

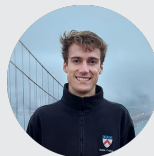
# Feed-Forward **SceneDINO** for Unsupervised Semantic Scene Completion



Aleksandar Jevtić\* <sup>1</sup>



**Christoph Reich**\* <sup>1,2,4,5</sup>



Felix Wimbauer<sup>1,4</sup>



Oliver Hahn<sup>2</sup>



Christian Rupprecht<sup>3</sup>



Stefan Roth<sup>2,5,6</sup>

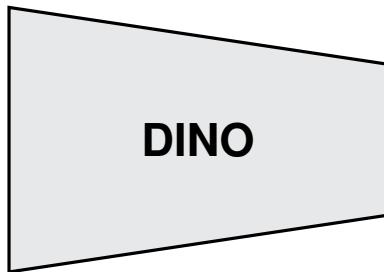


Daniel Cremers<sup>1,4,5</sup>

\*equal contribution



# Motivation

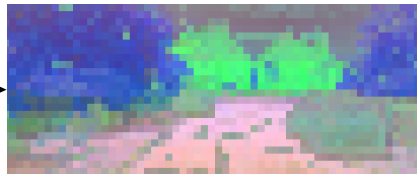
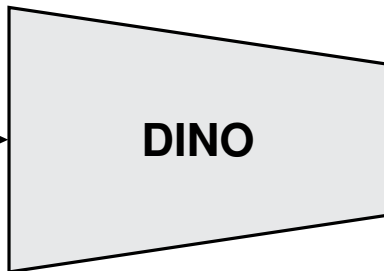


# Motivation



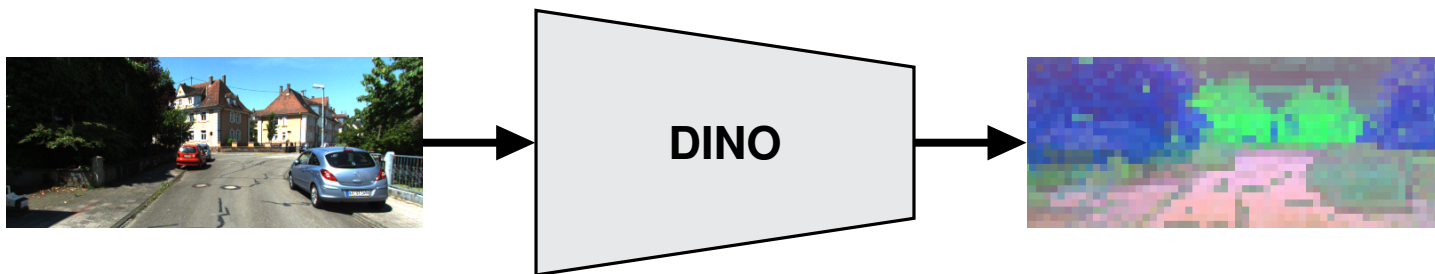
**DINO**

# Motivation



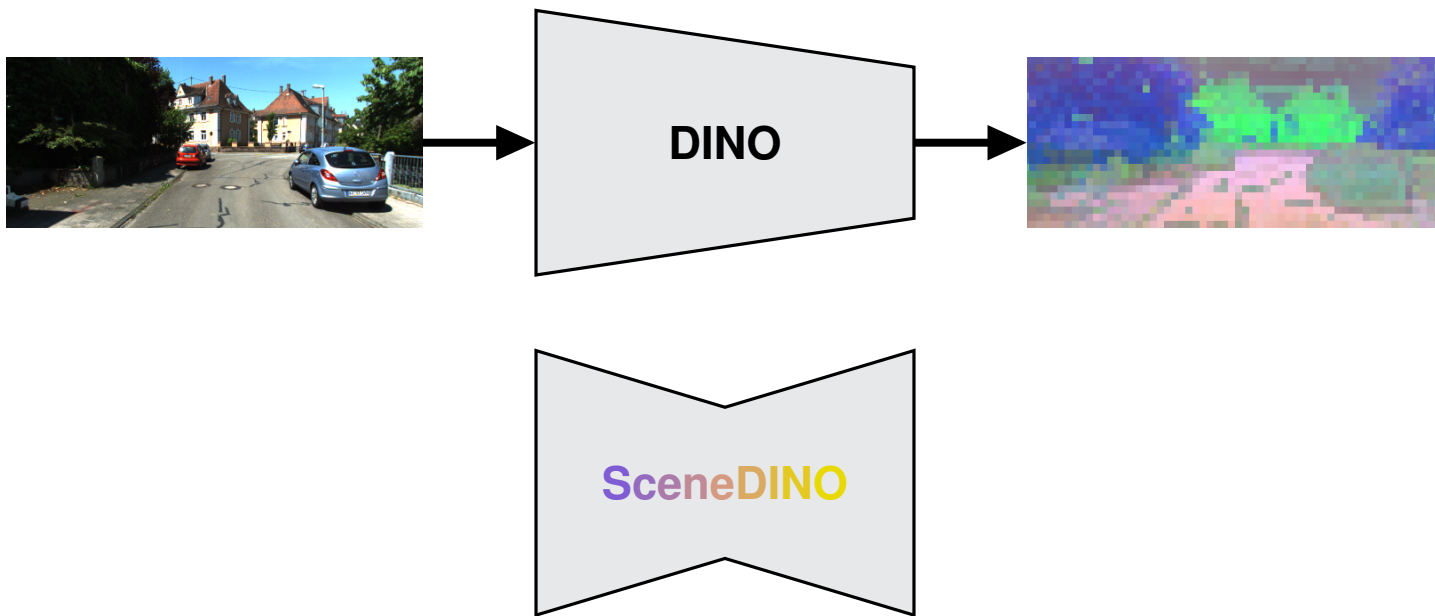


# Motivation



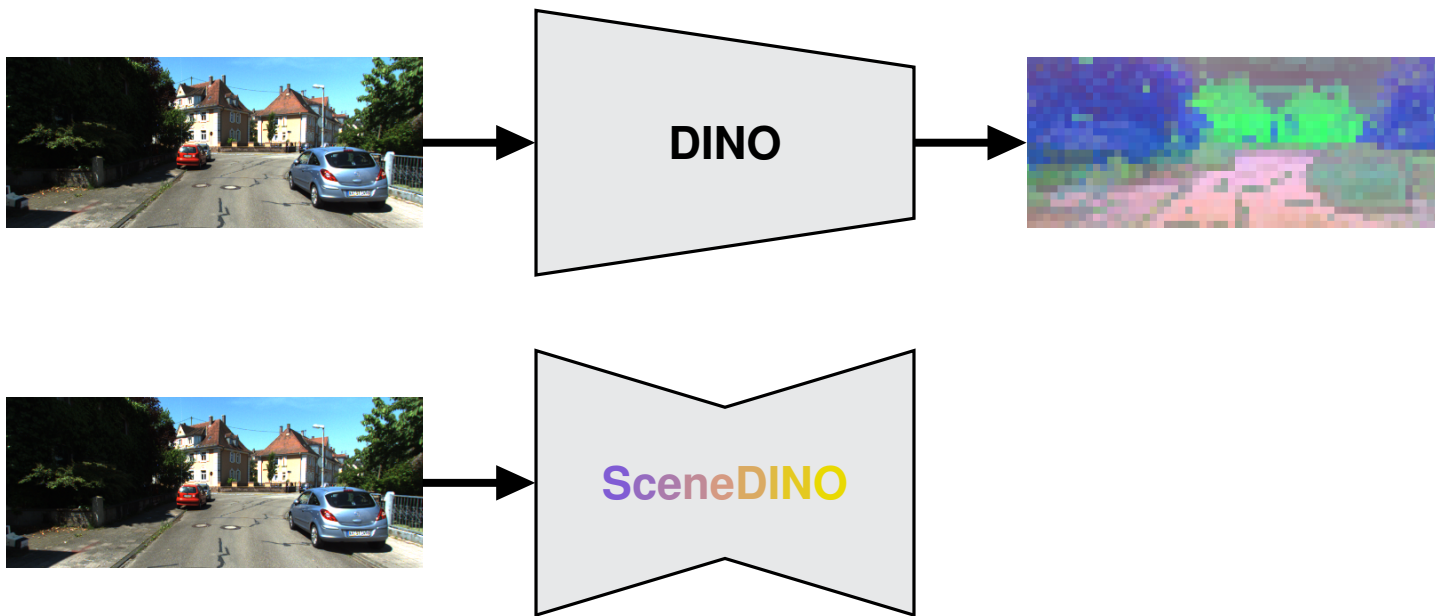
Bring DINO to 3D 🚀

# Motivation



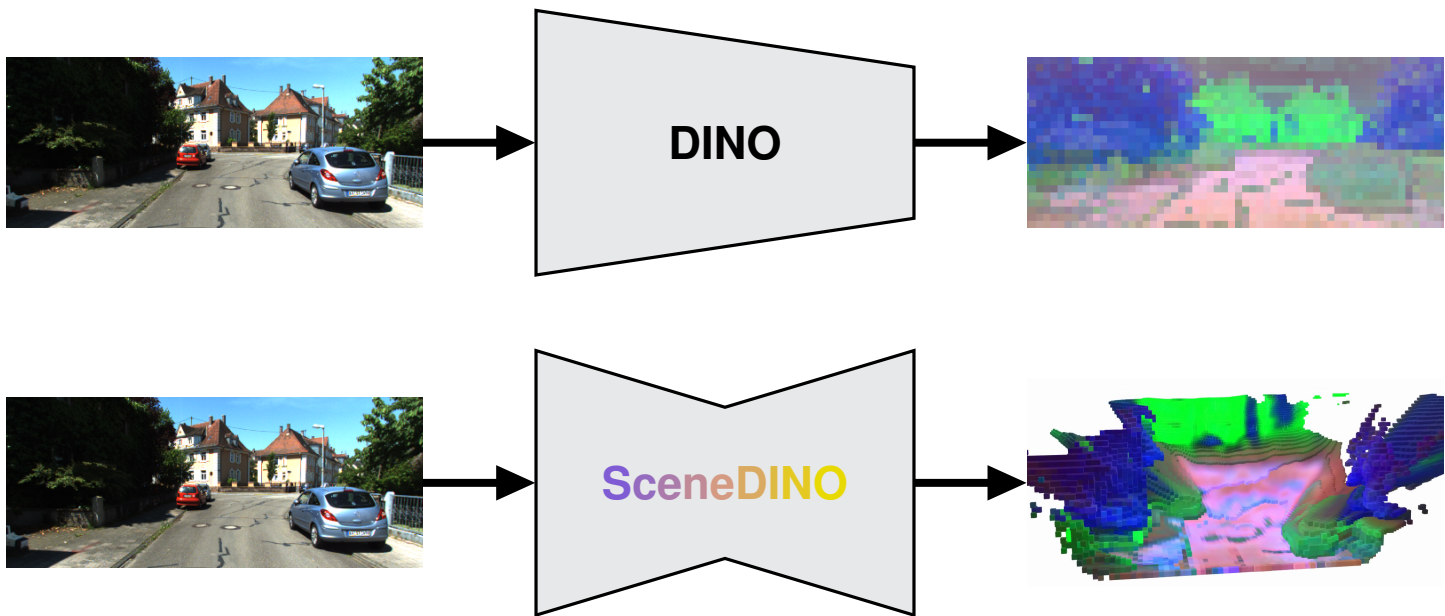
Bring DINO to 3D 🚀

# Motivation



Bring DINO to 3D 🚀

# Motivation



Bring DINO to 3D 🚀

# Semantic Scene Completion (SSC)

a.k.a. Semantic Occupancy Prediction

$n$  input images



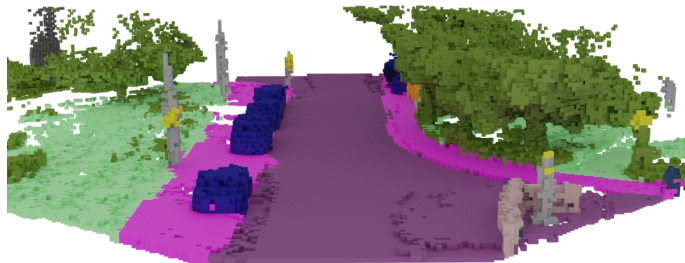
# Semantic Scene Completion (SSC)

