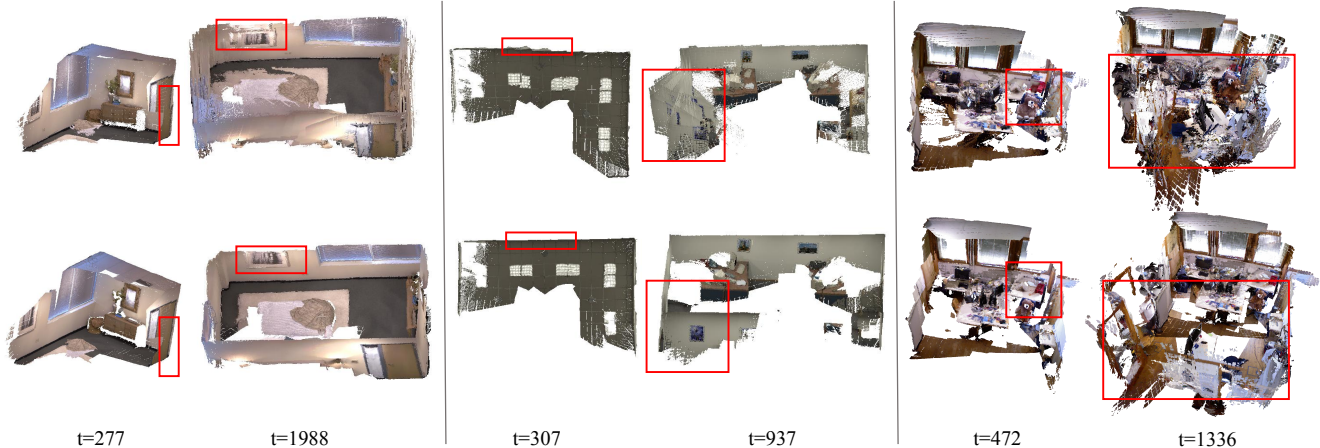


Figure 5. Qualitative reconstruction comparison with and without backend assistance. The first row displays the frontend-only results at the corresponding timestamps, while the second row shows the results with the assistance of backend KV caches.

to col of SLAM3R.

**7-scenes Reconstruction.** The dense reconstruction results on 7-Scenes are shown in Tab. 5. Our method exhibits a significant gap compared to other dense SLAM state-of-the-art methods. In terms of reconstruction quality, our method achieves the highest reconstruction accuracy of 0.017m, while the others are all above 0.05m. In terms of completeness and chamfer distance, our method achieves 0.037m and 0.027m, respectively, still outperforming all baselines by around 50%.

This consistently superior performance across all major reconstruction metrics is also demonstrated in our reconstruction demo Fig. 4. As shown in the first two rows (office and pumpkin), the baselines exhibit mismatched surfaces between frames within the red-windowed regions. In contrast, our SLAM-Former’s reconstruction consistently shows coherent and accurate structure.

**Replica Reconstruction.** The dense reconstruction results on the Replica dataset are listed in Tab. 4. Our method also achieves the best in terms of both accuracy and completion across all baselines. Specifically, our 2.09/1.56 Acc./Comp. show at least 1cm gap to the second best on both metrics, respectively.

We also demonstrate the reconstruction in the third row of Fig. 4. Here, StreamVGGT shows multiple layers of surface in the room, as highlighted in the red-windowed region. More severely, VGGT-SLAM exhibits layers with substantial scale discrepancies. While SLAM-Former closely matches the ground truth. The less dense point clouds of the baselines are caused by the mismatch of layers, as the test samples a constant number of points for evaluation.

Method	Sequence										Avg
	360	desk	desk2	floor	plant	room	rpy	teddy	xyz		
F.-only	0.137	0.041	0.045	0.264	0.053	0.547	0.036	0.073	0.013	0.134	
F. + EB.	0.073	0.021	0.026	0.078	0.022	0.104	0.018	0.027	0.011	0.042	
F. + MB.	0.067	0.018	0.025	0.081	0.023	0.082	0.017	0.031	0.011	<b>0.039</b>	
F. + MB. + EB.	<b>0.067</b>	<b>0.018</b>	<b>0.026</b>	<b>0.079</b>	<b>0.021</b>	<b>0.082</b>	<b>0.017</b>	<b>0.030</b>	<b>0.011</b>	<b>0.039</b>	

Table 6. Evaluation of module cooperation with RMSE of ATE on TUM RGB-D [31] (unit: m).

