



# CityDreamer

## Compositional Generative Model of Unbounded 3D Cities

Authors: Haozhe Xie et al.

Affiliation: Nanyang Technological University

Conference: CVPR 2024

Presenter: Zhiguo Liu  
Date: November 19



# Challenges: Previous Works



GANCraft [CVPR'21]



InfiniCity [ICCV'23]



SceneDreamer [arXiv 2303.01330]

# Challenges: Natural Data vs City Data



## Similarity

- Objects are similar in natural scenes



- Buildings are diverse in cities

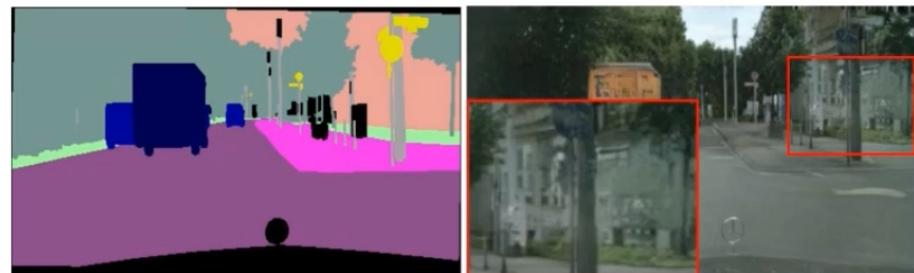


## Realistic

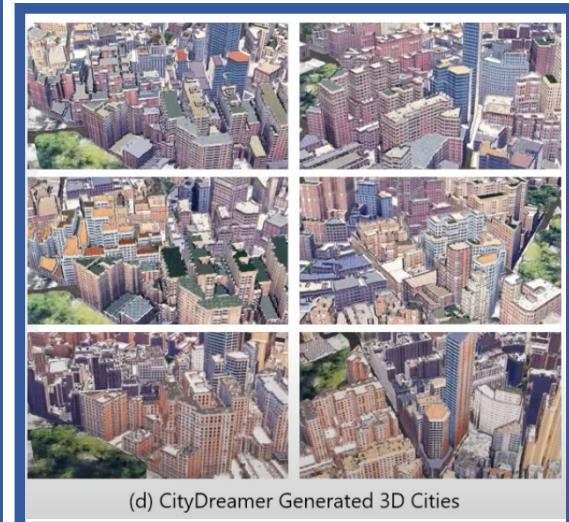
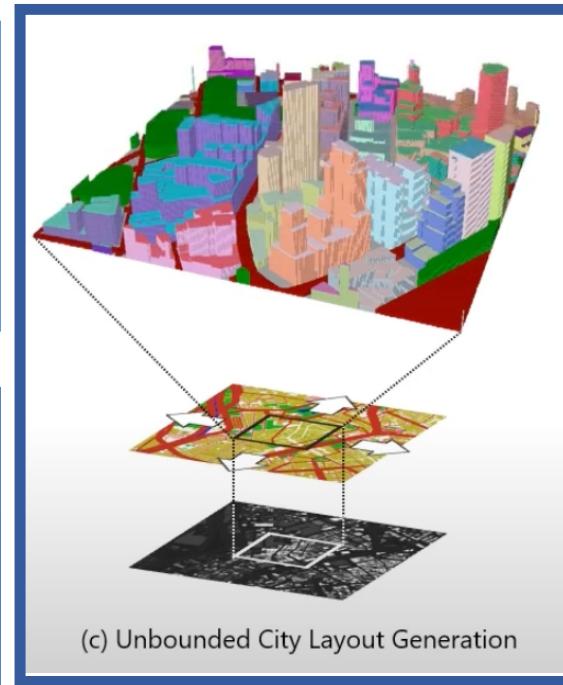
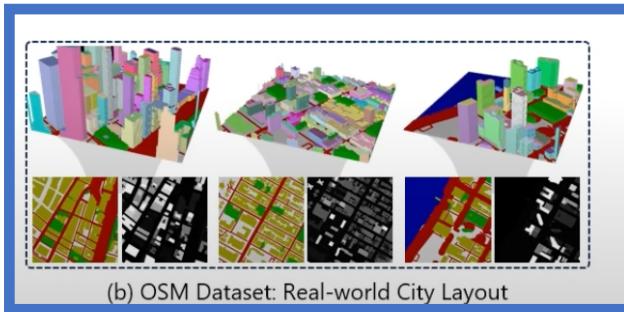
- Synthetic nature images are realistic



