



# Generative Zero-Shot Learning for Semantic Segmentation of 3D Point Clouds

Poster presentation - 3DV 2021

Bjoern Michele, Alexandre Boulch, Gilles Puy, Maxime Bucher, Renaud Marlet



# Zero-Shot learning

Training



Seen    Unseen



Test-Time

Zero-Shot learning (ZSL)



# Zero-Shot learning

Training



Seen Unseen



Test-Time

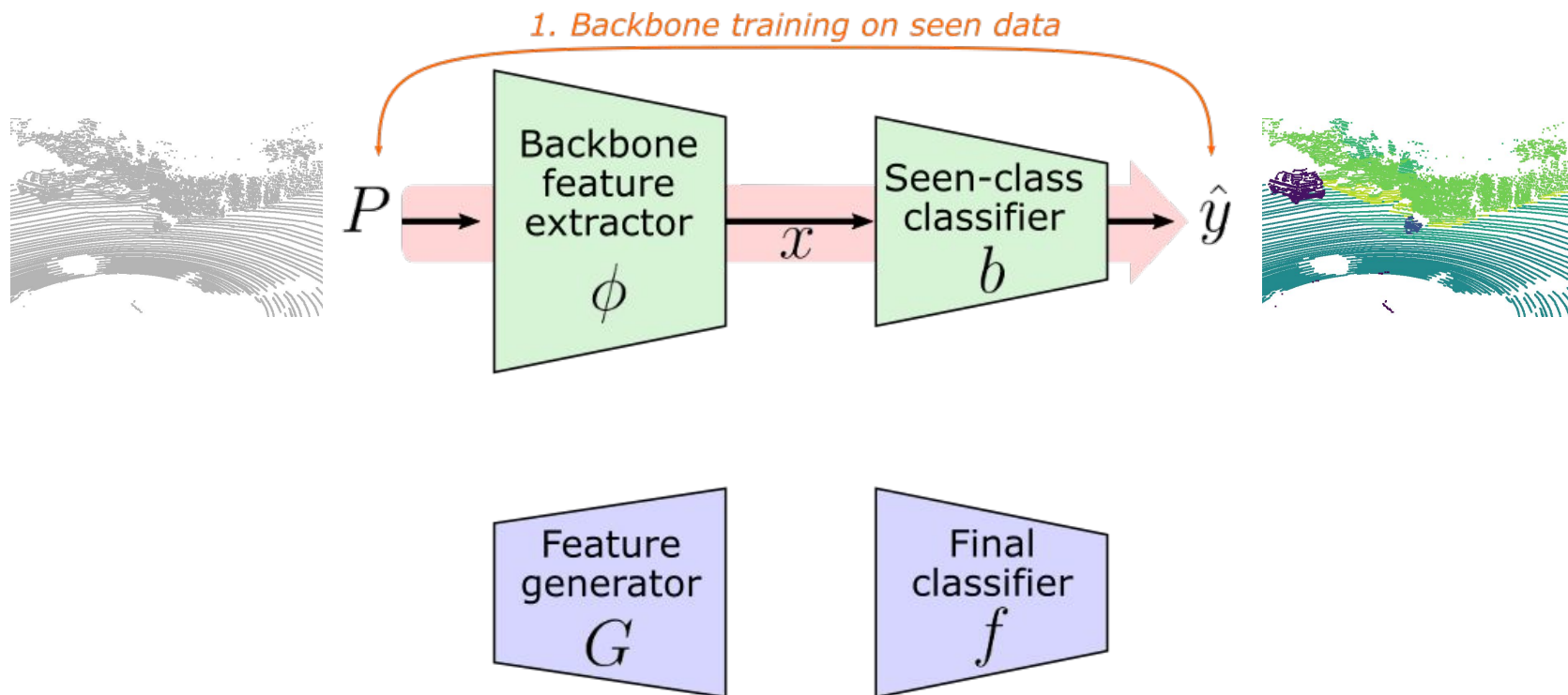
Zero-Shot learning (ZSL)



Generalized Zero-Shot learning (GZSL)

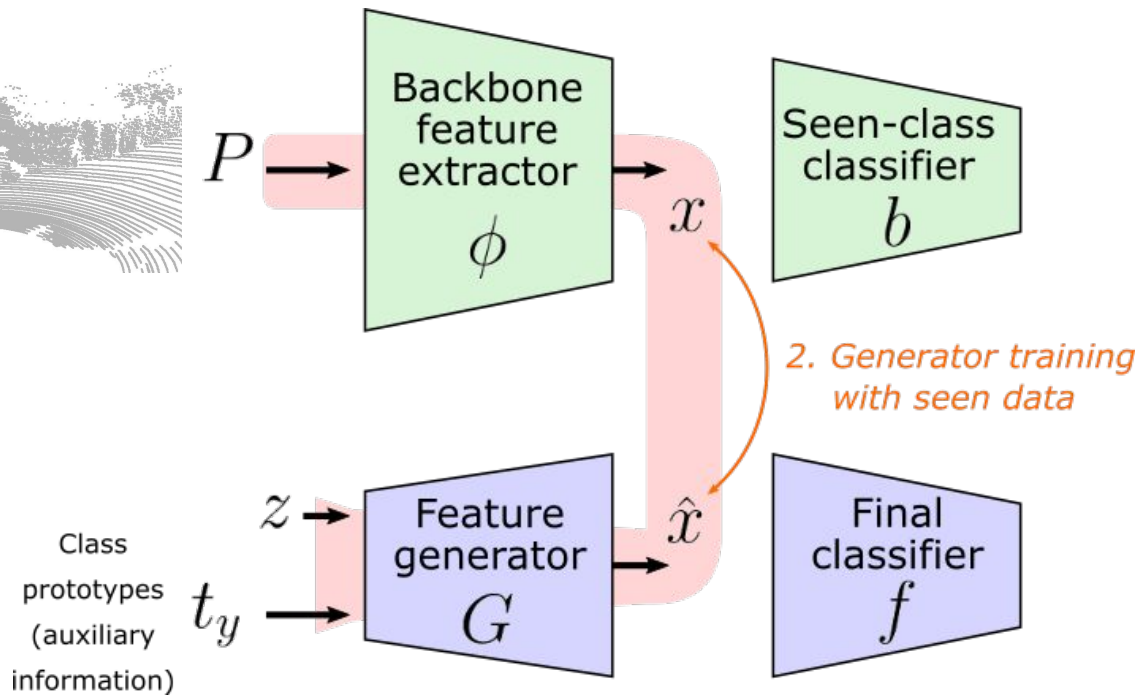
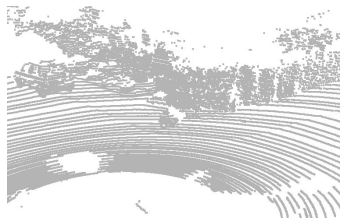


