

SinNeRF: Training Neural Radiance Fields on Complex Scenes from a Single Image

Dejia Xu*, Yifan Jiang*, Peihao Wang, Zhiwen Fan, Humphrey Shi, Zhangyang Wang
ECCV 2022

STRUCT Group Seminar
Presenter: Rundong Luo
2023.01.15

OUTLINE

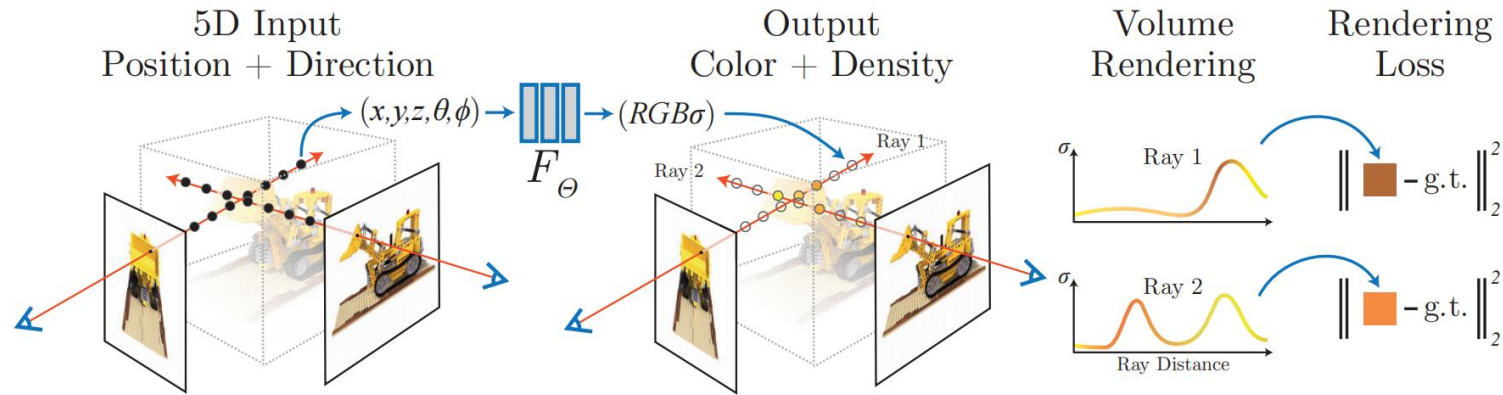
- Authorship
- Background
- Method
- Experiments
- Conclusion

OUTLINE

- Authorship
- **Background**
- Method
- Experiments
- Conclusion

BACKGROUND: NeRF

Overview



Volume rendering

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt, \text{ where } T(t) = \exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s)) ds\right).$$

BACKGROUND: NeRF

Critical points to consider when reading NeRF papers (personally)

- What's target task? Specifically, what's the input and output during inference?
- What's the training input? Multiview? Calibrated? etc.
- Is the model test-time optimization (one NeRF for each scene) or train-time optimization (one NeRF for multiple scenes)

The original NeRF

- Training input: many (>10) calibrated images of the same scene
- Training output: a NeRF for this specific scene
- Inference input: camera pose
- Inference output: an image taken from this pose
- **Test-time optimization**

BACKGROUND: NeRF

Fields of research on NeRF

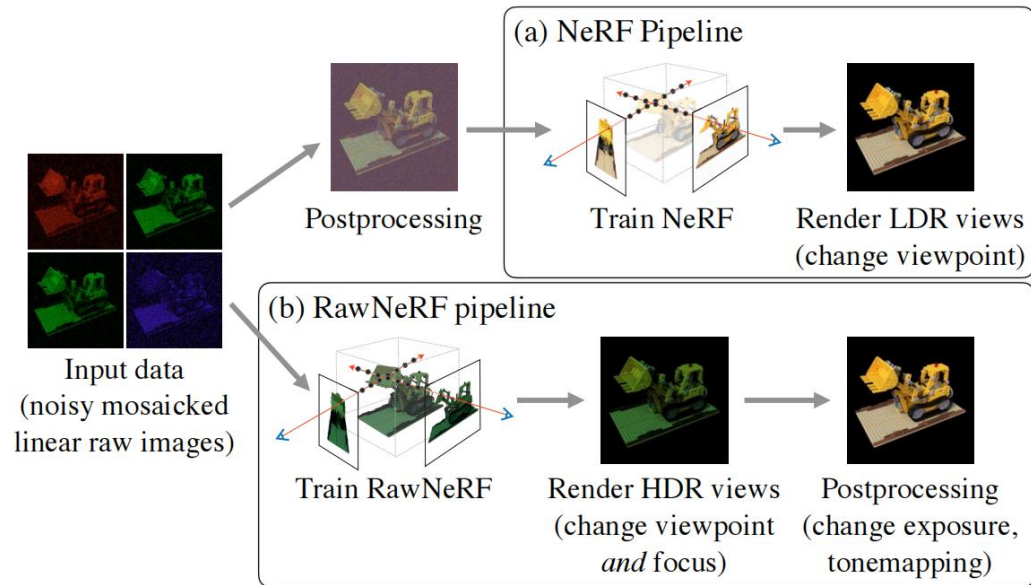
- Faster training & inference
- Artistic effect, e.g. HDR NeRF
- Deformable
- **Generalization**, e.g. RawNeRF, GRAF, PixelNeRF, SinNeRF
- **Compositionality**, e.g. GIRAFFE, uORF

