



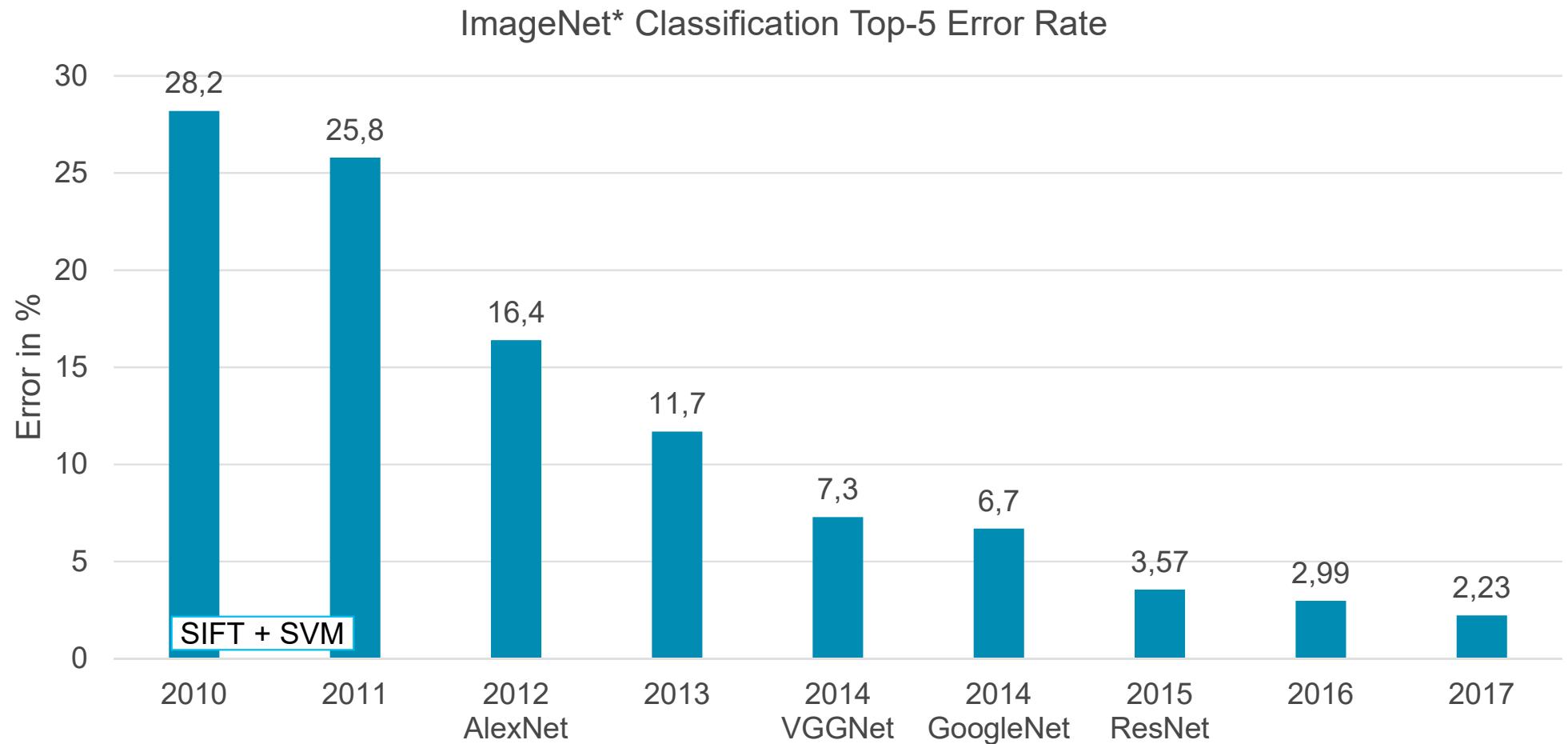
**University of Stuttgart**  
Institute for Photogrammetry



# Sparse 3D CNNs for citywide semantic mapping with ALS

Stefan Schmohl

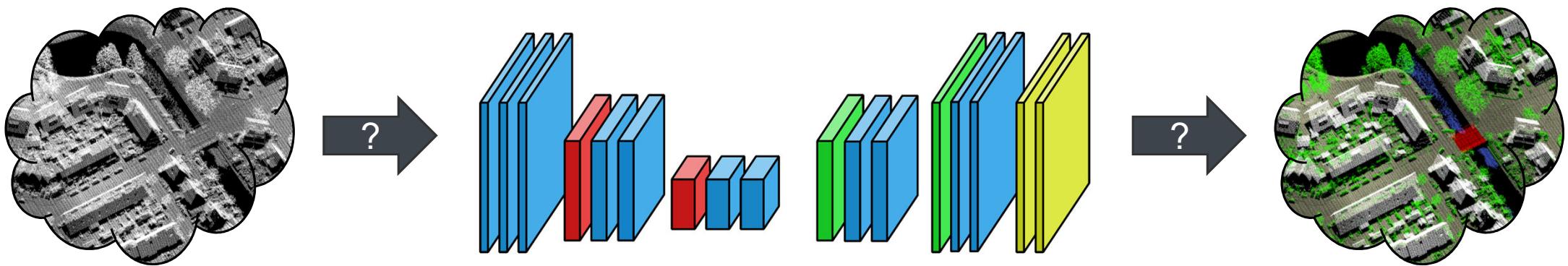
# Motivation: CNNs pushing forward the progress in vision tasks



<http://image-net.org>

Russakovsky et al., 2015: *ImageNet: Large Scale Visual Recognition Challenge*

# CNNs – how to apply them to point clouds?



## → Challenges:

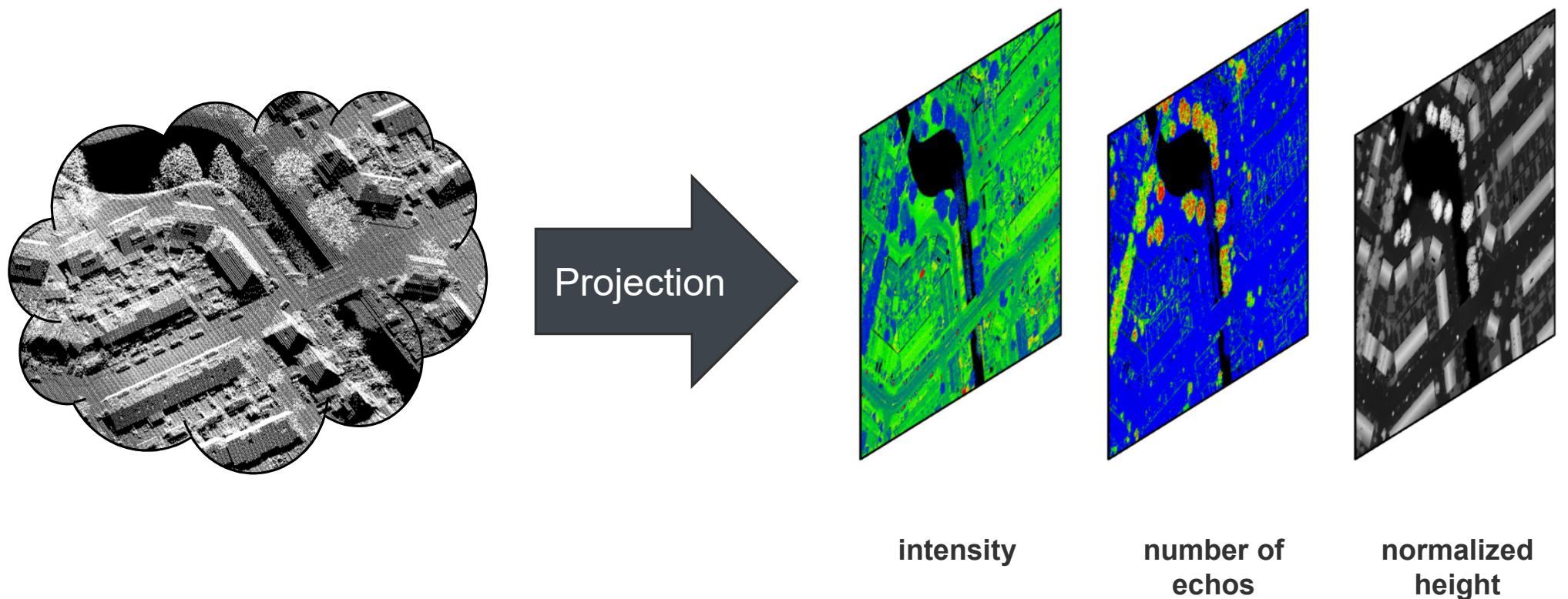
point clouds are...