a.k.a. Semantic Occupancy Prediction

$n$  input images



Dense 3D geometry & semantics



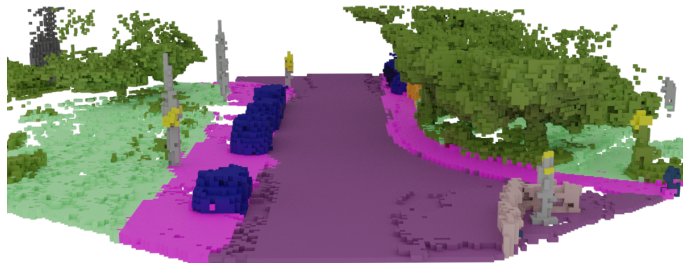
# Semantic Scene Completion (SSC)

a.k.a. Semantic Occupancy Prediction

Single input image



Dense 3D geometry & semantics



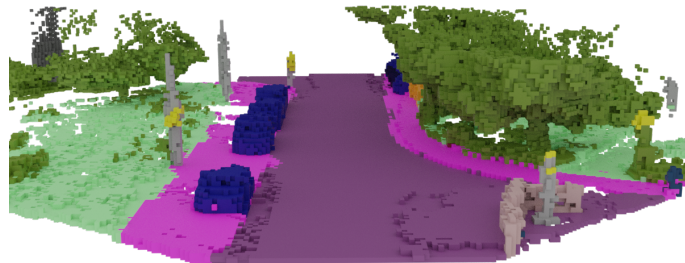
# Semantic Scene Completion (SSC)

a.k.a. Semantic Occupancy Prediction

Single input image



Dense 3D geometry & semantics



✓ Comprehensive 3D scene understanding task



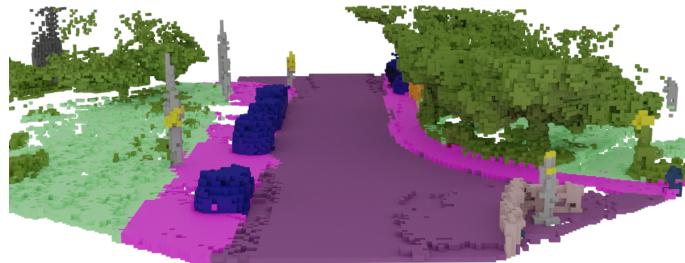
# Semantic Scene Completion (SSC)

a.k.a. Semantic Occupancy Prediction

Single input image



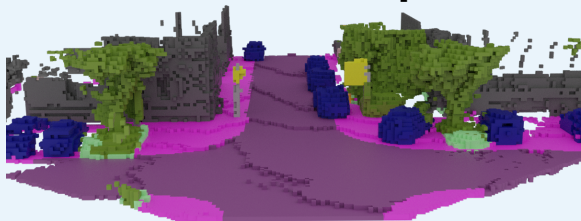
Dense 3D geometry & semantics



- ✓ Comprehensive 3D scene understanding task
- ✓ Applications in robotics, autonomous driving, medical image analysis, and civil engineering

# Related Work: SSC

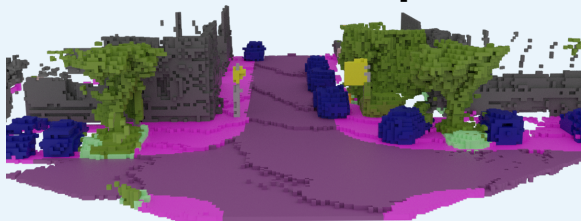
## Geometric & 3D semantic supervision (e.g., [1])



- [1] S. Song *et al.*, “Semantic scene completion from a single depth image,” in *CVPR*, 2017, pp. 190–198.
- [2] Y. Huang *et al.*, “SelfOcc: Self-supervised vision-based 3D occupancy prediction,” in *CVPR*, 2024, pp. 19 946–19 956.

# Related Work: SSC

## Geometric & 3D semantic supervision (e.g., [1])



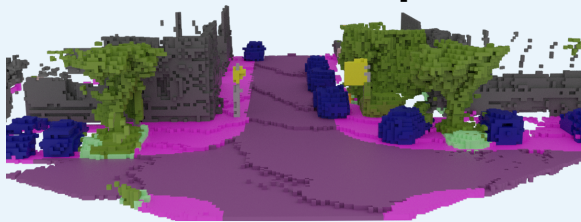
- Ground truth very expensive
- Infeasible to scale
- Special hardware needed

[1] S. Song *et al.*, “Semantic scene completion from a single depth image,” in *CVPR*, 2017, pp. 190–198.

[2] Y. Huang *et al.*, “SelfOcc: Self-supervised vision-based 3D occupancy prediction,” in *CVPR*, 2024, pp. 19 946–19 956.

# Related Work: SSC

## Geometric & 3D semantic supervision (e.g., [1])



- Ground truth very expensive
- Special hardware needed
- Infeasible to scale

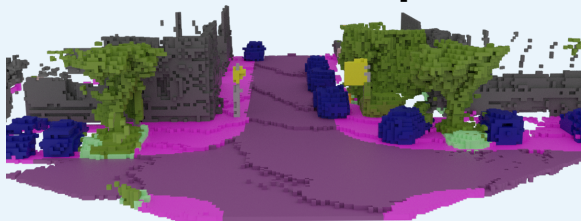
## 2D supervision (e.g., [2])



- [1] S. Song *et al.*, “Semantic scene completion from a single depth image,” in *CVPR*, 2017, pp. 190–198.  
[2] Y. Huang *et al.*, “SelfOcc: Self-supervised vision-based 3D occupancy prediction,” in *CVPR*, 2024, pp. 19 946–19 956.

# Related Work: SSC

## Geometric & 3D semantic supervision (e.g., [1])



- Ground truth very expensive
- Infeasible to scale
- Special hardware needed

## 2D supervision (e.g., [2])



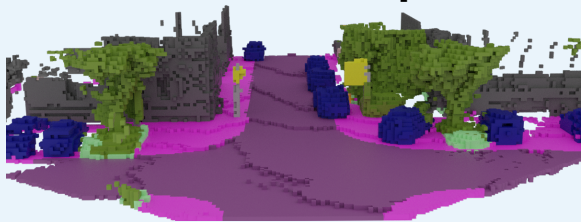
- Still, expensive to obtain
- Limited generalization

[1] S. Song *et al.*, “Semantic scene completion from a single depth image,” in *CVPR*, 2017, pp. 190–198.

[2] Y. Huang *et al.*, “SelfOcc: Self-supervised vision-based 3D occupancy prediction,” in *CVPR*, 2024, pp. 19 946–19 956.

# Related Work: SSC

## Geometric & 3D semantic supervision (e.g., [1])



- Ground truth very expensive
- Infeasible to scale
- Special hardware needed

## 2D supervision (e.g., [2])



- Still, expensive to obtain
- Limited generalization

**Large-scale SSC annotations infeasible → unsupervised SSC**

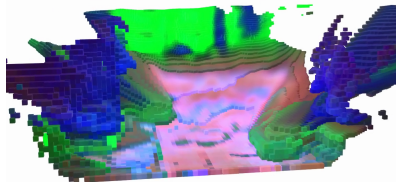
[1] S. Song *et al.*, “Semantic scene completion from a single depth image,” in *CVPR*, 2017, pp. 190–198.

[2] Y. Huang *et al.*, “SelfOcc: Self-supervised vision-based 3D occupancy prediction,” in *CVPR*, 2024, pp. 19 946–19 956.



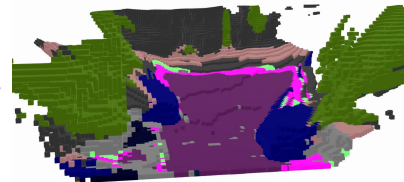
Single Input Image

SceneDINO



3D Feature Field

Distill & Cluster



SSC Prediction

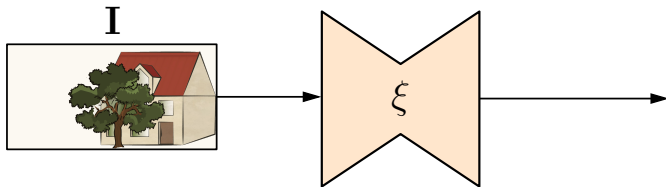
✓ Fully unsupervised

✓ Multi-view self-supervision

✓ Feed-forward inference

# Model Architecture

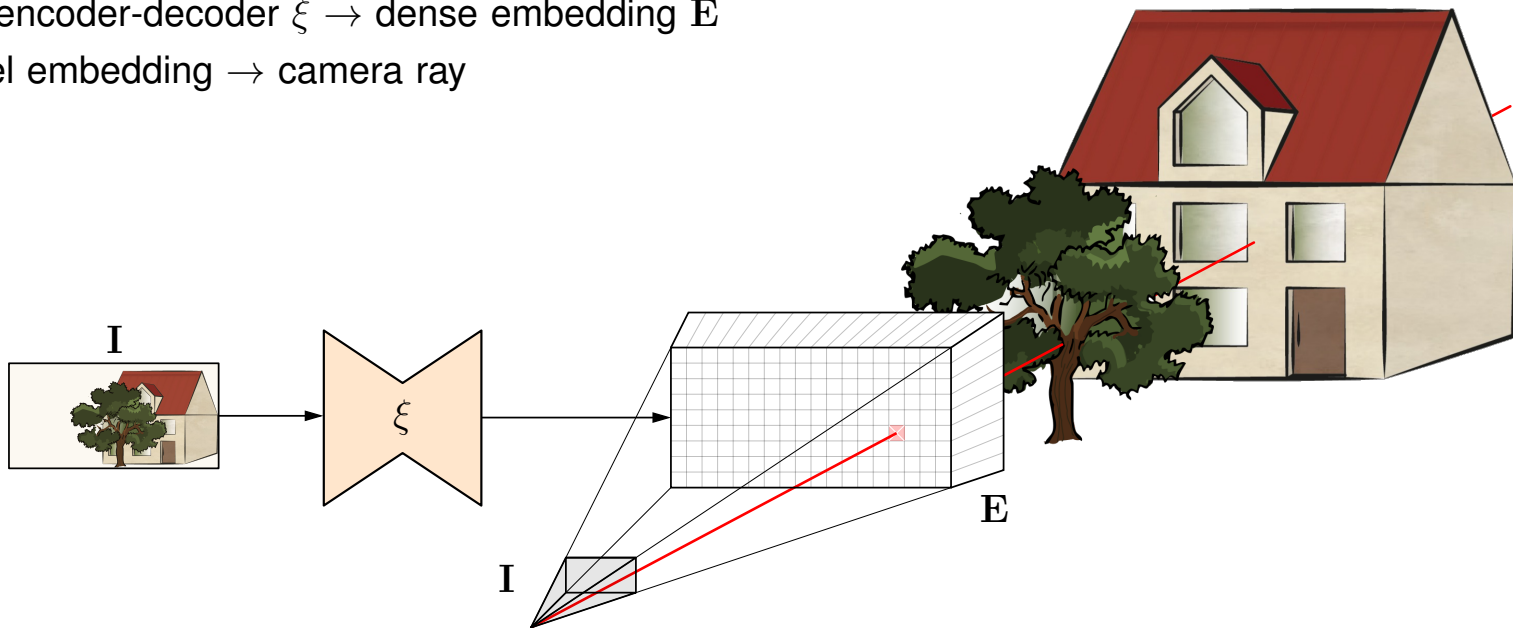
- Single input image  $I$
- 2D encoder-decoder  $\xi \rightarrow$  dense embedding  $E$