# Method<sup>1)</sup>

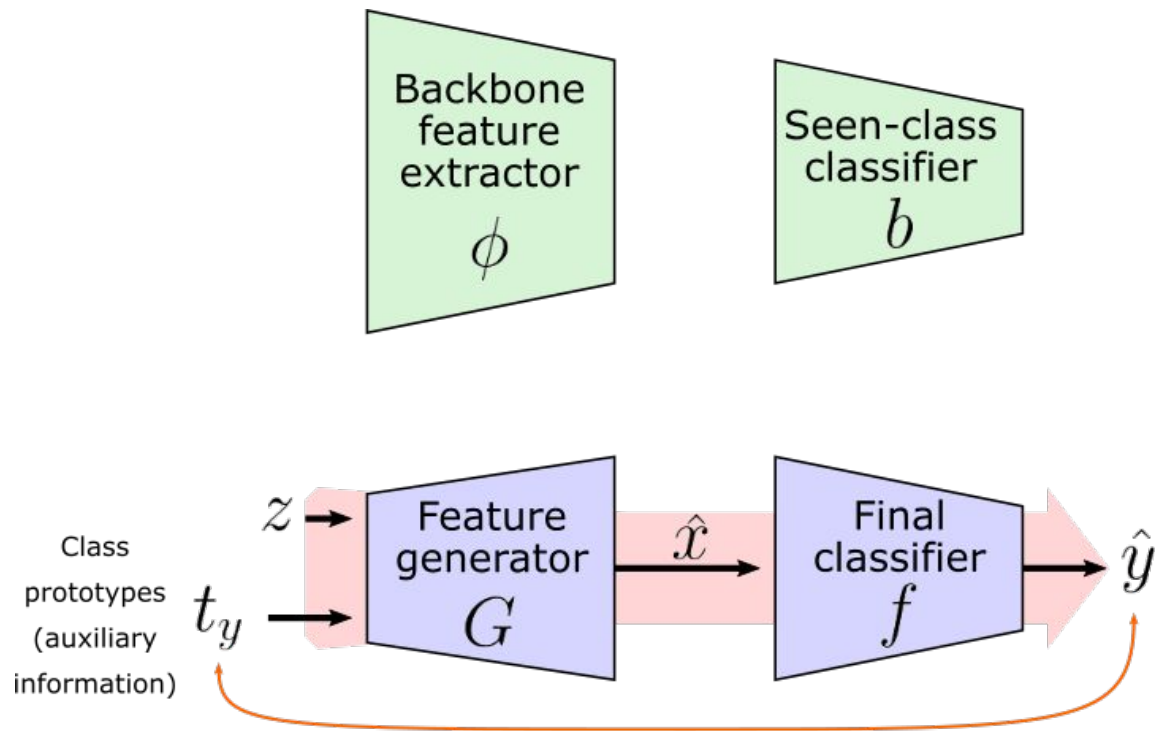


<sup>1)</sup> Adapted from:  
Bucher et al. Generating visual representations for zero-shot classification. ICCV, 2017.  
Bucher et al. Zero-shot semantic segmentation. NeurIPS, 2019.