- a) three dimensional
- b) irregular

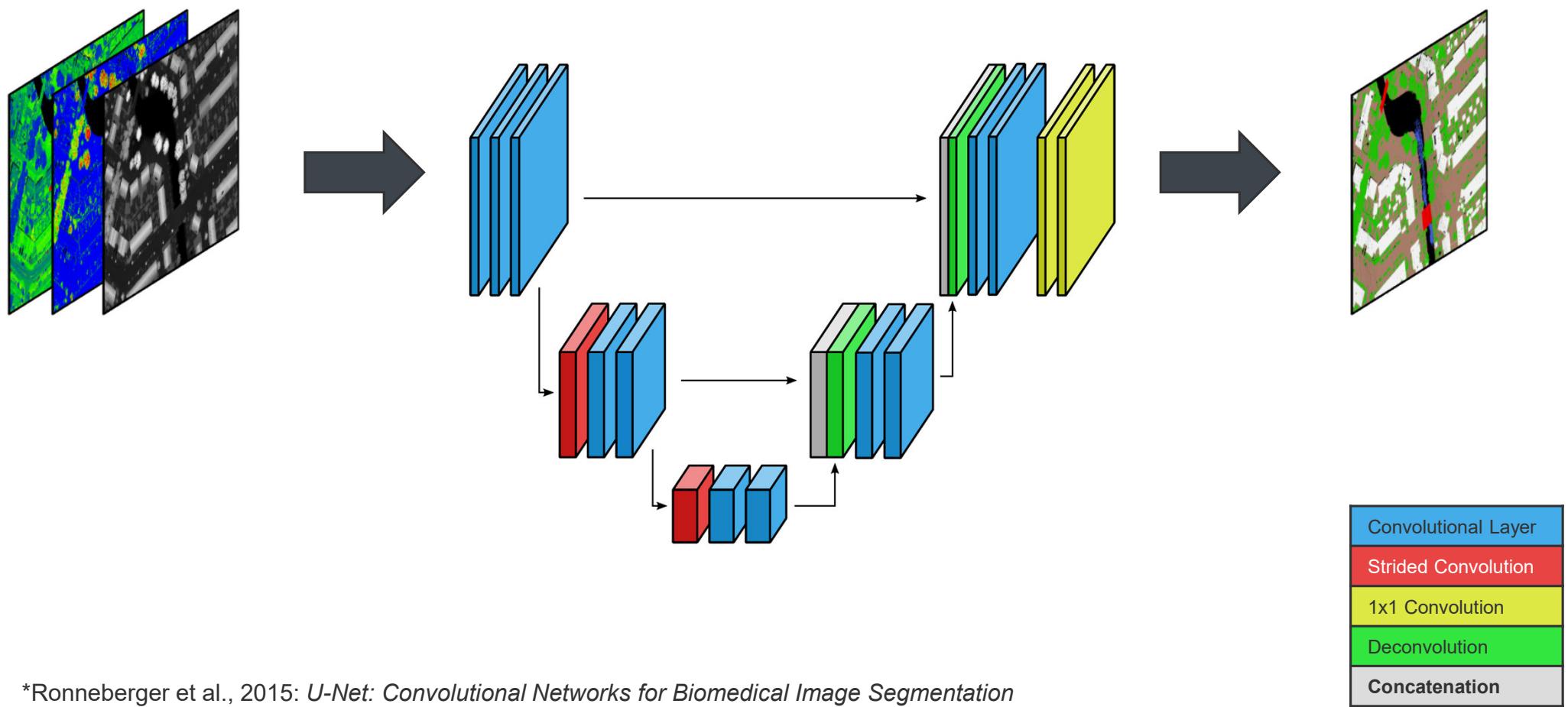
Convolutional Layer
Strided Convolution
1x1 Convolution
Deconvolution

