



西安电子科技大学  
XIDIAN UNIVERSITY



清华大学 智能产业研究院  
Institute for AI Industry Research, Tsinghua University

# Delving into Shape-aware Zero-shot Semantic Segmentation

**Xinyu Liu<sup>1,2</sup>, Beiwen Tian<sup>2,4</sup>, Zhen Wang<sup>3</sup>, Rui Wan<sup>3</sup>, Kehua Sheng<sup>3</sup>, Bo Zhang<sup>3</sup>,  
Hao Zhao<sup>2</sup>, Guyue Zhou<sup>2</sup>**

<sup>1</sup>Xidian University, China; <sup>2</sup>Institute for AI Industry Research (AIR), Tsinghua University, China; <sup>3</sup>Didi Chuxing, China; <sup>4</sup>Department of Computer Science and Technology, Tsinghua University, China

liuxinyu@stu.xidian.edu.cn, zhaohao@air.tsinghua.edu.cn

# Content

**1. Introduction**

**2. Method**

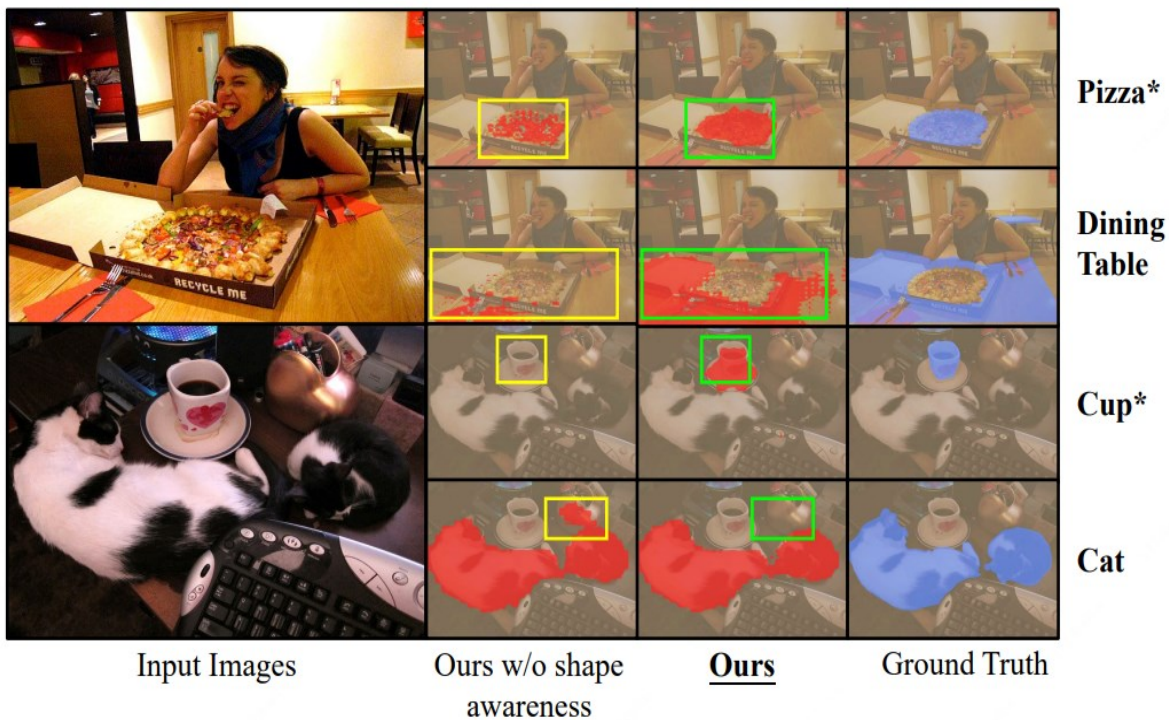
**3. Experiment**

**4. Conclusion**



# Introduction

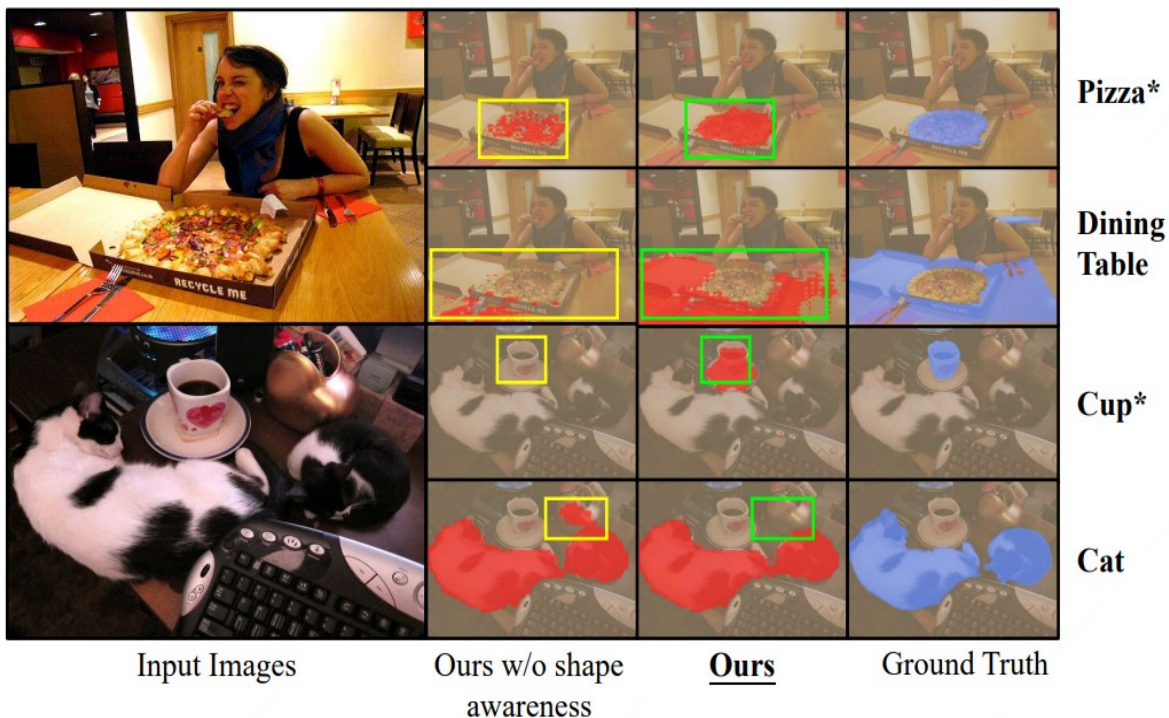
## Problems:



- In real-world, unknown objects may appear, hence there is a need to learn and achieve dense predictions for unknown object categories from limited samples.
- Dense prediction task requires not only accurate semantic understanding but also fine shape delineation.
- However, existing vision-language models are trained with image-level language descriptions.