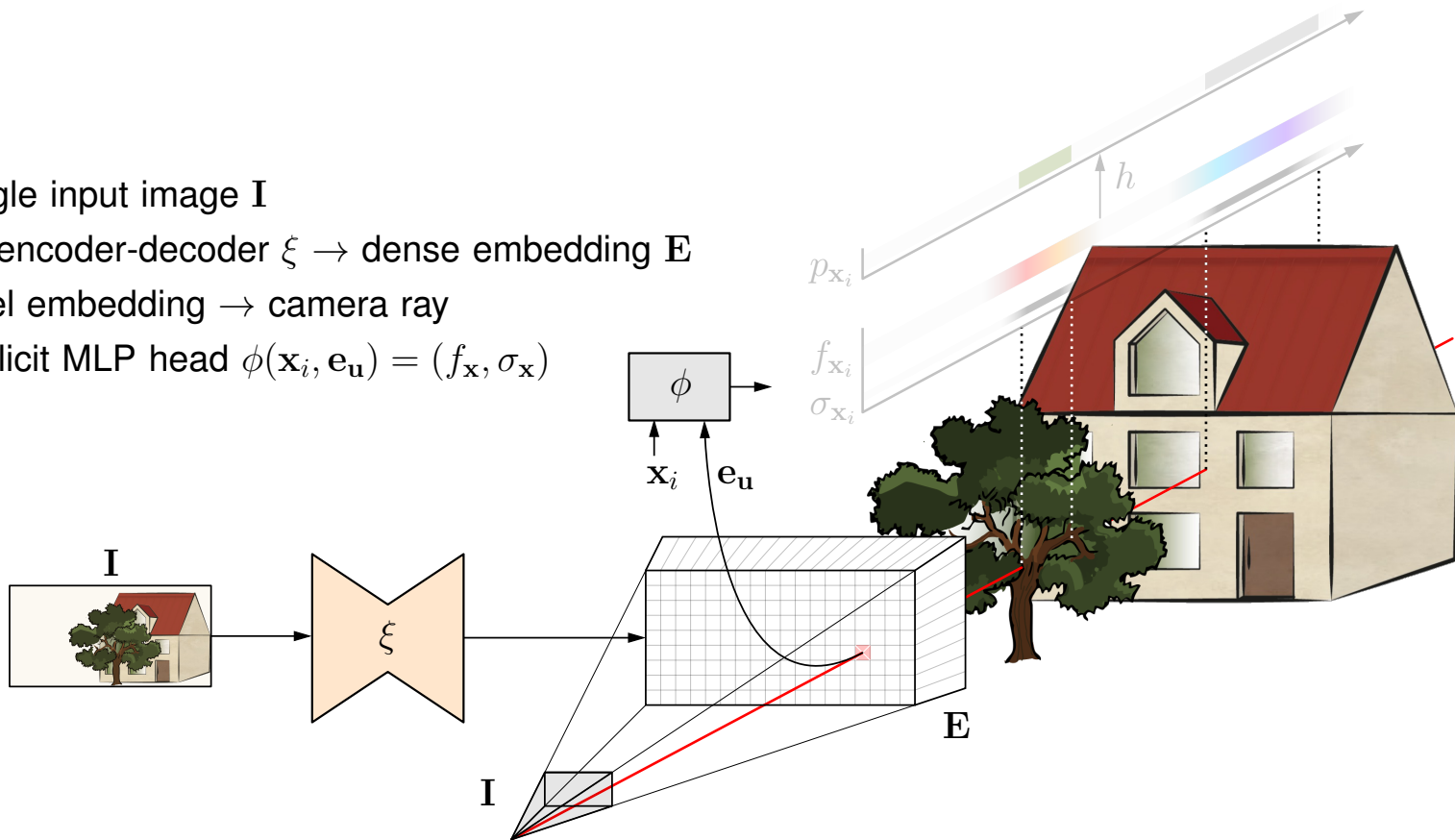
# Model Architecture

- Single input image  $I$
- 2D encoder-decoder  $\xi \rightarrow$  dense embedding  $E$
- Pixel embedding  $\rightarrow$  camera ray



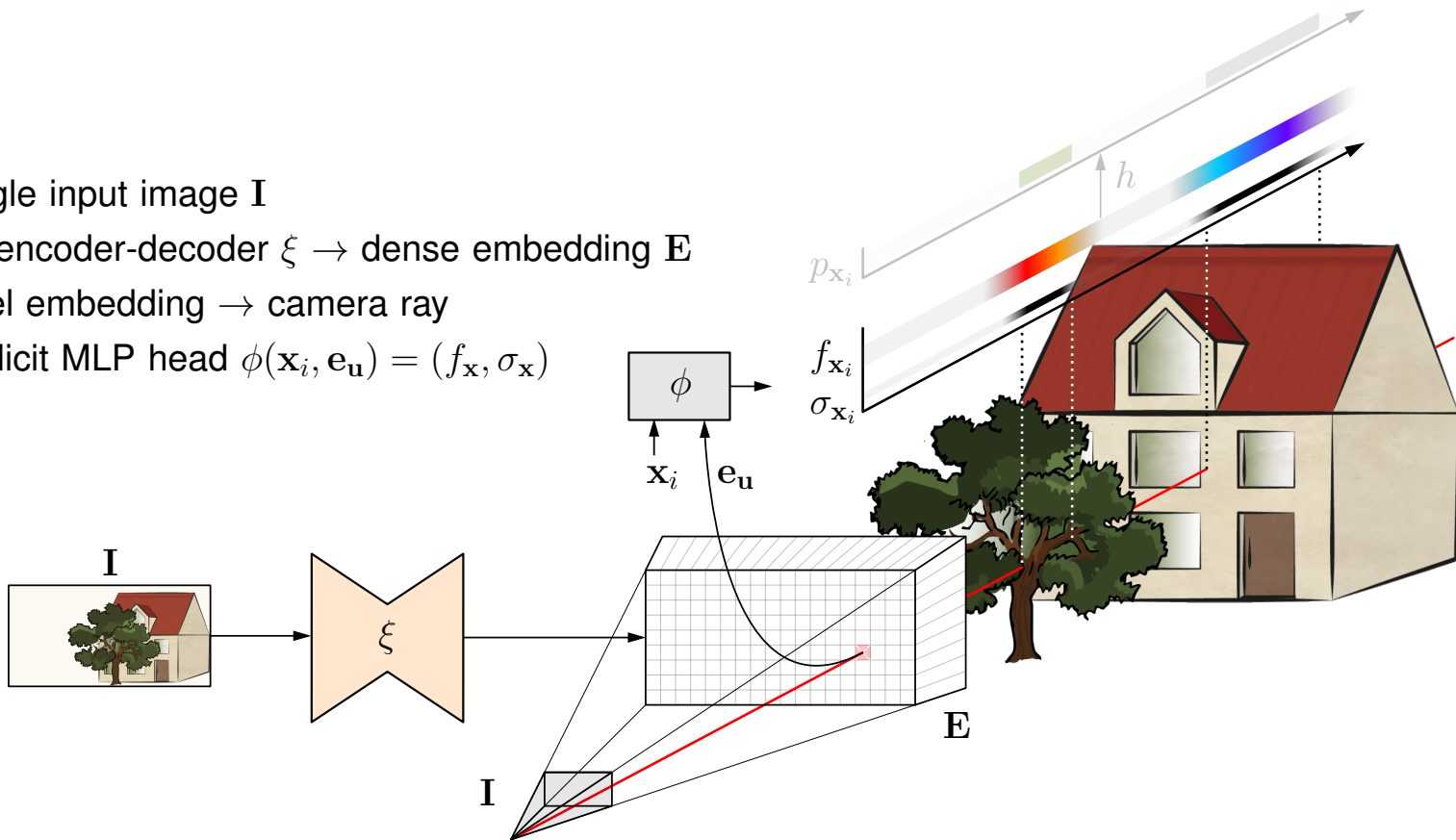
# Model Architecture

- Single input image  $I$
- 2D encoder-decoder  $\xi \rightarrow$  dense embedding  $E$
- Pixel embedding  $\rightarrow$  camera ray
- Implicit MLP head  $\phi(\mathbf{x}_i, \mathbf{e}_u) = (f_{\mathbf{x}}, \sigma_{\mathbf{x}})$



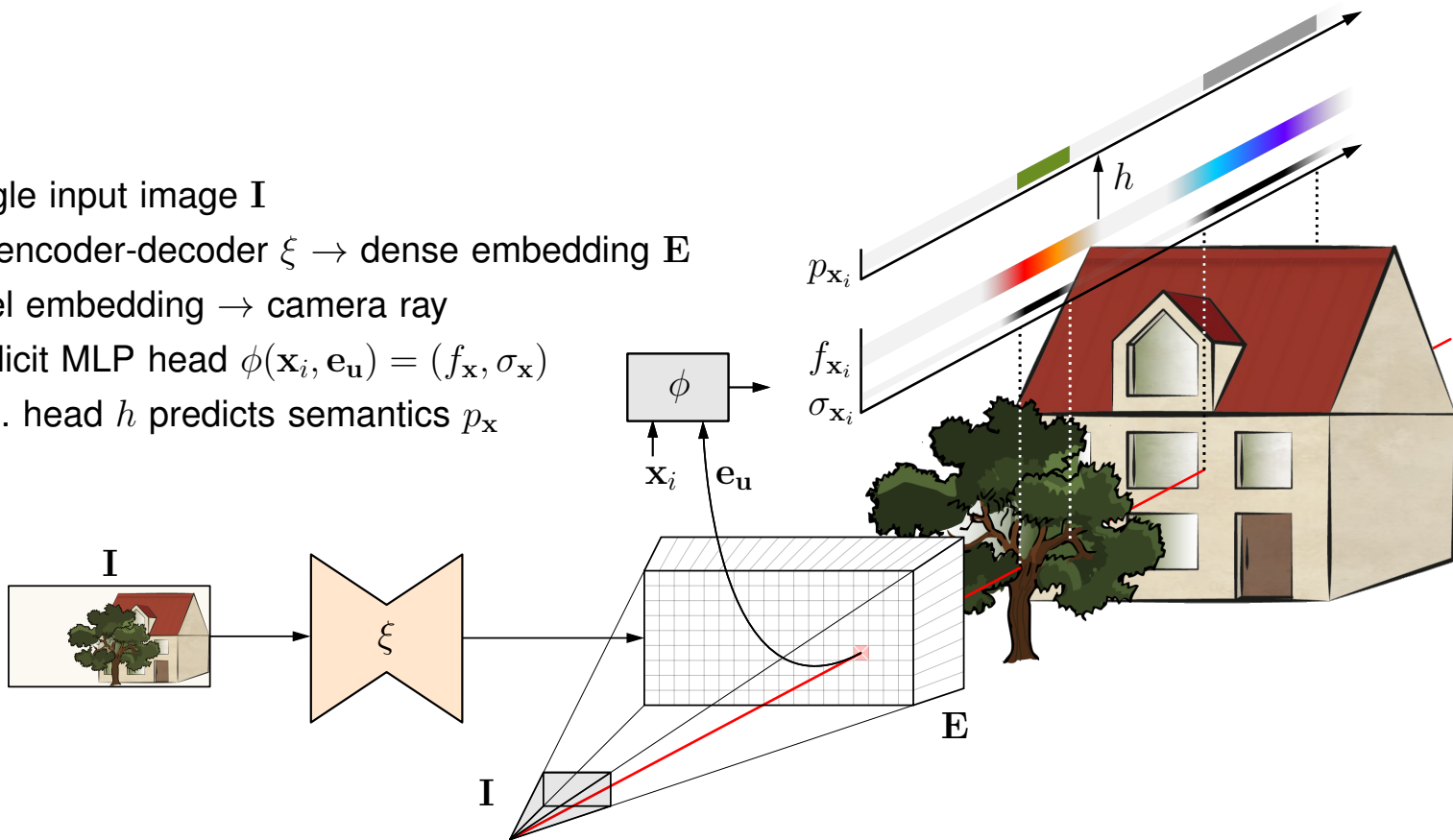
# Model Architecture

- Single input image  $I$
- 2D encoder-decoder  $\xi \rightarrow$  dense embedding  $E$
- Pixel embedding  $\rightarrow$  camera ray
- Implicit MLP head  $\phi(\mathbf{x}_i, \mathbf{e}_u) = (f_{\mathbf{x}}, \sigma_{\mathbf{x}})$



