# Occlusion-Free Scene Recovery via Neural Radiance Fields

Chengxuan Zhu[1,2]   Renjie Wan[3]   Yunkai Tang[1,2]   Boxin Shi[1,2]

[1]National Key Laboratory for Multimedia Information Processing, School of Computer Science, Peking University
[2]National Engineering Research Center of Visual Technology, School of Computer Science, Peking University
[3]Department of Computer Science, Hong Kong Baptist University

# Occlusion Removal

Background (desired)

Scribbles, fences, waterdrops…
(undesired foreground occlusion)

occluded

- Opaque occlusions can block useful information from reaching the camera
  - Structure-from-motion may fail
  - Occlusions cause trouble for downstream vision tasks

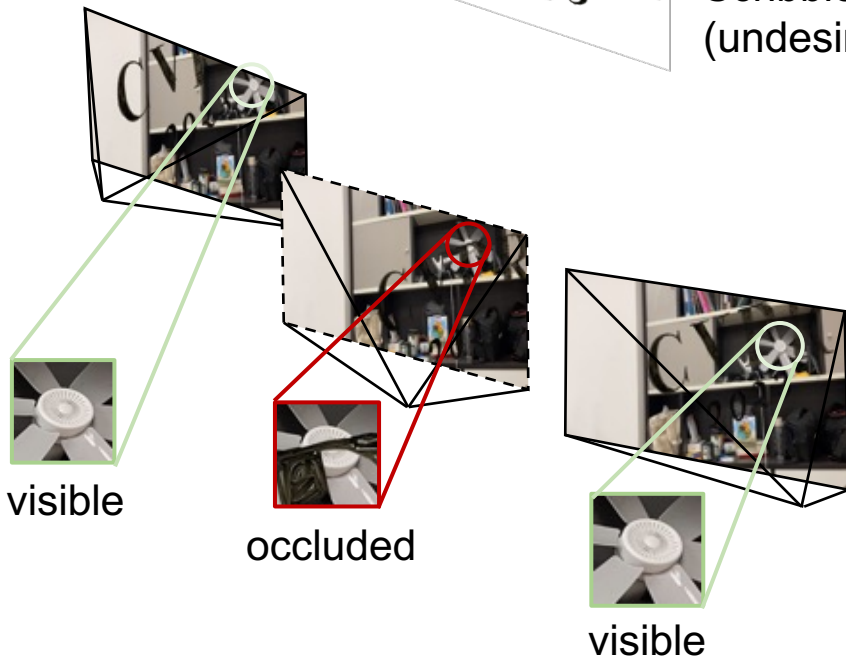Occlusion-free scene representation?

2

# Occlusion Removal

Background (desired)

Scribbles, fences, waterdrops…
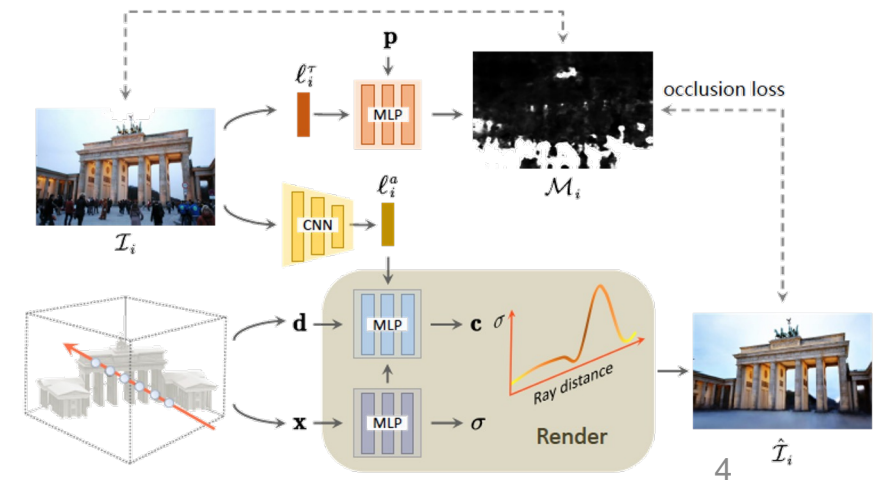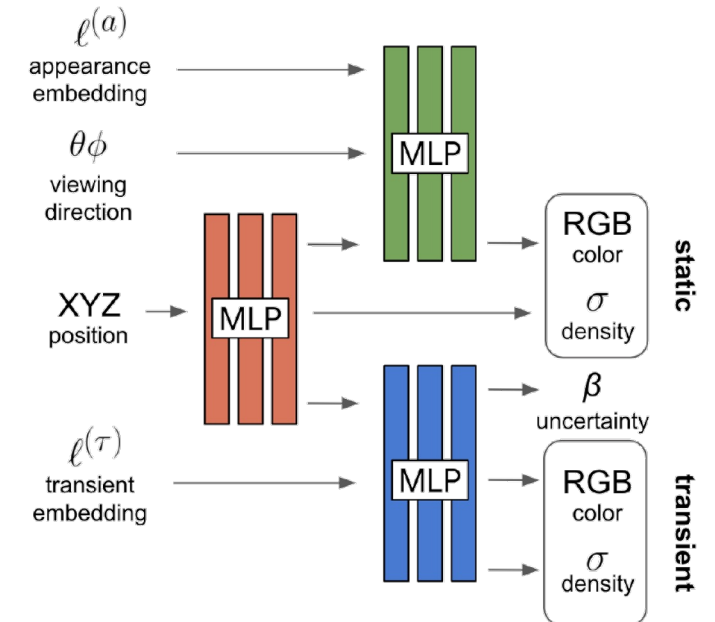(undesired foreground occlusion)