NeRF paper collection:

<https://github.com/awesome-NeRF/awesome-NeRF>

Previously on NeRF

RawNeRF (Mildenhall et al. CVPR 2022)

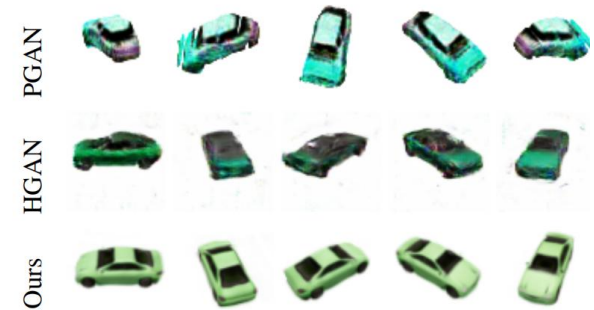
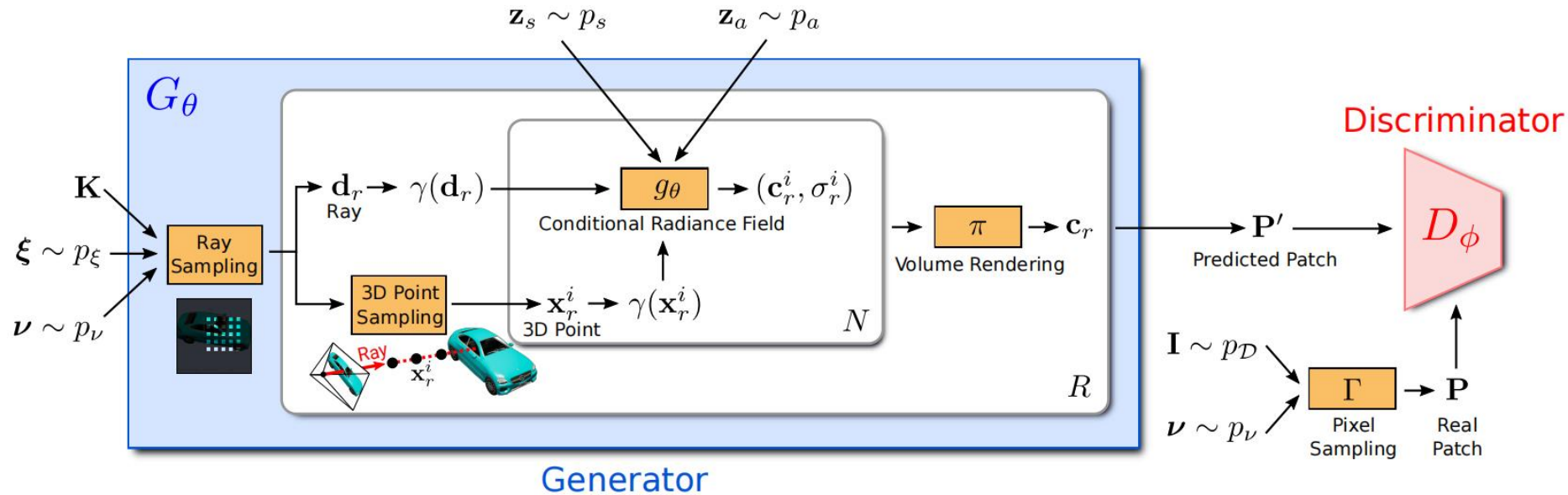


NeRF for denoising

- Training input identical to original NeRF
- Modified loss function & training inputs (multiple shutter speed)

Previously on NeRF

GRAF (Schwarz et al. NeurIPS 2020)