# Model Architecture

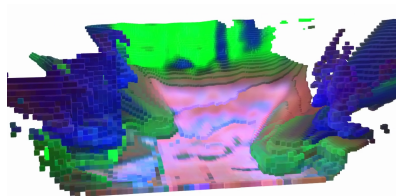
- Single input image  $I$
- 2D encoder-decoder  $\xi \rightarrow$  dense embedding  $E$
- Pixel embedding  $\rightarrow$  camera ray
- Implicit MLP head  $\phi(\mathbf{x}_i, \mathbf{e}_u) = (f_{\mathbf{x}}, \sigma_{\mathbf{x}})$
- Seg. head  $h$  predicts semantics  $p_{\mathbf{x}}$



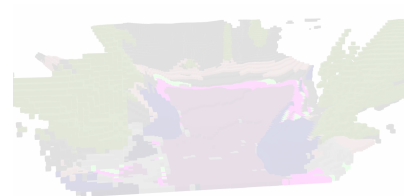
# SceneDINO Training



Single Input Image



3D Feature Field

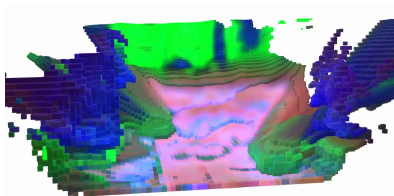


SSC Prediction

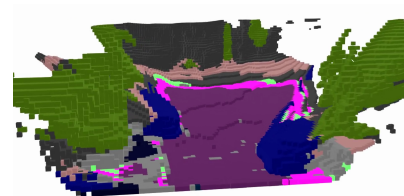
# SceneDINO Training



Single Input Image



3D Feature Field

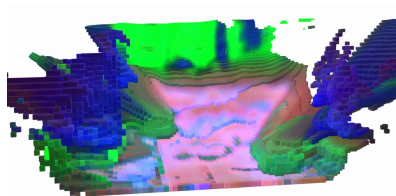


SSC Prediction

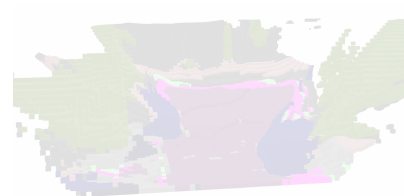
# SceneDINO Training



Single Input Image

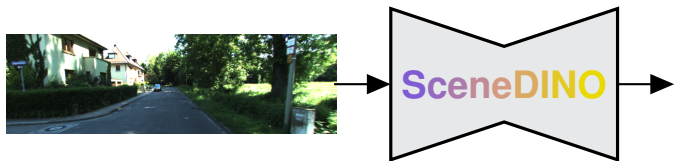


3D Feature Field



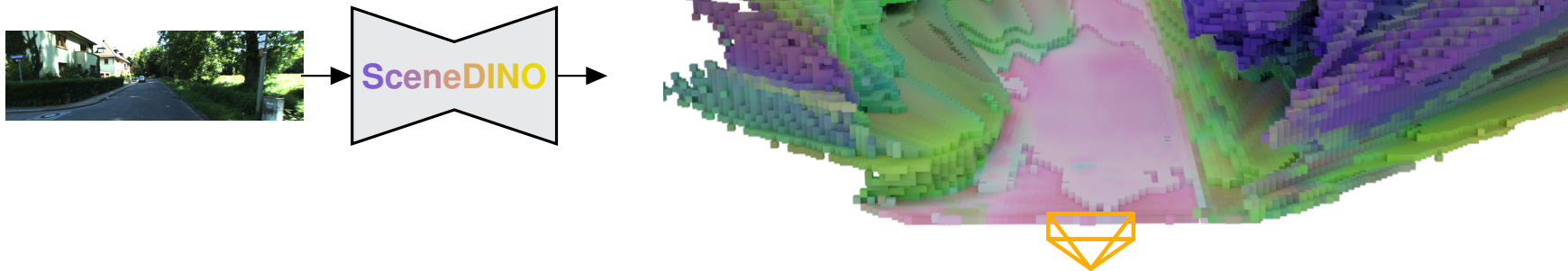
SSC Prediction

# Multi-View Self-Supervision

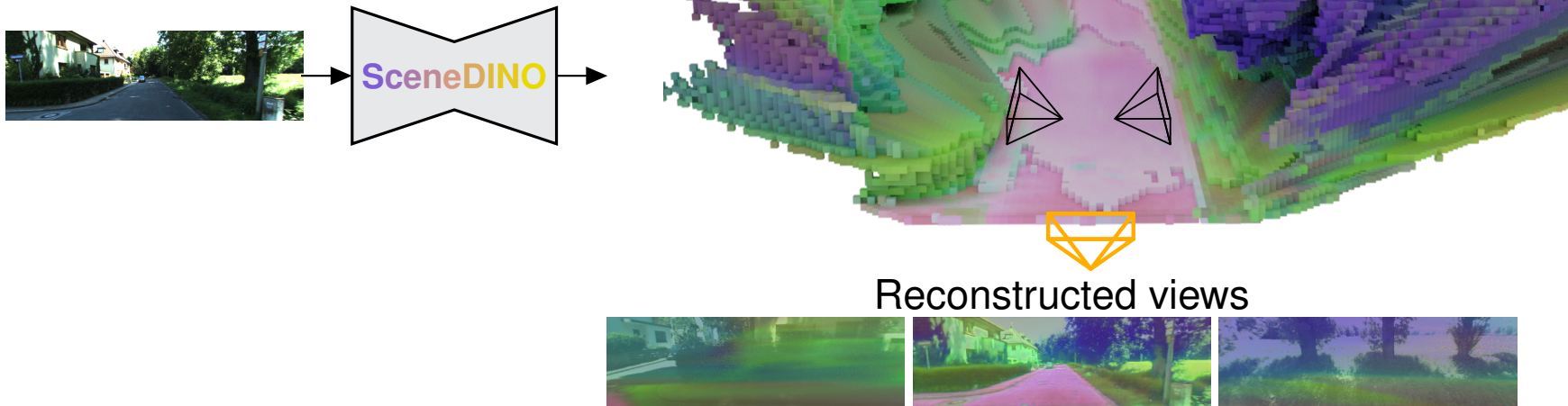




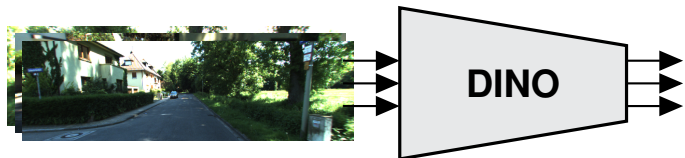
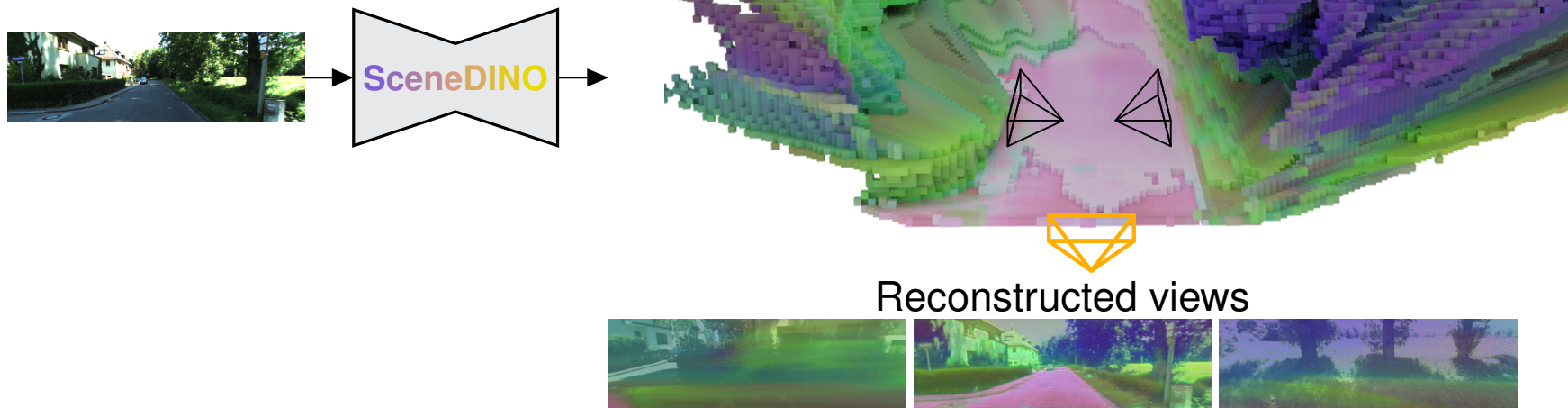
# Multi-View Self-Supervision



