

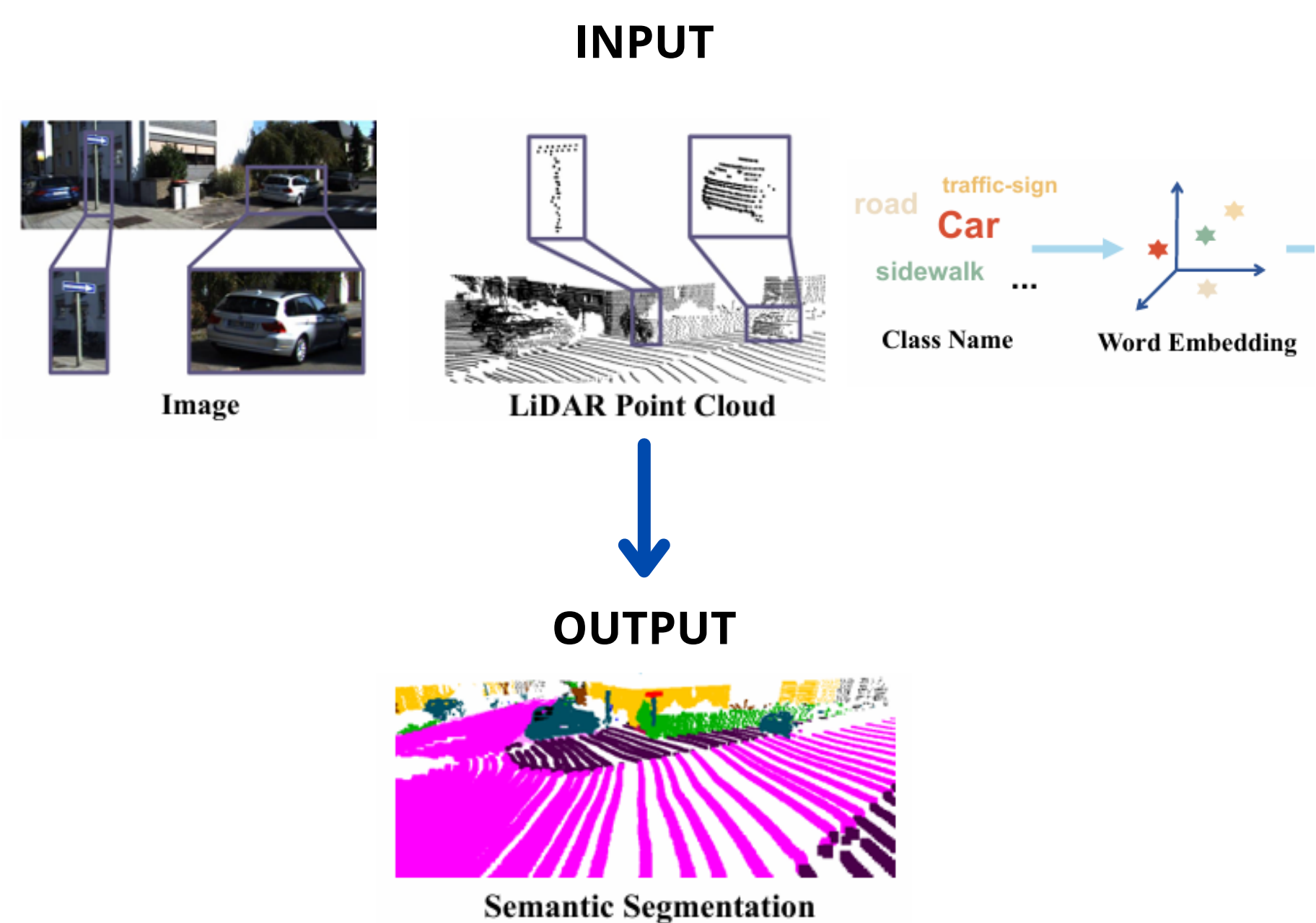
# See More and Know More: Zero-shot Point Cloud Segmentation via Multi-modal Visual Data

