



## Novel View Synthesis from only a 6-DoF Camera Pose by Twostage Networks

Xiang Guo, Bo Li, Yuchao Dai \*, Tongxin Zhang, Hui Deng

School of Electronics and Information, Northwestern Polytechnical University, Xi'an, China

# 

## Content

- Introduction
- Approach
- Results
- Ablation Studies
- Conclusion

#### The task of Novel View Synthesis (NVS)







Choi et al., ICCV'19



The Dependency of NVS in Synthesizing Process

Camera Pose

#### The Dependency of NVS in Synthesizing Process

- Camera Pose
- + Input Images for Reference



Zhou et al., ECCV'16



The Dependency of NVS in Synthesizing Process

- Camera Pose
- + Input Images for Reference
- + Depth



Choi et al., ICCV'19

- The Dependency of NVS
- Camera Pose
- + Input Images for Reference+ Depth
- +3D model



Pittaluga et al., CVPR'19



Mildenhall et al., ECCV'20



Riegler and Koltun, ECCV'20

- The Dependency of NVS
- Camera Pose
- + Input Images for Reference
  + Depth
  + 3D model



Mildenhall et al., ECCV'20 (b) Coarse Image (c) Refined Image (a) Camera Pose (d) Ground Truth

Only Pose as Input (Ours)

Treat rendering and scene representation as one task

- Camera Pose
- + Input Images for Reference
- + Depth
- +3D model

#### Limitation

- Accuracy
- Availability
- Memory Consumption
- Computation complexity





- Camera Pose
- + Input Images for Reference
- + Depth
- +3D model

VS



• Camera Pose Only!



(b) Coarse Image



(c) Refined Image





(a) Camera Pose

(d) Ground Truth

Only Pose as Input (Ours)



## Approach - Overview



Overview of our method. The scene information is embedded by the network weights during Training





GenNet network structure. Expand pose in channel dimension





- U-Net Style, follow pix2pix [1]
- Loss: L1 norm, Perceptual Loss, Adversarial Loss

$$G^* = \arg\min_{G} \max_{D} l_{cGAN}(G, D) + \lambda_1 l_{L1}(G) + \lambda_2 l_{style}(G) + \lambda_3 l_{content}(G)$$

[1] P. Isola, J. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2017, pp. 5967–5976.

#### Results



OldHospital



#### Examples on Cambridge Landmarks



ShopFacade



KingsCollege







## Results

#### Results on Cambridge Landmarks [1]

	GreatCourt			KingsCollege			OldHospital		
	Coarse	Refined (PL)	Refined (w/o PL)	Coarse	Refined (PL)	Refined (w/o PL)	Coarse	Refined (PL)	Refined (w/o PL)
SSIM	0.4361	0.3916	0.3903	0.2953	0.2397	0.2370	0.1897	0.1300	0.1308
PSNR	11.5266	11.3886	11.3648	12.9879	12.5296	12.5288	12.4578	11.5834	11.5353
L1	0.2023	0.2050	0.2055	0.1706	0.1803	0.1800	0.1785	0.1977	0.1987
Brenner	191.3672	466.3329	471.1897	232.5756	744.3240	752.2070	180.5834	1360.8733	1368.9296
	ShopFacade			StMarysChurch			Street		
		ShopFacad	le		StMarysChu	ırch		Street	
	Coarse	ShopFacad Refined (PL)	le Refined (w/o PL)	Coarse	StMarysChu Refined (PL)	Refined (w/o PL)	Coarse	Street Refined (PL)	Refined (w/o PL)
SSIM	Coarse 0.2495	ShopFacad Refined (PL) 0.1768	le Refined (w/o PL) 0.1372	Coarse 0.2962	StMarysChu Refined (PL) 0.2172	Refined (w/o PL) 0.2158	Coarse 0.2255	Street Refined (PL) 0.1303	Refined (w/o PL) 0.1372
SSIM PSNR	Coarse 0.2495 12.7433	ShopFacad           Refined (PL)           0.1768           11.9713	le Refined (w/o PL) 0.1372 11.9403	Coarse 0.2962 12.8790	StMarysChu Refined (PL) 0.2172 12.4334	rch Refined (w/o PL) 0.2158 12.3585	Coarse 0.2255 10.3989	Street           Refined (PL)           0.1303           9.7498	Refined (w/o PL) 0.1372 9.6893
SSIM PSNR L <sub>1</sub>	Coarse 0.2495 12.7433 0.1818	ShopFacad           Refined (PL)           0.1768           11.9713           0.1951	le Refined (w/o PL) 0.1372 11.9403 0.1955	Coarse 0.2962 12.8790 0.1777	StMarysChu           Refined (PL)           0.2172           12.4334           0.1855	rch Refined (w/o PL) 0.2158 12.3585 0.1867	Coarse 0.2255 10.3989 0.2456	Street           Refined (PL)           0.1303           9.7498           0.2633	Refined (w/o PL) 0.1372 9.6893 0.2653
SSIM PSNR L <sub>1</sub> Brenner	Coarse 0.2495 12.7433 0.1818 113.3604	ShopFacad           Refined (PL)           0.1768           11.9713           0.1951           692.7051	le Refined (w/o PL) 0.1372 11.9403 0.1955 697.9377	Coarse 0.2962 12.8790 0.1777 117.7310	StMarysChu           Refined (PL)           0.2172           12.4334           0.1855           538.1169	rch Refined (w/o PL) 0.2158 12.3585 0.1867 619.7828	Coarse 0.2255 10.3989 0.2456 109.2868	Street           Refined (PL)           0.1303           9.7498           0.2633           870.7560	Refined (w/o PL) 0.1372 9.6893 0.2653 997.1627

QUANTITATIVE EVALUATION OF THE SYNTHESIZED IMAGES QUALITY. THREE MEASURE METHODS WITH REFERENCE IMAGE: SSIM, PSNR,  $L_1$  norm and one method without reference image: Brenner. Coarse means the coarse images generated by GenNet and Refined means refined image by RefineNet with or without Perceptual Loss (PL).

[1] A. Kendall, M. Grimes, and R. Cipolla, "Posenet: A convolutional network for real-time 6-dof camera relocalization," in *Proc. IEEE Int. Conf. Comp. Vis.*, 2015.

# 

### Results

#### Examples on 7-Scenes [1]





[1] J. Shotton, B. Glocker, C. Zach, S. Izadi, A. Criminisi, and A. W. Fitzgibbon, "Scene coordinate regression forests for camera relocaliza- tion in RGB-D images," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2013, pp. 2930–2937.

## **Ablation Studies**

• Effect of RefineNet







PL

w/o PL

gt

ShopFacade





KingsCollege





GreatCourt

#### • Effect of Perceptual Loss









**ICPR**<sup>20</sup>

























### Conclusion



- A new problem configuration of NVS: take only camera pose as input
- A two-stage training strategy which is consisted of two consecutive networks: GenNet and RefineNet, utilizing GAN and perceptual loss.
- Experiments show promising results in generating visually pleasant images
- Limitations: should be trained for each scene; distortion and geometric disalignment



## Thank You!