

Tackling View-Dependent Semantics in 3D Language Gaussian Splatting

Jiazhong Cen¹, Xudong Zhou¹, Jiemin Fang^{2*}, Changsong Wen¹, Lingxi Xie², Xiaopeng Zhang², Wei Shen^{1*}, Qi Tian²

¹MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University

²Huawei Technologies Co., Ltd.



Introduction

Integrating vision-language features into 3D-GS is a common practice for open-vocabulary perception.

However, existing methods typically operate on rendered 2D feature maps and lack direct point-level understanding in 3D space.

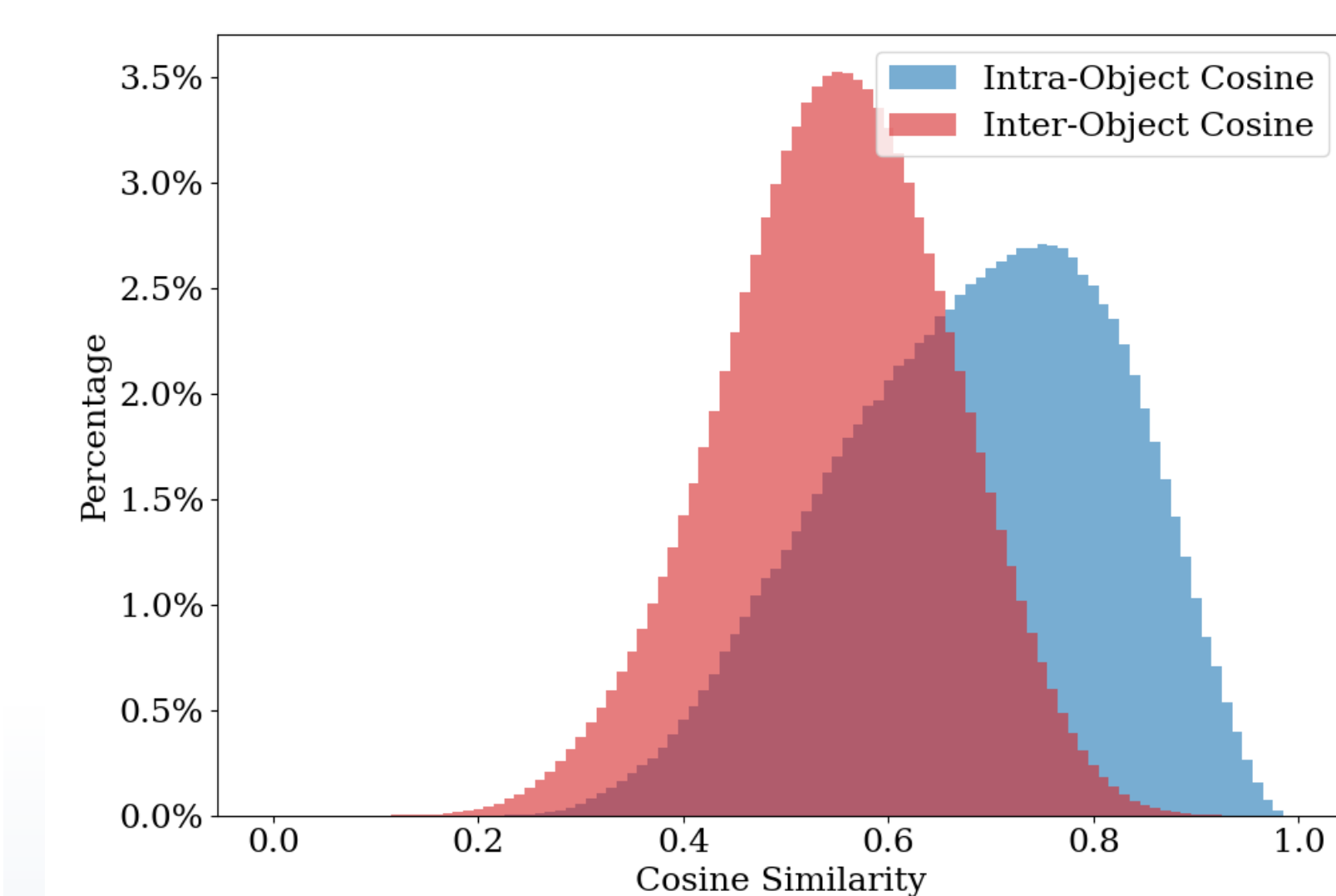
But why?