visible

occluded

visible

- Opaque occlusions can block useful information from reaching the camera
  - Structure-from-motion may fail
  - Occlusions cause trouble for downstream vision tasks

- Good news: We have multiple views!
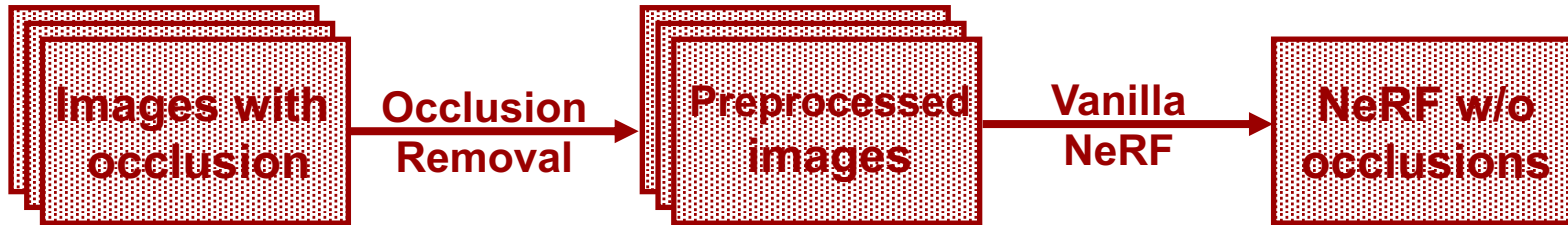
Occlusio Occlusion-free NeRF? tion?

# Related Methods

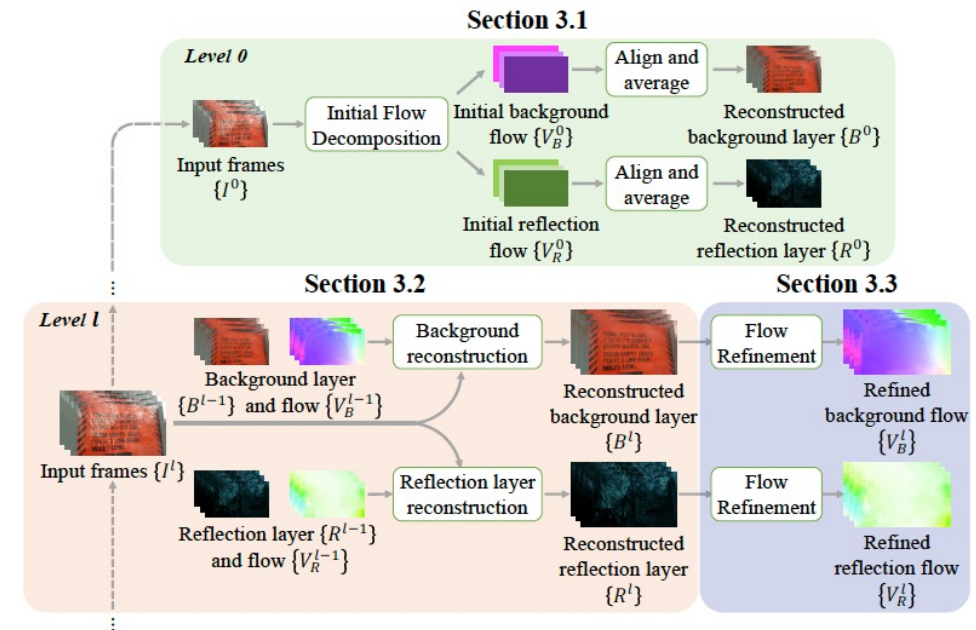**Images with occlusion** → Excluding uncertain parts → **NeRF w/o occlusions**

- Modeling a static scene and transient objects respectively (NeRF-W) [Martin-Brualla, CVPR'21]

- Exclude transient objects with a visible possibility map (Ha-NeRF) [Chen, CVPR'22]

- Problem: Reliant on pre-computed camera poses, and may not work with **static occlusions**



4

# Related Methods

Images with occlusion → **Occlusion Removal** → Preprocessed images → **Vanilla NeRF** → NeRF w/o occlusions

- Exploiting motion parallax estimated from multiple shots to separate different layers [Liu, CVPR'20]

- Inpainting based on single image is omitted due to its ill-posed nature and poor performance