- Synthetic city images are not realistic



# The Proposed Method: Datasets



# The Proposed Method: Structure



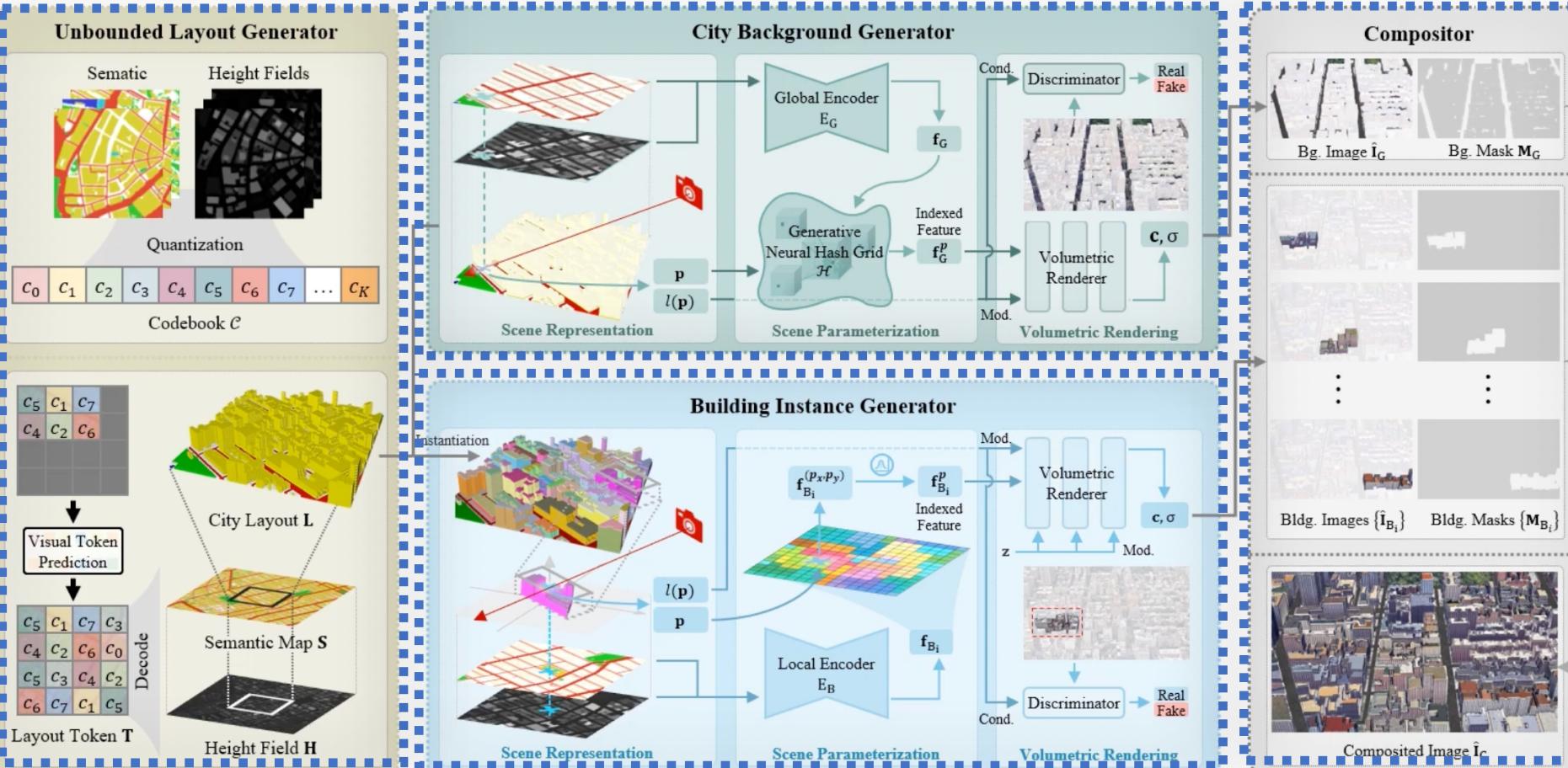
Unbounded Layout Generator

City Background Generator

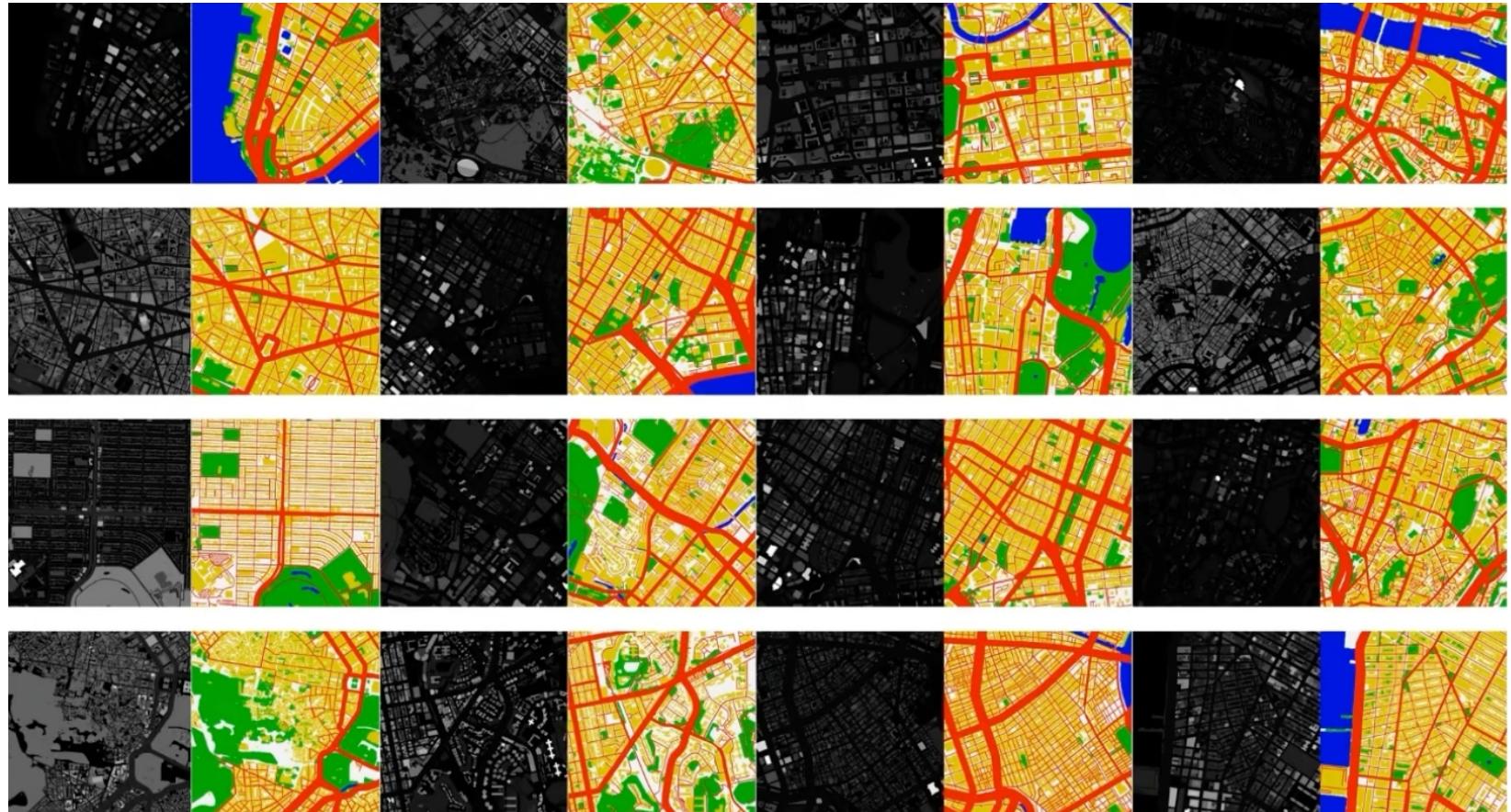
Compositor

Building Instance Generator

# The Proposed Method: Structure



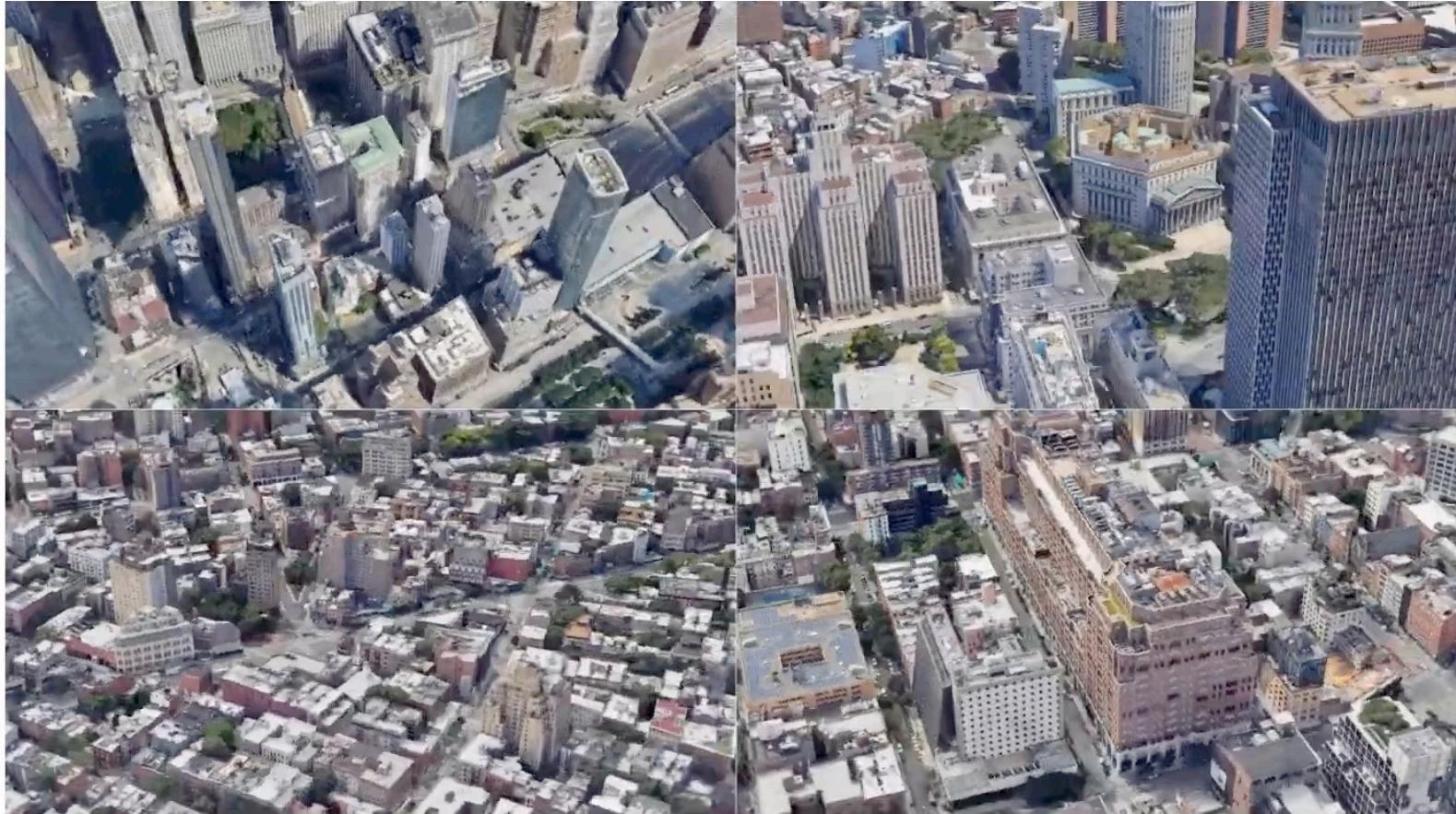
# The OSM Dateset



# The OSM Dateset



# The GoogleEarth Dateset



# Comparison to SOTA Methods



PersistentNature



SceneDreamer



CityDreamer



