

---

# Supplementary Materials: Video World Models with Long-term Spatial Memory

---

Tong Wu<sup>\*1</sup>, Shuai Yang<sup>\*2,4</sup>, Ryan Po<sup>1</sup>, Yinghao Xu<sup>1</sup>,  
Ziwei Liu<sup>5</sup>, Dahua Lin<sup>3,4</sup>, Gordon Wetzstein<sup>1</sup>

<sup>1</sup> Stanford University   <sup>2</sup> Shanghai Jiao Tong University   <sup>3</sup> The Chinese University of Hong Kong

<sup>4</sup> Shanghai Artificial Intelligence Laboratory   <sup>5</sup> S-Lab, Nanyang Technological University

<https://spmем.github.io/>

In the supplementary material, we provide additional implementation details in Section A, extended experimental results in Section B, and further discussions on related approaches in Section C. We also include more video examples in **video.mp4** to better illustrate the effectiveness and qualitative performance of our method.

## A Additional details in methodology

**Autoregressive point cloud fusion.** In autoregressive generation, the length of the generated video increases over time. At each step, a new static points map is produced and updated into our spatial memory. Since Mega-SAM performs reconstruction in the NDC coordinate system, it is impossible to directly merge results from different stages without alignment, and long video inference will also fail due to CUDA memory limitations. Therefore, unlike the precise reconstruction using Mega-SAM during the data construction stage, we employ CUT3R [5] for 4D reconstruction during the inference stage. CUT3R is a unified online 3D perception framework featuring a stateful recurrent model that incrementally updates a persistent internal state representation with each new observation. Given an image stream (video or unordered photos), the model simultaneously updates its state and predicts metric-scale pointmaps (per-pixel 3D points in a shared world coordinate system) and camera parameters for each input in an online manner. At each inference step, we save the state dict of the current CUT3R model and the parameters of the pose retriever to serve as initialization for the next inference step, ensuring that the reconstruction results of each step remain within the same coordinate system. Therefore, as shown in Figure S1, after tsdf-fusion, the filtered points cloud of the current step can be directly merged and aligned with the previous spatial static memory to achieve autoregressive point cloud fusion.

**Details in static point extraction for the dataset.** In Mega-SAM, we first resize the input video resolution to 384×672. Based on the initial results, we perform optical flow estimation to refine the estimated camera motion through pixel-level motion cues. Subsequently, a Covariance-based Variable Decomposition strategy is employed to further enhance the robustness and accuracy of the predicted results. In TSDF-fusion, we compute the initial grid dimensions based on the current voxel size and the scene bounds. If the maximum dimension exceeds a predefined threshold (1200), we proportionally scale the voxel size to ensure that the grid dimensions remain within the limit. The final adjusted voxel size is returned. The core principle is to control the grid resolution to prevent memory overflow.

## B Additional experimental results

**User study setup.** We conducted a user study to evaluate the perceptual quality of our generated videos. Specifically, we invited 13 graduate students who are actively engaged in video generation

\* denotes equal contribution.

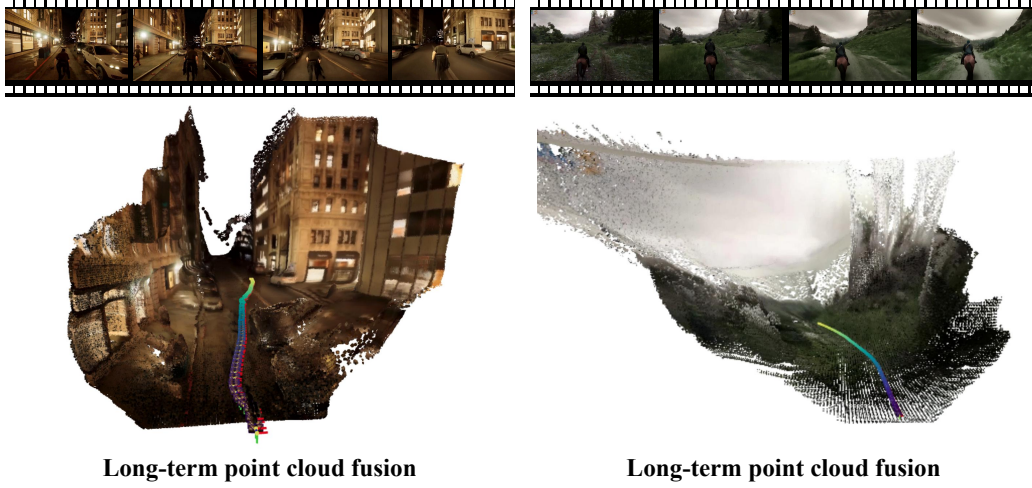


Figure S1: **Autoregressive point cloud fusion.** The system continuously updates spatial memory by integrating newly observed static maps. These maps are reconstructed online using CUT3R in a recurrent manner, while TSDF-Fusion filters out dynamic elements to maintain map consistency.