## 2D Rasterization

---



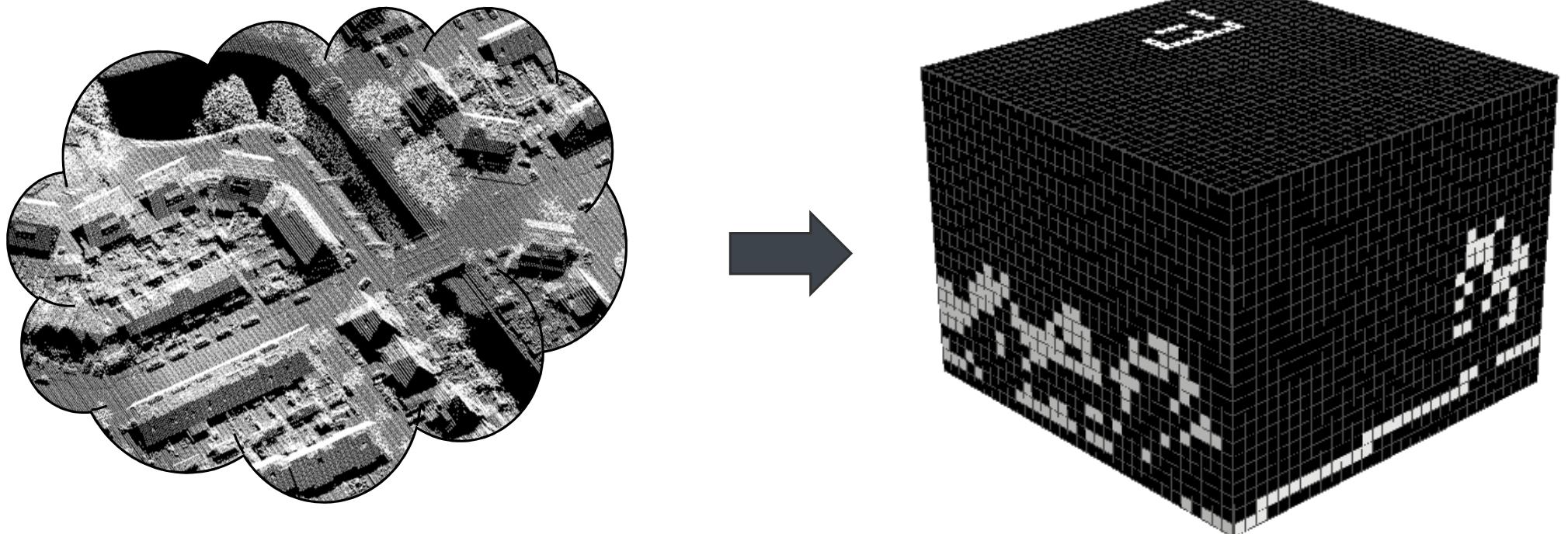
## 2D U-Net\*



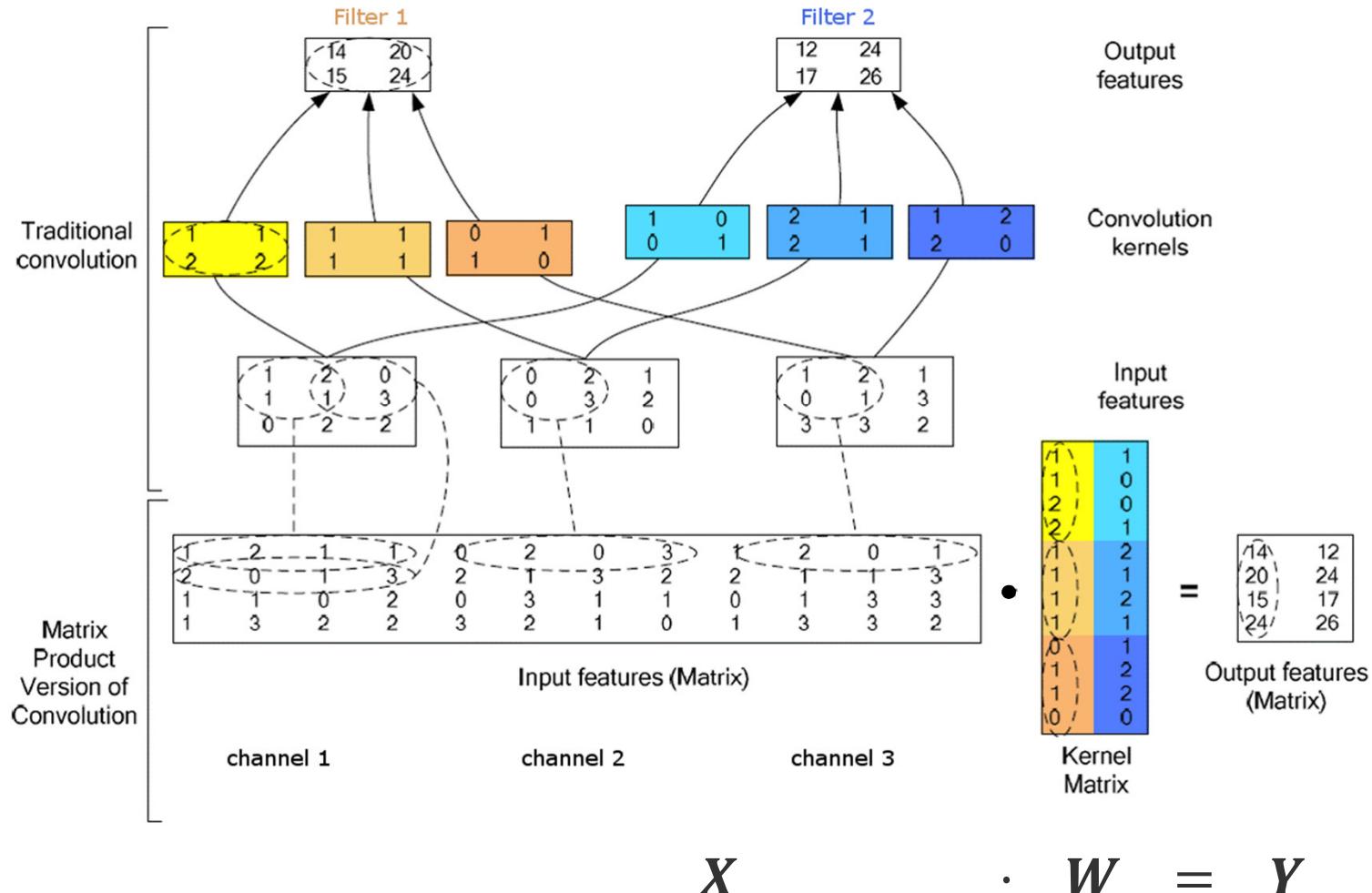
\*Ronneberger et al., 2015: *U-Net: Convolutional Networks for Biomedical Image Segmentation*

# 3D CNNs for voxelized point clouds?

---



# Convolution by GEneral Matrix Multiply (GEMM)



Chellapilla et al., 2006: High Performance Convolutional Neural Networks for Document Processing

# Convolution by GEneral Matrix Multiply (GEMM)

---

$$Y = X * W$$



One row per spatial kernel location →

$$Y = X \cdot W$$



One column per filter = output channel

---

Chellapilla et al., 2006: *High Performance Convolutional Neural Networks for Document Processing*

# Sparse 3D CNN

**Idea 1:** input data structure does not matter

=> construct  $X$  out of list of coordinates (i.e. “voxel cloud”)

**Idea 2:** sparse GEMM

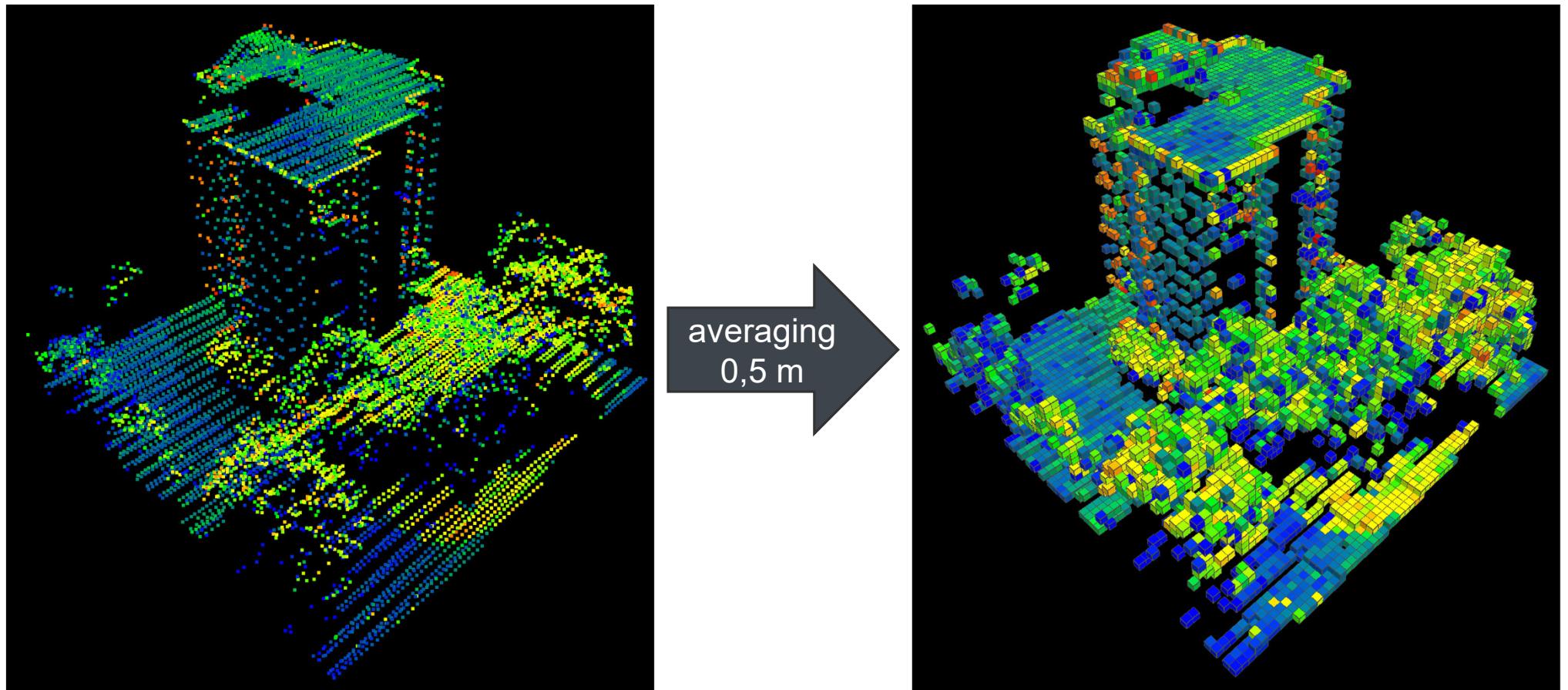
=> construct  $X$  without the “empty” locations (or rows, respectively).