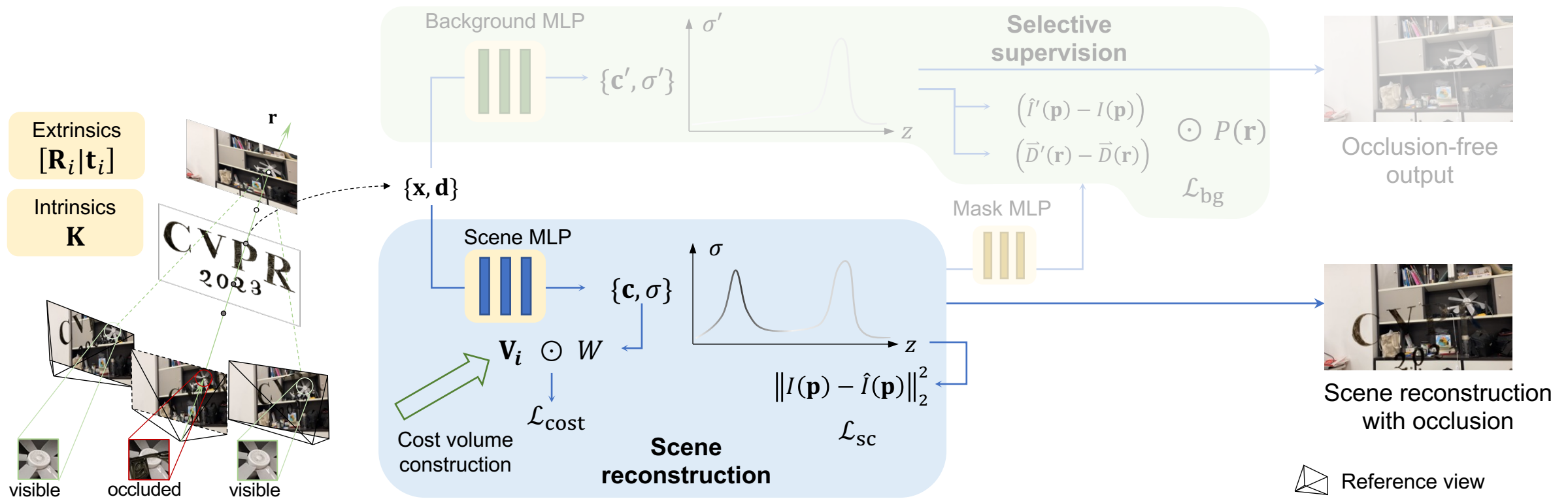- Problem: **Not 3D consistent** (NeRF may fail), and requires the input views to be close for flow estimation

# Occlusion Removal

- Joint optimization of **pose refinement** and **scene reconstruction** by effective multi-view feature fusion

- **Self-supervised** occlusion detection and occlusion-free scene recovery via NeRF

- Opaque occlusions can block useful information from reaching the camera
  - Structure-from-motion may fail
  - Occlusions cause trouble for downstream vision tasks
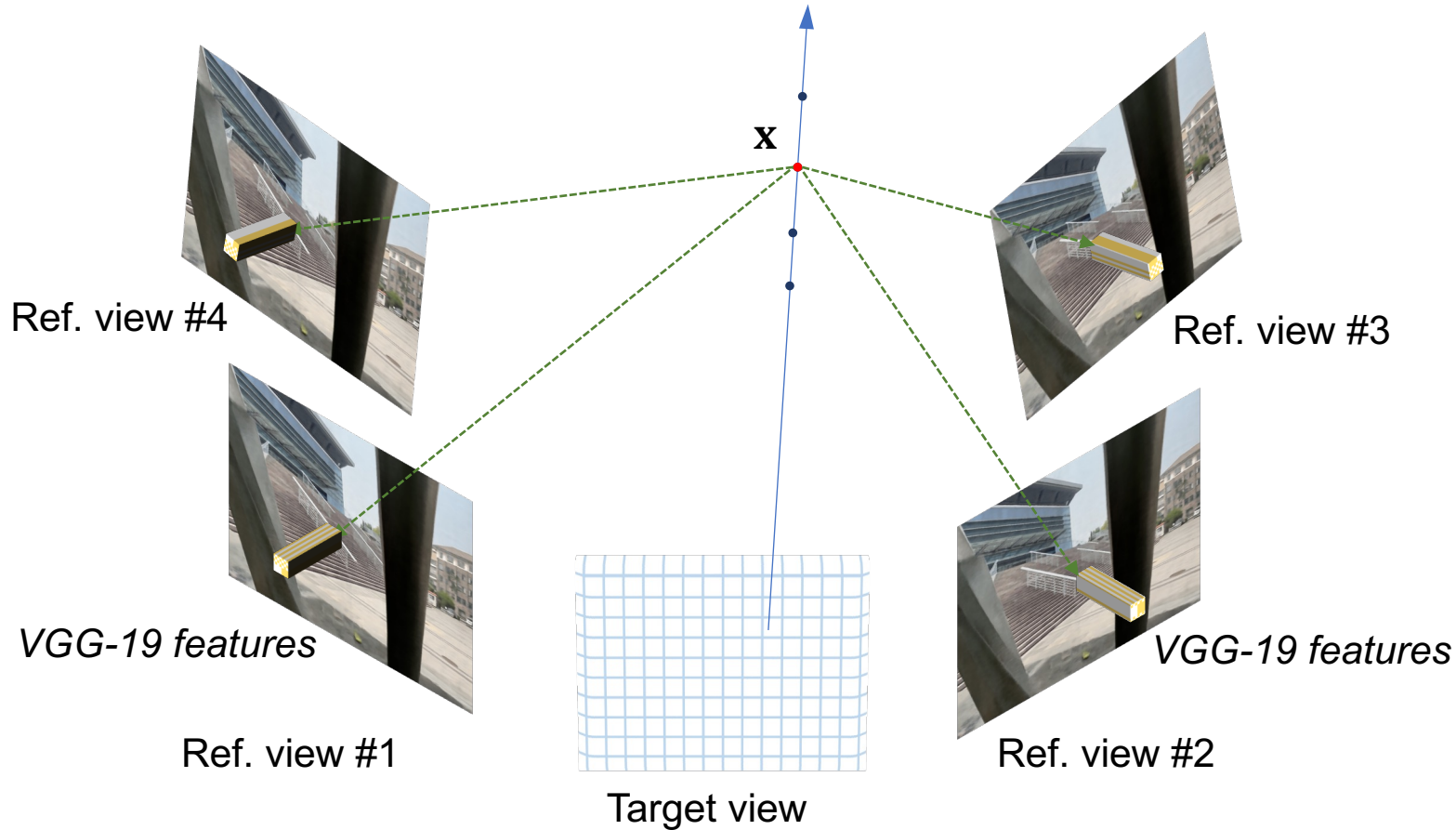
- Good news: We have multiple views!