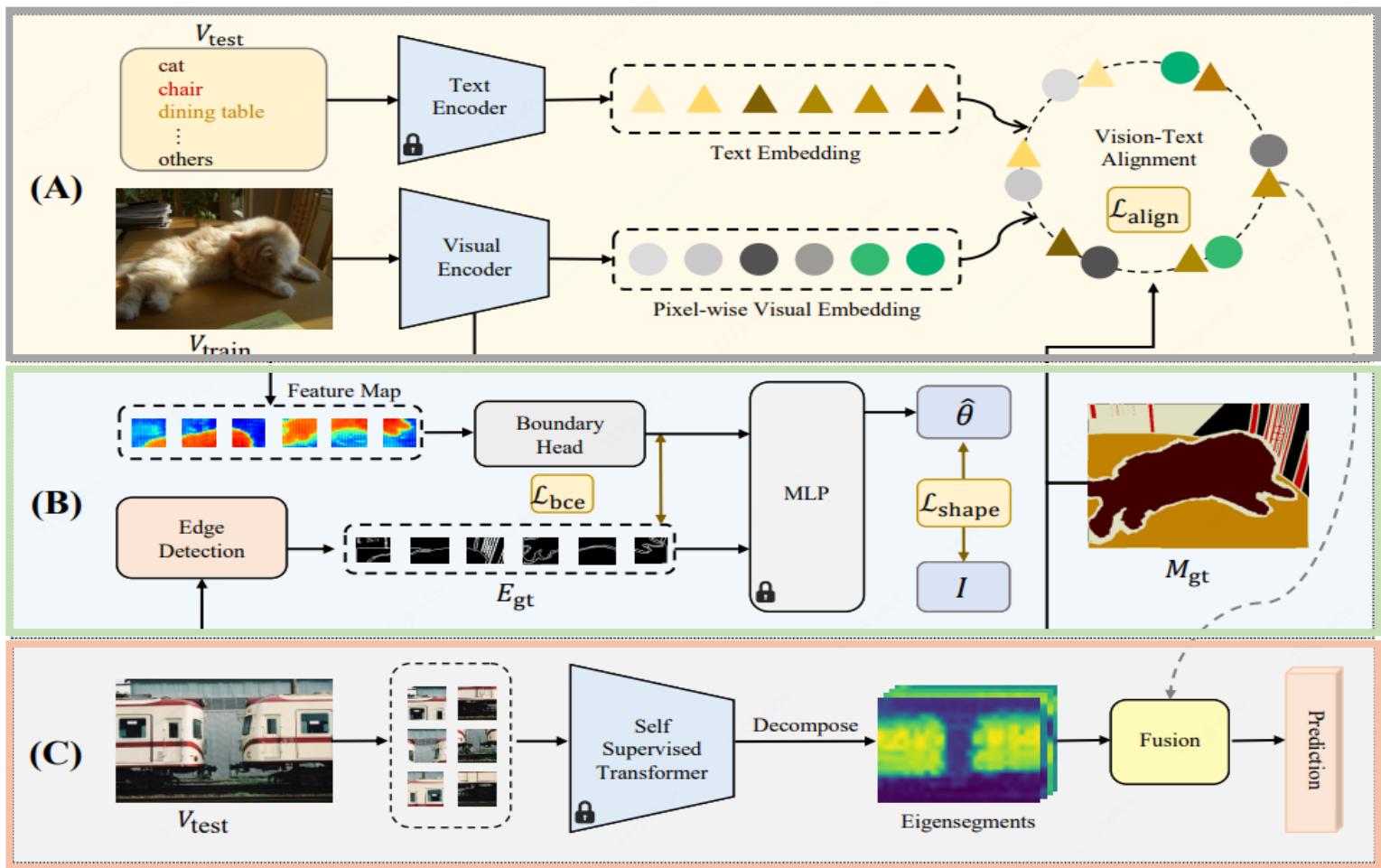
# Introduction

## Solutions:



- Leverage the eigen vectors of Laplacian matrices constructed with self-supervised pixelwise features to promote shape-awareness.
- Propose to jointly optimizing a boundary segment constraint that aligns both the boundary of predicted semantic regions and the ground truth regions.