=> **Submanifold Sparse Convolutional Networks\***

\*Graham et al., 2018: *Semantic Segmentation with Submanifold Sparse Convolutional Networks*

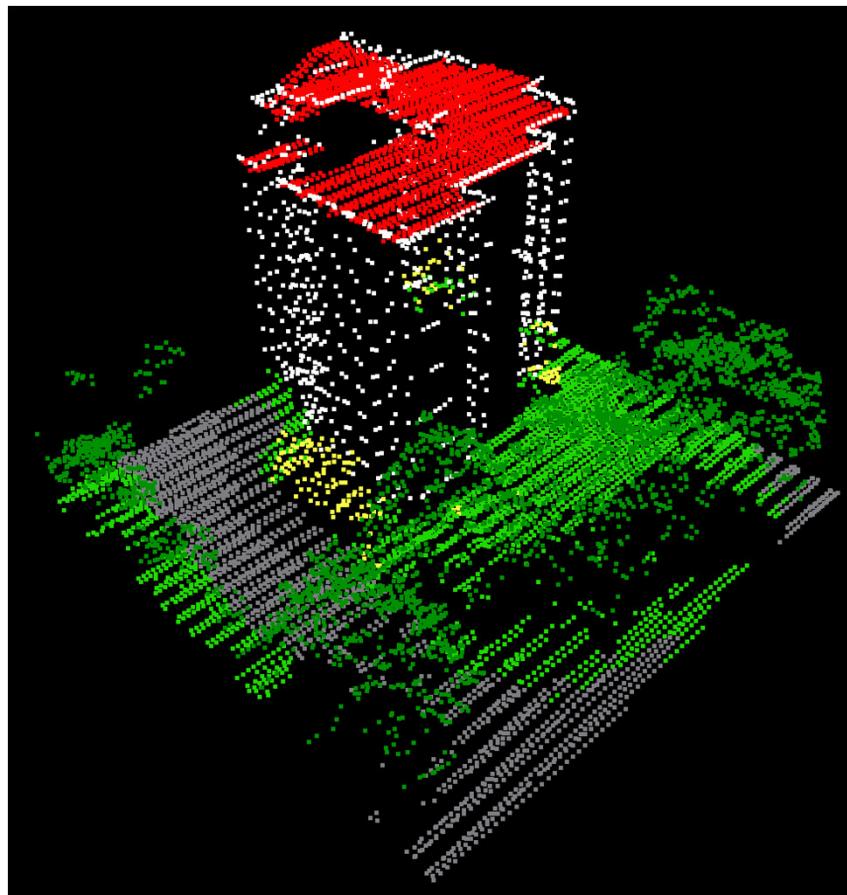
# Voxelization

---

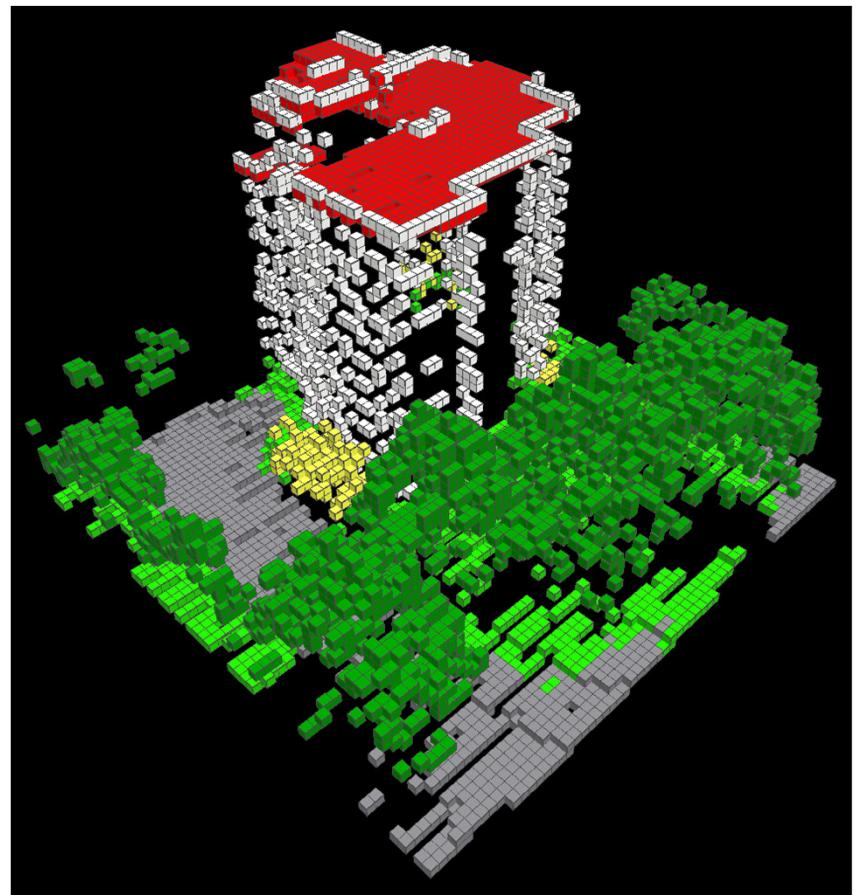


# Voxelization

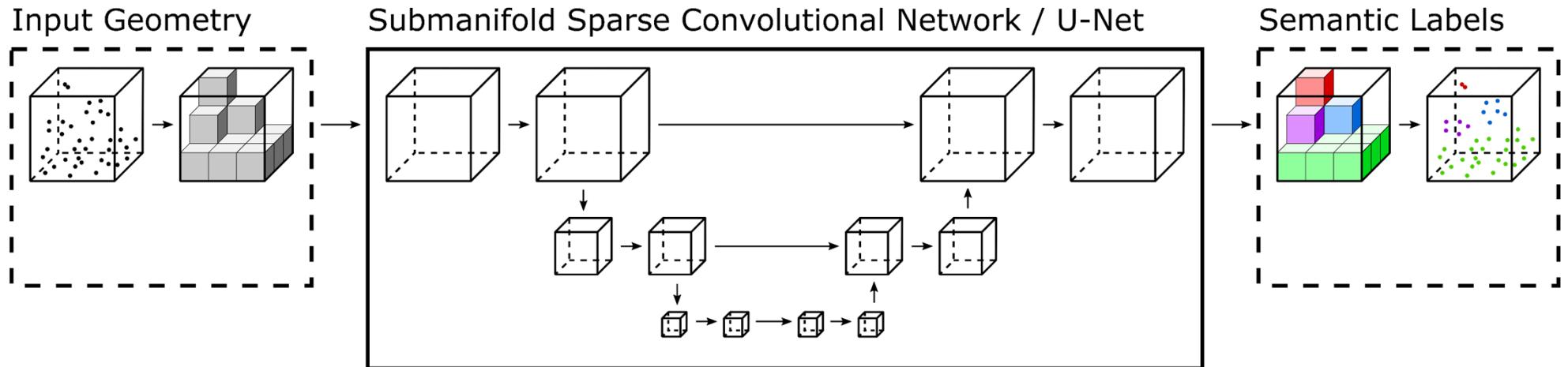
---



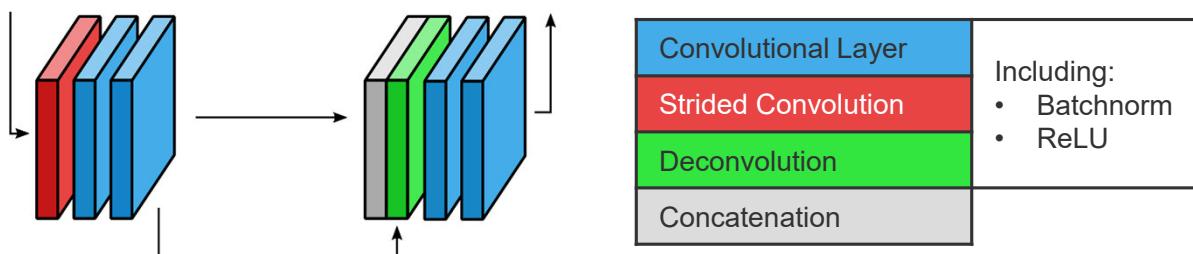
label  
transfer



# Network Architecture

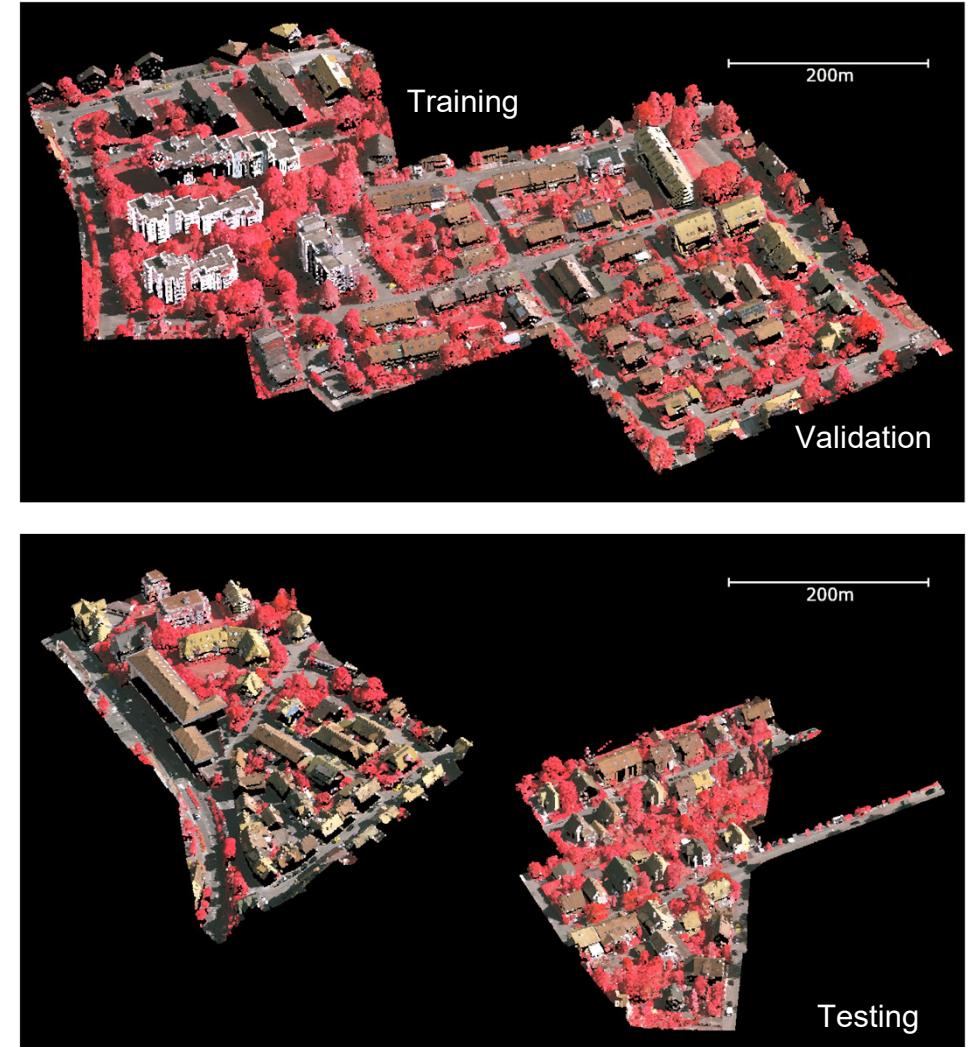


