# JacobiNeRF: NeRF Shaping with Mutual Information Gradients

Xiaomeng Xu[1,*], Yanchao Yang[2,3,*,†], Kaichun Mo[3,4], Boxiao Pan[3], Li Yi[1,5,6], Leonidas Guibas[3,7]

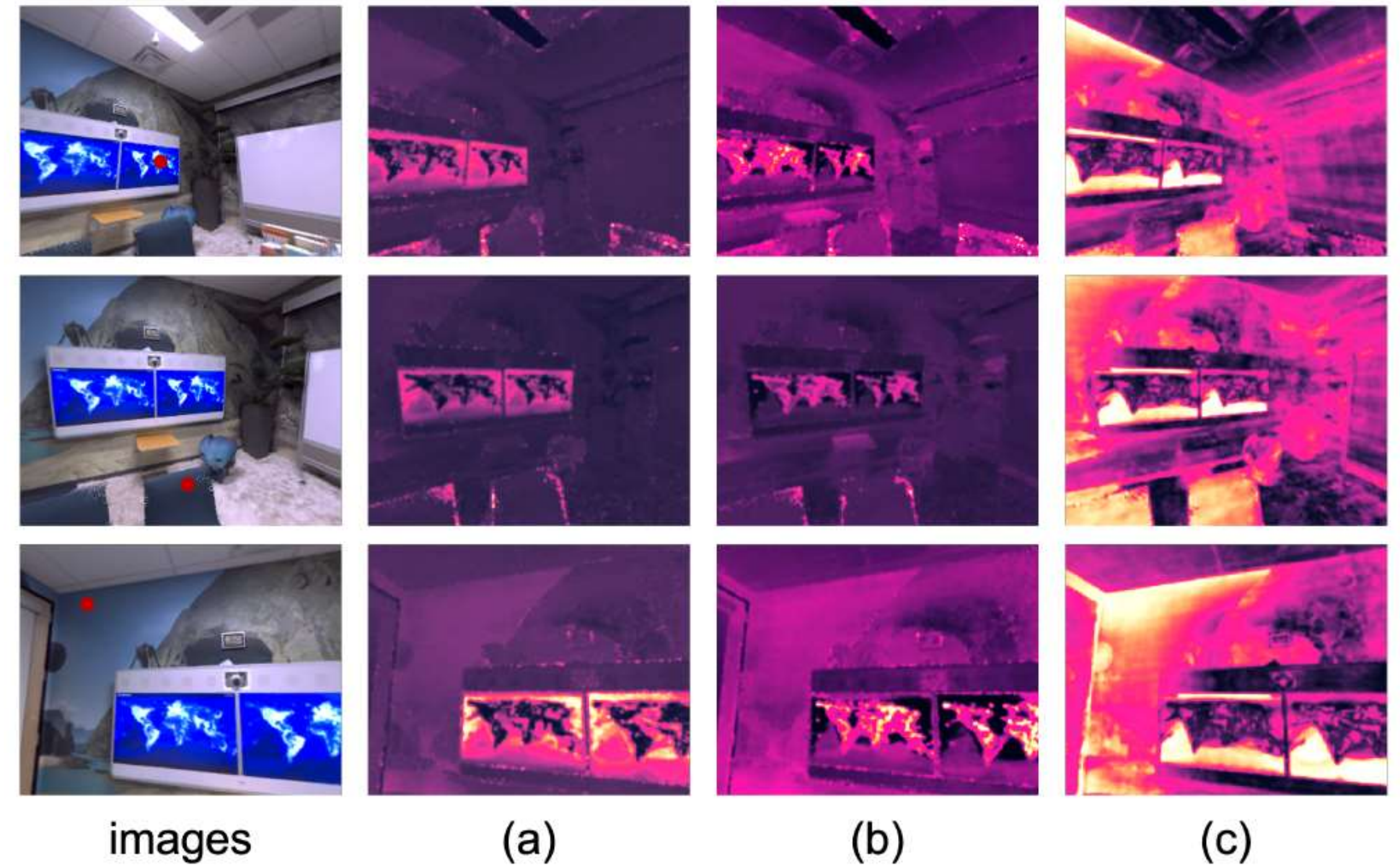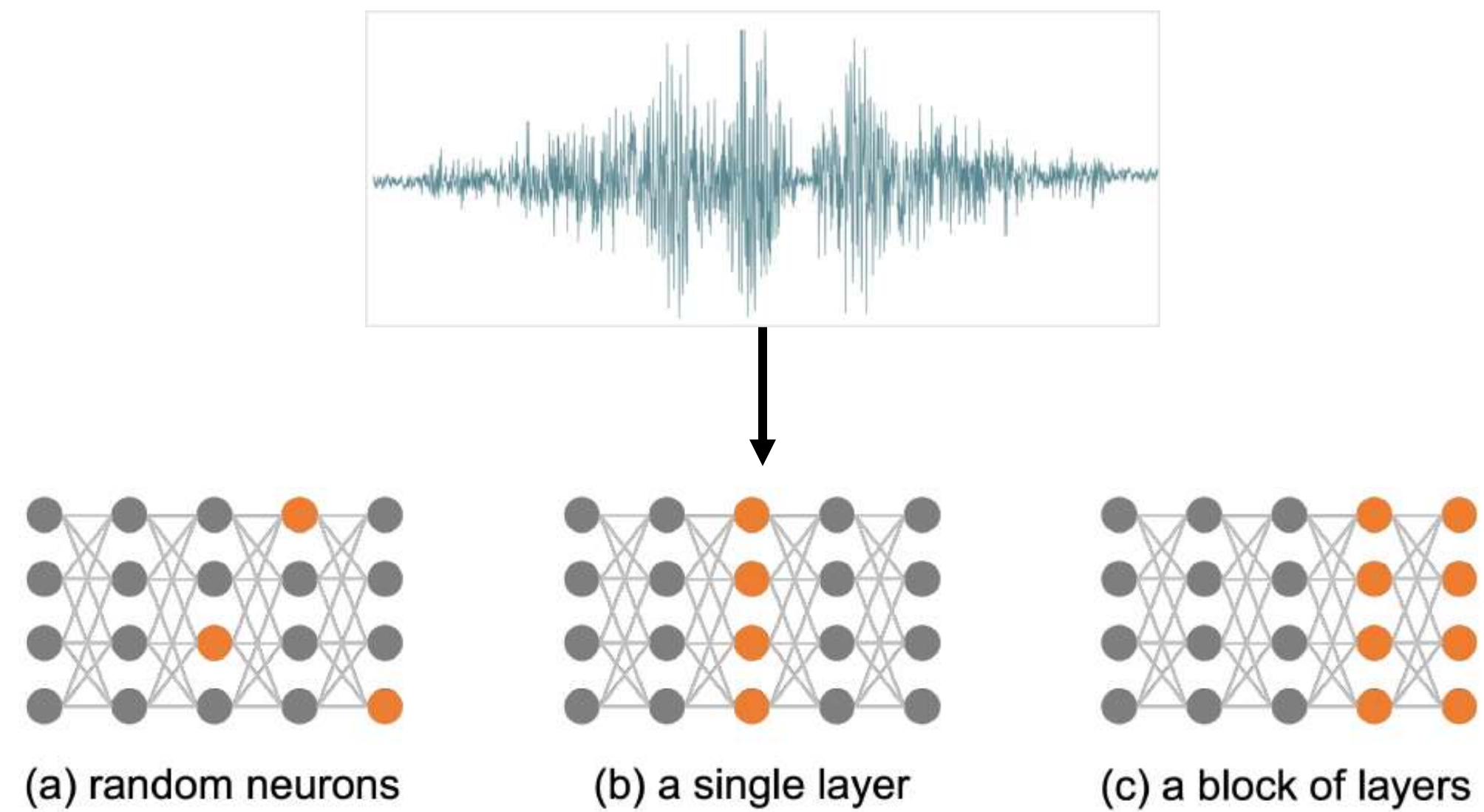[1]Tsinghua University, [2]The University of Hong Kong, [3]Stanford University, [4]NVIDIA Research, [5]Shanghai AI Laboratory, [6]Shanghai Qizhi Institute, [7]Google Research

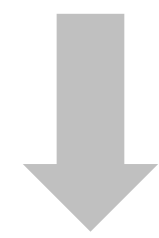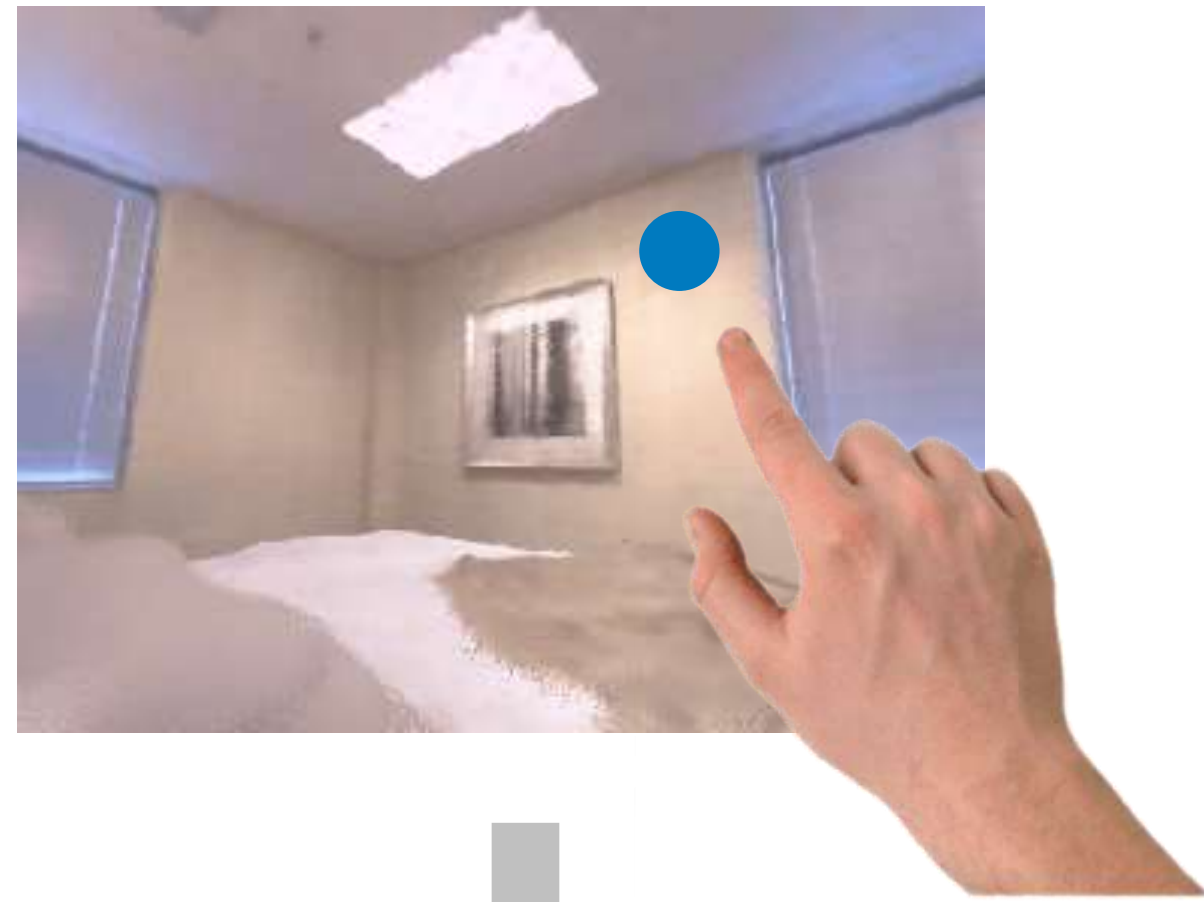*Equal Contributions, †Corresponding Author

Paper Tag: THU-AM-002

(a) random neurons    (b) a single layer    (c) a block of layers

images    (a)    (b)    (c)

