

SST: Real-time End-to-end Monocular 3D Reconstruction via Sparse Spatial-Temporal Guidance

Chenyanguang Zhang^{1,*}, Zhiqiang Lou^{1,*}, Yan Di², Federico Tombari^{2,3}, and Xiangyang Ji¹

¹ Tsinghua University ² Technical University of Munich ³ Google
{zcyg22,lzq20}@mails.tsinghua.edu.cn, xyji@tsinghua.edu.cn

Abstract—Real-time monocular 3D reconstruction is a challenging problem that remains unsolved. Although recent end-to-end methods demonstrate promising results, tiny structures and geometric boundaries are hardly captured due to their insufficient supervision neglecting spatial details and oversimplified feature fusion ignoring temporal cues. To address the problems, we propose an end-to-end 3D reconstruction network SST, which utilizes Sparse estimated points from visual SLAM system as additional Spatial guidance and fuses Temporal features via a cross-modal attention mechanism, achieving more detailed reconstruction results. We propose a Local Spatial-Temporal Fusion module to exploit more informative spatial-temporal cues from multi-view color information and sparse priors, as well a Global Spatial-Temporal Fusion module to refine the local TSDF volumes with the world-frame model from coarse to fine. Extensive experiments on ScanNet and 7-Scenes demonstrate that SST outperforms all state-of-the-art competitors, whilst keeping a high inference speed at 59 FPS, enabling real-world applications with real-time requirements.

Index Terms—3D reconstruction, real time, visual SLAM guidance

I. INTRODUCTION

Monocular 3D scene reconstruction has aroused wide research interest due to its promising applications in AR/VR, robotic manipulation, and scene understanding. Unlike Lidar or RGB-D methods, monocular 3D reconstruction refers to predicting the 3D scene model from several consecutive frames, without direct distance measurements. Thus, it is severely ill-posed and remains unsolved yet.

Traditional visual SLAM systems can estimate ego-motion accurately, but can only provide coarse sparse reconstruction results due to high computational consumption and matching ambiguity. Recently, deep-learning-based methods have been proposed to achieve satisfactory dense reconstruction performance, which can be divided into two-stage and end-to-end methods. Two-stage methods [5]–[8], occupying the mainstream, first estimate per-frame depth and then fuse the back-projected point clouds. End-to-end methods directly regress 3D Truncated Signed Distance Function (TSDF) volumes via neural networks.

This work was supported by the National Key R&D Program of China under Grant 2018AAA0102801. Authors with * have equal contributions.

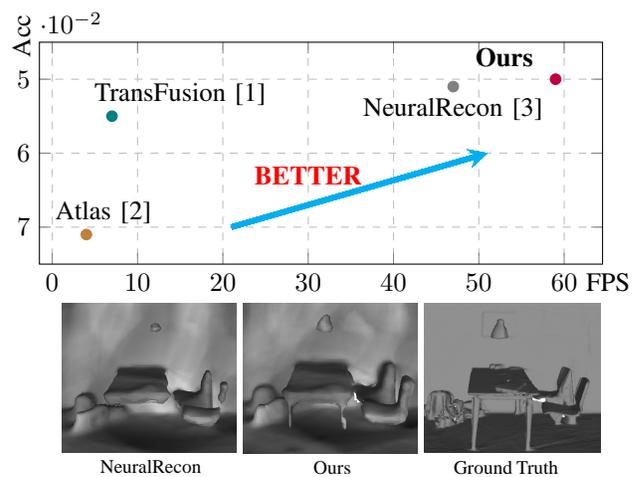


Fig. 1. Quantitative and qualitative comparison on ScanNet [4] benchmark. Compared to other competitors, our SST achieves best comprehensive performance, reconstructing more precise tiny structures and geometric boundaries.

For end-to-end methods, if well reconstructed geometric details are demanded, intolerable GPU memory cost from high-resolution geometry volumes emerges. Although sparse volume representation [3] has been proposed, current methods are still use insufficiently detailed voxel-based representation. Thus, they can only be coarsely supervised via ground-truth TSDF volume with large voxel size, losing many of geometry details compared to two-stage methods with dense depth supervision. To enhance the geometry information neglected by current end-to-end methods, we introduce the sparse depth input, which is a byproduct in ego-motion estimation process of visual SLAM system, and has been proved to contain heuristic geometric guidance [9], [10]. We harness a lightweight CNN to encode the sparse depth input, and then fuse this spatial guidance along with color features.