**6-7 U-Net levels:**



# Dataset I - ISPRS 3D Semantic Labeling (Vaihingen 3D)

- Training area: ~ 700.000 points
- Testing area: ~ 400.000 points
- Point density: ~ 4-8 points / m<sup>2</sup>
- Features:
  - intensity
  - echo number
  - number of echos
  - CIR
- 10 semantic classes
- Voxel size 0.5 m



# Results on V3D

## Features:

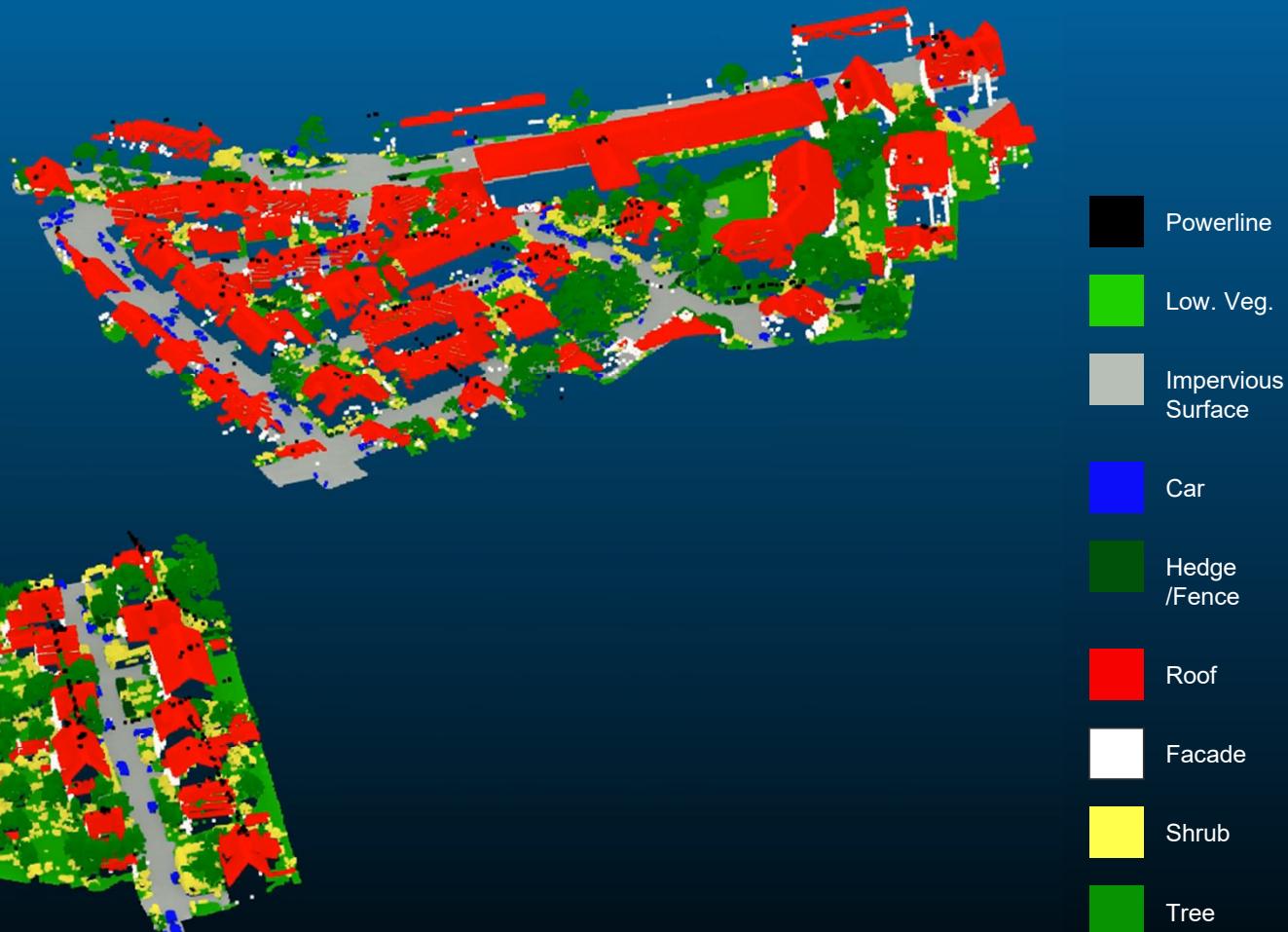
- intensity
- echo number
- # of echos
- CIR