# Method



Overall architecture

# Experiment

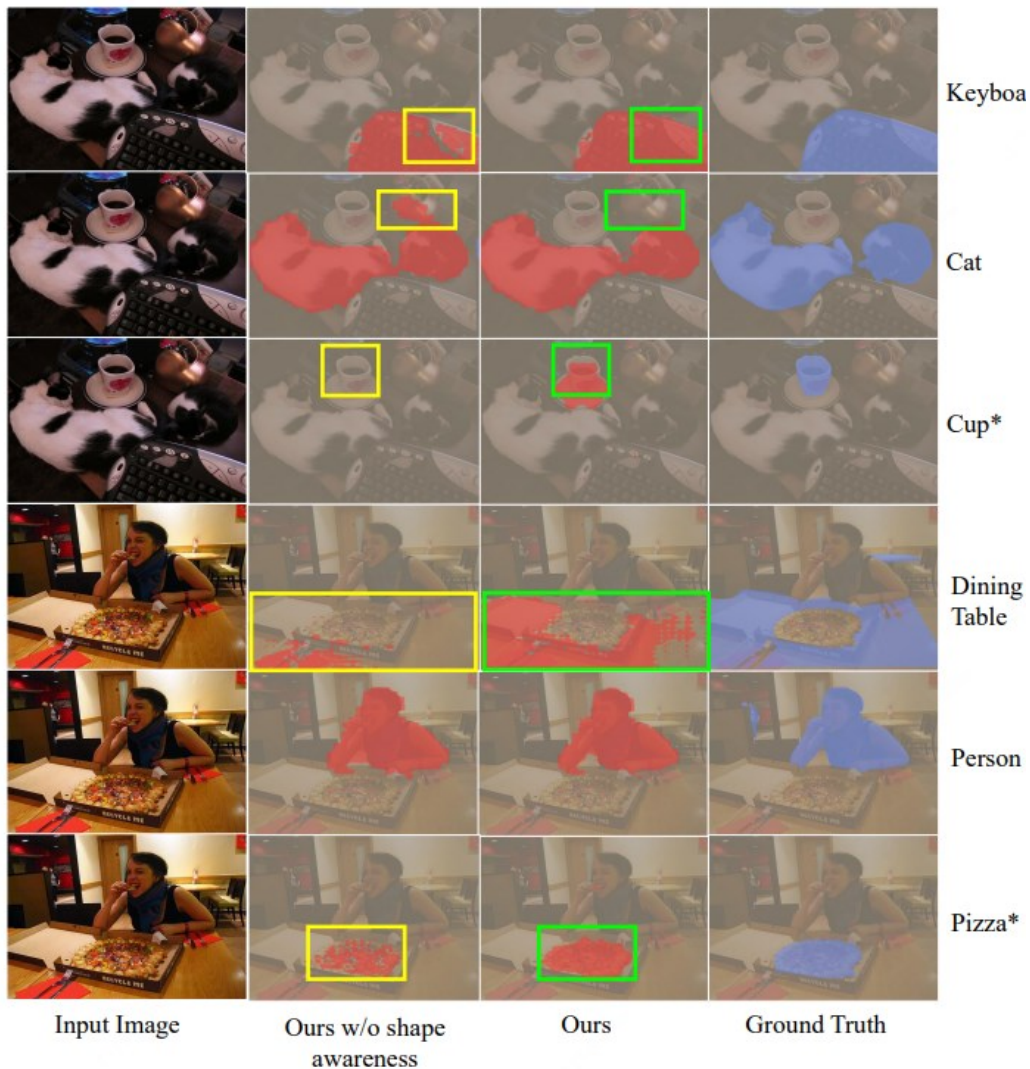
Our results on PASCAL-5<sup>i</sup> and COCO-20<sup>i</sup> datasets

