

# Scalable Cross-Modal 3D Scene Understanding

Sayan Deb Sarkar

**Microsoft Spatial AI Lab**  
July 18, 2025

# Who Am I?

## PhD student since 2024.09

- Advisor: Prof. Iro Armeni
- Gradient Spaces Research Group, part of Stanford Vision and Learning Lab (SVL)



## Computer Science MSc 2022.09 - 2024.08

- Advisor: Prof. Marc Pollefeys
- Computer Vision And Geometry Group (CVG)



Computer Vision  
and Geometry Lab



## At Industry

- Internships at Microsoft Spatial AI Lab & Qualcomm XR
- Computer Vision Engineer at Mercedes Benz R & D



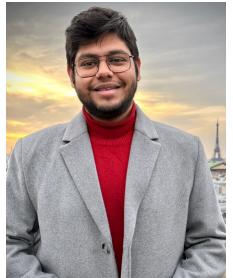
[sayands.github.io](https://sayands.github.io)



★★★ Highlight  
 **CrossOver**

# 3D Scene Cross-Modal Alignment

Sayan Deb Sarkar



Ondrej Miksik



Marc Pollefeys



Dániel Béla Baráth



Iro Armeni



**ETH** zürich

 Microsoft

  
Computer Vision  
and Geometry Lab

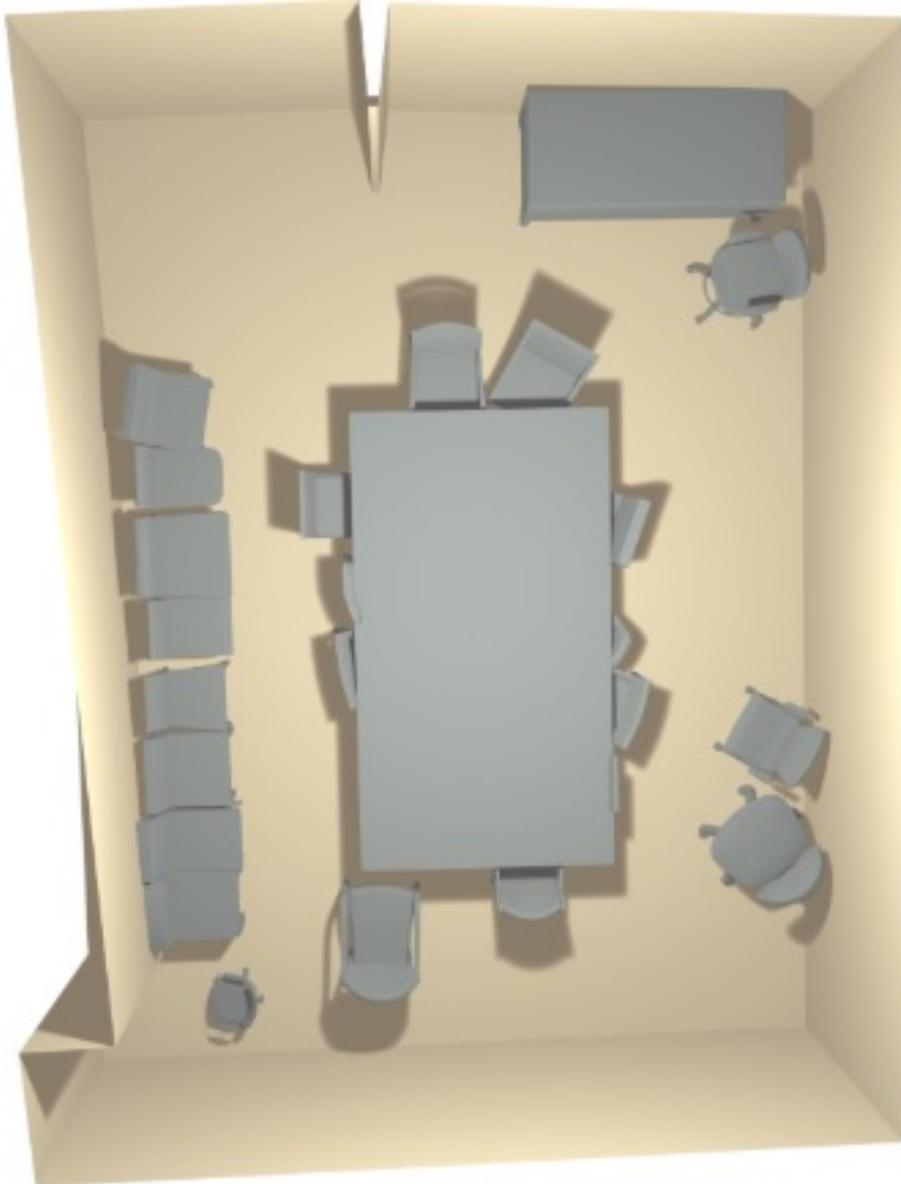
**Stanford**  
University

  
Gradient  
Spaces

**HUN**  
REN

 SZTAKI

# Application

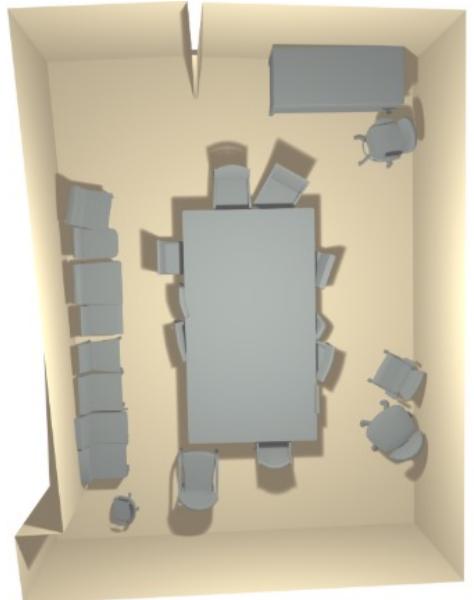


Virtual CAD Scene



Scene Database

# Application

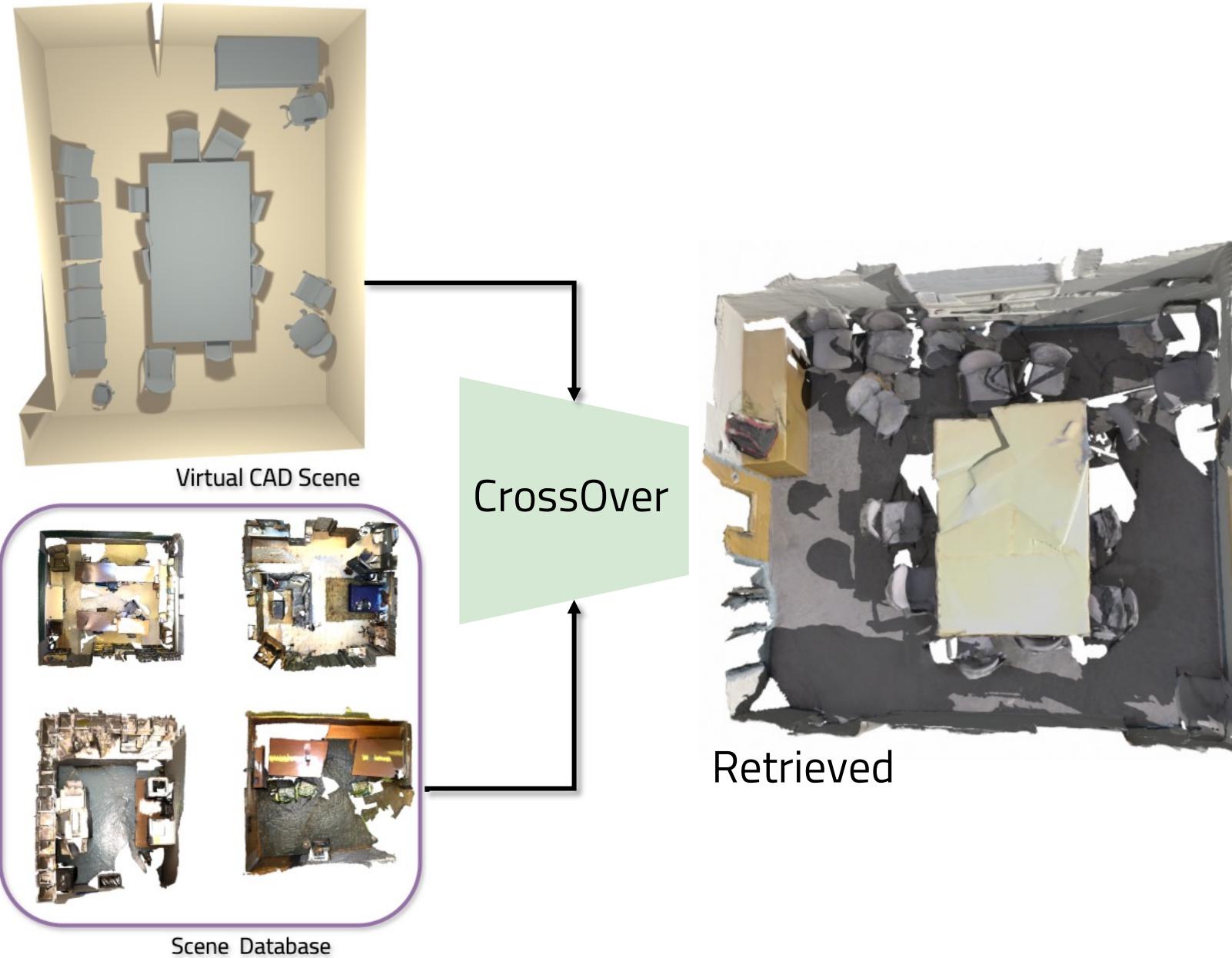


Virtual CAD Scene

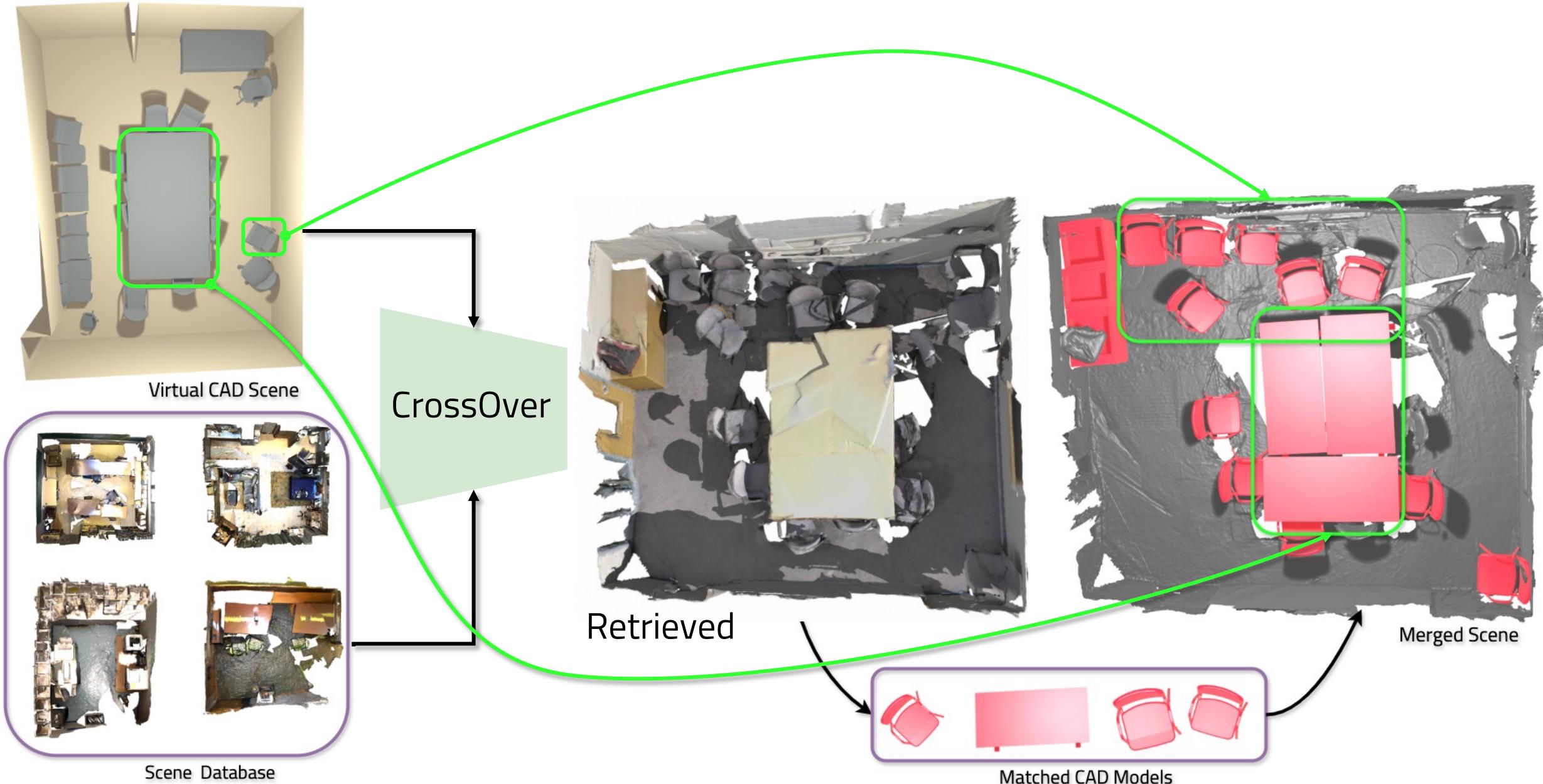


Scene Database

# Application



# Application

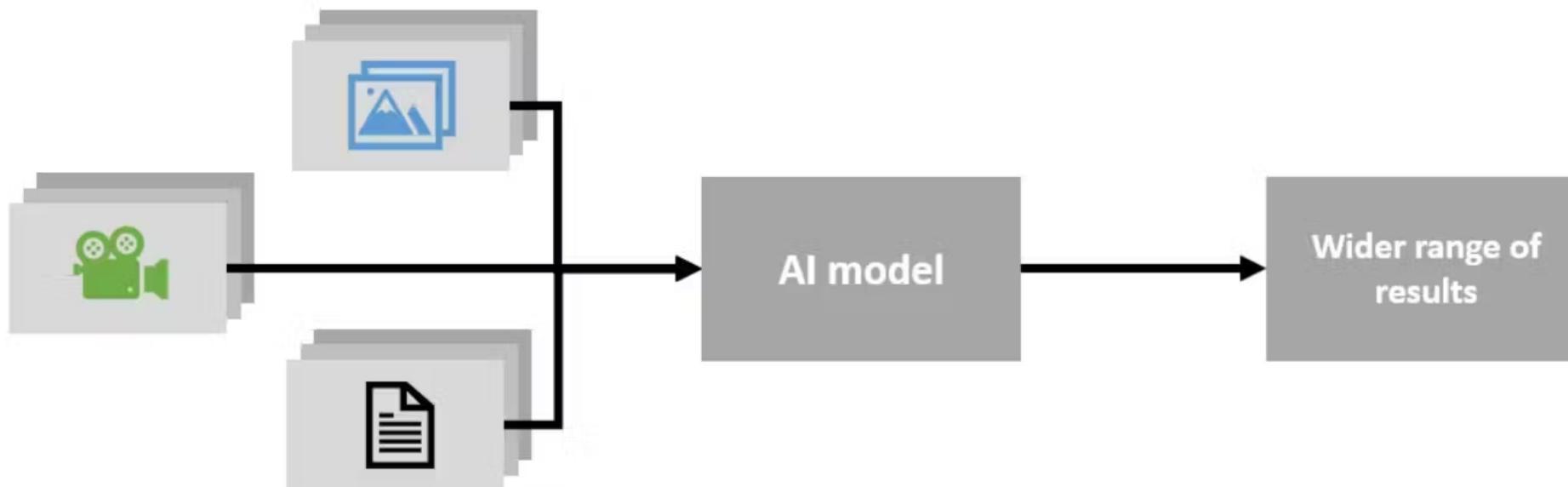


# Uni-modal vs Multi-modal

Unimodal AI model

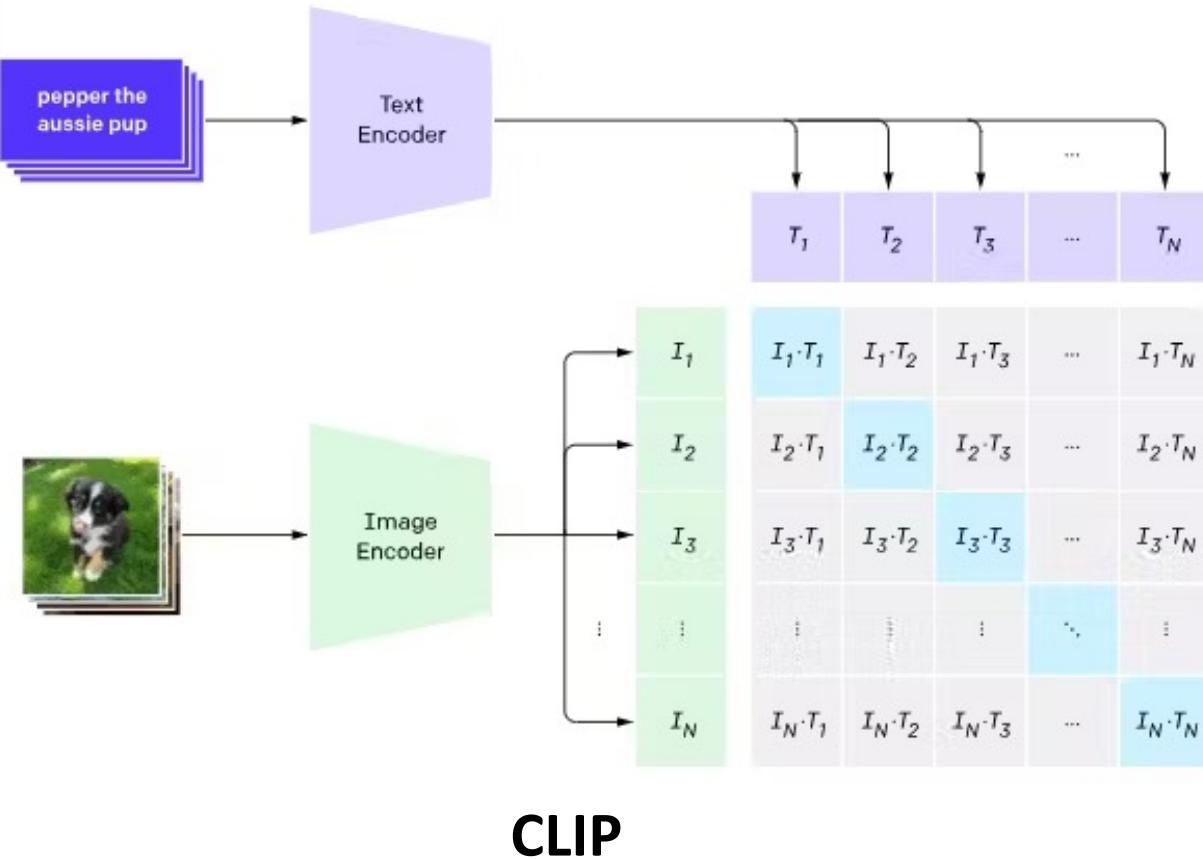