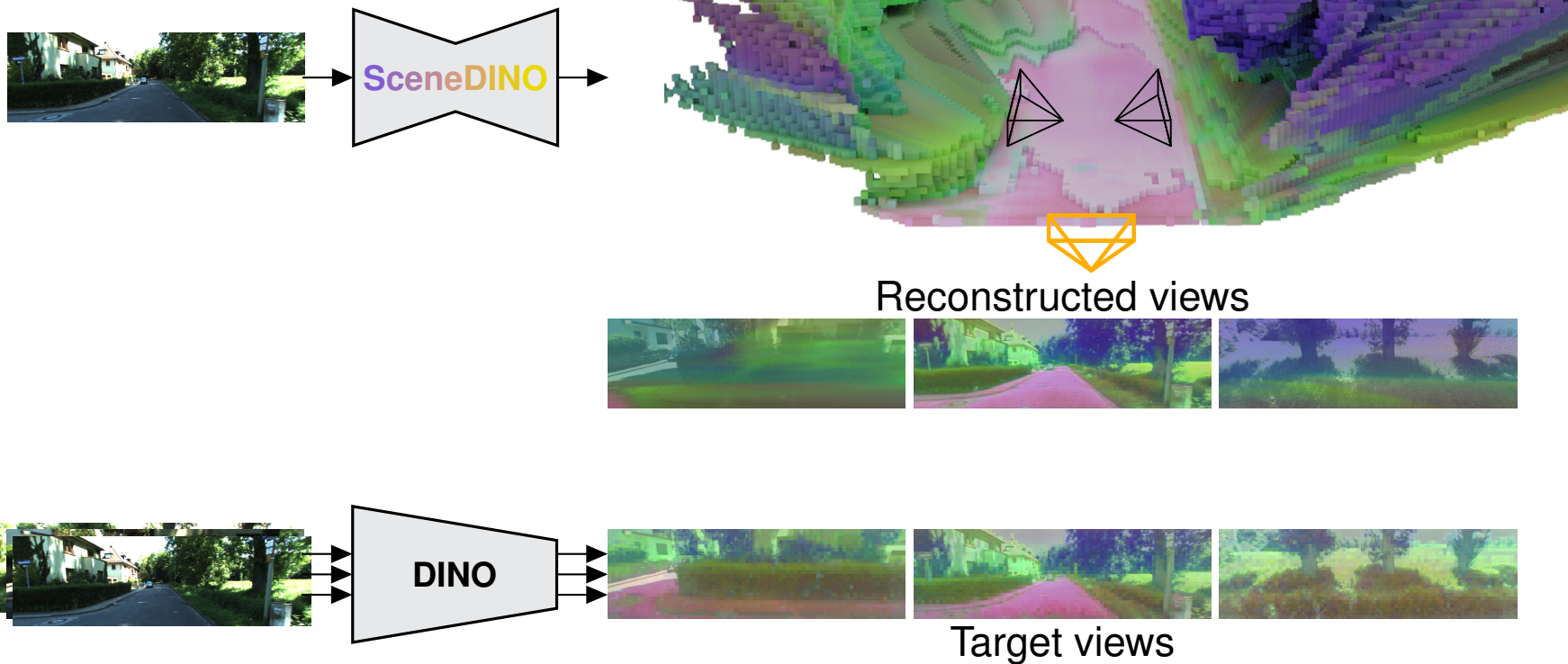
# Multi-View Self-Supervision



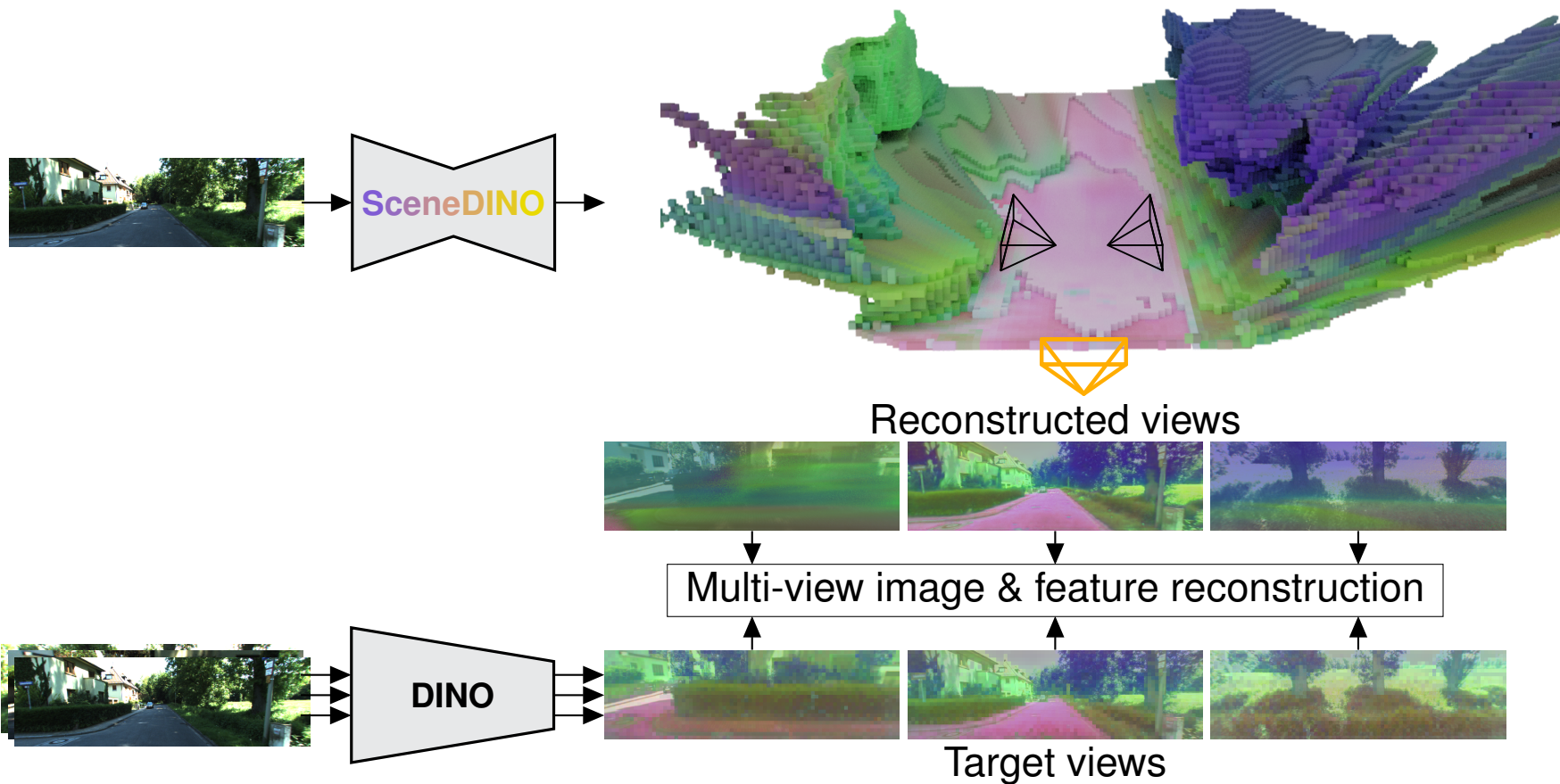
# Multi-View Self-Supervision



# Multi-View Self-Supervision



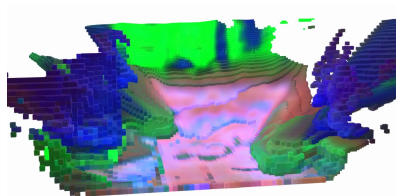
# Multi-View Self-Supervision



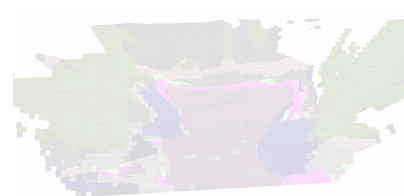
# SceneDINO Training



Single Input Image



3D Feature Field



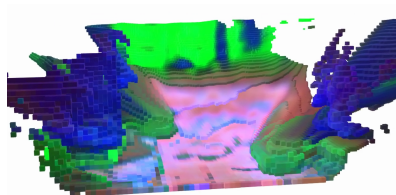
SSC Prediction



# SceneDINO Training

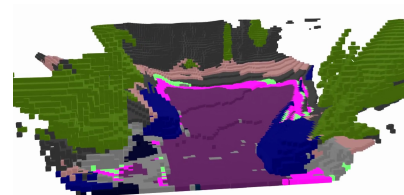


SceneDINO



3D Feature Field

Distill & Cluster



SSC Prediction

# Unsupervised Segmentation

- **Goal:** Learn unsupervised segmentation head

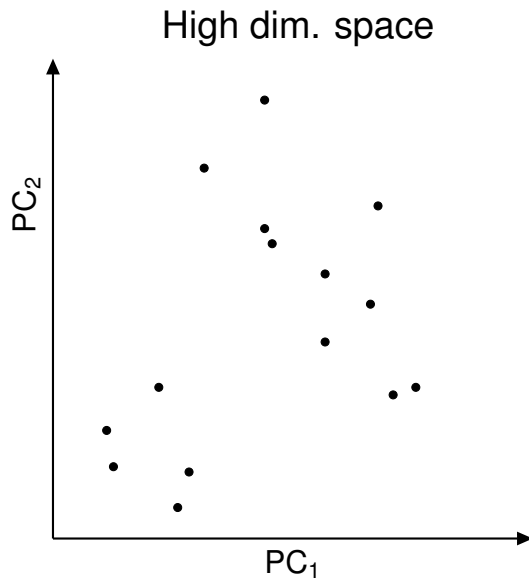


# Unsupervised Segmentation

- **Goal:** Learn unsupervised segmentation head
- **Idea:** Magnify semantic correspondence & cluster features

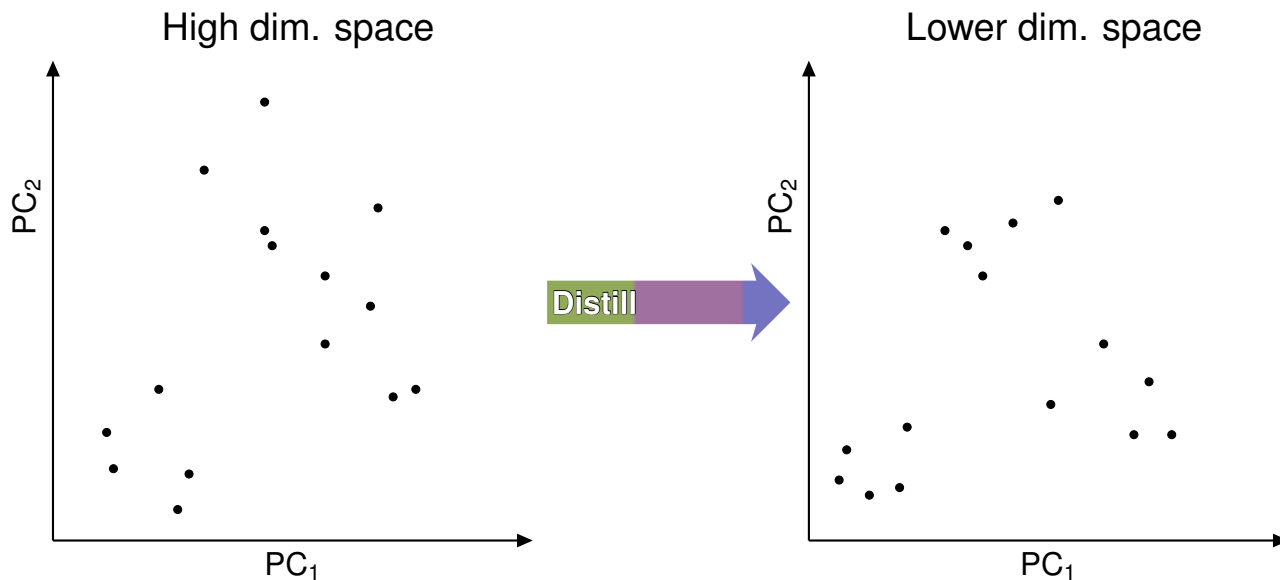
# Unsupervised Segmentation

- **Goal:** Learn unsupervised segmentation head
- **Idea:** Magnify semantic correspondence & cluster features