"Jiggling" the neuron weights of the NeRF results in diffuse perturbations that are not semantically aligned

selecting

editing

label propagation

original state

effect after perturbation

MI-shaping  $\theta$

$\theta$ (NeRF)

perturbing p

$\frac{\partial I(p)}{\partial \theta}$

$\theta'$

Semantic Seg

Instance Seg

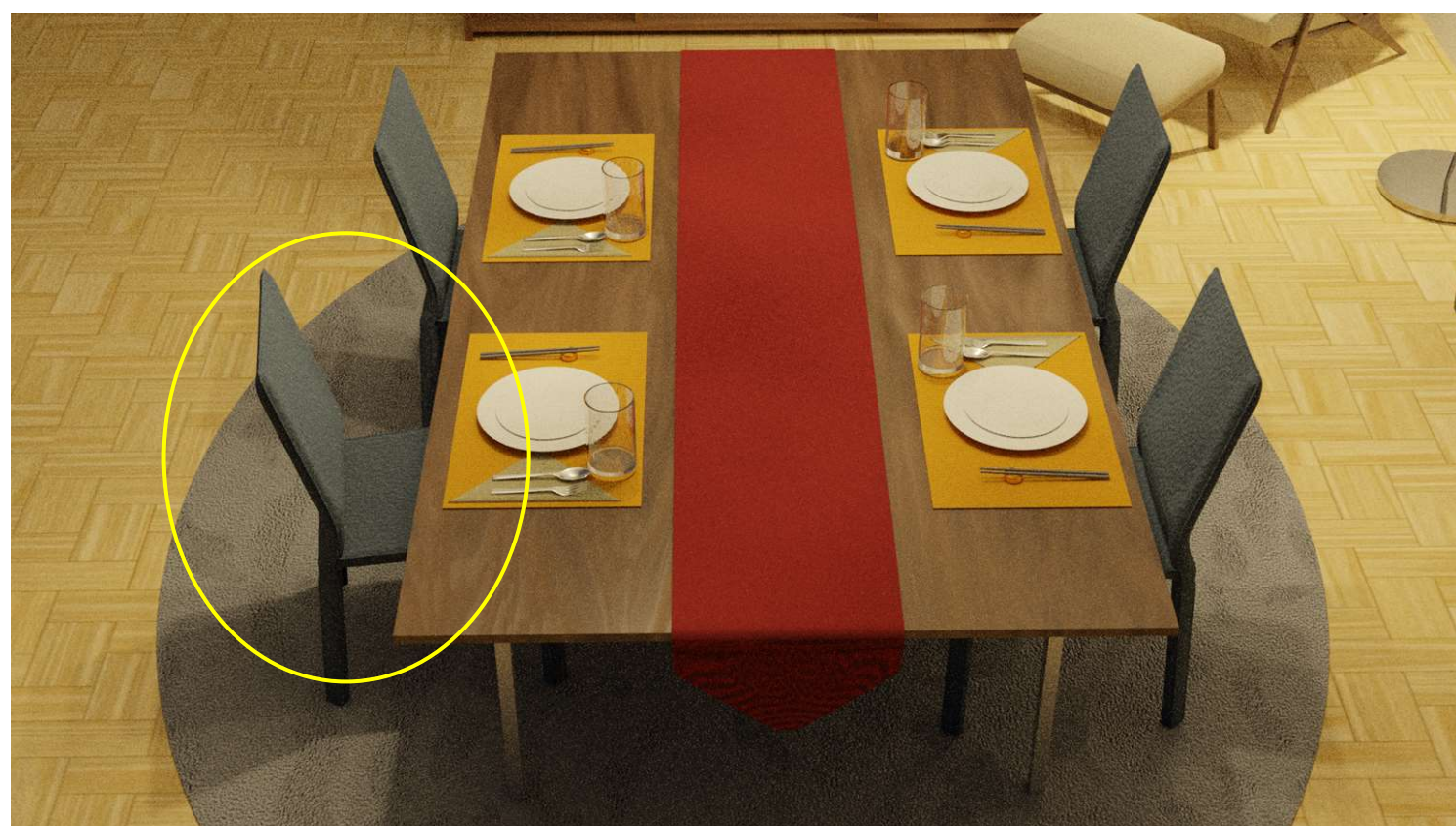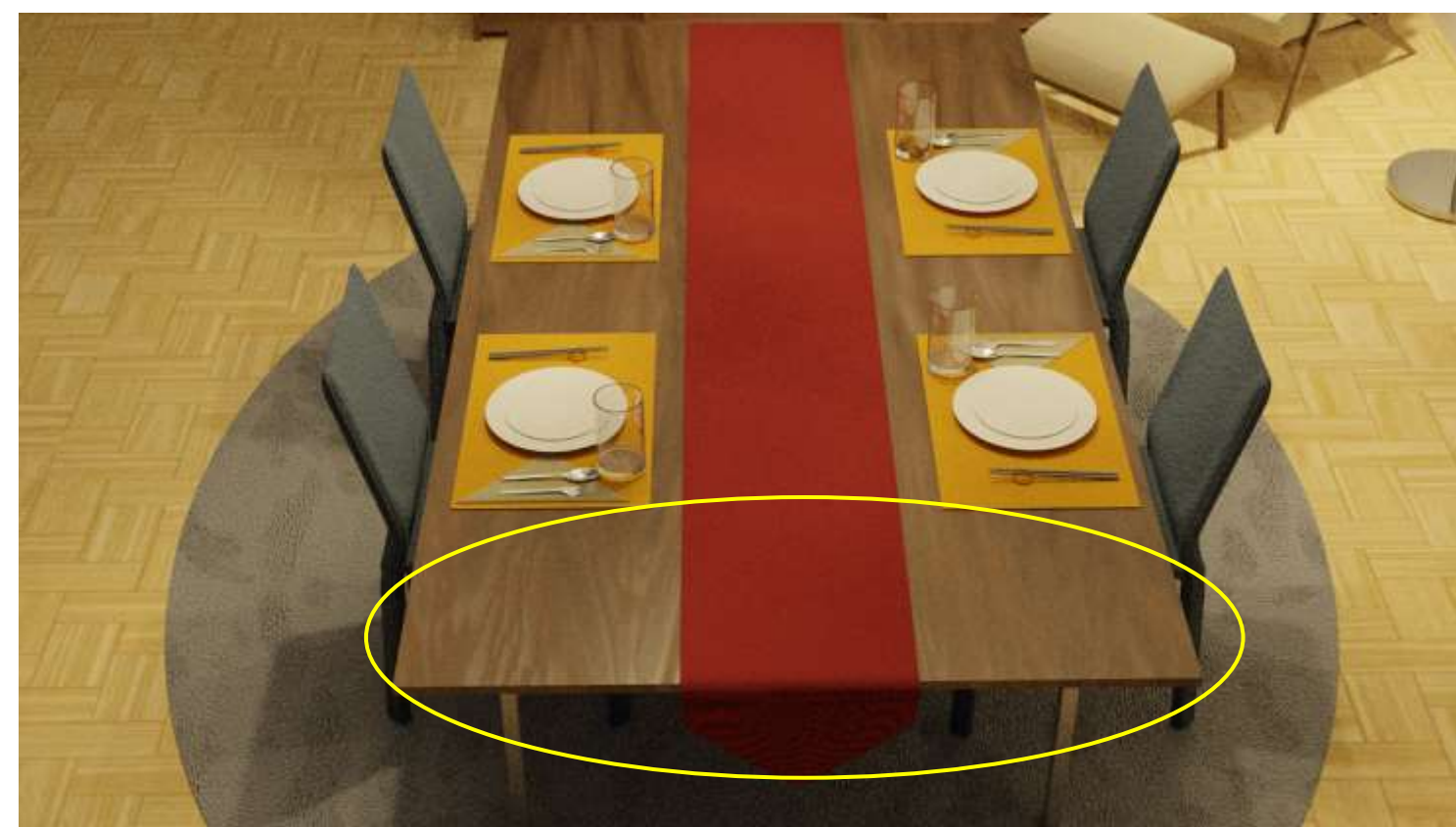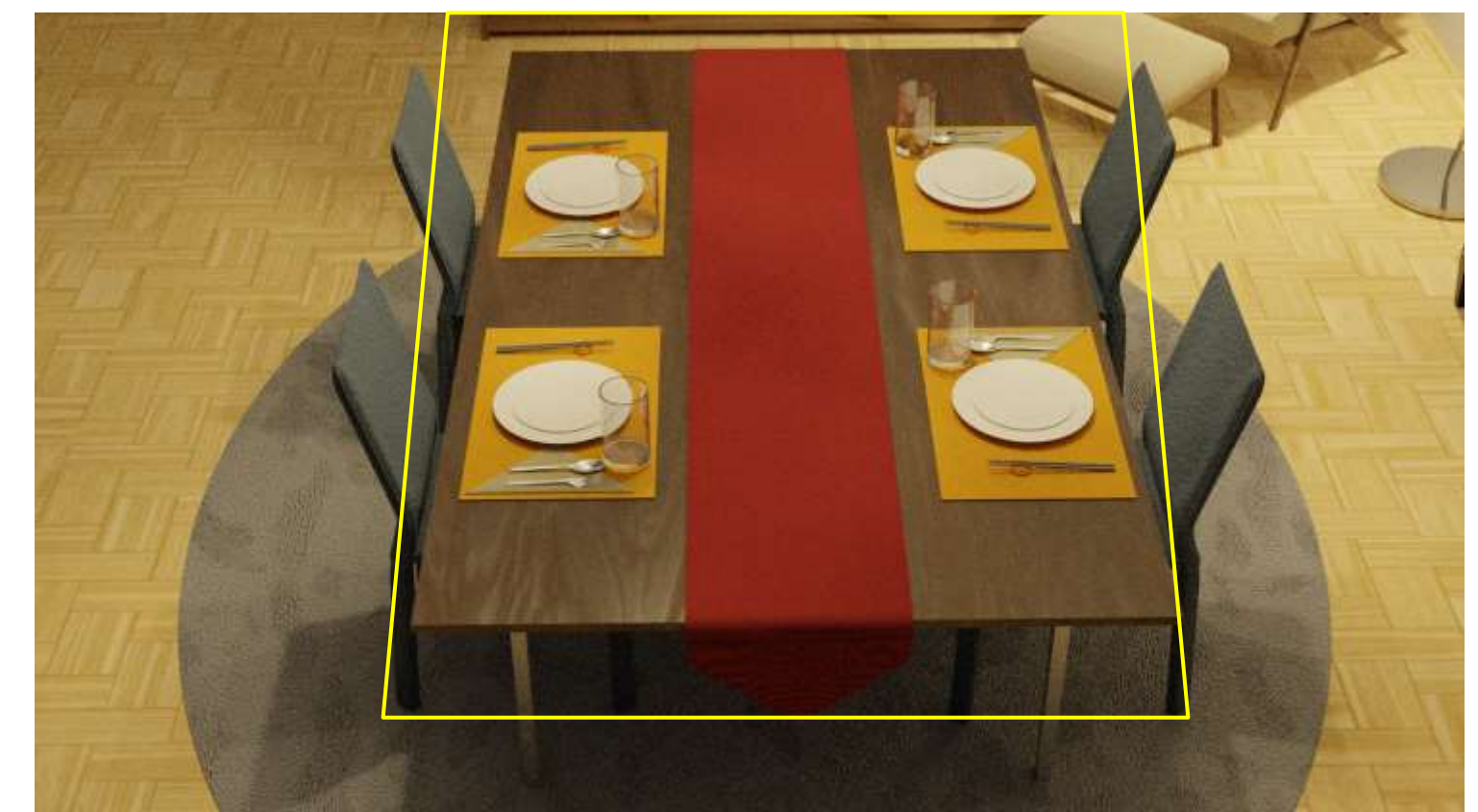**Example sparse annotation interaction**

# Semantic Structure of a Scene is Reflected in its Co-Variations



this chair has moved

the table became longer
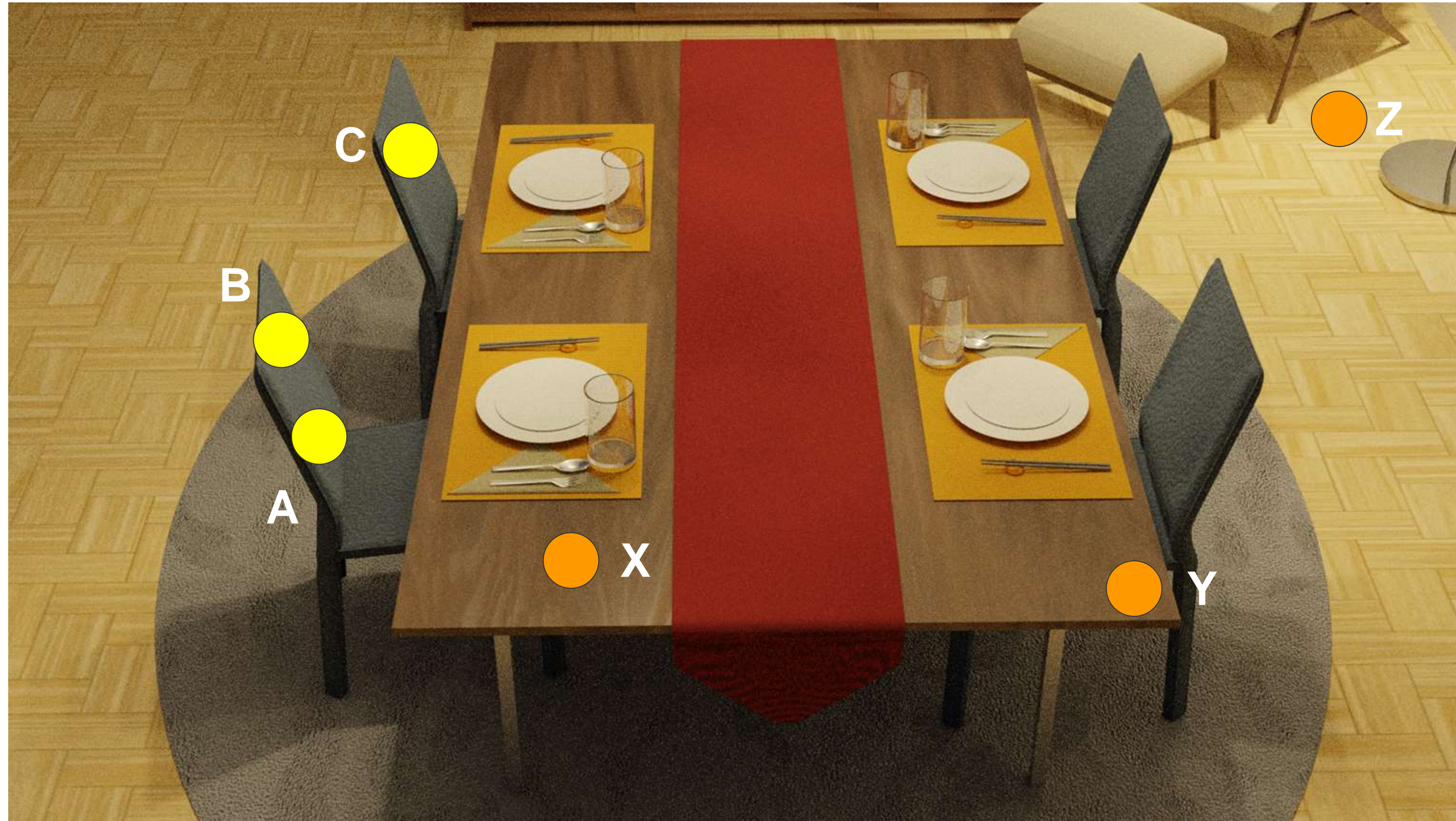
the table became darker

