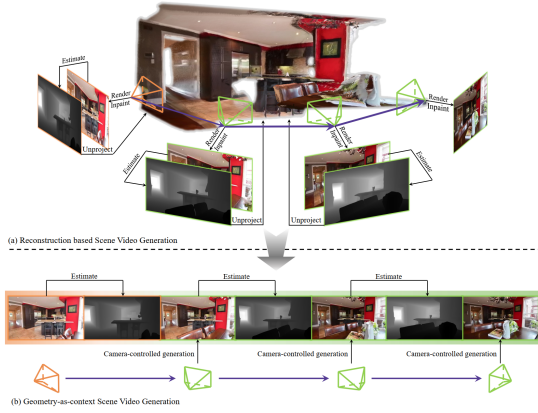


MILabCVPR-26

MILab **Geometry-as-ContextGaC** *Geometry-as-context: Modulating Explicit 3D in Scene-consistent Video Generation to Geometry Context CVPR2026*

Geometry-as-Context



Geometry-as-Context

Method	RealEstate10K [49]						Tanks-and-Temples [17]					
	PSNR ↑	SSIM ↑	LPIPS ↓	FID ↓	$R_{err} ↓$	$T_{err} ↓$	PSNR ↑	SSIM ↑	LPIPS ↓	FID ↓	$R_{err} ↓$	$T_{err} ↓$
CameraCtrl [9]	15.93	0.550	0.440	77.93	0.020	0.266	12.28	0.389	0.642	201.72	0.080	0.533
ViewCrafter [47]	16.72	0.585	0.417	80.47	0.022	0.327	12.59	0.438	0.549	115.16	0.058	0.399
GEN3C [27]	18.12	0.624	0.402	66.20	0.027	0.344	15.32	0.506	0.509	90.52	0.074	0.430
SEVA [48]	17.42	0.601	0.431	73.95	0.021	0.284	15.60	0.518	0.490	88.76	0.072	0.477
Voyager [14]	18.70	0.616	0.395	65.12	0.035	0.596	15.24	0.487	0.510	90.06	0.081	0.515
GaC (Ours)	19.01	0.656	0.354	55.76	0.024	0.270	15.77	0.532	0.507	93.29	0.072	0.442

Table 1. Quantitative results of scene video generation from single view with given camera trajectory.

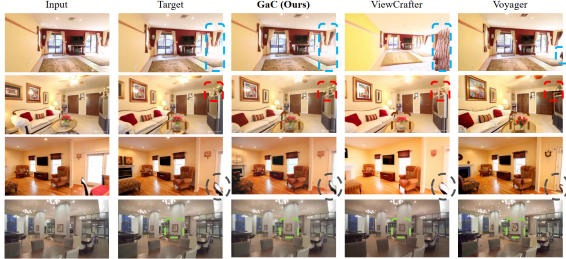
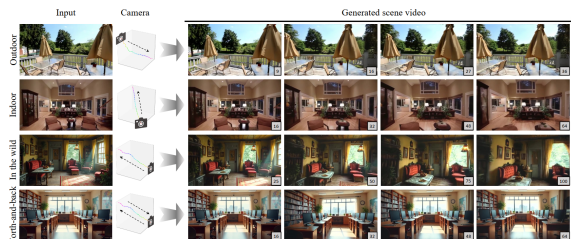


Figure 4. Qualitative results of scene video generation from single view. Compared to the baselines, our model generates more consistent novel views. These images are the 20-th frame of the generated video clip except the input one. For a clearer visualization, please zoom in.

Method	RealEstate10K [49]						Tanks-and-Temples [17]					
	PSNR ↑	SSIM ↑	LPIPS ↓	FID ↓	$R_{err} ↓$	$T_{err} ↓$	PSNR ↑	SSIM ↑	LPIPS ↓	FID ↓	$R_{err} ↓$	$T_{err} ↓$
CameraCtrl [9]	13.80	0.391	0.501	98.32	0.044	0.325	11.42	0.282	0.769	171.88	0.110	0.0766
ViewCrafter [47]	15.77	0.437	0.430	72.14	0.042	0.411	13.96	0.358	0.640	133.66	0.134	0.576
GEN3C [27]	15.28	0.524	0.431	80.03	0.081	0.454	13.80	0.321	0.648	129.91	0.098	0.778
SEVA [48]	15.12	0.528	0.472	85.39	0.048	0.422	14.00	0.349	0.620	101.62	0.100	0.662
Voyager [14]	15.80	0.521	0.409	79.81	0.077	0.638	14.06	0.337	0.652	89.37	0.120	0.650
GaC (Ours)	16.34	0.547	0.399	64.31	0.050	0.429	14.29	0.359	0.629	101.24	0.093	0.638



[Arxiv](#)