Yuhang Lu, Qi Jiang, Runnan Chen, Yuenan Hou, Xinge Zhu, Yuexin Ma



## GOAL

Point cloud semantic segmentation over:  
**seen & unseen objects**

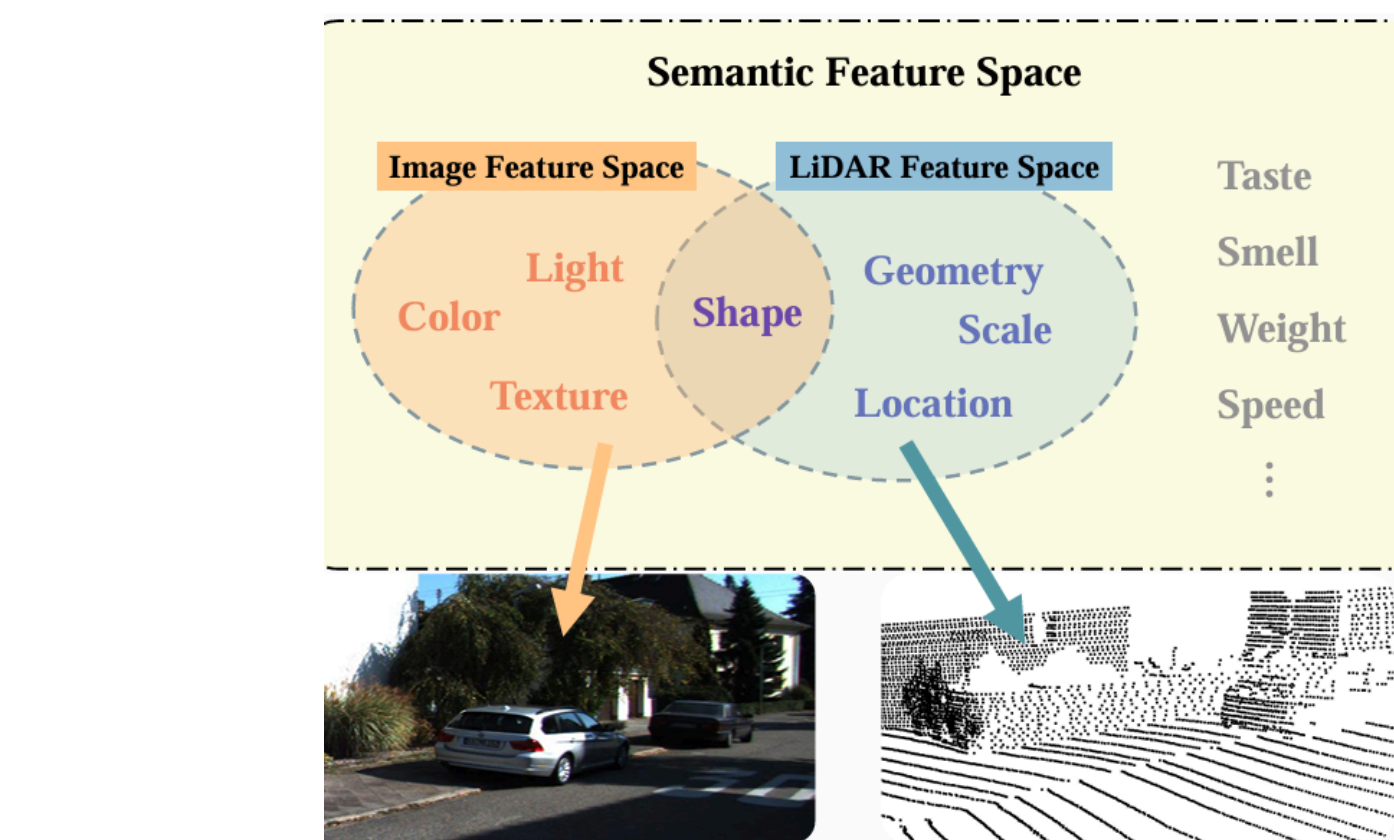


## INTRODUCTION

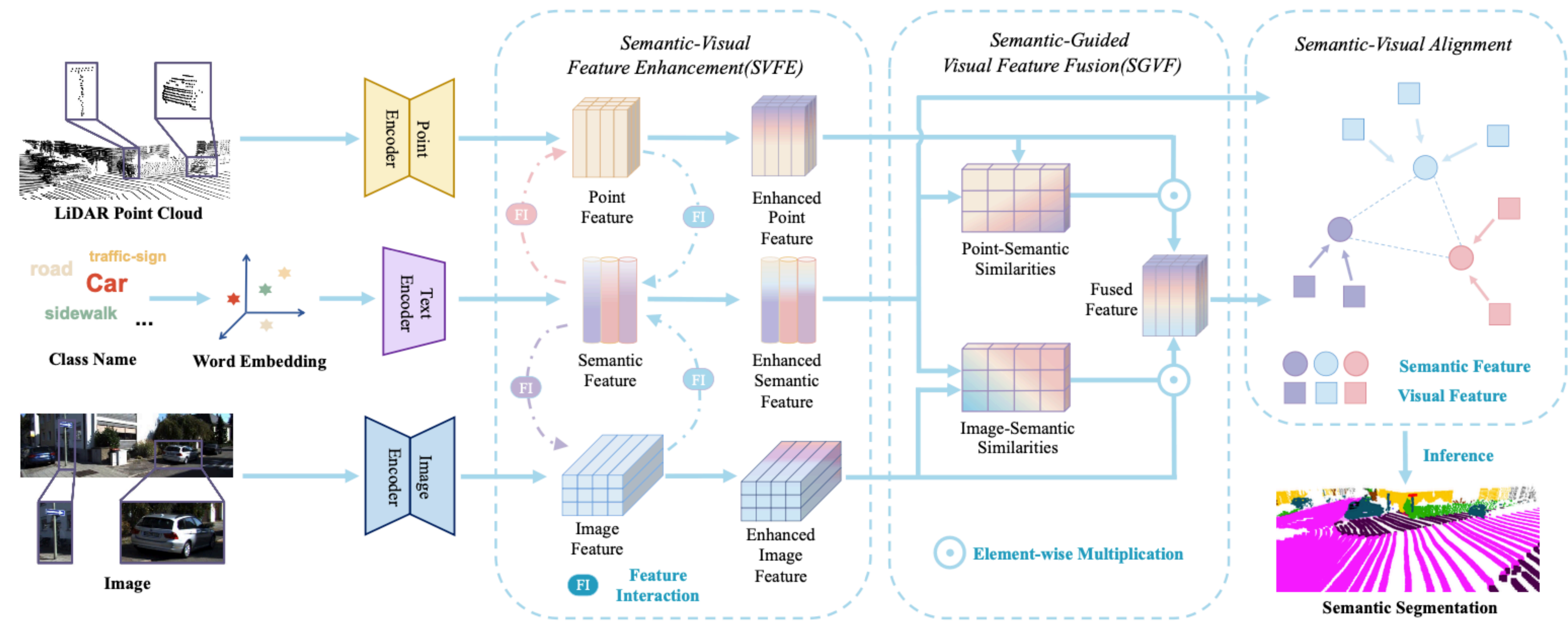
- Point cloud **semantic segmentation** over: seen & unseen objects
- Train a model** to recognize unseen objects
- Useful for **autonomous driving**

## CHALLENGES

- Hard to generalize** over unlabelled training set of 3D data
- Manual labelling** is infeasible
- Few** 3D semantic segmentation models
- Generative methods
  - fake features generator (3DGenZ)
- Projection-based methods alignment (TGF)



## MODEL



### Visual-Semantic Feature Extraction

- Cylinder3D**, to extract 3D features from point cloud
- ResUnet**, to extract 2D features from image
- W2V**, **Glove**, followed by **MLP** to map the word embedding to its semantic feature