#### 4.4. Frontend and backend co-operation

To investigate how the backend design of our SLAM-Former contributes to the overall system performance in a SLAM-like manner, we conducted a series of ablation experiments. The results are summarized in Tab. 6. Here, **F**, **MB.**, and **EB.** denote the Frontend, Middle Backend, and End Backend components of our architecture, respectively. All evaluations are conducted on the TUM RGB-D benchmark using the RMSE of ATE as the metric.

The results demonstrate that the inclusion of backend modules yields significantly improved accuracy over using the frontend alone, confirming the effectiveness of our proposed frontend-backend design.

##### How does the backend assist the frontend?

Furthermore, while the individual contributions of MB. and EB. appear comparable in terms of ATE, and the combined use of MB. and EB. does not show marked improvement on the overall metric, it is important to emphasize that MB. plays a crucial role beyond what is captured by the final ATE.

We demonstrate intermediate results on some of the most challenging sequences in datasets, including Replica-room1, ICLNUIM [44]-ofkt1, and TUM-room, all of which



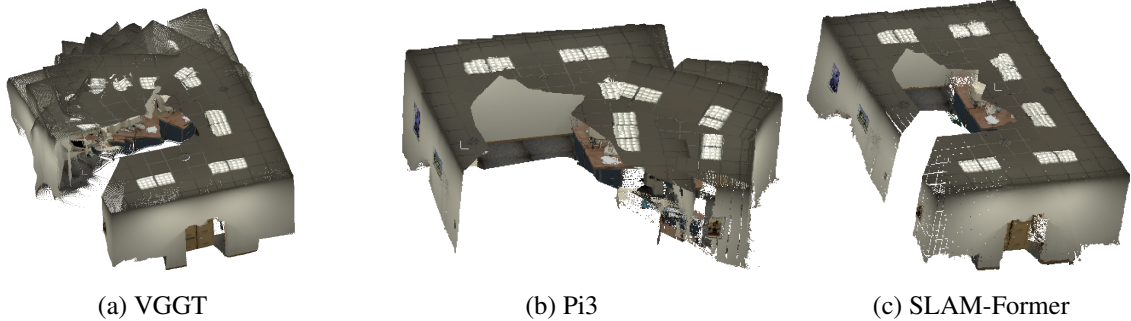


Figure 6. Qualitative reconstruction comparison on ICL-NUIM ofkt1. The results from VGGT, Pi3, and ours are displayed from left to right, respectively. Both VGGT and Pi3 suffer from pose drift, leading to geometric inaccuracies, while our method demonstrates consistent and accurate reconstruction.

are inside-out captures of an indoor environment as illustrated in Fig. 5. Initially, the errors in the frontend-only results are relatively small, as highlighted by the red windows. However, as time progresses, the frontend-only reconstruction becomes significantly distorted. This distortion arises because the frontend-only **accumulates errors over time**, leading to substantial inaccuracy in the later stages. In contrast, our model, which incorporates a backend, **maintains consistency throughout the entire process**, effectively mitigating these issues.

#### How does the backend benefit from the frontend?

The above test demonstrates that the frontend can benefit from the backend to achieve consistent performance over time.

However, how can the backend benefit from the frontend? Is it equivalent to simply running all keyframe images through a single pass of VGGT or Pi3?

To answer this question, we conduct a demonstration using the ICL-NUIM scene, ofkt0 sequence. As shown in Fig. 6, the left two images display the results of VGGT and Pi3 when all keyframe images are used as input, without any sequential information. The right image shows our result. It is evident that without the sequential information provided by our frontend, VGGT and Pi3 produce disordered reconstructions. In contrast, our backend leverages the implicit order provided by the frontend to achieve a more coherent and accurate reconstruction.

#### 4.5. Execution Speed

We also record the time cost of the method, as shown in Tab. 7. We demonstrate the module execution time on different datasets. The KF-detection and frontend run in less than 100ms on average, while the backend, though slower, executes less frequently. Note that the overall speed is  $> 10\text{Hz}$ , indicating that our method operates in real-time.

	KF-detection TPE(ms)	Frontend TPE(ms)	Backend TPE(ms)	Summary FPS
TUM: room	89	97	187	10.8
7Scene: chess	89	77	83	11.0
Replica: room1	87	76	113	11.3

Table 7. Time cost on Datasets. TPE denotes Time per Execution, and FPS represents Frames per Second.

## 5. Conclusion

In this work, we introduce SLAM-Former, putting the full SLAM capabilities into a single transformer. With alternating incremental frontend processing and global backend processing, SLAM-Former enables the frontend and backend to cooperate and enhance each other, resulting in an overall improvement. The results demonstrate that SLAM-Former significantly outperforms traditional geometry foundation-based SLAM methods in both tracking and reconstruction. Furthermore, it achieves highly competitive tracking performance and vastly superior reconstruction compared to traditional methods when tested with real-world data.