Method	Backbone	Setting	PASCAL-5 <sup>i</sup>						COCO-20 <sup>i</sup>					
			5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	mIoU	FBIoU	20 <sup>0</sup>	20 <sup>1</sup>	20 <sup>2</sup>	20 <sup>3</sup>	mIoU	FBIoU
FWB [38]	ResNet	1-shot	51.3	64.5	56.7	52.2	56.2	—	17.0	18.0	21.0	28.9	21.2	—
DAN [52]	ResNet	1-shot	54.7	68.6	57.8	51.6	58.2	71.9	—	—	—	—	24.4	62.3
PFENet [51]	ResNet	1-shot	60.5	69.4	54.4	55.9	60.1	72.9	36.8	41.8	38.7	36.7	38.5	63.0
HSNet [37]	ResNet	1-shot	67.3	72.3	62.0	63.1	66.2	77.6	37.2	44.1	42.4	41.3	41.2	69.1
SPNet [54]	ResNet	zero-shot	23.8	17.0	14.1	18.3	18.3	44.3	—	—	—	—	—	—
ZS3Net [4]	ResNet	zero-shot	40.8	39.4	39.3	33.6	38.3	57.7	18.8	20.1	24.8	20.5	21.1	55.1
LSeg [25]	ResNet	zero-shot	52.8	53.8	44.4	38.5	47.4	64.1	22.1	25.1	24.9	21.6	23.4	57.9
<b>Ours</b>	DRN	zero-shot	<b>57.3</b>	<b>60.3</b>	<b>58.4</b>	<b>45.9</b>	<b>55.5</b>	<b>66.4</b>	<b>34.2</b>	<b>36.5</b>	<b>34.6</b>	<b>35.6</b>	<b>35.2</b>	<b>58.4</b>
LSeg [25]	ViT-L	zero-shot	61.3	63.6	43.1	41.0	52.3	67.6	28.1	27.5	30.0	23.2	27.2	<b>59.9</b>
<b>Ours</b>	ViT-L	zero-shot	<b>62.7</b>	<b>64.3</b>	<b>60.6</b>	<b>50.2</b>	<b>59.4</b>	<b>69.0</b>	<b>33.8</b>	<b>38.1</b>	<b>34.4</b>	<b>35.0</b>	<b>35.3</b>	58.2

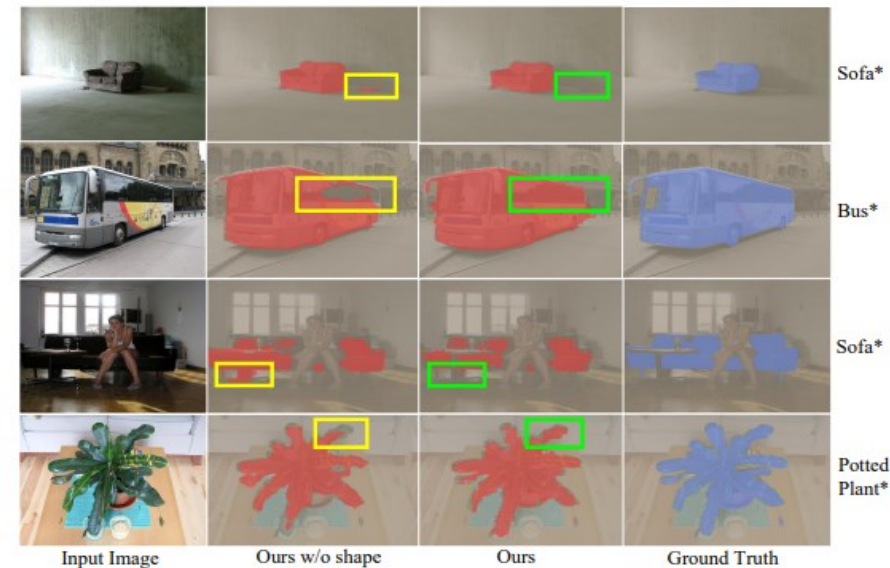
The cross dataset results

Model	Backbone	external dataset	target dataset	PASCAL-5 <sup>i</sup>
LSeg	ViT-L	✗	✓ (seen classes)	52.3
SPNet	ResNet	✗	✓ (seen classes)	18.3
ZS3Net	ResNet	✗	✓ (seen classes)	38.3
LSeg	ResNet	✗	✓ (seen classes)	47.4
LSeg+	ResNet	COCO	✗	59.0
OpenSeg [14]	ResNet	COCO	✗	60.0
<b>Ours</b>	DRN	COCO	✗	<b>62.7</b>