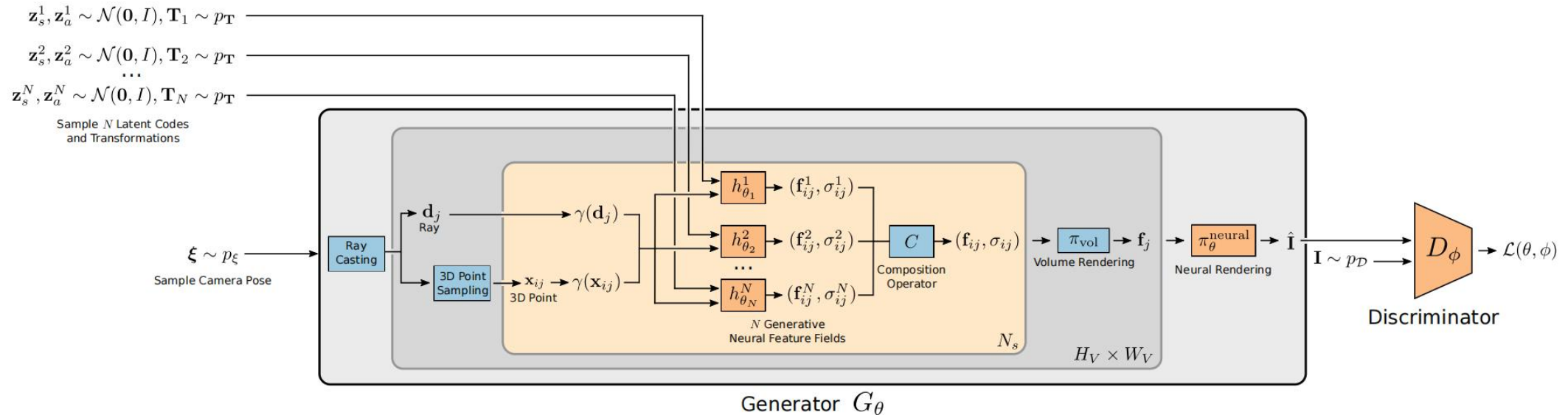


Generative Radiance Fields

- Training input: some real image patches of the same category (e.g. cars)
- Training output: NeRF for this type of input
- **Trained on GAN loss (instead of pixel loss in original NeRF)**
- Inference input: sampled from certain distributions, but no exact meaning
- Inference output: novel images of this type (conditioned on inputs)

Previously on NeRF

GIRAFFE (Niemeyer et al. CVPR 2021)

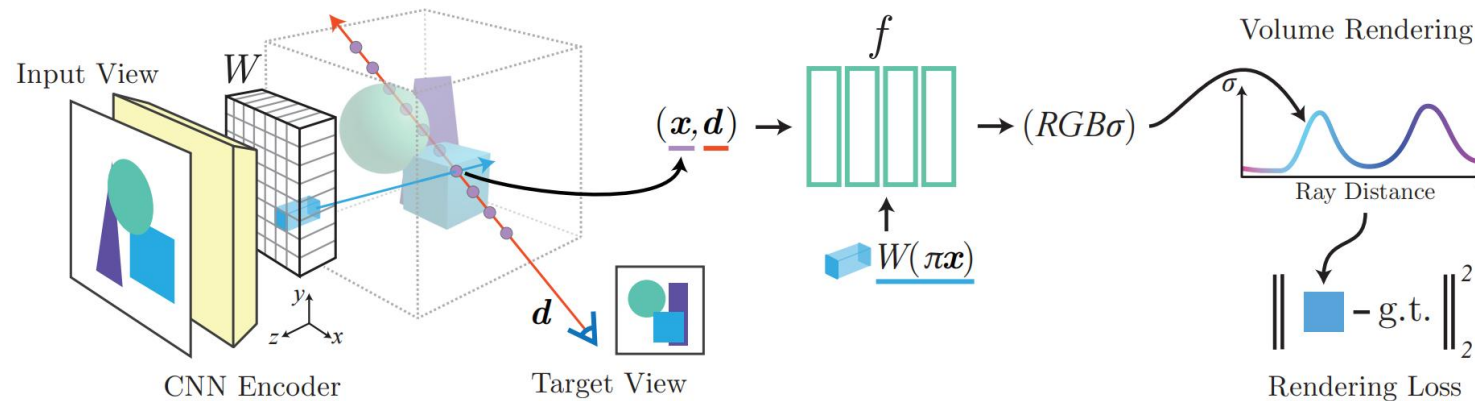


Generative Radiance Fields for Multiple Objects

- Training input: some real scenes with one or multiple objects
- Training output: encoder, NeRF **in scene space**
- Inference input: sampled from certain distributions. Inference latents, then stack together
- Inference output: novel images

Previously on NeRF

PixelNeRF (Yu et al. CVPR 2021)