**Limitations** SLAM-Former still has several limitations. First of all, we utilize full attention in the backend to replace the loop-detection & optimization in traditional SLAM. However, using full attention causes problems due to the  $\mathcal{O}(n^2)$  time complexity. We believe this could be solved in either SLAM ways (using a sparse graph) or through transformer techniques (such as sparse attention and token merge) in future work. Secondly, SLAM-Former does not support a frontend mode locally; all previous KV caches should be fed into the model during inference.

## References

- [1] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, “Orb-slam: A versatile and accurate monocular slam system,” *IEEE Transactions on Robotics*, vol. 31, p. 1147–1163, Oct. 2015.

- [2] J. Engel, T. Schöps, and D. Cremers, “Lsd-slam: Large-scale direct monocular slam,” in *European conference on computer vision*, pp. 834–849, Springer, 2014.
- [3] J. Elseberg, D. Borrmann, and A. Nüchter, “One billion points in the cloud—an octree for efficient processing of 3d laser scans,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 76, pp. 76–88, 2013.
- [4] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon, “Kinectfusion: Real-time dense surface mapping and tracking,” in *2011 10th IEEE International Symposium on Mixed and Augmented Reality*, pp. 127–136, 2011.
- [5] Z. Teed and J. Deng, “Droid-slam: Deep visual slam for monocular, stereo, and rgb-d cameras,” in *Advances in Neural Information Processing Systems* (M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, eds.), vol. 34, pp. 16558–16569, Curran Associates, Inc., 2021.
- [6] Y. Yuan, M. Bleier, and A. Nüchter, “Scenefactory: A workflow-centric and unified framework for incremental scene modeling,” 2025.
- [7] R. Murai, E. Dexheimer, and A. J. Davison, “Mast3r-slam: Real-time dense slam with 3d reconstruction priors,” 2025.
- [8] D. Maggio, H. Lim, and L. Carlone, “Vggt-slam: Dense rgb slam optimized on the  $sl(4)$  manifold,” 2025.
- [9] S. Wang, V. Leroy, Y. Cabon, B. Chidlovskii, and J. Revaud, “Dust3r: Geometric 3d vision made easy,” 2024.
- [10] J. Wang, M. Chen, N. Karaev, A. Vedaldi, C. Rupprecht, and D. Novotny, “Vggt: Visual geometry grounded transformer,” 2025.
- [11] D. Zhuo, W. Zheng, J. Guo, Y. Wu, J. Zhou, and J. Lu, “Streaming 4d visual geometry transformer,” 2025.
- [12] Y. Lan, Y. Luo, F. Hong, S. Zhou, H. Chen, Z. Lyu, S. Yang, B. Dai, C. C. Loy, and X. Pan, “Stream3r: Scalable sequential 3d reconstruction with causal transformer,” 2025.
- [13] Z. Zhu, S. Peng, V. Larsson, Z. Cui, M. R. Oswald, A. Geiger, and M. Pollefeys, “Nicer-slam: Neural implicit scene encoding for rgb slam,” 2023.
- [14] M. Bloesch, J. Czarowski, R. Clark, S. Leutenegger, and A. J. Davison, “Codeslam - learning a compact, optimisable representation for dense visual slam,” 2019.
- [15] J. Czarowski, T. Laidlow, R. Clark, and A. J. Davison, “Deepfactors: Real-time probabilistic dense monocular slam,” *IEEE Robotics and Automation Letters*, vol. 5, p. 721–728, Apr. 2020.
- [16] Y. Yao, Z. Luo, S. Li, T. Fang, and L. Quan, “Mvsnet: Depth inference for unstructured multi-view stereo,” 2018.
- [17] L. Koestler, N. Yang, N. Zeller, and D. Cremers, “Tandem: Tracking and dense mapping in real-time using deep multi-view stereo,” in *Conference on Robot Learning*, pp. 34–45, PMLR, 2022.
- [18] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” 2020.
- [19] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, “3d gaussian splatting for real-time radiance field rendering,” 2023.
- [20] A. Rosinol, J. J. Leonard, and L. Carlone, “Nerf-slam: Real-time dense monocular slam with neural radiance fields,” 2022.
- [21] C. Yan, D. Qu, D. Xu, B. Zhao, Z. Wang, D. Wang, and X. Li, “Gs-slam: Dense visual slam with 3d gaussian splatting,” 2024.
- [22] V. Leroy, Y. Cabon, and J. Revaud, “Grounding image matching in 3d with mast3r,” 2024.
- [23] J. Yang, A. Sax, K. J. Liang, M. Henaff, H. Tang, A. Cao, J. Chai, F. Meier, and M. Feiszli, “Fast3r: Towards 3d reconstruction of 1000+ images in one forward pass,” 2025.
- [24] Y. Wang, J. Zhou, H. Zhu, W. Chang, Y. Zhou, Z. Li, J. Chen, J. Pang, C. Shen, and T. He, “ $\pi^3$ : Scalable permutation-equivariant visual geometry learning,” 2025.
- [25] H. Wang and L. Agapito, “3d reconstruction with spatial memory,” 2024.
- [26] Q. Wang, Y. Zhang, A. Holynski, A. A. Efros, and A. Kanazawa, “Continuous 3d perception model with persistent state,” 2025.
- [27] Z. Chen, M. Qin, T. Yuan, Z. Liu, and H. Zhao, “Long3r: Long sequence streaming 3d reconstruction,” 2025.
- [28] Z. Teed and J. Deng, “Deepv2d: Video to depth with differentiable structure from motion,” 2020.
- [29] L. Lipson, Z. Teed, and J. Deng, “Deep patch visual slam,” 2024.
- [30] Y. Zhang, F. Tosi, S. Mattoccia, and M. Poggi, “Go-slam: Global optimization for consistent 3d instant reconstruction,” 2023.
- [31] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, “A benchmark for the evaluation of rgb-d slam systems,” in *Proc. of the International Conference on Intelligent Robot Systems (IROS)*, Oct. 2012.
- [32] B. Glocker, S. Izadi, J. Shotton, and A. Criminisi, “Real-time rgb-d camera relocalization,” in *International Symposium on Mixed and Augmented Reality (ISMAR)*, IEEE, October 2013.
- [33] Y. Liu, S. Dong, S. Wang, Y. Yin, Y. Yang, Q. Fan, and B. Chen, “Slam3r: Real-time dense scene reconstruction from monocular rgb videos,” 2025.
- [34] J. Straub, T. Whelan, L. Ma, Y. Chen, E. Wijmans, S. Green, J. J. Engel, R. Mur-Artal, C. Ren, S. Verma, A. Clarkson, M. Yan, B. Budge, Y. Yan, X. Pan, J. Yon, Y. Zou, K. Leon, N. Carter, J. Briales, T. Gillingham, E. Mueggler, L. Pesqueira, M. Savva, D. Batra, H. M. Strasdat, R. D. Nardi, M. Goesele, S. Lovegrove, and R. Newcombe, “The replica dataset: A digital replica of indoor spaces,” 2019.
- [35] G. Baruch, Z. Chen, A. Dehghan, T. Dimry, Y. Feigin, P. Fu, T. Gebauer, B. Joffe, D. Kurz, A. Schwartz, and E. Shulman, “Arkitscenes: A diverse real-world dataset for 3d indoor scene understanding using mobile rgb-d data,” 2022.
- [36] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner, “ScanNet: Richly-annotated 3d reconstructions of indoor scenes,” in *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2017.
- [37] C. Yeshwanth, Y.-C. Liu, M. Nießner, and A. Dai, “ScanNet++: A high-fidelity dataset of 3d indoor scenes,” in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2023.

