# DeepDR: Deep Structure-Aware RGB-D Inpainting for Diminished Reality

Christina Gsaxner<sup>1</sup>, Shohei Mori<sup>1</sup>, Dieter Schmalstieg<sup>1,2</sup>, Jan Egger<sup>1,3</sup>, Gerhard Paar<sup>4</sup>, Werner Bailer<sup>4</sup> and Denis Kalkofen<sup>1,5</sup>.

<sup>1</sup>Graz University of Technology, <sup>2</sup>University of Stuttgart, <sup>3</sup>University of Duisburg-Essen, <sup>4</sup>Joanneum Research, <sup>5</sup>Flinders University

# Introduction

Diminished reality (DR) involves virtually removing real objects from the environment using inpainting techniques. However, existing methods struggle with maintaining coherent structure and 3D geometry, particularly for advanced tasks like 3D scene editing. In response, we introduce DeepDR, a real-time RGB-D inpainting framework tailored for DR, ensuring both realistic image and depth inpainting with minimal artifacts, achieved through a structure-aware generative network explicitly conditioned on scene semantics.

# Network architecture



Input images / and depth maps D are encoded separately, then fused on a higher dimension for joint completion. Our semantics-aware decoder comprises up blocks conditioning features on semantic information. Outputs  $I_{o}^{t}$  and  $D_{o}^{t}$  serve as auxiliary inputs in a recurrent feedback loop for subsequent time steps.



DeepDR is a framework for **Diminished Reality that removes** objects from the environment via real-time, structure-aware image and depth inpainting.

#### Results

We show that DeepDR can outperform related methods on three datasets qualitatively and quantitatively. DeepFillV2 **PanoDR** E2FGVI DeepDR (Ours)



![](_page_0_Figure_14.jpeg)

	Method	$LPIPS\downarrow$	$FID\downarrow$	Depth RMSE $\downarrow$
InteriorNet [5]	DeepFillV2 [1]	0.0150	0.448	0.572
	PanoDR [2]	<u>0.0128</u>	0.606	0.564
	E2FGVI [3]	0.0131	<u>0.363</u>	<u>0.563</u>
	DeepDR (Ours)	0.0104	0.218	0.278
DynaFill [4]	DeepFillV2 [1]	0.0238	4.122	7.92
	PanoDR [2]	0.0250	5.579	8.12
	E2FGVI [3]	<u>0.0169</u>	2.826	7.83
	DynaFill [4]	0.0197	<u>2.665</u>	<u>7.78</u>
	DeepDR (Ours)	0.0168	2.415	4.51
ScanNet [6]	DeepFillV2 [1]	0.0208	0.693	0.508
	PanoDR [2]	0.0119	0.348	0.536
	E2FGVI [3]	<u>0.0110</u>	<u>0.295</u>	0.512
	DeepDR (Ours)	0.0108	0.292	0.484

### Structure-aware RGB-D decoder

![](_page_0_Figure_17.jpeg)

Our up blocks predict a semantic segmentation map  $S_i$  from the up-sampled image feature  $i_i$  via a pyramid pooling module. This segmentation guides the image and depth feature generation consistently.

# Training

![](_page_0_Picture_20.jpeg)

We employ loss terms and adversarial learning to enhance the accuracy and realism of inpainted RGB and depth data ( $\mathcal{L}I$ ,  $\mathcal{L}D$ ). Temporal consistency is ensured through the computation of a temporal loss  $\mathcal{L}t$  using optical flow. We enforce structural learning by incorporating semantic supervision via  $\mathcal{L}seg$ .

In CVPR, 2022. In CVPR, 2017.

![](_page_0_Picture_23.jpeg)

Acknowledgement: The work was funded by the FFG project "TRIP" (BRIDGE 883658) and the FWF project "enFaced 2.0" (KLI 1044). We thank xCAD Solutions GmbH for their continuous support.

[1] Yu et al., "Free-form image inpainting with gated convolution". In ICCV, 2019. [2] Giktas et al., "PanoDR: Spherical panorama diminished reality for indoor scenes". In CVPR Workshops, 2021.

[3] Li et al., "Towards an end-to-end framework for flow-guided video inpainting".

[4] Bešić et al., "Dynamic object removal and spatio-temporal RGB-D inpainting via geometry-aware adversarial learning". IEEE trans. intell. veh., 2022.

[5] Li et al., "InteriorNet: Mega-scale Multi-sensor Photo-realistic Indoor Scenes Dataset". In BMVC, 2018.

[6] Dai et al., "Scannet: Richly-annotated 3d reconstructions of indoor scenes".