Multimodal AI model

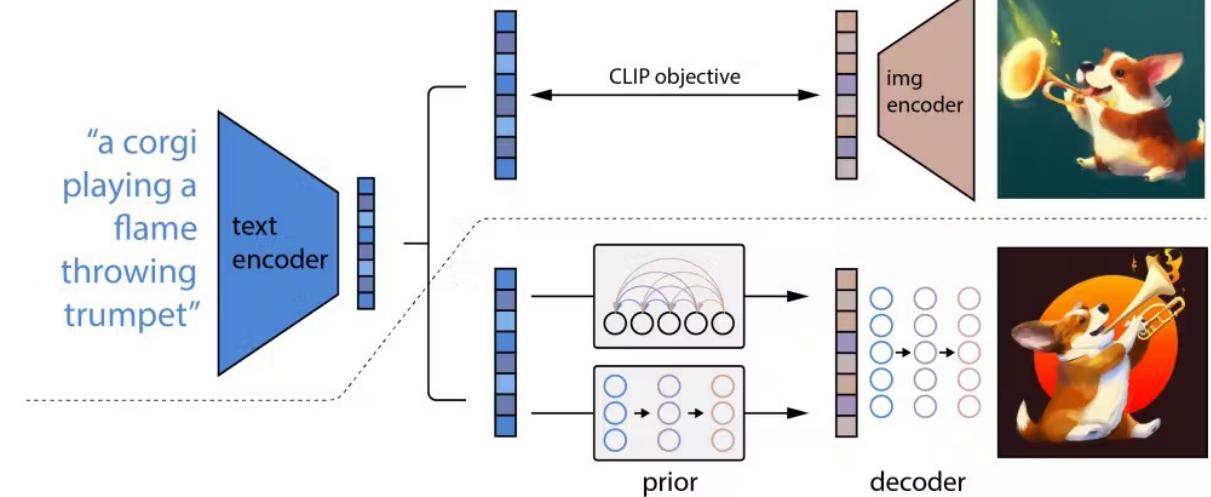


They enable knowledge and experience transfer across modalities.

# Multimodal Models Today in 2D



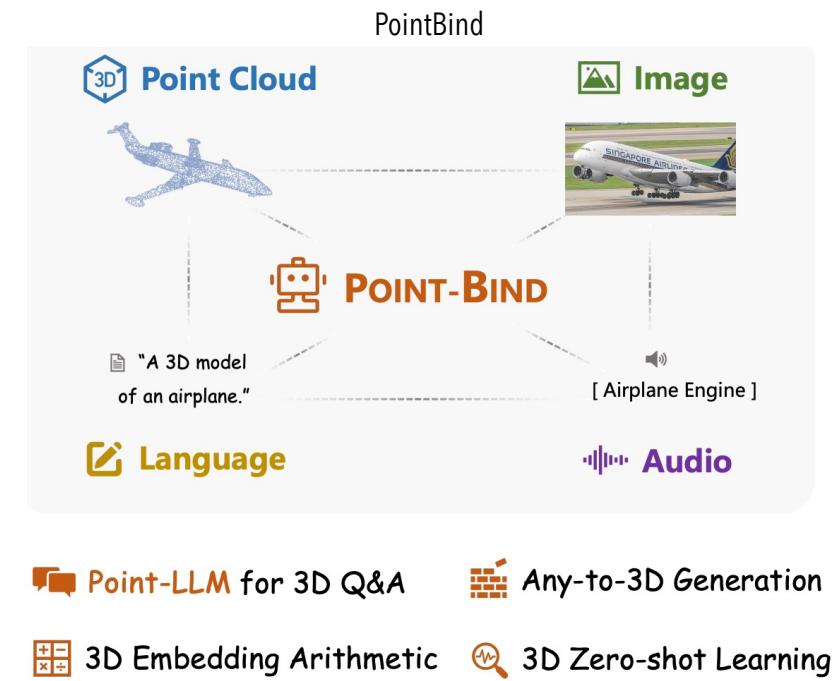
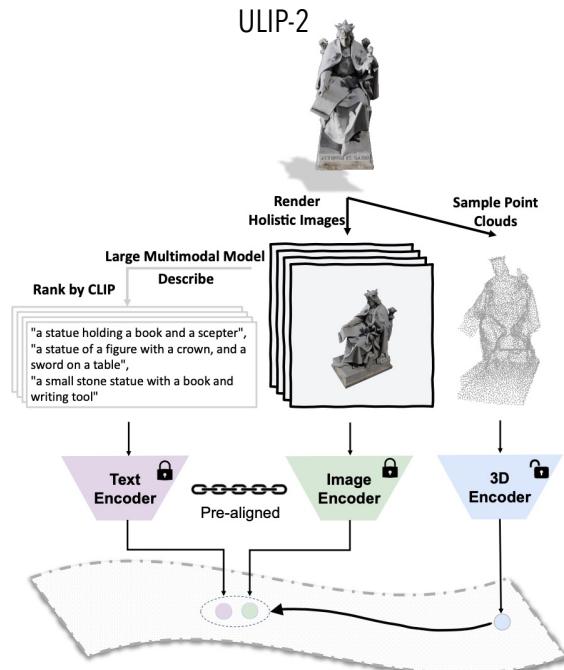
**CLIP**



**DALL-E**

Main focus is on dual-modality, connecting images with text.

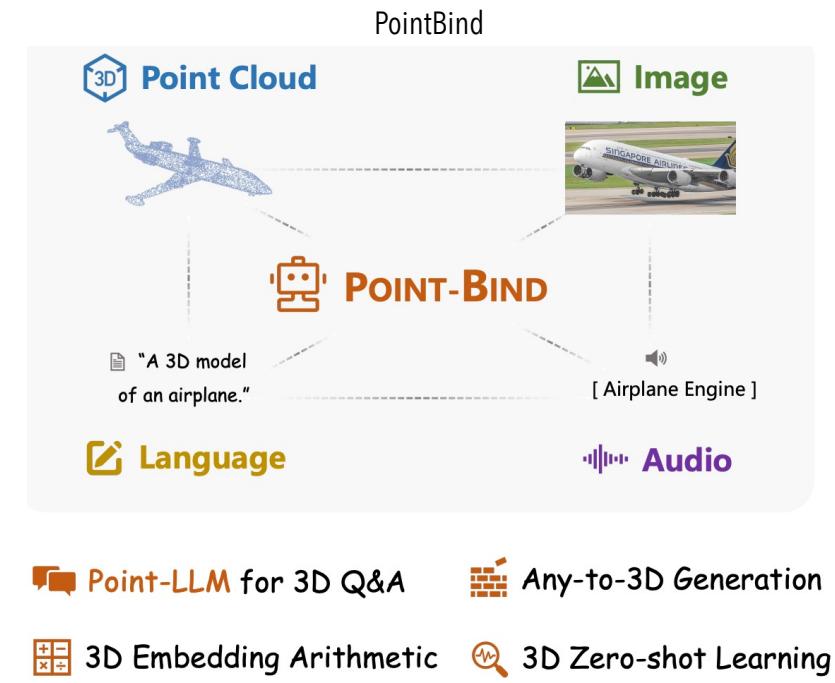
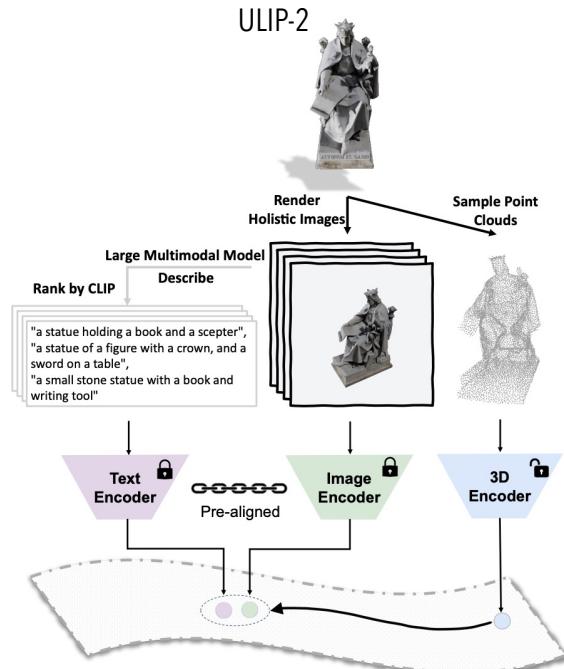
# Motivation



[1] Xue et al, *ULIP-2: Towards Scalable Multimodal Pre-training for 3D Understanding*, CVPR 2024

[2] Guo et al, *Point-Bind & Point-LLM: Aligning 3D with Multi-modality*, arXiv 2024

# Motivation

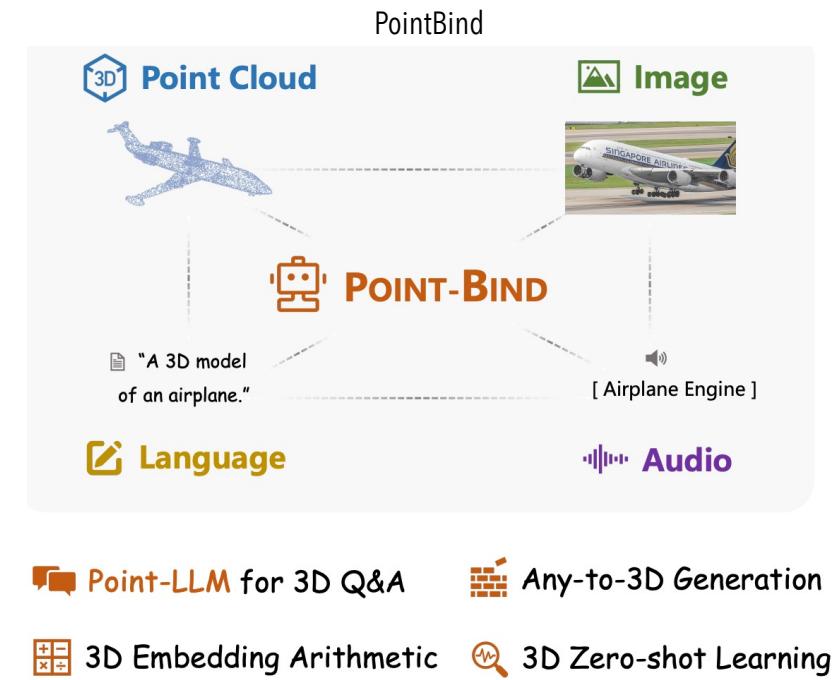
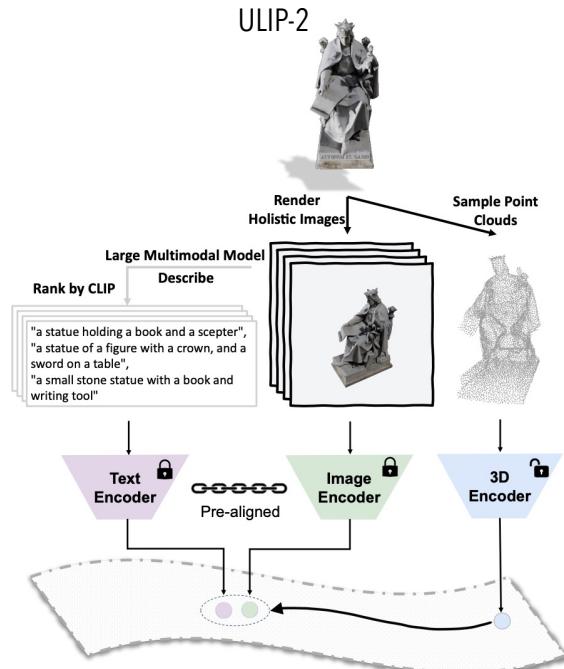


- Perfect Modality Alignments + Complete Object Data

[1] Xue et al, *ULIP-2: Towards Scalable Multimodal Pre-training for 3D Understanding*, CVPR 2024

[2] Guo et al, *Point-Bind & Point-LLM: Aligning 3D with Multi-modality*, arXiv 2024

# Motivation

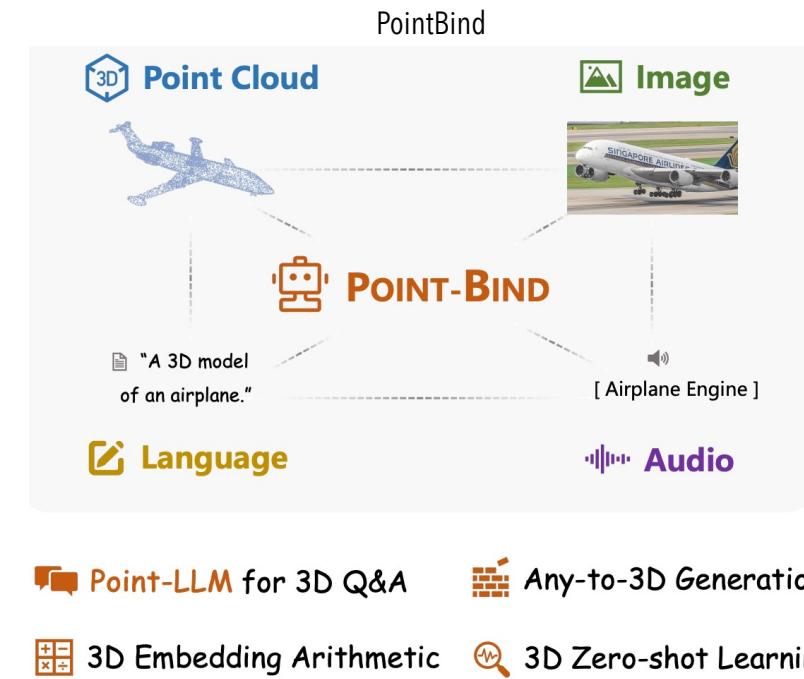
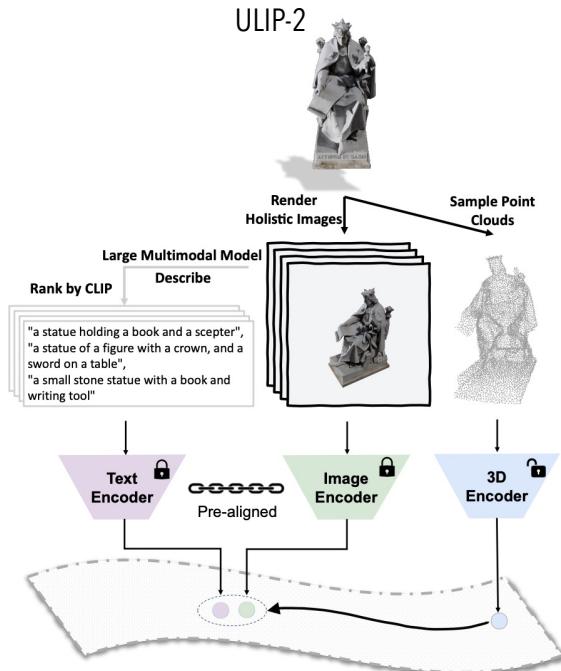


- Perfect Modality Alignments + Complete Object Data
- No Context Within A Scene

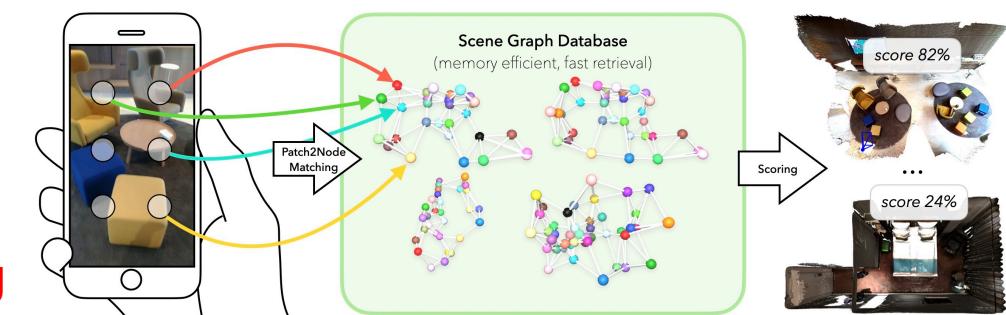
[1] Xue et al, *ULIP-2: Towards Scalable Multimodal Pre-training for 3D Understanding*, CVPR 2024

[2] Guo et al, *Point-Bind & Point-LLM: Aligning 3D with Multi-modality*, arXiv 2024