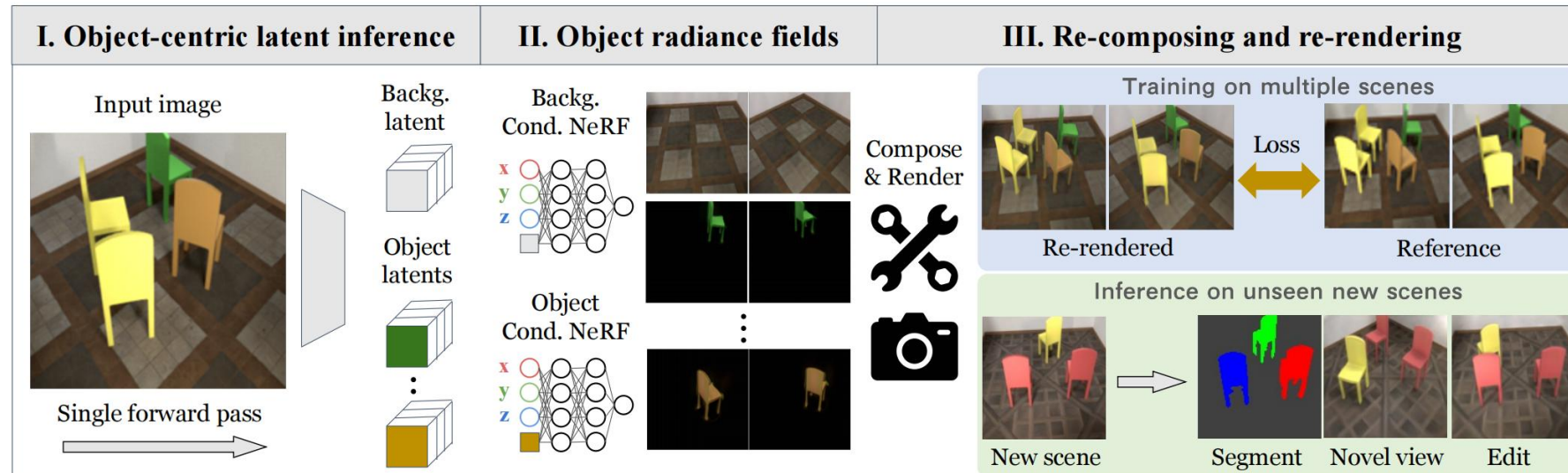


Single image NeRF for generating novel views

- **Train-time optimization**
- **Feature encoder for synthesis**
- Training input: **multi-view** (calibrated) scenes
- Training output: NeRF in view space, takes the reference view feature as an additional input
- Inference input: a single reference image and the desired camera pose (or relative pose to reference image)
- Inference output: novel images

Previously on NeRF

uORF (Yu et al. ICLR 2022)



Single image NeRF for scene editing & synthesis

- **Train-time optimization**
- Training input: **multi-view** (calibrated) scenes
- Training output: background NeRF in cononical space, object NeRF in view space
- Inference input: a single reference image
- Inference task: novel view synthesis, scene editing, scene segmentation

OUTLINE

- Authorship
- Background
- Method
- Experiments
- Conclusion

METHOD

Task description: generate novel views with only one image (test-time optimization)

Inputs

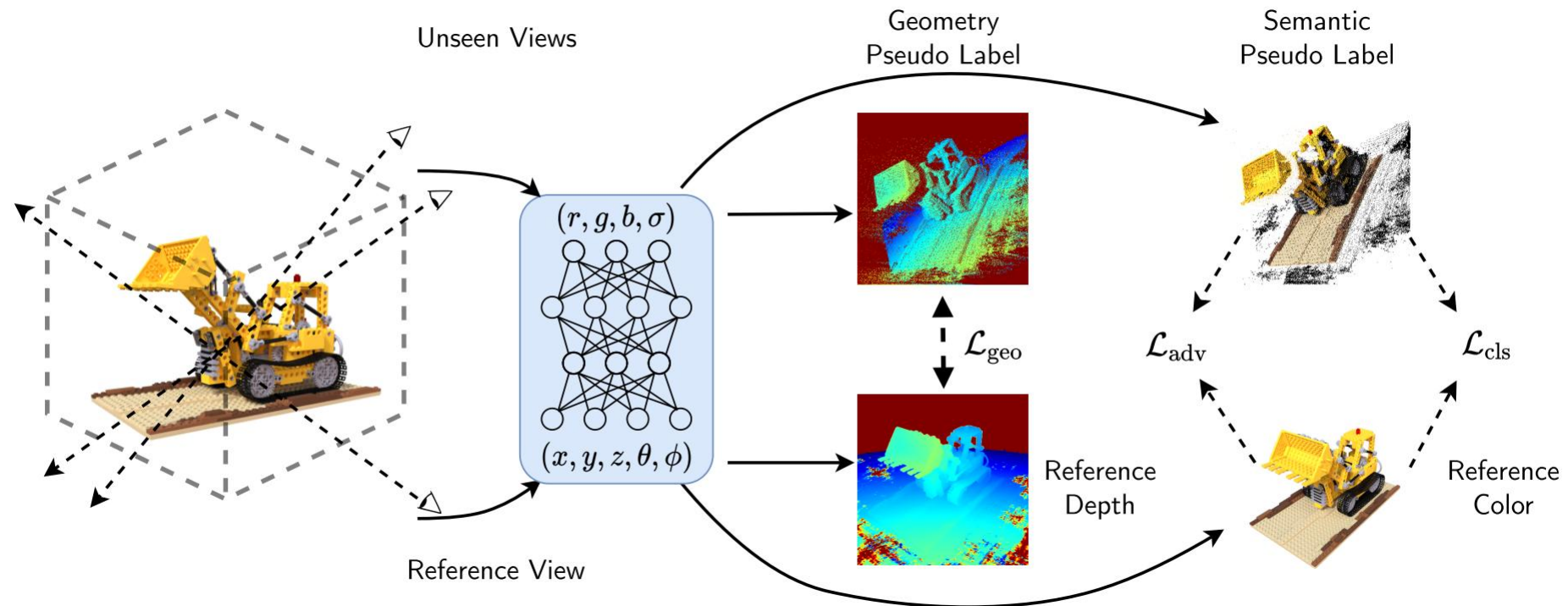
- Single RGB image
- Calibration (intrinsics & extrinsics)
- Depth

