



**Robotics and
Visual Intelligence Lab**

PixelSynth: Generating a 3D-Consistent Experience from a Single Image

ICCV 2021

Chris Rockwell

David F. Fouhey
University of Michigan

Justin Johnson

Gyeongsu Cho @ UNIST

Lab Seminar

2022.03.03 (Thr)

Introduction

Network

Experiments

Conclusion

Introduction

View Synthesis?

Introduction

View Synthesis?



Introduction

View Synthesis?



Introduction

Single Image View Synthesis

Introduction

Single Image View Synthesis



Introduction

Single Image View Synthesis



Translation: Step Forward
Rotation: Turn left

Introduction

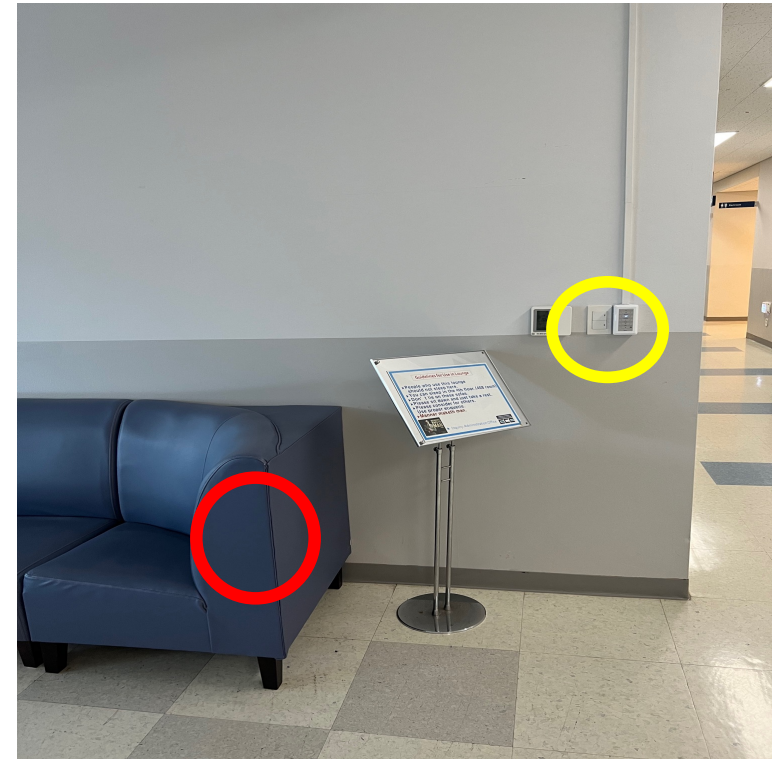
Single Image View Synthesis



Translation: Step Forward
Rotation: Turn left

Introduction

Challenges of Single Image View Synthesis



Challenge 1: Need to know depth

Introduction

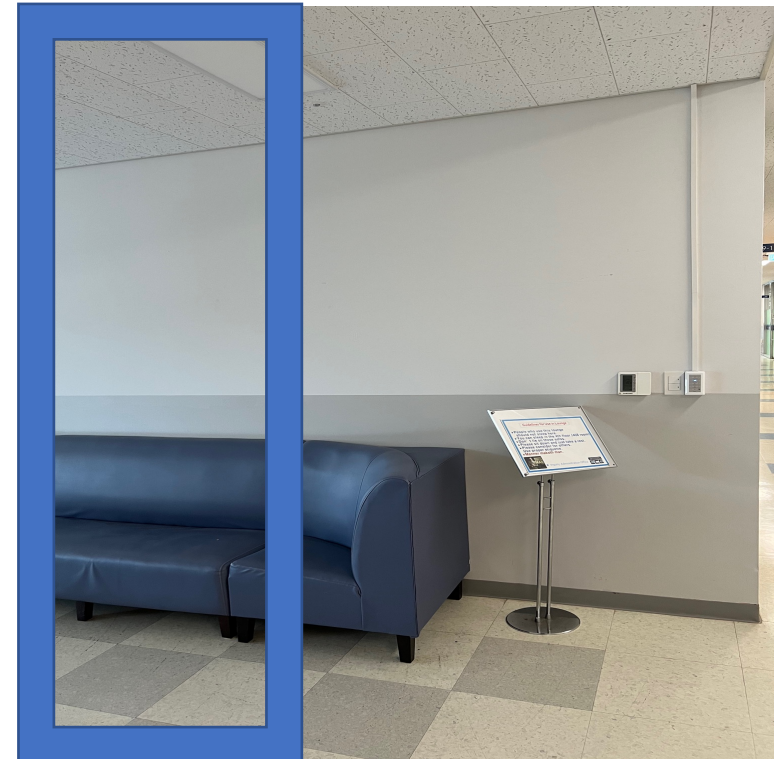
Challenges of Single Image View Synthesis



Challenge 2: Inpainting missing regions

Introduction

Challenges of Single Image View Synthesis



Challenge 3: Outpainting missing regions consistently

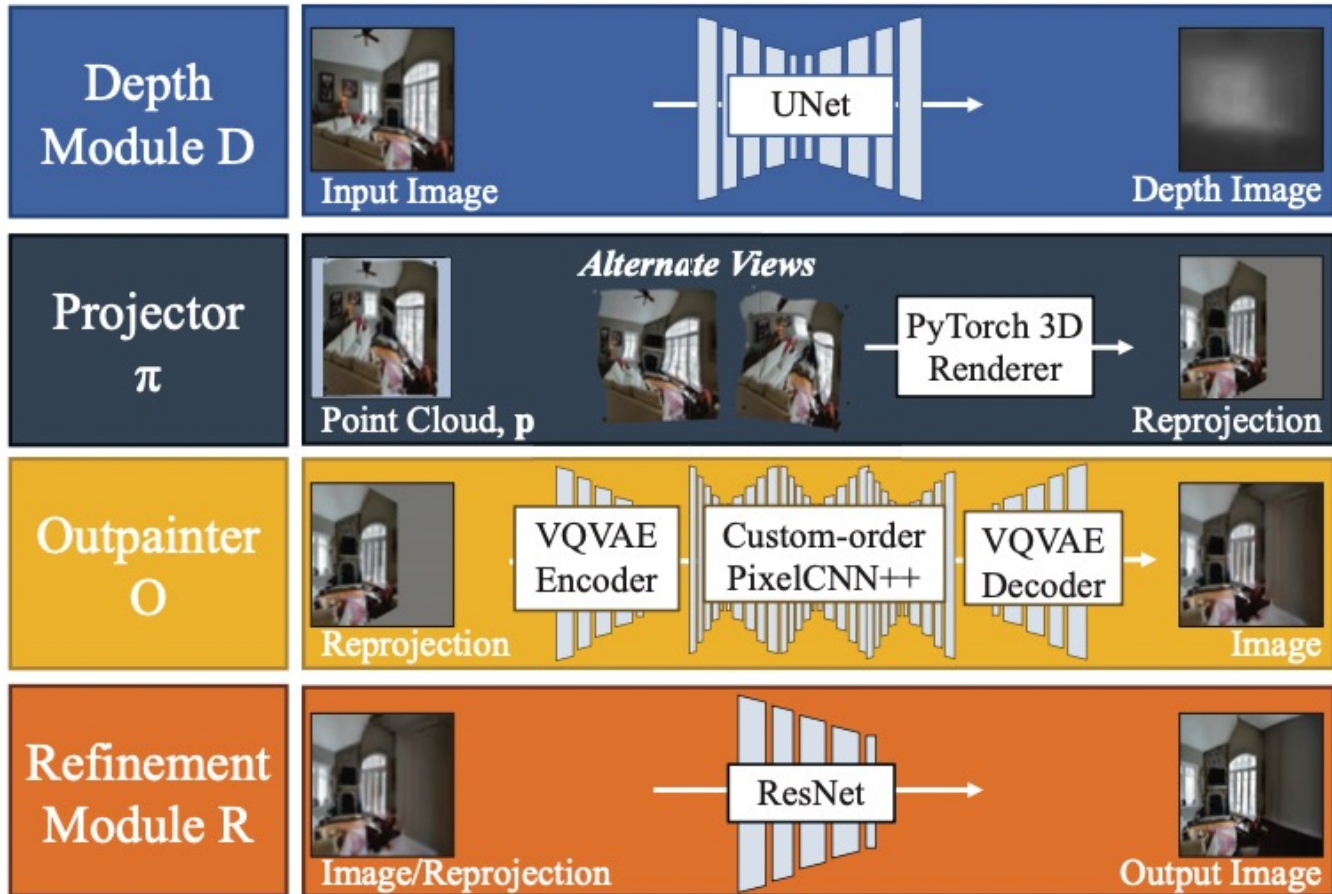
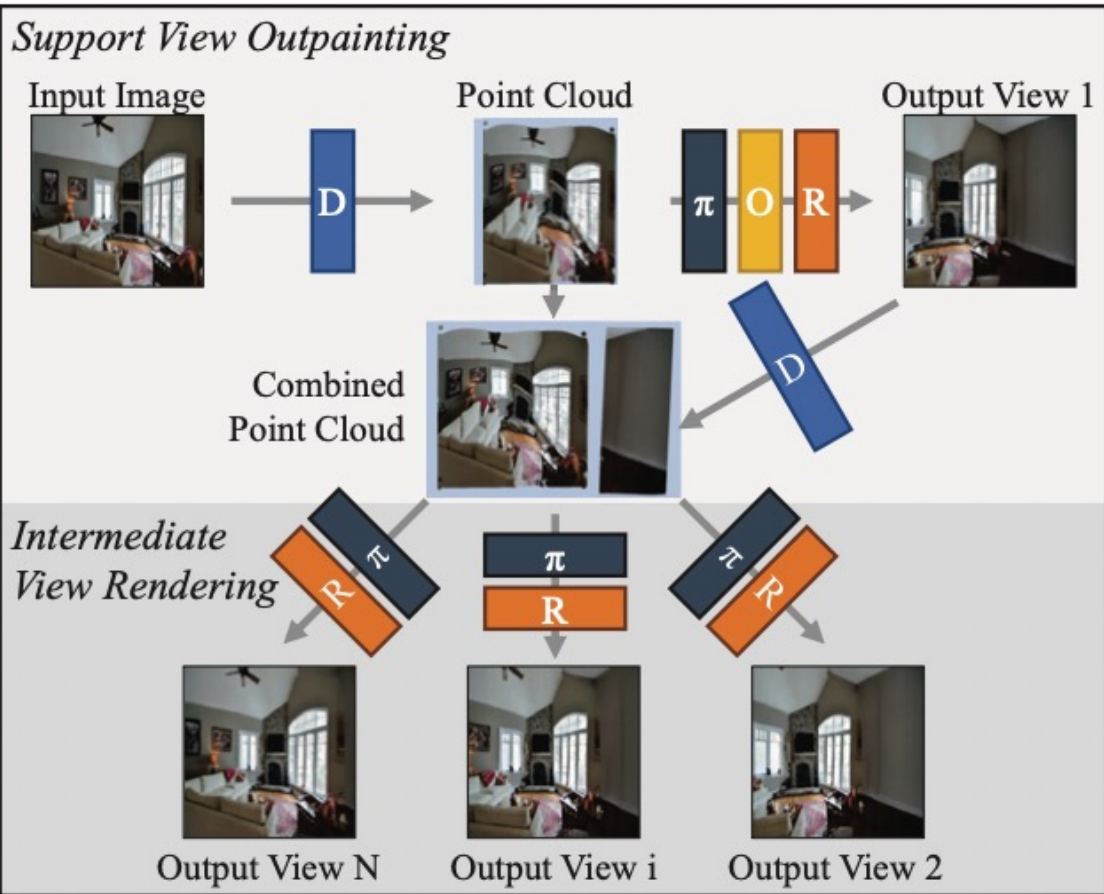
Introduction

