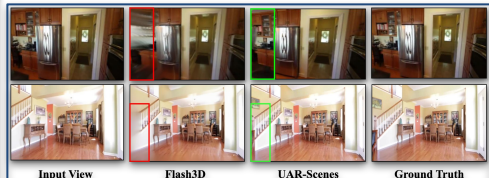


Problem Overview

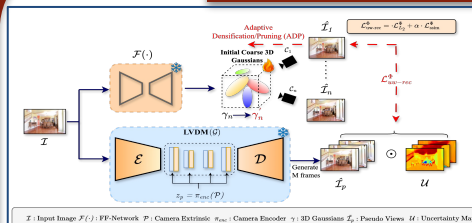


How to refine scenes so that it retains high fidelity in regions it has seen and generates consistent plausible completions for the unseen regions?

Contributions

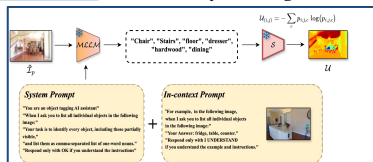
- ✓ Camera Controllable Video Diffusion-Guided Refinement of 3D Scenes for occluded and out-of-view regions.
- ✓ Uncertainty Estimation by distilling entropy with segmentation and MLLM guidance.
- ✓ Texture alignment on-the-fly with for matching the color of regressed and generated images.

Method Overview



$$\mathcal{L}_{uw-rec} = \left\| (1 - U_{(i,j)}) \odot (\tilde{I} - \tilde{I}_p) \right\|_2 + \alpha \text{SSIM}(\tilde{I}, \tilde{I}_p),$$

Simple Distillation of semantics with a MLLM driven prior is enough to estimate uncertainty!

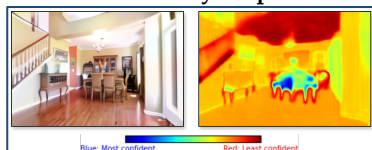


Building Blocks:

1. Obtain coarse scene gaussians obtained from feed-forward model.
2. Generate pseudo-views using camera controllable video diffusion model.
3. Improve regions using uncertainty mask guidance.

Results

Uncertainty Map U



RealEstate-10K (In-Domain)

Model	5 frames			10 frames			[0, 30, 30]		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
SynSis [52]	-	-	-	24.40	0.812	-	22.30	0.740	-
SV-MPI [53]	27.10	0.870	-	24.52	0.818	-	23.52	0.785	-
BTS [16]	24.15	0.897	0.113	25.00	0.855	0.194	24.00	0.755	0.184
Splitter Image [19]	24.15	0.894	0.110	25.60	0.760	0.240	23.10	0.730	0.290
MINE [54]	28.45	0.897	0.100	25.90	0.850	0.190	24.75	0.820	0.170
Flash3D [1]	28.46	0.899	0.100	25.94	0.857	0.133	24.93	0.833	0.160
UAR-Scenes	28.67	0.902	0.095	26.54	0.861	0.112	27.81	0.887	0.107

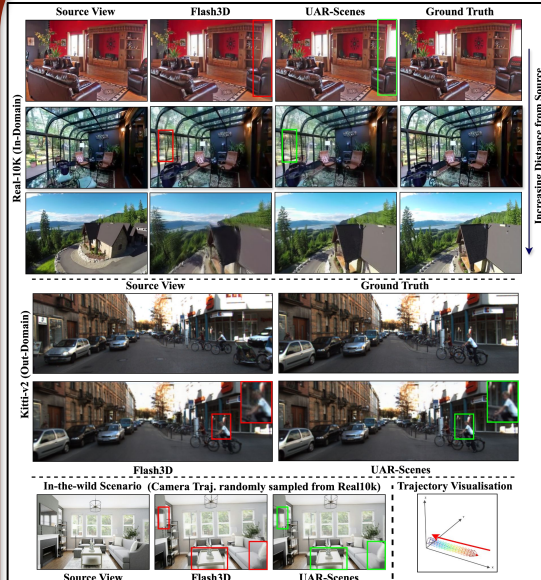
$$o_t = \mathcal{S}(\tilde{I}_p, O) \in \mathbb{R}^{n \times (N+1)},$$

$$U_{(i,j)} = - \sum_c p_{i,j,c} \log(p_{i,j,c}),$$

KITTI-v2 (Out-Domain)

Method	KITTI		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
LDI [56]	16.50	0.572	-
SV-MPI [53]	19.50	0.733	-
BTS [16]	20.10	0.761	0.144
MINE [54]	21.90	0.828	0.112
Flash3D [1]	21.96	0.826	0.132
UAR-Scenes	22.31	0.844	0.128

Novel View Synthesis



Qualitative results on benchmark datasets and in-the-wild scenarios

Acknowledgements: This work was supported by USDA NRI grant 021-67022-33453, UC MRPI grant A21-0101-S003, and NSF grant CMMI-2326309

Arxiv