- [38] M. Roberts, J. Ramapuram, A. Ranjan, A. Kumar, M. A. Bautista, N. Paczan, R. Webb, and J. M. Susskind, “Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding,” in *International Conference on Computer Vision (ICCV) 2021*, 2021.
- [39] Y. Yao, Z. Luo, S. Li, J. Zhang, Y. Ren, L. Zhou, T. Fang, and L. Quan, “Blendedmvs: A large-scale dataset for generalized multi-view stereo networks,” 2020.
- [40] Z. Li and N. Snavely, “Megadepth: Learning single-view depth prediction from internet photos,” in *Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [41] P.-H. Huang, K. Matzen, J. Kopf, N. Ahuja, and J.-B. Huang, “Deepmvs: Learning multi-view stereopsis,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [42] C. Campos, R. Elvira, J. J. G. Rodriguez, J. M. M. Montiel, and J. D. Tardos, “Orb-slam3: An accurate open-source library for visual, visual–inertial, and multimap slam,” *IEEE Transactions on Robotics*, vol. 37, p. 1874–1890, Dec. 2021.
- [43] H. Li, X. Gu, W. Yuan, L. Yang, Z. Dong, and P. Tan, “Dense rgb slam with neural implicit maps,” 2023.
- [44] A. Handa, T. Whelan, J. McDonald, and A. J. Davison, “A benchmark for rgb-d visual odometry, 3d reconstruction and slam,” in *2014 IEEE international conference on Robotics and automation (ICRA)*, pp. 1524–1531, IEEE, 2014.