Supervision

- Rendering (pixel loss)
- Geometry
- Semantic

METHOD

Overview



$$\mathcal{L}_{total} = \mathcal{L}_{pix} + \lambda_1 \mathcal{L}_{geo} + \lambda_2 \mathcal{L}_{adv} + \lambda_3 \mathcal{L}_{cls}$$

METHOD

Geometry supervision

- Reconstruct 3D geometry through warping

Warp from reference view i to unseen view j

$$p_j = K_{\text{unseen}} T(K_{\text{ref}}^{-1} Z_i p_i)$$

Multi-view geometry revisited

成像原理: (以下都是 Homogeneous Coordinates)

$$\begin{pmatrix} \lambda x \\ \lambda y \\ \lambda \end{pmatrix} = \begin{pmatrix} \alpha f & 0 & P_x \\ 0 & \alpha f & P_y \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} R_{3 \times 3} & t_{3 \times 1} \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix}$$

↑
图像中像素点的位置

↑
Intrinsics

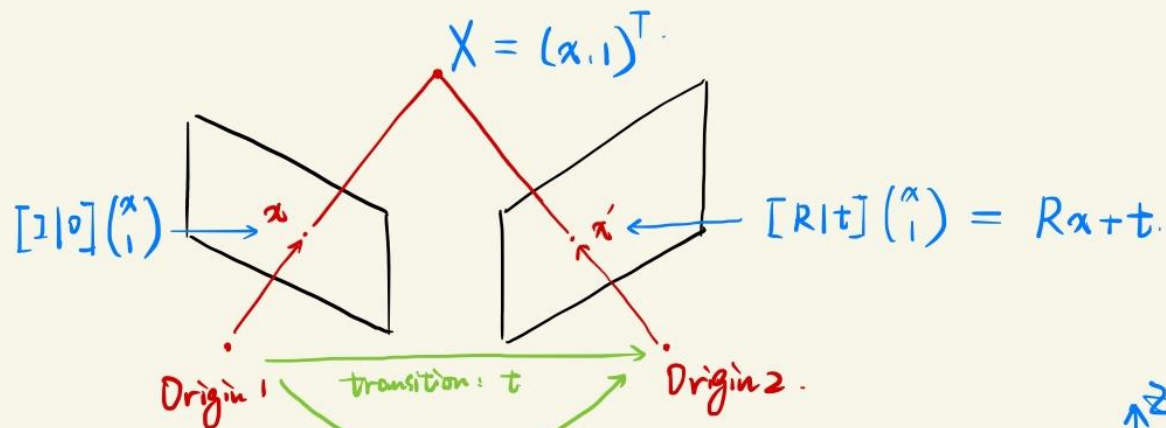
↑
Projection

↑
Extrinsics
(世界坐标系到相机坐标系)

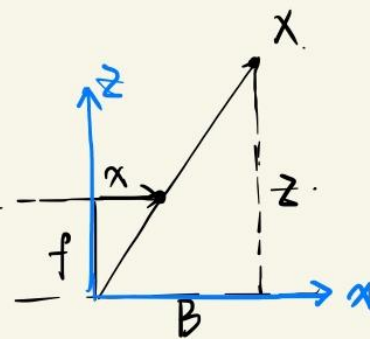
↑
三维坐标

Multi-view geometry revisited

以第一个相机的相机坐标系作为世界坐标系



从二维 $\begin{pmatrix} x \\ y \\ 1 \end{pmatrix}$ 到三维点坐标 $\begin{pmatrix} z \cdot x \\ z \cdot y \\ z \end{pmatrix}$ Image Plane



nominal unit. ← real units

再乘以内参 $K_i^{-1} \cdot z \cdot p$ 得到相机坐标系中坐标。(图中的 x)

计算 $R\alpha + t$ 得到第二个相机坐标系下的坐标。

再用第二个相机的外参得到最终坐标。

Multi-view geometry revisited

文中的公式: $P_j = K_{unseen} T \cdot (K_{ref}^{-1} z_i \cdot P_i)$

实际版本:

① $T_{unseen} T_{ref}^{-1} \begin{pmatrix} K_{ref}^{-1} z_i \cdot \begin{pmatrix} P_i \\ 1 \end{pmatrix} \\ 1 \end{pmatrix}_{4 \times 3}$

② 丢弃最后一维, 再左乘 K_{unseen} (3×1)

③ 从 Homogeneous coordinates 得到 real coordinates
同时获得了 depth

METHOD

Geometry supervision: reconstruct 3D geometry through warping

- Warp from reference view i to unseen view j

$$p_j = K_{\text{unseen}} T(K_{\text{ref}}^{-1} Z_i p_i)$$

- Smoothness constraint

$$\mathcal{L}_{\text{smooth}}(d_i) = e^{-\nabla^2 \mathcal{I}(\mathbf{x}_i)} (|\partial_{xx} d_i| + |\partial_{xy} d_i| + |\partial_{yy} d_i|)$$

- Consistency constraint

$$\mathcal{L}_{\text{geo}} = \mathcal{L}_1(d_1, f(d_2)) + \mathcal{L}_1(f(d_1), d_2) + \lambda_4 \mathcal{L}_{\text{smooth}}$$