Besides the drawback of coarse supervision due to voxel scale restriction, current end-to-end methods [2], [3] cannot construct fine 3D feature volumes since they conduct simple average feature fusion, which neglects temporal connections from different observation views. Thus, considering simultaneously utilizing the sparse cross-modal guidance and the

multi-view temporal cues, we carefully design a novel sparse cross-modal attention mechanism. Integrating sparse spatial-temporal guidance, it is capable to extract more informative 3D feature volumes for accurate reconstruction.

We dedicate ourselves to a real-time monocular 3D reconstruction pipeline SST with the state-of-the-art performance and less computational resource consumption. The proposed SST follows a economic local-to-global pipeline consisting of three main modules: Feature Extraction (FE), Local Spatial-Temporal Fusion (LSTF) and Global Spatial-Temporal Fusion (GSTF). Given a fragment of images, corresponding camera poses and sparse depth measurements from real-time visual SLAM system, SST first extracts pixel-wise cross-modal features from images and sparse depth points respectively, via the FE module. Then, after back-projecting features into 3D space, the novel LSTF module distills and aggregates multi-view features locally using learned attention-aware weights to construct effective 3D feature volumes. Finally, for globally consistent reconstruction, we design a GSTF module to refine and fuse local TSDF volumes with global 3D model. We design a well-organized lightweight recurrent unit for GSTF, which plays an indispensable role to accelerate the whole pipeline. SST follows a coarse-to-fine manner, and the final reconstructed TSDF volumes are outputted at the highest level.

The main contributions are summarized as follows,

- We propose a novel end-to-end network SST for real-time monocular 3D scene reconstruction, which outperforms all real-time competitors and achieves comparable results to the state-of-the-art method but runs 8X faster at 59FPS on ScanNet [4] and 7-scenes [11] benchmarks.
- We investigate the shortcomings of previous end-to-end methods and present a novel LSTF module with proposed sparse cross-modal attention mechanism that allows for adaptive feature aggregation for color information and sparse point-based geometry priors, distilling more informative cues for accurate reconstruction.
- We design a lightweight recurrent unit for our GSTF module to efficiently fuse local TSDF volumes with global 3D model, reducing time and storage consumption considerably.

II. RELATED WORKS

A. 3D Scene Reconstruction

There are two streams of methods handling 3D scene reconstruction problem, two-stage and end-to-end methods. Two-stage methods follow a long-standing reconstruction pipeline that first estimates depth map of each frame and then fuses the back-projected point clouds. [12]–[14] are popular traditional methods with a patch matching based pipeline. Learning-based approaches [5], [6], [8], [15], [16] often build a shared 3D cost volume in the reference camera space using feature averaging. However, these two-stage methods still confront with spatial inconsistency and redundant computation problems. Atlas [2] is the first end-to-end scene reconstruction approach, proposing a global averaging feature volume, but with low efficiency

and not incrementally. TransformerFusion [1] and NeuralRecon [3] are two impressive incremental methods. [1] constructs a global feature volume, gaining superior results but suffering from more time and storage consumption. [3] extracts feature volumes locally and then fuses to the global. Different from them, our SST fuses multi-frame spatial-temporal features in local coordinate system and then converts them to the global, which is demonstrated to be efficient for 3D reconstruction.

B. Sparse Geometry Feature Fusion

Some works [9], [10] illustrate noisy sparse points can enhance monocular depth estimation results. They design up-sampling networks to fuse sparse geometry features and RGB features, aiming at producing a completed dense depth map. Other related works concentrate on 3D object detection [17]–[20] or end-to-end autopilot control task [21]. [17]–[19] focus on unifying representation for multi-modal feature, whether voxel-based or point-based. We have not found any mature pipeline with respect to fusing sparse geometry information for end-to-end 3D scene reconstruction task. Our proposed method fills in the gaps of this field and gains impressive results considering both accuracy and speed.

III. METHOD

A. Overview

As illustrated in Fig. 2, given a monocular image sequence $\{I_t\}$, SST first utilizes a visual SLAM system [22] to generate noisy sparse geometry priors $\{G_t\}$ and corresponding camera poses $\{T_t\}$ as the network inputs. $\{G_t\}$ consist of projected sparse depth maps $\{SD_t\}$ and corresponding reprojection error maps $\{E_t\}$. We convert $\{E_t\}$ to confidence maps $\{CO_t\}$ via $CO_t = \exp(-\lambda E_t)$, and then concatenate $\{CO_t\}$ and $\{SD_t\}$ as final geometric priors $\{G_t\}$.