Network

Experiments

Conclusion

Network



Network

Depth
Module D



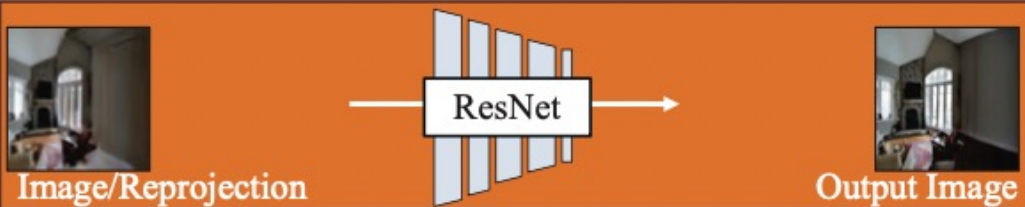
Projector
 π



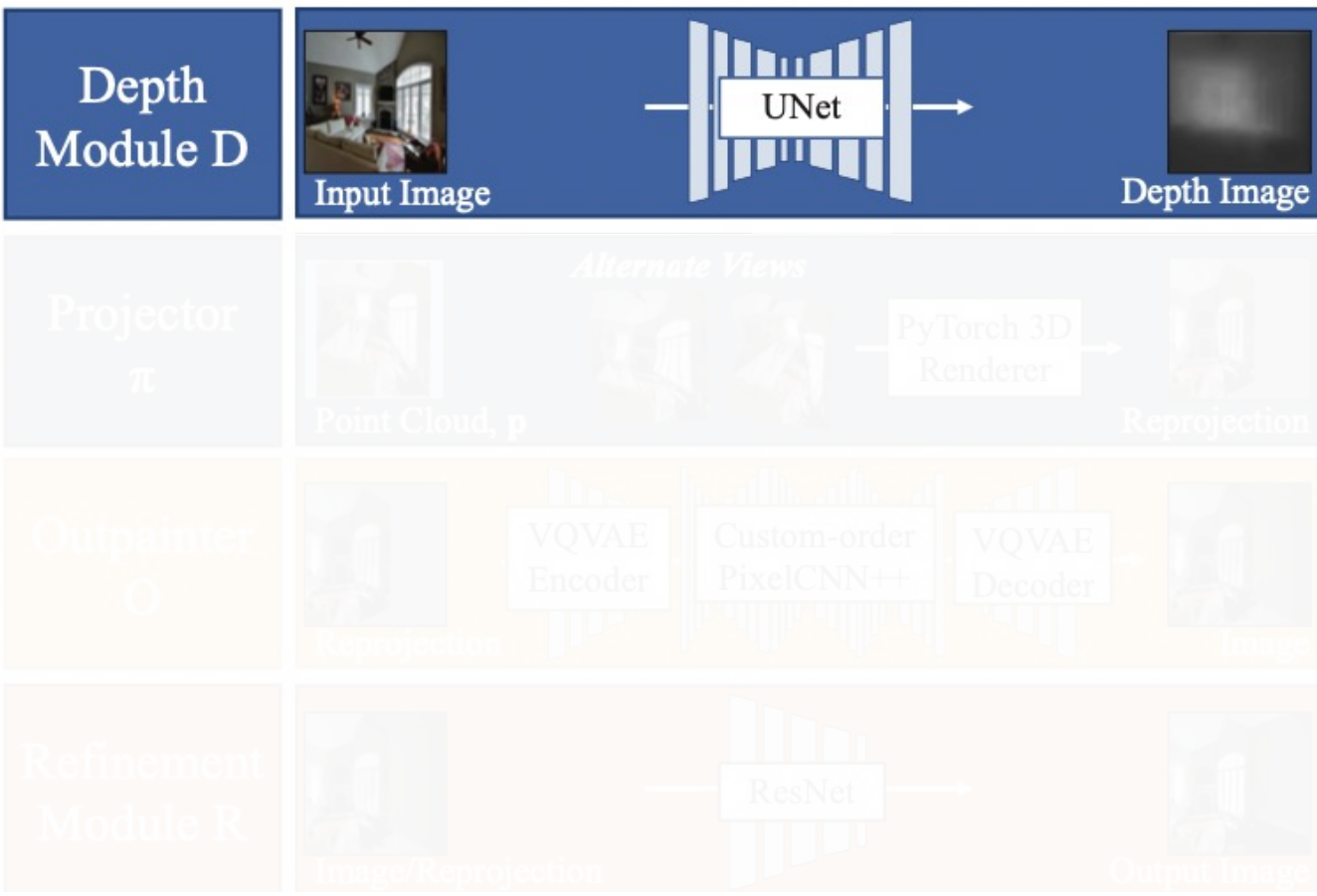
Outpainter
O



Refinement
Module R



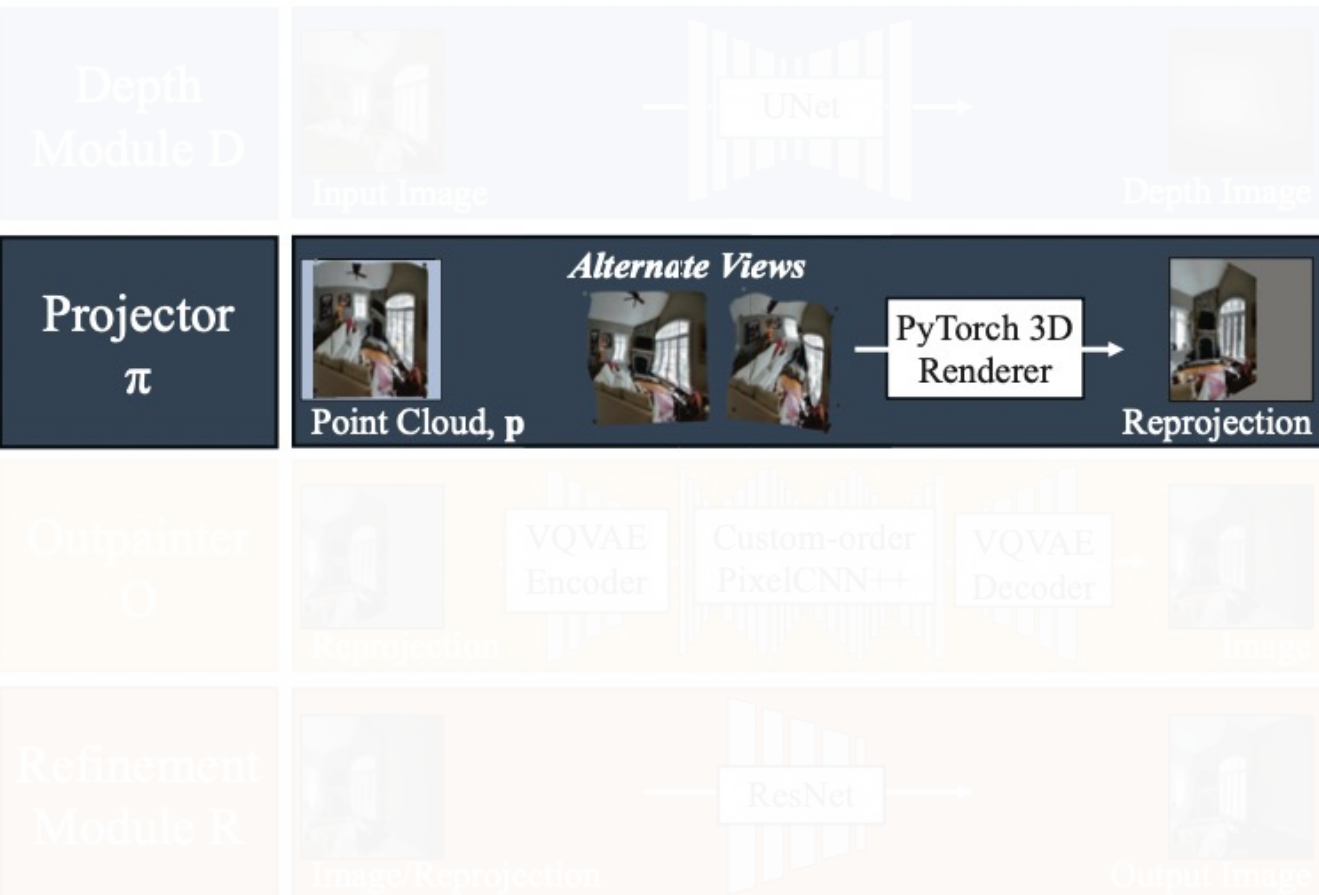
Network



$$C = D(I)$$

Point Cloud = Depth Module (Image)