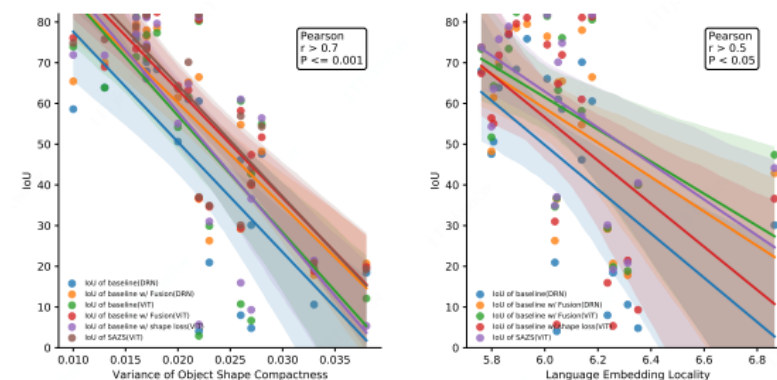
# Experiment



Our results on COCO-20<sup>i</sup> dataset



Our results on PASCAL-5<sup>i</sup> dataset



Correlation of CO variance and mean language embedding locality on IoU

# Experiment

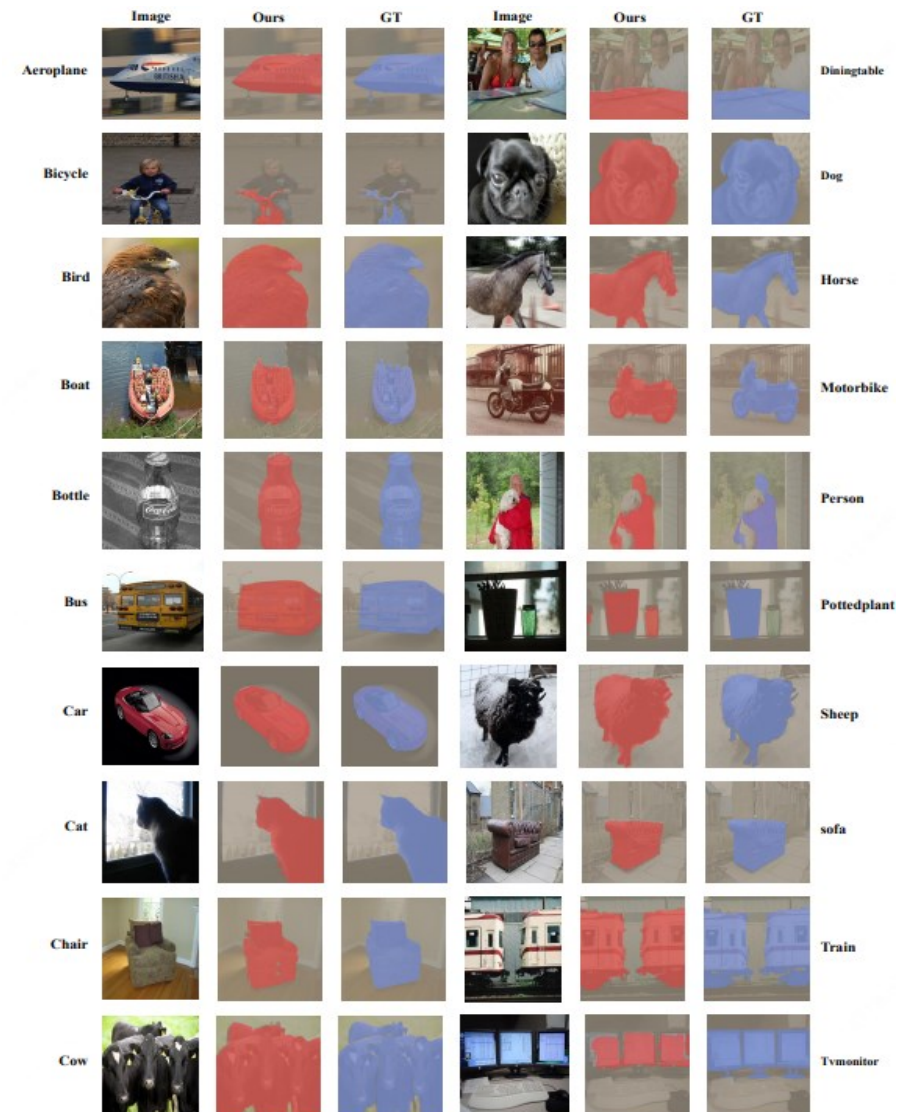


清华大学 智能产业研究院

Institute for AI Industry Research, Tsinghua University



Our results on COCO-20<sup>i</sup> dataset



Our results on PASCAL-5<sup>i</sup> dataset



# Experiment

Ablation study on COCO-20<sup>i</sup> (ViT backbone)

Model	Fusion	$\mathcal{L}_{\text{shape}}$	20 <sup>0</sup>	20 <sup>1</sup>	20 <sup>2</sup>	20 <sup>3</sup>	mIoU
SAZS	✓	✓	<b>34.2</b>	36.5	<b>34.6</b>	<b>35.6</b>	<b>35.2</b>
SAZS	✓		33.7	<b>38.2</b>	33.4	35.5	35.2
SAZS		✓	28.4	27.6	25.4	25.1	26.6
SAZS			24.2	28.5	24.4	23.3	25.1
LSeg [25]			22.1	25.1	24.9	21.6	23.4

Ablation study on COCO-20<sup>i</sup> (DRN backbone)

Model	Fusion	$\mathcal{L}_{\text{shape}}$	20 <sup>0</sup>	20 <sup>1</sup>	20 <sup>2</sup>	20 <sup>3</sup>	mIoU
SAZS	✓	✓	<b>33.8</b>	38.1	<b>34.4</b>	<b>35.0</b>	<b>35.3</b>
SAZS	✓		33.3	<b>39.0</b>	33.9	32.7	34.7
SAZS		✓	30.0	30.4	27.5	28.5	29.1
SAZS			26.3	32.0	26.2	26.2	27.7
LSeg [25]			28.1	27.5	30.0	23.2	27.2

Ablation study on PASCAL-5<sup>i</sup> (ViT backbone)