Occlusion-free NeRF?

# Our Pipeline

Extrinsics $[\mathbf{R}_i | \mathbf{t}_i]$

Intrinsics $\mathbf{K}$

$\mathbf{r}$

$\{\mathbf{x}, \mathbf{d}\}$

Background MLP

$\{\mathbf{c}', \sigma'\}$

$\sigma'$

$z$

**Selective supervision**

$\left(\hat{I}'(\mathbf{p}) - I(\mathbf{p})\right)$

$\left(\vec{D}'(\mathbf{r}) - \vec{D}(\mathbf{r})\right)$

$\odot \; P(\mathbf{r})$

$\mathcal{L}_{\mathrm{bg}}$

Occlusion-free output

Scene MLP

$\{\mathbf{c}, \sigma\}$

$\sigma$

$z$

Mask MLP

$\mathbf{V}_i \odot W$

$\mathcal{L}_{\mathrm{cost}}$

$\left\| I(\mathbf{p}) - \hat{I}(\mathbf{p}) \right\|_2^2$

$\mathcal{L}_{\mathrm{sc}}$

Cost volume construction

**Scene reconstruction**

Scene reconstruction with occlusion

visible    occluded    visible

Reference view

Target view

Trainable components

$\mathbf{V}_i$ Cost volume

$W$ Weight

# Joint Optimization
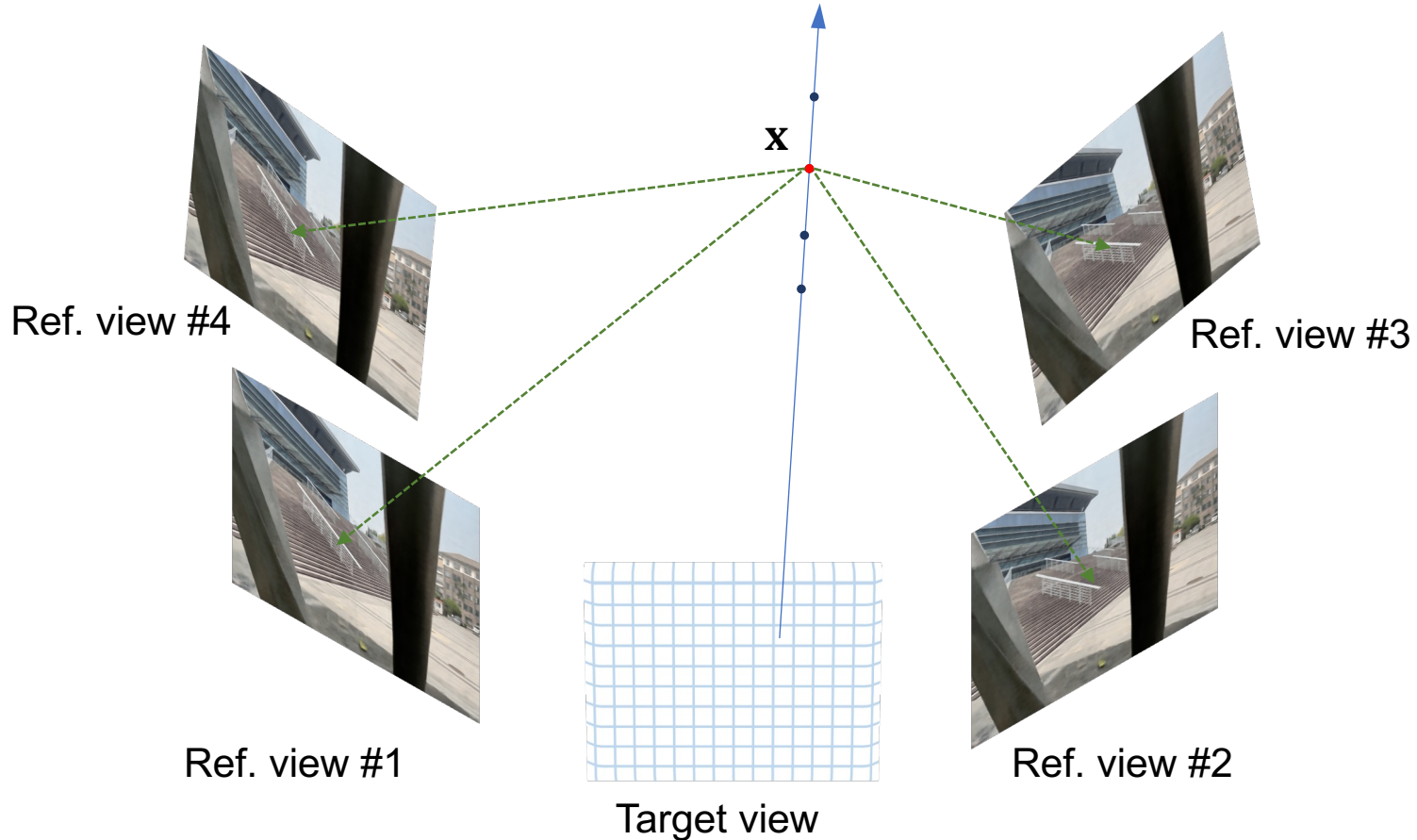
- Guided by **cost volume**



Ref. view #4

Ref. view #3

**x**

VGG-19 features

VGG-19 features

Ref. view #1

Target view

Ref. view #2

# Joint Optimization

- Guided by **cost volume**



Ref. view #4

Ref. view #1

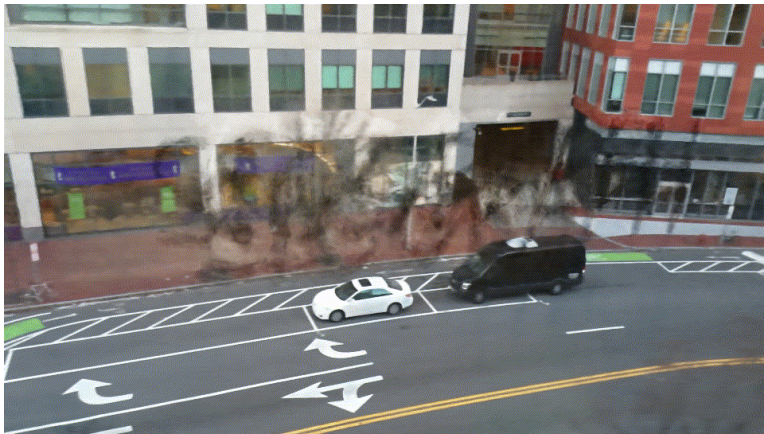Target view

Ref. view #3

Ref. view #2

**Intuition about cost volume**

The larger Variance ( ),
the less probable a visible point is
located at **x**

We use a "scene MLP" to jointly
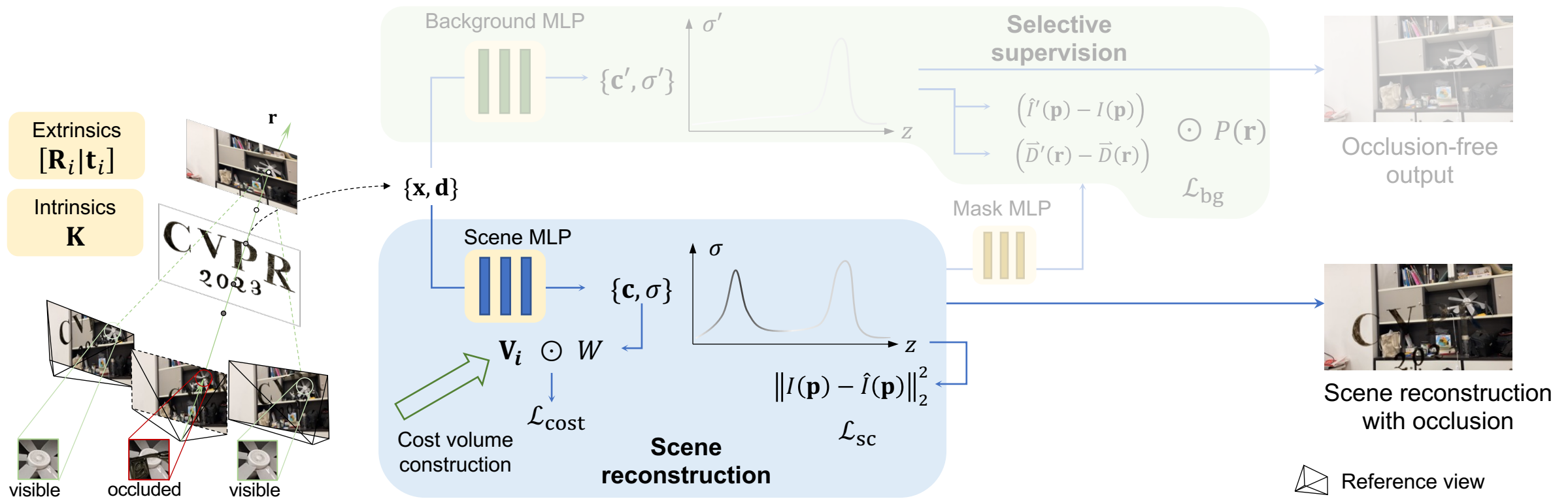reconstruct the scene and optimize
camera parameters