A is more correlated with B than with C $\qquad$ $\mathbb{I}(A, B) > \mathbb{I}(A, C)$

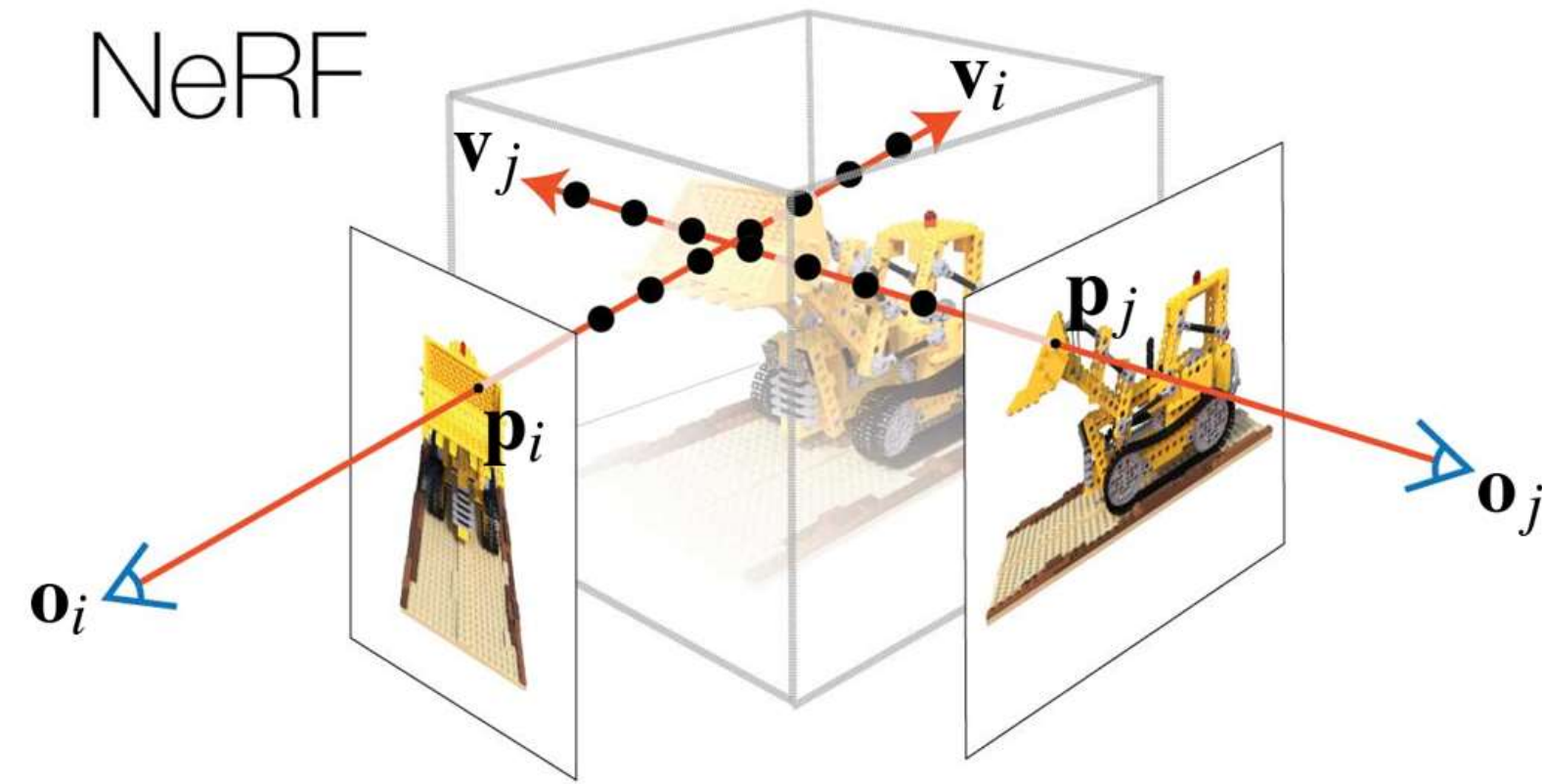X is more correlated with Y than with Z $\qquad$ $\mathbb{I}(X, Y) > \mathbb{I}(X, Z)$

Can we shape NeRFs to better reflect mutual correlations in the scene?

The scene tangent space



$\mathbb{I}(A,B) > \mathbb{I}(A,C)$
$\mathbb{I}(X,Y) > \mathbb{I}(X,Z)$

# Mutual Information via NeRF Gradients



$$I(\mathbf{p}_i) = \Phi(\mathbf{o}_i, \mathbf{v}_i; \boldsymbol{\theta})$$

$$I(\mathbf{p}_j) = \Phi(\mathbf{o}_j, \mathbf{v}_j; \boldsymbol{\theta})$$

$$\hat{I}(\mathbf{p}_i) = \Phi(\mathbf{o}_i, \mathbf{v}_i; \boldsymbol{\theta}^D + \mathbf{n})$$

$$\hat{I}(\mathbf{p}_j) = \Phi(\mathbf{o}_j, \mathbf{v}_j; \boldsymbol{\theta}^D + \mathbf{n})$$

Mutual information $\mathbb{I}$

$$\mathbb{I}(\hat{I}(\mathbf{p}_i), \hat{I}(\mathbf{p}_j)) \approx \left| \cos\left( \frac{\partial \Phi_i}{\partial \theta^D}, \frac{\partial \Phi_j}{\partial \theta^D} \right) \right|$$
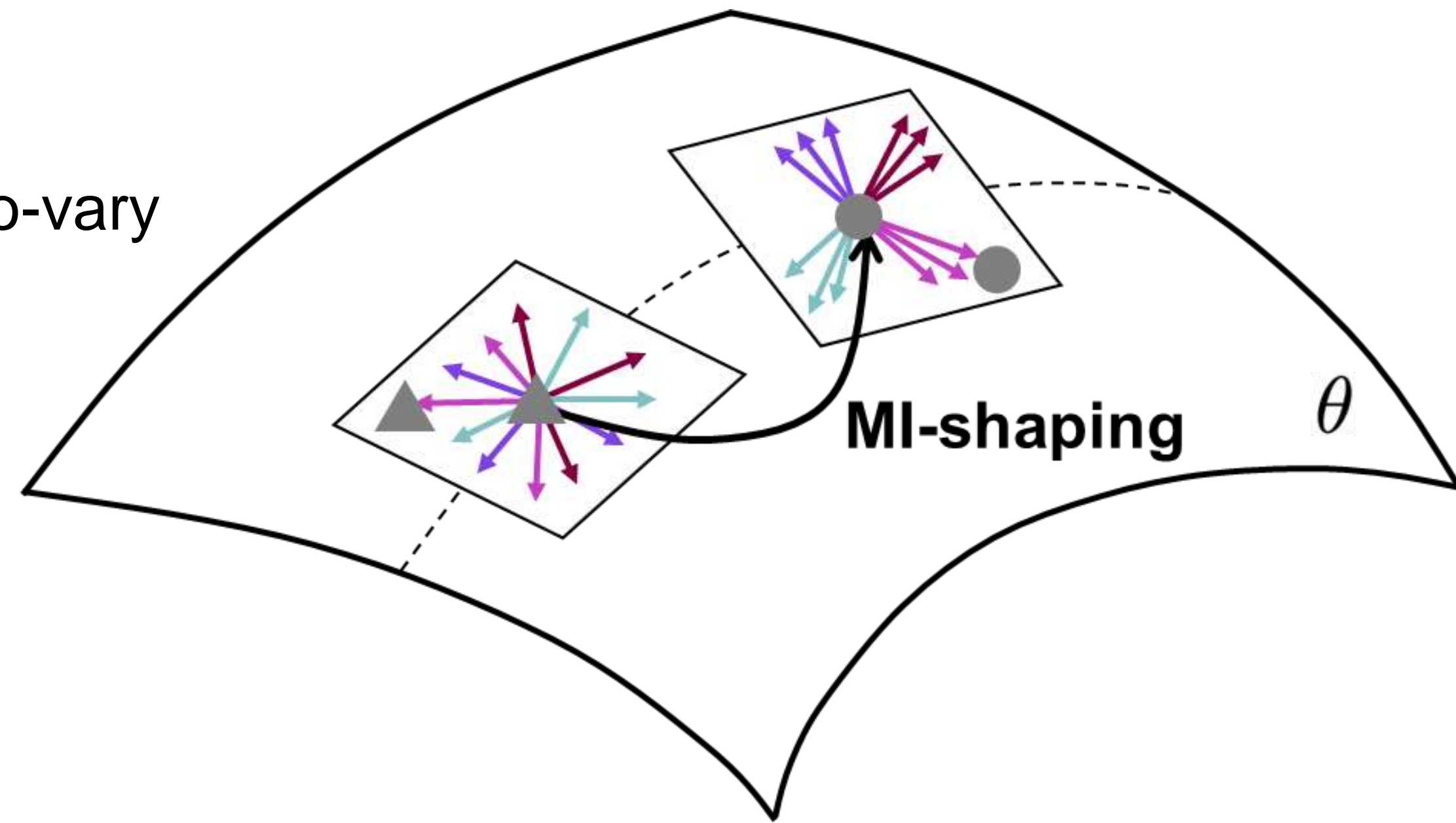
*Inter-pixel correlations are captured by cosine similarity of the NeRF Jacobians*

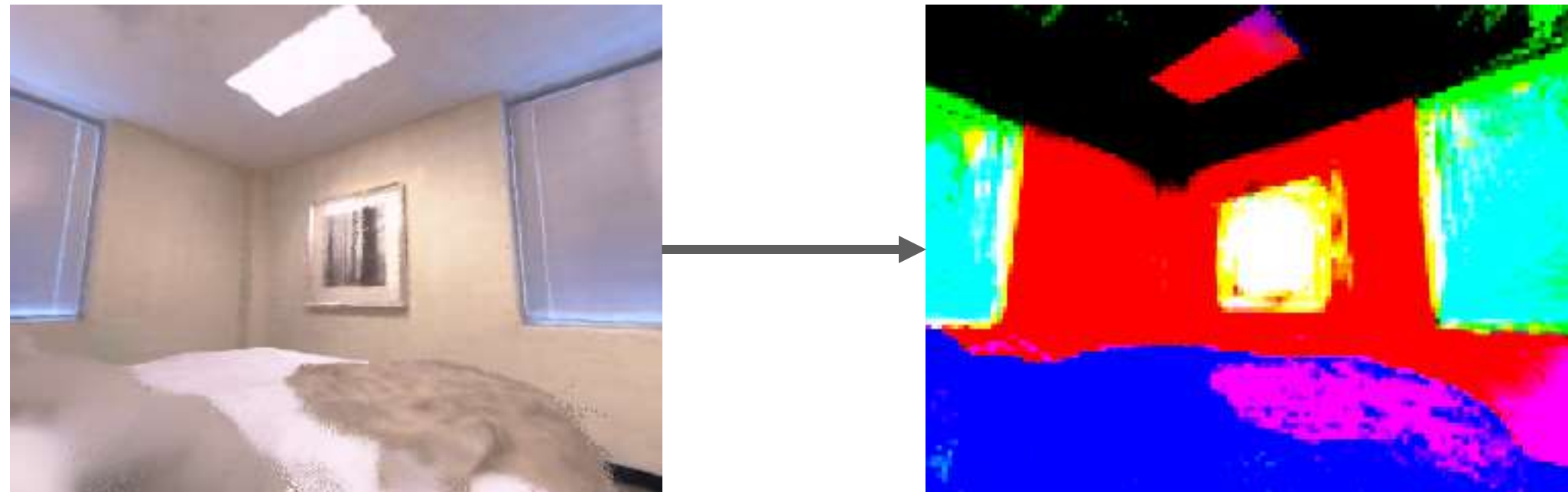# Setting up Semantic "Neuronal Resonances" via Aligning Gradients