“横看成岭侧成峰
远近高低各不同”

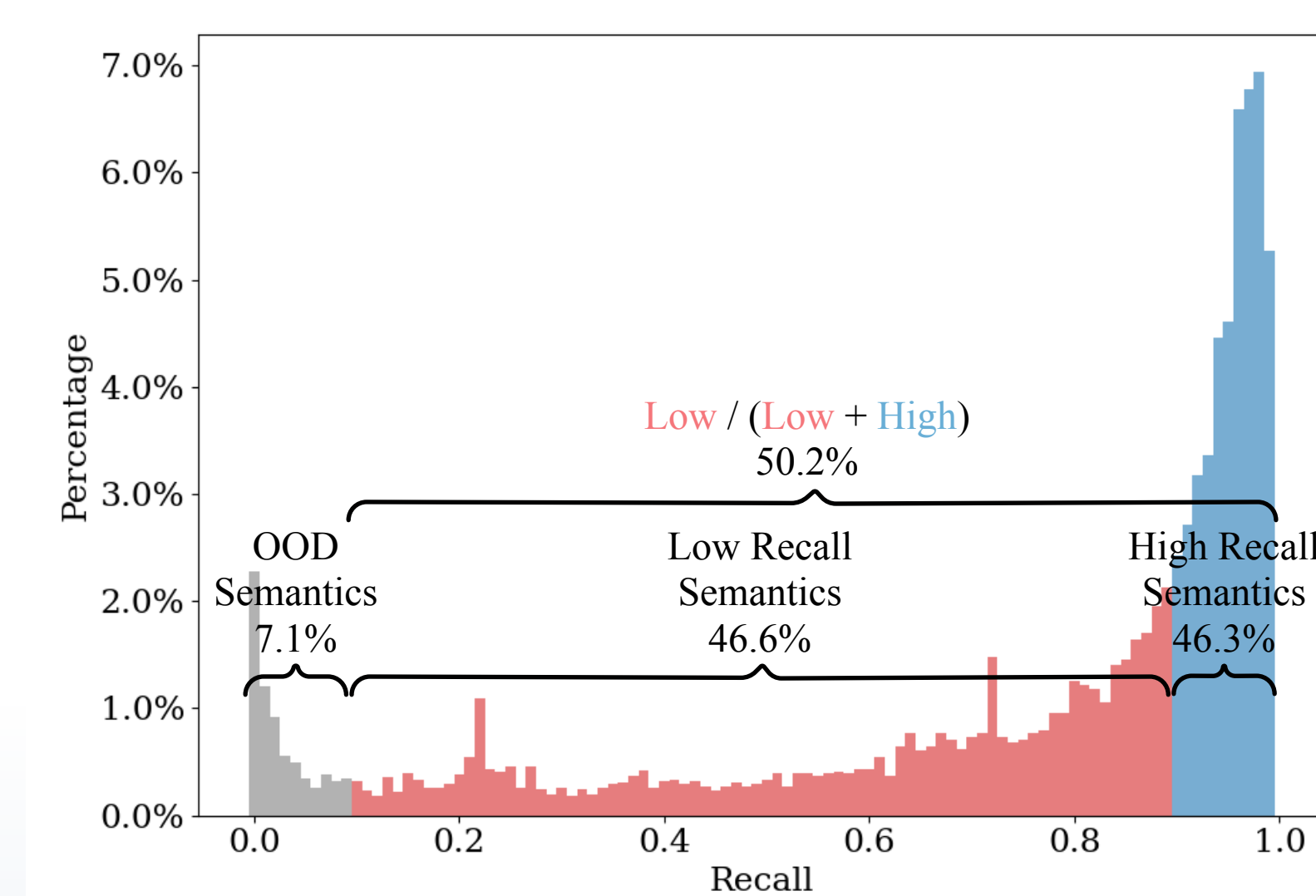
Consider a book, we can only recognize its name from limited viewpoints, indicating its semantics is **view-dependent**.

