

# Incorporating depth information into few-shot semantic segmentation

*Yifei Zhang<sup>1\*</sup>, Désiré Sidibé<sup>2</sup>, Olivier Morel<sup>1</sup>, Fabrice Meriaudeau<sup>1</sup>*

*<sup>1</sup>ERL VIBOT CNRS 6000, ImViA, Université Bourgogne Franche Comté, 71200, Le Creusot, France*

*<sup>2</sup>Université Paris-Saclay, Univ Evry, IBISC, 91020, Evry, France*

*[\\*Yifei.Zhang@u-bourgogne.fr](mailto:Yifei.Zhang@u-bourgogne.fr)*

ICPR 2020, Milan, Italy



- Introduction on multimodal few-shot segmentation
- Proposed model
- Dataset
- Experimental results
- Conclusion and future work

# Introduction on multimodal few-shot segmentation

Few-shot semantic segmentation presents a significant challenge for semantic scene understanding under limited supervision. Namely, this task targets at generalizing the segmentation ability of the model to new categories given a few samples. In order to obtain complete scene understanding, we extend the RGB-centric methods to take advantage of complementary depth information.

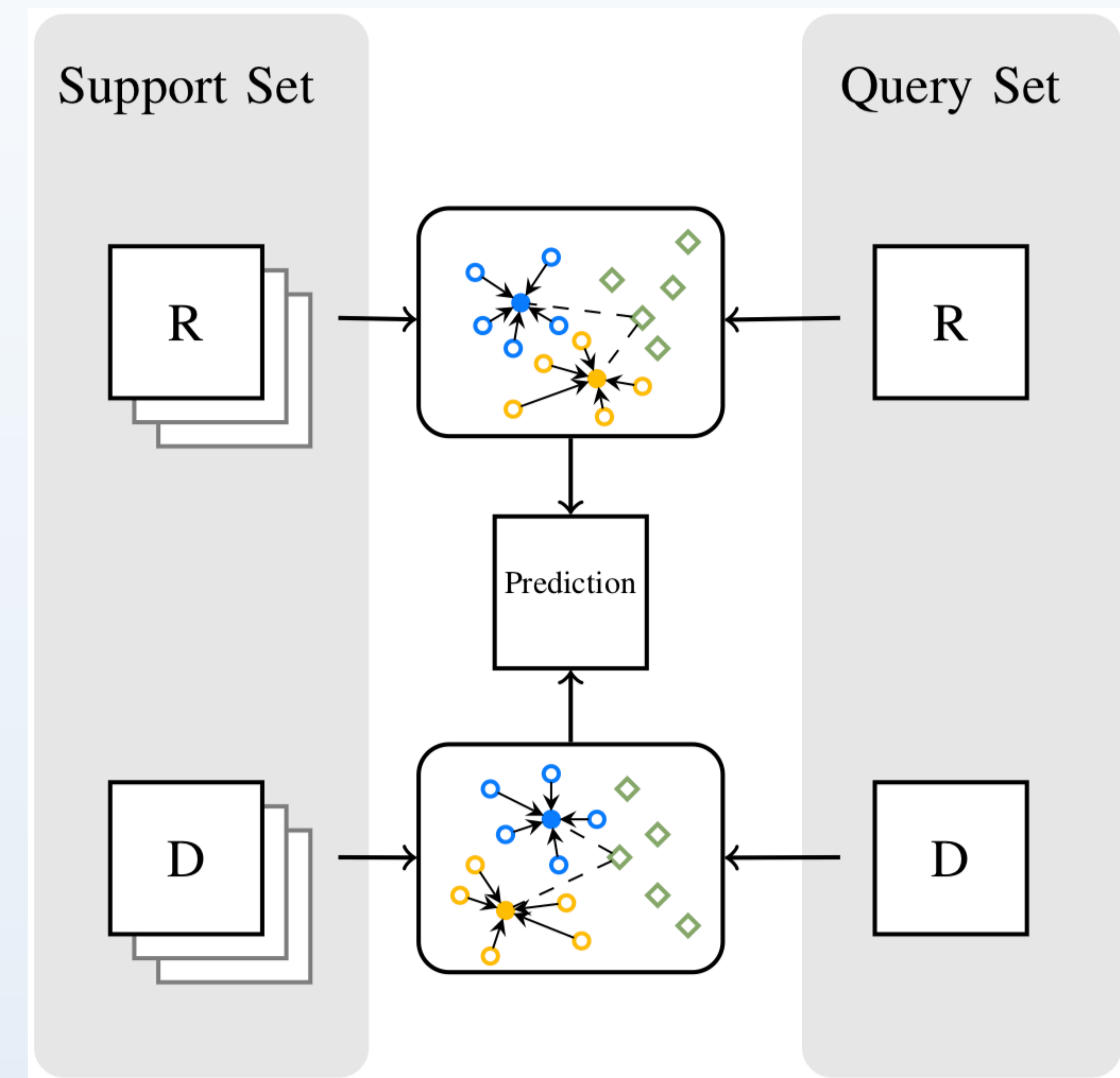


Figure 1: An overview of the proposed method (RDNet).

- Introduction on multimodal few-shot segmentation
- **Proposed model**
- Dataset
- Experimental results
- Conclusion and future work