## Voxel size:

0.5 m

⇒ OA = 85.0 %

⇒ Time: 11 s



# Memory (and Runtime) on V3D

Graphics memory	Voxel size				
	<u>2.0 m</u>	<u>1.0 m</u>	<u>0.5 m</u>	<u>0.25 m</u>	<u>0.125 m</u>
dense	1.5 GB	7.7 GB	-	-	-
sparse	0.9 GB	1.5 GB	<b>2.2 GB</b>	4.9 GB	7.9 GB

## Dataset II - Dutch National Height Model (AHN3)

- Training area: ~ 20.000.000 points
- Testing area: ~ 40.000.000 points
- Point density: ~ 9-16 points / m<sup>2</sup>
- Features:
  - intensity
  - echo number
  - number of echos
  - scan angle
- 5 semantic classes
- Voxel size 0.25 m



# Results on AHN3

## Features:

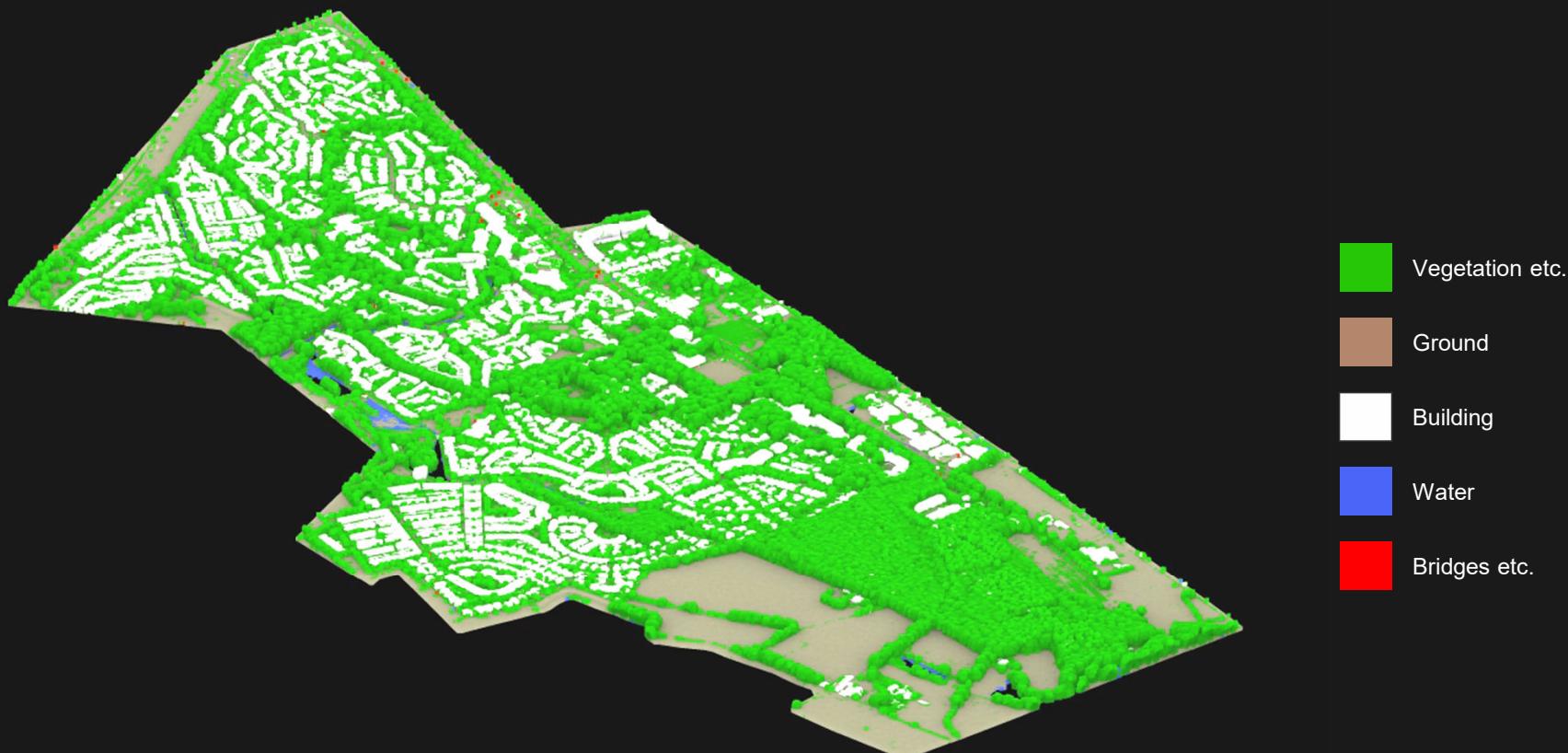
- intensity
- echo number
- # of echos
- scan angle

## Voxel size:

0.25 m

⇒ OA = **96 %**

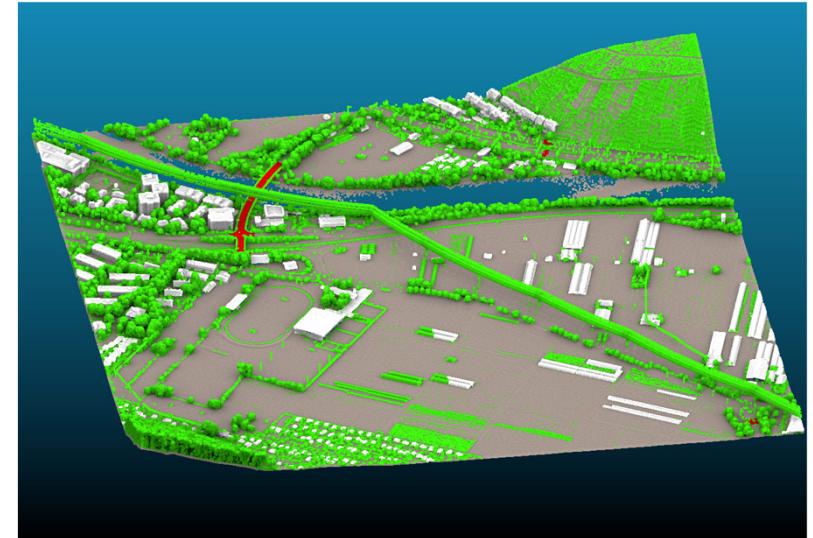
⇒ Time: **108 s**



## Dataset III - Stuttgart

Let's go even bigger!

- Point cloud properties & ground truth similar to AHN3.
  - however not the data domain, obviously.
- Coverage of Stuttgart:
  - 337 tiles, each 1 km<sup>2</sup>
- Goal:
  - Semantic classification of Stuttgart (with more than ground truth classes)
- Idea:
  - Find appropriate training data.