Our SST, consisting of Feature Extraction (FE), Local Spatial-Temporal Fusion (LSTF) and Global Spatial-Temporal Fusion (GSTF) modules, incrementally receives $\{I_t\}$, $\{G_t\}$ and $\{T_t\}$, and reconstructs accurate dense 3D geometry structures in real time. Specifically, the input sequence will be split into several fragments F_l and each fragment contains N frames, predefined as 9 during our implementation. FE extracts and back-projects features from I_t and G_t for each key frame in the sequence to obtain raw feature volumes FV_t^{IG} . Then LSTF collects features $\{FV_t^{IG}\}_N$ in the current fragment F_l and fuses them into a local feature volume FV_l^{frag} via proposed sparse cross-modal attention mechanism, which strengthens the mixed features by the adaptive weights considering both temporal and spatial guidance. Finally, FV_l^{frag} will be efficiently updated to the global feature volume FV^{global} in GSTF with the lightweight recurrent unit, followed by simple feed-forward layers to output final TSDF volumes. Inspired by [3], our network repeats the above inference pipeline in a coarse-to-fine manner with three scales and utilizes coarse predictions to guide the high-resolution 3D feature volume for less computation cost.

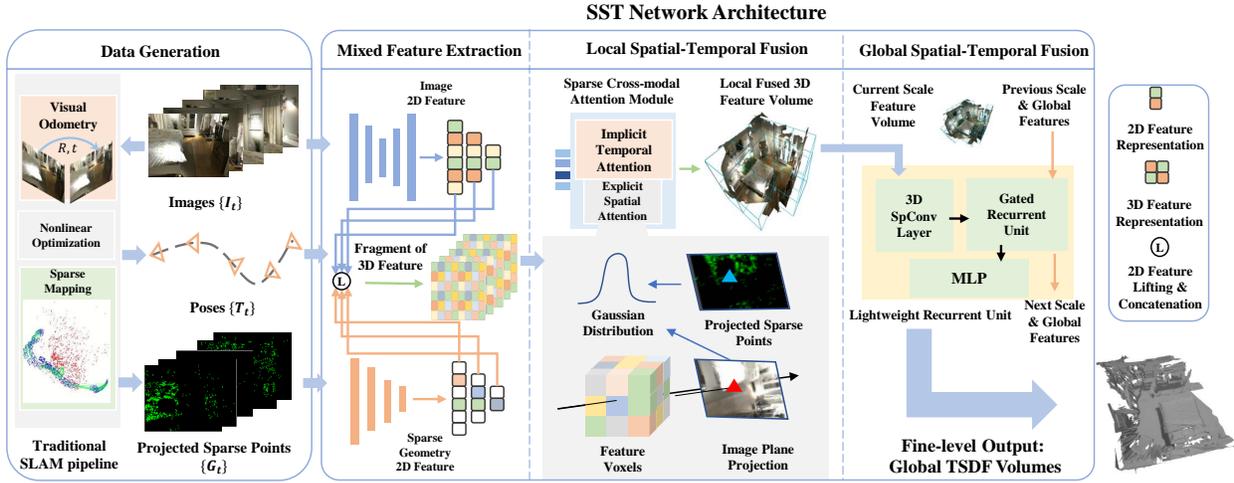


Fig. 2. Pipeline Overview of SST.

B. Feature Extraction

Given a color image I_t and corresponding geometry prior G_t , the FE module generates mixed 3D feature volume FV_t^{IG} which contains rich semantic cues from I_t and spatial information from G_t . According to [23], we select different feature extraction backbones for I_t and G_t respectively because of their different modalities. Concretely, a lightweight MnasNet [24] is applied for image feature $F_{t,color}$ extraction. For informative geometry prior input G_t , we design a lightweight CNN composed of 4 cascaded convolution blocks to output implicit geometry feature $F_{t,geo}$ with 8 channels.

We then back-project $F_{t,geo}$ and $F_{t,color}$ into 3D space to construct feature volume FV_t^{IG} for key frame I_t ,

$$FV_t^{IG} = \text{Backproj}(\text{Cat}(F_{t,color}, F_{t,geo})) \quad (1)$$

where Cat means channel-wise concatenation and Backproj denotes 2D-3D back-projection operation from the image plane to the local frame coordinate. Note that FV_t^{IG} , $F_{t,color}$ and $F_{t,geo}$ are all calculated at 3 resolution scales for our coarse-to-fine inference framework.

C. Local Spatial-Temporal Fusion

Given l -th fragment's 3D features $\{FV_t^{IG}\}_N$ from N consecutive frames $\{I_t\}_N$, the task of LSTF is to output an effective 3D feature volume FV_l^{frag} for current local fragment, similar to local window optimization in traditional SLAM pipelines.