# Method

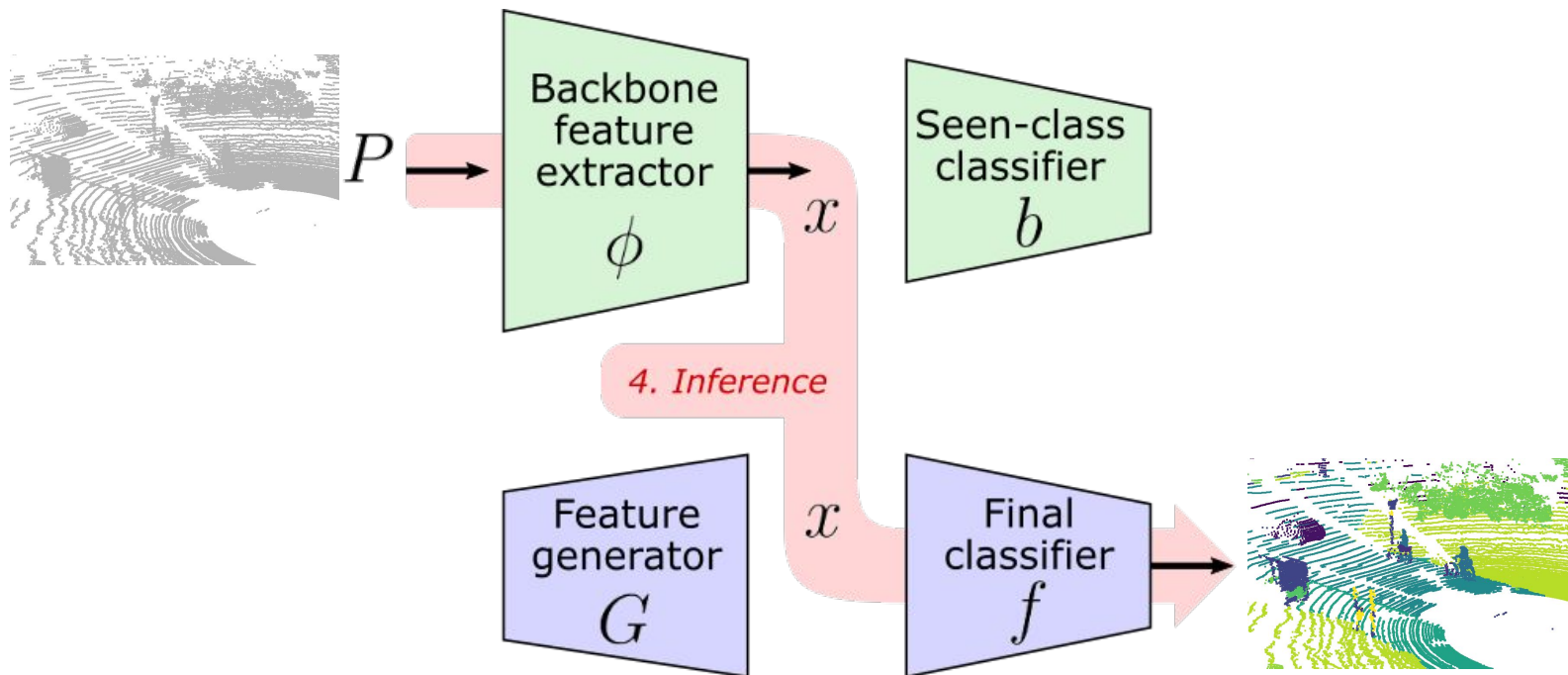


# Method



*3. Classifier training on generated unseen data  
(+ seen for GZSL)*

# Method



## Experiments - Classification ZSL

Method	Generative	Full supervision Acc.	ZSL		
			W2V Acc.	GloVe Acc.	Glove + W2V Acc.
PointNet		89.2			
f-CLSWGAN <sup>1)</sup> *	✓		20.7	-	-
CADA-VAE <sup>1)</sup> *	✓		23.0	-	-
ZSLPC <sup>2)</sup>			28.0	20.9	20.5
MHPC <sup>3)</sup>			<b>33.9</b>	28.7	-
3DGenZ (ours)	✓		28.6	<b>29.3</b>	<b>36.8</b>

\*: adaptation of 2D methods to 3D point clouds, implemented in 1).

1) Cheraghian et al. Transductive zero-shot learning for 3D point cloud classification. WACV, 2020

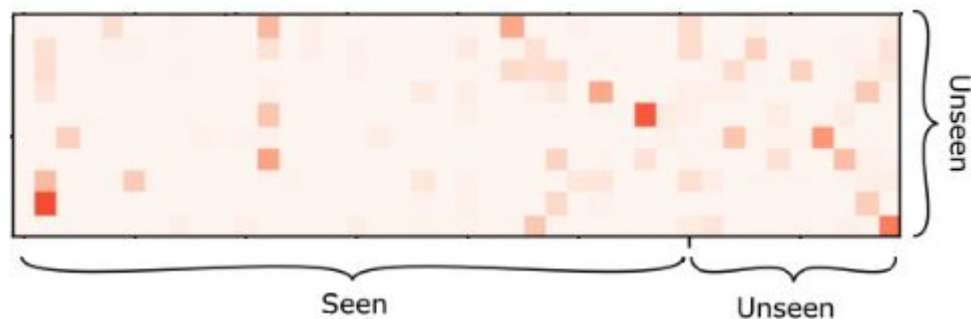
2) Cheraghian et al. Zero-shot learning of 3d point cloud objects. MVA, 2019.

3) Cheraghian et al. Mitigating the hubness problem for zero-shot learning of 3d objects. 2019.



# Bias problem in Generalized Zero-Shot learning (GZSL)

Confusion matrix for unseen classes (GZSL on Modelnet40)