Existing methods ignore this characteristic, simply adopting 3D-GS to fit multi-view 2D semantic features! Thus leading to sub-optimal results.

Quantifying View-Dependency



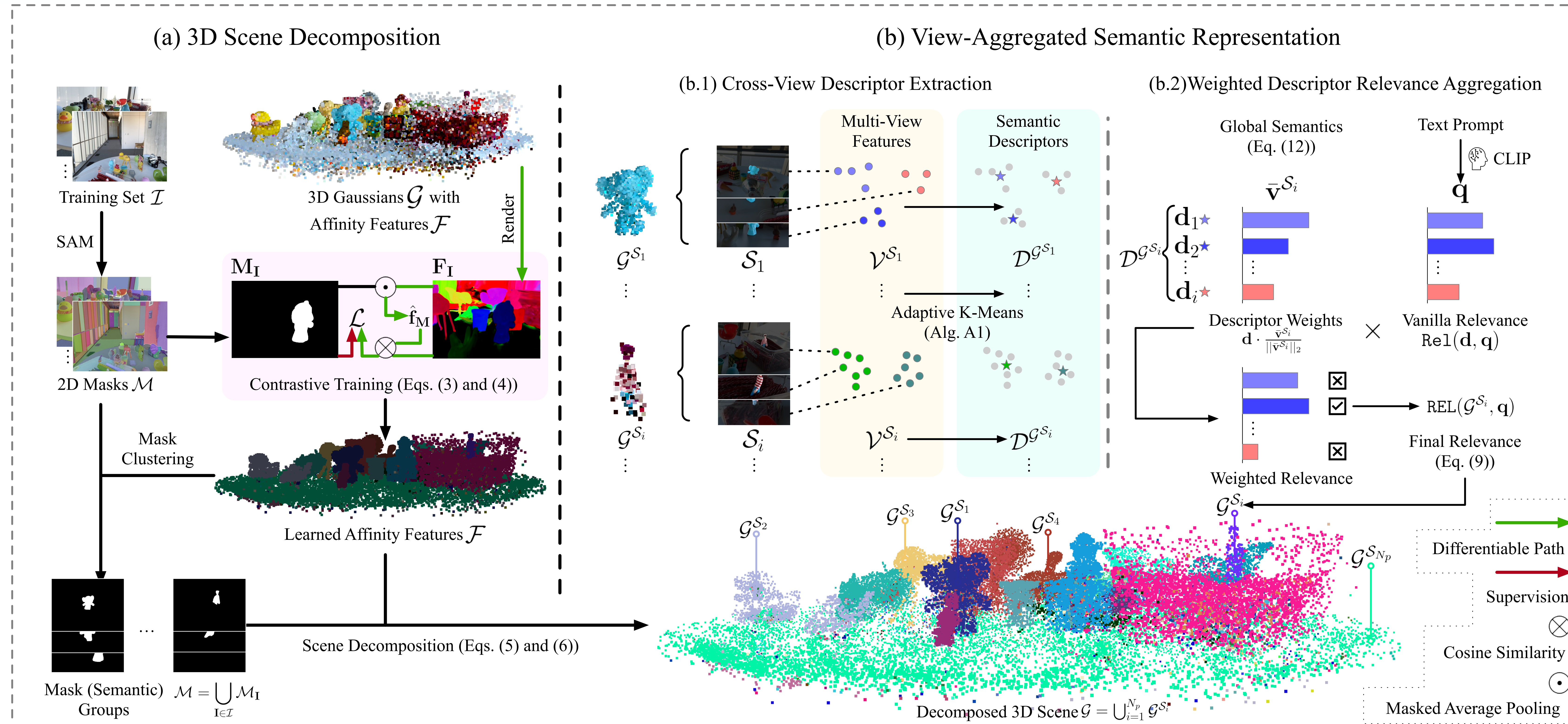
Cosine similarity of multi-view semantic features within the same object and across different objects: A substantial overlap between the inter- and intra-object distributions, suggesting inconsistent multi-view semantics of each object.