For voxel v_i in local frame coordinate system, previous methods [2], [3] directly generate the voxel feature FV_{l,v_i}^{frag} via simple averaging as $\frac{1}{N} \sum_{t=1}^N FV_{t,v_i}^{IG}$. This simple averaging operation loses most spatial information from multi-view features since it cannot learn to attend to important parts for surface reconstruction, i.e. corners, vertical lines, planes etc. Also, the averaging operation treats each frame in a fragment equally, neglecting temporal correlation among multi-view features. Furthermore, considering the input geometry

priors, if using direct averaging, $F_{t,geo}$ and $F_{t,color}$ will be processed separately without any interaction, which increases the difficulty of subsequent module to learn informative fusion features for accurate reconstruction. To handle these shortcomings, as shown in Fig. 3, we propose LSTF that utilizes the proposed sparse cross-modal attention mechanism, enabling both adaptive weighted feature fusion in temporal dimension and channel-wise multi-modal feature interaction for spatial feature fusion. The inference procedure is shown as,

$$\begin{aligned} A_{in} &= [FV_{1,v_i}^{IG}, \dots, FV_{N,v_i}^{IG}] \\ Q &= W_q A_{in}, \quad K = W_k A_{in}, \quad V = W_v A_{in} \\ \omega_{l,im} &= \text{Softmax}(QK^T), \quad A_{out} = \omega_{l,im} \omega_{l,ex} V \end{aligned} \quad (2)$$

where A_{in}, A_{out}, W_* denote the input, output, and learnable parameters of the attention module. The final attention of the sparse cross-modal attention mechanism is the mixture of implicit temporal attention $\omega_{l,im}$ and explicit spatial attention $\omega_{l,ex}$. To compute $\omega_{l,ex}$, for the weight $\omega_{l,i,t}$ of voxel v_i in the t -th view of this fragment, we first project v_i on the image plane of the t -th view to get the projected depth d_{v_i} and pixel position p_{v_i} , and $\omega_{l,i,t}$ is computed as

$$\omega_{l,i,t} = \begin{cases} \text{Gauss}_{\sigma_{Exp}(p_{v_i})}(\|SD_t(p_{v_i}) - d_{v_i}\|) & SD_t(p_{v_i}) \geq 0 \\ 1 & \text{otherwise} \end{cases} \quad (3)$$

where $SD_t(p_{v_i})$ means corresponding sparse depth value of p_{v_i} in SD_t . $\omega_{l,ex}$ tends to assign higher weight for voxels near sparse points and the learned adaptive weight $\omega_{l,im}$ can attenuate the influence of unnecessary features and persevere important cues, together assisting in constructing effective feature volumes. The value weight matrix W_v further mixes features channels from geometry priors and color images, integrating sparse spatial information into final output features, conducting spatial fusion process. The fragment feature FV_{l,v_i}^{frag} of voxel v_i is computed as $FV_{l,v_i}^{frag} = \text{LSTF}(FV_{1,v_i}^{IG}, \dots, FV_{N,v_i}^{IG})$.

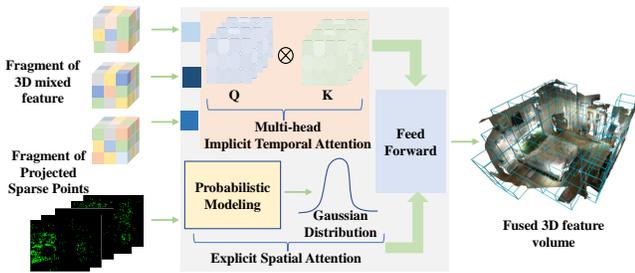


Fig. 3. Details of the proposed LSTF.

D. Global Spatial-Temporal Fusion

To enable consistent incremental reconstruction, we adopt a proposed lightweight recurrent unit for fusing current fragment feature volume FV_l^{frag} into previous established global feature volume, conceptually similar to global optimization in traditional SLAM pipelines. Specifically, the lightweight recurrent unit is a 3D sparse convolutional variant of Gated Recurrent Unit (GRU). It maintains a global hidden state volume $H_{l,t-1}^{global}$, and updates it as follows. First, FV_l^{frag} is input to 3D sparse convolution layers for extracting its corresponding surface geometry feature $S_{l,t}^{frag}$. Then hidden state volume $H_{l,t-1}^{frag}$ in local coordinate system is generated from $H_{l,t-1}^{global}$ and fused with $S_{l,t}^{frag}$, resulting in updated local hidden state volume $H_{l,t}^{frag}$ then transformed into the global volume $H_{l,t}^{global}$. Two simple MLP layers are applied to $H_{l,t}^{global}$ to estimate occupancy grid O_l^{global} and TSDF volume $TSDF_l^{global}$ respectively. The whole pipeline is shown in Fig. 4 (c).