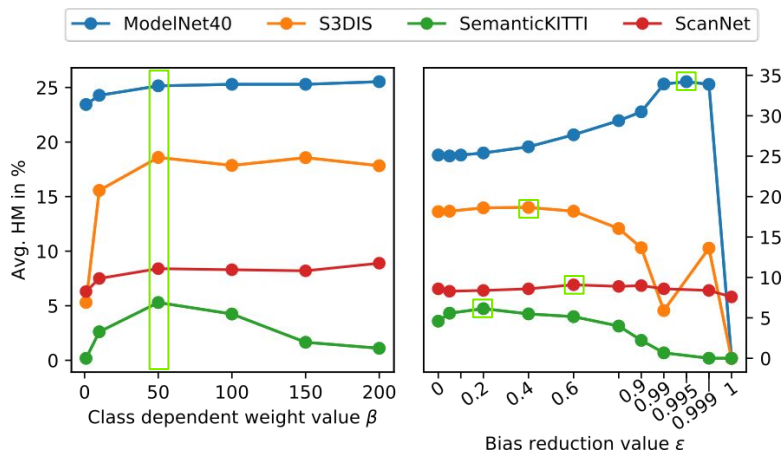


Semantic Segmentation is naturally in GZSL

# Bias problem and reduction

## Additional bias reduction techniques

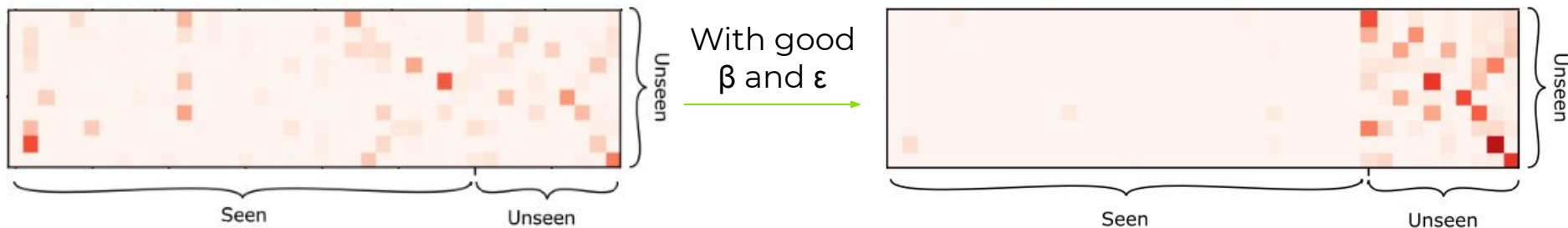
- **Class-dependant weighting:** Loss for unseen classes is weighted (factor  $\beta > 1$ ) in classifier training
- **Calibrated Stacking:** Subtracting a small value  $\epsilon$  from the seen-class score (after softmax) at test time
- $\beta$  and  $\epsilon$  are Hyperparameters



# Bias problem and reduction

## Additional bias reduction techniques

- **Class-dependant weighting:** Loss for unseen classes is weighted (factor  $\beta > 1$ ) in classifier training
- **Calibrated Stacking<sup>1)</sup>:** Subtracting a small value  $\epsilon$  from the seen-class score (after softmax) at test time
- $\beta$  and  $\epsilon$  are Hyperparameters



1) Chao et al. An empirical study and analysis of generalized zero-shot learning for object recognition in the wild. ECCV, 2016.

# Experiments - Classification GZSL

Method		Full	GZSL									
		super- vision Acc.	Bias reduct.	W2V			GloVe			GloVe + W2V		
	Gener- ative			Acc. $\mathcal{S}$	Acc. $\mathcal{U}$	HM	Acc. $\mathcal{S}$	Acc. $\mathcal{U}$	HM	Acc. $\mathcal{S}$	Acc. $\mathcal{U}$	HM
PointNet		89.2										
f-CLSWGAN <sup>1)</sup> *	✓			76.3	3.7	7.0	-	-	-	-	-	-
CADA-VAE <sup>1)</sup> *	✓			84.7	1.3	2.6	-	-	-	-	-	-
ZSLPC <sup>2)</sup>				40.1	22.5	28.8	49.2	18.2	26.6	-	-	-
MHPC <sup>3)</sup>			✓	<b>53.8</b>	26.2	35.2	<b>53.8</b>	25.7	<b>34.8</b>	-	-	-
3DGenZ (ours)	✓		✓	48.8	<b>29.3</b>	<b>36.6</b>	44.7	<b>28.4</b>	34.7	<b>47.8</b>	<b>36.5</b>	<b>41.3</b>

\*: adaptation of 2D methods to 3D point clouds, implemented in 1).

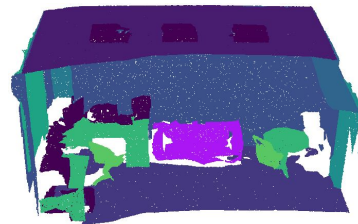
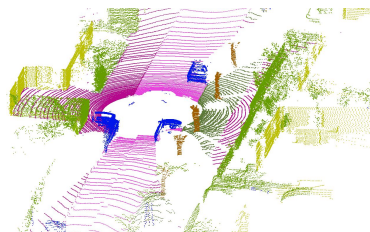
1) Cheraghian et al. Transductive zero-shot learning for 3D point cloud classification. WACV, 2020.

2) Cheraghian et al. Zero-shot learning of 3d point cloud objects. MVA, 2019.

3) Cheraghian et al. Mitigating the hubness problem for zero-shot learning of 3d objects. 2019.

# Semantic Segmentation - Datasets and baselines

- 3 Datasets
  - Outdoor
    - SemanticKITTI<sup>1)</sup> (4 Unseen, 15 Seen)
  - Indoor
    - S3DIS<sup>2)</sup> (4 Unseen, 9 Seen)
    - ScanNet<sup>3)</sup> (4 Unseen, 16 Seen)
- 2 Baselines
  - Based on DeVISE<sup>4)</sup> (2D) and ZSLPC<sup>5)</sup> (3D classification)
  - Additional bias reduction



1) Behley et al. SemanticKITTI: A Dataset for Semantic Scene Understanding of LiDAR Sequences. ICCV, 2019.

2) Armeni et al. 3D semantic parsing of large-scale indoor spaces. CVPR, 2016.

3) Dai et al. Scannet: Richly-annotated 3d reconstructions of indoor scenes. CVPR, 2017.

4) Frome et al. DeVISE: A deep visual-semantic embedding model. NIPS, 2013.

5) Cheraghian et al. Zero-shot learning of 3d point cloud objects. MVA, 2019.

# Experiments Semantic Segmentation - SemanticKITTI

15 seen, 4 unseen classes. Outdoor LiDAR dataset.