# Details

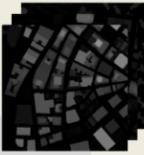


## Unbounded Layout Generator (§3.1)

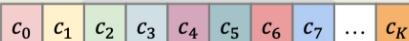
Sematic Maps



Height Fields



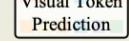
Quantization



Codebook  $\mathcal{C}$



Visual Token Prediction



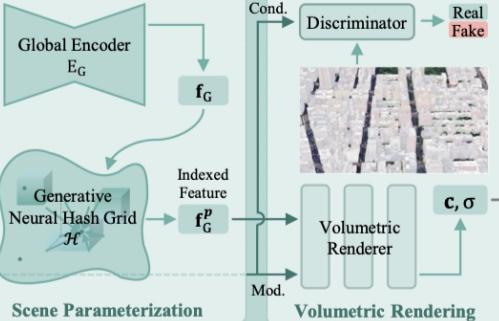
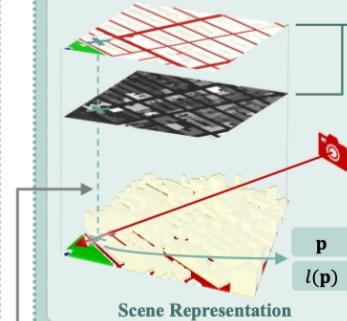
Decode



Layout Token  $T$

Height Field  $H$

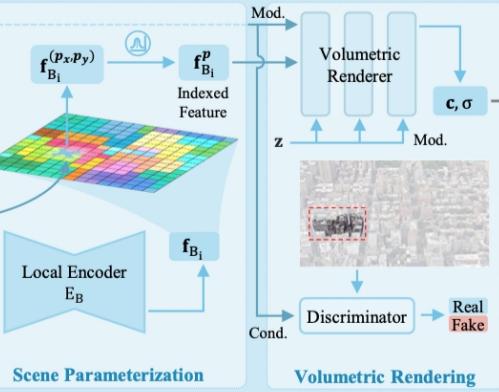
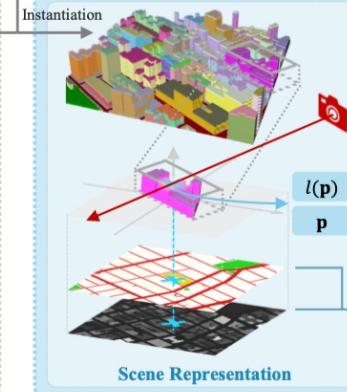
## City Background Generator (§3.2)