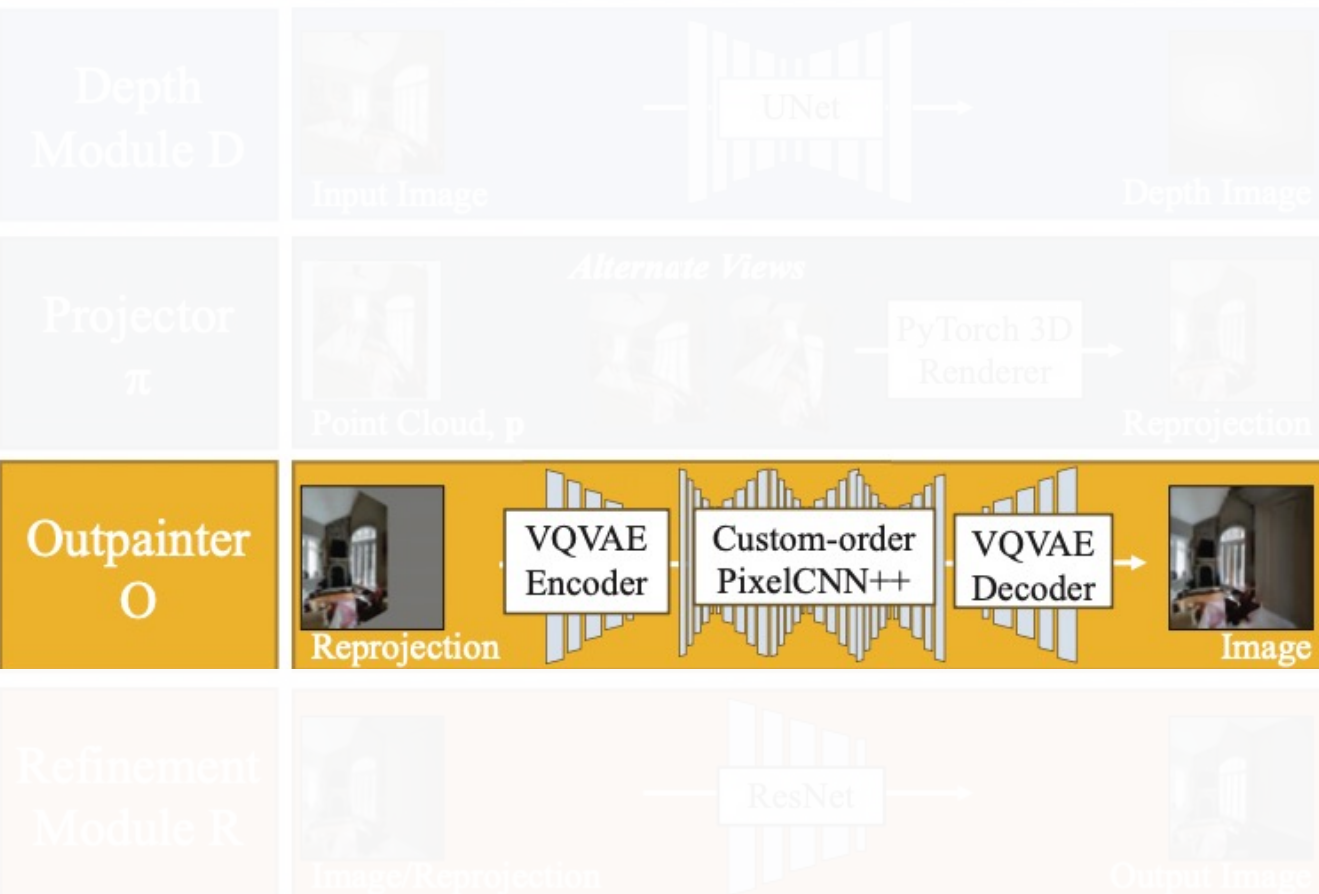
Network



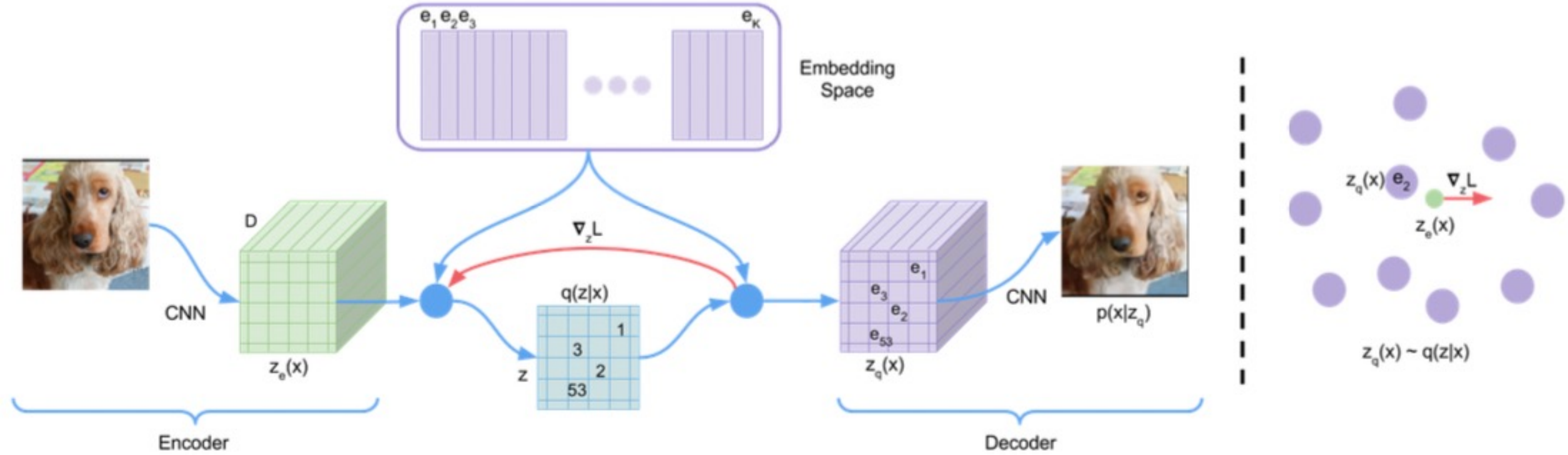
$$I = \pi(C, \mathbf{p})$$

Image Projector Point Cloud Target pose

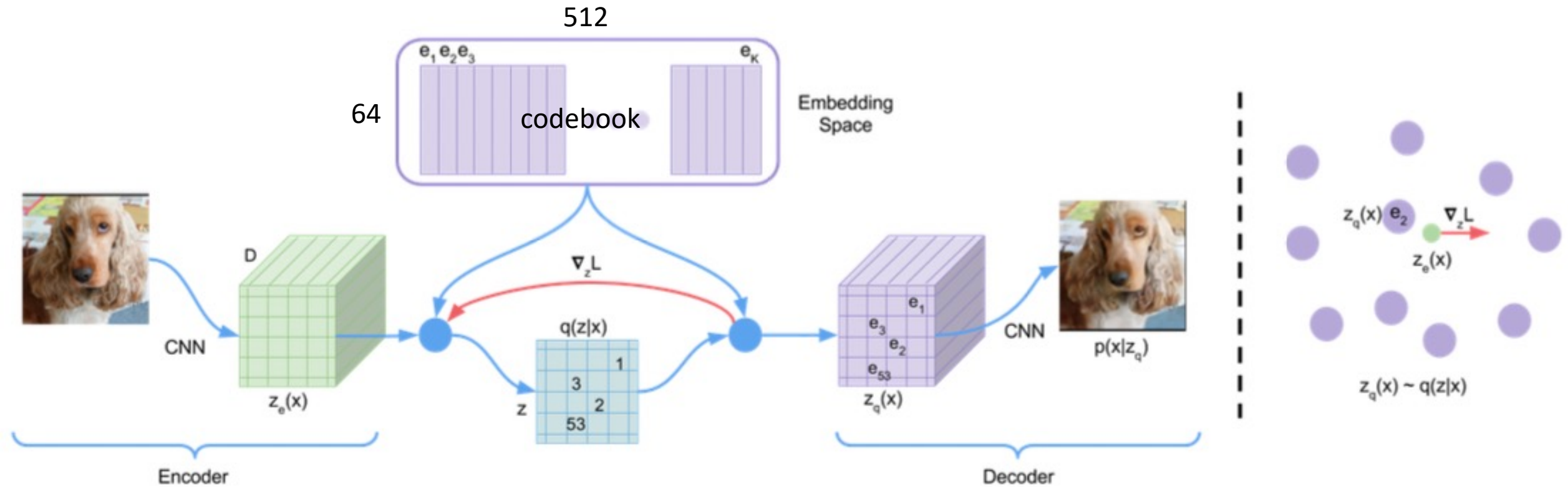
Network



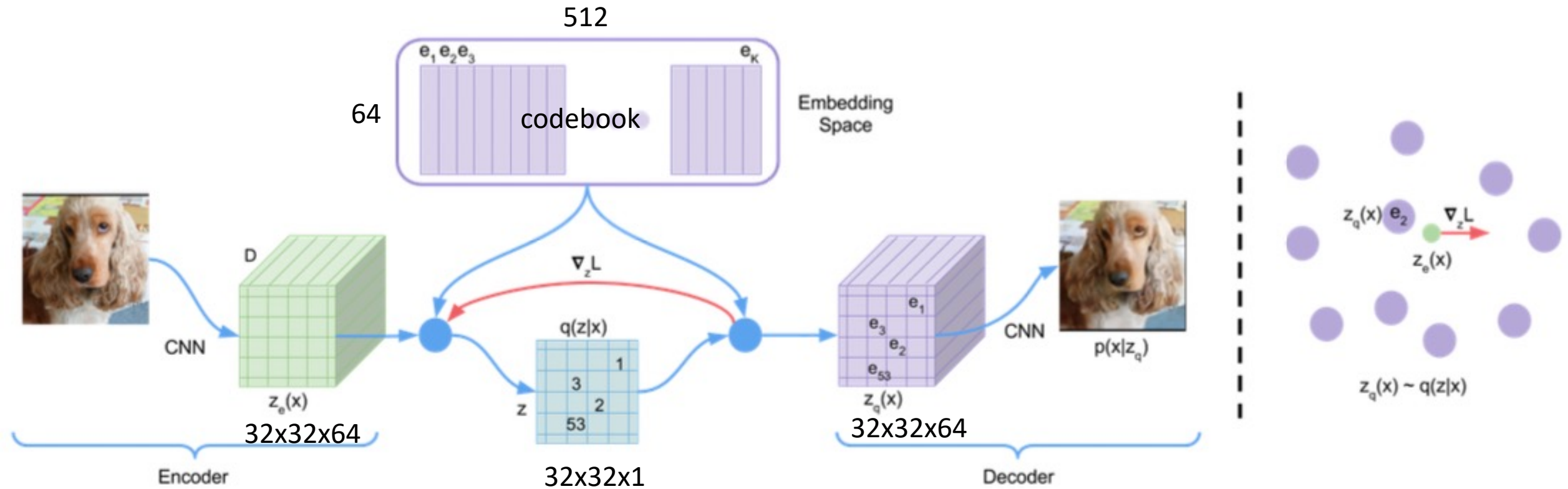
Network-VQVAE



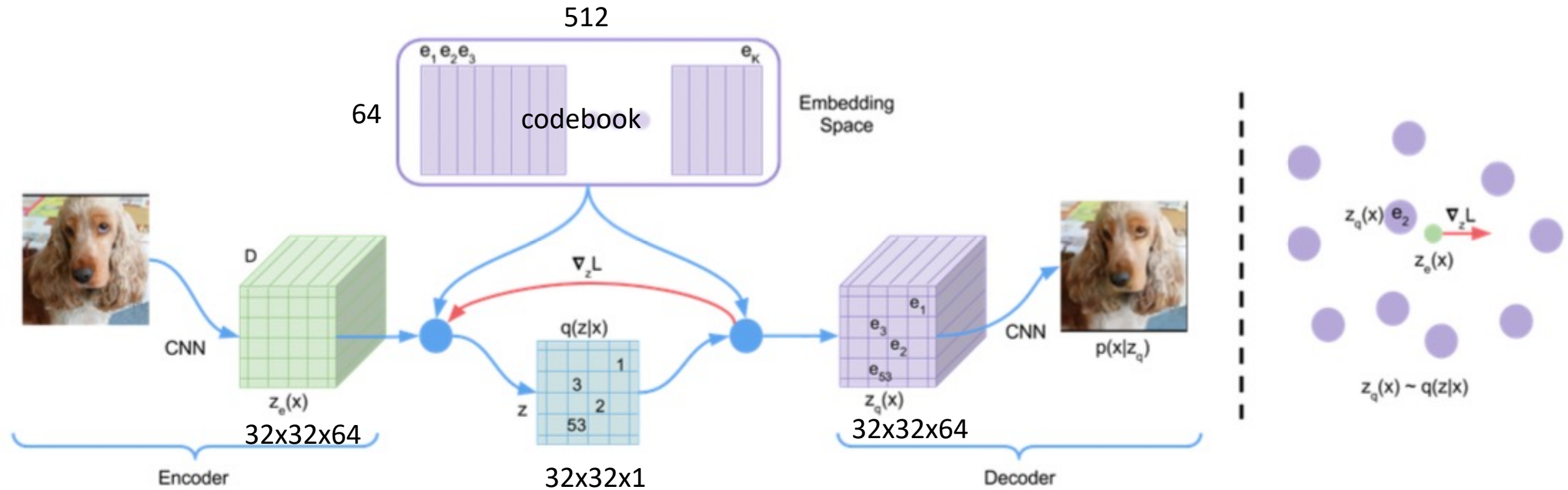
Network-VQVAE



Network-VQVAE

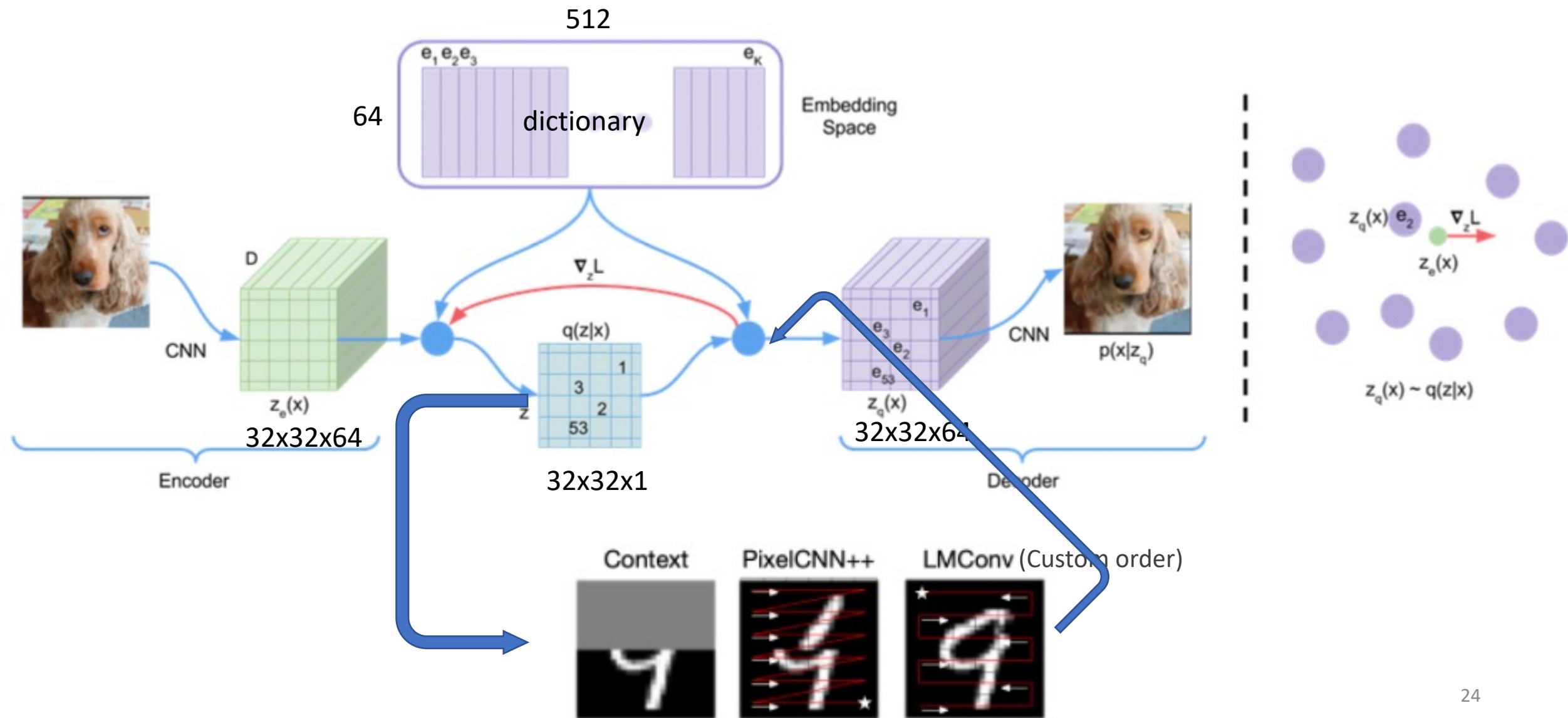


Network-VQVAE