	Training set		SemanticKITTI		
	Back- bone	Classi- fier	mIoU		HmIoU
			$\mathcal{S}$	$\mathcal{U}$	
<i>Supervised methods with different levels of supervision</i>					
Full supervision	$\mathcal{S} \cup \mathcal{U}$	$\mathcal{S} \cup \mathcal{U}$	59.4	50.3	54.5
ZSL backbone	$\mathcal{S}$	$\mathcal{S} \cup \mathcal{U}$	52.9	13.2	21.2
ZSL-trivial	$\mathcal{S}$	$\mathcal{S}$	55.8	0.0	0.0
<i>Generalized zero-shot-learning methods</i>					
ZSLPC-Seg <sup>1)</sup> *	$\mathcal{S}$	$\mathcal{U}$	49.1	0.0	0.0
DeViSe-3DSeg <sup>2)</sup> *	$\mathcal{S}$	$\mathcal{U}$	49.7	0.0	0.0
ZSLPC-Seg <sup>1)</sup>	$\mathcal{S}$	$\mathcal{U}$	26.4	10.2	14.7
DeViSe-3DSeg <sup>2)</sup>	$\mathcal{S}$	$\mathcal{U}$	42.9	4.2	7.5
3DGenZ (ours)	$\mathcal{S}$	$\mathcal{S} \cup \hat{\mathcal{U}}$	<b>41.4</b>	<b>10.8</b>	<b>17.1</b>

Upper bounds

\*Direct, unrepaired (failing) adaptation.

1) Cheraghian et al. Zero-shot learning of 3d point cloud objects. MVA, 2019.

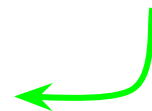
2) Frome et al. DeVISE: A deep visual-semantic embedding model. NIPS, 2013.

# Experiments Semantic Segmentation - SemanticKITTI

15 seen, 4 unseen classes. Outdoor LiDAR dataset.

	Training set		SemanticKITTI		
	Back- bone	Classi- fier	mIoU		HmIoU
			$\mathcal{S}$	$\mathcal{U}$	
<i>Supervised methods with different levels of supervision</i>					
Full supervision	$\mathcal{S} \cup \mathcal{U}$	$\mathcal{S} \cup \mathcal{U}$	59.4	50.3	54.5
ZSL backbone	$\mathcal{S}$	$\mathcal{S} \cup \mathcal{U}$	52.9	13.2	21.2
ZSL-trivial	$\mathcal{S}$	$\mathcal{S}$	55.8	0.0	0.0
<i>Generalized zero-shot-learning methods</i>					
ZSLPC-Seg <sup>1</sup> *	$\mathcal{S}$	$\mathcal{U}$	49.1	0.0	0.0
DeViSe-3DSeg <sup>2</sup> *	$\mathcal{S}$	$\mathcal{U}$	49.7	0.0	0.0
ZSLPC-Seg <sup>1</sup>	$\mathcal{S}$	$\mathcal{U}$	26.4	10.2	14.7
DeViSe-3DSeg <sup>2</sup>	$\mathcal{S}$	$\mathcal{U}$	42.9	4.2	7.5
3DGenZ (ours)	$\mathcal{S}$	$\mathcal{S} \cup \hat{\mathcal{U}}$	41.4	10.8	17.1

Baselines



\*Direct, unrepaired (failing) adaptation.

- 1) Cheraghian et al. Zero-shot learning of 3d point cloud objects. MVA, 2019.  
2) Frome et al. DeVISE: A deep visual-semantic embedding model. NIPS, 2013.

# Experiments Semantic Segmentation - SemanticKITTI

15 seen, 4 unseen classes. Outdoor LiDAR dataset.

	Training set		SemanticKITTI		
	Back- bone	Classi- fier	mIoU		HmIoU
			$\mathcal{S}$	$\mathcal{U}$	
<i>Supervised methods with different levels of supervision</i>					
Full supervision	$\mathcal{S} \cup \mathcal{U}$	$\mathcal{S} \cup \mathcal{U}$	59.4	50.3	54.5
ZSL backbone	$\mathcal{S}$	$\mathcal{S} \cup \mathcal{U}$	52.9	13.2	21.2
ZSL-trivial	$\mathcal{S}$	$\mathcal{S}$	55.8	0.0	0.0
<i>Generalized zero-shot-learning methods</i>					
ZSLPC-Seg <sup>1)</sup> *	$\mathcal{S}$	$\mathcal{U}$	49.1	0.0	0.0
DeViSe-3DSeg <sup>2)</sup> *	$\mathcal{S}$	$\mathcal{U}$	49.7	0.0	0.0
ZSLPC-Seg <sup>1)</sup>	$\mathcal{S}$	$\mathcal{U}$	26.4	10.2	14.7
DeViSe-3DSeg <sup>2)</sup>	$\mathcal{S}$	$\mathcal{U}$	42.9	4.2	7.5
3DGenZ (ours)	$\mathcal{S}$	$\mathcal{S} \cup \hat{\mathcal{U}}$	41.4	10.8	17.1

Ours

\*Direct, unrepaired (failing) adaptation.

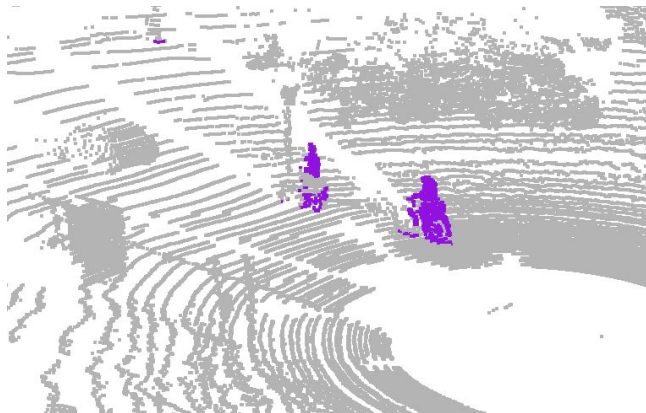
- 1) Cheraghian et al. Zero-shot learning of 3d point cloud objects. MVA, 2019.  
2) Frome et al. DeVISE: A deep visual-semantic embedding model. NIPS, 2013.