# Motivation



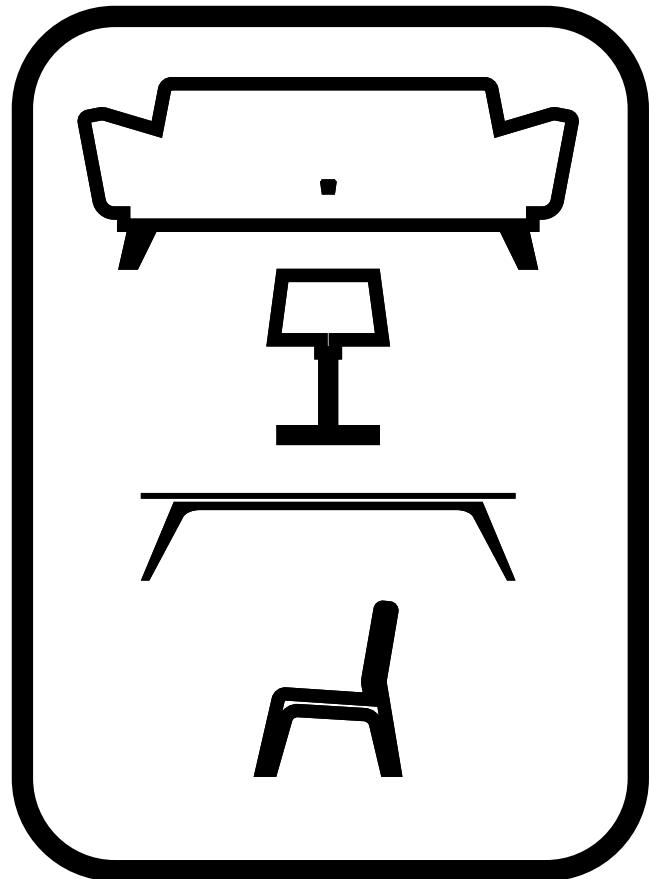
- Perfect Modality Alignments + Complete Object Data
- No Context Within A Scene
- Need For Semantic Annotations/Explicit Scene Graphs for 3D Scene Understanding



[1] Xue et al, *ULIP-2: Towards Scalable Multimodal Pre-training for 3D Understanding*, CVPR 2024

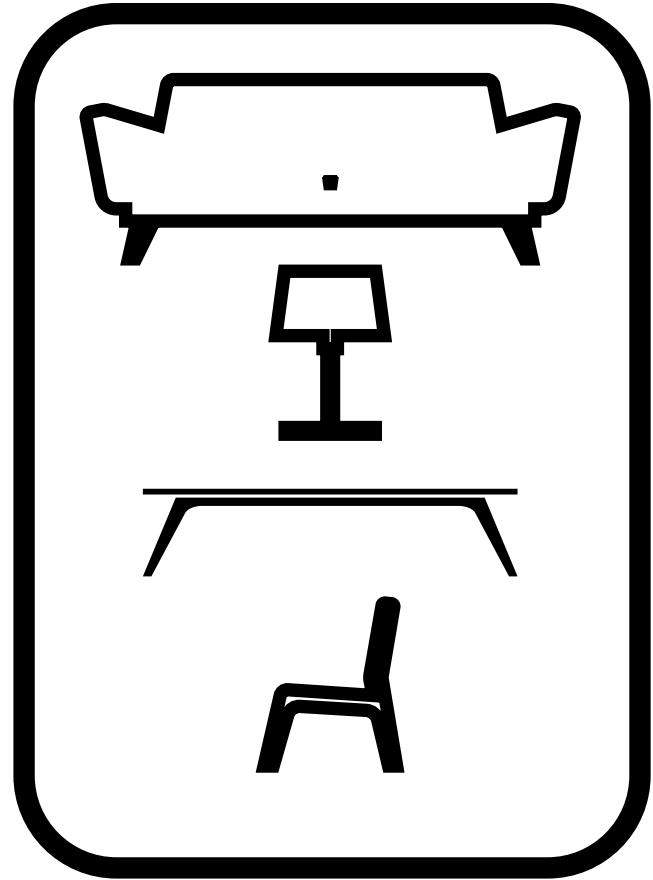
[2] Guo et al, *Point-Bind & Point-LLM: Aligning 3D with Multi-modality*, arXiv 2024

[2] Miao et al, *SceneGraphLoc: Cross-Modal Coarse Visual Localization on 3D Scene Graphs*, ECCV 2024