During implementation, we find the inference time bottleneck lies in the structure of 3D sparse convolution layers. Thus, we design a new lightweight structure with much higher inference speed compared to [3], in which the 3D sparse convolution layers are redundant and less organized. We refer to the implementation adopted by [25] but significantly modify the network structure and reduce the number of parameters. The concrete structure of the 3D sparse convolution (3D SpConv) unit design is exhibited in Fig. 4 (a). Fig. 4 (b) shows the detail of the 3D SpConv residual block, which is simplified by us compared with [25]. In practice, the novel lightweight recurrent unit maintains the quality of reconstruction but distinctively reduces inference time, playing a pivotal role in accelerating SST network.

IV. EXPERIMENTS

A. Experimental Setup

Datasets and Baselines. The experiments are performed on two common benchmarks, ScanNet [4] and 7-Scenes [11]. We use the same training and validation splits with previous methods [2] [3] [1]. Following [26] [3], we directly apply our model trained on ScanNet for evaluation on 7-Scenes to validate the generalization ability. We compare our method with two-stage MVS methods and end-to-end reconstruction

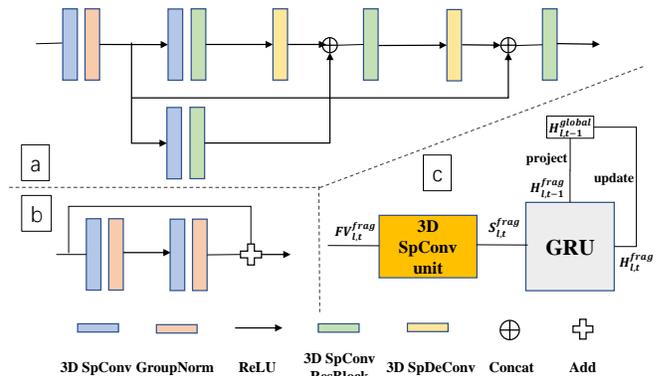


Fig. 4. Details of the GSTF module.

methods. The two-stage methods contain MVDNet [7], GP-MVS [16], DPSNet [8] and COLMAP [12], while the end-to-end methods include non-real-time Atlas [2], TransFusion [1] and real-time NeuralRecon [3].

Metrics. Current real-time state-of-the-art method [3] and its non-real-time counterpart [1] follow different 3D metrics to evaluate their methods on ScanNet. Thereby, we need to make comparison to them under two corresponding 3D metrics respectively. The 3D metrics exhibited in the top block of Tab. I follow [3]. In the bottom block of Tab. I, we present metrics following [1]. For evaluation under 2D depth metrics, we render the reconstructed mesh to the image plane and obtain depth estimations. All 2D depth metrics for ScanNet evaluation in Tab. II are defined in [3] [27]. For metrics in the generalization validation on 7-Scenes, we also follow metrics in [3] as presented in Tab. III.

Implementation Details. In the FE module, the MnasNet extracts color feature with 80, 40 and 24 channels for three coarse-to-fine scales. We initialize the MnasNet via weights pretrained from ImageNet [28]. In the LSTF module, after the calculation for sparse cross-modal attention, the used feed-forward layer remains the channels of output local fragment feature volume with 80, 40 and 24. The lightweight 3D sparse convolution utilized in GSTF is implemented by torchsparse [25] due to its best synthetic performance. For sparse depth input, we ignore pixels with depth deeper than 3.0m due to their large noise. For the voxel output, we predict the occupancy via a Sigmoid layer and the TSDF value via a full-connection layer from the global feature volume. As in previous end-to-end methods [1]–[3], the voxel size of the finest level is 4cm and TSDF truncation distance λ is set to 12cm. We adopt the binary cross-entropy loss for the predicted occupancy score and l_1 loss for the log-transform results of regressed TSDF results. Both loss functions are applied in the three coarse-to-fine levels.

B. Evaluation Results

ScanNet. 2D depth metrics and 3D geometry metrics are evaluated on the ScanNet [4] dataset. Tab. I exhibits quantita-