9

# Joint Optimization

NeRF

NeRF--
[Wang, ArXiv'21]

Ours, scene MLP

# Our Pipeline

Extrinsics $[\mathbf{R}_i|\mathbf{t}_i]$

Intrinsics $\mathbf{K}$

$\{\mathbf{x}, \mathbf{d}\}$

**Background MLP**

$\{\mathbf{c}', \sigma'\}$

$\sigma'$

$z$

**Selective supervision**

$\left(\hat{I}'(\mathbf{p}) - I(\mathbf{p})\right)$

$\left(\vec{D}'(\mathbf{r}) - \vec{D}(\mathbf{r})\right)$

$\odot \; P(\mathbf{r})$

$\mathcal{L}_{\mathrm{bg}}$

Occlusion-free output

Mask MLP

Scene MLP

$\{\mathbf{c}, \sigma\}$

$\sigma$

$z$

$\mathbf{V}_i \odot W$

$\mathcal{L}_{\mathrm{cost}}$

Cost volume construction

$\left\| I(\mathbf{p}) - \hat{I}(\mathbf{p}) \right\|_2^2$

$\mathcal{L}_{\mathrm{sc}}$

**Scene reconstruction**

Scene reconstruction with occlusion

visible    occluded    visible

Reference view

Target view

Trainable components

$\mathbf{V}_i$  Cost volume

$W$  Weight

11

# Our Pipeline

Extrinsics
$[\mathbf{R}_i | \mathbf{t}_i]$

Intrinsics
$\mathbf{K}$

Background MLP

$\{\mathbf{c}', \sigma'\}$

$\sigma'$

$z$

**Selective supervision**

$\left(\hat{I}'(\mathbf{p}) - I(\mathbf{p})\right)$

$\left(\vec{D}'(\mathbf{r}) - \vec{D}(\mathbf{r})\right)$

$\odot P(\mathbf{r})$

$\mathcal{L}_{\mathrm{bg}}$

Occlusion-free
output

$\{\mathbf{x}, \mathbf{d}\}$

Mask MLP

Scene MLP

$\{\mathbf{c}, \sigma\}$

$\sigma$

$z$

$\mathbf{V}_i \odot W$

$\mathcal{L}_{\mathrm{cost}}$

Cost volume
construction

**Scene
reconstruction**

$\|I(\mathbf{p}) - \hat{I}(\mathbf{p})\|_2^2$

$\mathcal{L}_{\mathrm{sc}}$

Scene reconstruction
with occlusion

visible     occluded     visible

Reference view

Target view

Trainable components

$\mathbf{V}_i$   Cost volume

$W$   Weight

# Selective Supervision

- Guided by **bidirectional depth inconsistency**

Density

Depth

Infinite far

The ray terminates here

Reverse ray terminates here

When bidirectional depth inconsistency is *small*, there is no occlusion.

# Selective Supervision

- Guided by **bidirectional depth inconsistency**



Density

Depth

Infinite far

The ray terminates here

Reverse ray terminates here

When bidirectional depth inconsistency is *large*, there is occlusion.

# Selective Supervision

- Guided by **bidirectional depth inconsistency**

Density

Depth

Mask MLP

→ Probability of the ray not hitting occlusion

(supervised by bidirectional depth inconsistency)

# Selective Supervision

- Supervise the background MLP only where there is no occlusion



Input view      Mask MLP output      Ours, background MLP

# Our pipeline

Background MLP

$\sigma'$

Selective supervision

$\{\mathbf{c}', \sigma'\}$

$z$

$\left(\hat{I}'(\mathbf{p}) - I(\mathbf{p})\right)$

$\left(\vec{D}'(\mathbf{r}) - \vec{D}(\mathbf{r})\right)$

$\odot \, P(\mathbf{r})$

$\mathcal{L}_{\text{bg}}$

Extrinsics
$[\mathbf{R}_i | \mathbf{t}_i]$

Intrinsics
$\mathbf{K}$

$\mathbf{r}$

$\{\mathbf{x}, \mathbf{d}\}$

Mask MLP

Occlusion-free output

Scene MLP

$\sigma$

$\{\mathbf{c}, \sigma\}$

$\mathbf{V}_i \odot W$

$z$

$\mathcal{L}_{\text{cost}}$

$\|I(\mathbf{p}) - \hat{I}(\mathbf{p})\|_2^2$

$\mathcal{L}_{\text{sc}}$

Cost volume construction

**Scene reconstruction**

visible    occluded    visible

Scene reconstruction with occlusion

Reference view

Target view

Trainable components

$\mathbf{V}_i$   Cost volume

$W$   Weight

17

# Dataset

- An evaluation dataset containing 10 different scenes
- Covering various types of occlusions