These pixels should co-vary

But these should not
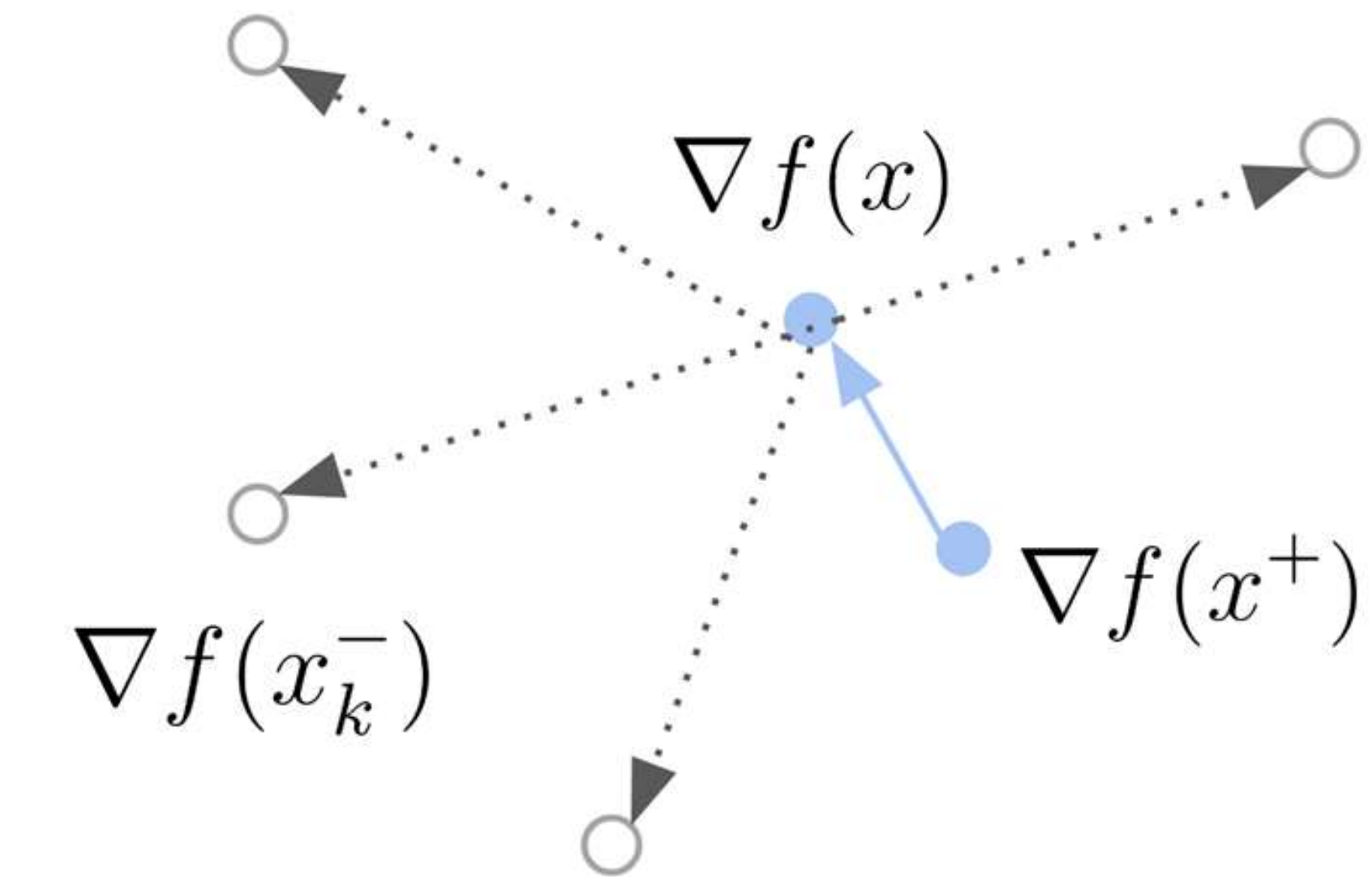
MI-shaping

$\theta$

**Shaping co-aligns gradients of correlated points (here points of the same semantic class)**

# NeRF MLP Shaping via Mutual Information Gradients
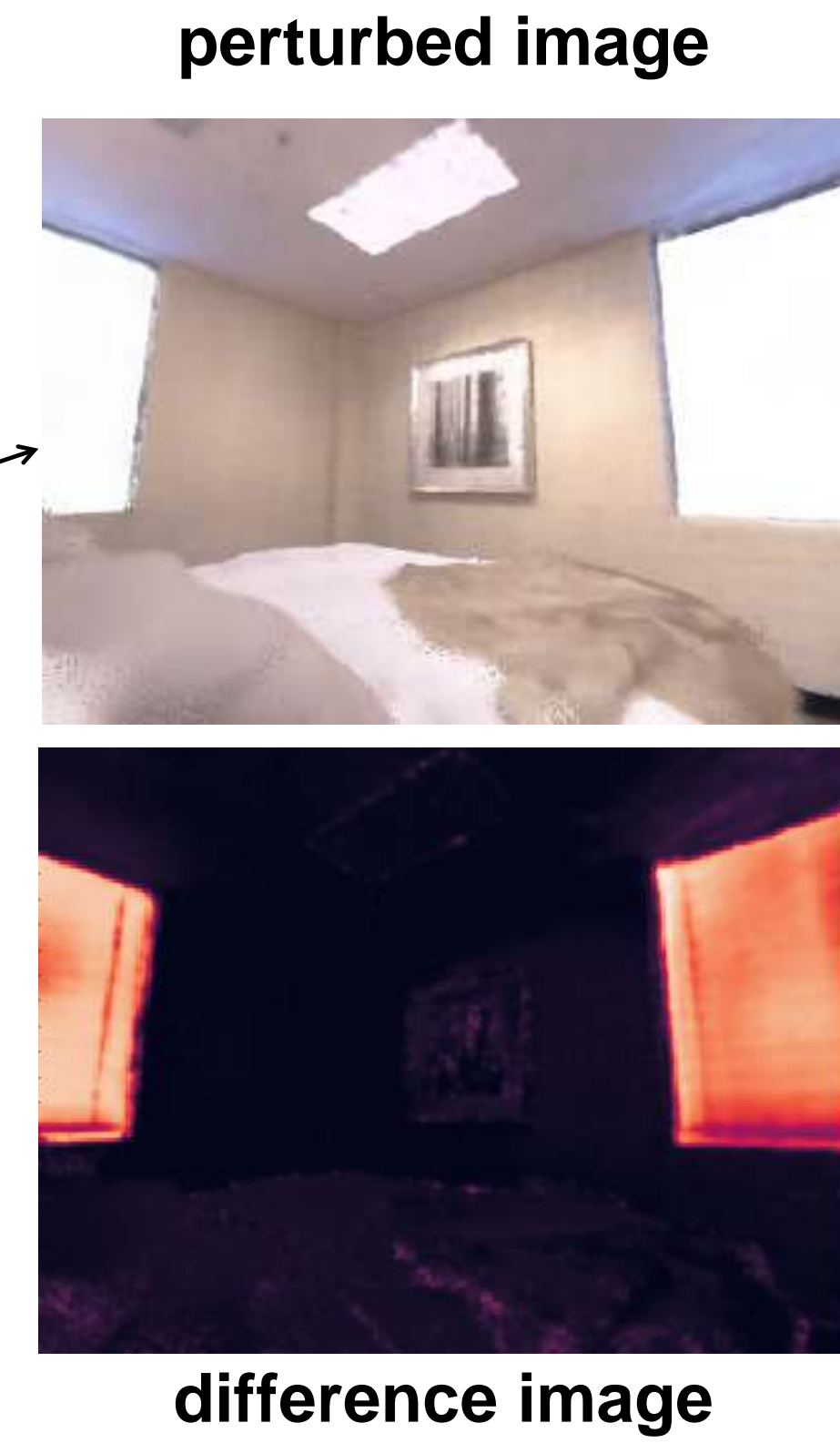


General purpose
image feature: DINO



InfoNCE contrastive loss on
gradients + reconstruction loss

**Obtain some source of semantic affinity**

**Contrastive training
aligning gradients**

# NeRF Shaping Causes Gradient Alignment



perturbed image

Before shaping

difference image

volume rendering

MI-shaping

$\theta$

After shaping

perturbed image