Using single view semantics to retrieve the corresponding object:

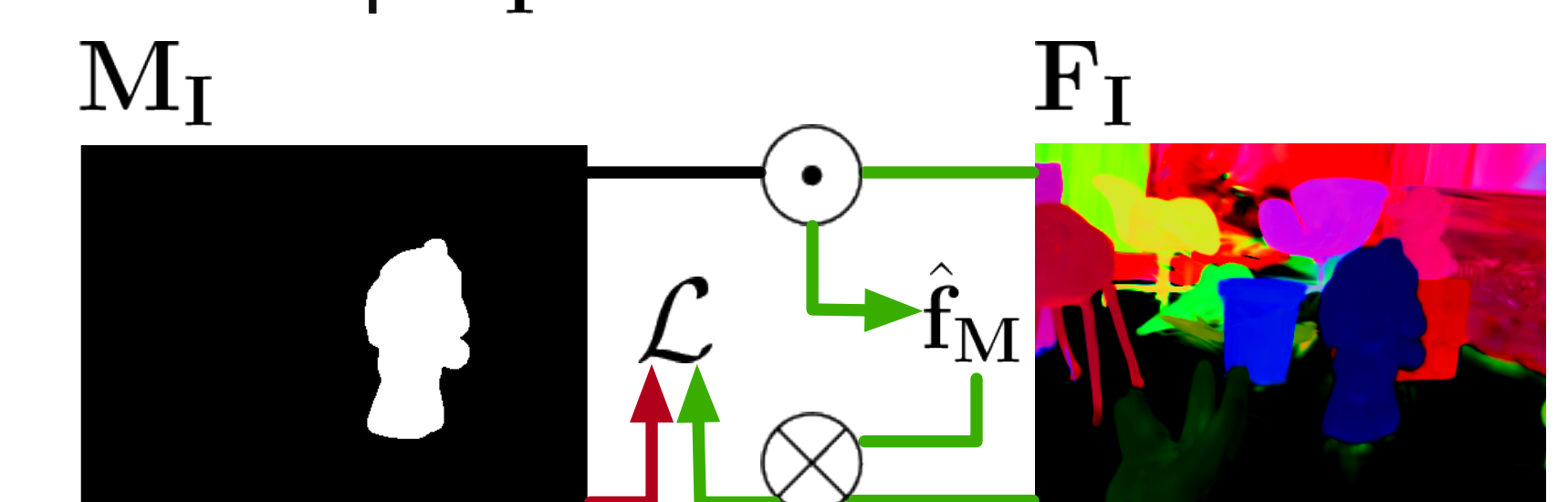
More than a half of semantic features can not completely retrieve its corresponding object.

Methodology



3D Scene Decomposition

LaGa first assign an affinity feature to each 3D Gaussian. During training, it renders the feature map \mathbf{F}_I



$$\hat{\mathbf{f}}_{M^I} = \text{MAP}(\mathbf{M}^I \mathbf{F}_I) = \frac{1}{\sum_{p \in \delta(I)} \mathbf{M}^I(p)} \sum_{p \in \delta(I)} \mathbf{M}^I(p) \mathbf{F}_I(p).$$

$$\mathcal{L} = \sum_{I \in \mathcal{I}} \sum_{M \in \mathcal{M}^I} \sum_{p \in \delta(I)} (1 - 2M(p)) \max(\langle \hat{\mathbf{f}}_M, \mathbf{F}_I(p) \rangle, 0).$$