METHOD

Semantic supervision: local texture guidance & global structure guidance

- Discriminate between reference image patches and NeRF outputs

Differentiable augmentation for collapse avoiding

$$\mathcal{L}_D = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [f_D(-D(T(\mathbf{x})))] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [f_D(D(T(G(\mathbf{z}))))]$$

$$\mathcal{L}_G = \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [f_G(-D(T(G(\mathbf{z}))))],$$

$$\mathcal{L}_{\text{adv}} = \mathcal{L}_D + \mathcal{L}_G,$$

- Extract global structure information through ViT

$$\mathcal{L}_{\text{cls}} = \|f_{\text{vit}}(A) - f_{\text{vit}}(B)\|^2$$

METHOD

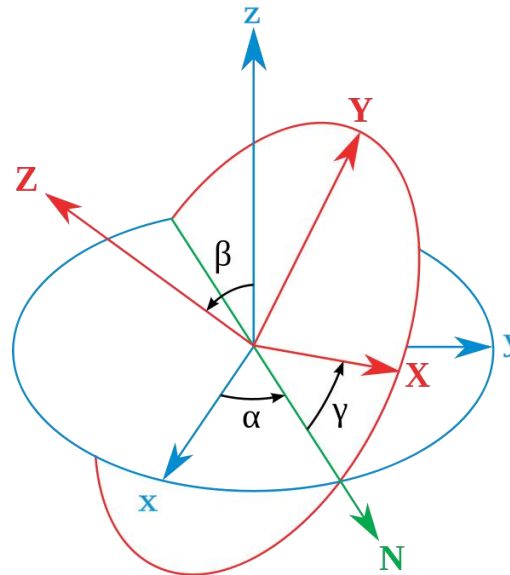
Training Techniques

- Progressive Strided Ray Sampling (from large stride to small stride)

$$\mathcal{P}(u, v, s) = \{(u + sx, v + sy) \mid x, y \in \{0, \dots, K\}\}$$

- Progressive Gaussian Pose Sampling (from small distortion to large distortion)

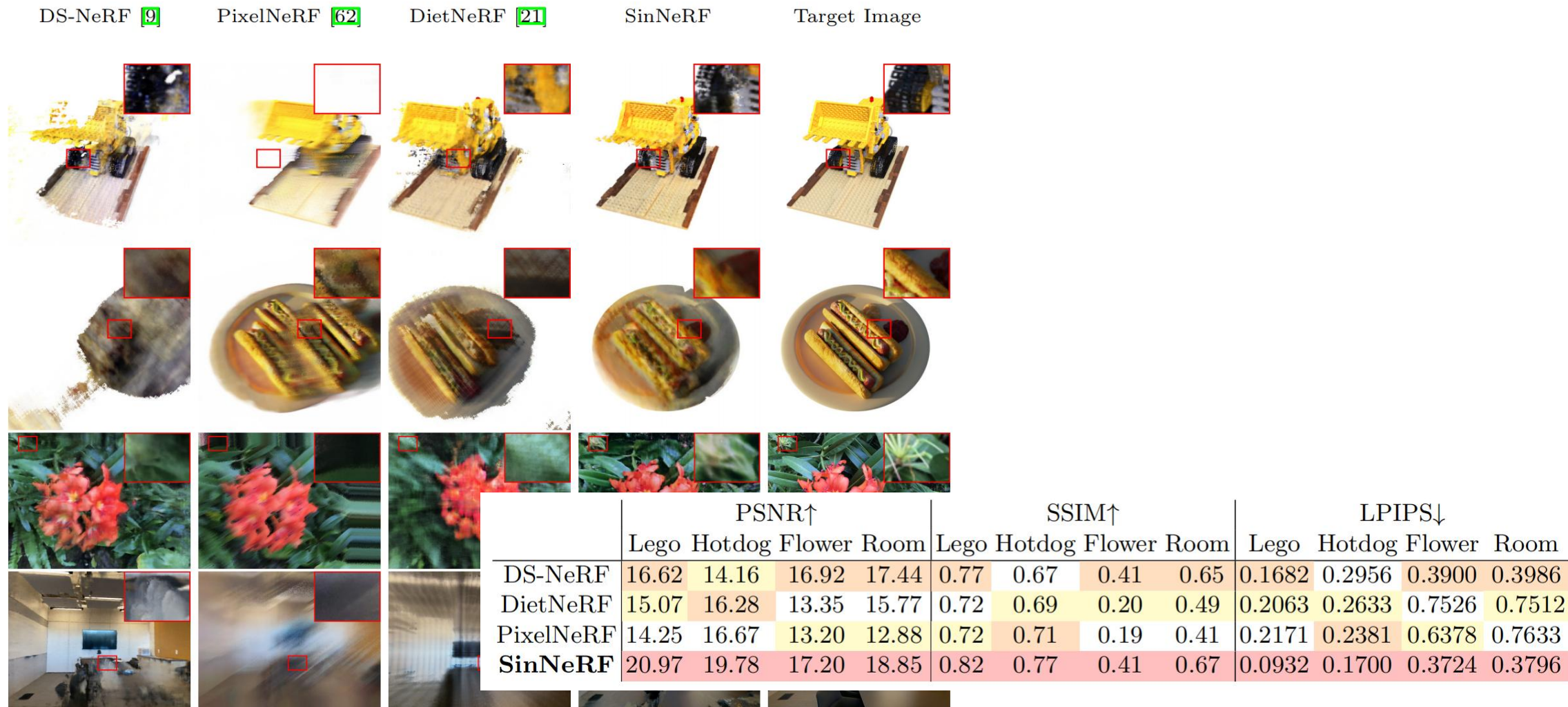
$$(\alpha, \beta, \gamma) \sim \mathcal{N}(0, \omega^2)$$