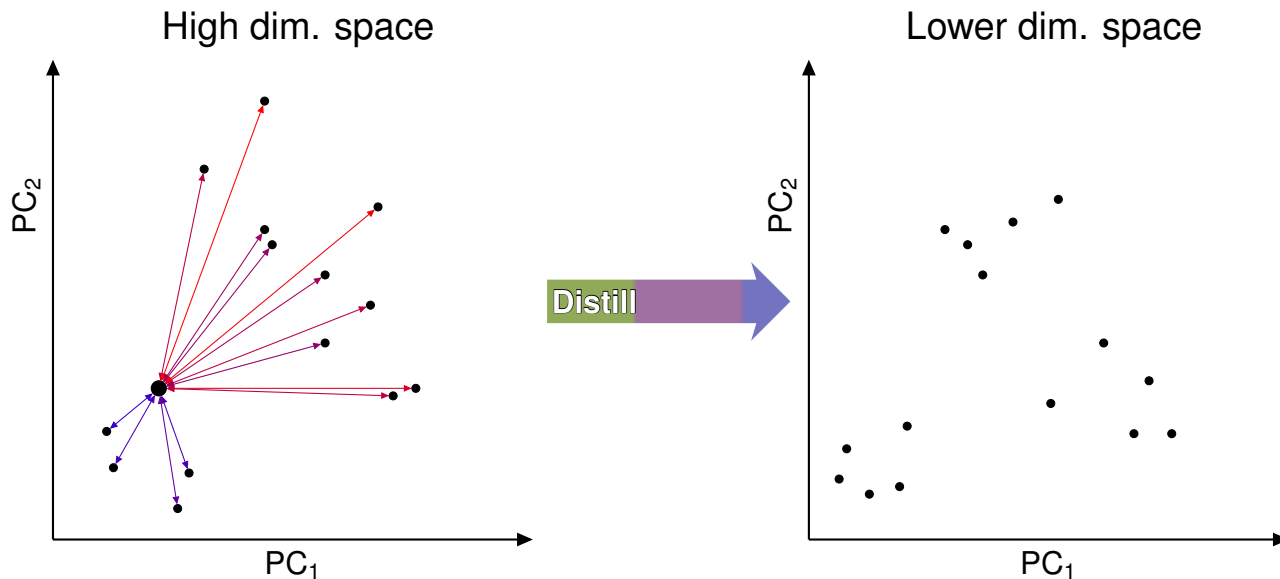
# Unsupervised Segmentation

- **Goal:** Learn unsupervised segmentation head
- **Idea:** Magnify semantic correspondence & cluster features



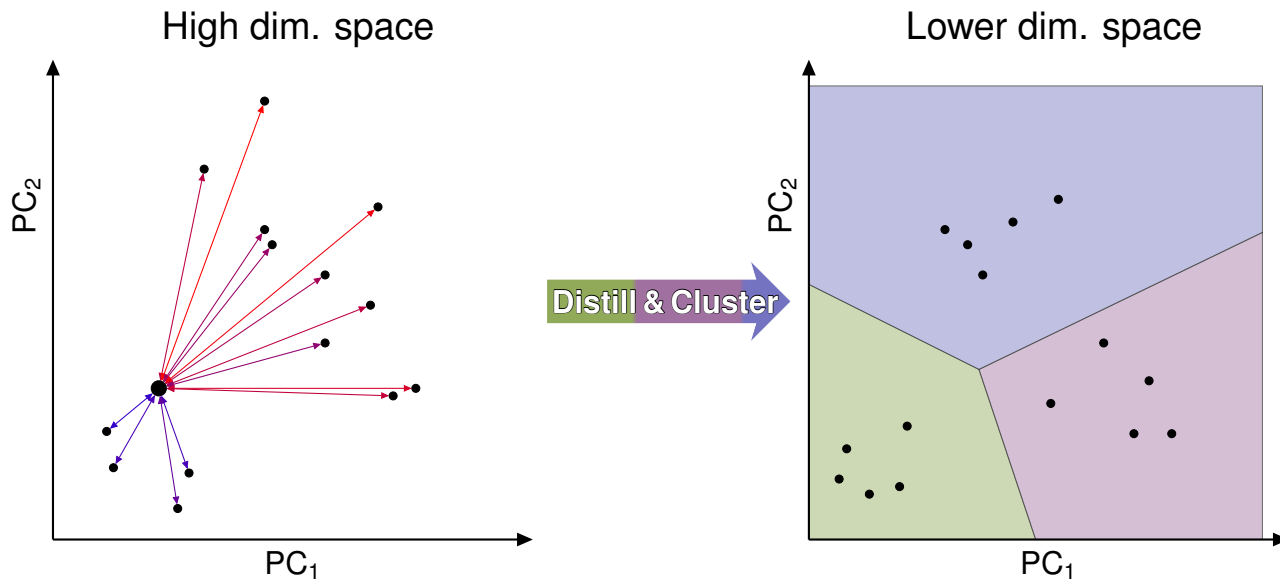
# Unsupervised Segmentation

- **Goal:** Learn unsupervised segmentation head
- **Idea:** Magnify semantic correspondence & cluster features



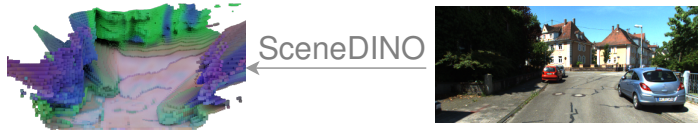
# Unsupervised Segmentation

- **Goal:** Learn unsupervised segmentation head
- **Idea:** Magnify semantic correspondence & cluster features



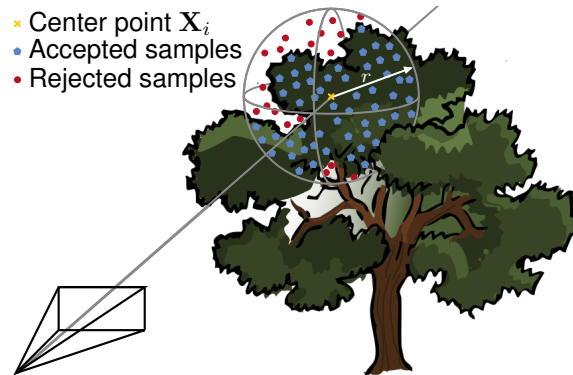
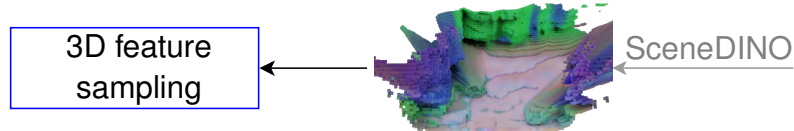
# Unsupervised Segmentation

- **Goal:** Learn unsupervised segmentation head
- **Idea:** Magnify semantic correspondence & cluster features



# Unsupervised Segmentation

- **Goal:** Learn unsupervised segmentation head
- **Idea:** Magnify semantic correspondence & cluster features



# Unsupervised Segmentation

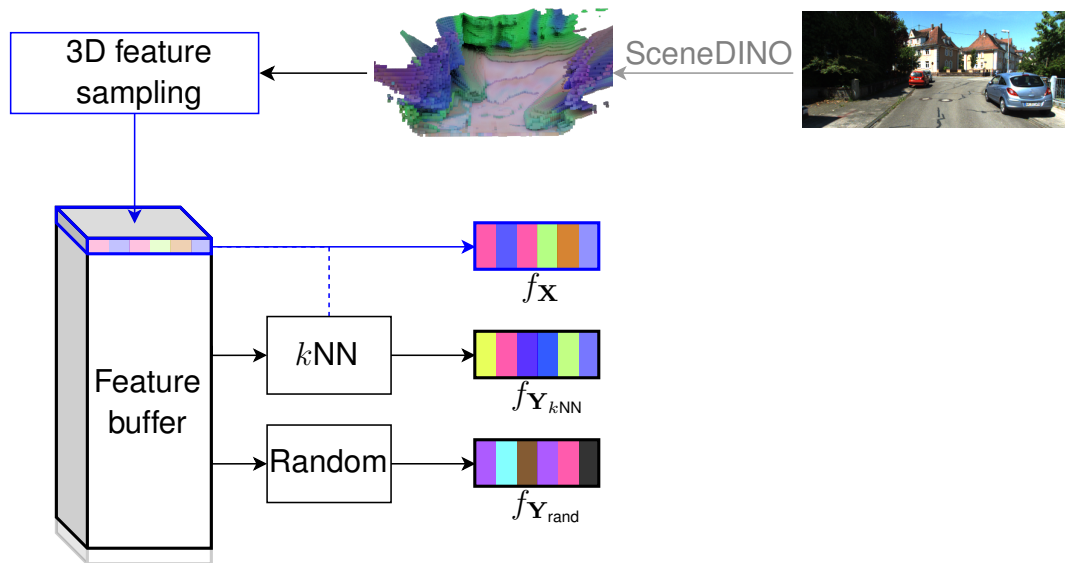
- **Goal:** Learn unsupervised segmentation head
- **Idea:** Magnify semantic correspondence & cluster features





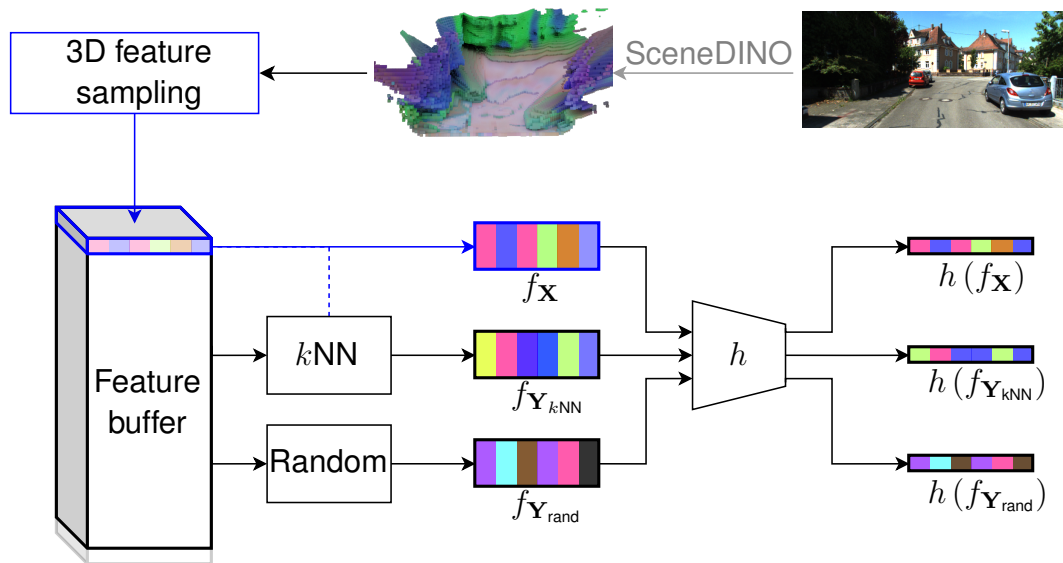
# Unsupervised Segmentation

- **Goal:** Learn unsupervised segmentation head
- **Idea:** Magnify semantic correspondence & cluster features