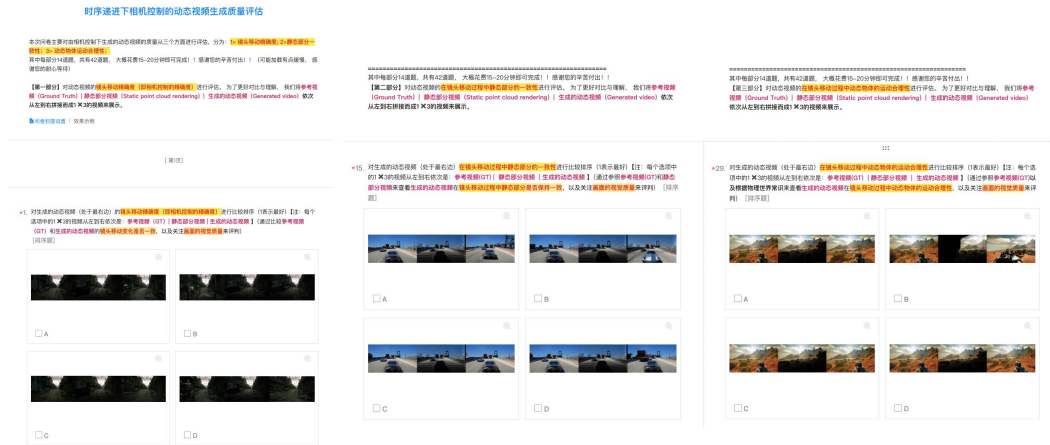


Figure S2: **User study questionnaire screenshots.**

research to participate in the evaluation. Each participant was presented with pairs of ground-truth and generated videos, and asked to rank them based on visual realism and temporal consistency. We report the average human ranking scores derived from their responses to reflect overall perceptual preference. To ensure transparency and clarity, we provide Figure S2, which shows partial screenshots of the user study questionnaire, and Figure S3, which presents partial ranking data for reference.

**Camera accuracy measurement.** We further evaluate camera pose accuracy along the reverse trajectory. Specifically, in our setting, each pose is visited twice during the sequence: once in the forward pass and once in the reverse, where ideally, the estimated camera poses should match. However, the presence of dynamic elements in the scene poses challenges for traditional SfM methods such as COLMAP. Therefore, we employ CUT3R to estimate the camera trajectories. To quantify the accuracy of the reverse trajectory, we compute both the rotation and translation errors between corresponding camera poses using the following metrics:

$$\text{RotErr} = \frac{1}{n} \frac{180}{\pi} \cdot \sum_{i=1}^n \arccos \left( \frac{\text{tr}(\mathbf{R}_0^i \mathbf{R}_{\text{rev}}^{i\top}) - 1}{2} \right), \quad (1)$$



without iterative optimization. To further process dynamic scenes, more recent advances [72, 5] further extend learning-based 3D reconstruction to dynamic scenes by incorporating recurrent models or persistent state, which improves efficiency and robustness in handling moving objects and changing environments. Some of these advances make it feasible to perform online estimation of camera poses and point clouds during video capture or generation.

In this work, we use Mega-SAM during dataset construction to obtain more stable and accurate point cloud and camera pose estimations. However, due to concerns about time efficiency and global alignment, we adopt CUT3R in our iterative video generation process. CUT3R jointly estimates per-frame point clouds while incorporating a dynamic-static disentanglement mechanism to preserve long-term static scene memory. Experimental results indicate that the difference in models does not lead to a significant performance gap.

## References

- [1] Jianhong Bai, Menghan Xia, Xiao Fu, Xintao Wang, Lianrui Mu, Jinwen Cao, Zuozhu Liu, Haoji Hu, Xiang Bai, Pengfei Wan, and Di Zhang. Recammaster: Camera-controlled generative rendering from a single video, 2025. URL <https://arxiv.org/abs/2503.11647>.
- [2] Hao He, Yinghao Xu, Yuwei Guo, Gordon Wetzstein, Bo Dai, Hongsheng Li, and Ceyuan Yang. Cameractrl: Enabling camera control for video diffusion models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=Z4ev0UYrk7>.
- [3] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [4] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny. Vggg: Visual geometry grounded transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025.
- [5] Qianqian Wang\*, Yifei Zhang\*, Aleksander Holynski, Alexei A. Efros, and Angjoo Kanazawa. Continuous 3d perception model with persistent state. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025.
- [6] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20697–20709, June 2024.
- [7] Junyi Zhang, Charles Herrmann, Junhwa Hur, Varun Jampani, Trevor Darrell, Forrester Cole, Deqing Sun, and Ming-Hsuan Yang. MonST3r: A simple approach for estimating geometry in the presence of motion. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=1JpqxFgWCM>.
- [8] Shangzhan Zhang, Jianyuan Wang, Yinghao Xu, Nan Xue, Christian Rupprecht, Xiaowei Zhou, Yujun Shen, and Gordon Wetzstein. Flare: Feed-forward geometry, appearance and camera estimation from uncalibrated sparse views. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025.