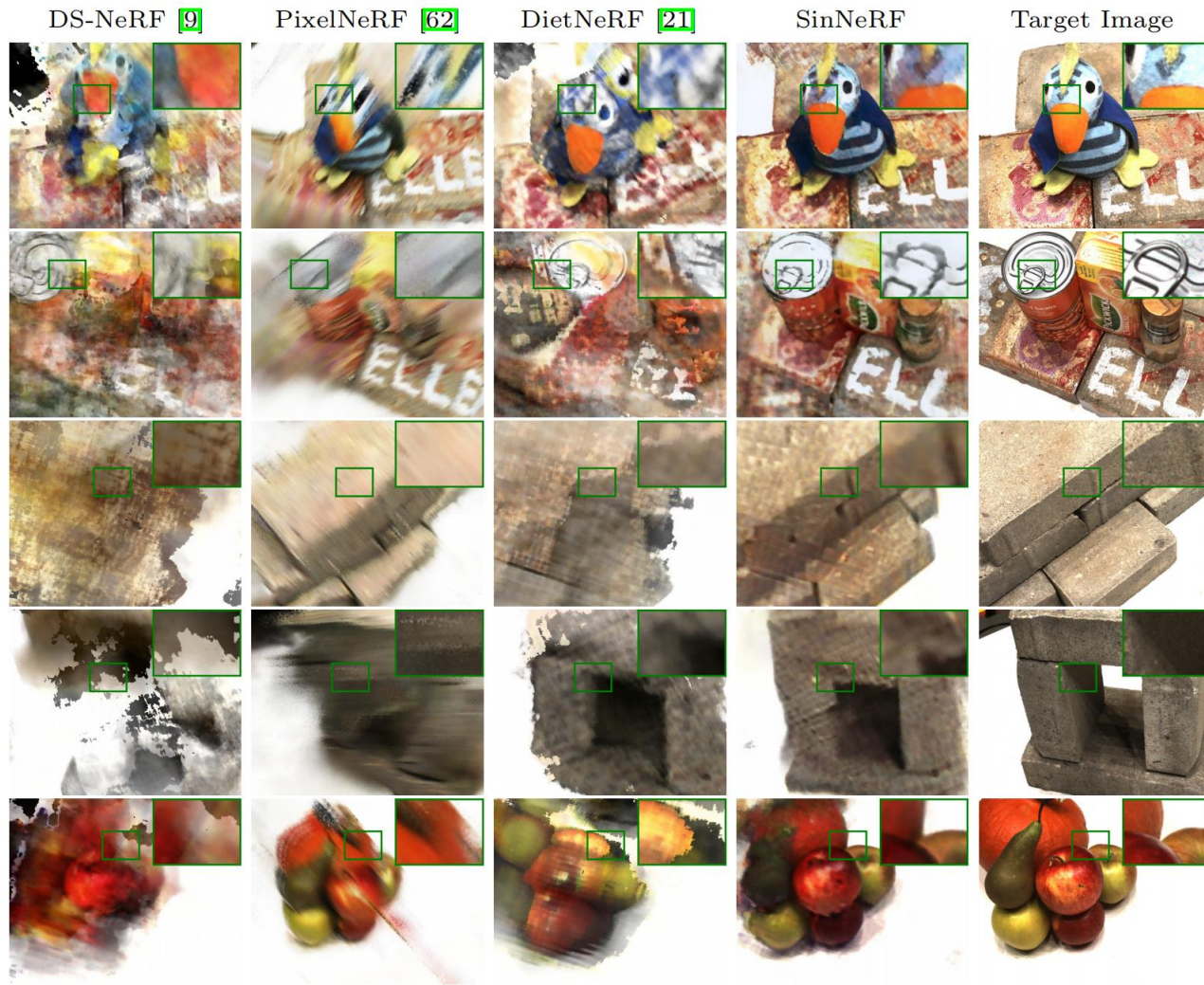
OUTLINE

- Authorship
- Background
- Method
- Experiments
- Conclusion

EXPERIMENTS



EXPERIMENTS



EXPERIMENTS

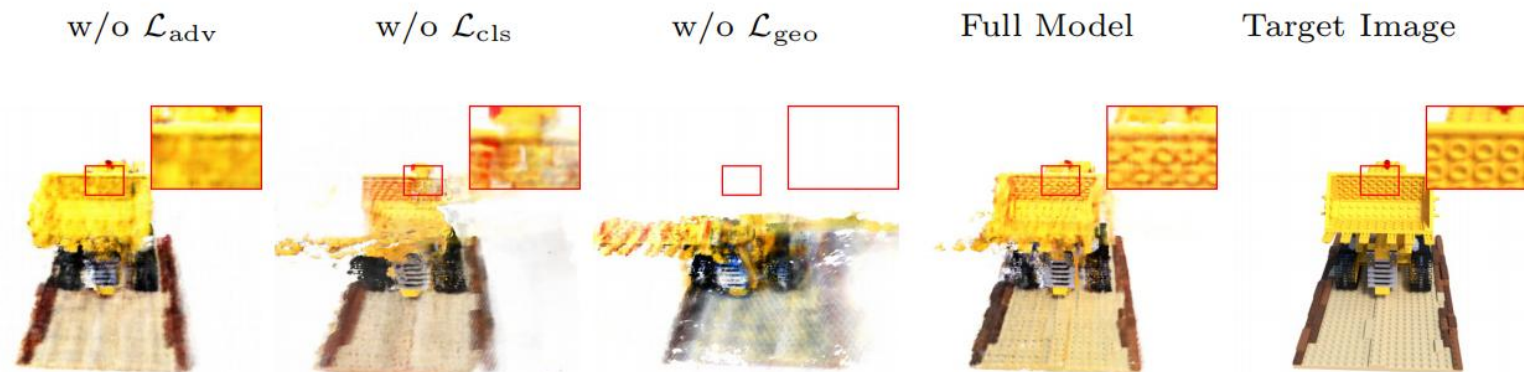


Fig. 5: Novel view synthesis from different variants of our proposed model.

Methods	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
w/o \mathcal{L}_{geo}	16.11 (-4.86)	0.74 (-0.08)	0.1919 (+0.0987)
w/o \mathcal{L}_{cls}	18.20 (-2.77)	0.76 (-0.06)	0.1348 (+0.0146)
w/o \mathcal{L}_{adv}	20.20 (-0.77)	0.79 (-0.03)	0.1306 (+0.0294)