## Compositor (§3.4)

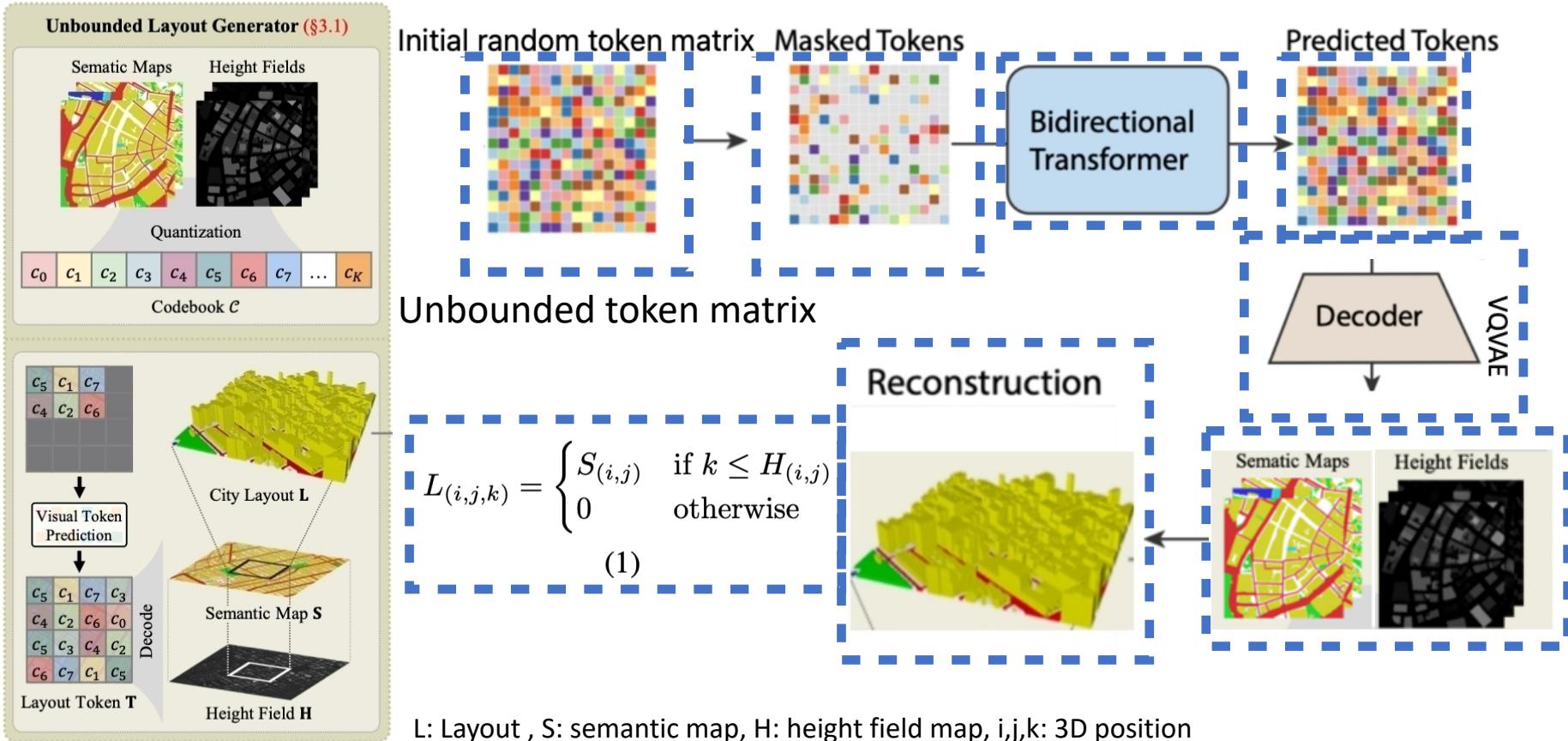


## Building Instance Generator (§3.3)

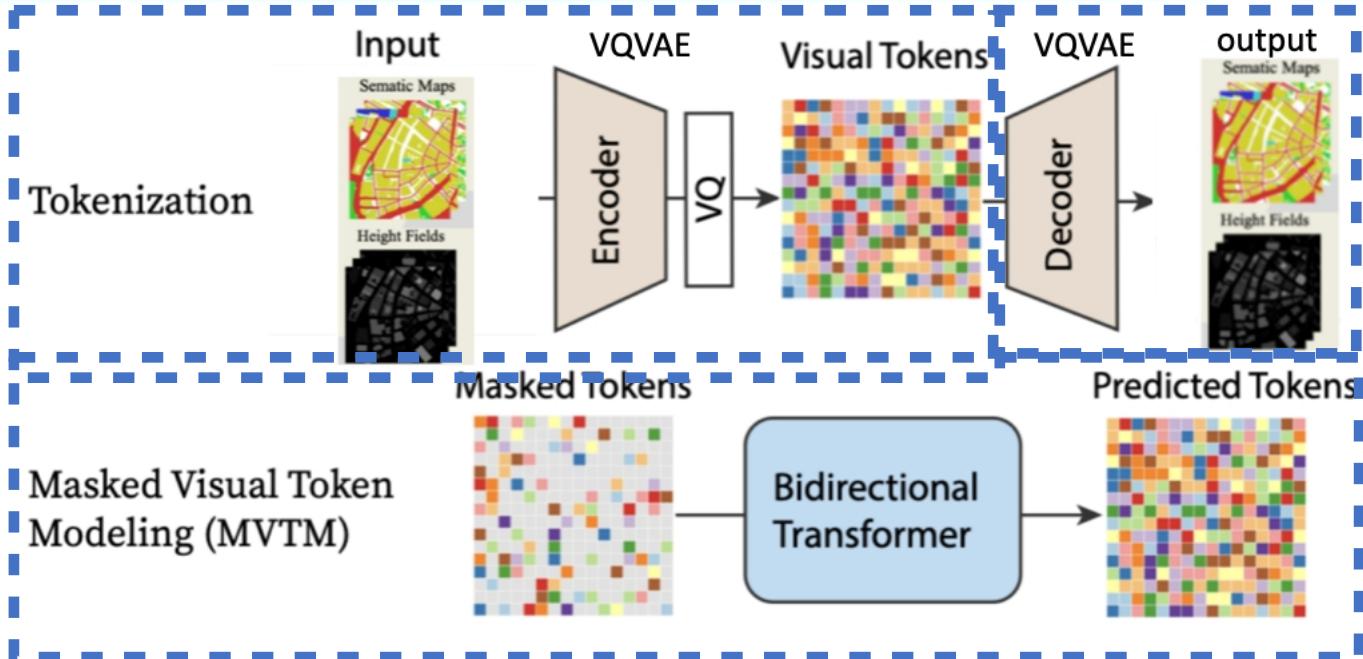
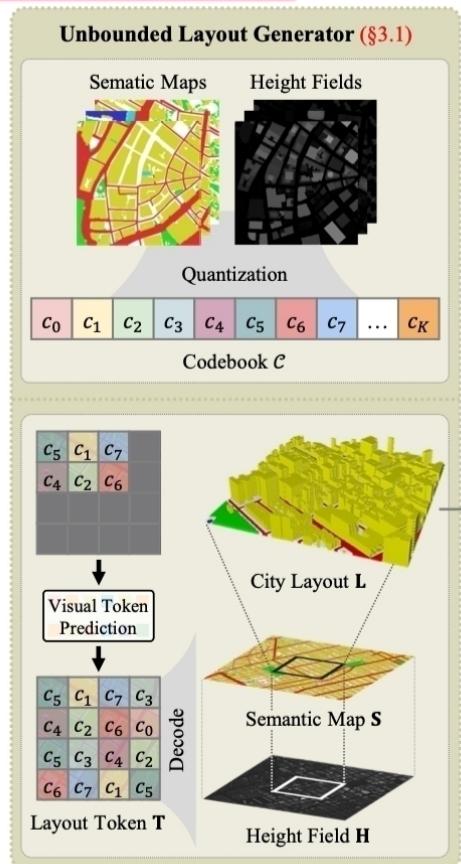


Composed Image  $\hat{I}_C$

# Unbounded Layout Generator: Maskgit Step



# Unbounded Layout Generator: Maskgit Model



$$\text{Reweighted\_ELBO} = w_1 \text{ELBO} + w_2 D_{KL}$$

ELBO: Evidence Lower Bound loss,  $D_{KL}$ : Kullback-Leibler Divergence  $w_1, w_2$ : weight parameter

# Unbounded Layout Generator—VQVAE



Unbounded Layout Generator (§3.1)

Sematic Maps Height Fields



Quantization

$c_0 \ c_1 \ c_2 \ c_3 \ c_4 \ c_5 \ c_6 \ c_7 \ \dots \ c_K$

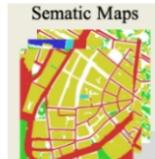
Codebook  $\mathcal{C}$

$$Loss = L_1 Loss + SmoothnessLossS + CrossEntropyLossE$$

$$\ell_{\text{VQ}} = \lambda_R \|\hat{\mathbf{H}}_p - \mathbf{H}_p\| + \lambda_S \mathcal{S}(\hat{\mathbf{H}}_p, \mathbf{H}_p) + \lambda_E \mathcal{E}(\hat{\mathbf{S}}_p, \mathbf{S}_p) \quad (2)$$

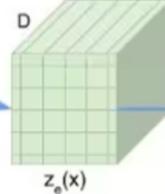
$\ell_{\text{VQ}}$ : VQVAE loss function,  $\lambda_R, \lambda_S, \lambda_E$ : weight parameter

$\hat{\mathbf{H}}_p, \hat{\mathbf{S}}_p$ : Generated Height field and semantic map patches,  $\mathbf{H}_p, \mathbf{S}_p$ : corresponding ground truth.