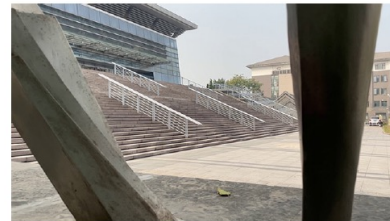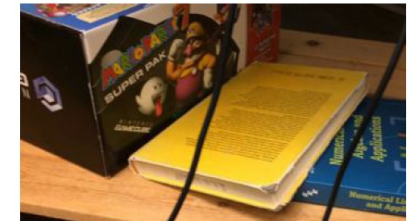(a) FENCE1    (b) FENCE2    (c) FENCE3    (d) SCRIBBLE1 (e) SCRIBBLE2

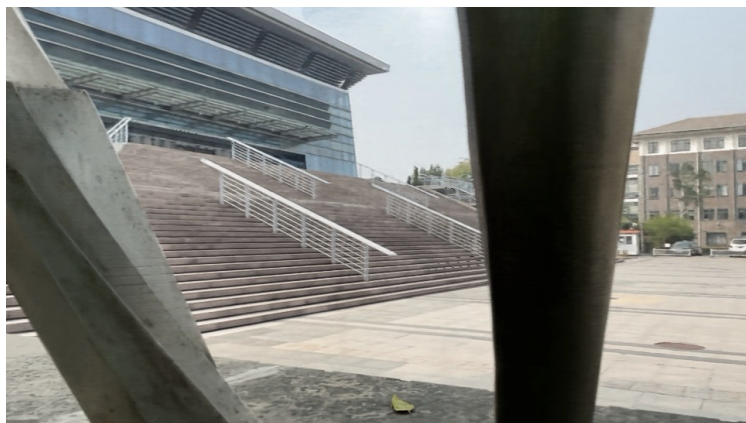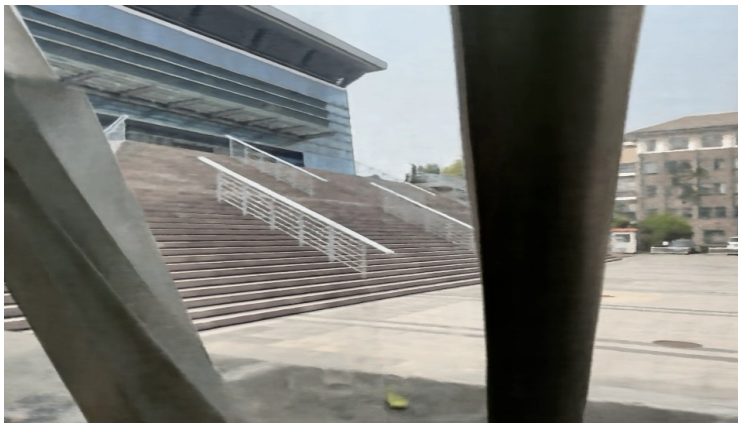(f) SCRIBBLE3 (g) RAINDROP    (h) STATUE    (i) WIRE1    (j) WIRE2
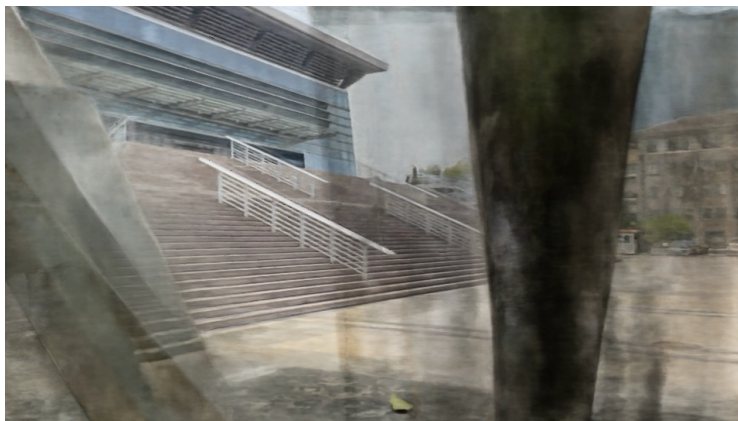
# Results