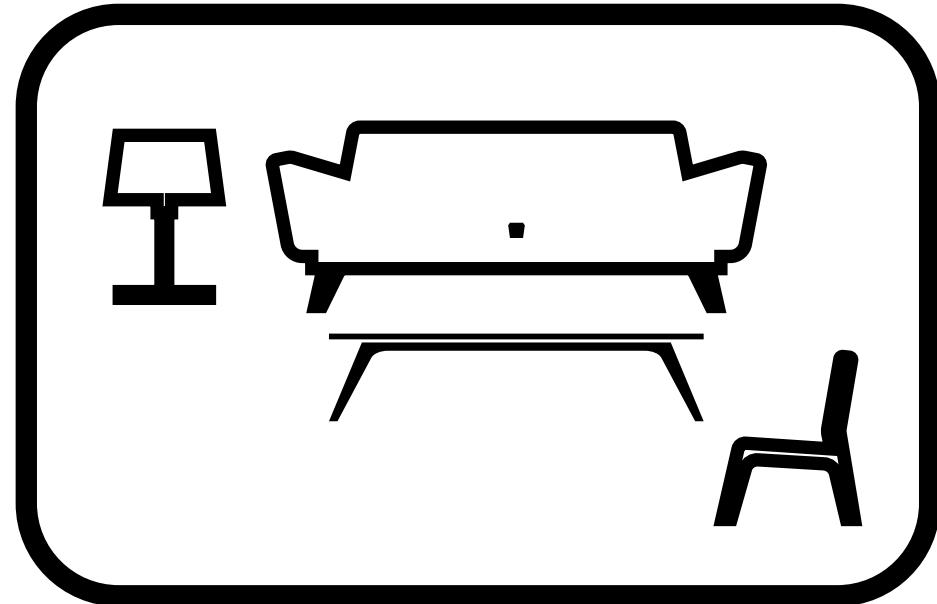


## I. Instance-Level Alignment

*to leverage large pre-trained models*

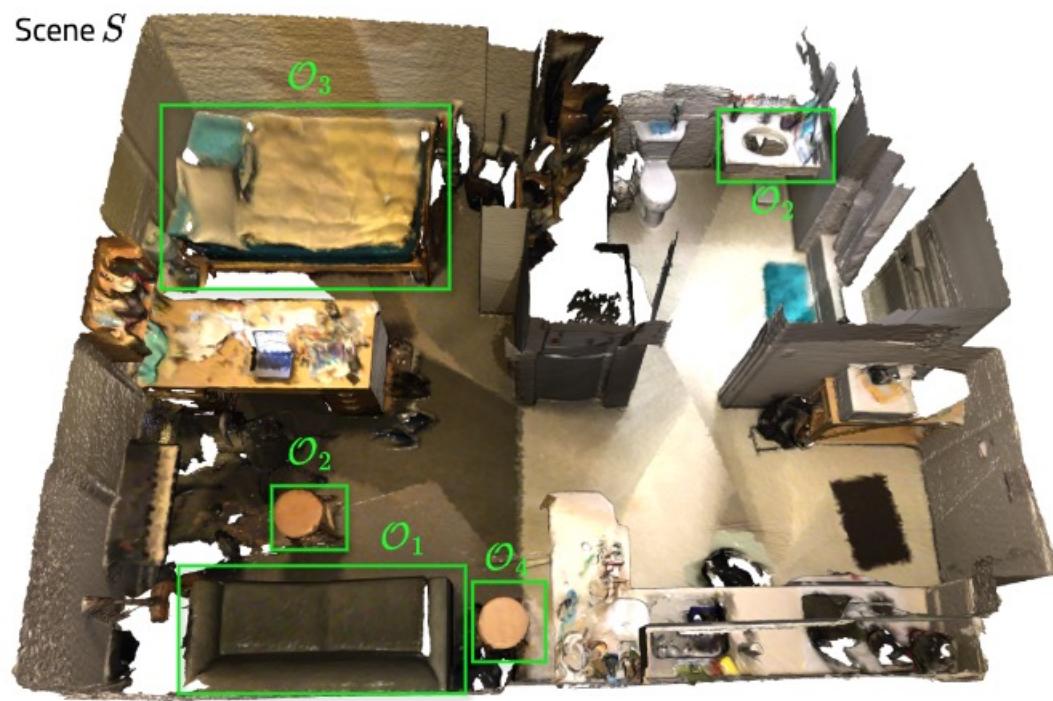


**I. Instance-Level Alignment**  
*to leverage large pre-trained models*

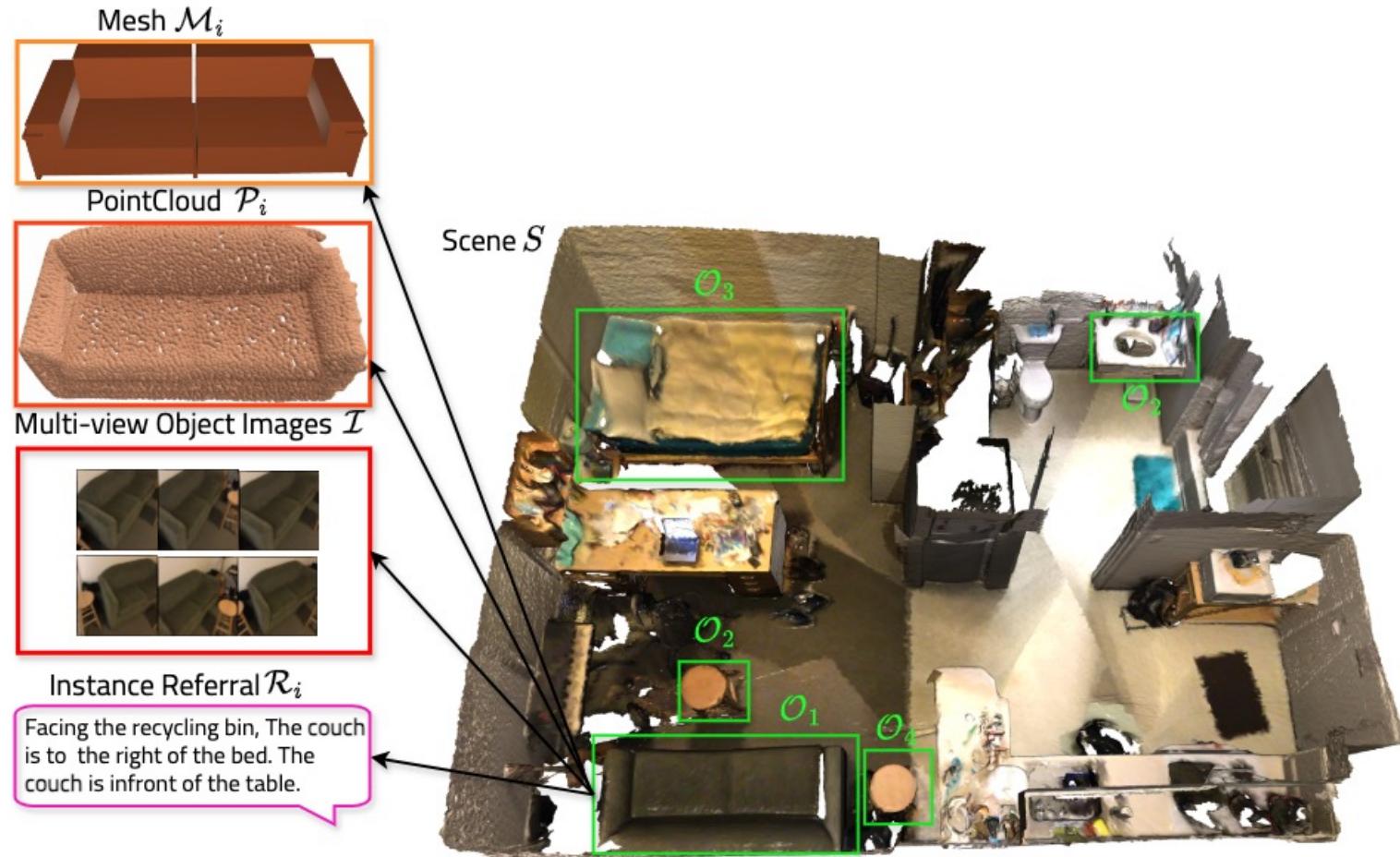


**II. Scene-Level Alignment**

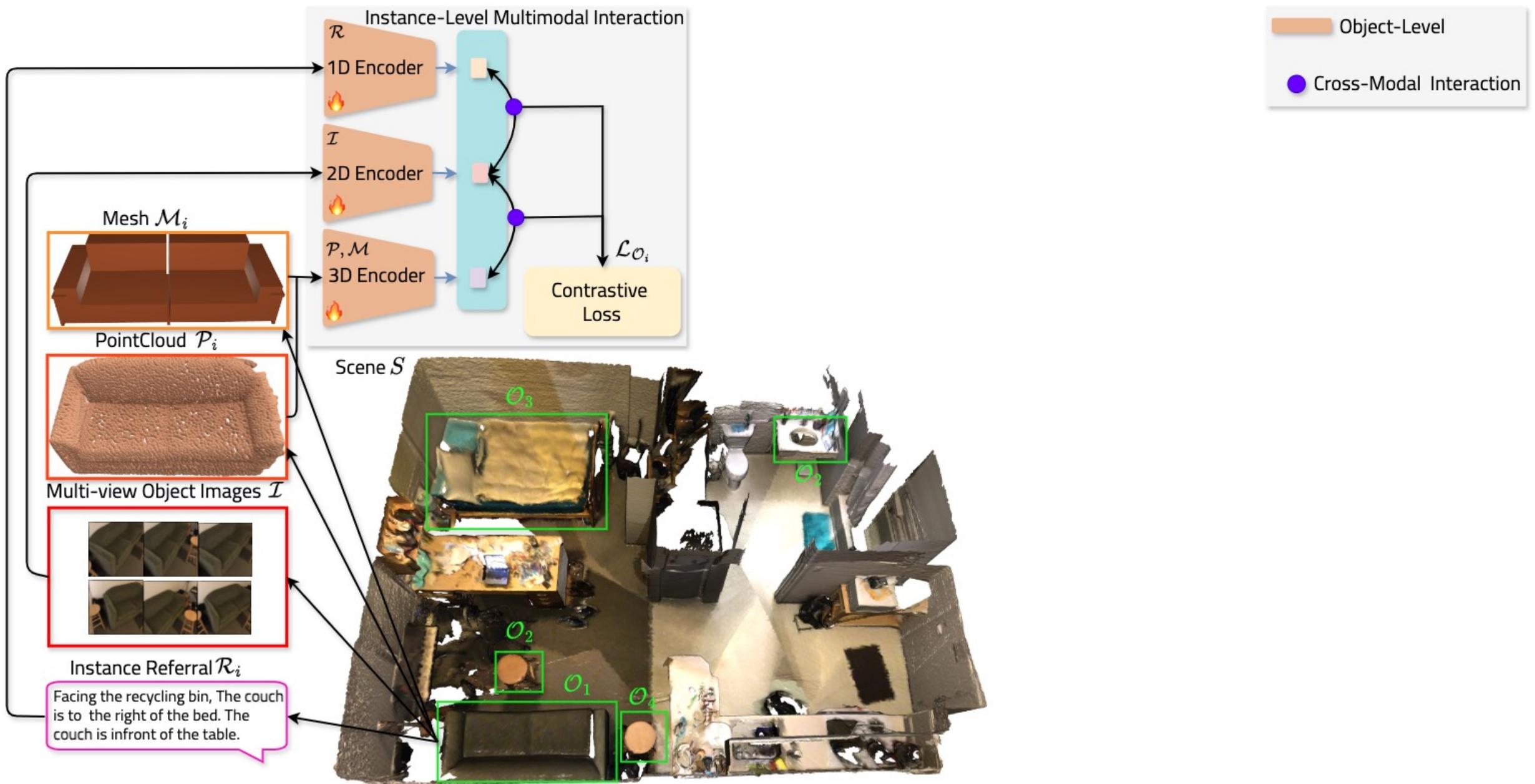
# Overview of CrossOver



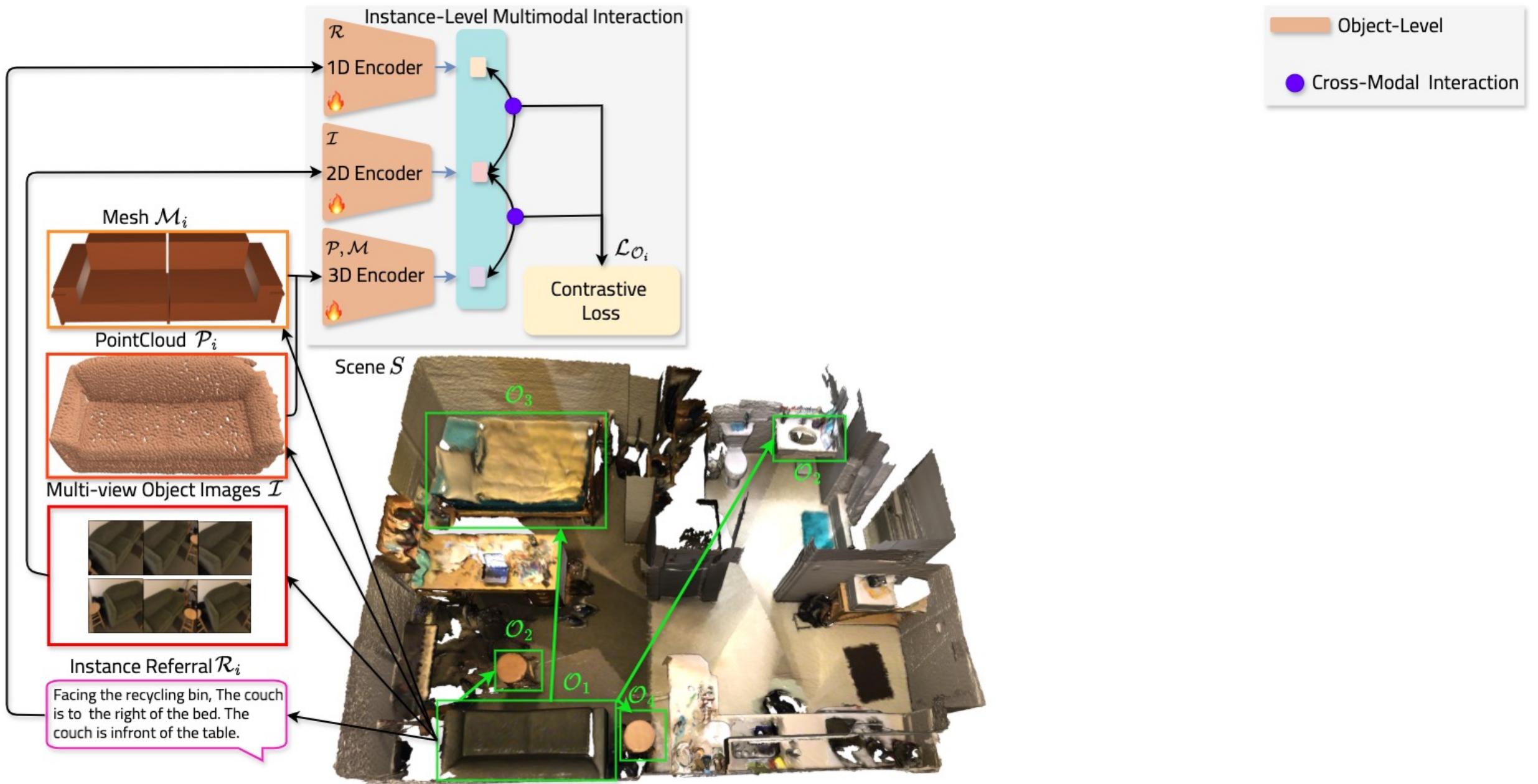
# Overview of CrossOver



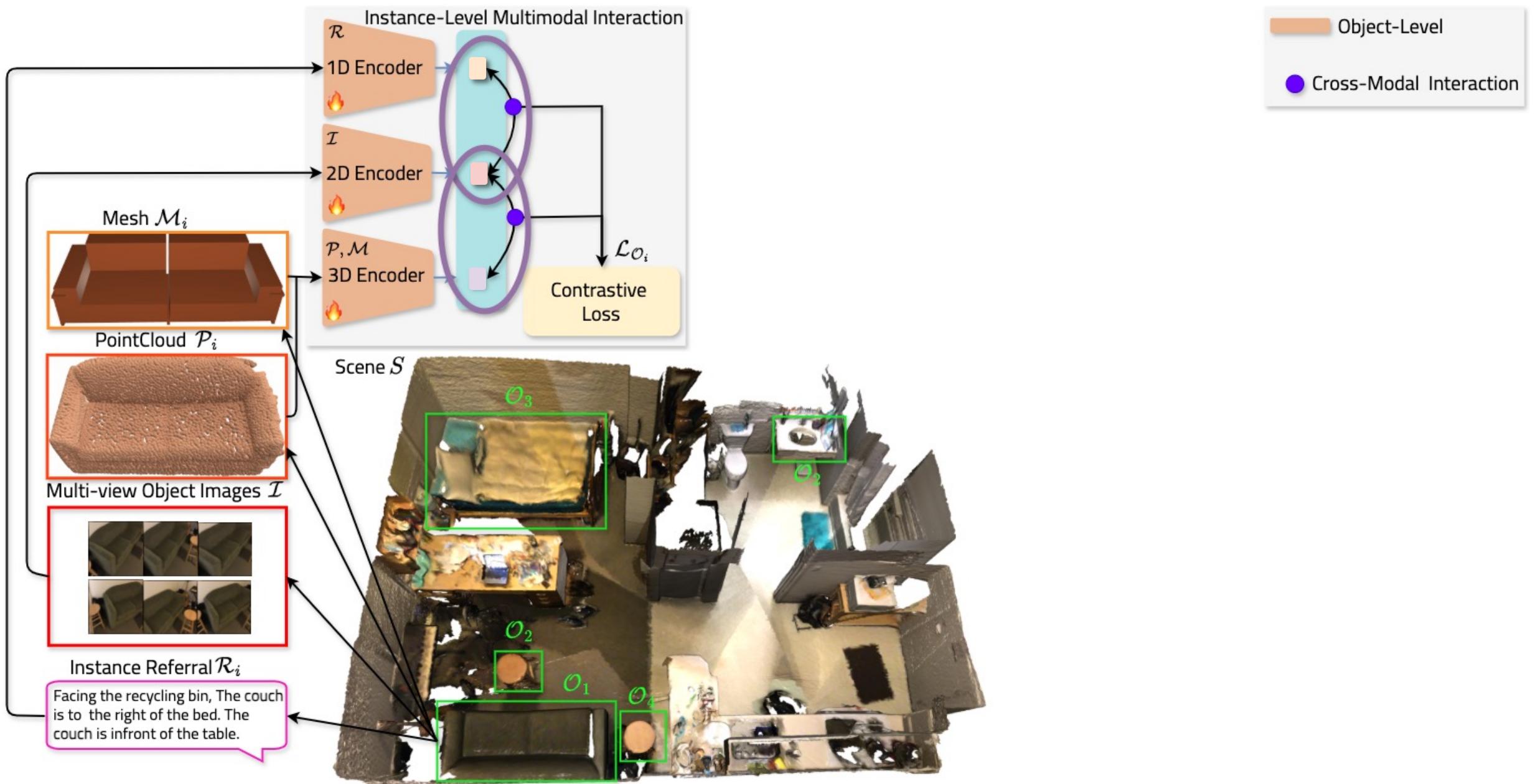
# Overview of CrossOver



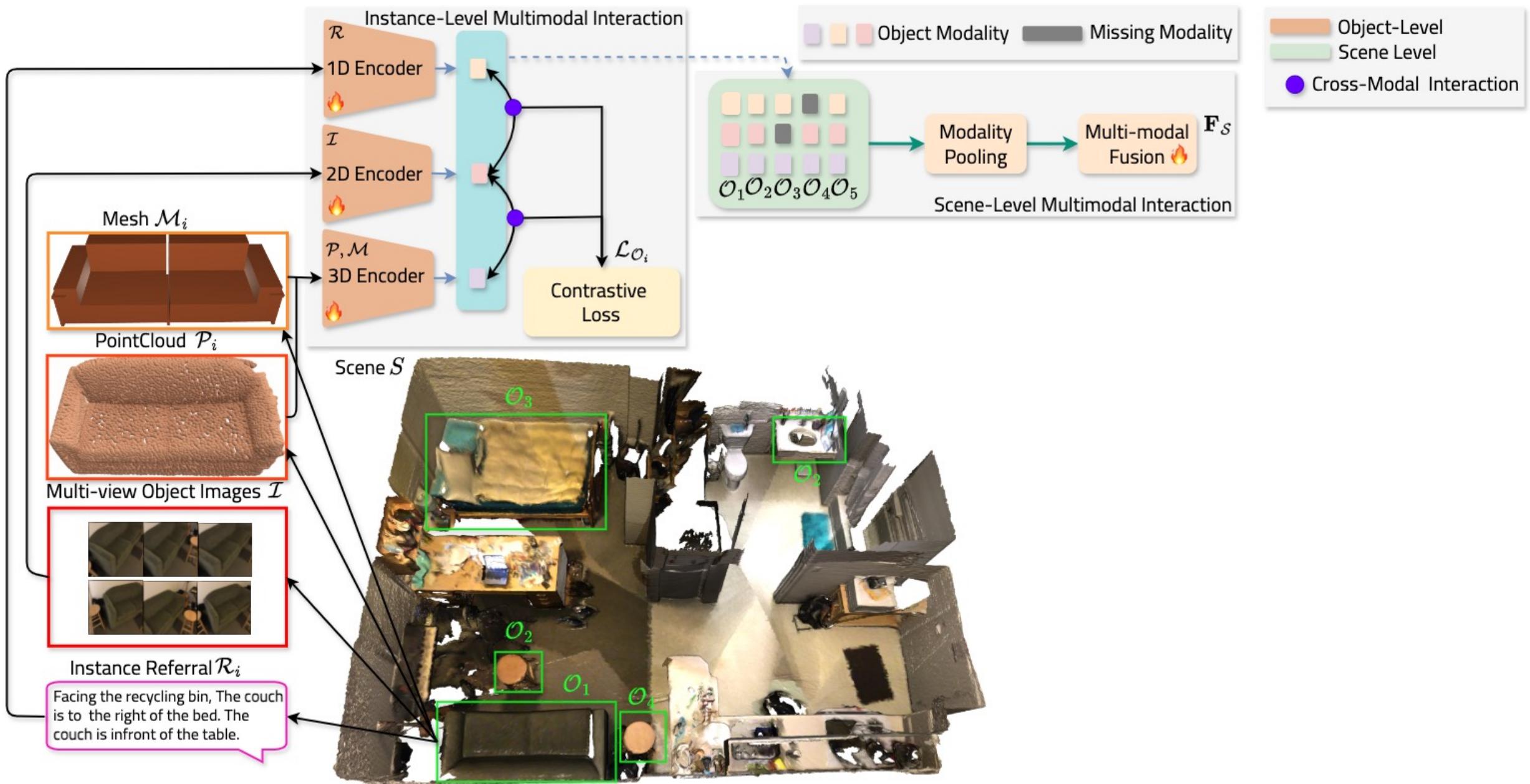
# Overview of CrossOver



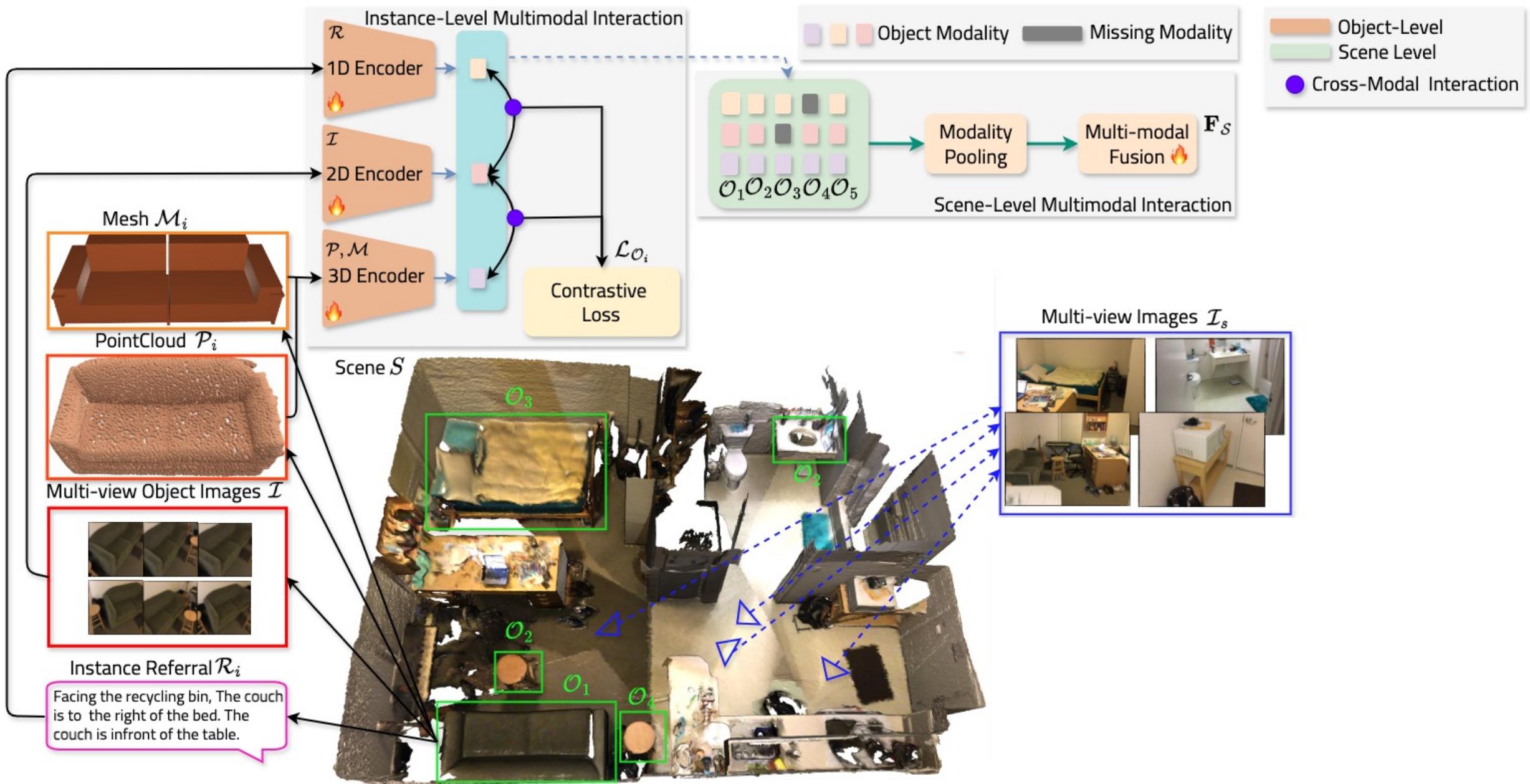
# Overview of CrossOver



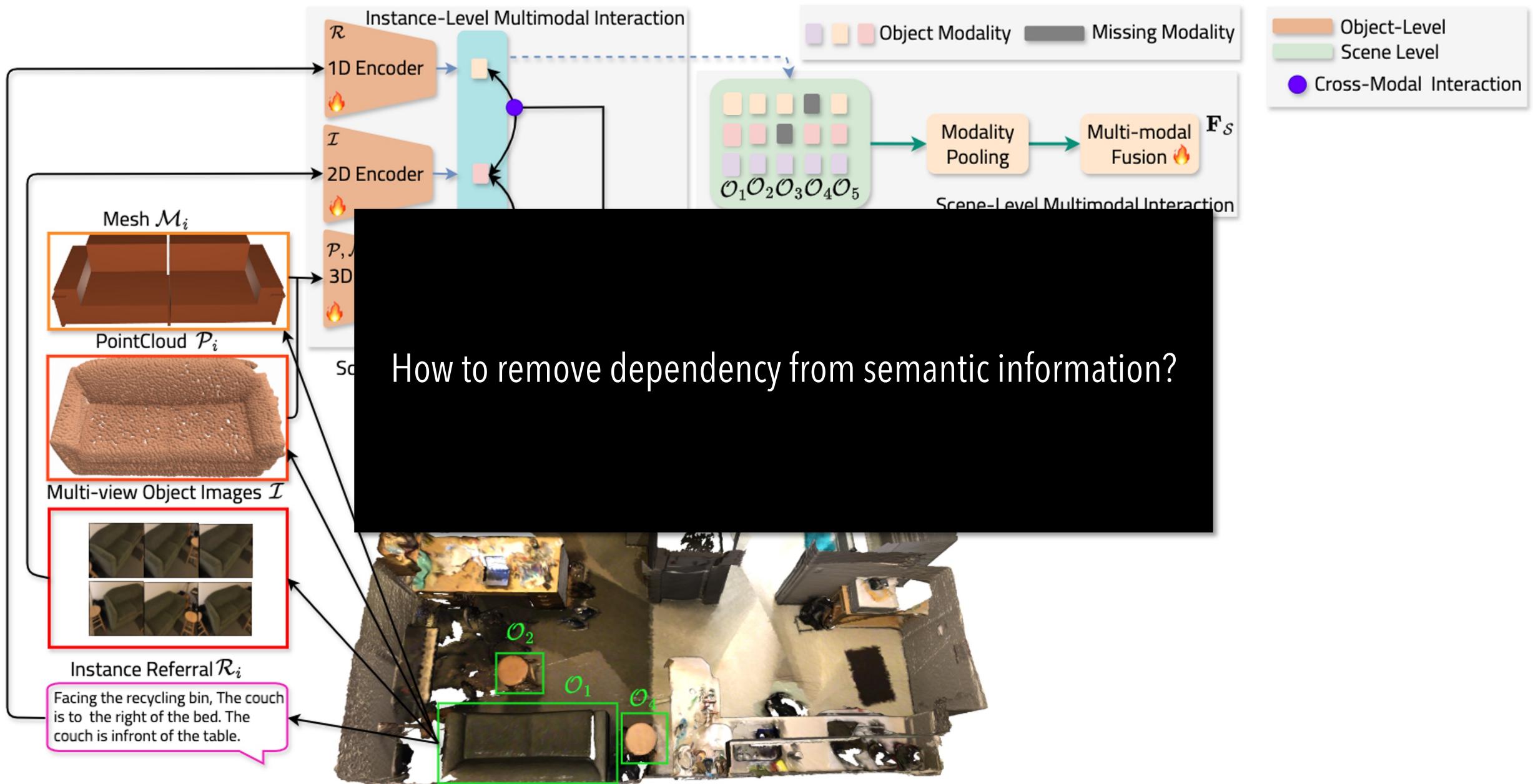
# Overview of CrossOver



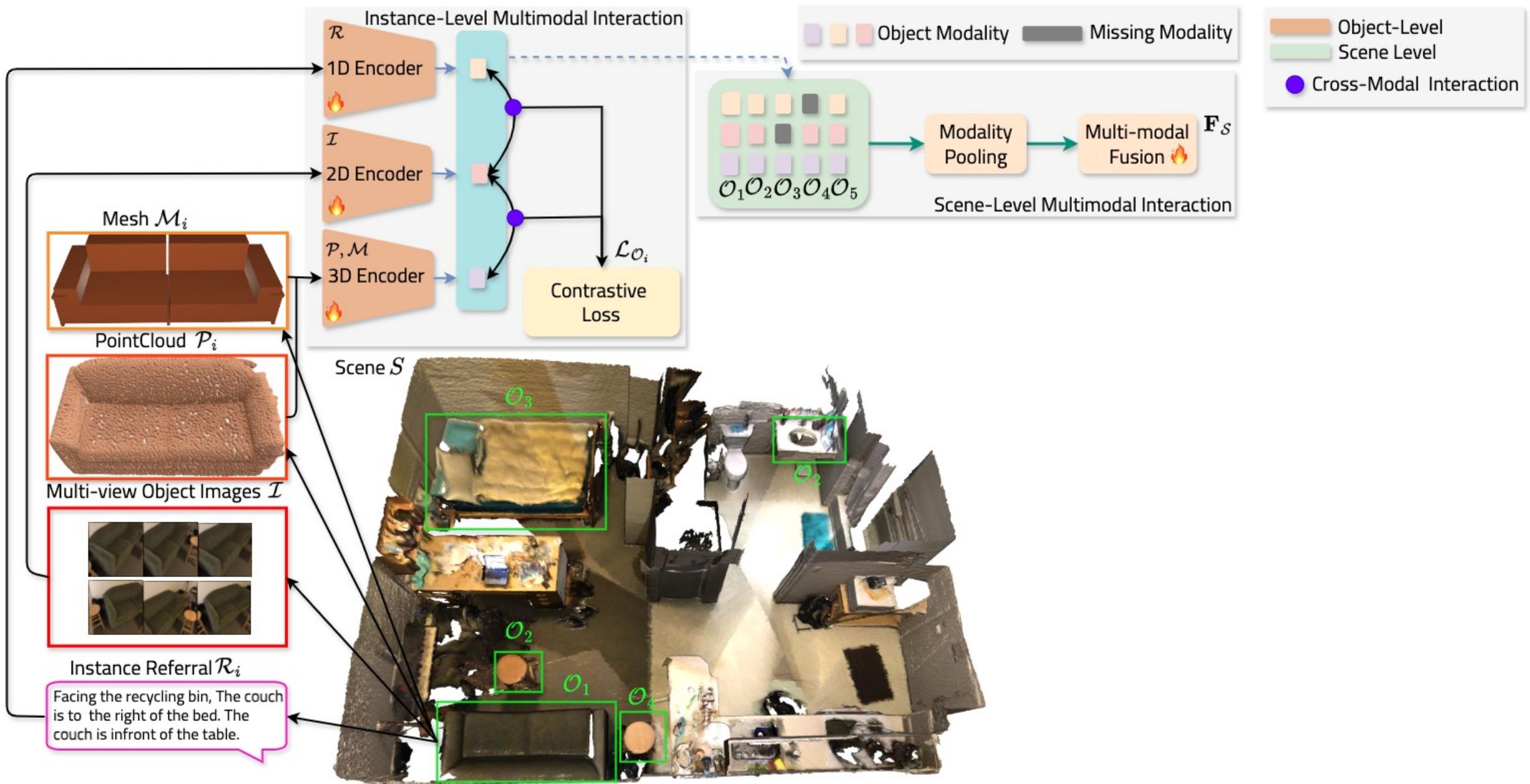
# Overview of CrossOver



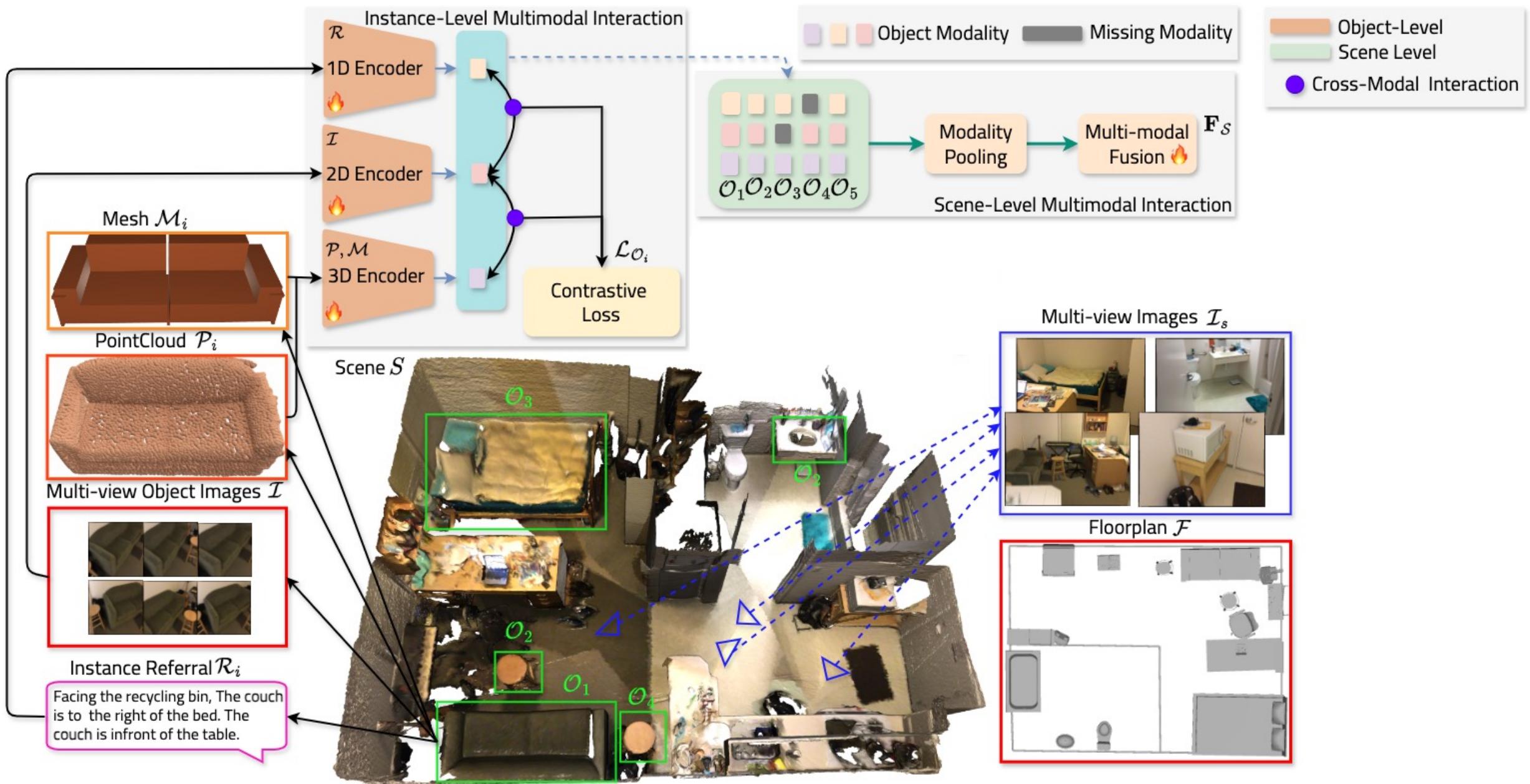
# Overview of CrossOver



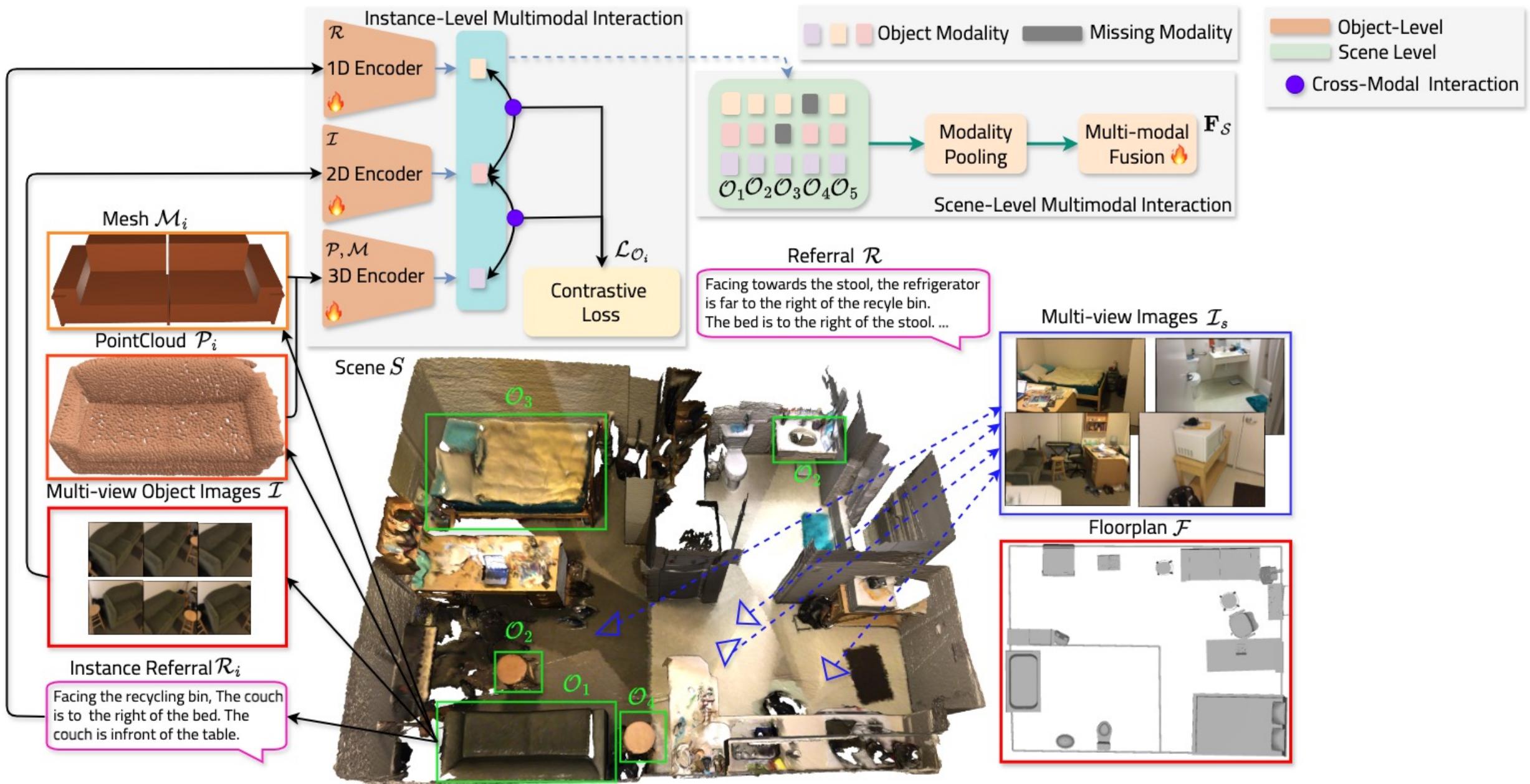
# Overview of CrossOver



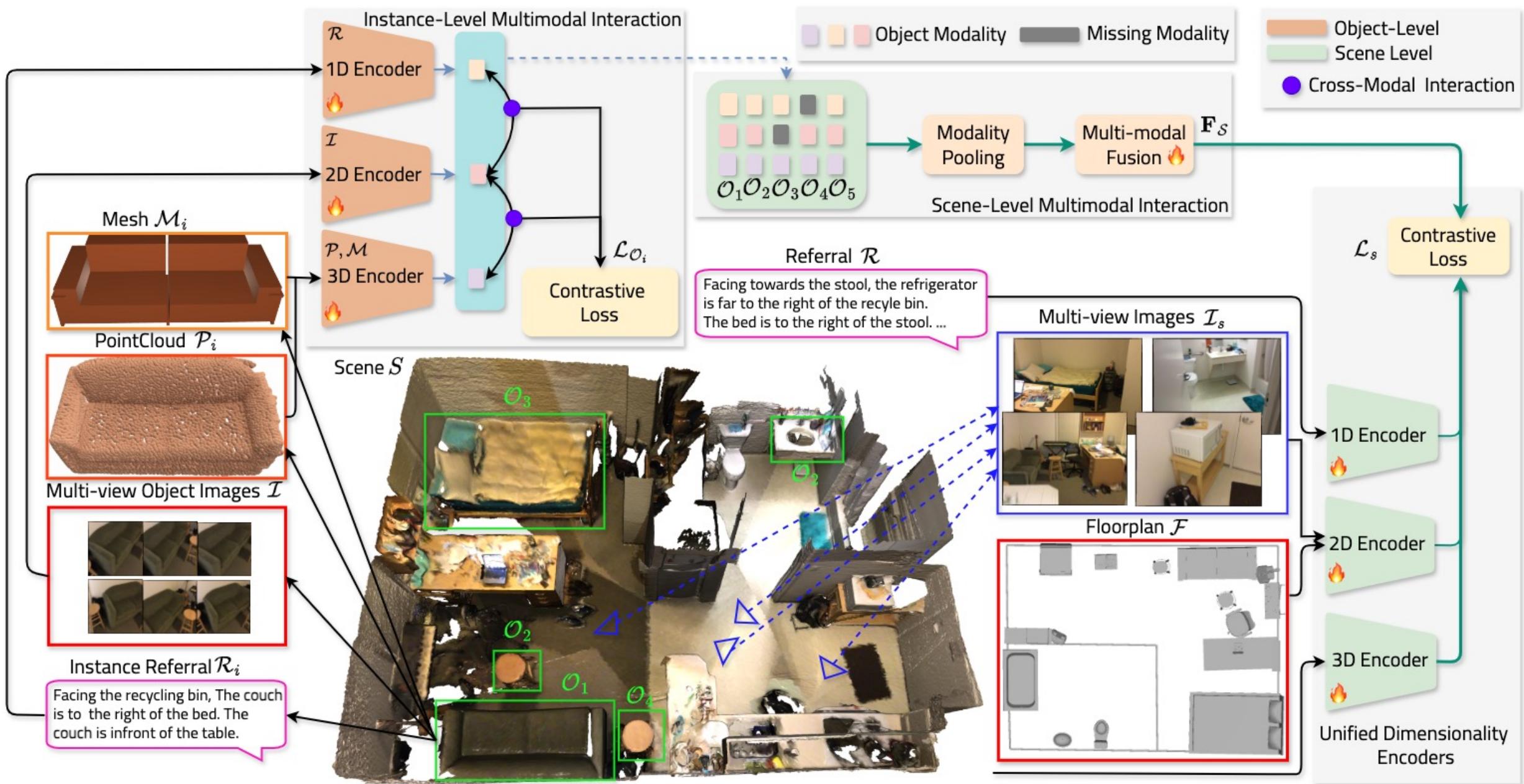