# Visualisations SemanticKITTI

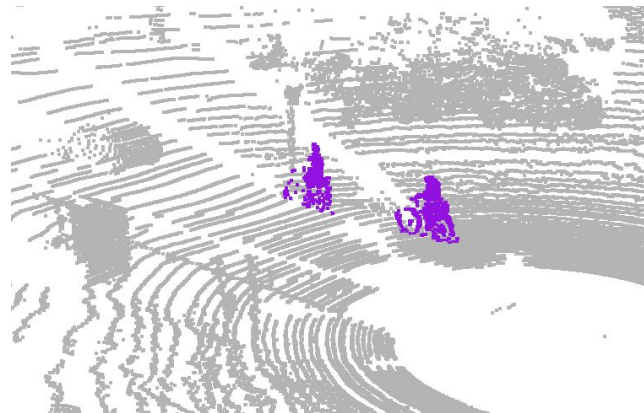
4 unseen, 15 seen classes

Predicted



Unseen Class  
Bicyclist

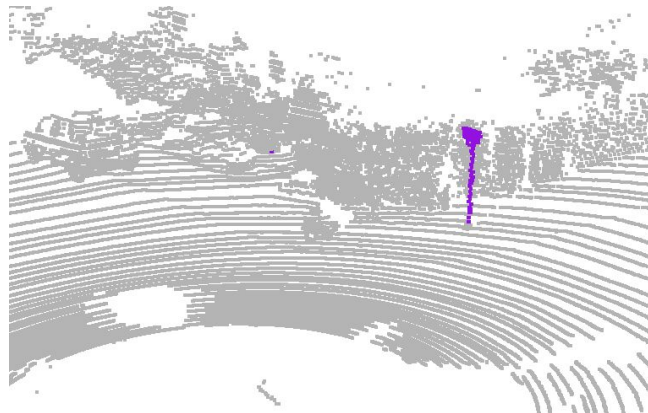
Ground-Truth



# Visualisations SemanticKITTI

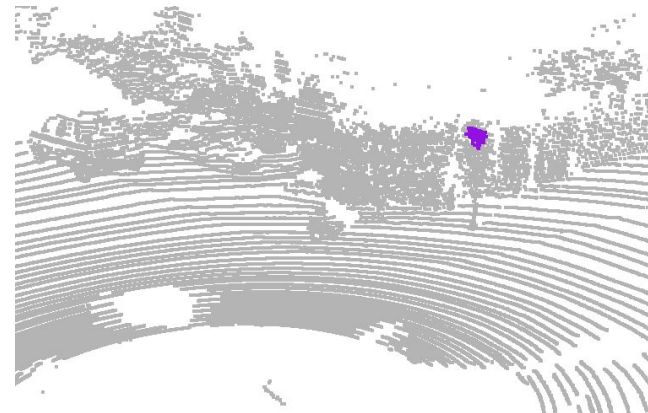
4 unseen, 15 seen classes

Predicted



Unseen Class  
Traffic-Sign

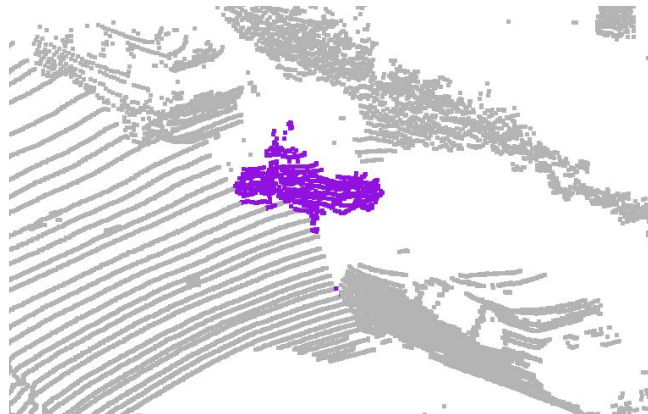
Ground-Truth



# Visualisations SemanticKITTI

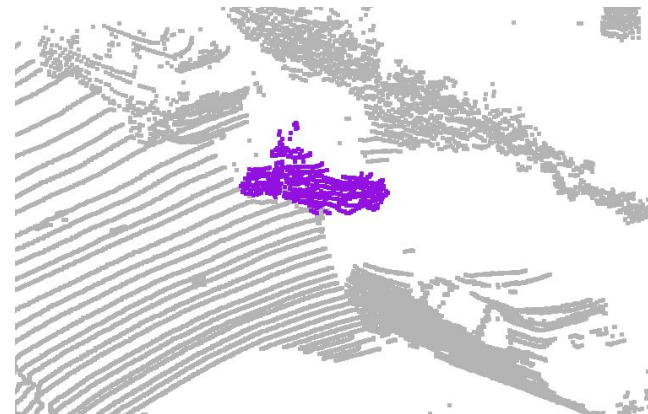
4 unseen, 15 seen classes

Predicted



Unseen Class  
Motorcycle

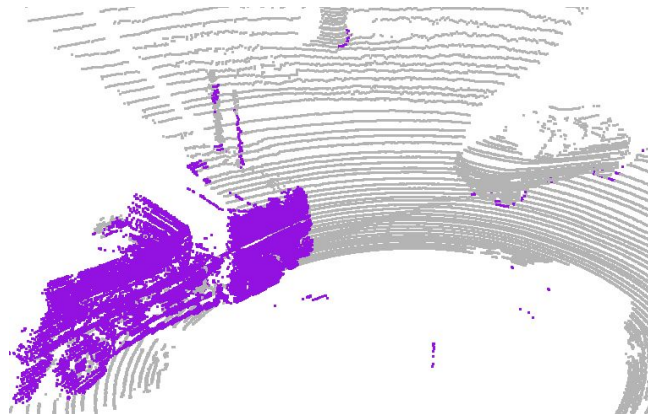
Ground-Truth



# Visualisations SemanticKITTI

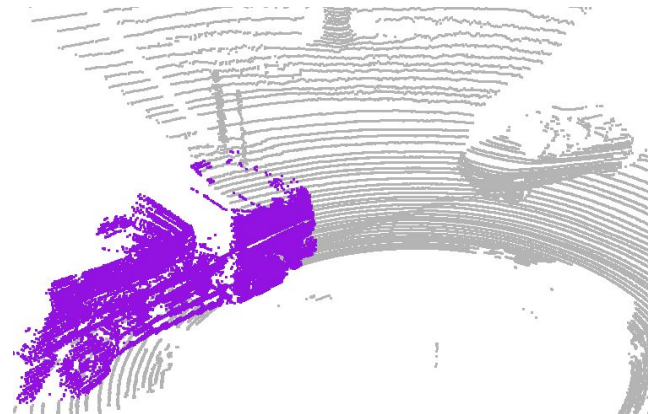
4 unseen, 15 seen classes

Predicted



Unseen Class  
Truck

Ground-Truth



# Experiments Semantic Segmentation - S3DIS, ScanNet

	Training set		S3DIS			ScanNet		
	Back- bone	Classi- fier	mIoU		HmIoU	mIoU		HmIoU
	$\mathcal{S}$	$\mathcal{U}$	$\mathcal{S}$	$\mathcal{U}$		$\mathcal{S}$	$\mathcal{U}$	
<i>Supervised methods with different levels of supervision</i>								
Full supervision	$\mathcal{S} \cup \mathcal{U}$	$\mathcal{S} \cup \mathcal{U}$	74.0	50.0	59.6	43.3	51.9	47.2
ZSL backbone	$\mathcal{S}$	$\mathcal{S} \cup \mathcal{U}$	60.9	21.5	31.8	41.5	39.2	40.3
ZSL-trivial	$\mathcal{S}$	$\mathcal{S}$	70.2	0.0	0.0	39.2	0.0	0.0
<i>Generalized zero-shot-learning methods</i>								
ZSLPC-Seg <sup>1</sup> *)	$\mathcal{S}$	$\mathcal{U}$	65.5	0.0	0.0	28.2	0.0	0.0
DeViSe-3DSeg <sup>2</sup> *)	$\mathcal{S}$	$\mathcal{U}$	70.2	0.0	0.0	20.0	0.0	0.0
ZSLPC-Seg <sup>1</sup> )	$\mathcal{S}$	$\mathcal{U}$	5.2	1.3	2.1	16.4	4.2	6.7