NeRF



Ours, Scene MLP



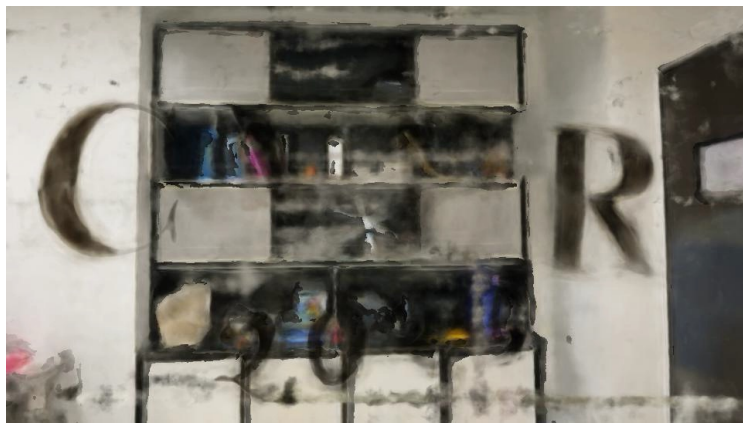Ours, Background MLP



Ha-NeRF



PWC-Net + NeRF



NeRF-W

# Results
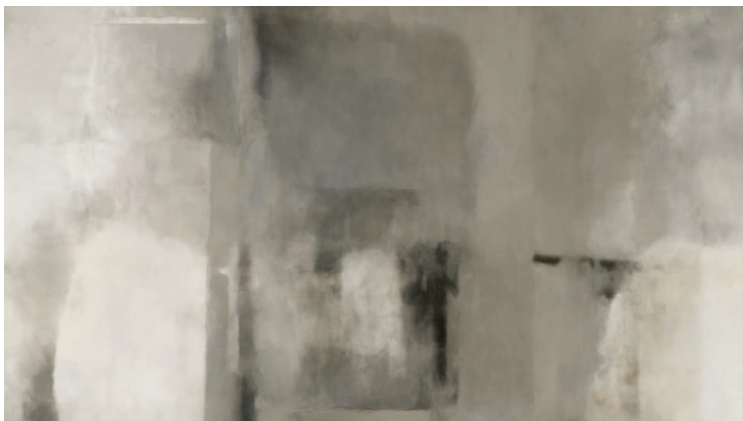
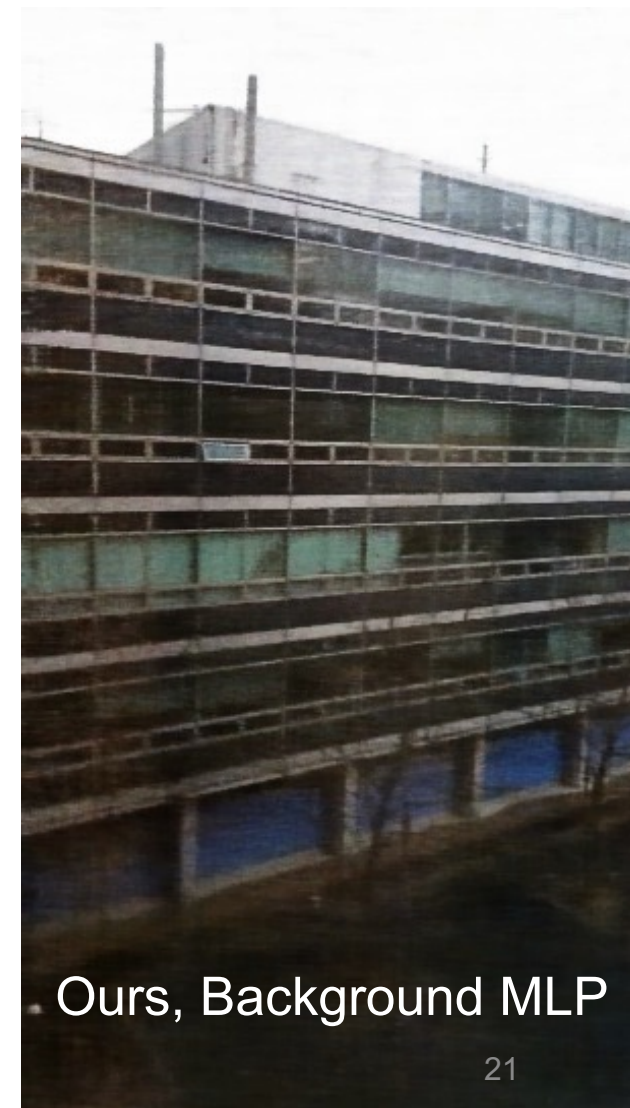NeRF



Ours, Scene MLP



Ours, Background MLP



Ha-NeRF



PWC-Net + NeRF



NeRF-W

# Results

* other baselines are omitted due to **COLMAP failure**



Input sample

Ours, Scene MLP

Ours, Background MLP

# Occlusion-Free Scene Recovery via Neural Radiance Fields

Chengxuan Zhu[1,2]   Renjie Wan[3]   Yunkai Tang[1,2]   Boxin Shi[1,2]

[1]National Key Laboratory for Multimedia Information Processing, School of Computer Science, Peking University
[2]National Engineering Research Center of Visual Technology, School of Computer Science, Peking University
[3]Department of Computer Science, Hong Kong Baptist University

22