# Overview of CrossOver



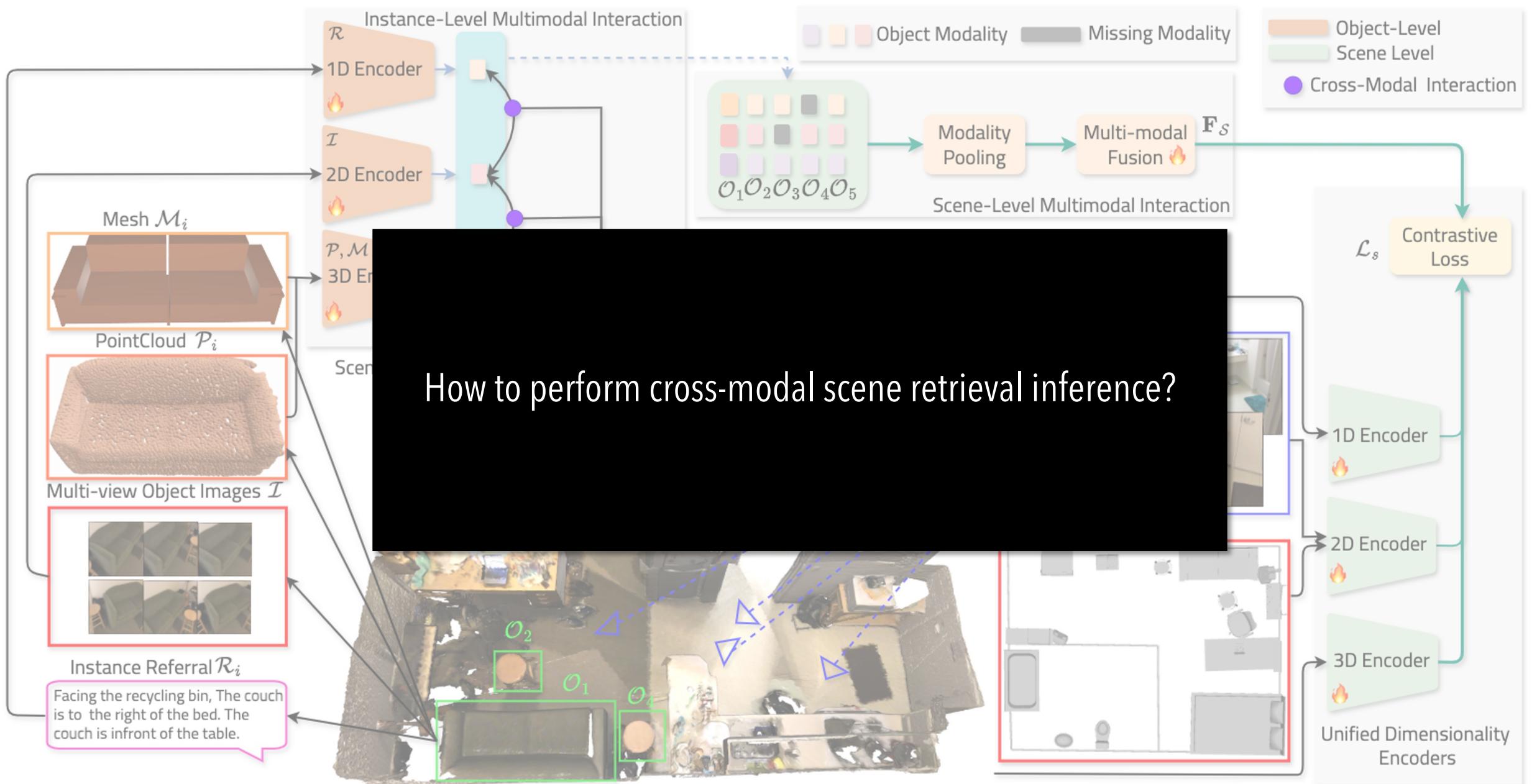
# Overview of CrossOver



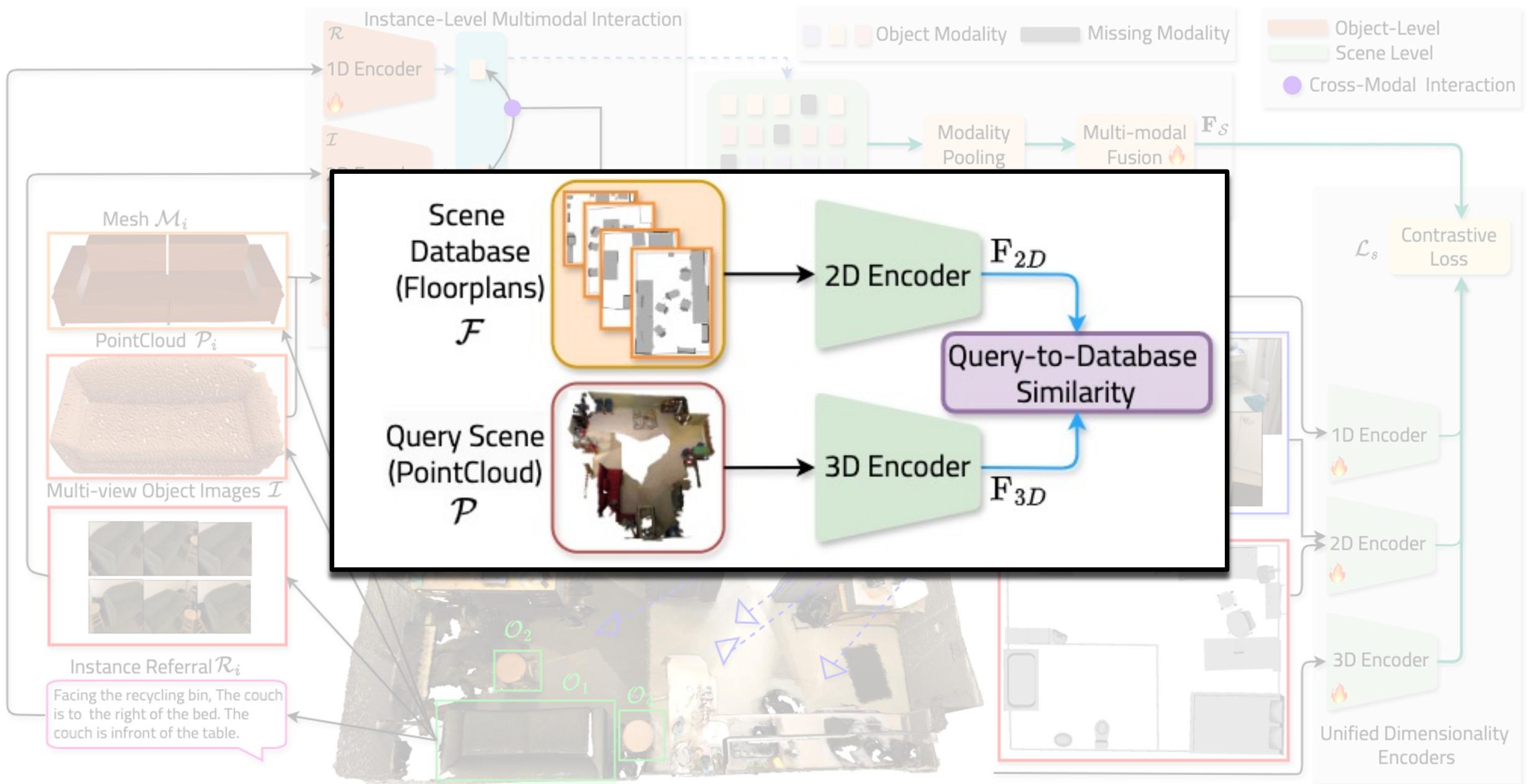
# Overview of CrossOver



# Overview of CrossOver



# Scene Retrieval Inference Pipeline



# **Experimental Results**

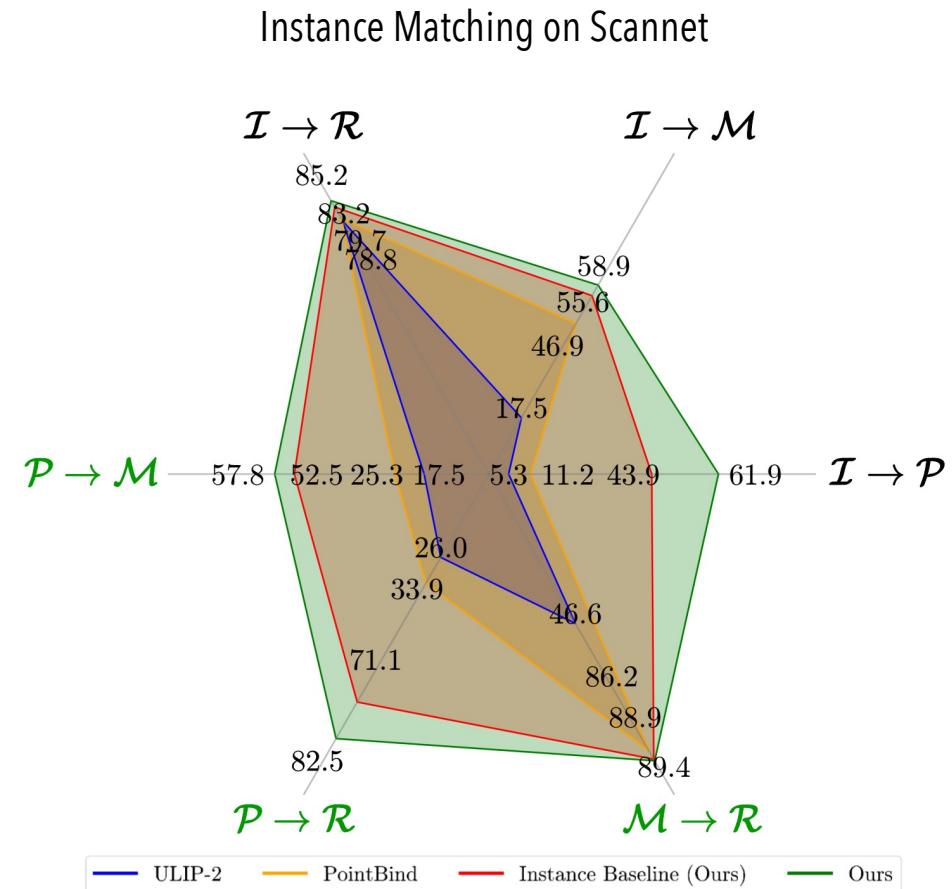
## **Cross-Modal Instance Matching**

# **Experimental Results**

## **Cross-Modal Instance Matching**

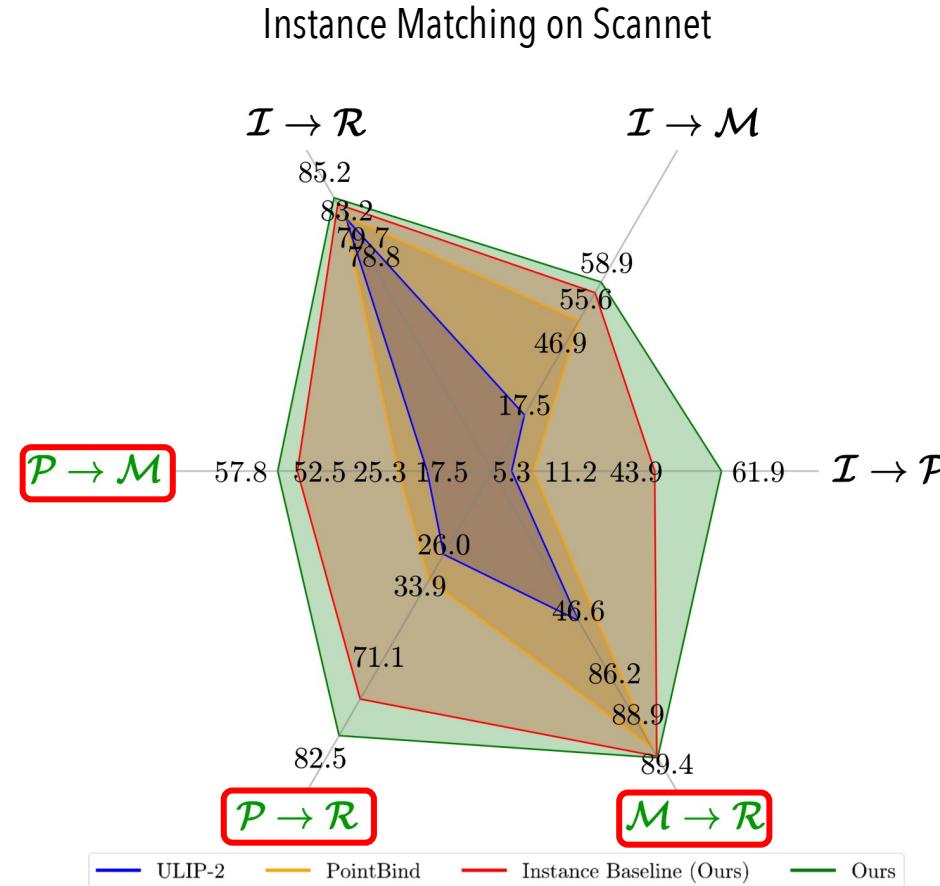
# Experimental Results

## Cross-Modal Instance Matching



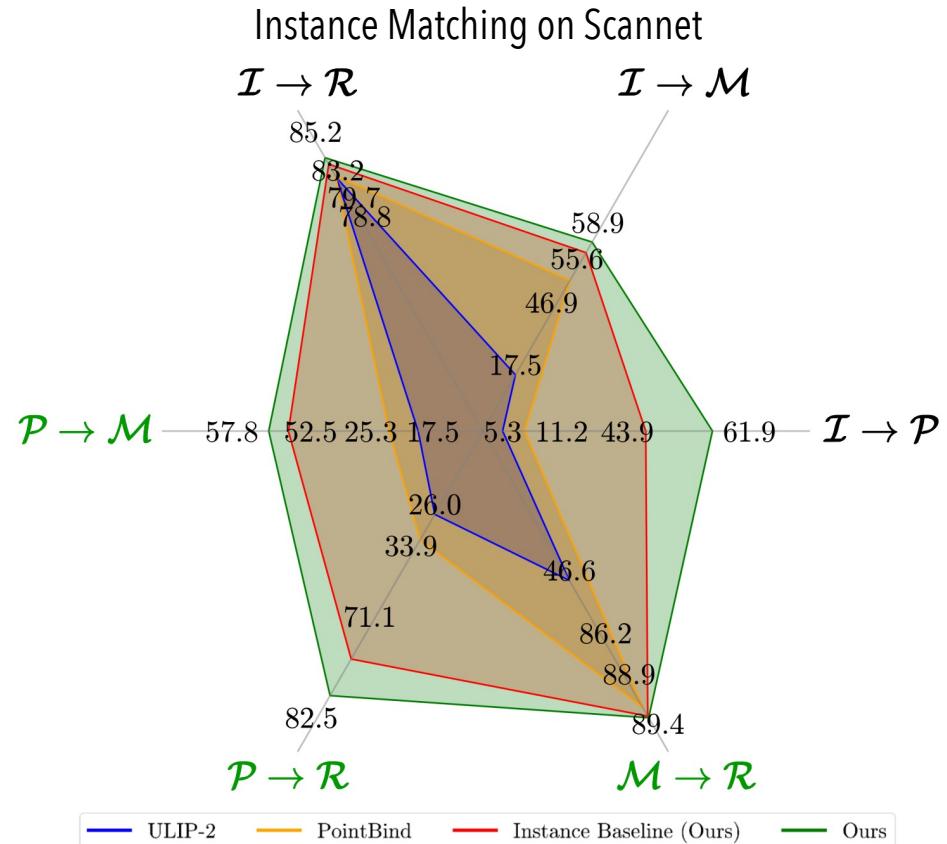
# Experimental Results

## Cross-Modal Instance Matching



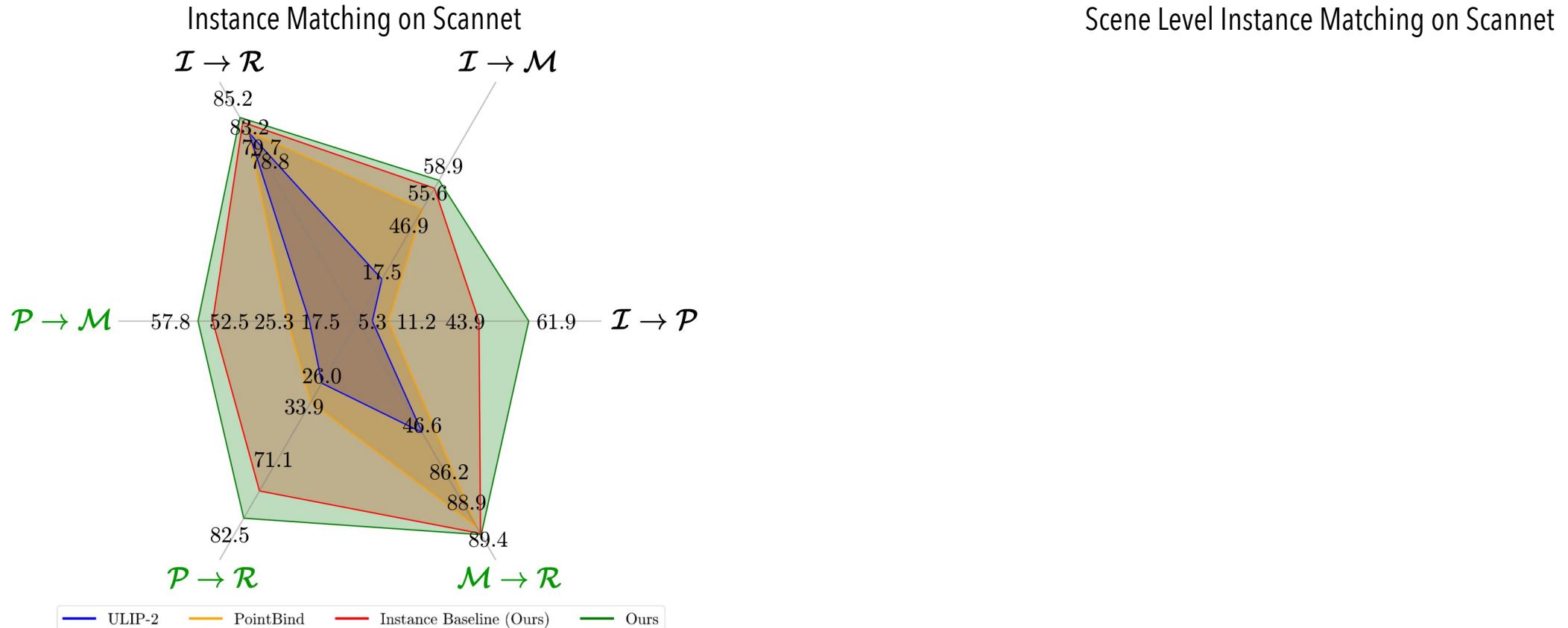
# Experimental Results

## Cross-Modal Instance Matching



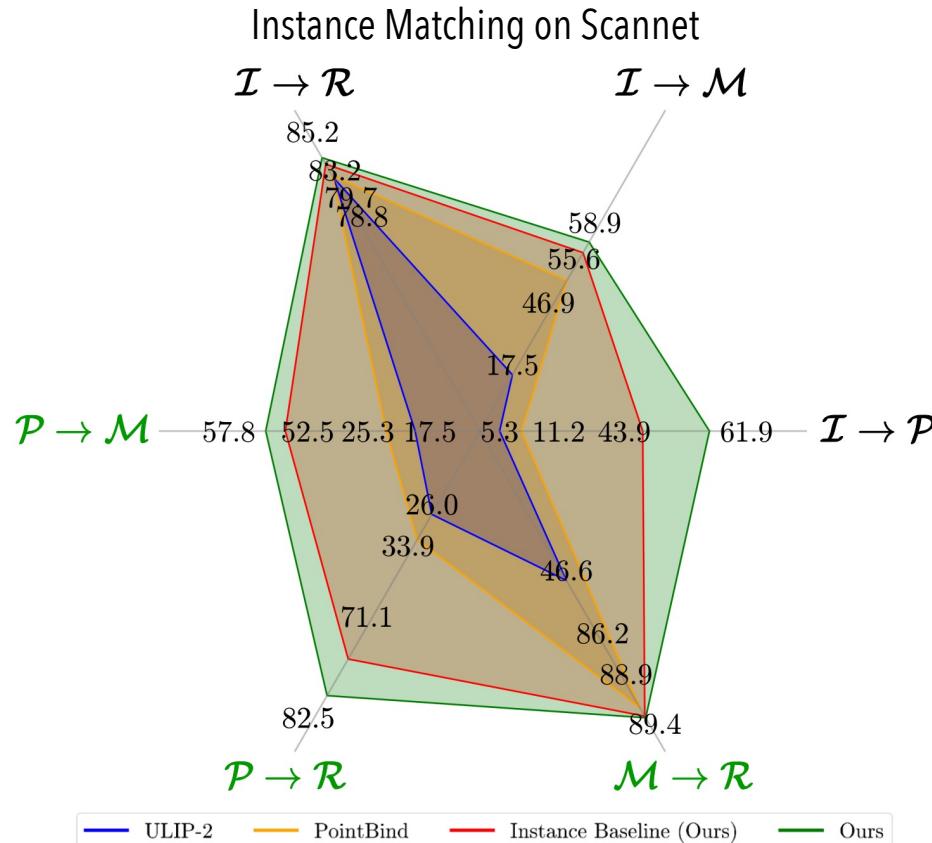
# Experimental Results

## Cross-Modal Instance Matching



# Experimental Results

## Cross-Modal Instance Matching

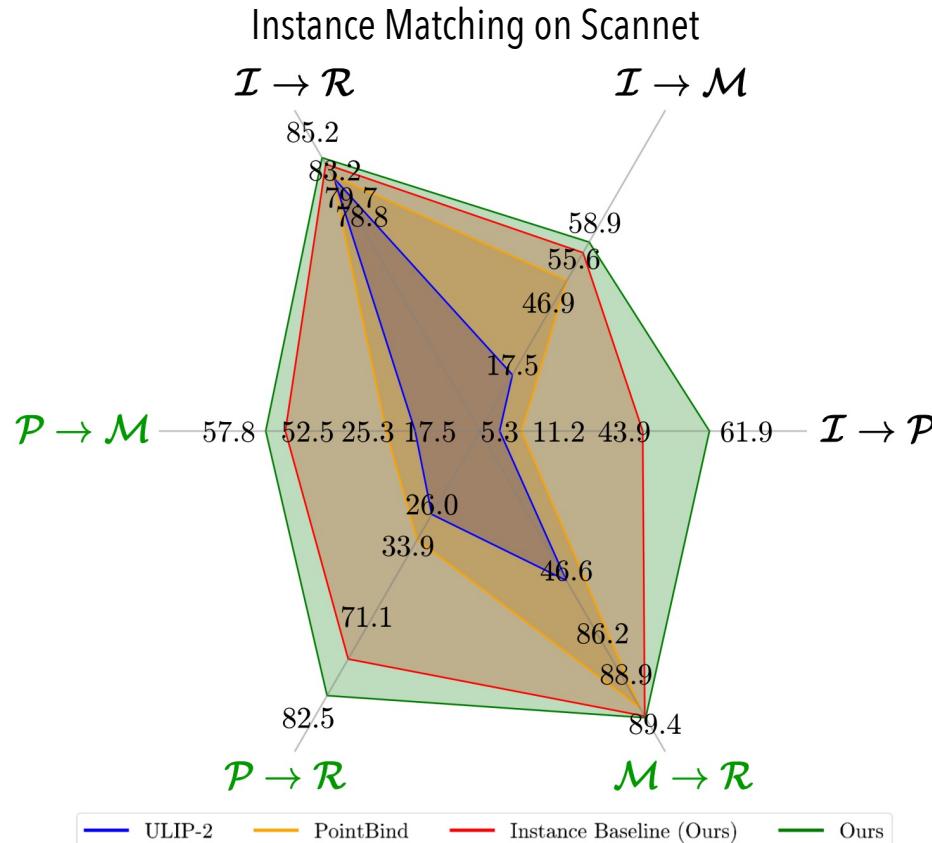


Scene Level Instance Matching on Scannet