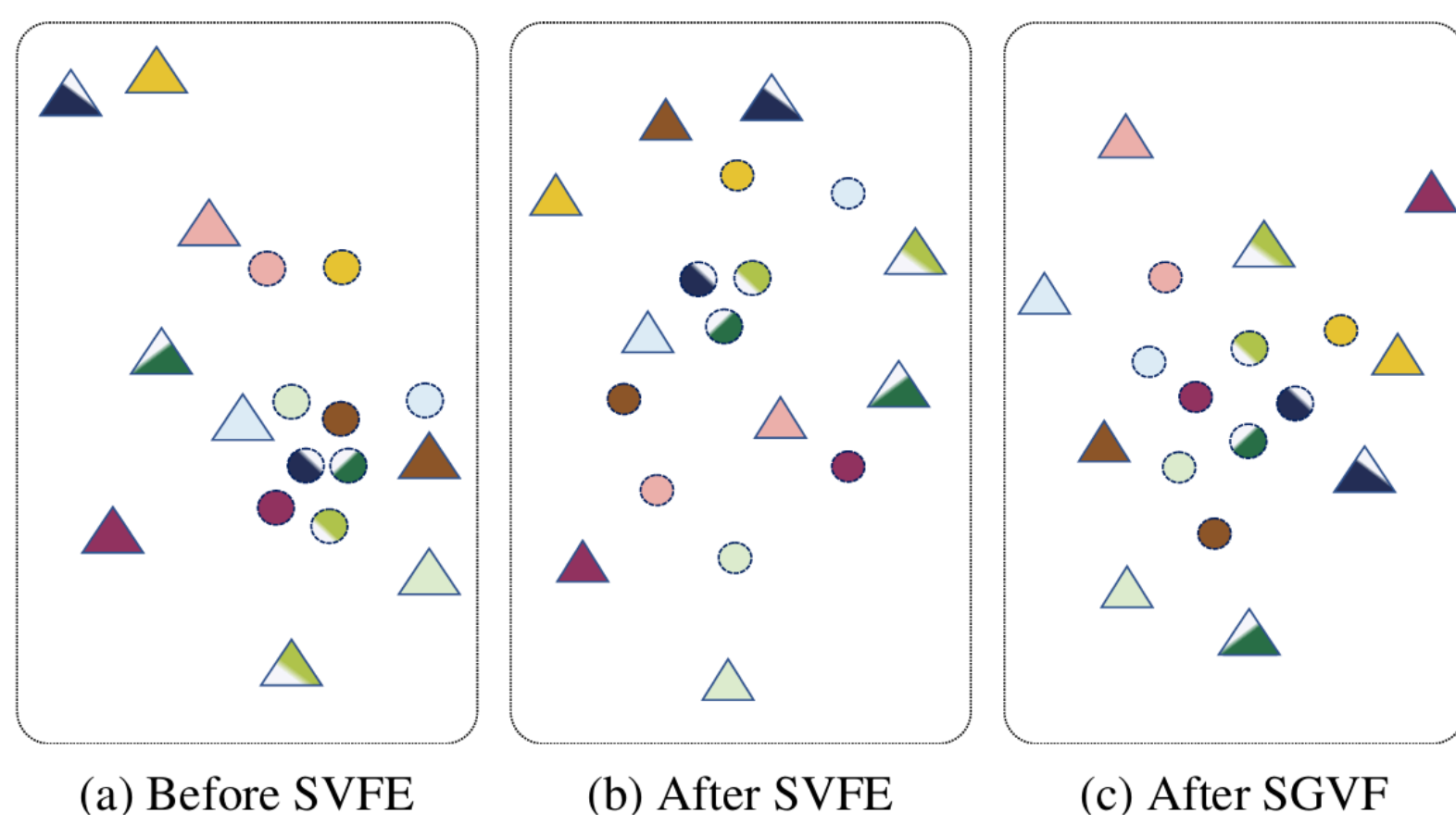
### Semantic-Visual Feature Enhancement

- Transformer Decoder with cross-attention layer to decrease the gap between visual and semantic features
- Semantic Feature Enhancement**, semantic features as query and visual features as key and value
  - Visual Feature Enhancement**, visual features as query and semantic features as key and value

### Semantic-Guided Visual Feature Fusion

Combine effectively visual features under the guide of the semantic features

- Weight matrices**, multihead-attention of enhanced visual features and enhanced semantic features
- Element-wise multiplication** between weight matrix and enhanced visual feature
- MLP with Softmax**



### Semantic-Visual Alignment

- Align semantic and visual feature spaces under seen classes supervision
- Loss function for seen classes, to have compact distribution within classes and distinguishable one between classes
  - Loss function for unseen classes, to push unseen classes' features apart to seen ones (seen bias)

### Inference

The class of each point is determined by the similarity between its fused visual feature and semantic features of all classes (seen and unseen)

## EXPERIMENTS

Datasets:

- SemanticKITTI
- nuScenes

Evaluation metrics:

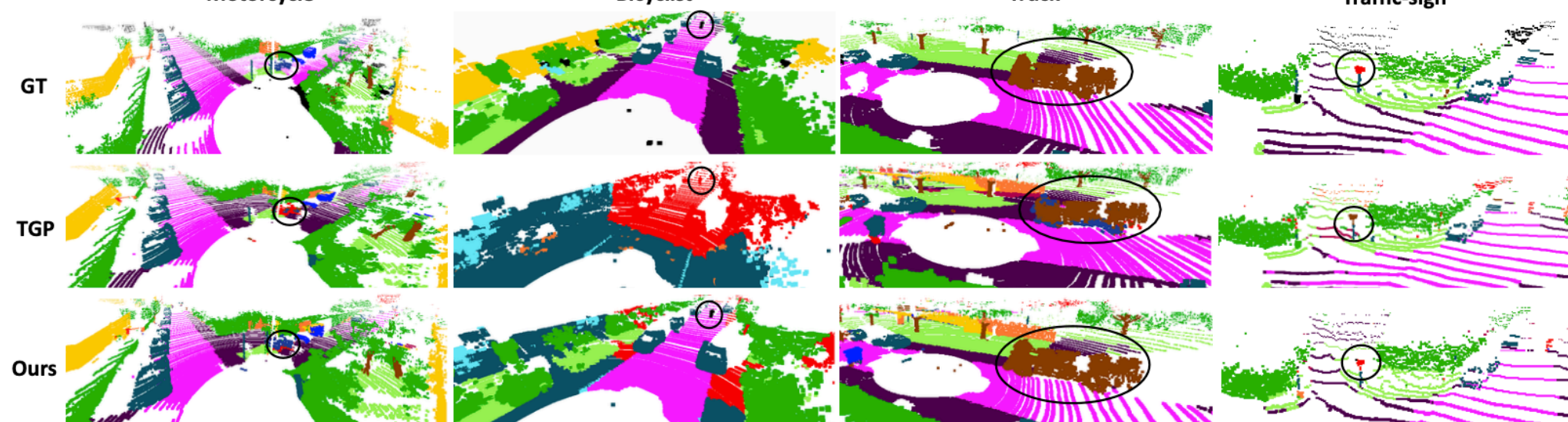
- harmonic mean IoU

$$hIoU = \frac{2 \times mIoU_{seen} \times mIoU_{unseen}}{mIoU_{seen} + mIoU_{unseen}}$$

## RESULTS

Our model outperforms SOTA with **more than 50% improvement** of unseen category mIOU. Furthermore, it yields real-time performance (0.097 s/f)

Setting	Model	SemanticKITTI				nuScenes			
		Seen mIoU	Unseen mIoU	Improvement	Overall mIoU	Seen mIoU	Unseen mIoU	Improvement	Overall mIoU
0	TGP[15]	-	-	-	59.1	-	-	-	67.9
	Ours	-	-	-	<b>62.6</b>	-	-	-	<b>69.1</b>
2	3DGenZ[45]	40.9	12.4	-	37.9	19.0	67.8	4.2	59.9
	TGP[15]	58.3	28.8	+3.5%	55.2	38.6	58.9	26.9	+25.3%
4	Ours	<b>59.5</b>	<b>29.8</b>	-	<b>56.4</b>	<b>39.7</b>	59.4	<b>33.7</b>	56.2
	Supervised	61.5	71.8	-	62.6	66.3	70.1	61.9	69.1
6	3DGenZ[45]	41.4	10.8	-	35.0	17.1	67.2	3.1	51.2
	TGP[15]	54.6	17.3	+54.9%	46.7	26.3	65.7	14.8	53.0
8	Ours	<b>58.8</b>	<b>26.8</b>	-	<b>52.1</b>	<b>36.8</b>	66.4	<b>23.1</b>	<b>55.6</b>
	Supervised	60.3	71.2	-	62.6	65.3	71.9	60.6	69.1
10	3DGenZ[45]	40.3	6.5	-	29.6	11.2	53.8	3.2	34.8
	TGP[15]	53.6	13.3	+79.7%	40.9	21.3	68.8	14.1	48.3
12	Ours	<b>56.6</b>	<b>23.9</b>	-	<b>46.3</b>	<b>33.6</b>	66.8	<b>22.1</b>	<b>50.0</b>
	Supervised	56.8	75.3	-	62.6	64.8	74.5	60.1	69.1
14	3DGenZ[45]	38.3	1.3	-	22.7	2.5	36.5	2.1	19.3
	TGP[15]	<b>53.2</b>	<b>8.6</b>	+70.9%	<b>34.4</b>	<b>14.8</b>	<b>68.4</b>	<b>13.7</b>	<b>41.1</b>
16	Ours	46.0	<b>14.7</b>	-	32.8	<b>22.3</b>	68.2	<b>21.5</b>	<b>44.9</b>
	Supervised	52.1	77.1	-	62.6	62.2	73.5	64.7	69.1



## ABLATION

Model	Seen mIoU	Unseen mIoU	Overall	
			mIoU	hIoU
Ours	58.8	<b>26.8</b>	<b>52.1</b>	<b>36.8</b>
Ours w/o SGVF	58.8	23.4	51.3	33.5
Ours w/o SVFE	<b>59.0</b>	19.9	50.8	29.8
Ours w/o Image	58.3	20.0	50.2	29.8