# Proposed model

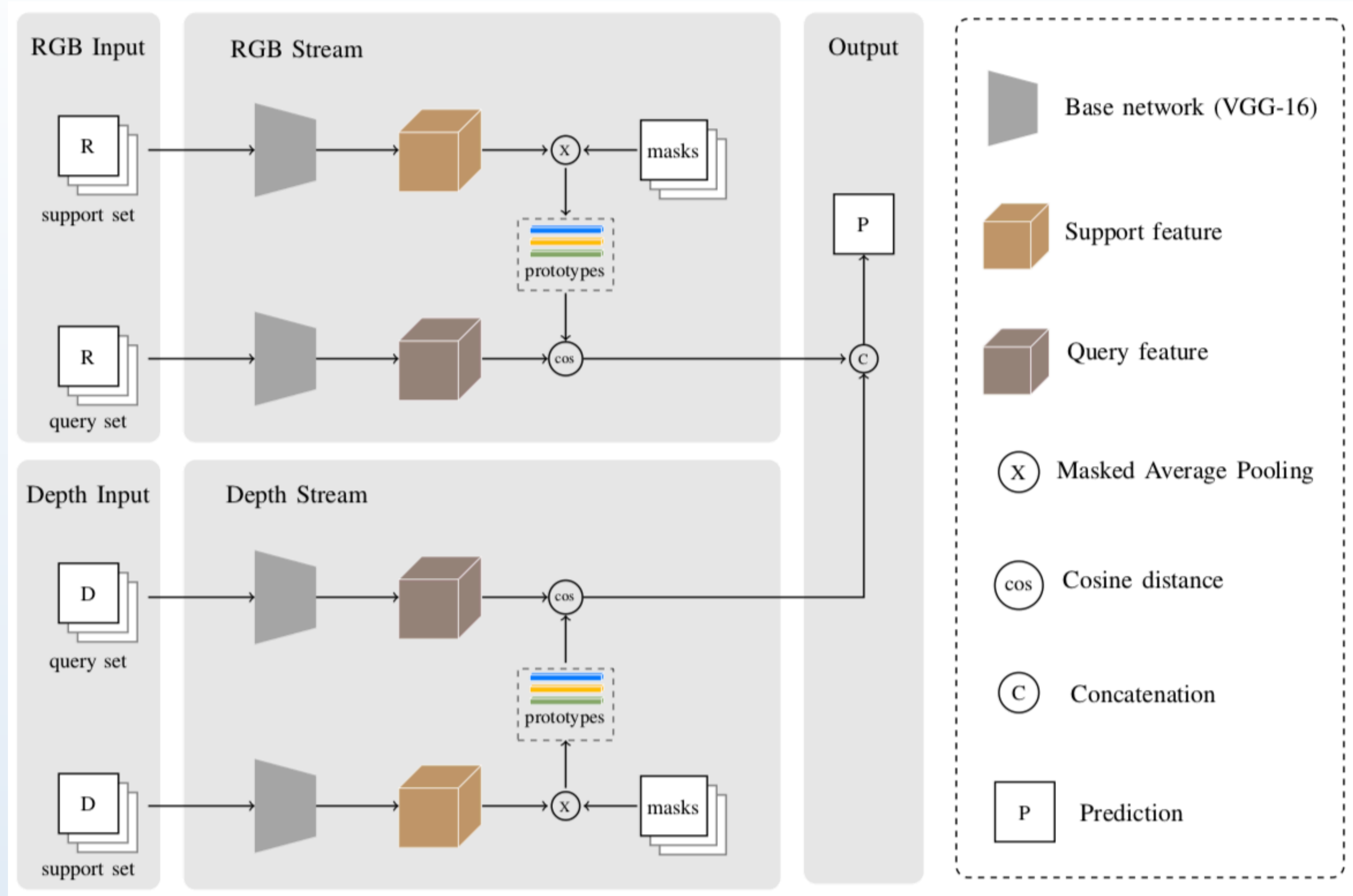


Figure 2: Illustration of the proposed method (RDNet)

- Introduction on multimodal few-shot segmentation
- Proposed model
- **Dataset**
- Experimental results
- Conclusion and future work

# Dataset

Dataset	Test classes
Cityscapes-3 <sup>0</sup>	road, sidewalk, bus
Cityscapes-3 <sup>1</sup>	vegetation, terrain, sky
Cityscapes-3 <sup>2</sup>	human, car, building

Figure 3: Training and evaluation on Cityscapes-3 dataset using 3-fold cross-validation



Figure 4: Example image from the PASCAL VOC dataset.



Figure 5: Example image from the Cityscapes dataset.

- Introduction on multimodal few-shot segmentation
- Proposed model
- Dataset
- Experimental results
- Conclusion and future work