Scene-level Recall $\uparrow$	Scannet [11]			3RScan [38]		
	R@25%	R@50%	R@75%	R@25%	R@50%	R@75%
$\mathcal{I} \rightarrow \mathcal{P}$						
ULIP-2 [43]	1.28	0.64	0.24	1.91	0.40	0.28
PointBind [18]	6.73	0.96	0.32	3.18	0.64	0.01
Inst. Baseline (Ours)	88.46	37.82	1.92	93.63	35.03	3.82
Ours	<b>98.08</b>	<b>76.92</b>	<b>23.40</b>	<b>99.36</b>	<b>79.62</b>	<b>22.93</b>
$\mathcal{I} \rightarrow \mathcal{R}$						
ULIP-2 [43]	98.12	96.21	60.34	98.66	85.91	36.91
PointBind [18]	98.22	95.17	62.07	<b>100</b>	87.25	41.61
Inst. Baseline (Ours)	99.31	97.59	71.13	<b>100</b>	92.62	55.03
Ours	<b>99.66</b>	<b>98.28</b>	<b>76.29</b>	<b>100</b>	<b>97.32</b>	<b>67.79</b>
$\mathcal{P} \rightarrow \mathcal{R}$						
ULIP-2 [43]	37.24	16.90	8.62	16.78	6.04	1.34
PointBind [18]	54.83	27.93	11.72	21.48	6.04	2.01
Inst. Baseline (Ours)	98.63	83.85	46.74	92.62	60.40	20.81
Ours	<b>99.31</b>	<b>96.56</b>	<b>70.10</b>	<b>100</b>	<b>89.26</b>	<b>50.34</b>