difference image

▲ MLP w/o shaping

● MLP with shaping

● point to perturb

⟶ gradient of the point
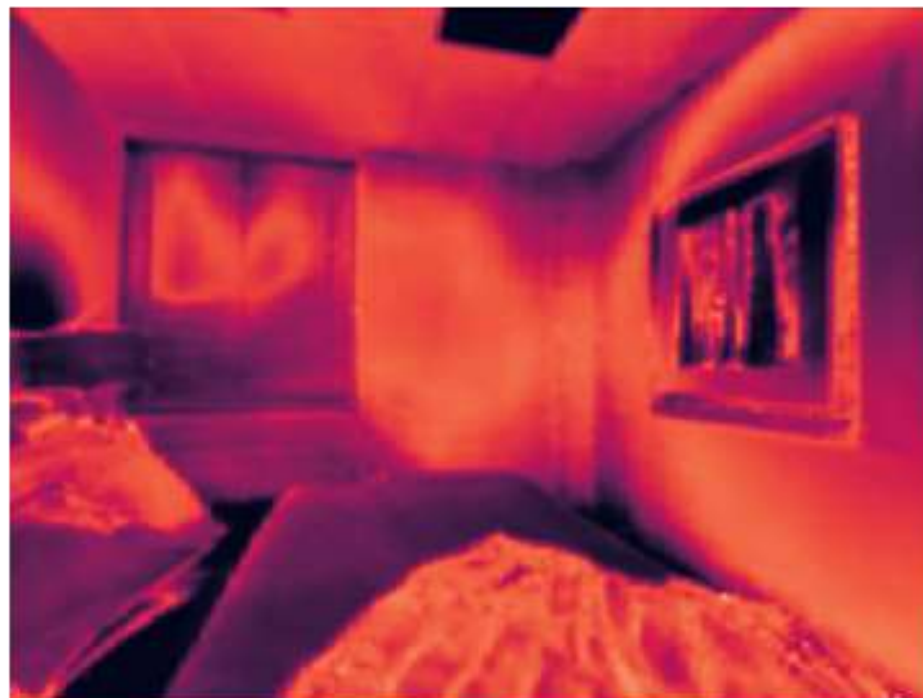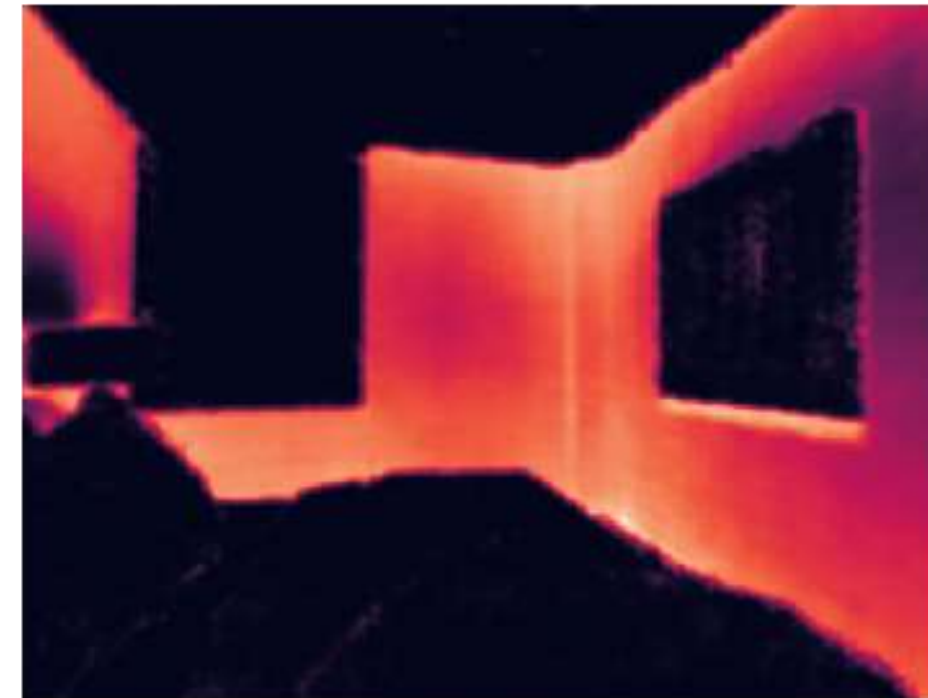
# Application: Entity Selection
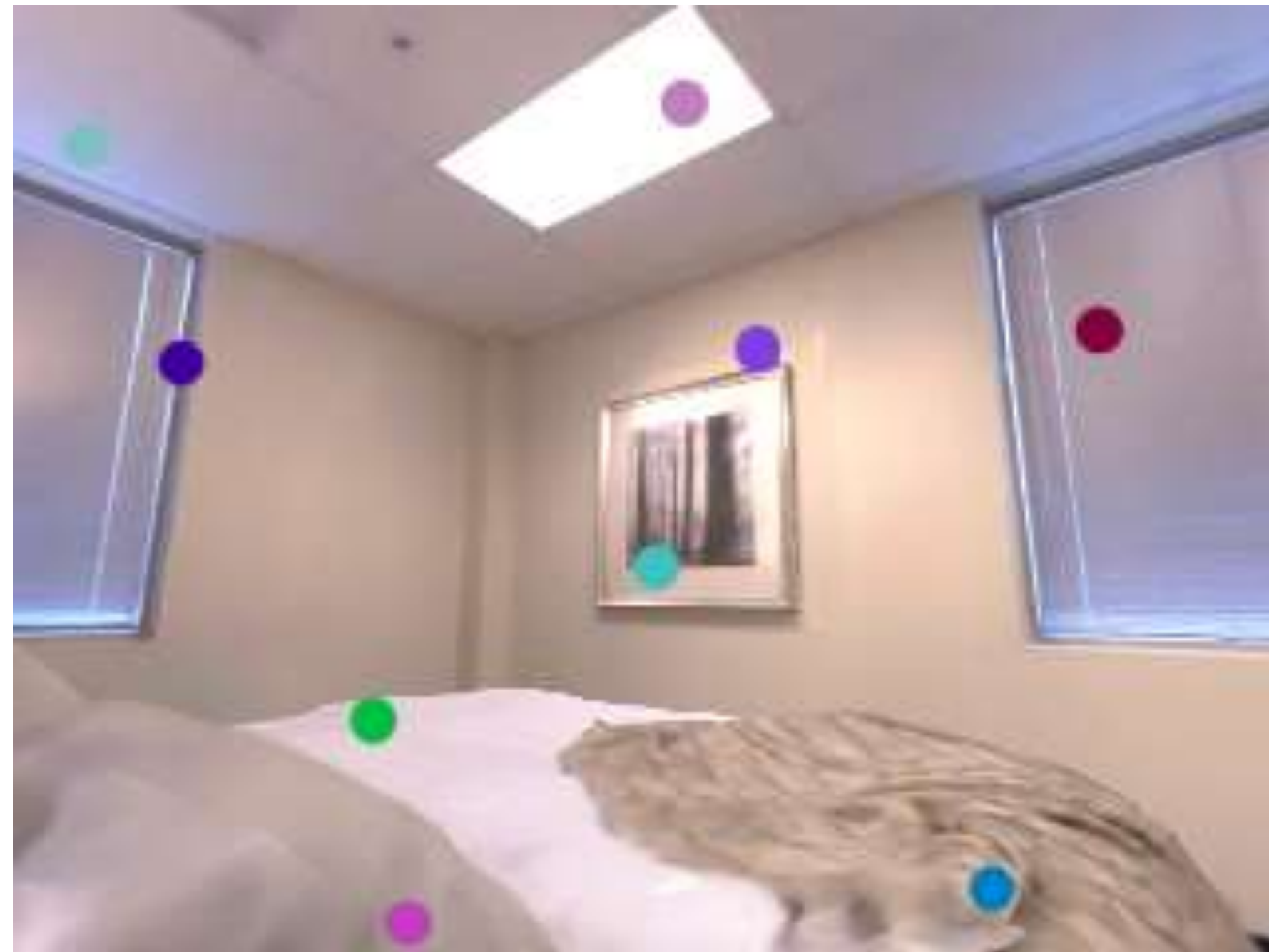


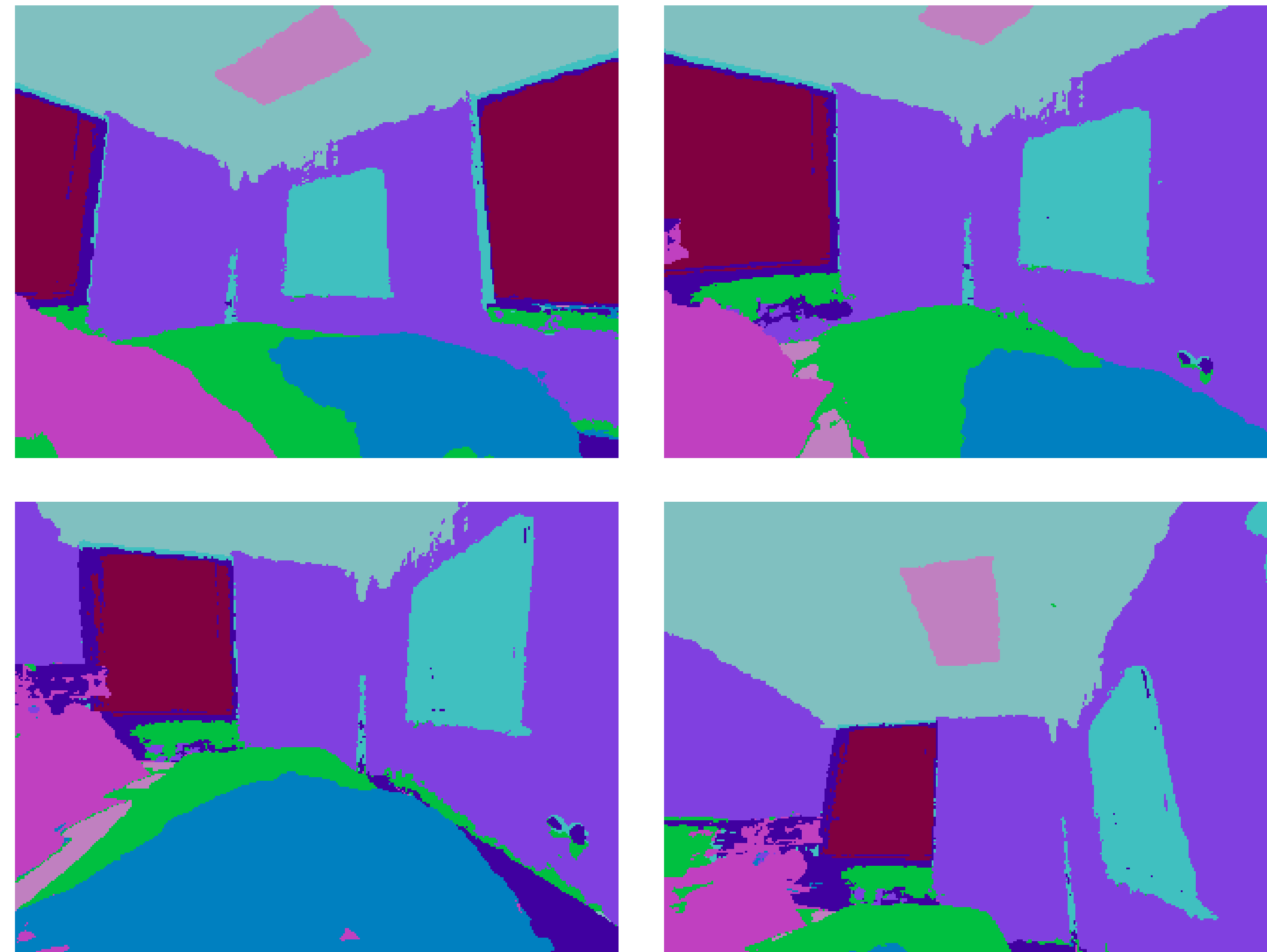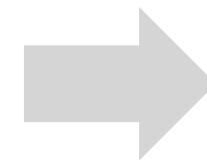image        NeRF        JacobiNeRF

**From a single point we can select an entire semantic entity.**

# Application: Label Propagation

Acquire dense labels of a scene given sparse annotations.



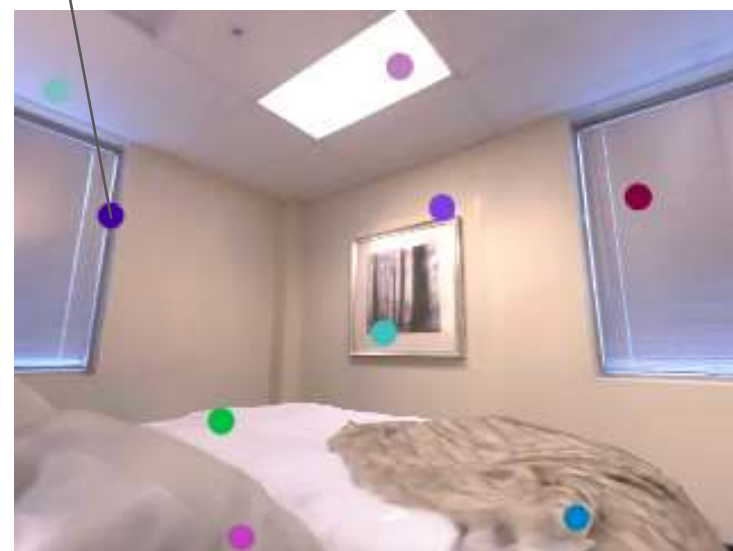Label one pixel for each class
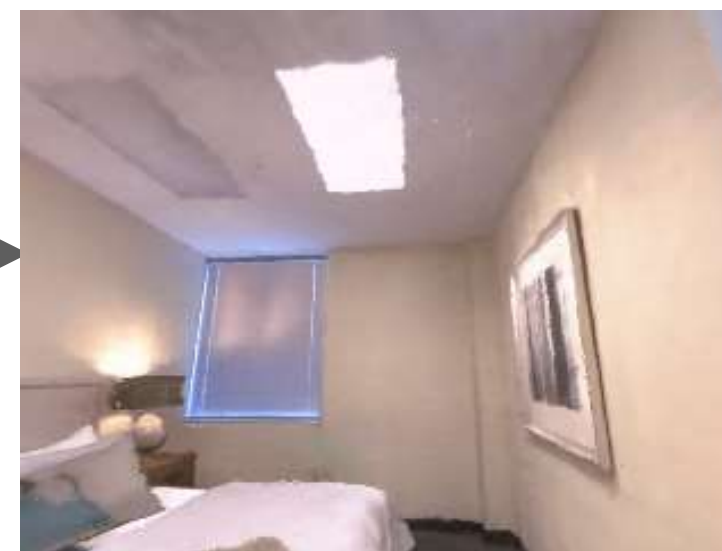from one view