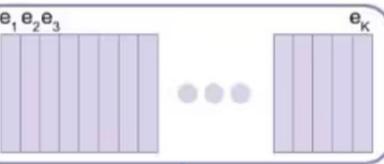


Height Fields

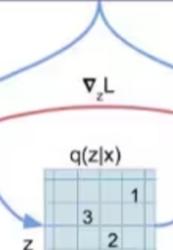
CNN



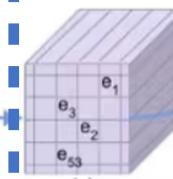
Encoder



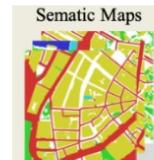
Embedding Space



Quantization



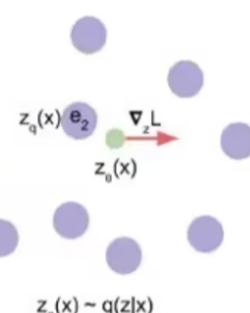
Decoder



Height Fields

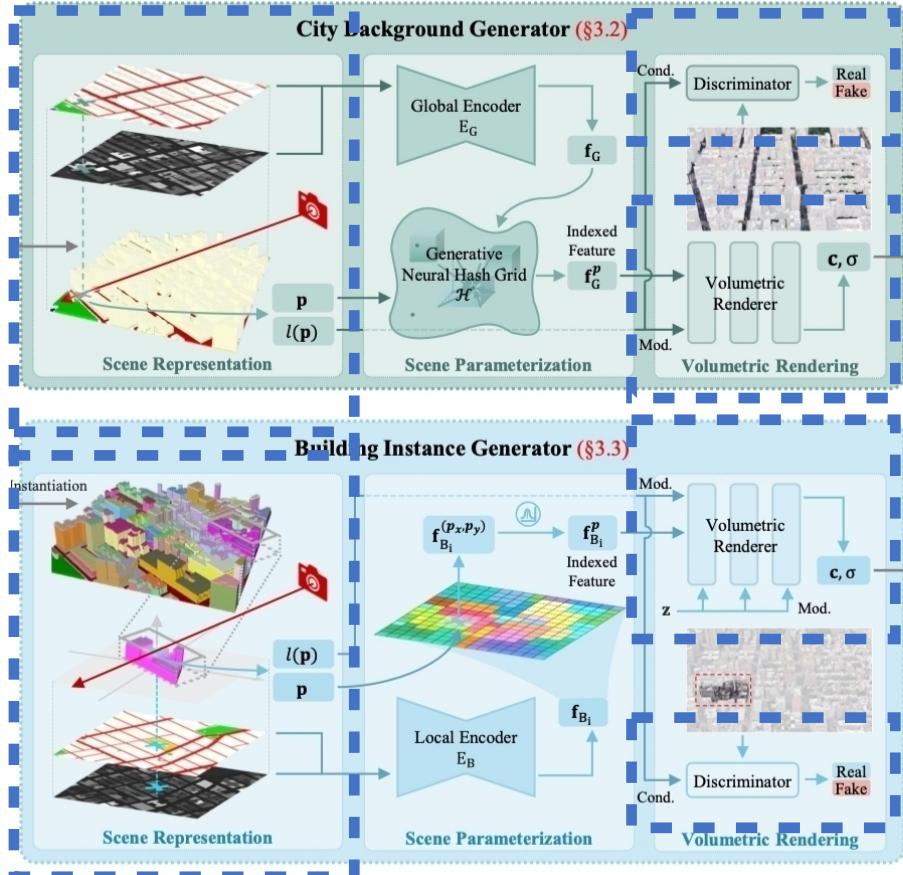
CNN

$p(x|z_q)$



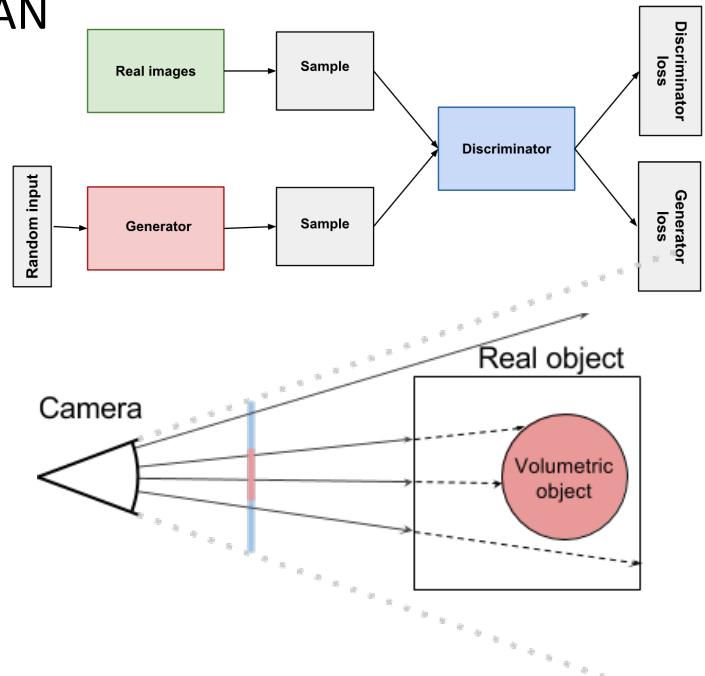
codebook

# City Background vs Building Generator: Same

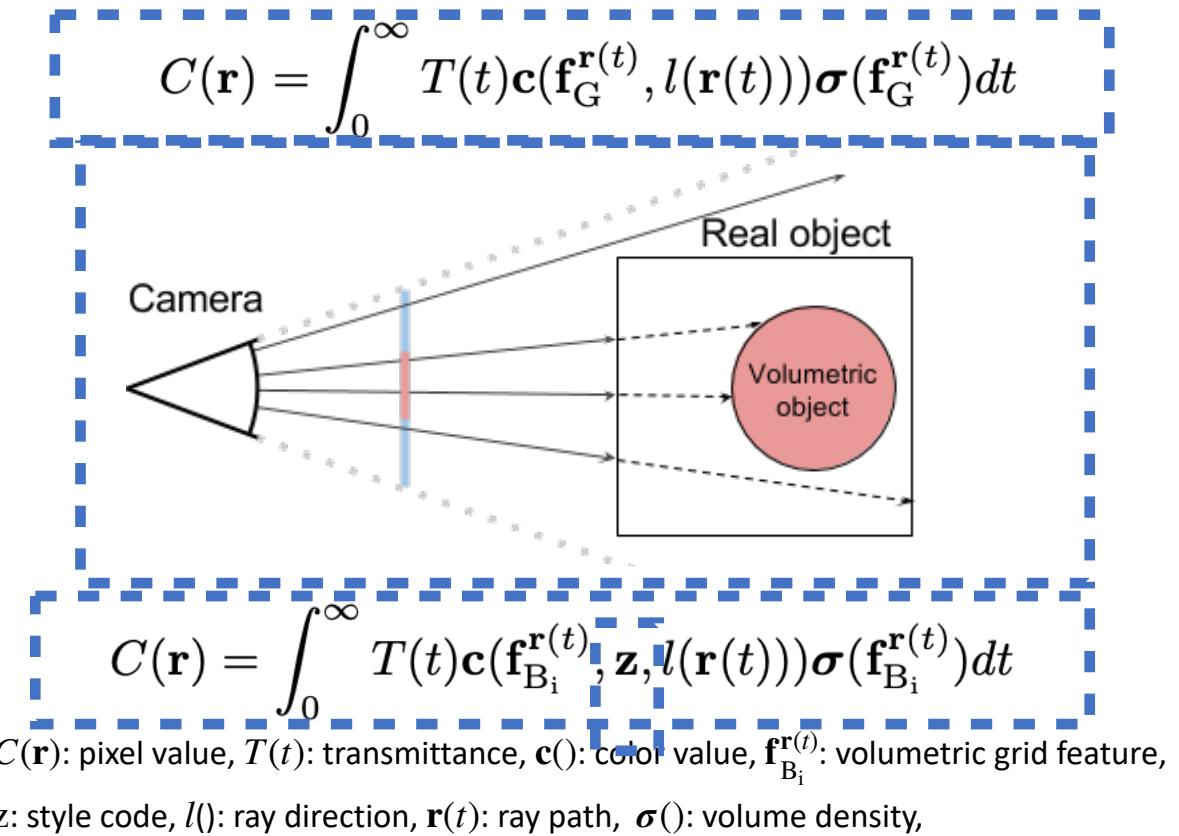
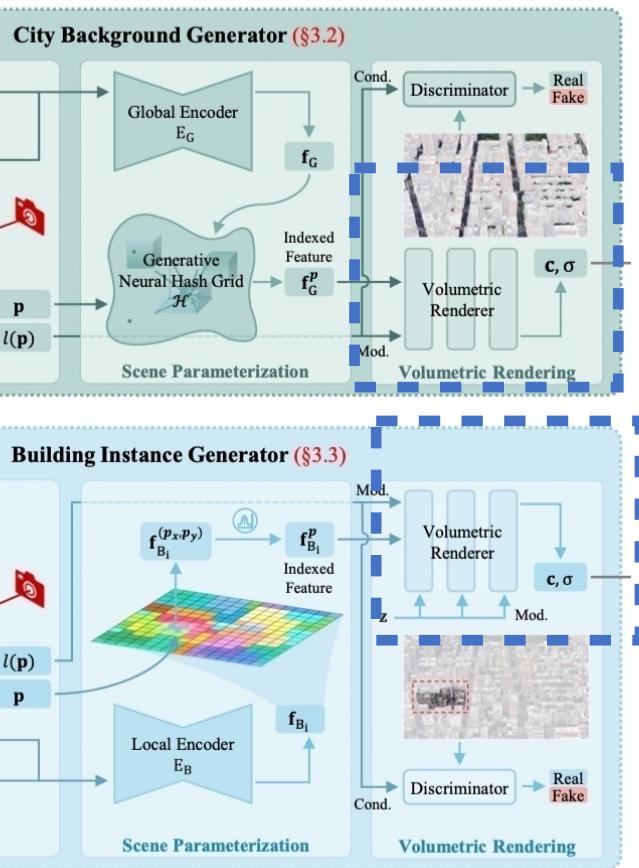