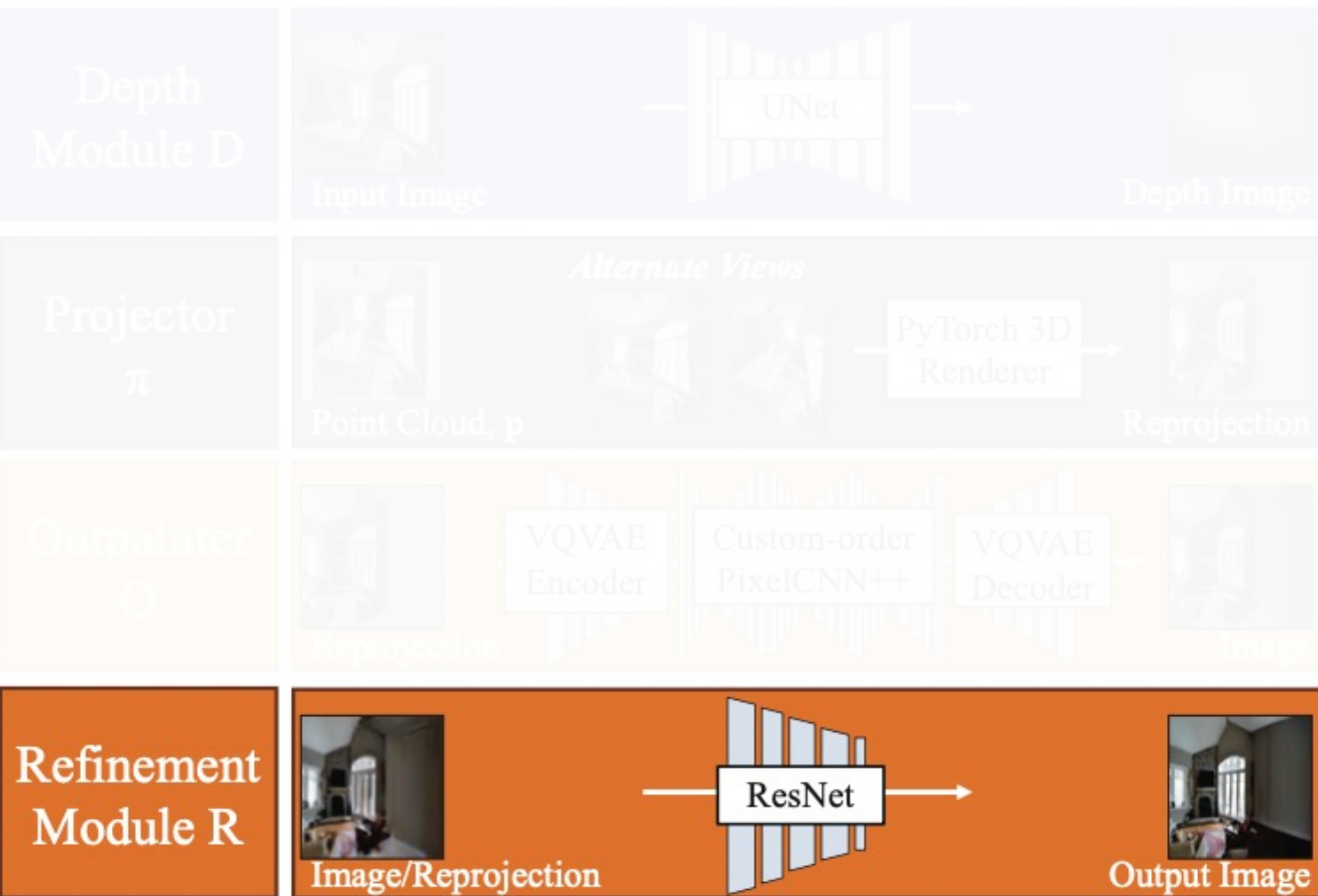


$$\begin{aligned}
 L &= \log p(x|z_q(x)) + \|sg[z_e(x)] - e\|_2^2 + \beta \|z_e(x) - sg[e]\|_2^2, \\
 &= \underbrace{\|\mathbf{x} - D(\mathbf{e})\|_2^2}_{\text{reconstruction}} + \underbrace{\|sg[E(\mathbf{x})] - \mathbf{e}\|_2^2}_{\text{codebook}} + \underbrace{\beta \|sg[\mathbf{e}] - E(\mathbf{x})\|_2^2}_{\text{commitment}}
 \end{aligned}$$