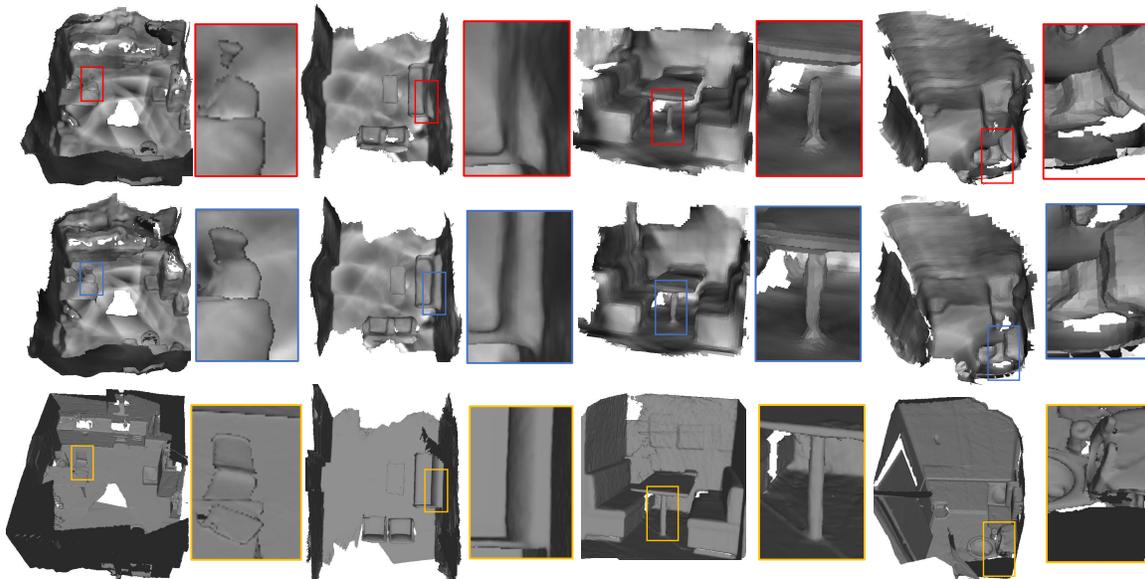


Fig. 5. Qualitative comparison on ScanNet. From top to down are the results of NeuralRecon [3], our results and ground truth respectively. Compared to real-time state-of-the-art [3], our SST presents more accurate and coherent reconstruction results.

TABLE I
3D GEOMETRY METRICS ON SCANNET.

Method	Comp	Acc	Recall	Prec	F-score	FPS
MVDNet [7]	0.040	0.240	0.831	0.208	0.329	28
GPMVS [16]	0.031	0.879	0.871	0.188	0.304	27
DPSNet [8]	0.045	0.284	0.793	0.223	0.344	4
COLMAP [12]	0.069	0.135	0.634	0.505	0.558	0.4
NeuralRecon [3]	0.138	0.053	0.472	0.687	0.559	47
Ours	0.124	0.053	0.505	0.695	0.584	59
TransFusion [1]	0.082	0.055	0.600	0.728	0.655	7
Atlas [2]	0.076	0.071	0.605	0.675	0.636	4
NeuralRecon [3]	0.075	0.051	0.556	0.706	0.621	47
Ours	0.071	0.050	0.584	0.714	0.643	59

TABLE II
2D DEPTH METRICS ON SCANNET.

Method	Abs.Rel.	Abs.Diff.	Sq.Rel.	RMSE	$\delta < 1.25$
MVDNet [7]	0.098	0.191	0.061	0.293	89.6
GPMVS [16]	0.130	0.239	0.339	0.472	90.6
DPSNet [8]	0.087	0.158	0.035	0.232	92.5
COLMAP [12]	0.137	0.264	0.138	0.502	83.4
Atlas [2]	0.065	0.123	0.045	0.251	93.6
NeuralRecon [3]	0.065	0.099	0.034	0.197	93.7
Ours	0.060	0.092	0.034	0.185	94.0

tive comparison under two categories of 3D geometry metrics followed [3] and [1] respectively.

For 3D metrics consistent with [3], our SST yields best Accuracy, Precision and F-score. Compared with recent state-of-the-art real-time method [3], SST surpasses in all 3D metrics, exceeding [3] in F-score, Precision and Recall by 2.5%, 0.8% and 3.3% respectively. Compared with two-stage methods, SST surpasses the best two-stage method [12] in Accuracy by nearly 61%, with the promotion of Precision by 19.0%. For 3D metrics consistent with [1], SST yields best

Completeness and Accuracy as well as second best Precision and F-score. We surpass the state-of-the-art real-time method [3] (47 FPS) in all 3D metrics. SST exceeds [3] in F-score, Precision and Recall by 2.2%, 0.8% and 2.8% respectively. Compared with non-real-time method [2] (4 FPS), we get a slightly inferior Recall since it is an off-line method with advantage of having a global context to complete their previous TSDF predictions. For Precision and F-score, SST surpasses [2] by 3.9% and 0.7% respectively. Compared with state-of-the-art non-real-time method [1] (7 FPS), we surpass it in Completeness and Accuracy by 13.4% and 9.1%, while get slightly inferior Precision and Recall by only 1.6% and 1.4%. However, [1] depends on many cascaded transformer blocks, resulting in great storage and non-real-time performance (7 FPS), while our SST maintains lightweight architecture and real-time inference performance (59 FPS).