Same

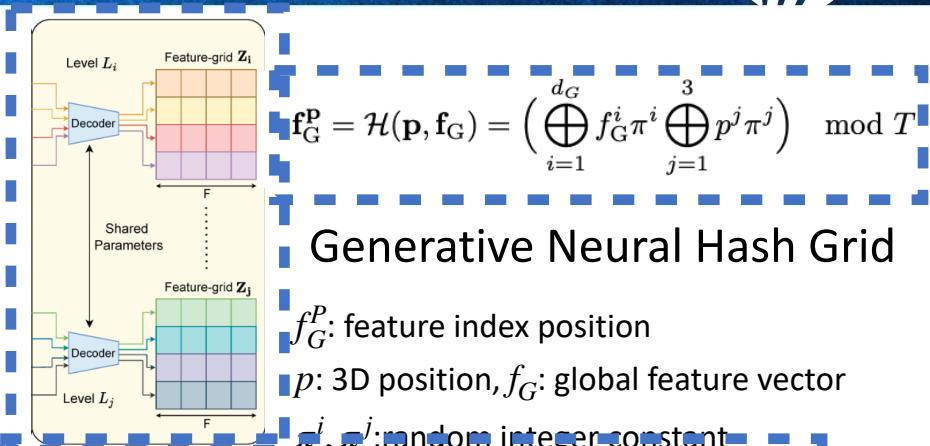
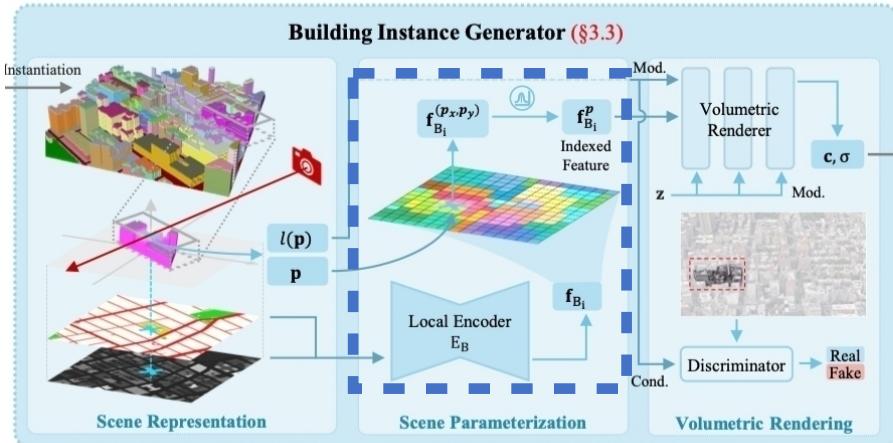
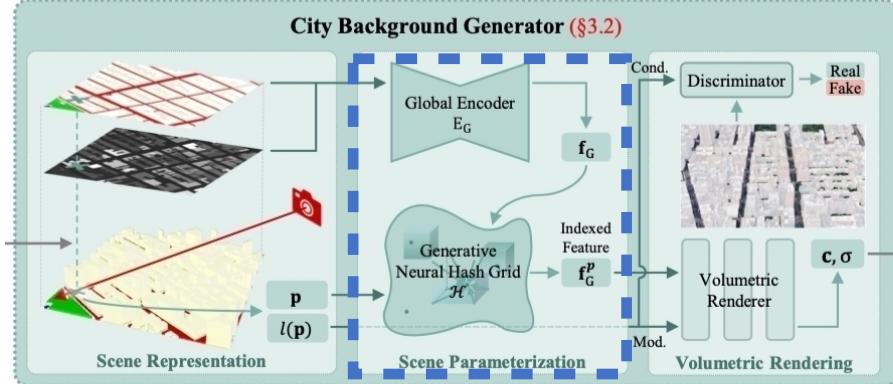
- Bird'd Eye View
- Volumetric Renderer
- GAN



