

Deep Realistic Novel View Generation for City-Scale Aerial Images





Overview

- Introduced a novel end-to-end framework for generation of large scale synthetic aerial image sequences with associated precise ground truth camera metadata
- Proposed a novel deep learning network with an explicit edgemap processing stream to remove image artifacts



ARNet-Edge: Edge-Augmented Artifact Removal Network

Custom U-Net like encoder-decoder network. Learns to map 2D raw synthetic images to realistic-looking synthetic images with reduced artifact

- Encoder: two parallel series of Squeeze and Excitation Resnet (SE-Resnet Blocks). Stream-1: extracts features from RGB input; Stream-2: extracts features from edge-map
- Decoder: deconvolution layers followed by two convolution layers attempting to reconstruct the original image from encoders features gi))
- Combined loss: mean square error + structural similarity

$$L_{ssim}(gt, gi) = \frac{1}{N} \sum_{i=1}^{N} (1 - SSIM(gt, gi))$$
$$L_{mse}(gt, gi) = \frac{1}{N} \sum_{i=1}^{N} (gt - gi)^2$$



Experimental Results

1. Image Quality

Model	SSIM		PSNR	
	ABQ	LA	ABQ	LA
Raw Synthetic	0.573	0.651	20.65	23.53
REDNet	0.629	0.686	22.29	24.35
Deep Image Prior	0.619	0.710	21.12	25.21
ARNet	0.638	0.704	22.40	25.15
ARNet-Edge	0.665	0.721	23.21	25.89



2. Dense 3D Point Cloud Reconstruction

- Reconstructed a dense point cloud from the denoised ABQ image sequence (PC_d) and from the raw ABQ image sequence (PC_r) ; point cloud generated from the original aerial image sequence as the ground truth (PC_{o})
- Computed cloud-to-cloud distance



 PC_d vs PC_o

 PC_r vs PC_o