Dense label for any view

## 2D version (JacobiNeRF-2D)



**Given m labels**

**Perturb along gradients**
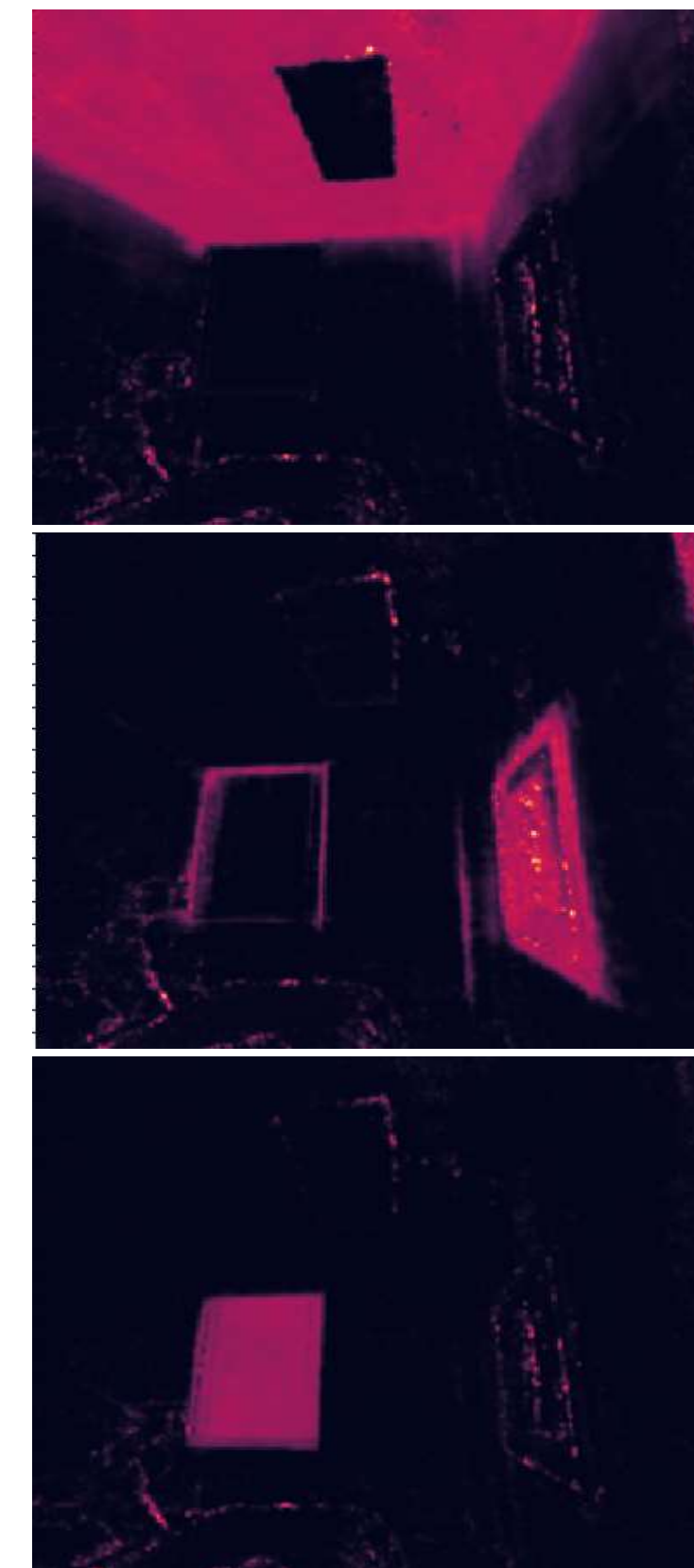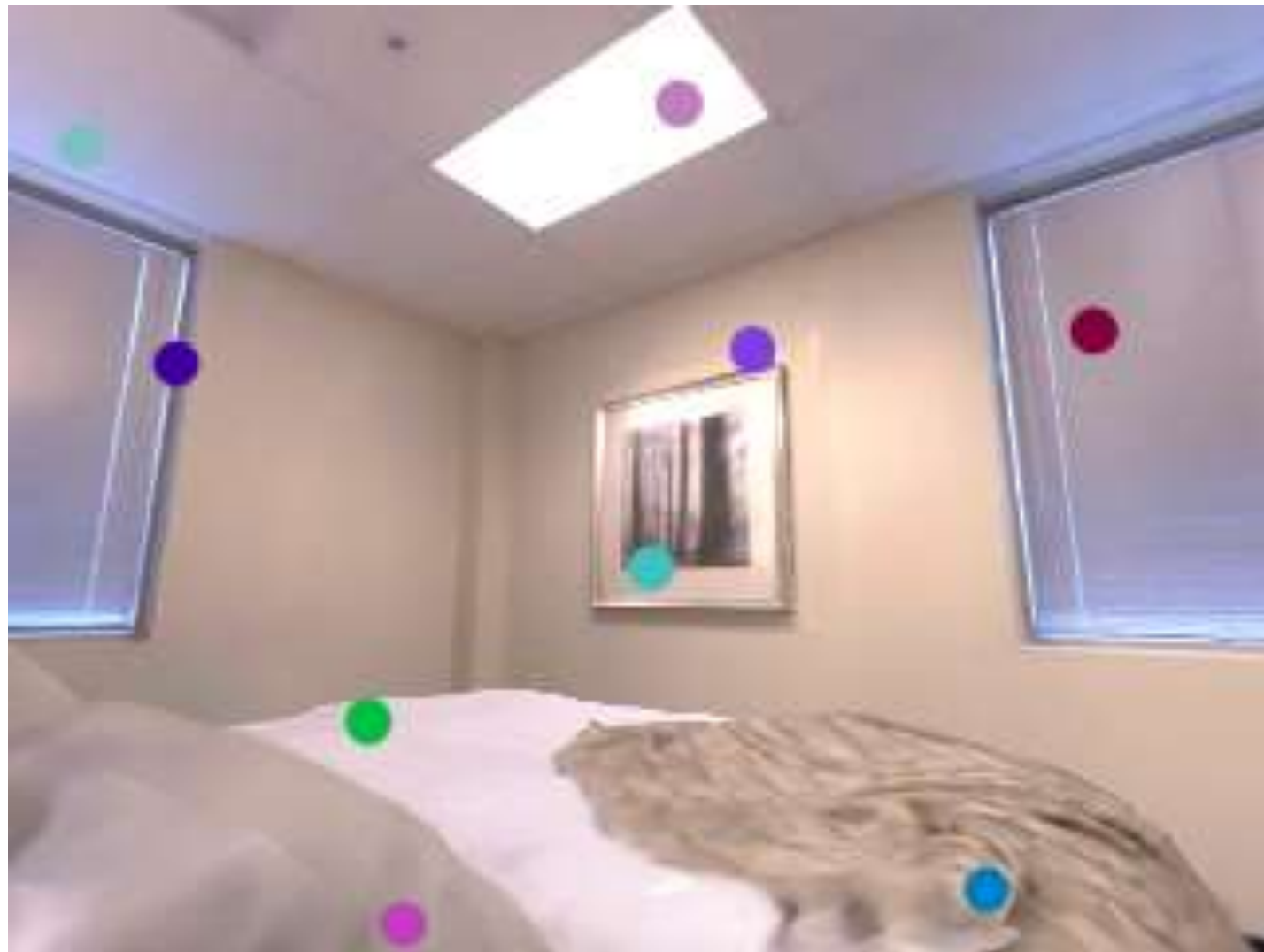
**argmax**

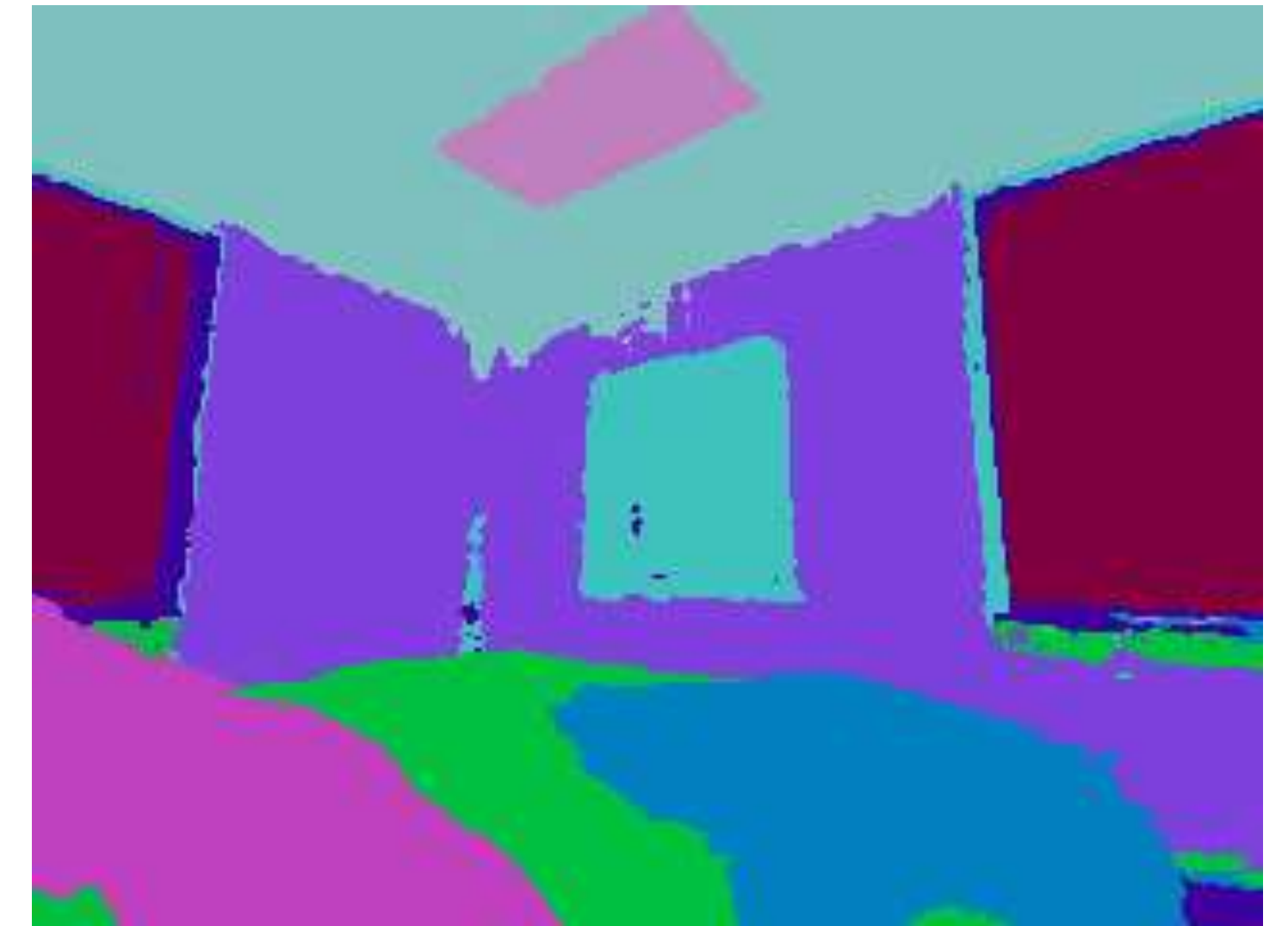**Source view**

**Target view**

**……**

**m responses**

**Label**

# Semantic Segmentation (sparse 1pix/class, Replica)



Given label

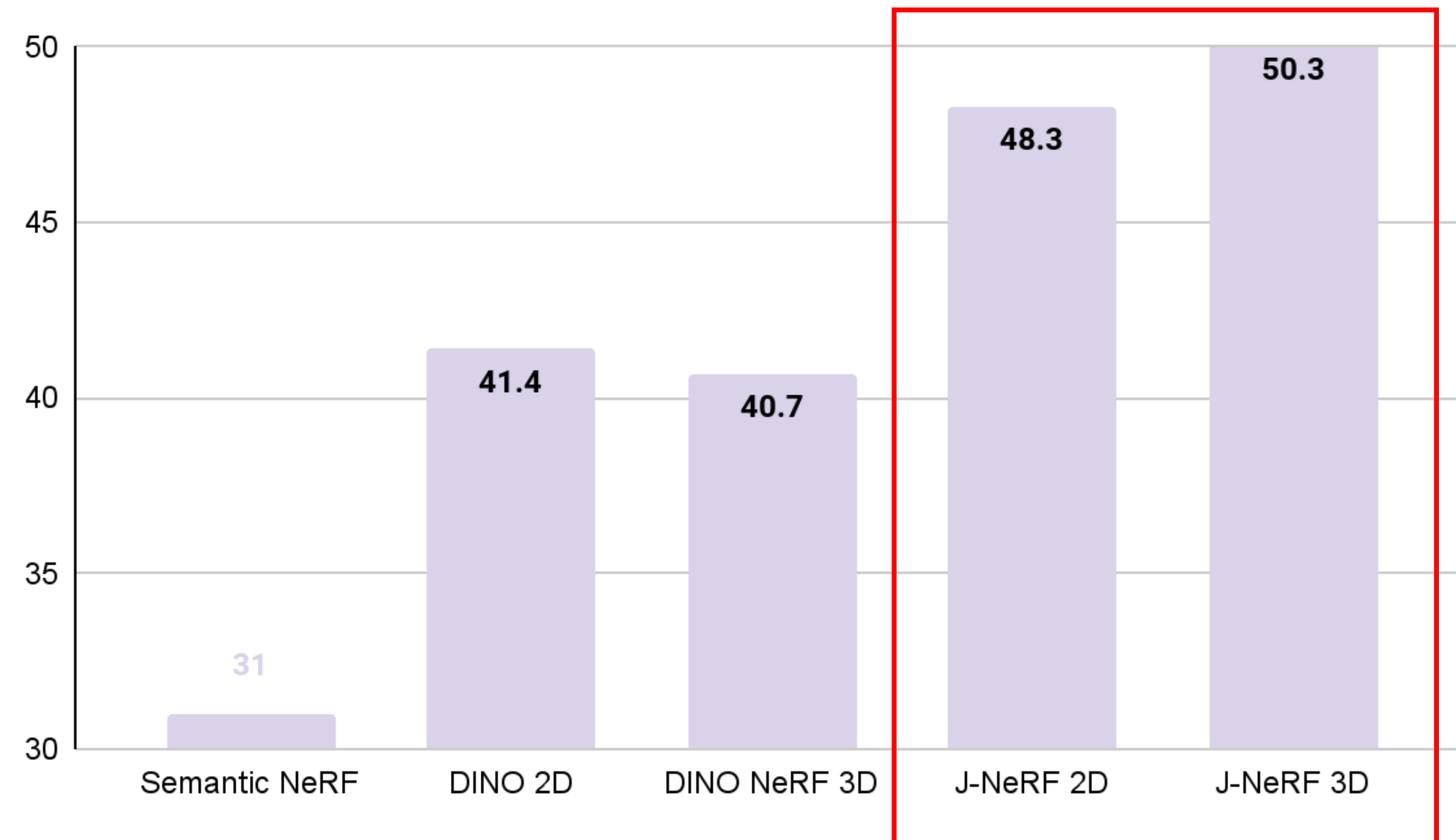J-NeRF 3D

# Semantic Segmentation (sparse 1pix/class, Replica)



1 pix/class 1 view

**mIoU**

| Semantic NeRF | DINO 2D | DINO NeRF 3D | J-NeRF 2D | J-NeRF 3D |
|---|---|---|---|---|
| 18.7 | 18.1 | 25.3 | 26.3 | 28.3 |