# Experimental results

Methods	Modality	1-way 1-shot				1-way 2-shot			
		Cityscapes-3 <sup>0</sup>	Cityscapes-3 <sup>1</sup>	Cityscapes-3 <sup>2</sup>	Mean	Cityscapes-3 <sup>0</sup>	Cityscapes-3 <sup>1</sup>	Cityscapes-3 <sup>2</sup>	Mean
PANet	RGB	35.2	19.7	32.1	29.0	37.2	23.2	36.7	32.4
RDNet-R		35.7	22.3	32.6	30.2	36.7	24.1	37.5	32.8
PANet	Depth	32.6	14.5	19.3	22.1	34.2	15.8	22.5	24.2
RDNet-D		35.1	15.8	21.0	24.0	33.7	17.3	25.3	25.4
RDNet-concat	RGB-D	33.8	15.7	20.7	23.4	34.3	17.9	26.9	26.4
RDNet (ours)		<b>36.8</b>	<b>23.5</b>	<b>33.3</b>	<b>31.2</b>	<b>37.3</b>	<b>26.1</b>	<b>37.6</b>	<b>33.7</b>

Table 1: Results of 1-way 1-shot and 1-way 2-shot semantic segmentation on Cityscapes-3<sup>i</sup> using mean-IoU (%) metric.