Network-Custom order pixelcnn++



Network

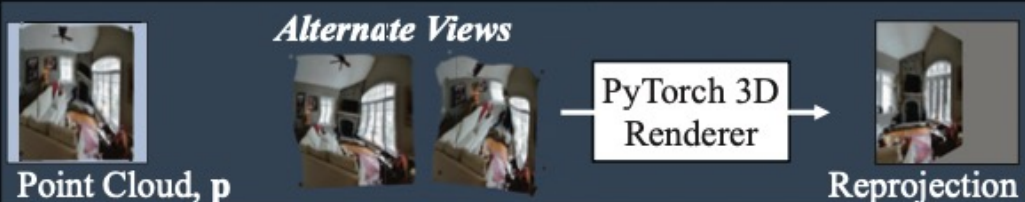


Network-Training

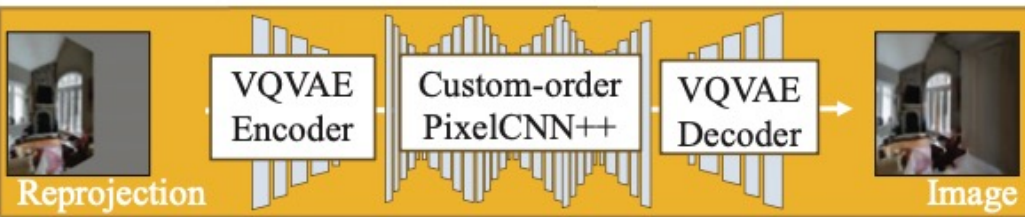
Depth
Module D



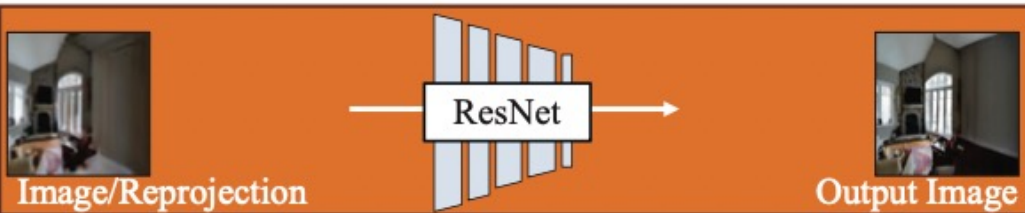
Projector
 π



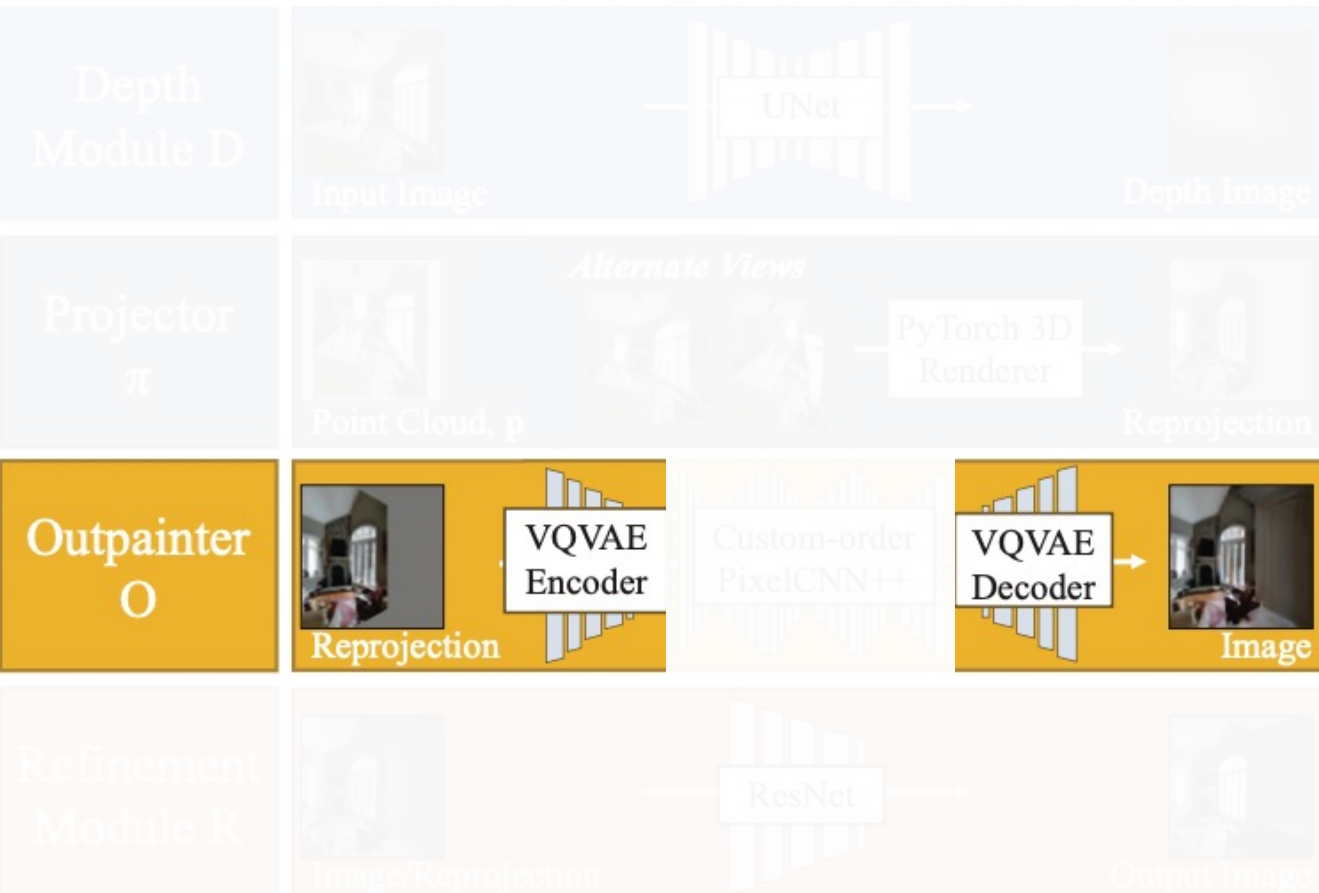
Outpainter
O



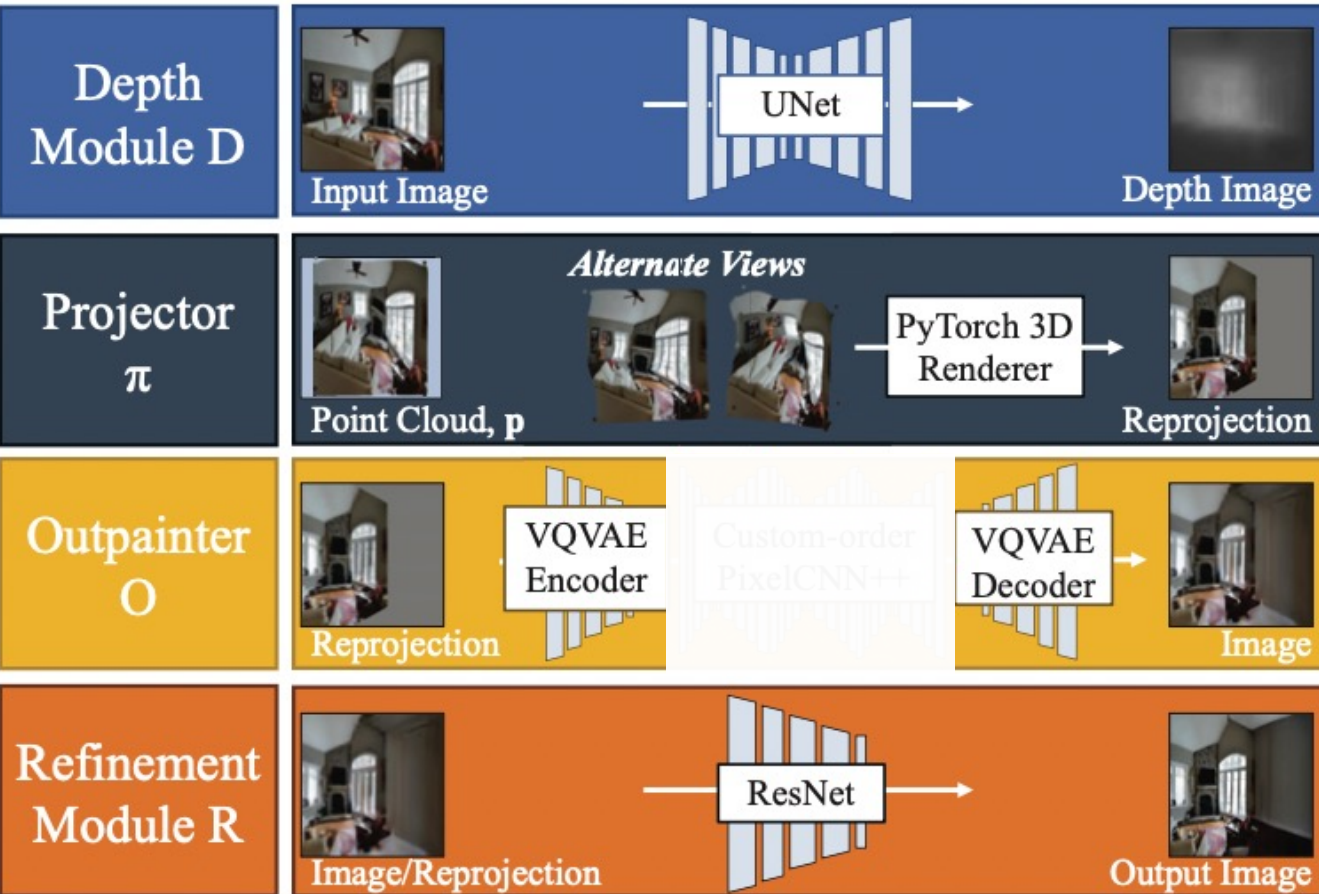
Refinement
Module R



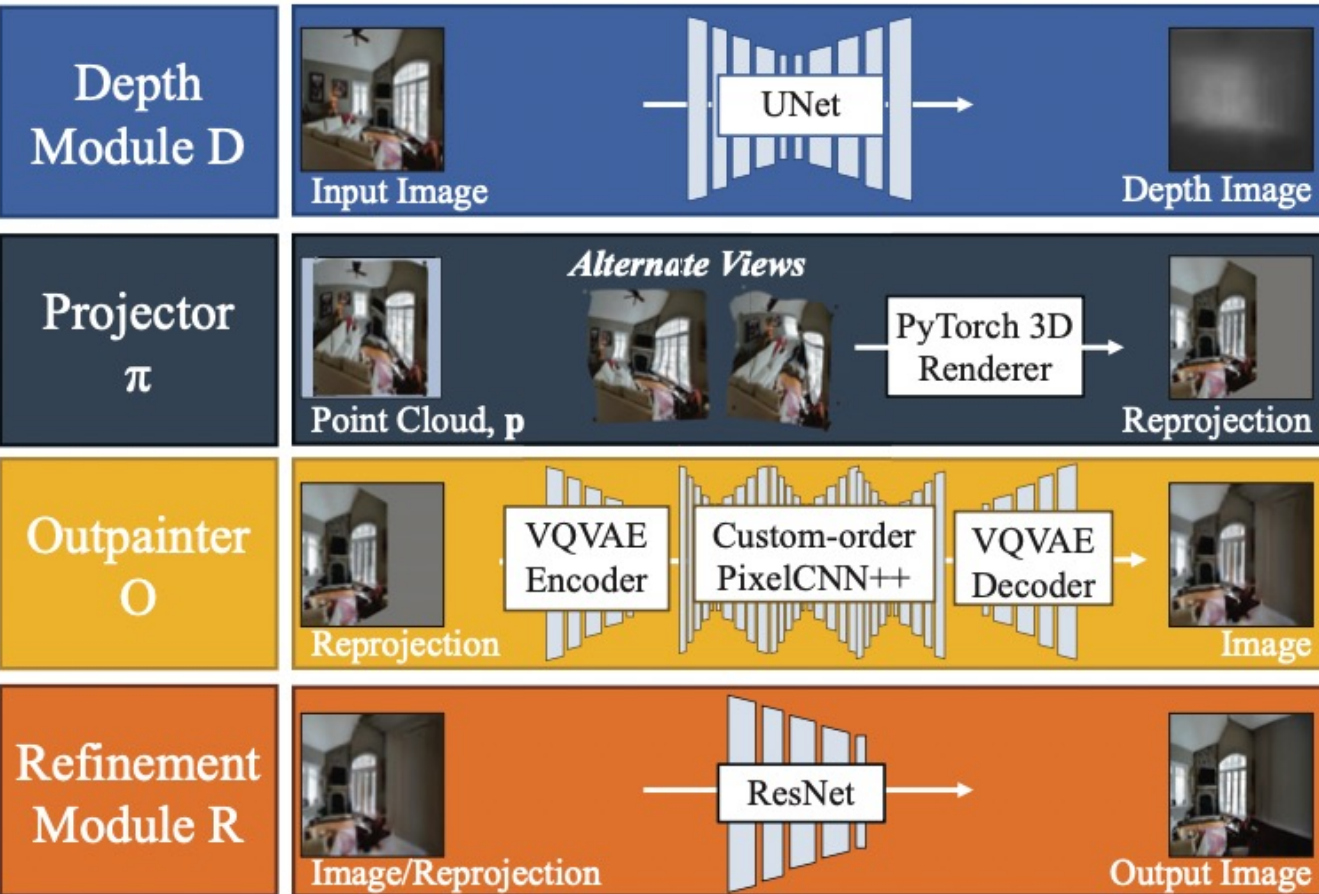
Network-Training



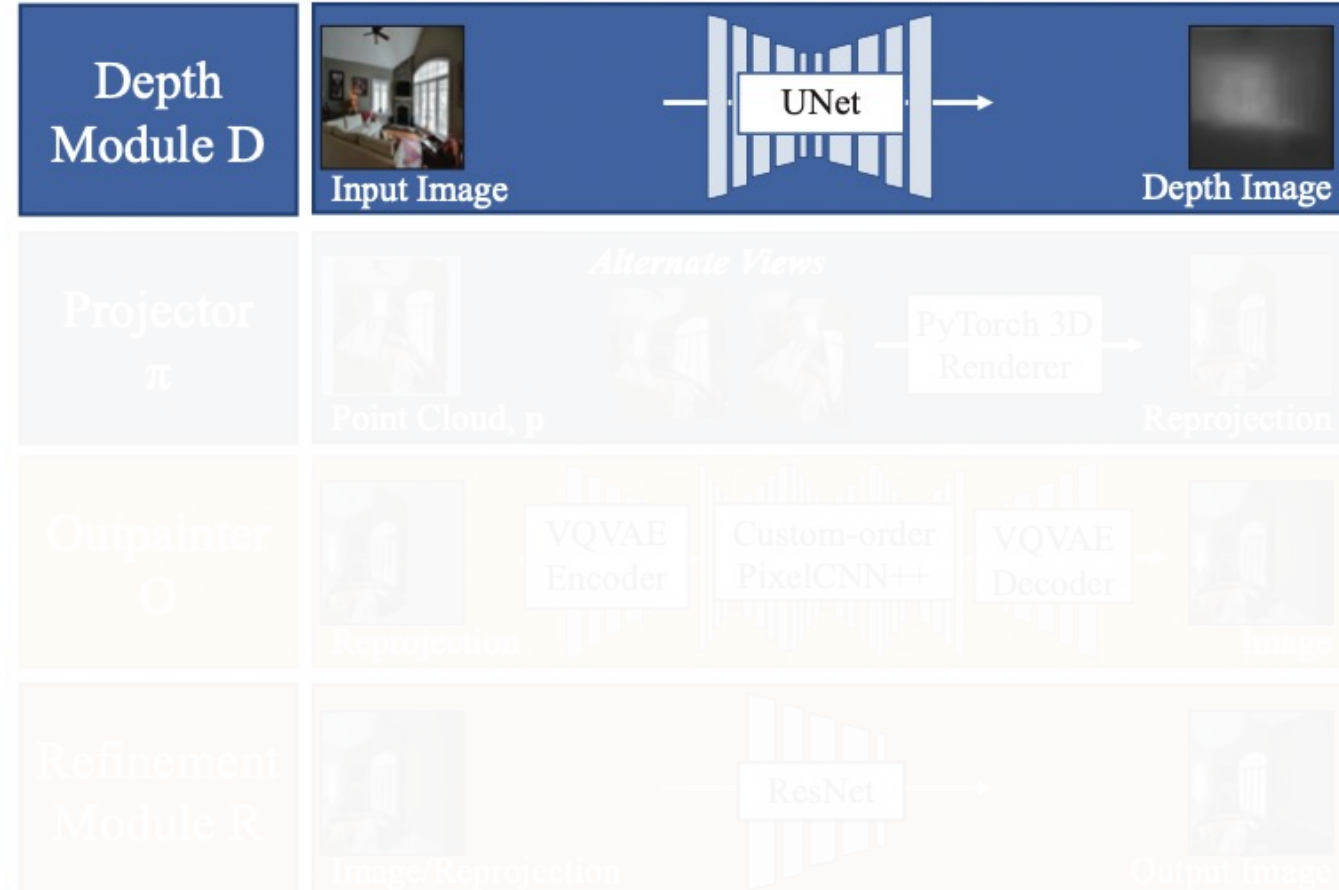
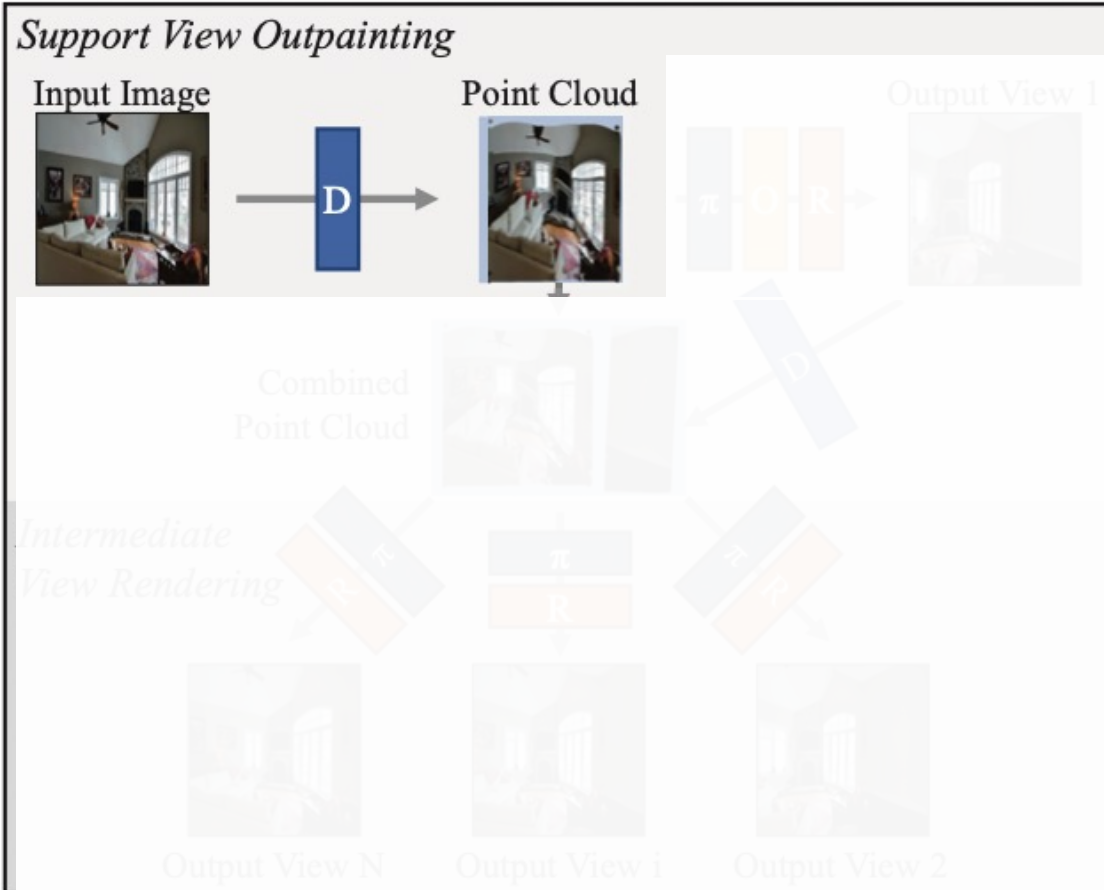
Network-Training