1 pix/class 1 view

**Acc**

| Semantic NeRF | DINO 2D | DINO NeRF 3D | J-NeRF 2D | J-NeRF 3D |
|---|---|---|---|---|
| 31 | 41.4 | 40.7 | 48.3 | 50.3 |

Average results on 7 scenes, 180 test views for each scene

Given label

J-NeRF 3D

# Semantic Segmentation (dense 1view, Replica)

Dense label 1 view



mIoU

Dense label 1 view



Acc

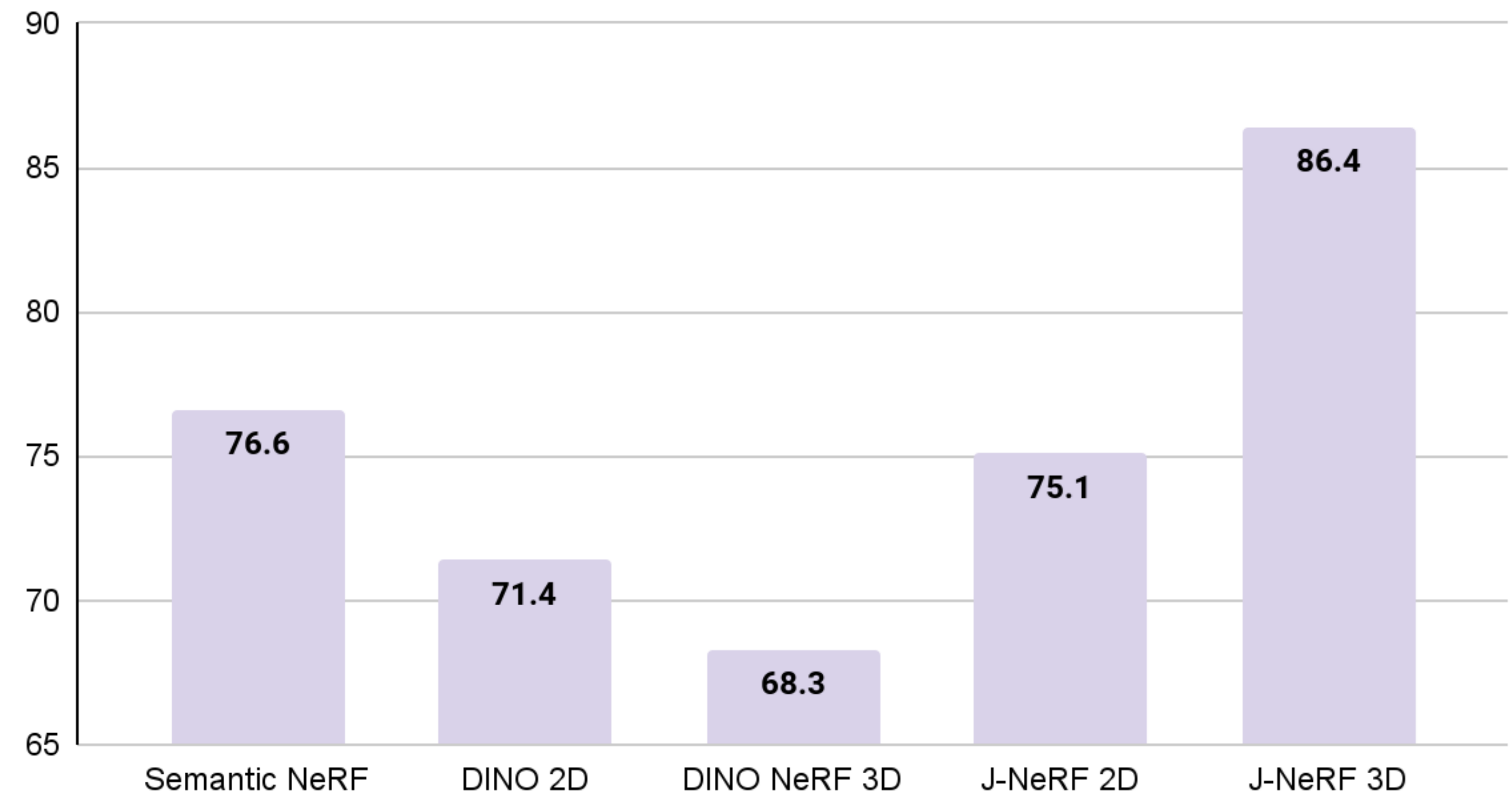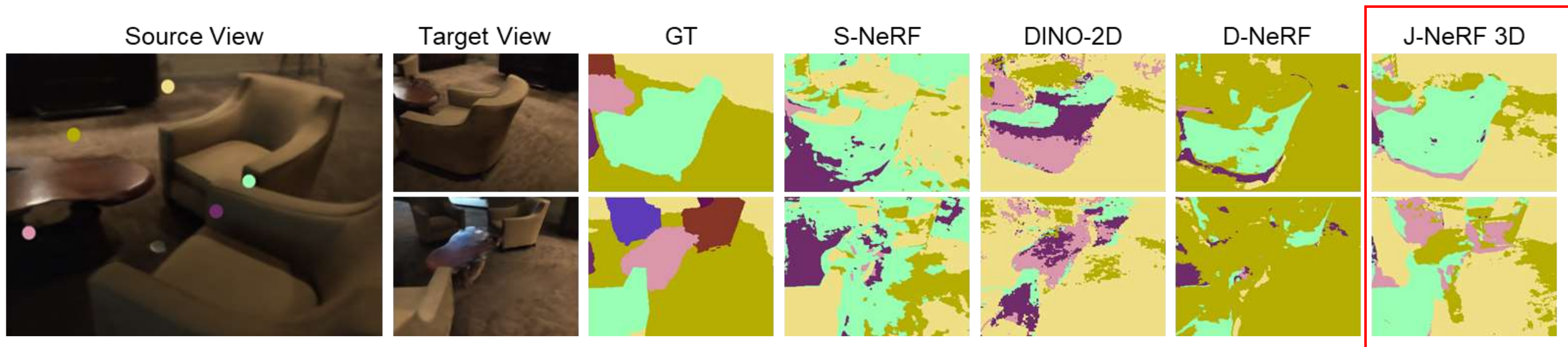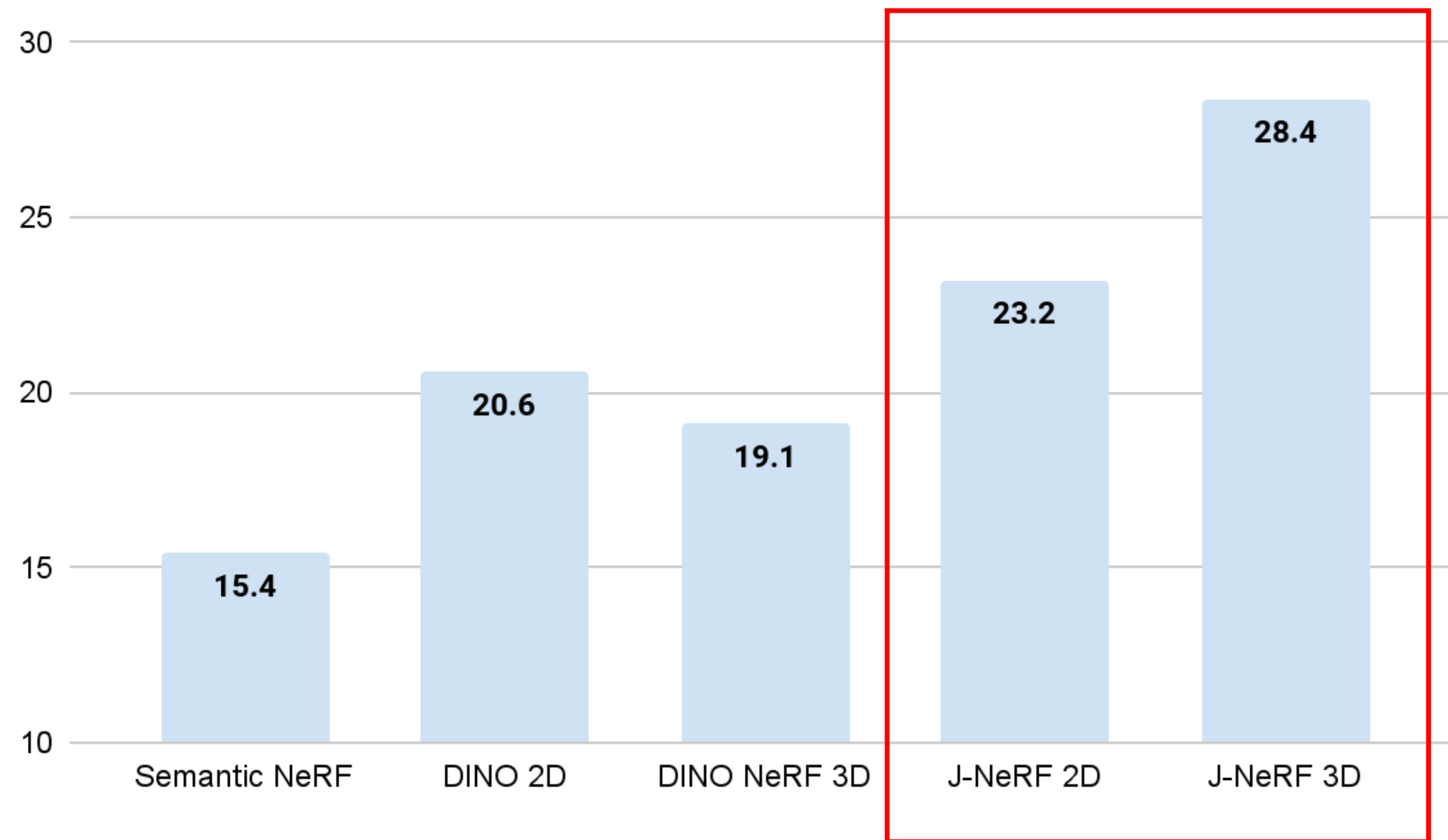Average results on 7 scenes, 180 test views for each scene
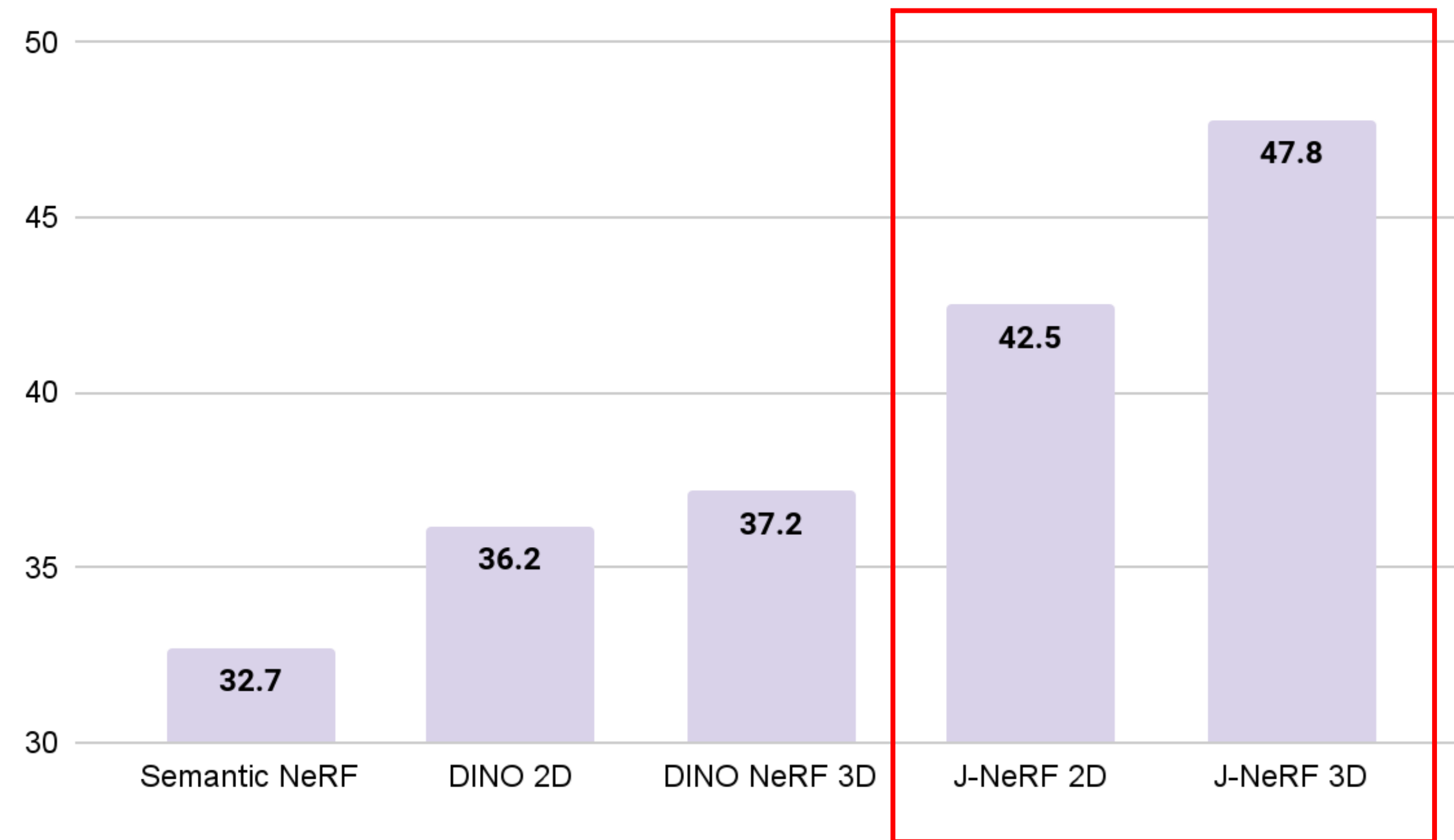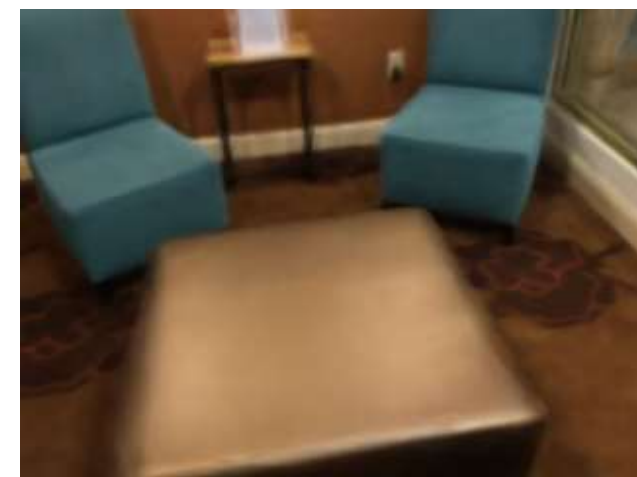
# Instance segmentation (sparse, 1pix/instance, ScanNet)