LaGa conducts clustering on the learned mask prototypes $\hat{\mathbf{f}}_M$. The masks of each 3D object will be grouped into a same cluster S_i

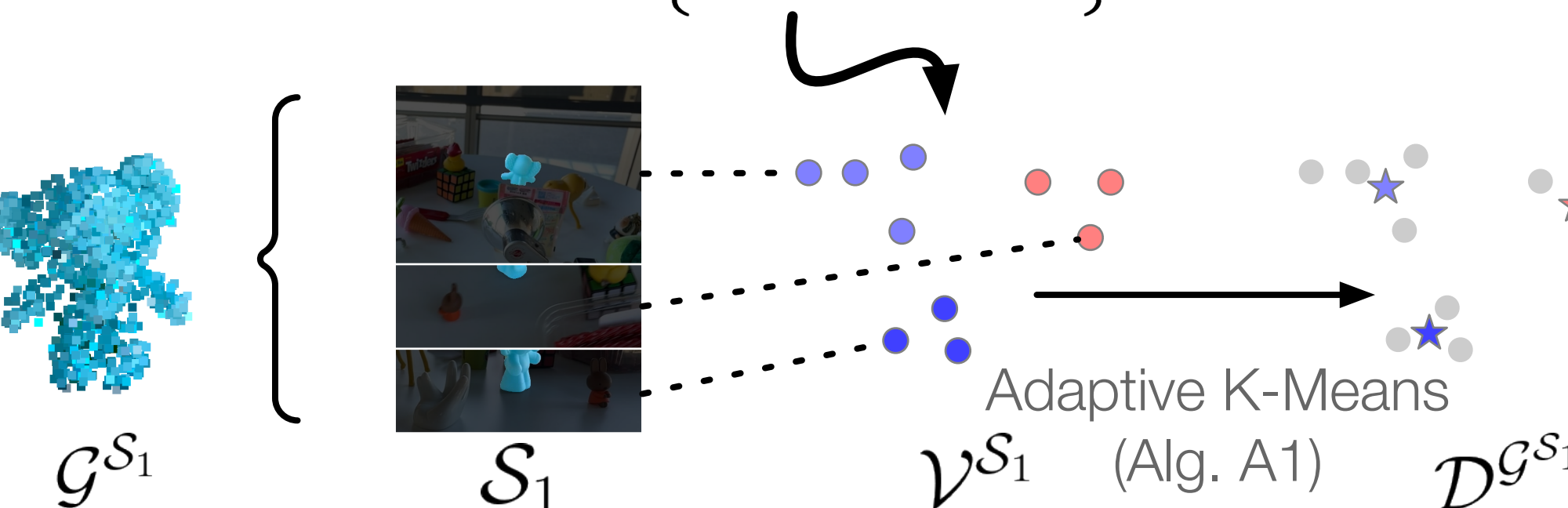
To decompose the 3D scene, LaGa derives classifier t^{S_i} from the mask clusters. The 3D Gaussian with \mathbf{f}_g belongs to the i^* -th object.

$$t^{S_i} = \frac{1}{|S_i|} \sum_{M \in S_i} \hat{\mathbf{f}}_M. \quad i^* = \arg \max_i \langle \mathbf{f}_g, t^{S_i} \rangle.$$

Cross-View Descriptor Extraction

For each 3D object, LaGa uses CLIP to extract its multi-view semantics:

$$\mathcal{V}^{S_i} = \{ \mathbf{v}^M \mid M \in S_i \}.$$



The multi-view semantics of different 3D objects exhibit varying levels of complexity.

To account for this, LaGa employs an **adaptive** K-Means algorithm (Alg. A1) on each object's semantic features to adaptively extract a set of descriptors:

$$\mathcal{D}^{\mathcal{G}^{S_i}} = \{ \mathbf{d}_i \in \mathbb{R}^{C'} \mid i \in \{1, \dots, N^{\mathcal{G}^{S_i}}\} \}.$$

Note: C' denotes the number of channels of affinity features and C' for the CLIP semantic features.