For 2D depth metrics shown in Tab. II, SST outperforms all end-to-end and two-stage competitors. Compared with the state-of-the-art real-time end-to-end method [3], SST exceeds in Abs. Rel. by 7.7%, Abs. Diff. by 7.1% and RMSE by 6.1%. When compared with two-stage methods, our method also has significantly better performance.

We also compare visualization results in Fig. 5, illustrating our method can reconstruct more accurate and consistent surface geometry than previous state-of-the-art real-time method NeuralRecon [3], especially at regions of tiny structures, e.g. the unbroken table leg and the vertical back of chair.

7-Scenes. To demonstrate the generalization ability of SST, we utilize 7-Scenes dataset to evaluate our model trained on ScanNet. Our method still achieves outperforming performance to the state-of-the-art real-time method [3] (Tab. III).

Efficiency. In Tab. I, we report average running frames-per-second (FPS) measured on an NVIDIA RTX 3090 GPU

TABLE III
3D METRICS ON 7-SCENES.

Method	Comp	Acc	Recall	Prec	F-score	FPS
NeuralRecon [3]	0.228	0.100	0.228	0.389	0.282	47
Ours	0.225	0.104	0.242	0.392	0.298	59

TABLE IV
LSTF ,FE AND GSTF ARCHITECTURE ABLATIONS ON SCANNet UNDER
3D METRICS ALONG WITH THE TOP BLOCK OF TAB. I.

	SCAM	Geometry Priors	LWRU	Recall	Prec	F-score	FPS
a	×	×	×	0.472	0.687	0.559	47
b	✓	×	×	0.494	0.690	0.574	41
c	✓	×	✓	0.495	0.695	0.576	62
d	✓	✓	×	0.496	0.695	0.579	38
e	✓	✓	✓	0.505	0.695	0.584	59

for all methods. Our method achieves the highest inference speed of 59 FPS among all competitors (8X faster than the state-of-the-art method [1]), enabling applications with real-time requirement, e.g. AR on mobile devices.

C. Ablation Studies

In this section, we conduct ablation experiments on ScanNet to demonstrate the effectiveness of our proposed architecture for making full use of spatial-temporal information for 3D reconstruction. The proposed sparse cross-modal attention mechanism (SCAM) in LSTF, sparse geometry priors encoder in FE and the lightweight recurrent unit (LWRU) in GSTF are three main concerns of ablations we conduct. Results shown in Tab. IV demonstrate significant effectiveness of three modules above.

Comparing (a) (b), our pipeline with the proposed SCAM yields higher 3D geometry metrics, promoting F-score by 1.7%. The geometry priors encoder plays another important role in exploiting the sparse spatial guidance. Concretely, comparing (c) (e), we find that encoding the sparse geometry priors by our proposed CNN promotes the F-score by 1.0%. The proposed LWRU plays a pivotal role in inference speed enhancement of SST, as is illustrated in (d) and (e). With the design of our LWRU, the promotion of FPS reaches near 50%, and the F-score is slightly enhanced simultaneously.

V. CONCLUSIONS

We introduce a real-time monocular 3D reconstruction network SST, yielding accurate reconstruction results whilst keeping high inference speed. The key idea is to fully utilize spatial-temporal guidance provided by multi-view features and sparse geometry priors. We design a LSTF module to treat multi-view feature attentively, simultaneously fusing sparse spatial information, and a efficient GSTF module to maintain global consistency of reconstruction. Extensive experiments demonstrate that our SST network achieves comparable results to the state-of-the-art method but runs 8X faster. We believe that our local-to-global method with sparse spatial-temporal

guidance can inspire future researches aiming at end-to-end 3D reconstruction task.