# V3D → Stuttgart

## Features:

- intensity
- echo number
- # of echos

## Voxel size:

0.5 m



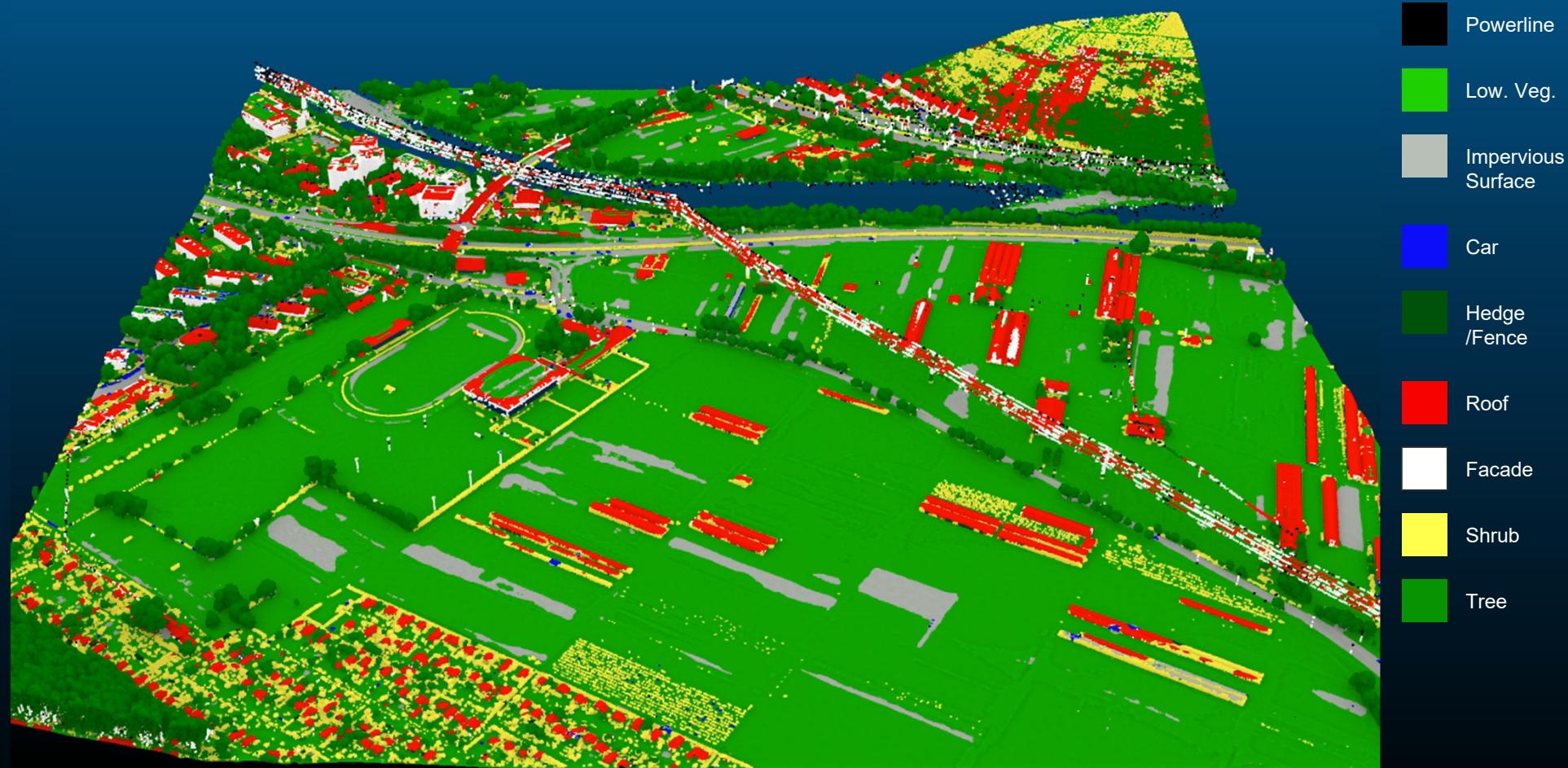
# V3D → Stuttgart

## Features:

- intensity
- echo number
- # of echos

## Voxel size:

0.5 m



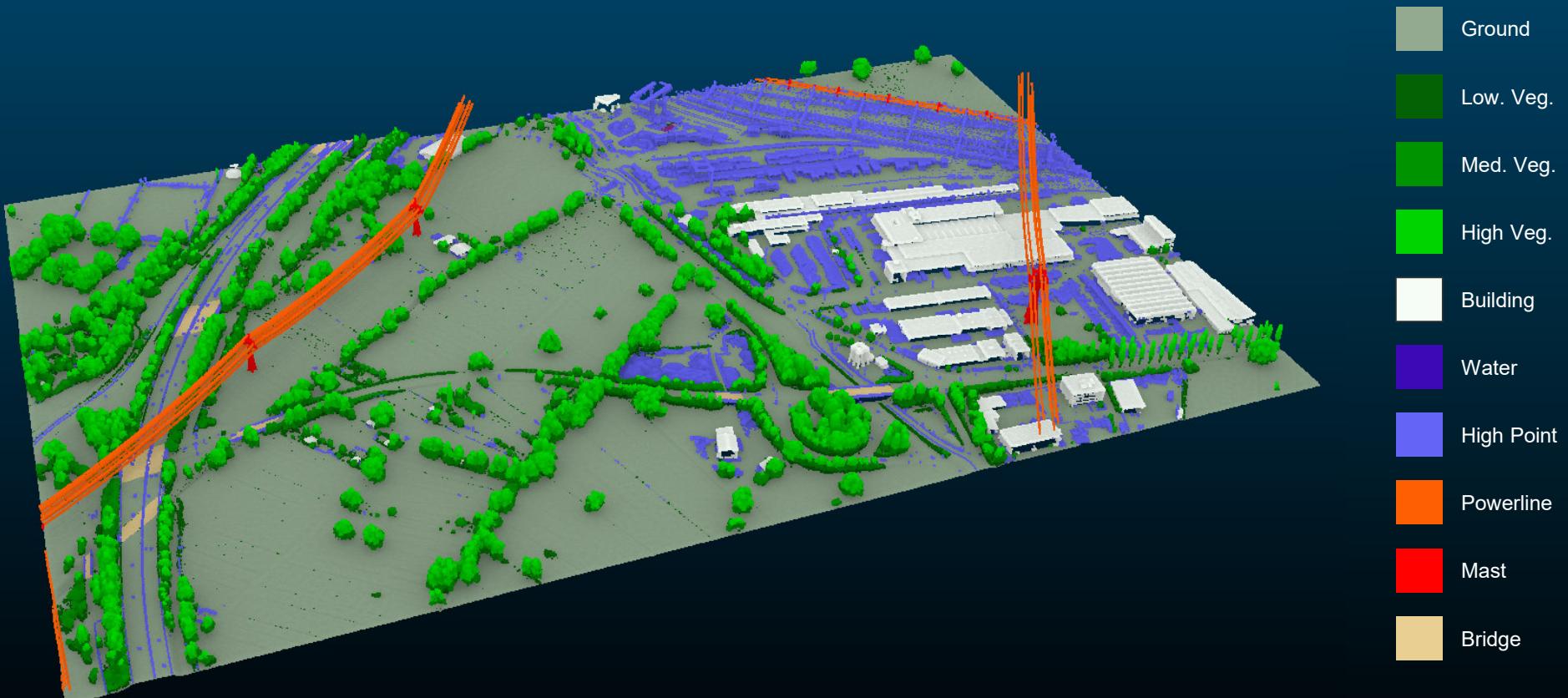
## Dataset IV – Vorarlberg, Austria

- Point cloud properties & [data domain](#)  
similar to Stuttgart.
  - But higher class diversity!
- Coverage over entire Vorarlberg:
  - 536 tiles, each 6,25 km<sup>2</sup>
- Selected training area:
  - Parts of the city of Bregenz & vicinity
  - ~ 63 M points
  - ~ 3 km<sup>2</sup>