# Experimental Results

## Cross-Modal Instance Matching



Scene Level Instance Matching on Scannet

Scene-level Recall ↑	Scannet [11]			3RScan [38]		
	R@25%	R@50%	R@75%	R@25%	R@50%	R@75%
$\mathcal{I} \rightarrow \mathcal{P}$						
ULIP-2 [43]	1.28	0.64	0.24	1.91	0.40	0.28
PointBind [18]	6.73	0.96	0.32	3.18	0.64	0.01
Inst. Baseline (Ours)	88.46	37.82	1.92	93.63	35.03	3.82
Ours	<b>98.08</b>	<b>76.92</b>	<b>23.40</b>	<b>99.36</b>	<b>79.62</b>	<b>22.93</b>
$\mathcal{I} \rightarrow \mathcal{R}$						
ULIP-2 [43]	98.12	96.21	60.34	98.66	85.91	36.91
PointBind [18]	98.22	95.17	62.07	<b>100</b>	87.25	41.61
Inst. Baseline (Ours)	99.31	97.59	71.13	<b>100</b>	92.62	55.03
Ours	<b>99.66</b>	<b>98.28</b>	<b>76.29</b>	<b>100</b>	<b>97.32</b>	<b>67.79</b>
$\mathcal{P} \rightarrow \mathcal{R}$						
ULIP-2 [43]	37.24	16.90	8.62	16.78	6.04	1.34
PointBind [18]	54.83	27.93	11.72	21.48	6.04	2.01
Inst. Baseline (Ours)	98.63	83.85	46.74	92.62	60.40	20.81
Ours	<b>99.31</b>	<b>96.56</b>	<b>70.10</b>	<b>100</b>	<b>89.26</b>	<b>50.34</b>

- Outperforms all baselines in all datasets
- Emergent behaviour within the embedding space, highlighting robustness of learning cross-modal interactions

# **Experimental Results**

## **Temporal Instance Matching**

# Experimental Results

## Temporal Instance Matching

Method	Scene-level Recall ↑		
	R@25%	R@50%	R@75%
<i>same-modal (<math>\mathcal{P} \rightarrow \mathcal{P}</math>)</i>			
MendNet [15]	80.68	64.77	37.50
VN-DGCNN <sub>cls</sub> [13]	72.32	53.41	29.55
VN-ONet <sub>recon</sub> [13]	86.36	71.59	44.32
LivingScenes [47]	87.50	78.41	50.00
Ours	<b>92.31</b>	<b>84.62</b>	<b>57.69</b>
<i>cross-modal (ours)</i>			
$\mathcal{I} \rightarrow \mathcal{P}$	89.74	73.08	42.31
$\mathcal{I} \rightarrow \mathcal{R}$	62.33	38.96	18.18
$\mathcal{P} \rightarrow \mathcal{R}$	68.83	40.26	22.08

# Experimental Results

## Temporal Instance Matching

Method	Scene-level Recall ↑		
	R@25%	R@50%	R@75%
<i>same-modal (<math>\mathcal{P} \rightarrow \mathcal{P}</math>)</i>			
MendNet [15]	80.68	64.77	37.50
VN-DGCNN <sub>cls</sub> [13]	72.32	53.41	29.55
VN-ONet <sub>recon</sub> [13]	86.36	71.59	44.32
LivingScenes [47]	87.50	78.41	50.00
Ours	<b>92.31</b>	<b>84.62</b>	<b>57.69</b>
<i>cross-modal (ours)</i>			
$\mathcal{I} \rightarrow \mathcal{P}$	89.74	73.08	42.31
$\mathcal{I} \rightarrow \mathcal{R}$	62.33	38.96	18.18
$\mathcal{P} \rightarrow \mathcal{R}$	68.83	40.26	22.08

- Better performance in the same-modal task compared to baselines, despite **not being specifically trained on temporal data!**

# **Experimental Results**

## **Cross-Modal Scene Retrieval**

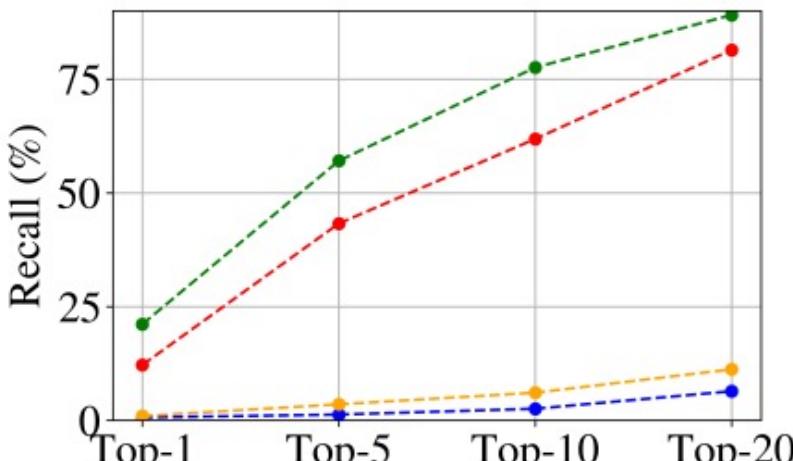
# Experimental Results

## Cross-Modal Scene Retrieval

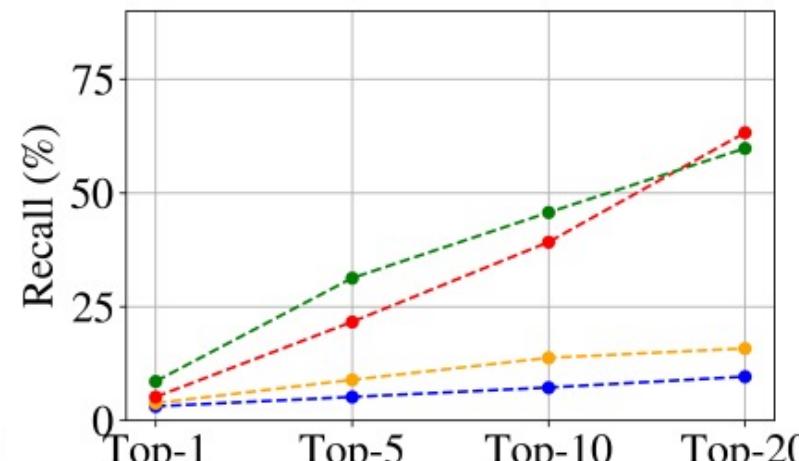
Scene matching recall of different methods on three modality pairs

top-K out of 312 validation scans

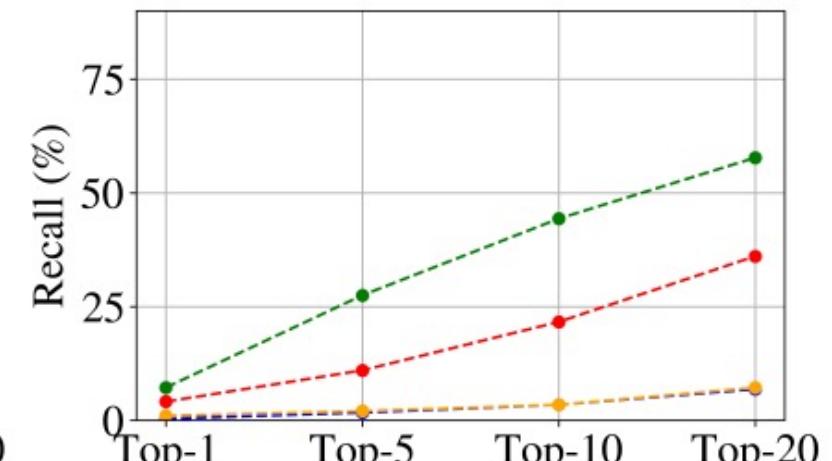
$\mathcal{I} \rightarrow \mathcal{P}$



$\mathcal{I} \rightarrow \mathcal{R}$



$\mathcal{P} \rightarrow \mathcal{R}$



— ULIP-2 — PointBind — Instance Baseline (Ours) — Ours

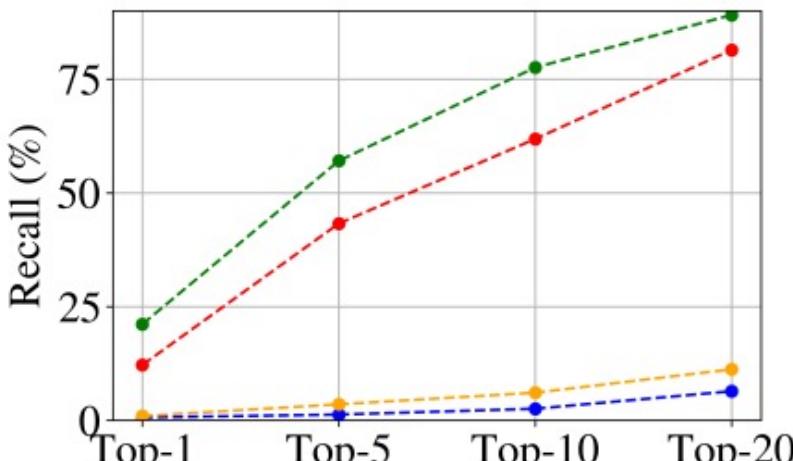
# Experimental Results

## Cross-Modal Scene Retrieval

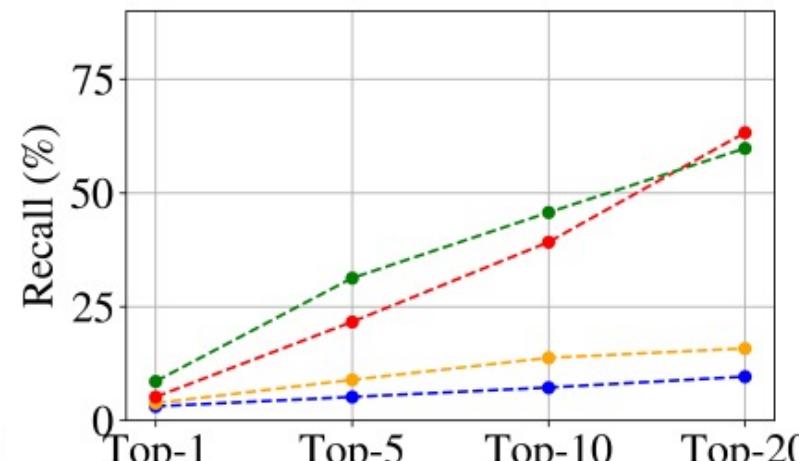
Scene matching recall of different methods on three modality pairs

top-K out of 312 validation scans

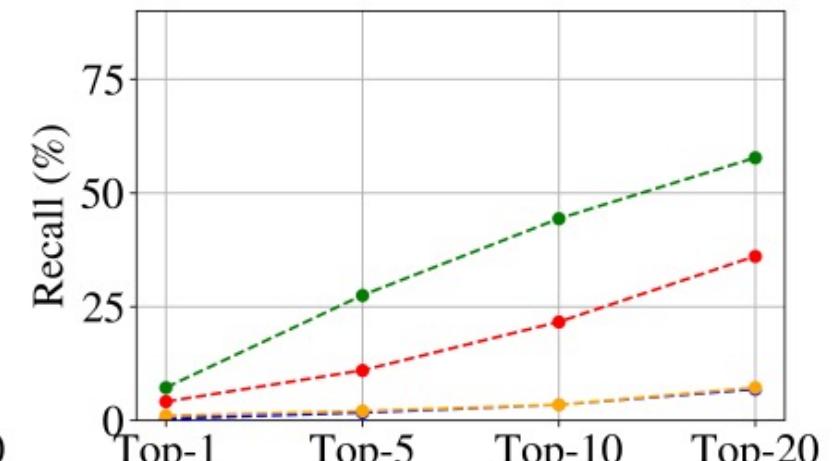
$\mathcal{I} \rightarrow \mathcal{P}$



$\mathcal{I} \rightarrow \mathcal{R}$



$\mathcal{P} \rightarrow \mathcal{R}$



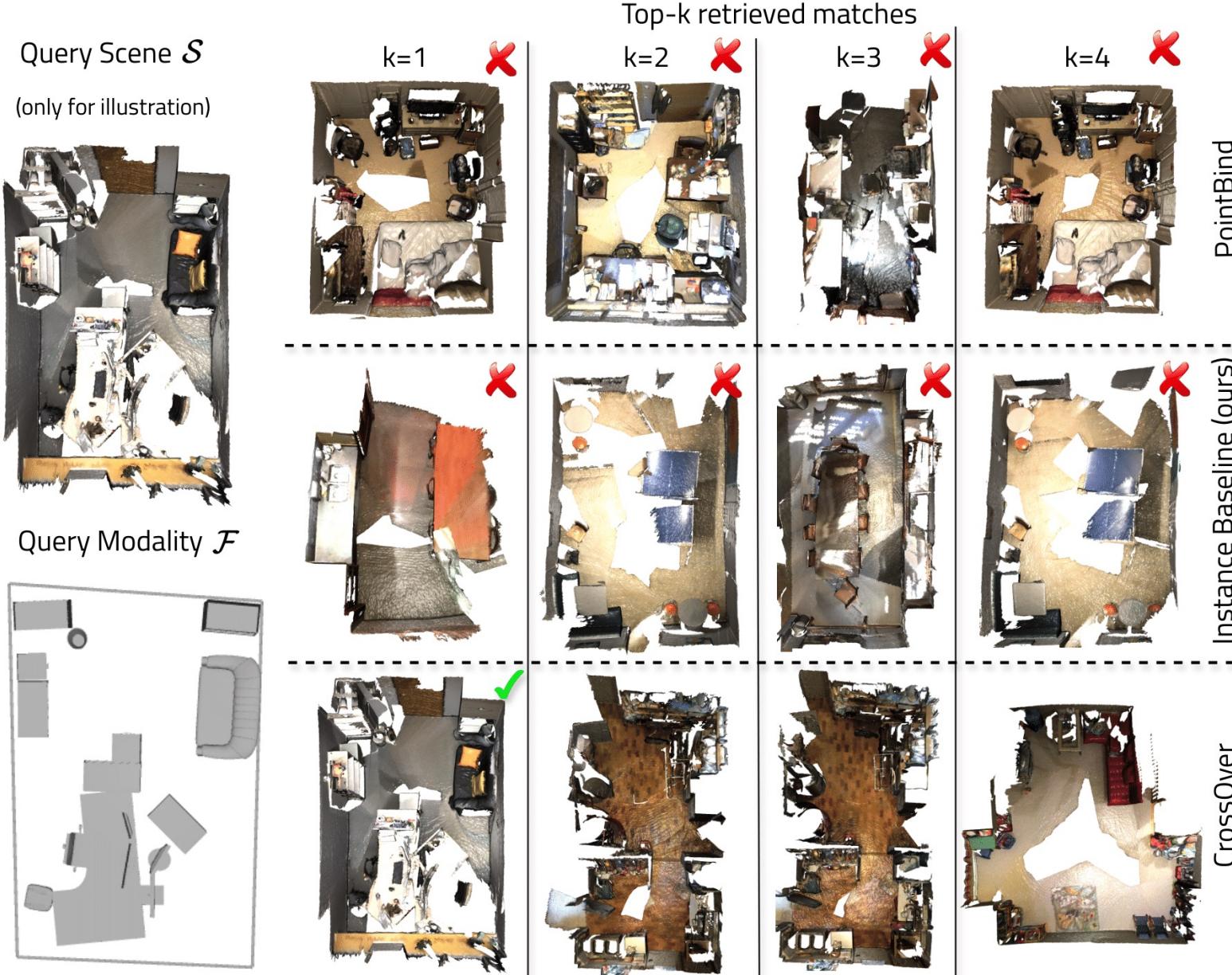
- Unified dimensionality encoders: do not rely on semantics, consistently outperform prior methods in all pairwise modalities