REFERENCES

- [1] A. Bozic, P. Palafox, J. Thies *et al.*, “Transformerfusion: Monocular rgb scene reconstruction using transformers,” *NeurIPS*, vol. 34, 2021.
- [2] Z. Murez, T. van As, J. Bartolozzi *et al.*, “Atlas: End-to-end 3d scene reconstruction from posed images,” in *ECCV*, 2020, pp. 414–431.
- [3] J. Sun, Y. Xie, L. Chen *et al.*, “Neuralrecon: Real-time coherent 3d reconstruction from monocular video,” in *CVPR*, 2021, pp. 15598–15607.
- [4] A. Dai, A. X. Chang, M. Savva *et al.*, “ScanNet: Richly-annotated 3d reconstructions of indoor scenes,” in *CVPR*, 2017, pp. 5828–5839.
- [5] A. Duzceker, S. Galliani, C. Vogel *et al.*, “Deepvideomvs: Multi-view stereo on video with recurrent spatio-temporal fusion,” in *CVPR*, 2021, pp. 15324–15333.
- [6] X. Long, L. Liu, W. Li *et al.*, “Multi-view depth estimation using epipolar spatio-temporal networks,” in *CVPR*, 2021, pp. 8258–8267.
- [7] K. Wang and S. Shen, “Mvdepthnet: Real-time multiview depth estimation neural network,” in *3DV*, 2018, pp. 248–257.
- [8] S. Im, H. Jeon, S. Lin *et al.*, “Dpsnet: End-to-end deep plane sweep stereo,” *arXiv preprint arXiv:1905.00538*, 2019.
- [9] F. Ma and S. Karaman, “Sparse-to-dense: Depth prediction from sparse depth samples and a single image,” in *ICRA*, 2018, pp. 4796–4803.
- [10] Y. Huang, Y. Liu, T. Wu *et al.*, “S3: Learnable sparse signal superdensity for guided depth estimation,” in *CVPR*, 2021, pp. 16706–16716.
- [11] J. Shotton, B. Glocker, C. Zach *et al.*, “Scene coordinate regression forests for camera relocalization in rgb-d images,” in *CVPR*, 2013, pp. 2930–2937.
- [12] J. Schönberger, E. Zheng, J. Frahm *et al.*, “Pixelwise view selection for unstructured multi-view stereo,” in *ECCV*, 2016, pp. 501–518.
- [13] Y. Di, H. Morimitsu, S. Gao, and X. Ji, “Monocular piecewise depth estimation in dynamic scenes by exploiting superpixel relations,” in *ICCV*, 2019, pp. 4363–4372.
- [14] Y. Di, H. Morimitsu, Z. Lou, and X. Ji, “A unified framework for piecewise semantic reconstruction in dynamic scenes via exploiting superpixel relations,” in *ICRA*, 2020, pp. 10737–10743.
- [15] C. Liu, J. Gu, K. Kim *et al.*, “Neural rgb (r) d sensing: Depth and uncertainty from a video camera,” in *CVPR*, 2019, pp. 10986–10995.
- [16] Y. Hou, J. Kannala, and A. Solin, “Multi-view stereo by temporal nonparametric fusion,” in *ICCV*, 2019, pp. 2651–2660.
- [17] J. Yoo, Y. Kim, J. Kim *et al.*, “3d-cvf: Generating joint camera and lidar features using cross-view spatial feature fusion for 3d object detection,” in *ECCV*, 2020, pp. 720–736.
- [18] T. Huang, Z. Liu, X. Chen *et al.*, “Epnnet: Enhancing point features with image semantics for 3d object detection,” in *ECCV*, 2020, pp. 35–52.
- [19] L. Xie, C. Xiang, Z. Yu *et al.*, “Pi-rcnn: An efficient multi-sensor 3d object detector with point-based attentive cont-conv fusion module,” in *AAAI*, vol. 34, no. 07, 2020, pp. 12460–12467.
- [20] Y. Su, Y. Di, G. Zhai, F. Manhardt, J. Rambach, B. Busam, D. Stricker, and F. Tombari, “Opa-3d: Occlusion-aware pixel-wise aggregation for monocular 3d object detection,” *IEEE Robotics and Automation Letters*, 2023.
- [21] A. Prakash, K. Chitta, and A. Geiger, “Multi-modal fusion transformer for end-to-end autonomous driving,” in *CVPR*, 2021, pp. 7077–7087.
- [22] R. Mur-Artal and J. D. Tardós, “Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras,” *IEEE transactions on robotics*, vol. 33, no. 5, pp. 1255–1262, 2017.
- [23] N. Merrill, P. Geneva, and G. Huang, “Robust monocular visual-inertial depth completion for embedded systems,” in *ICRA*, 2021, pp. 5713–5719.
- [24] M. Tan, B. Chen, R. Pang *et al.*, “Mnasnet: Platform-aware neural architecture search for mobile,” in *CVPR*, 2019, pp. 2820–2828.
- [25] H. Tang, Z. Liu, S. Zhao *et al.*, “Searching efficient 3d architectures with sparse point-voxel convolution,” in *ECCV*, 2020, pp. 685–702.
- [26] X. Long, L. Liu, C. Theobalt *et al.*, “Occlusion-aware depth estimation with adaptive normal constraints,” in *ECCV*, 2020, pp. 640–657.
- [27] D. Eigen, C. Puhrsch, and R. Fergus, “Depth map prediction from a single image using a multi-scale deep network,” *NeurIPS*, vol. 27, 2014.
- [28] J. Deng, W. Dong, R. Socher *et al.*, “Imagenet: A large-scale hierarchical image database,” in *CVPR*, 2009, pp. 248–255.