DeViSe-3DSeg <sup>2</sup> )	$\mathcal{S}$	$\mathcal{U}$	3.6	1.4	2.0	12.8	3.0	4.8
3DGenZ (ours)	$\mathcal{S}$	$\mathcal{S} \cup \hat{\mathcal{U}}$	<b>53.1</b>	<b>7.3</b>	<b>12.9</b>	<b>32.8</b>	<b>7.7</b>	<b>12.5</b>

Baselines



\*Direct, unrepaired (failing) adaptation.

- 1) Cheraghian et al. Zero-shot learning of 3d point cloud objects. MVA, 2019.  
 2) Frome et al. DeVISE: A deep visual-semantic embedding model. NIPS, 2013.

# Experiments Semantic Segmentation - S3DIS, ScanNet

	Training set		S3DIS			ScanNet		
	Back- bone	Classi- fier	mIoU		HmIoU	mIoU		HmIoU
	$\mathcal{S}$	$\mathcal{U}$	$\mathcal{S}$	$\mathcal{U}$		$\mathcal{S}$	$\mathcal{U}$	
<i>Supervised methods with different levels of supervision</i>								
Full supervision	$\mathcal{S} \cup \mathcal{U}$	$\mathcal{S} \cup \mathcal{U}$	74.0	50.0	59.6	43.3	51.9	47.2
ZSL backbone	$\mathcal{S}$	$\mathcal{S} \cup \mathcal{U}$	60.9	21.5	31.8	41.5	39.2	40.3
ZSL-trivial	$\mathcal{S}$	$\mathcal{S}$	70.2	0.0	0.0	39.2	0.0	0.0
<i>Generalized zero-shot-learning methods</i>								
ZSLPC-Seg <sup>1</sup> *)	$\mathcal{S}$	$\mathcal{U}$	65.5	0.0	0.0	28.2	0.0	0.0
DeViSe-3DSeg <sup>2</sup> *)	$\mathcal{S}$	$\mathcal{U}$	70.2	0.0	0.0	20.0	0.0	0.0
ZSLPC-Seg <sup>1</sup> )	$\mathcal{S}$	$\mathcal{U}$	5.2	1.3	2.1	16.4	4.2	6.7
DeViSe-3DSeg <sup>2</sup> )	$\mathcal{S}$	$\mathcal{U}$	3.6	1.4	2.0	12.8	3.0	4.8
3DGenZ (ours)	$\mathcal{S}$	$\mathcal{S} \cup \hat{\mathcal{U}}$	<b>53.1</b>	<b>7.3</b>	<b>12.9</b>	<b>32.8</b>	<b>7.7</b>	<b>12.5</b>

Ours

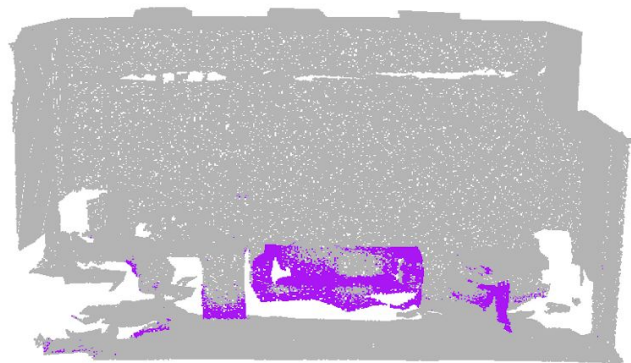
\*Direct, unrepaired (failing) adaptation.