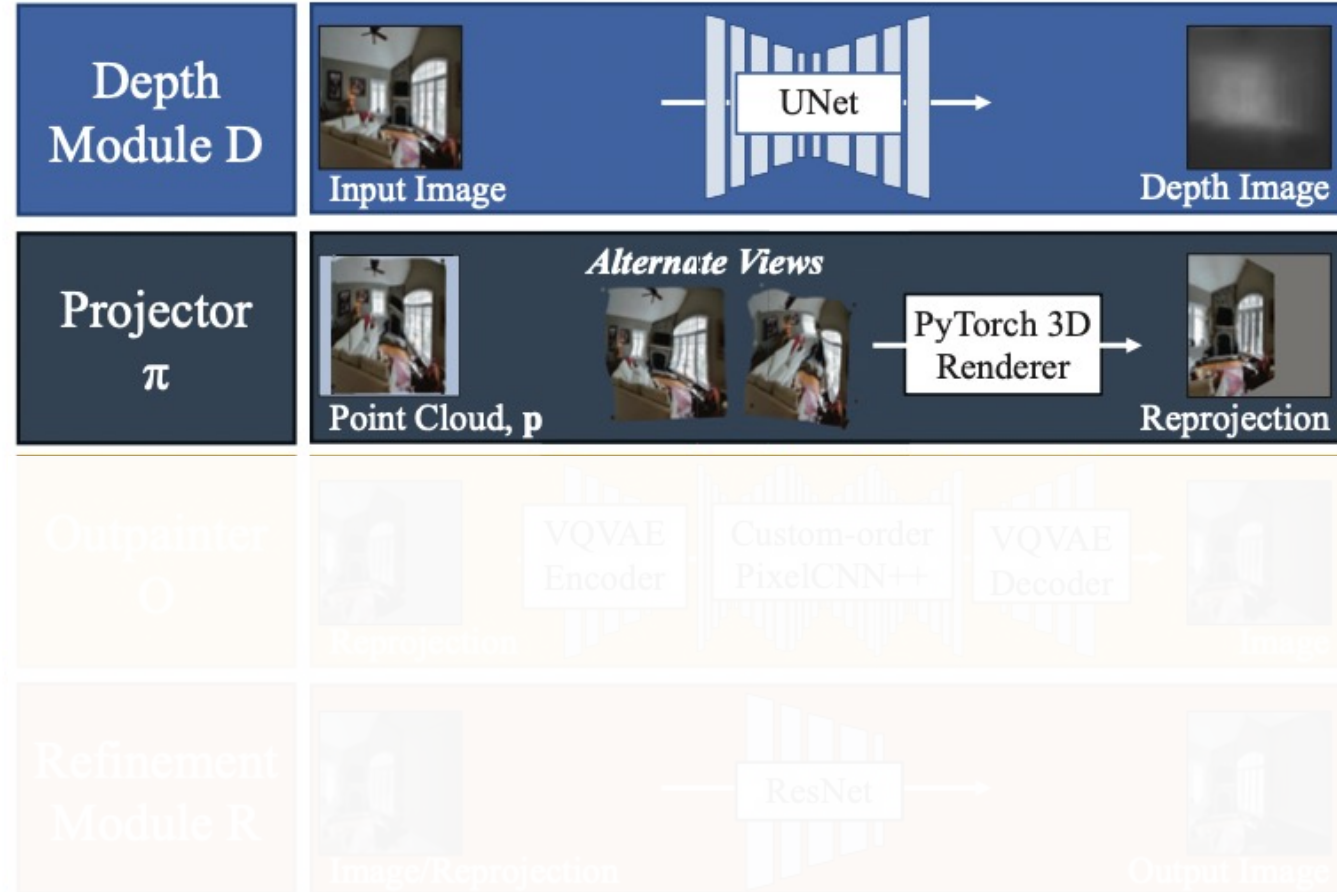
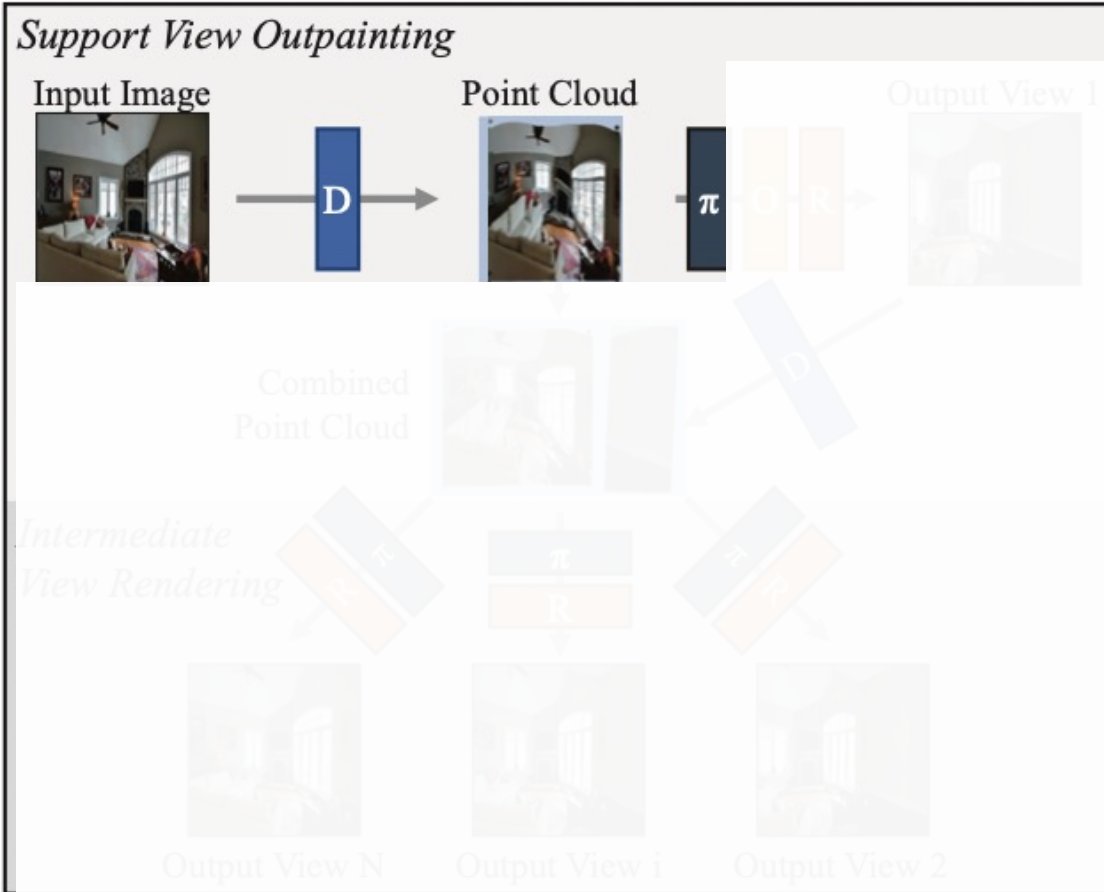
Network-Training



Network-Inference

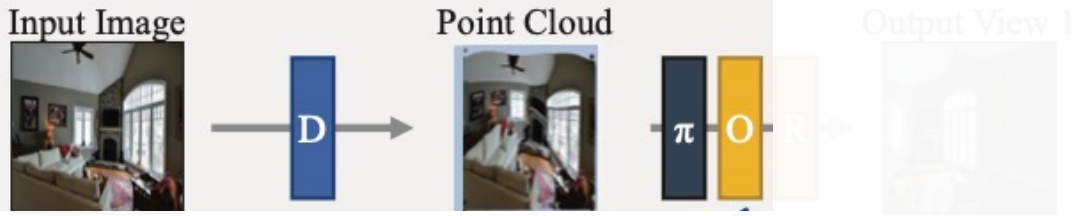


Network-Inference



Network-Inference

Support View Outpainting



Combined Point Cloud

Intermediate View Rendering



Depth Module D



Projector π



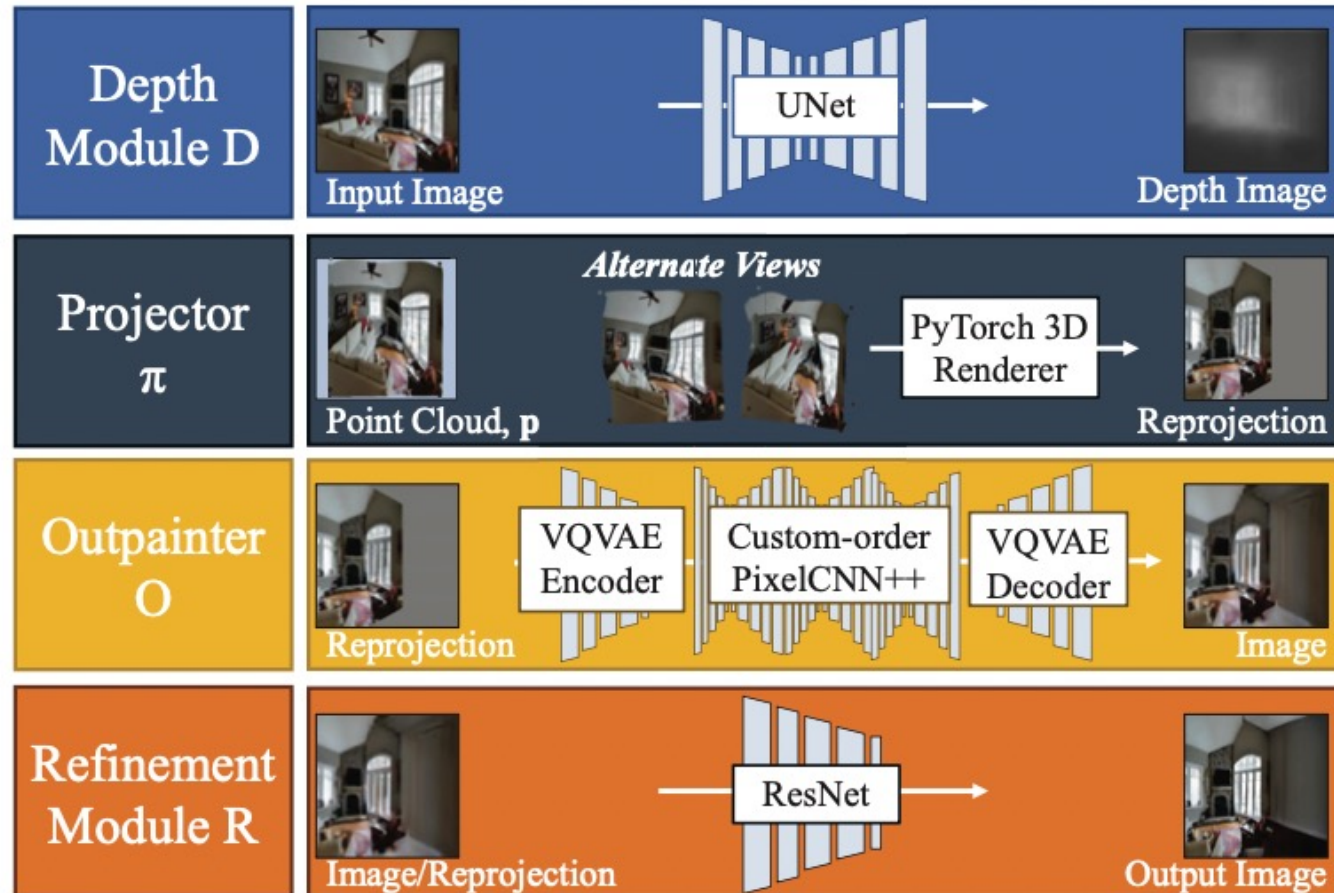
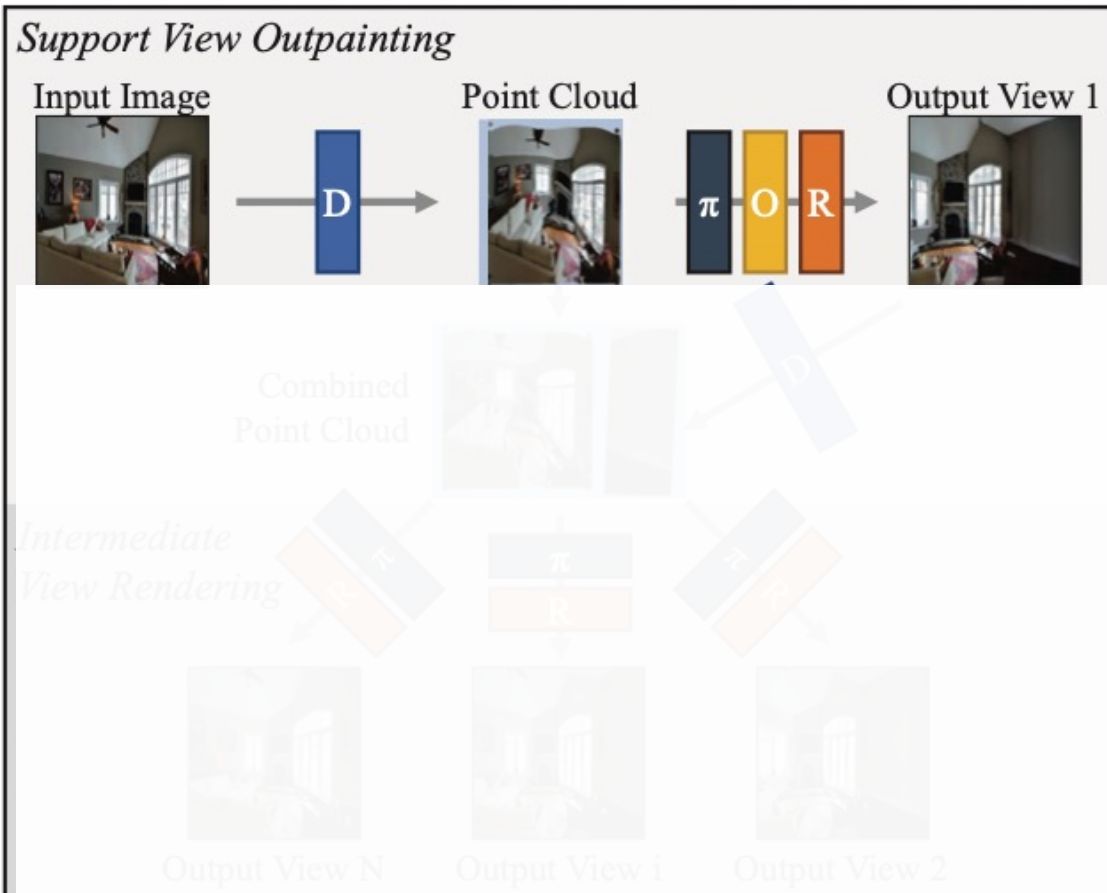
Outpainter O