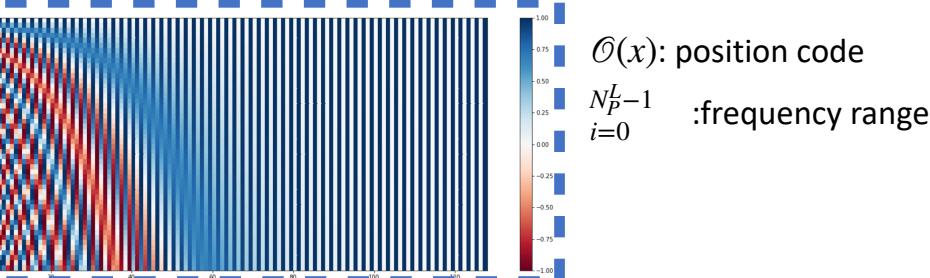
# City Background vs Building Generator: Volumetric Renderer



# City Background & Building Generator: different encode strategy

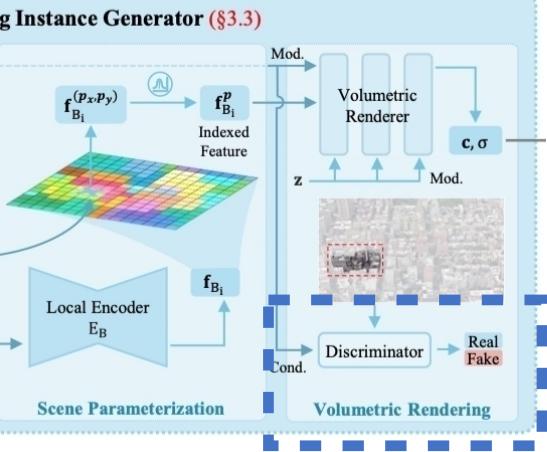
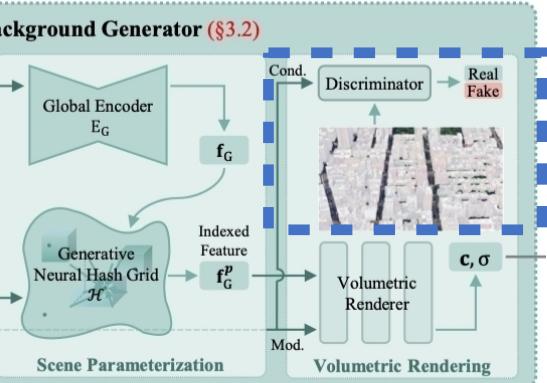


$$\mathcal{O}(x) = \{\sin(2^i \pi x), \cos(2^i \pi x)\}_{i=0}^{N_P^L - 1} \quad (9)$$

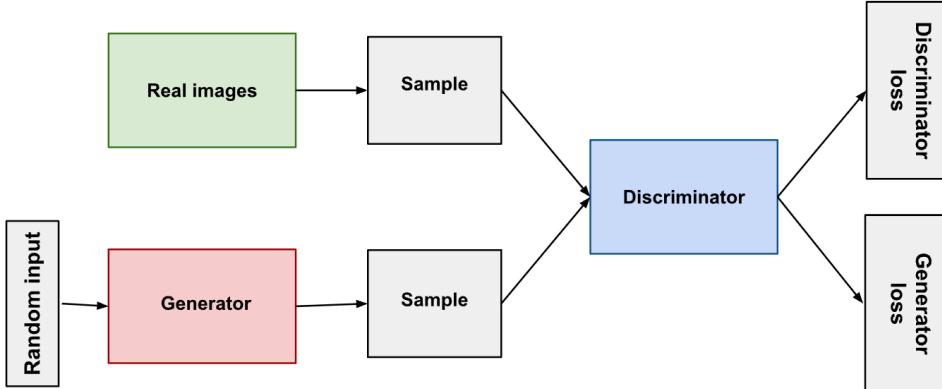


Position encode

# City Background vs Building Generator - Diff GAN Loss



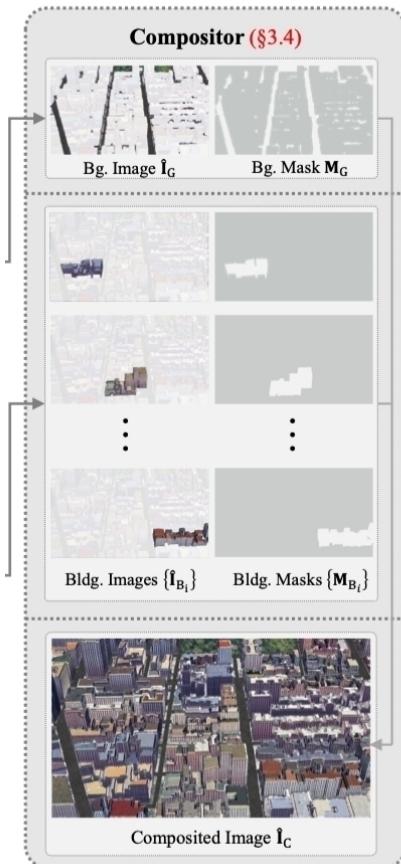
$$\ell_G = \lambda_{L1} \|\hat{\mathbf{I}}_G - \mathbf{I}_G\| + \lambda_P \mathcal{P}(\hat{\mathbf{I}}_G, \mathbf{I}_G) + \lambda_G \mathcal{G}(\hat{\mathbf{I}}_G, \mathbf{S}_G) \quad (6)$$



$$\ell_B = \mathcal{G}(\hat{\mathbf{I}}_{B_i}, \mathbf{S}_{B_i}) \quad (11)$$

$\ell_G$ : Total generation loss,  $\lambda_{L1}, \lambda_P, \lambda_G$ : Weight factor,  $\hat{\mathbf{I}}_G$ : Generated image,  $\mathbf{I}_G$ : Real image,  $\mathbf{S}_G$ : Semantic map,  $\|\cdot\|$ : L1 Loss, P: perceptual loss,  $\mathcal{G}$ : GAN loss

# Compositor



$$\mathbf{I}_C = \hat{\mathbf{I}}_G \mathbf{M}_G + \sum_{i=1}^n \hat{\mathbf{I}}_{B_i} \mathbf{M}_{B_i} \quad (12)$$

$\mathbf{I}_C$ : Generated complete city image,  $\hat{\mathbf{I}}_G$ : Background image,  $\mathbf{M}_G$ :Background semantic mask  
 $\hat{\mathbf{I}}_{B_i}$ : Building instance image,  $\mathbf{M}_{B_i}$ : Semantic mask of the building

# Thoughts



Q&A



Q&A