- 1) Cheraghian et al. Zero-shot learning of 3d point cloud objects. MVA, 2019.  
 2) Frome et al. DeViSE: A deep visual-semantic embedding model. NIPS, 2013.

# Visualisations S3DIS

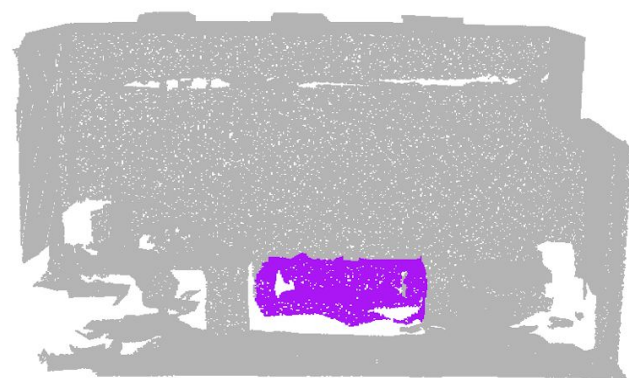
4 unseen, 9 seen classes

Predicted



Unseen Class  
Sofa

Ground-Truth

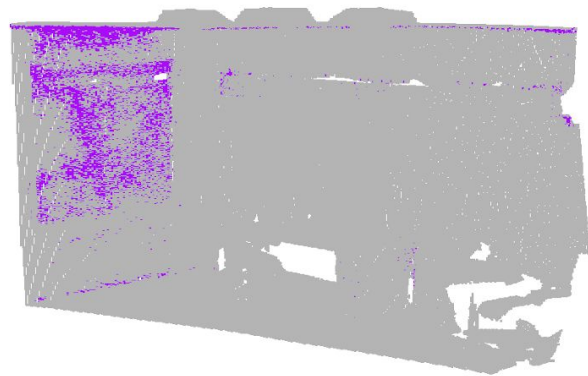




# Visualisations S3DIS

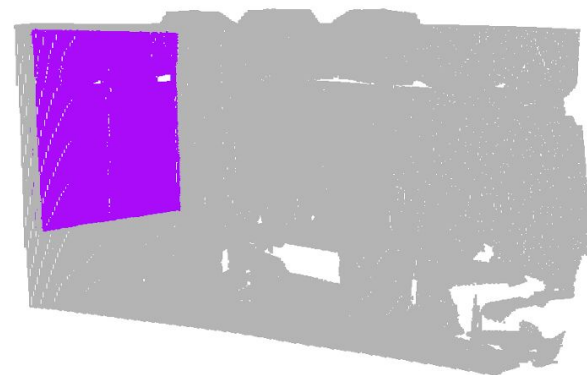
4 unseen, 9 seen classes

Predicted



Unseen Class  
Window

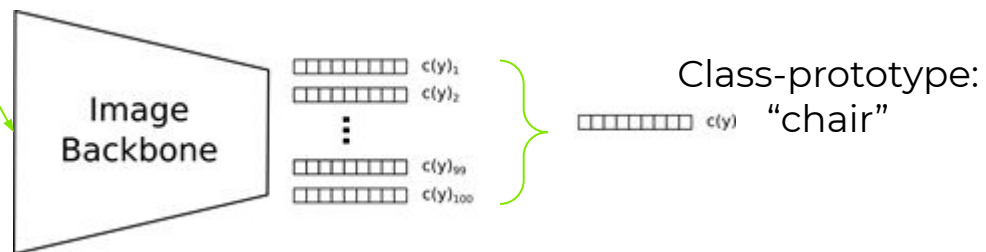
Ground-Truth





# Image based 3D Zero-shot

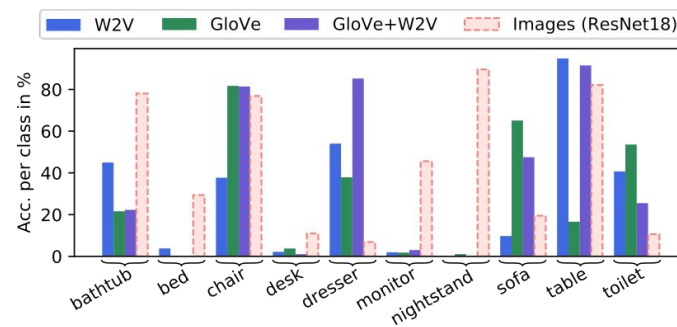
- Image-based class prototypes



# Image based 3D zero-shot - Results

Representation	Classif.		Segmentation	
	ModelNet40		ScanNet	KITTI
	ZSL	GZSL	HmIoU	
W2V+GloVe (self-sup.)	36.8	<b>41.3</b>	12.5	<b>17.1</b>
ResNet-18 (sup.)	<b>43.6</b>	40.0	13.9	3.6
ResNet-50 (self-sup.)	37.0	36.5	<b>15.5</b>	5.3

Only image  
based



# Conclusion and outview

## Conclusion

1. Reaching state of the art in **Classification** with additional bias reduction
2. GZSL for **Semantic Segmentation** on 3D PCs for 3 datasets:
  - a. Improving over natural baselines
  - b. Creation of a benchmark
3. **Image-based class prototypes** can outperform text-based ones

## Outview

- Multi-Modal prototypes
- Phrasal (multi-word) embeddings to discover complex corner cases



SMART TECHNOLOGY  
FOR SMARTER MOBILITY