# **Experimental Results**

**Cross-Modal Scene Retrieval Visualization: Success**

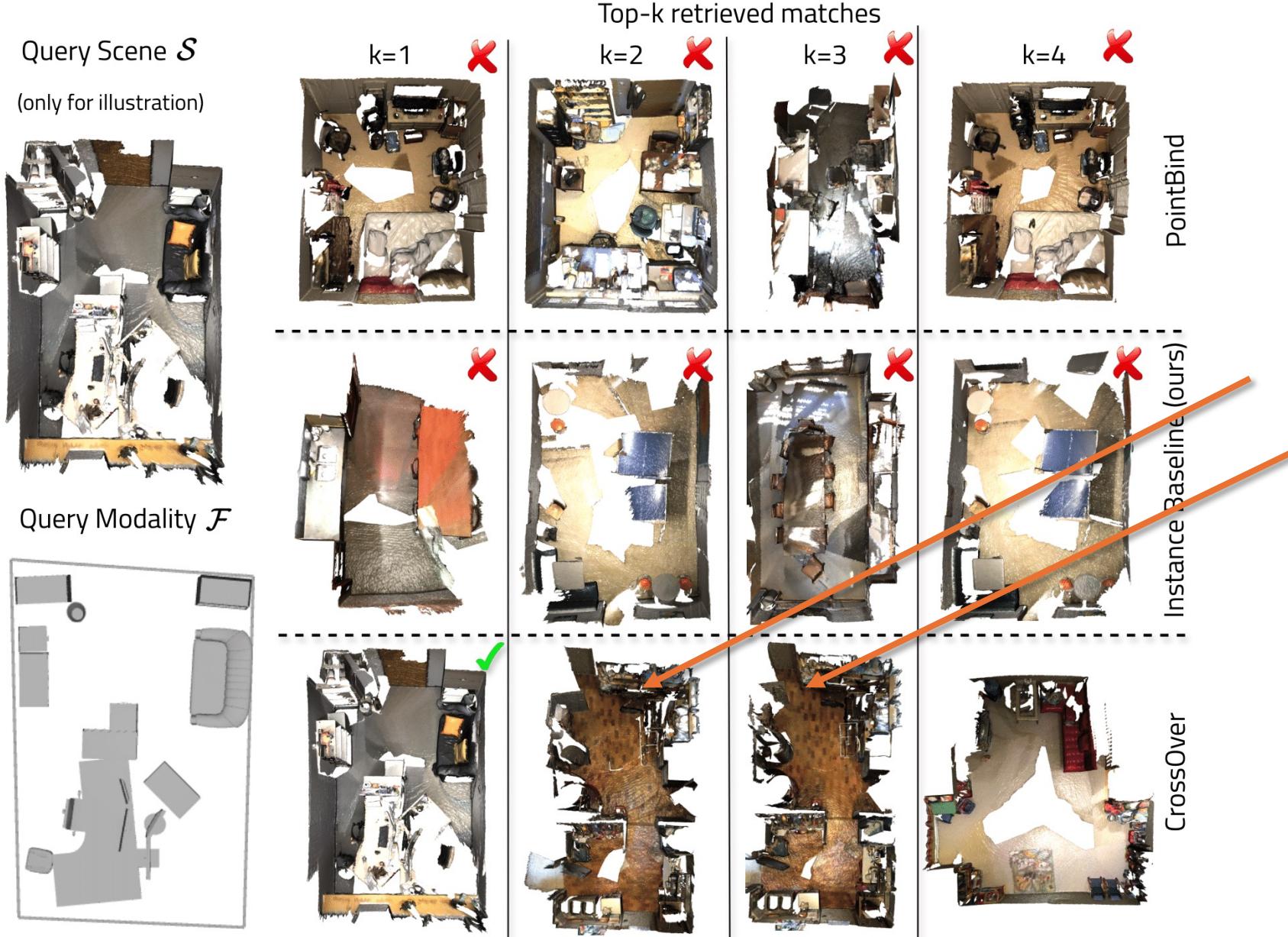
# Experimental Results

## Cross-Modal Scene Retrieval Visualization: Success



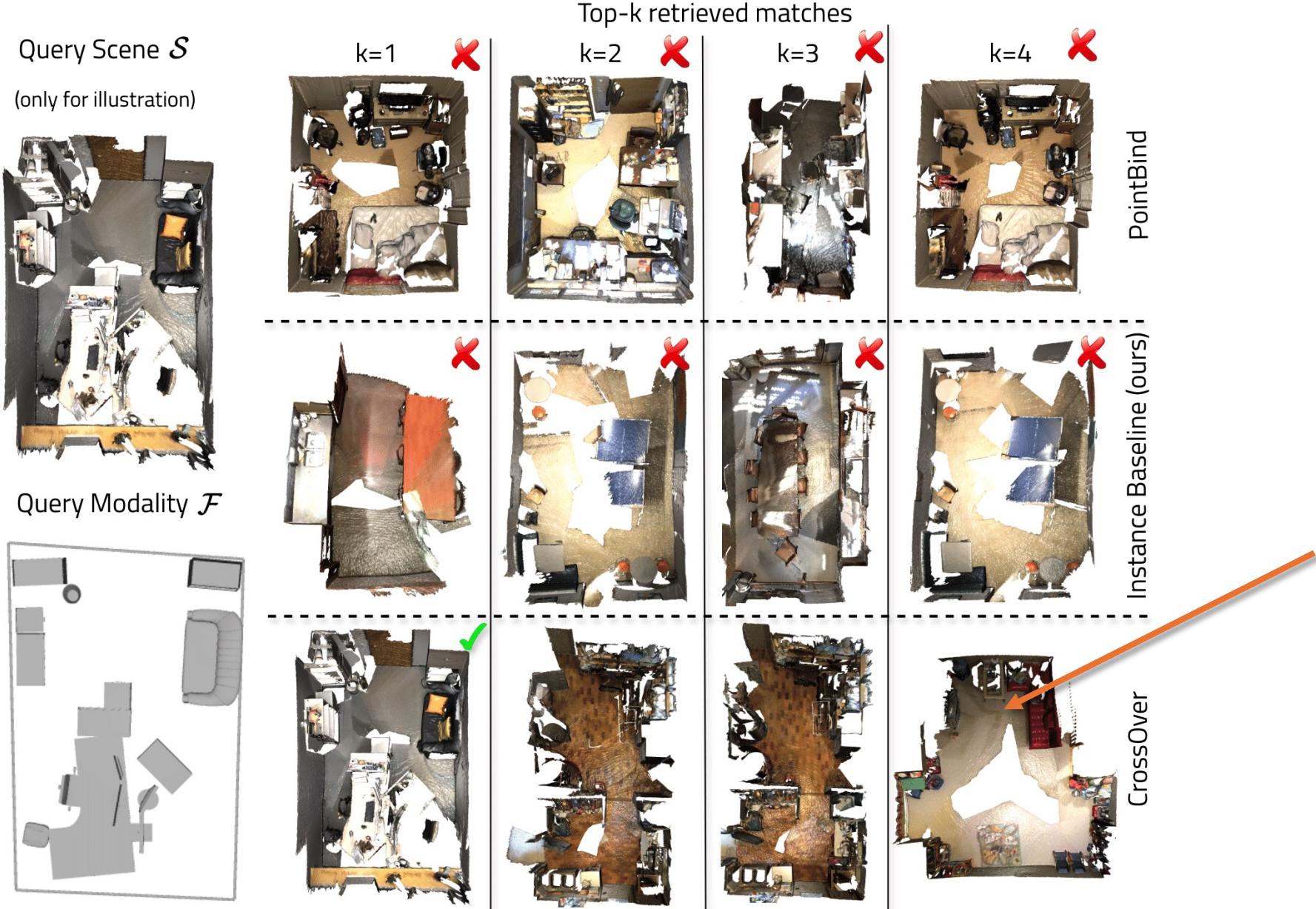
# Experimental Results

## Cross-Modal Scene Retrieval Visualization: Success



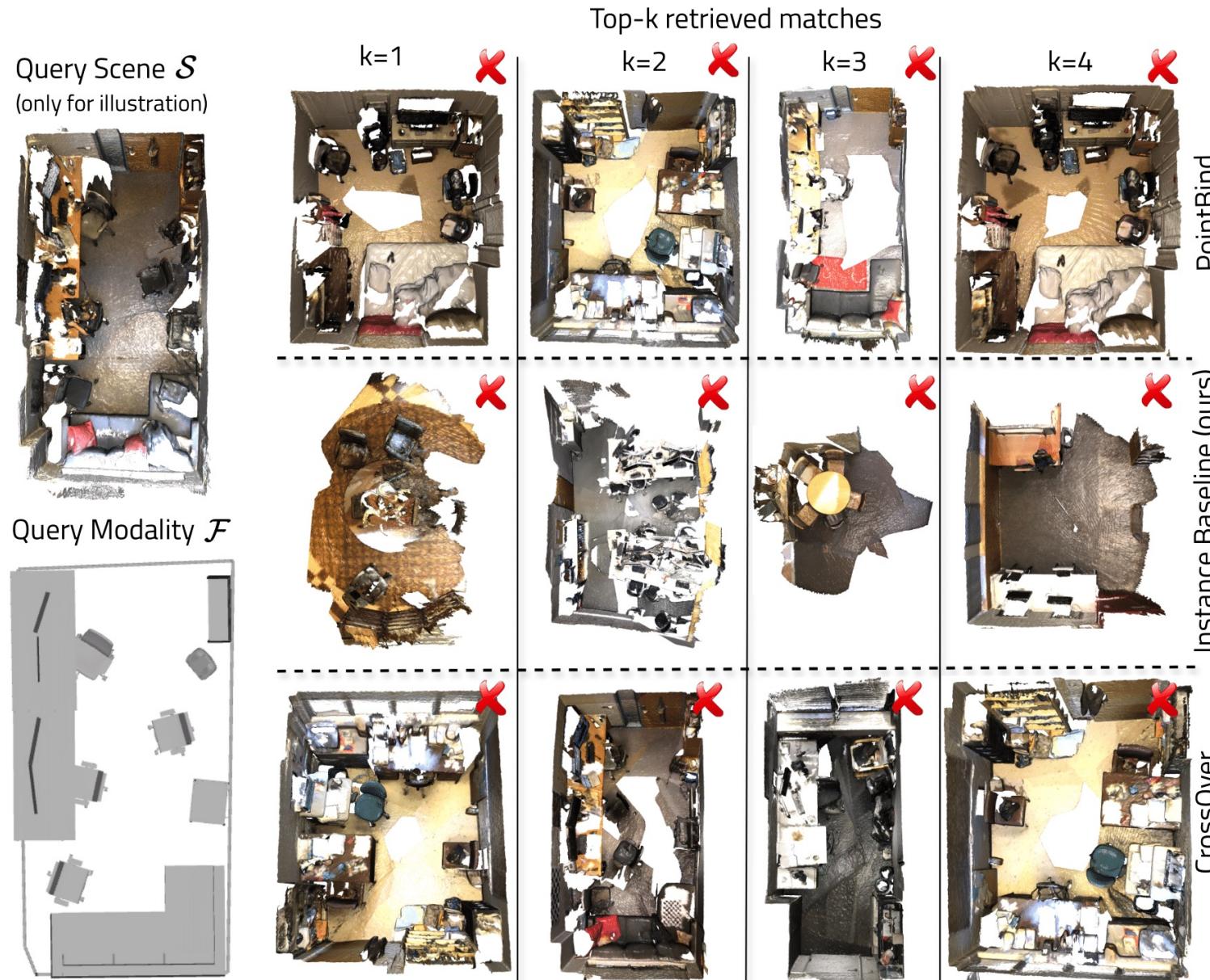
# Experimental Results

## Cross-Modal Scene Retrieval Visualization: Success



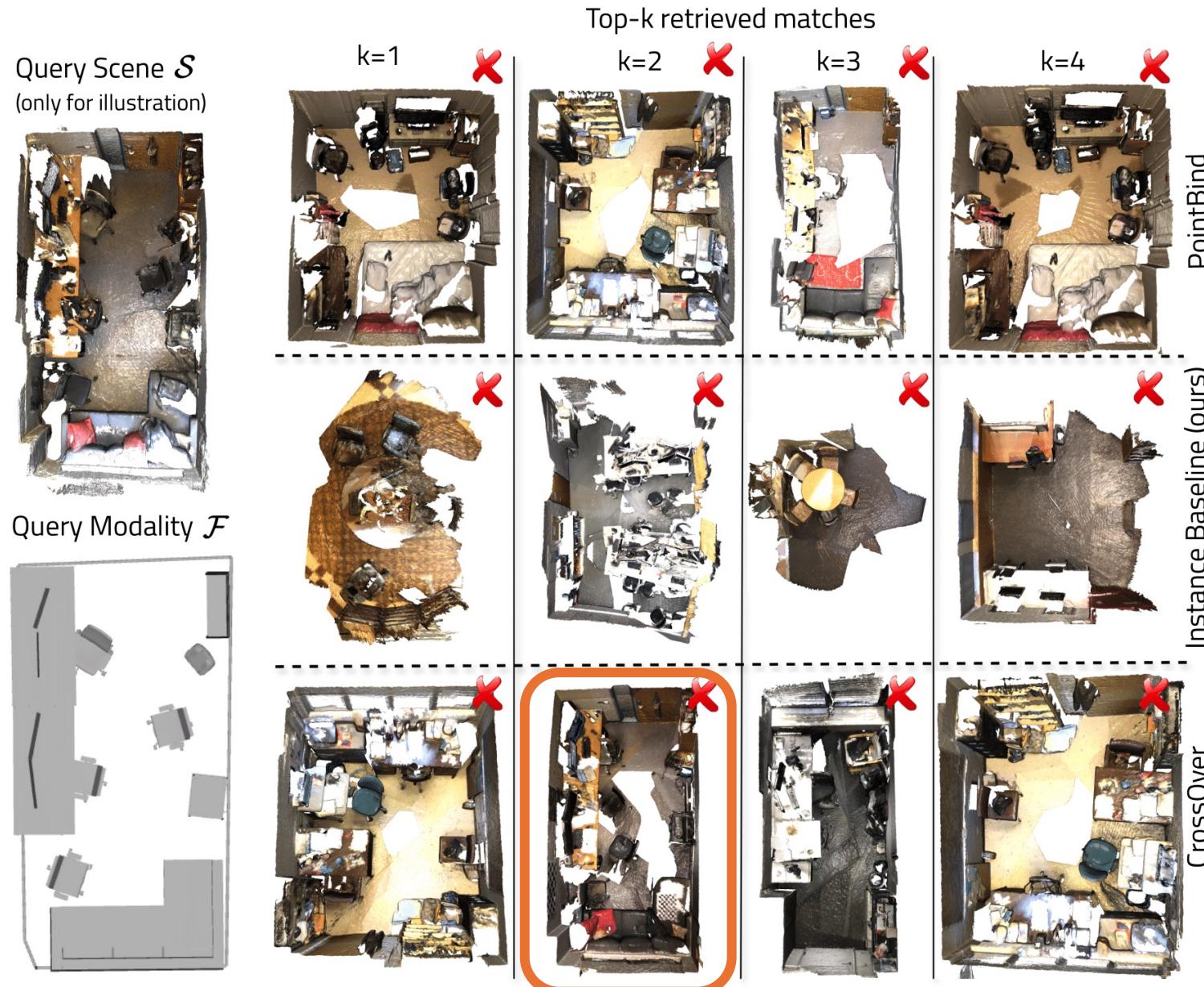
# Experimental Results

## Cross-Modal Scene Retrieval Visualization: Failure



# Experimental Results

## Cross-Modal Scene Retrieval Visualization: Failure



# **Experimental Results**

## **Missing Modalities**

# Experimental Results

## Missing Modalities

Ablation on instance matching on ScanNet with non-overlapping data per modality pair

Available Data		Instance Matching Recall ↑			
$\mathcal{I} \rightarrow \mathcal{P}$ (%)	$\mathcal{I} \rightarrow \mathcal{M}$ (%)	same	diff	top-1	top-3
25	75	86.32	<b>73.38</b>	55.46	79.73
50	50	<b>87.46</b>	70.02	57.49	79.94
75	25	87.35	67.65	54.99	79.45
100	100	87.44	72.46	<b>59.88</b>	<b>80.81</b>

- Although partial data availability decreases recall, our matching only decreases by 3% even when using 25% data
- Common in [real-world applications](#), where certain modalities might be [scarce!](#)

# **Future Work**

**Key Takeaways:**

# Future Work

## Key Takeaways:

- **Cross-modal alignment method** for 3D scenes that learns a unified, modality-agnostic embedding space, enabling a range of scene understanding tasks.
- CrossOver leverages a **unified embedding space** centered on image features, allowing it to generalize across **unpaired modalities** and **outperform existing methods** on real-world datasets

# Future Work

## Key Takeaways:

- **Cross-modal alignment method** for 3D scenes that learns a unified, modality-agnostic embedding space, enabling a range of scene understanding tasks.
- CrossOver leverages a **unified embedding space** centered on image features, allowing it to generalize across **unpaired modalities** and **outperform existing methods** on real-world datasets

## Current Directions:

# Future Work

## Key Takeaways:

- **Cross-modal alignment method** for 3D scenes that learns a unified, modality-agnostic embedding space, enabling a range of scene understanding tasks.
- CrossOver leverages a **unified embedding space** centered on image features, allowing it to generalize across **unpaired modalities** and **outperform existing methods** on real-world datasets

## Current Directions:

- How can we scale up cross-modal approaches on a scene level? Multiple challenges: noisy data, bad annotations, etc
- How can we **enable traditional 3D scene understanding** with CrossOver, eg, segmentation, 3D question answering & visual grounding?
- How our approach can be applied to **dynamic scene reconstruction and real-time navigation**, leading to interactive and immersive mixed-reality experiences?

# Future Work

## Key Takeaways:

- **Cross-modal alignment method** for 3D scenes that learns a unified, modality-agnostic embedding space, enabling a range of scene understanding tasks.
- CrossOver leverages a **unified embedding space** centered on image features, allowing it to generalize across **unpaired modalities** and **outperform existing methods** on real-world datasets

## Current Directions:

- How can we scale up cross-modal approaches on a scene level? Multiple challenges: noisy data, bad annotations, etc
- How can we **enable traditional 3D scene understanding** with CrossOver, eg, segmentation, 3D question answering & visual grounding?
- How our approach can be applied to **dynamic scene reconstruction and real-time navigation**, leading to interactive and immersive mixed-reality experiences?

## Stay Tuned!

- New models trained on much bigger datasets with zero-shot capabilities, releasing soon!

**Thank You!**