Full Model	20.97	0.82	0.0932

OUTLINE

- Authorship
- Background
- Method
- Experiments
- Conclusion

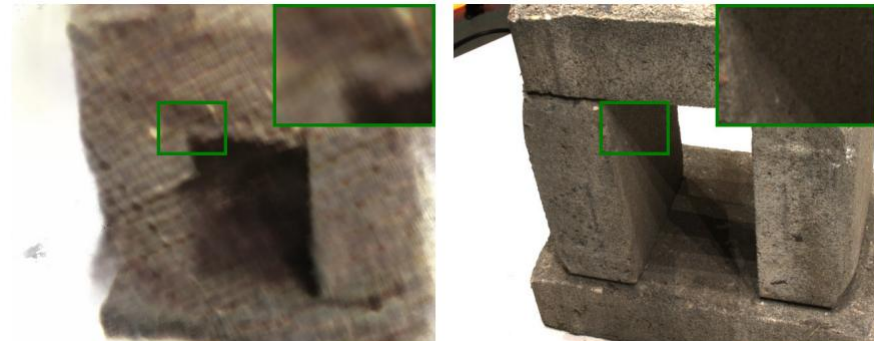
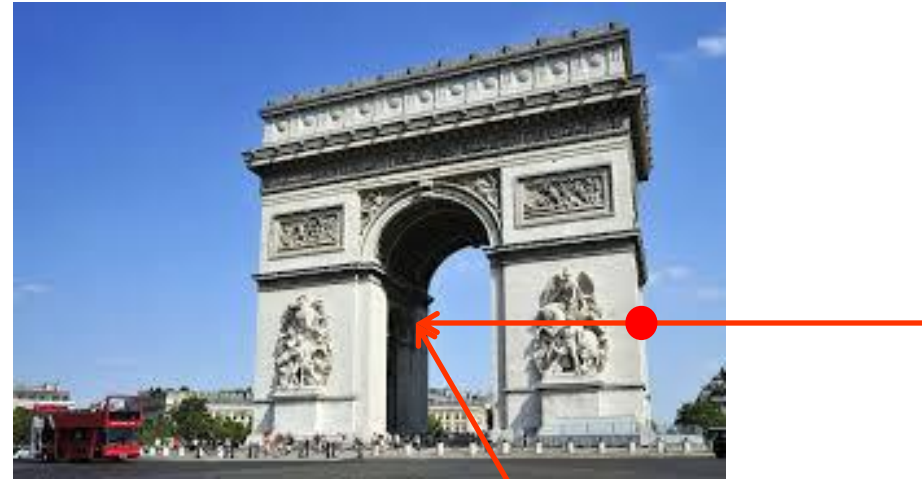
CONCLUSION

Pros

- The authors proposed SinNeRF, a novel view synthesis approach which only requires one calibrated view with depth
- SoTA on multiple benchmarks

Cons

- Test time optimization (same as original NeRF), i.e., one model per scene
- Simple scene
- Very complex design
- Cannot handle occlusion



Thanks for listening!