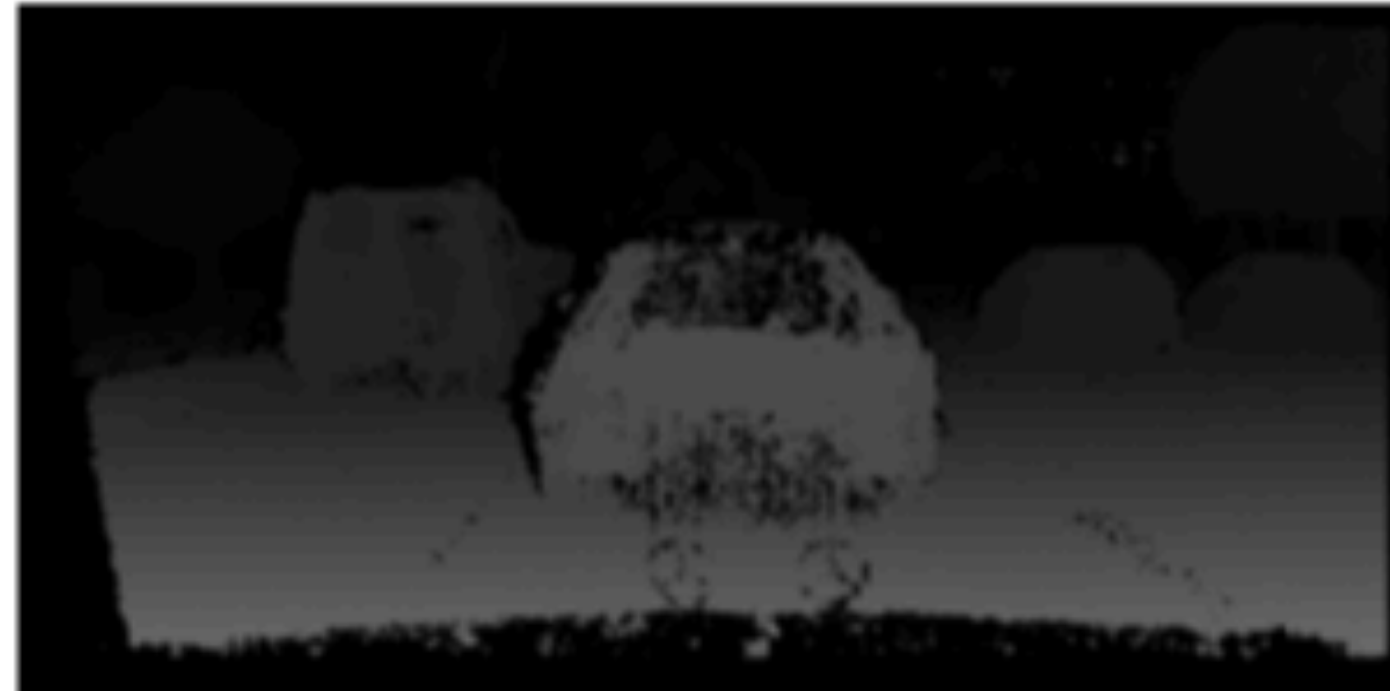

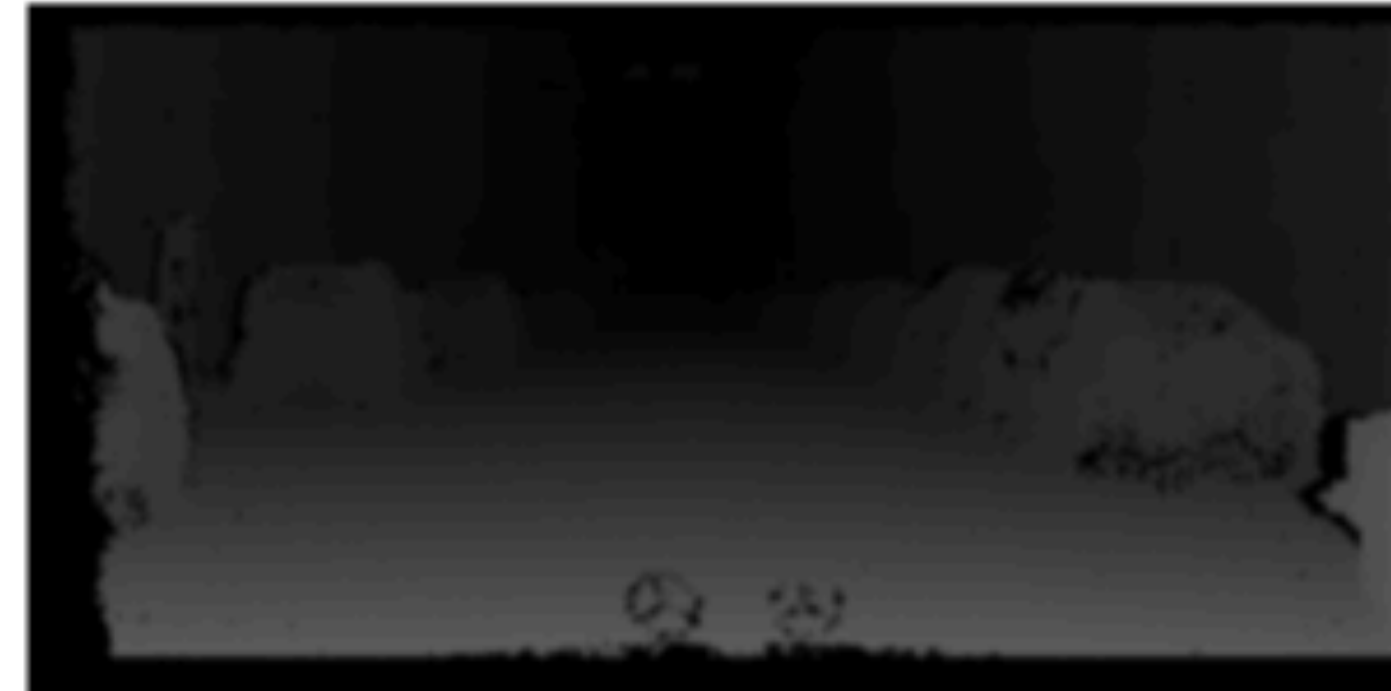

Class	Support RGB	Support depth	Query GT	Query depth	Prediction
Road					
Car					

Table 2: Qualitative results of our method for 1-way 1-shot semantic segmentation on Cityscapes-3<sup>i</sup>.

# Experimental results

Class	RDNet	RDNet-R	RDNet-D
Mean	<b>31.2</b>	30.2	24.0
Road	83.0	80.9	<b>84.4</b>
Sidewalk	<b>17.8</b>	15.7	15.7
Bus	9.5	<b>10.6</b>	5.3
Vegetation	<b>43.1</b>	40.2	26.9
Terrain	8.3	<b>10.1</b>	6.8
Sky	<b>19.1</b>	16.7	13.7
Human	<b>47.8</b>	46.6	36.9
Car	<b>12.1</b>	12.1	5.0
Building	<b>39.9</b>	39.2	21.1

Table 3: Per-class mean-IoU (%) comparison of ablation studies for 1-way 1-shot semantic segmentation.

Mehtods	Modality	binary IoU	Runtime
PANet	RGB	55.0	71ms
RDNet-R		56.5	65ms
RDNet-concat	RGB-D	51.9	67ms
RDNet (ours)		57.9	135ms

Table 4: Results of 1-way 1-shot semantic segmentation using binary-IoU and the runtime.

# Experimental results

Visualization using T-SNE:

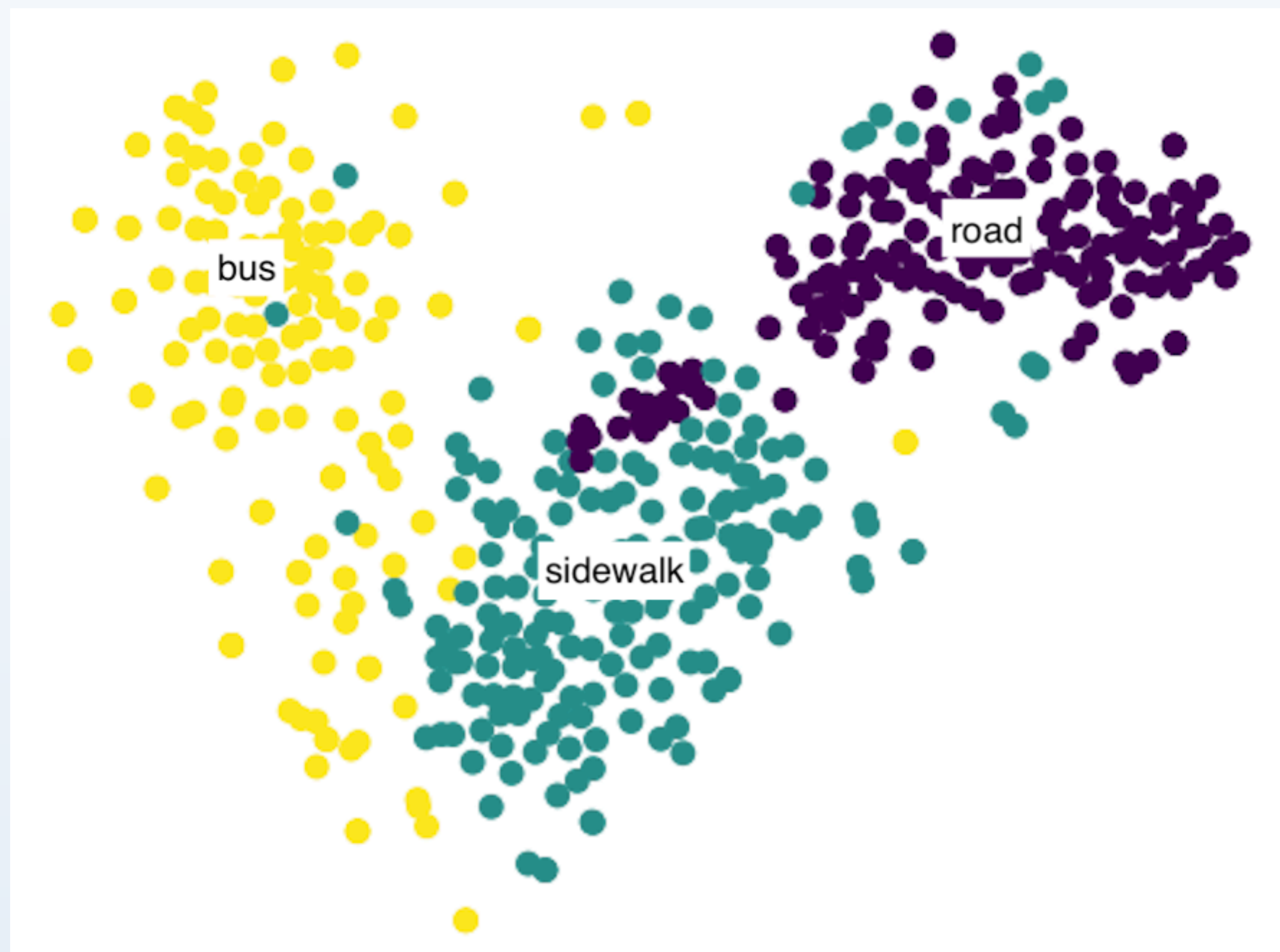


Figure 6: RGB embeddings in Cityscapes-30.



Figure 7: Depth embeddings in Cityscapes-30.

- Introduction on multimodal few-shot segmentation
- Proposed model
- Dataset
- Experimental results
- Conclusion and future work

# Conclusion and future work

## Conclusion:

- Comprehensive experiments and ablation studies on Cityscapes-3d dataset demonstrate the improved generalizability and discriminating ability of our method.
- The proposed method is simple yet effective, and explore the positive use of depth information in few-shot segmentation tasks.

## Future work:

- The integration of different multimodal information in few-shot learning tasks

**Thank you!**