Introduction

Novel view synthesis is a challenging problem in computer vision and robotics. Different from the existing works, which need the reference images or 3D models of the scene to generate images under novel views, we propose a novel paradigm to this problem. That is, we synthesize the novel view from only a 6-DoF camera pose directly. Although this setting is the most straightforward way, there are few works addressing it. While, our experiments demonstrate that, with a concise CNN, we could get a meaningful parametric model that could reconstruct the correct scenery images only from the 6-DoF pose. To this end, we propose a two-stage learning strategy,

which consists of two **CNNs**: consecutive GenNet and RefineNet. GenNet generates a coarse image from a camera pose. RefineNet is a generative adversarial network that refines the coarse image.



(a) Camera Pose Takes only the camera pose as input to predict the image under input view



Method

Overview

Overview of the whole pipeline. A two-stage network consist of GenNet and RefineNet. The scene information is embedded into model parameters during training.

Novel View Synthesis from only a 6-DoF Camera Pose by Two-stage Networks Xiang Guo, Bo Li, Yuchao Dai *, Tongxin Zhang, Hui Deng School of Electronics and Information, Northwestern Polytechnical University, Xi'an, China



We propose a two stage framework consist of GenNet and RefineNet that connect in a sequence to synthesize an image from a camera pose. In the first stage, we use GenNet to get a coarse from pose P. In the second stage, we use RefineNet to refine the coarse image to get a fine detailed image. GenNet



Structure of GenNet. Use two fully connect layer to expand pose alone channel dimension

RefineNet

We use U-Net structure, following pix2pix. The loss function includes L1 norm, Perceptual Loss and Adversarial Loss.

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





Examples of results on Cambridge Landmarks. For each scene, the first row contains synthesized images, and second row contains corresponding ground truth images

Great		GreatCou	rt	KingsCollege			OldHospital		
SSIM	Coarse 0.4361	Ketined (PL) 0.3916	Refined (w/o PL) 0.3903	Coarse 0.2953	Ketined (PL) 0.2397	Ketined (w/o PL) 0.2370	Coarse 0.1897	Ketined (PL) 0.1300	Ketined (w/o PL) 0.1308
PSNR	11.5266	11.3886	11.3648	12.9879	12.5296	12.5288	12.4578	11.5834	11.5353
L1 Brenner	0.2023	0.2050 466.3329	0.2055 471.1897	0.1706	0.1803 744.3240	0.1800 752.2070	0.1785 180.5834	0.1977 1360.8733	0.1987 1368.9296
	~	ShopFacad	le	~	StMarysChu	irch	~	Street	
SSIM	Coarse 0.2495	Refined (PL) 0.1768	Refined (w/o PL) 0.1372	Coarse 0.2962	Refined (PL) 0.2172	Refined (w/o PL) 0.2158	Coarse 0.2255	Refined (PL) 0.1303	Refined (w/o PL) 0.1372
PSNR	12.7433	11.9713	11.9403	12.8790	12.4334	12.3585	10.3989	9.7498	9.6893
L_1	0.1818	0.1951	0.1955	0.1777	0.1855	0.1867	0.2456	0.2633	0.2653
NTITATIV D ONE ME	E EVALUAT	ION OF THE SY HOUT REFEREN REF	NTHESIZED IMAGE ICE IMAGE: BRENN FINED IMAGE BY R	ES QUALITY VER. COARS EFINENET V	TABLE I TABLE MEAS THREE MEAS SE MEANS THE WITH OR WITH	URE METHODS WI COARSE IMAGES C OUT PERCEPTUAL	TH REFERENTED LOSS (PL)	NCE IMAGE: SS BY GENNET A	SIM, PSNR, L ₁ NG
	n 7-i	REF	FINED IMAGE. DREIN FINED IMAGE BY R BRANCE FINED IMAGE BY R BRANCE FINED IMAGE. DREIN FINED IMAGE. DREIN FINED IMAGE. DREIN FINED IMAGE BY R BRANCE FINED IMAGE BY R BRANCE FINED IMAGE BY R FINED IMAGE BY R FINE	efineNet v	WITH OR WITH	COARSE IMAGES COUT PERCEPTUAL	Loss (PL)		
		Chess		Fire	RedKitchen	Heads		Office	
A	Examp bla	ples of resu images, an tion S	ults on 7Scener nd second roy Studies	es. For e v contai S	each scene, ns corresp	, the first row onding groun	contain d truth i	s synthesiz mages	zed
	Image: set of the s	Series of the se	Image: series of the series of th	ade	PL				
	Image: Constraint of the second se	College	Image: Second	urt	The eff image a The eff remove	Fect of Refine and second refect of Perce cunrealistic a	eNet. The second	he first ro fined imagoss. PL co (top)	w is coarse ge (Left) ould help to
			<u>C</u>	on	<u>clusi</u>	on			
Ve 1	oron	ose a	new	proł	olem	config	urat	ion o	of NVS
ke am	onl ewo	y the ork co	came onsist	of	pose two	as inp consec	out. Cutiv	A two ve ne	o-stag tworks
rom nag	nct nisin ges. ork	and lg real Ther for	sults i e are each s	n g also	eneration e: dig	itation	yisua s: n	ally j eed rest	pleasan to trai

limited generalization ability.