Weighted Descriptor Relevance Aggregation

After assigning a set of descriptors to each 3D object, LaGa adjust the importance of descriptors according to two metrics:

- 1. Directional Consistency:** reliable descriptors should have higher cosine similarity with the objects global feature.
- 2. Internal Compactness:** good descriptors should be generated by compact multi-view l2-normalized features. Its l2-norm indicates this compactness.

$$\omega^d = \mathbf{d} \cdot \frac{\bar{\mathbf{v}}^{S_i}}{\|\bar{\mathbf{v}}^{S_i}\|_2}$$

$$= \underbrace{\frac{\mathbf{d}}{\|\mathbf{d}\|_2} \cdot \frac{\bar{\mathbf{v}}^{S_i}}{\|\bar{\mathbf{v}}^{S_i}\|_2}}_{(i) \text{ Directional Consistency}} \times \underbrace{\|\mathbf{d}\|_2}_{(ii) \text{ Internal Compactness}}$$

During inference, given the text prompt \mathbf{q} , the highest weighted relevance score of the object is regarded as its final score:

$$\text{REL}(\mathcal{G}^{S_i}, \mathbf{q}) = \max_{d \in \mathcal{D}^{\mathcal{G}^{S_i}}} \omega^d \cdot \text{Rel}(\mathbf{d}, \mathbf{q}).$$

Quantitative Results

Results on the LERF-OVS dataset
(* denotes concurrent preprints)

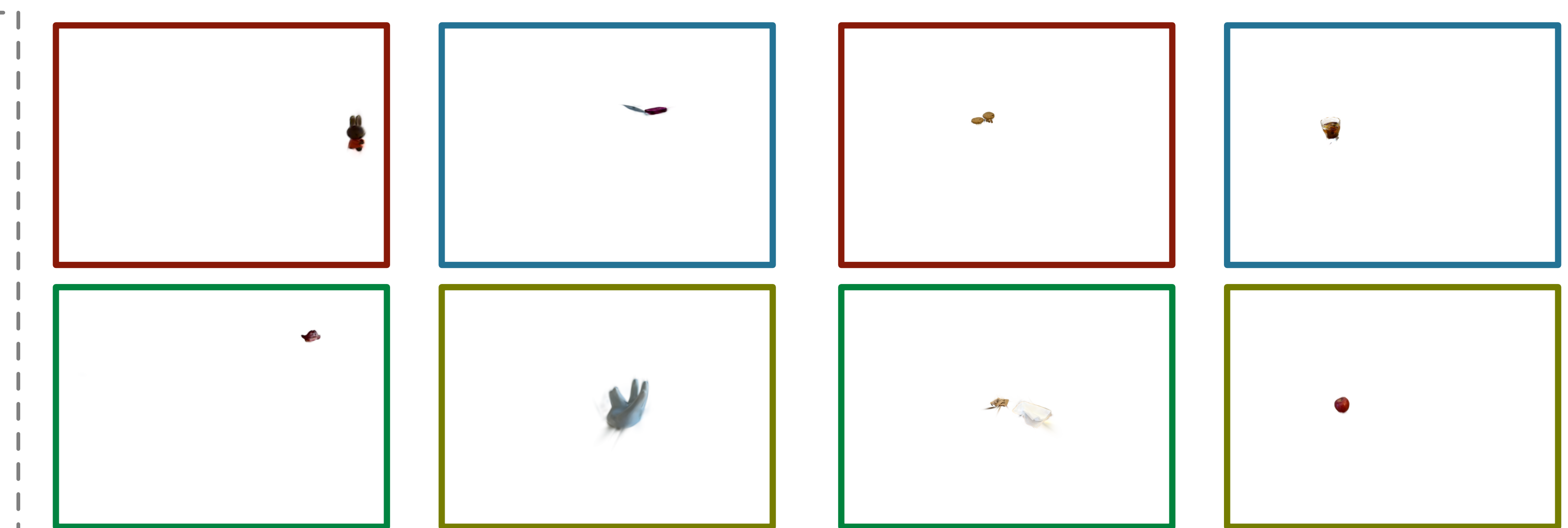
Under the same experiment setting, LaGa surpasses previous method for **+18.7% mIoU**.