Refinement Module R

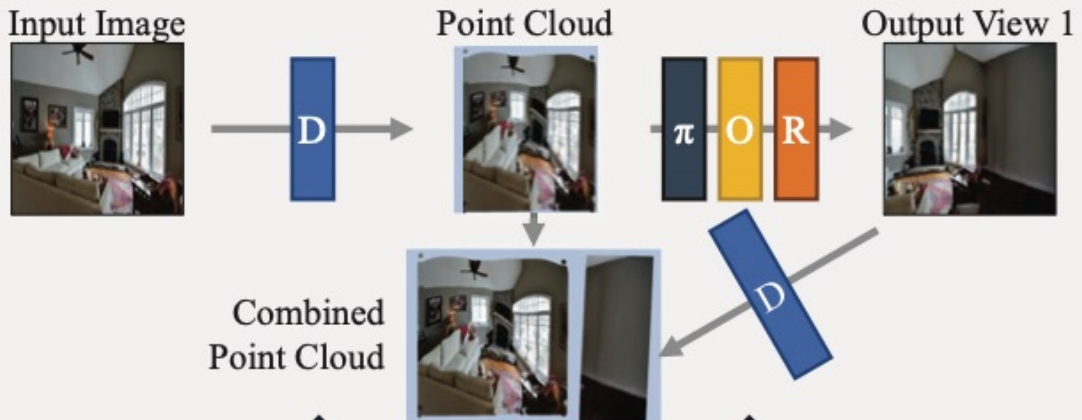


Network-Inference

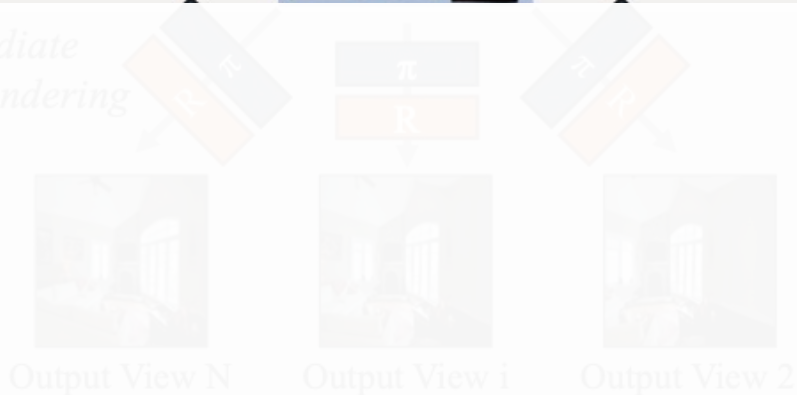


Network-Inference

Support View Outpainting



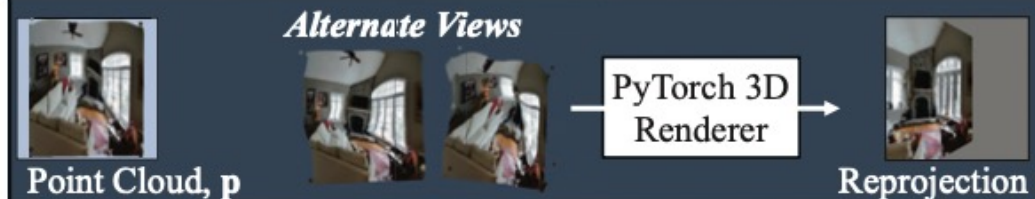
Intermediate View Rendering



Depth Module D



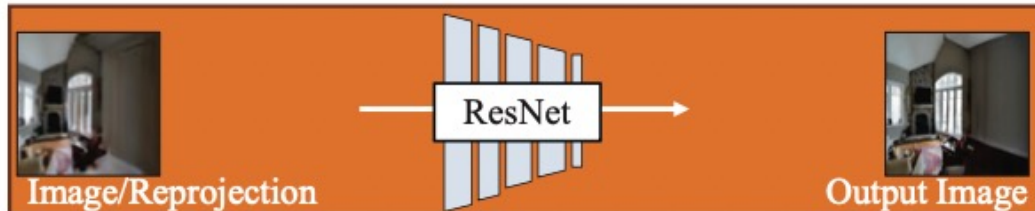
Projector π



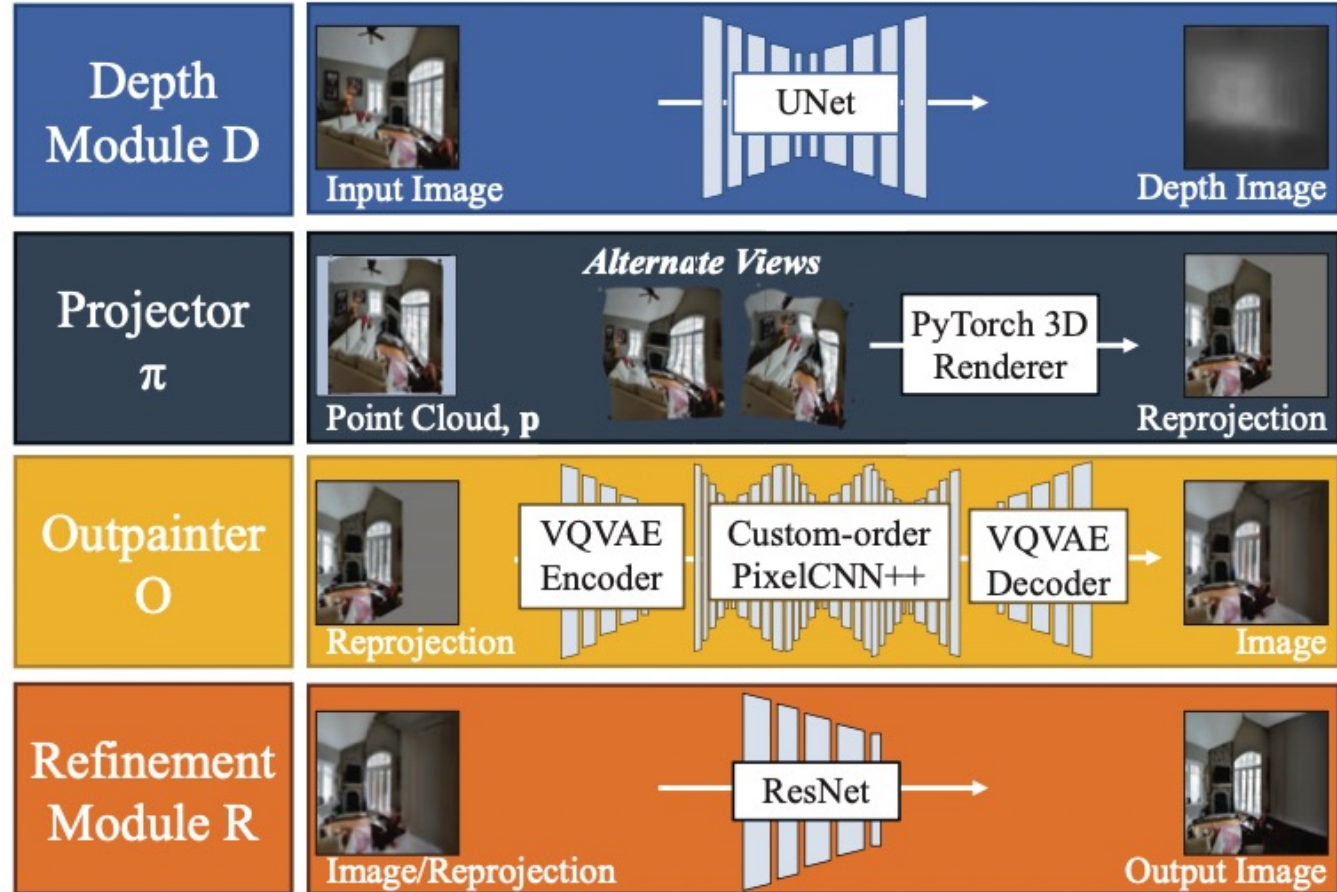
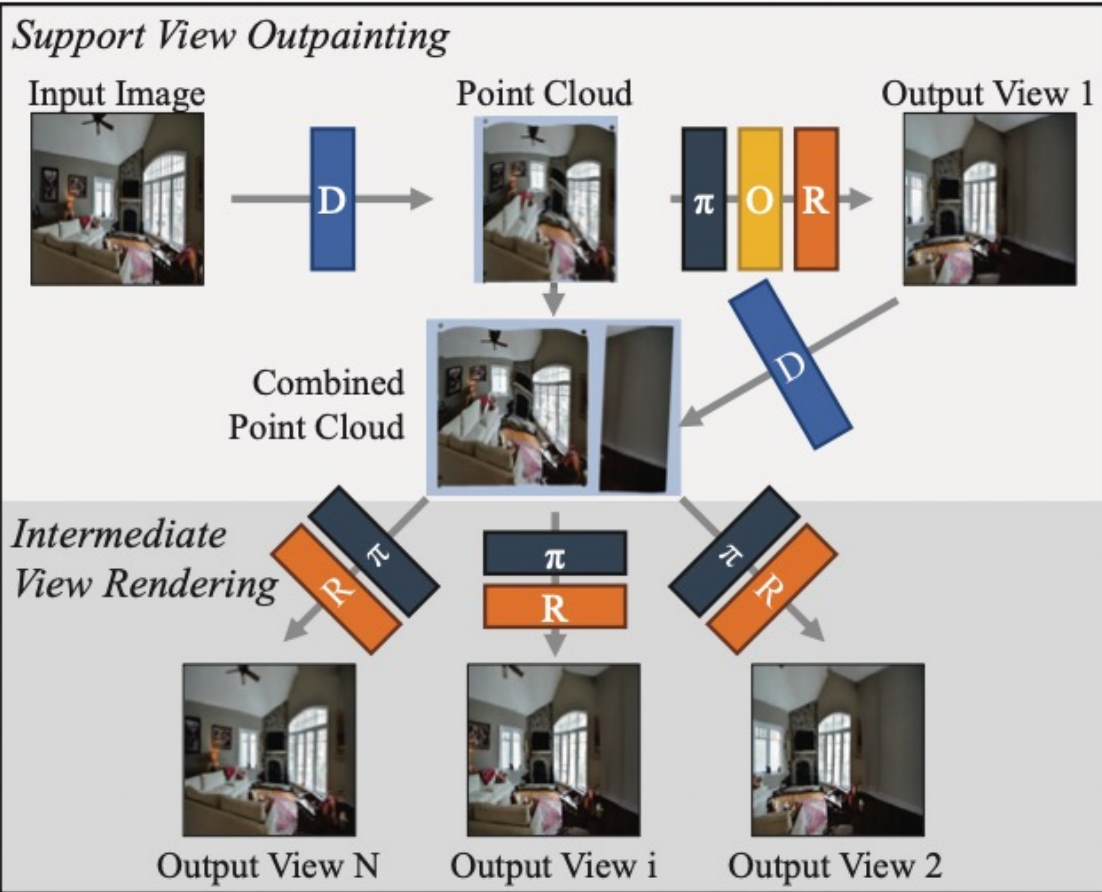
Outpainter O