Source View | Target View | GT | S-NeRF | DINO-2D | D-NeRF | J-NeRF 3D

# Instance segmentation (sparse, 1pix/instance, ScanNet)



1 pix/class 1 view

| | mIoU |
|---|---|
| Semantic NeRF | 15.4 |
| DINO 2D | 20.6 |
| DINO NeRF 3D | 19.1 |
| J-NeRF 2D | 23.2 |
| J-NeRF 3D | 28.4 |

1 pix/class 1 view

| | Acc |
|---|---|
| Semantic NeRF | 32.7 |
| DINO 2D | 36.2 |
| DINO NeRF 3D | 37.2 |
| J-NeRF 2D | 42.5 |
| J-NeRF 3D | 47.8 |

Average results on 4 scenes, ~180 test views for each scene
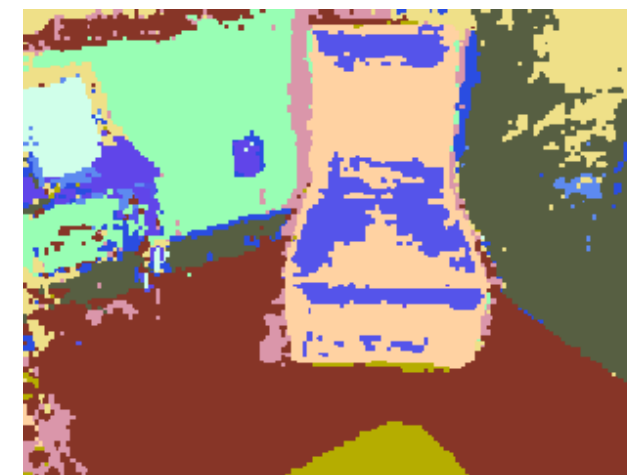
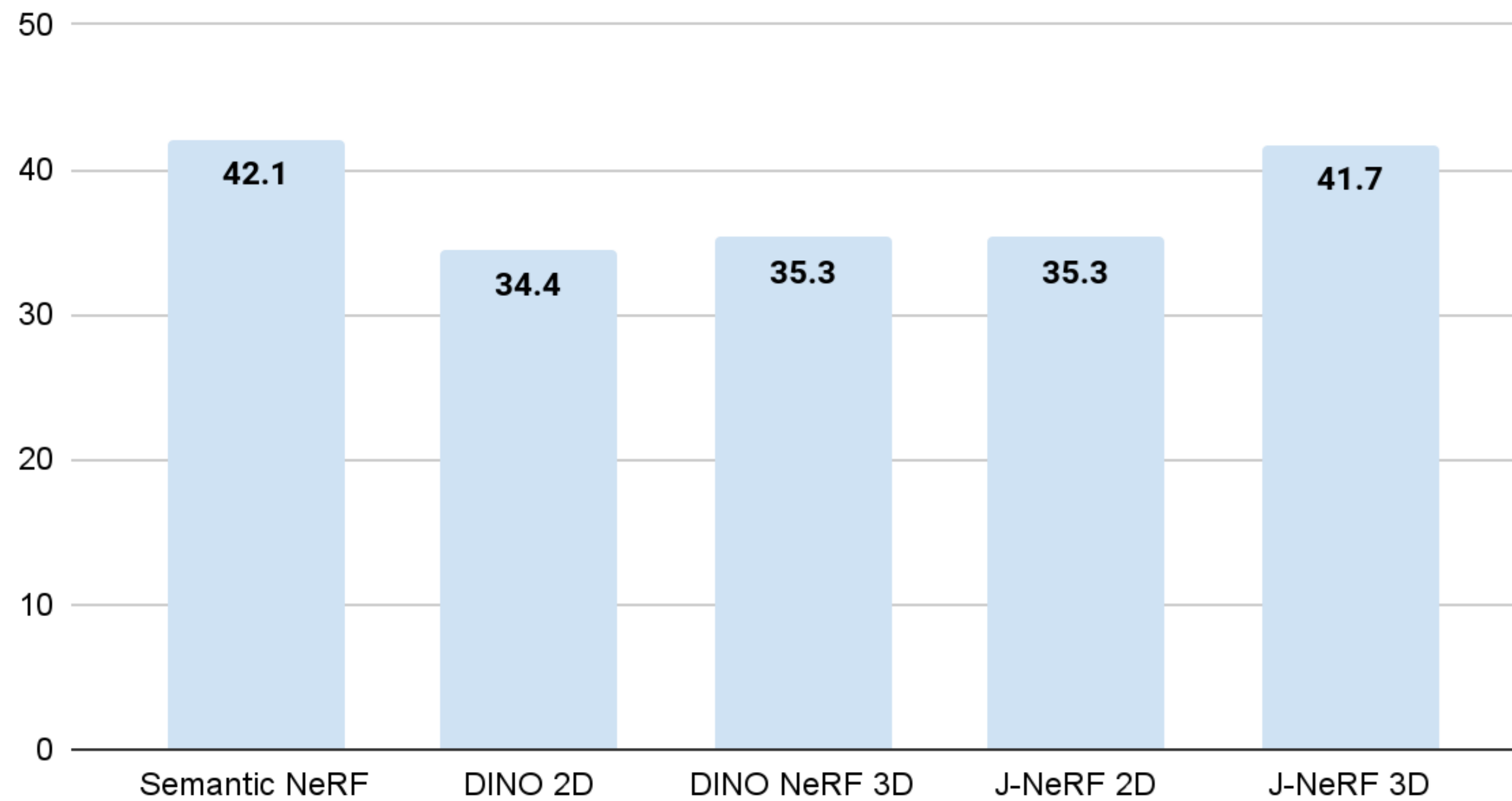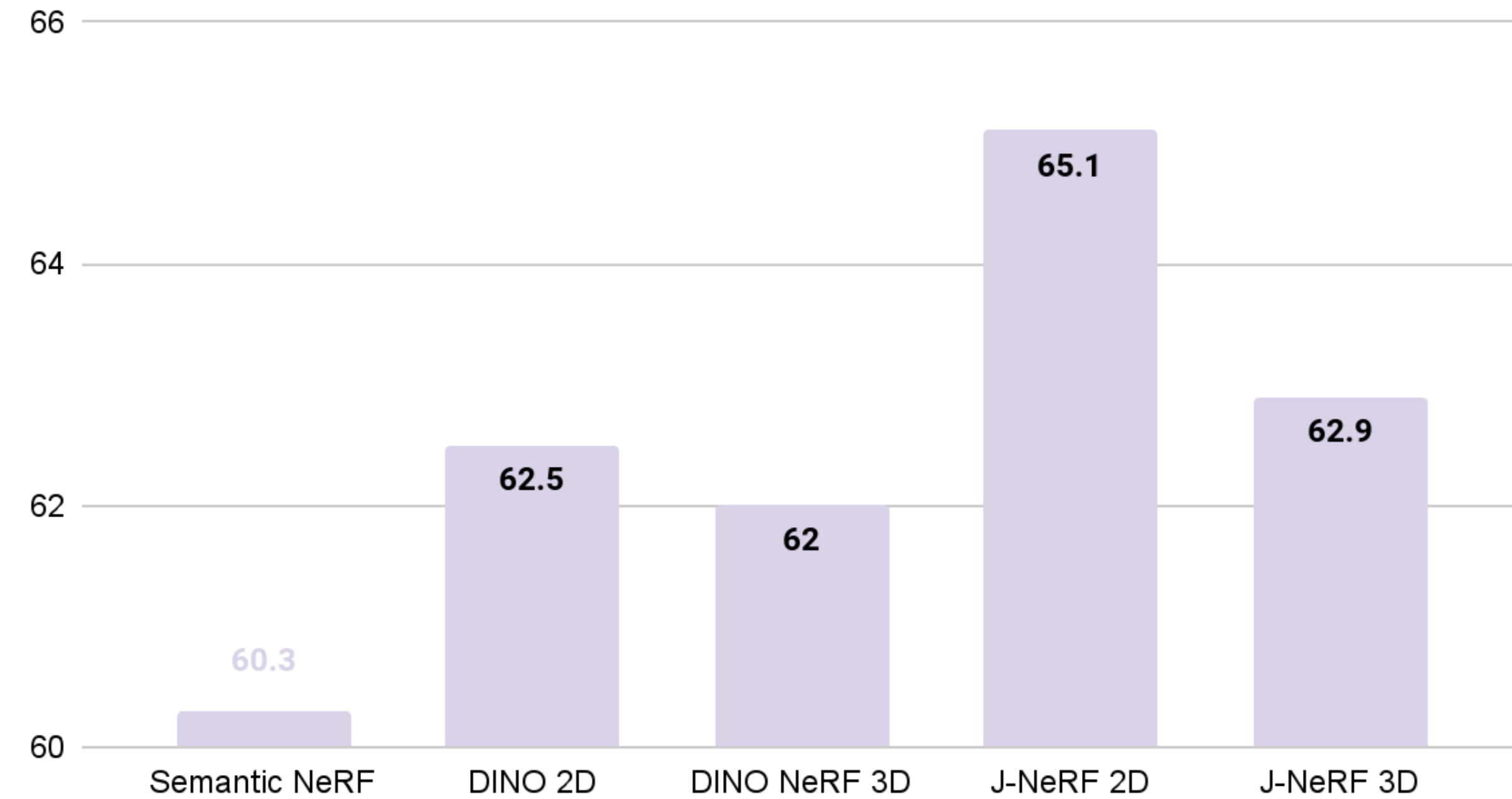| Source view | Target view | GT | DINO-2D | D-NeRF | J-NeRF 3D |

# Instance segmentation (1 view, dense, ScanNet)



Dense label 1 view

mIoU

Dense label 1 view

Acc

Average results on 4 scenes, ~180 test views for each scene

# GitHub Repo

**https://github.com/xxm19/jacobinerf**