	METHODS	F.	T.	R.	W.	MEAN
2D	LSEG	7.6	21.7	7.0	29.9	16.6
	LERF	38.6	45.0	28.2	37.9	37.4
	LEGAUSSIANS	60.3	44.5	52.6	41.4	46.9
	LANGSPLAT	44.7	65.1	51.2	44.5	51.4
	N2F2	47.0	69.2	56.6	47.9	54.4
	OCCAMLGS*	58.6	70.2	51.0	65.3	61.3
3D	VLGS*	58.1	73.5	61.4	54.8	62.0
	OPENGAUSSIAN [†]	39.3	60.4	31.0	22.7	38.4
	SAGA [‡]	36.2	19.3	53.1	14.4	30.7
	LANGSPLAT [‡]	25.9	35.6	29.3	33.5	31.1
	LEGAUSSIANS [‡]	31.2	34.5	17.6	17.3	25.2
	OPENGAUSSIAN [‡]	61.1	59.1	29.2	31.9	45.3
	SUPERGSEG*	43.7	55.3	18.1	26.7	35.9
	LAGA (OURS)	64.1	70.9	55.6	65.6	64.0

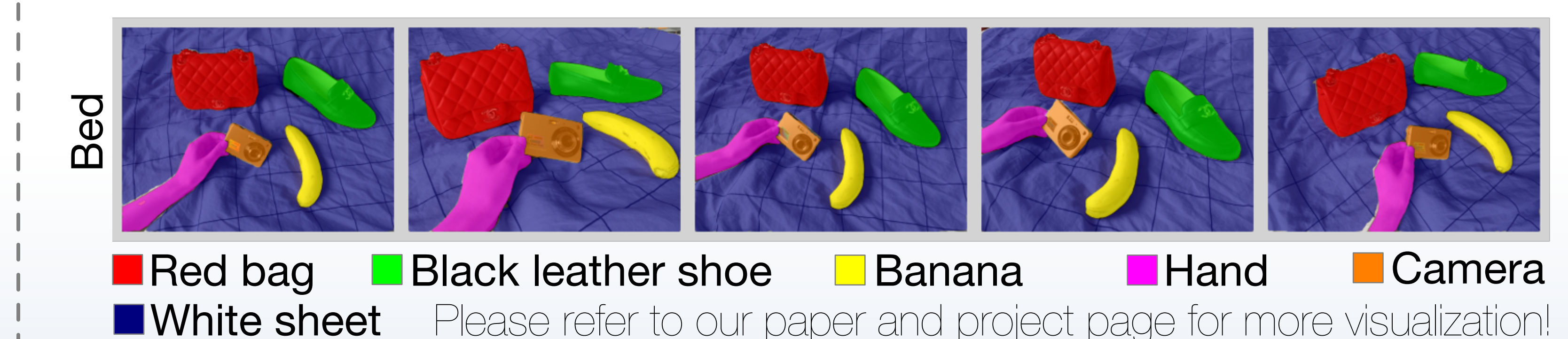
'2D' denotes conducting segmentation on rendered feature map. '3D' denotes conduct 3D segmentation then rendering the results to different views

Qualitative Results

Results on the LERF-OVS dataset



Results on the 3D-OVS dataset



Conclusion

This paper investigate the **view-dependency of multi-view semantics** in 3D objects—an issue ignored by prior methods. We propose **LaGa**, which **decomposes scenes into 3D objects** and **gathers their multi-view semantics adaptively**. Extensive experiments demonstrate its effectiveness.