Refinement Module R



Network-Inference



Introduction

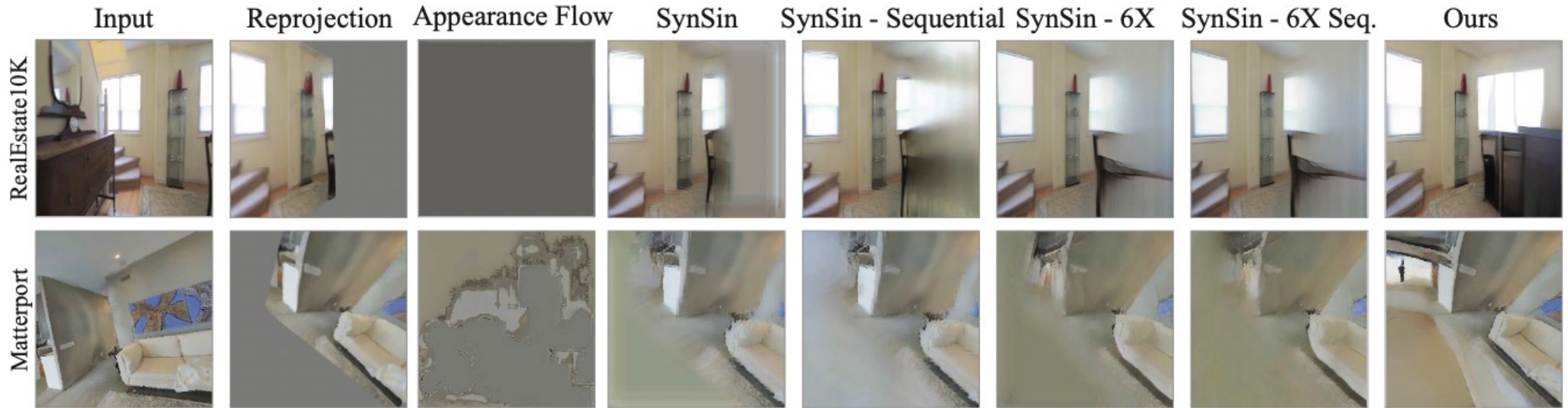
Network

Experiments

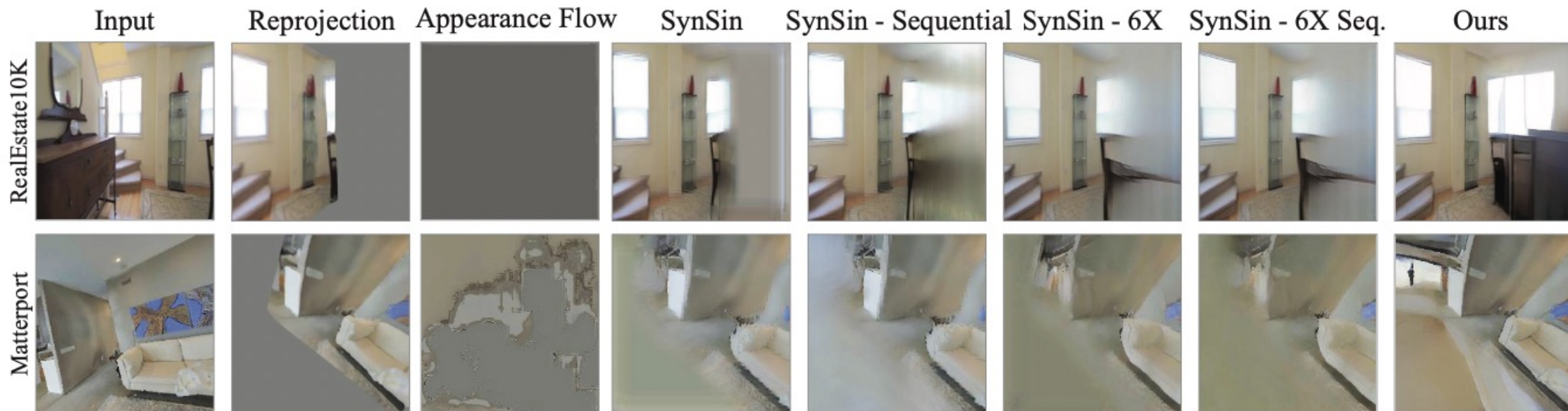
Conclusion



Experiments

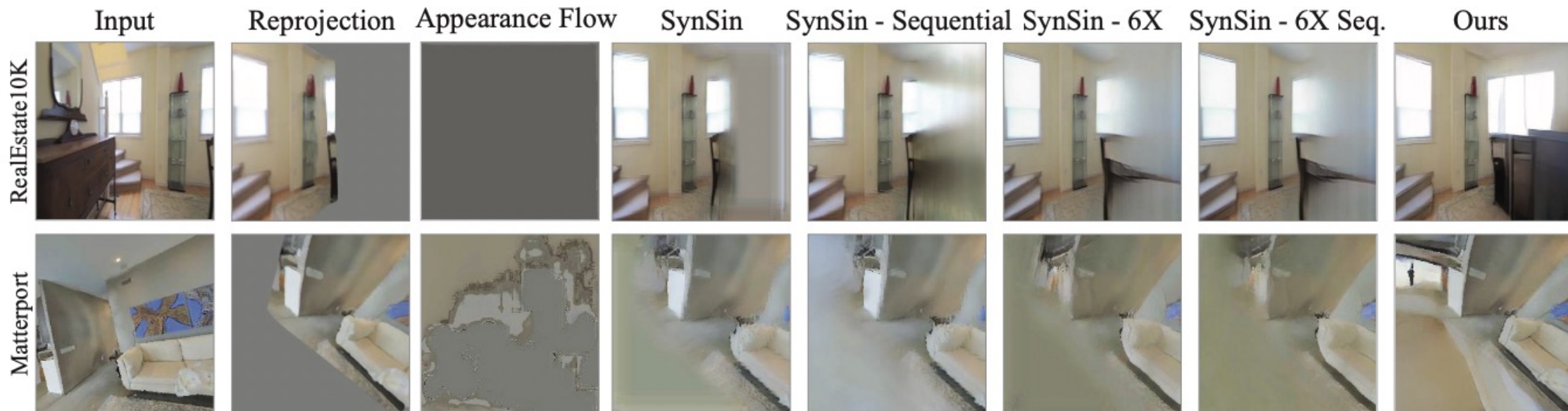


Experiments



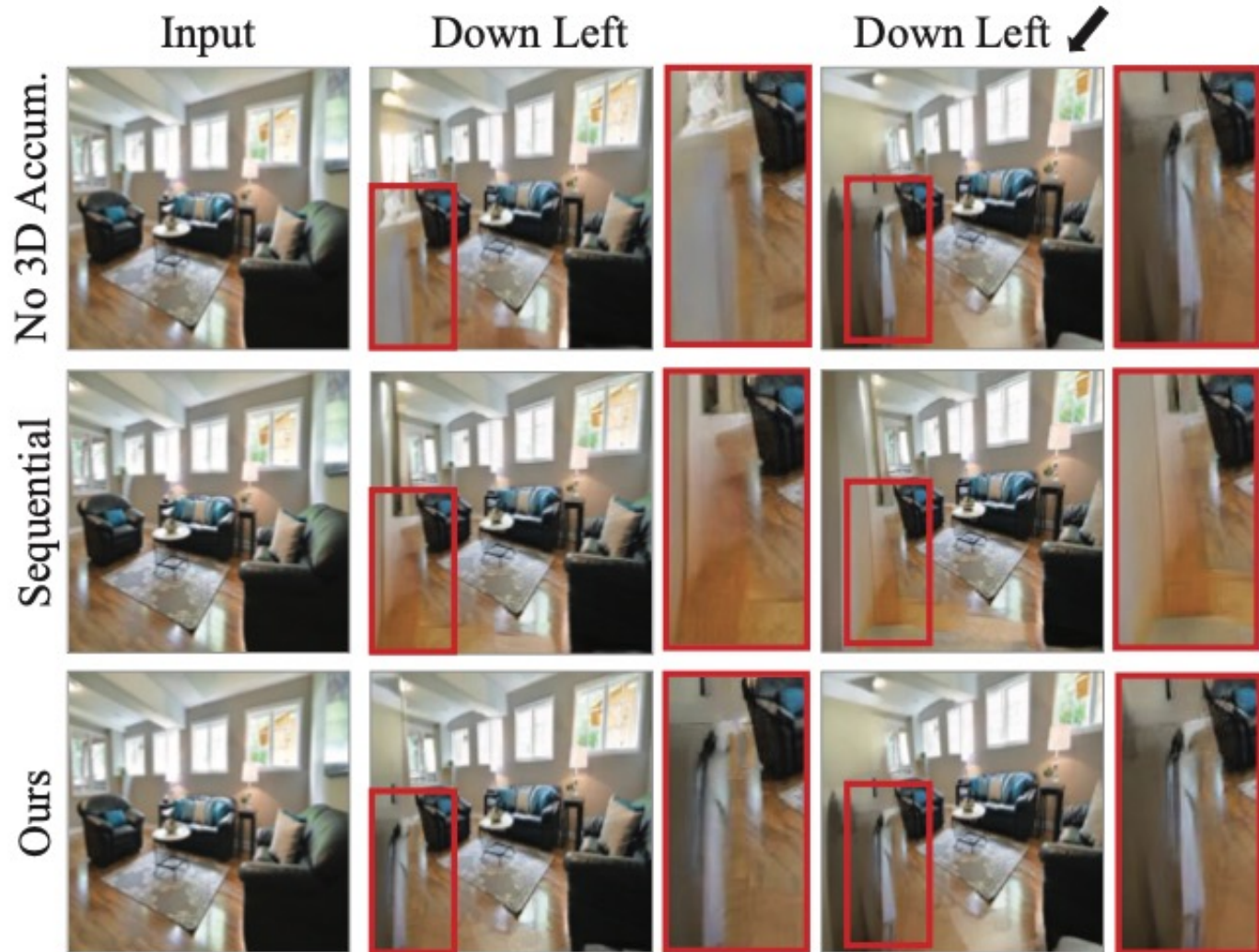
Method	Matterport		RealEstate	
	A/B ↑	FID ↓	A/B ↑	FID ↓
Tatarchenko <i>et al.</i> [45]	0.0%	427.0	0.0%	256.6
Appearance Flow [64]	19.8%	95.8	1.9%	248.3
Single-View MPI [47]	-	-	2.7%	74.8
SynSin [53]	14.8%	72.0	5.8%	34.7
SynSin - Sequential	19.5%	77.8	11.5%	34.9
SynSin - 6X	27.3%	70.4	22.0%	27.9
SynSin - 6X, Sequential	21.2%	79.3	14.4%	33.1
Ours	-	56.4	-	25.5

Experiments



Method	Matterport		RealEstate10K	
	PSNR \uparrow	Perc Sim \downarrow	PSNR \uparrow	Perc Sim \downarrow
Tatarchenko <i>et al.</i> [45]	13.72	3.82	10.63	3.98
Appearance Flow [64]	13.16	3.68	11.95	3.95
Single-View MPI [47]	-	-	12.73	3.45
SynSin - 6X, Sequential	15.61	3.17	14.21	2.73
Ours	14.60	3.17	13.10	2.88

Experiments



Method	A/B vs. Ours ↑	
	Matterport	RealEstate10K
No 3D Accumulation	22.6%	7.5%
Sequential Generation	44.0%	36.2%
Ours	-	-

Introduction

Network

Experiments

Conclusion

Conclusion