# Unsupervised Segmentation

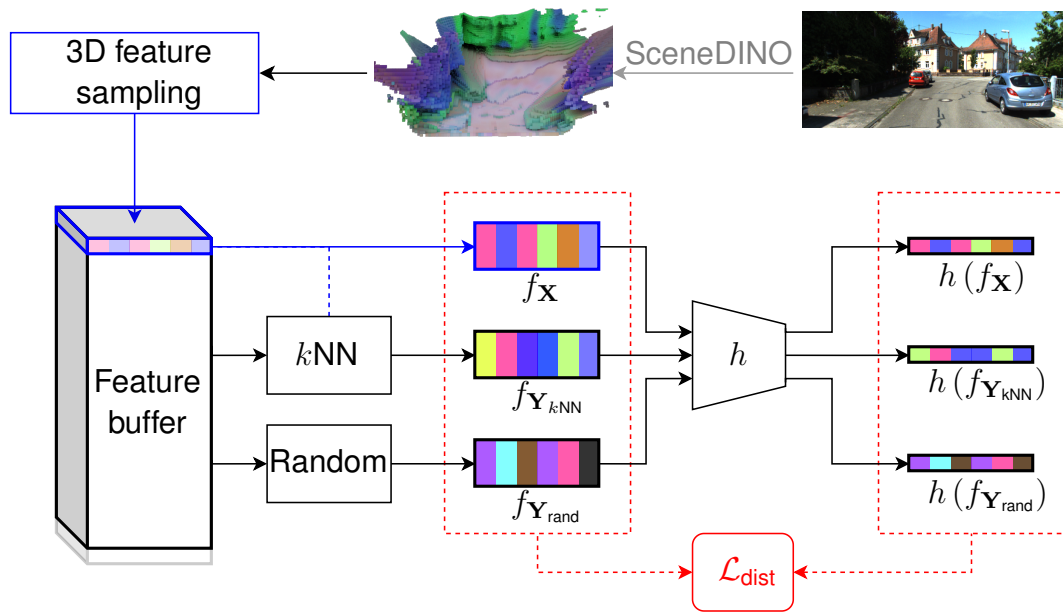
- **Goal:** Learn unsupervised segmentation head
- **Idea:** Magnify semantic correspondence & cluster features



- Head projects features down

# Unsupervised Segmentation

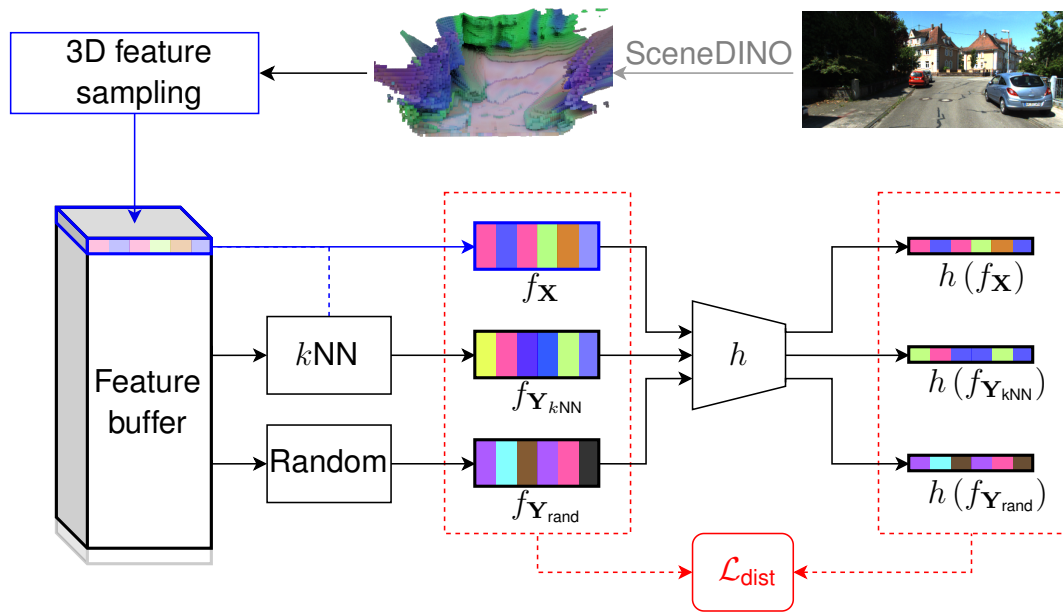
- **Goal:** Learn unsupervised segmentation head
- **Idea:** Magnify semantic correspondence & cluster features



- Head projects features down
- $\mathcal{L}_{dist}$  aligns correspondences

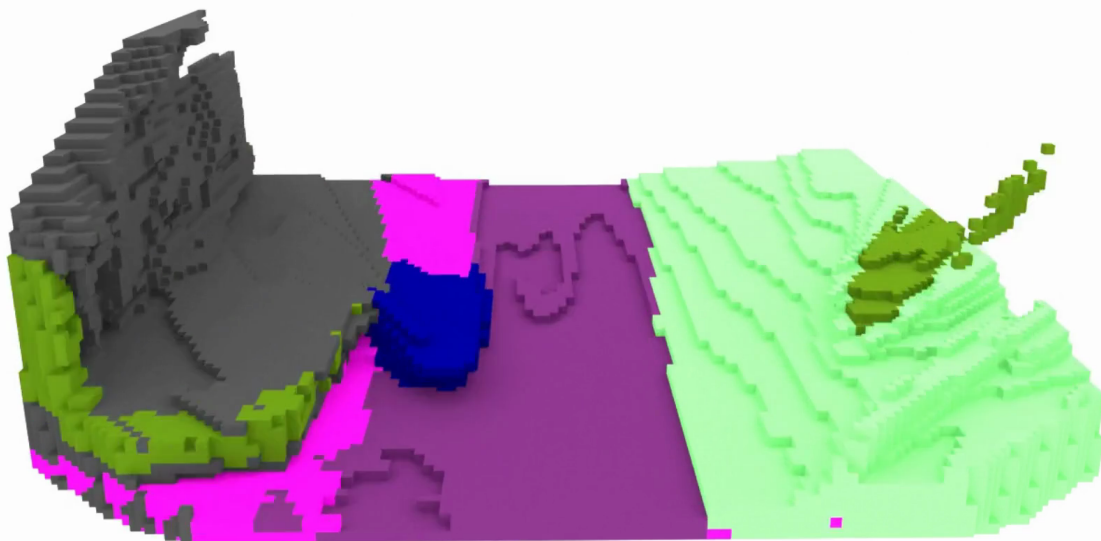
# Unsupervised Segmentation

- **Goal:** Learn unsupervised segmentation head
- **Idea:** Magnify semantic correspondence & cluster features



- Head projects features down
- $\mathcal{L}_{\text{dist}}$  aligns correspondences
$$f_{\mathbf{X}} \longleftrightarrow f_{\mathbf{X}}$$
$$f_{\mathbf{X}} \longleftrightarrow f_{\mathbf{Y}_{\text{kNN}}}$$
$$f_{\mathbf{X}} \longleftrightarrow f_{\mathbf{Y}_{\text{rand}}}$$
- $k$ -means cluster distilled features

# Results: SceneDINO



# Results: Unsupervised SSC

- KITTI-360-SSCBench experiments (full range 51.2 m validation)

Method	Unsupervised	Target features	mIoU (in %, $\uparrow$ )
S4C [3] (2D supervised)	✗	n/a	10.19
S4C [3] + STEGO [4]	✓	DINO	6.60

[3] A. Hayler *et al.*, “S4C: Self-supervised semantic scene completion with neural fields,” in *3DV*, 2024.

[4] M. Hamilton *et al.*, “Unsupervised semantic segmentation by distilling feature correspondences,” in *ICLR*, 2022.

# Results: Unsupervised SSC

- KITTI-360-SSCBench experiments (full range 51.2 m validation)

Method	Unsupervised	Target features	mIoU (in %, $\uparrow$ )
S4C [3] (2D supervised)	$\times$	n/a	10.19
S4C [3] + STEGO [4]	$\checkmark$	DINO	6.60
SceneDINO (Ours)	$\checkmark$	DINO	<b>8.00</b>

[3] A. Hayler *et al.*, “S4C: Self-supervised semantic scene completion with neural fields,” in *3DV*, 2024.

[4] M. Hamilton *et al.*, “Unsupervised semantic segmentation by distilling feature correspondences,” in *ICLR*, 2022.

# Results: Unsupervised SSC

- KITTI-360-SSCBench experiments (full range 51.2 m validation)

Method	Unsupervised	Target features	mIoU (in %, $\uparrow$ )
S4C [3] (2D supervised)	$\times$	n/a	10.19
S4C [3] + STEGO [4]	$\checkmark$	DINO	6.60
SceneDINO (Ours)	$\checkmark$	DINO	8.00
SceneDINO (Ours)	$\checkmark$	DINOv2	<b>9.08</b>

[3] A. Hayler *et al.*, “S4C: Self-supervised semantic scene completion with neural fields,” in *3DV*, 2024.

[4] M. Hamilton *et al.*, “Unsupervised semantic segmentation by distilling feature correspondences,” in *ICLR*, 2022.



# Results: Unsupervised SSC

- KITTI-360-SSCBench experiments (full range 51.2 m validation)

Method	Unsupervised	Target features	mIoU (in %, $\uparrow$ )
S4C [3] (2D supervised)	$\times$	n/a	10.19
S4C [3] + STEGO [4]	$\checkmark$	DINO	6.60
SceneDINO (Ours)	$\checkmark$	DINO	8.00
SceneDINO (Ours)	$\checkmark$	DINOv2	<b>9.08</b>

**State-of-the-art unsupervised semantic scene completion accuracy**

[3] A. Hayler *et al.*, “S4C: Self-supervised semantic scene completion with neural fields,” in *3DV*, 2024.

[4] M. Hamilton *et al.*, “Unsupervised semantic segmentation by distilling feature correspondences,” in *ICLR*, 2022.

# Results: Unsupervised SSC

- KITTI-360-SSCBench experiments (full range 51.2 m validation)

Method	Unsupervised	Target features	mIoU (in %, $\uparrow$ )
S4C [3] (2D supervised)	$\times$	n/a	10.19
S4C [3] + STEGO [4]	$\checkmark$	DINO	6.60
SceneDINO (Ours)	$\checkmark$	DINO	8.00
SceneDINO (Ours)	$\checkmark$	DINOv2	9.08
SceneDINO (Ours)	$\times$ (linear)	DINOv2	10.57

[3] A. Hayler *et al.*, “S4C: Self-supervised semantic scene completion with neural fields,” in *3DV*, 2024.

[4] M. Hamilton *et al.*, “Unsupervised semantic segmentation by distilling feature correspondences,” in *ICLR*, 2022.

# Results: Unsupervised SSC

- KITTI-360-SSCBench experiments (full range 51.2 m validation)

Method	Unsupervised	Target features	mIoU (in %, $\uparrow$ )
S4C [3] (2D supervised)	✗	n/a	10.19
S4C [3] + STEGO [4]	✓	DINO	6.60
SceneDINO (Ours)	✓	DINO	8.00
SceneDINO (Ours)	✓	DINOv2	9.08
SceneDINO (Ours)	✗ (linear)	DINOv2	10.57

**Linear probing outperforms 2D supervised S4C**

[3] A. Hayler *et al.*, “S4C: Self-supervised semantic scene completion with neural fields,” in *3DV*, 2024.

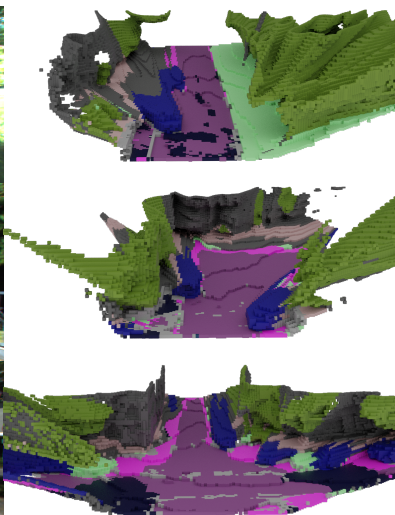
[4] M. Hamilton *et al.*, “Unsupervised semantic segmentation by distilling feature correspondences,” in *ICLR*, 2022.