Model	Fusion	$\mathcal{L}_{\text{shape}}$	5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	mIoU
SAZS	✓	✓	62.7	<b>64.3</b>	<b>60.6</b>	<b>50.2</b>	<b>59.4</b>
SAZS	✓		<b>63.1</b>	62.4	59.0	49.2	58.4
SAZS		✓	59.7	63.4	44.3	42.2	52.4
SAZS			59.2	61.9	43.8	41.9	51.7
LSeg [25]			61.3	63.6	43.1	41.0	52.3

Impact of  $Z_{\text{shape}}$  and  $Z_{\text{sem}}$  of fusion module

Model	external dataset	$Z_{\text{shape}}$	$Z_{\text{sem}}$	PASCAL-5 <sup>i</sup>
SAZS	COCO			58.4
SAZS	COCO	✓		58.6
SAZS	COCO		✓	62.7

## Our contributions:

- We propose a novel methodology for shape-aware zero-shot semantic segmentation models from language supervision, named SAZS, and set new state-of-the-art performance on PASCAL and COCO datasets.
- We propose to leverage the eigen vectors of Laplacian matrices constructed with self-supervised pixelwise features to promote shape-awareness.
- We propose to jointly optimizing a boundary segment constraint that aligns both the boundary of predicted semantic regions and the ground truth regions while narrowing the visual-language embedding space.
- We draw several interesting and conclusive observations: the benefits of promoting shape-awareness highly relates to objects' compactness and language embedding locality.