# Results: Unsupervised SSC

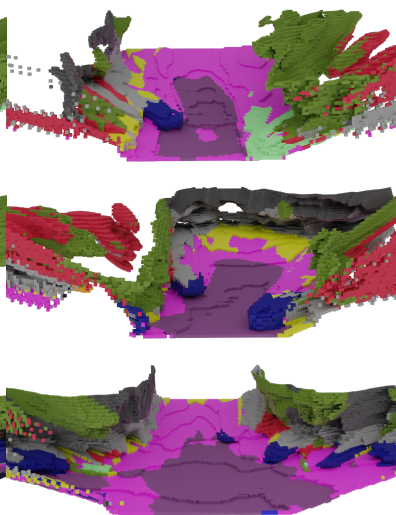
Input Image



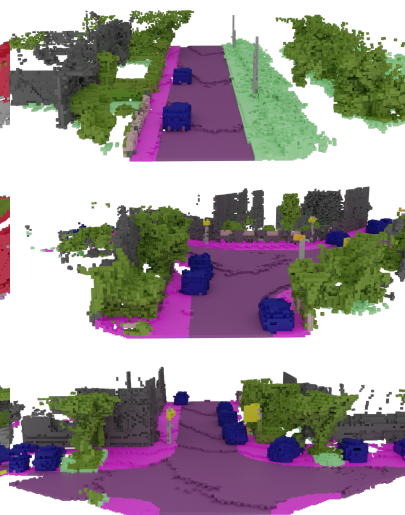
SceneDINO



S4C + STEGO



Ground Truth



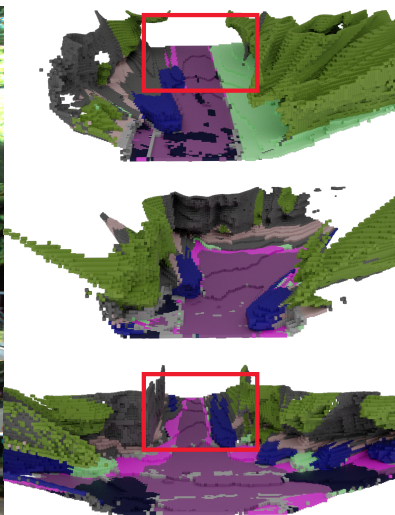
Road Sidewalk Building Fence Pole Other Object Traffic Sign Vegetation Terrain Person Car Other Vehicle Motorcycle Bicycle

# Results: Unsupervised SSC

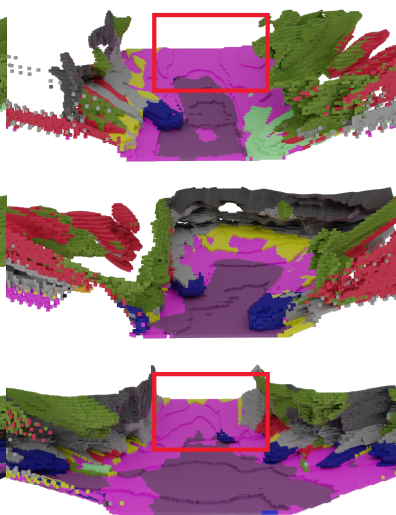
Input Image



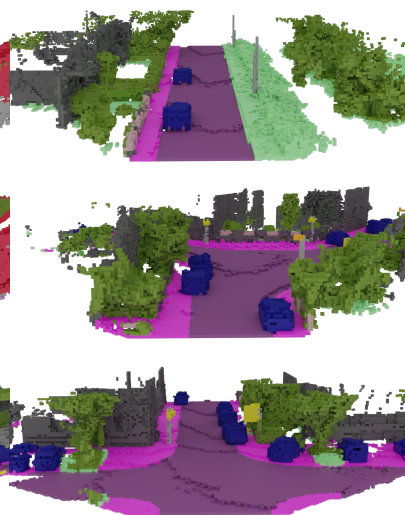
SceneDINO



S4C + STEGO



Ground Truth



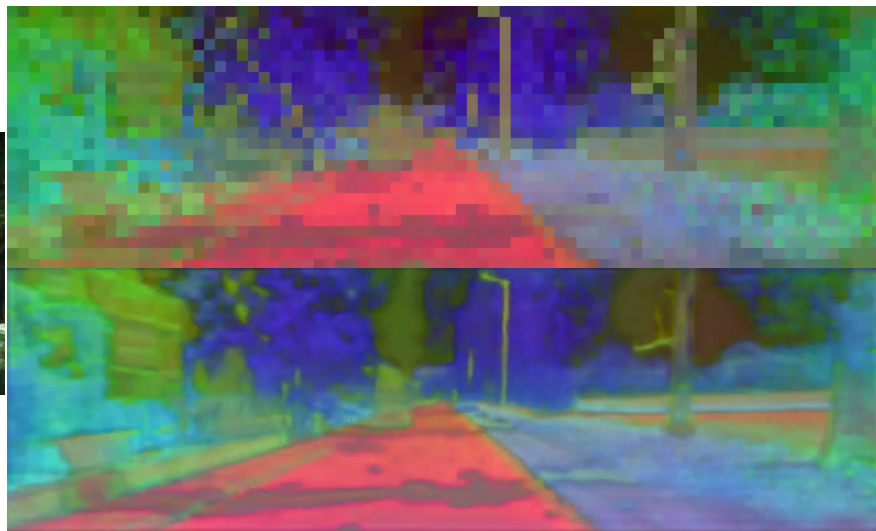
Road Sidewalk Building Fence Pole Other Object Traffic Sign Vegetation Terrain Person Car Other Vehicle Motorcycle Bicycle

# Results: SceneDINO in 2D

Input Image



DINO



SceneDINO

# Results: Multi-View Feature Consistency

- Multi-view consistency results using optical flow alignment

Method	KITTI-360		RealEstate10K	
	$L_1$ ( $\downarrow$ )	cos-sim ( $\uparrow$ )	$L_1$ ( $\downarrow$ )	cos-sim ( $\uparrow$ )
DINO [5]	16.06	0.70	14.41	0.75
SceneDINO (w/ DINO)	<b>6.45</b>	<b>0.93</b>	<b>5.87</b>	<b>0.95</b>
DINOv2 [6]	15.83	0.70	14.20	0.75
FiT3D [7]	7.02	0.93	5.67	0.95
SceneDINO (w/ DINOv2)	<b>5.24</b>	<b>0.96</b>	<b>4.87</b>	<b>0.97</b>

- [5] M. Caron *et al.*, “Emerging properties in self-supervised vision transformers,” in *ICCV*, 2021.  
[6] M. Oquab *et al.*, “DINOv2: Learning robust visual features without supervision,” *TMLR*, 2024.  
[7] Y. Yue *et al.*, “Improving 2D feature representations by 3D-aware fine-tuning,” in *ECCV*, 2024.



# Results: Multi-View Feature Consistency

- Multi-view consistency results using optical flow alignment

Method	KITTI-360		RealEstate10K	
	$L_1$ ( $\downarrow$ )	cos-sim ( $\uparrow$ )	$L_1$ ( $\downarrow$ )	cos-sim ( $\uparrow$ )
DINO [5]	16.06	0.70	14.41	0.75
SceneDINO (w/ DINO)	<b>6.45</b>	<b>0.93</b>	<b>5.87</b>	<b>0.95</b>
DINOv2 [6]	15.83	0.70	14.20	0.75
FiT3D [7]	7.02	0.93	5.67	0.95
SceneDINO (w/ DINOv2)	<b>5.24</b>	<b>0.96</b>	<b>4.87</b>	<b>0.97</b>

**SceneDINO's features are significantly more multi-view consistent**

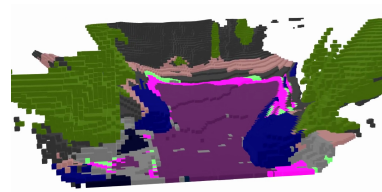
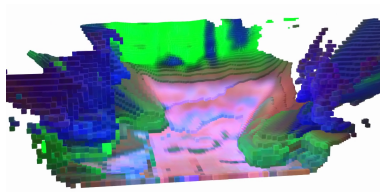
- [5] M. Caron *et al.*, “Emerging properties in self-supervised vision transformers,” in *ICCV*, 2021.  
[6] M. Oquab *et al.*, “DINOv2: Learning robust visual features without supervision,” *TMLR*, 2024.  
[7] Y. Yue *et al.*, “Improving 2D feature representations by 3D-aware fine-tuning,” in *ECCV*, 2024.



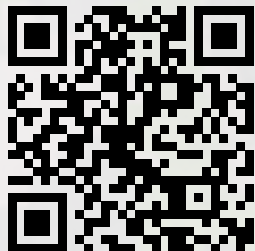
# Conclusion

We presented **SceneDINO** for unsupervised SSC

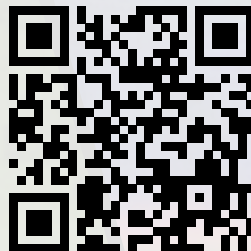
- **Multi-view self-supervision** effective for 3D scene understanding
- Single image  $\rightarrow$  **3D geometry & expressive features**
- Distilling & clustering leads to **SoTA accuracy** in unsupervised SSC
- Strong **linear probing, multi-view consistency**, and **domain generalization**



### Paper



### Project Page



### Code & Weights



<https://visinf.github.io/scenedino/>

This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No. 866008), the ERC Advanced Grant SIMULACRON, the DFG project CR 250/26-1 "4D-YouTube", the GNI Project "AICC", and the State of Hesse within the LOEWE emergenCITY center. This work was partially supported by the Deutsche Forschungsgemeinschaft (German Research Foundation, DFG) under Germany's Excellence Strategy (EXC 3066/1 "The Adaptive Mind", Project No. 533717223). C. Reich is supported by the Konrad Zuse School of Excellence in Learning and Intelligent Systems. C. Rupprecht is supported by an Amazon Research Award.





European Research Council  
Established by the European Commission





# Results: SceneDINO Analysis

- Analysing camera poses and target features

$\Delta$ mIoU	mIoU	Configuration
-0.12 	7.88	w/ estimated ORB-SLAM3 poses
—	8.00	Full framework (SceneDINO)
 +1.08	9.08	DINOv2 target features (vs. DINO)

# Results: SceneDINO Analysis



- Analysing camera poses and target features

$\Delta$ mIoU	mIoU	Configuration
-0.12 	7.88	w/ estimated ORB-SLAM3 poses
—	8.00	Full framework (SceneDINO)
 +1.08	9.08	DINOv2 target features (vs. DINO)

**SceneDINO can benefit from better target features**

# Results: SceneDINO Analysis

- Analysing camera poses and target features

$\Delta$ mIoU	mIoU	Configuration
-0.12 	7.88	w/ estimated ORB-SLAM3 poses
—	8.00	Full framework (SceneDINO)
 +1.08	9.08	DINOv2 target features (vs. DINO)



- Linear probing features (w/ 2D sem. GT)

Probing approach	Target features	mIoU
Linear	DINO	9.34
	DINOv2	<b>10.57</b>
S4C (full training)	n/a	10.19

**SceneDINO can benefit from better target features**

# Results: SceneDINO Analysis

- Analysing camera poses and target features

$\Delta$ mIoU	mIoU	Configuration
-0.12 	7.88	w/ estimated ORB-SLAM3 poses
—	8.00	Full framework (SceneDINO)
 +1.08	9.08	DINOv2 target features (vs. DINO)

**SceneDINO can benefit from better target features**

- Linear probing features (w/ 2D sem. GT)

Probing approach	Target features	mIoU
Linear	DINO	9.34
	DINOv2	<b>10.57</b>
S4C (full training)	n/a	10.19

**Linear probing outperforms 2D supervised S4C**