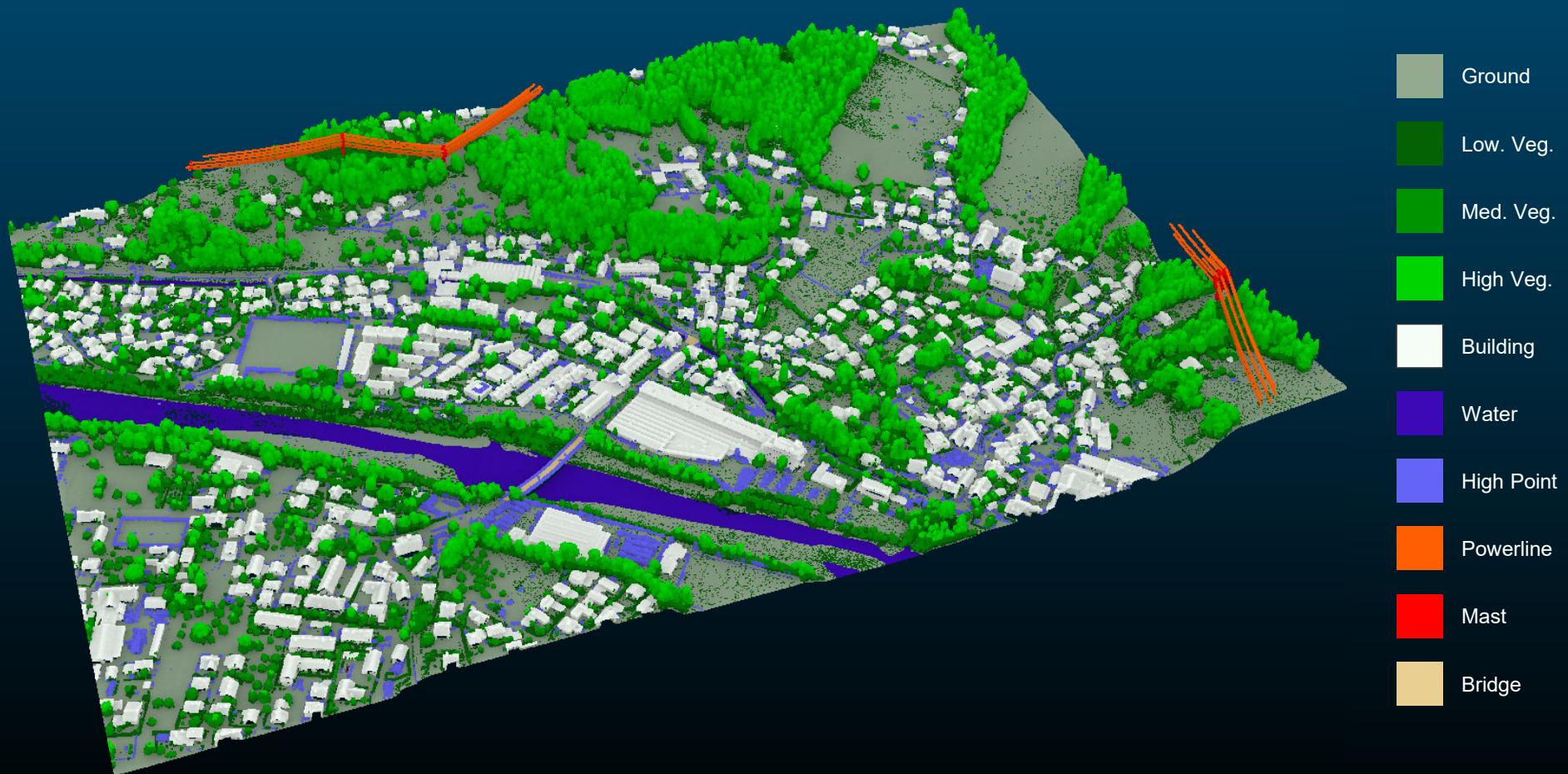


<https://de.wikipedia.org/wiki/Vorarlberg>

# Vorarlberg - Training part 1



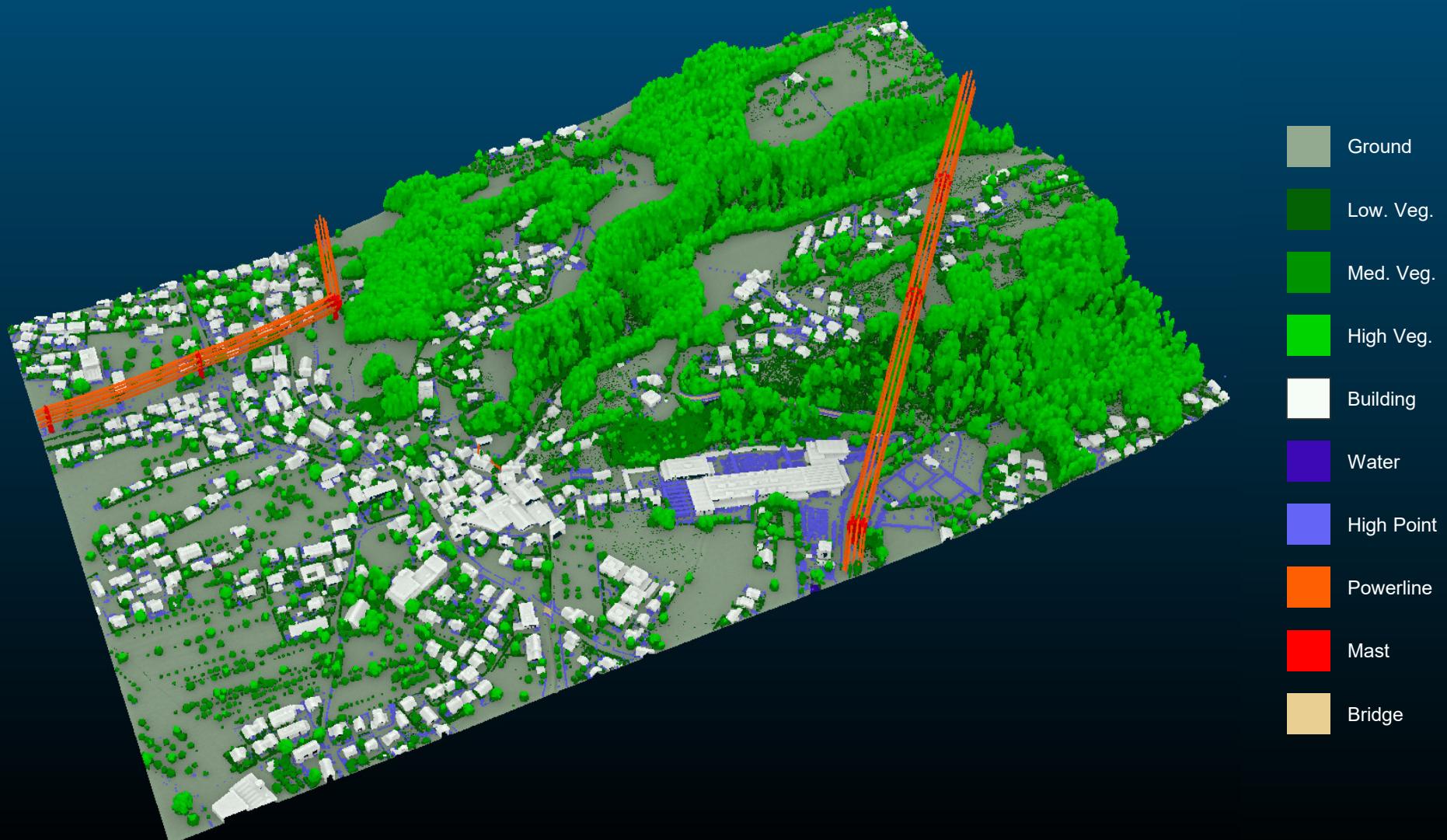
# Vorarlberg - Training part 2



# Vorarlberg - Training part 3



# Vorarlberg - Validation



# City wide mapping? Better buy a new graphics card!



**12 GB**, 3840 cores



**24 GB**, 4608 + 576 cores

=> updating drivers, cuda, PyTorch & sparse CNN framework

**=> 3x slower!!**



[WWW.PHDCOMICS.COM](http://WWW.PHDCOMICS.COM)

# Vorarlberg → Stuttgart: First results

## Features:

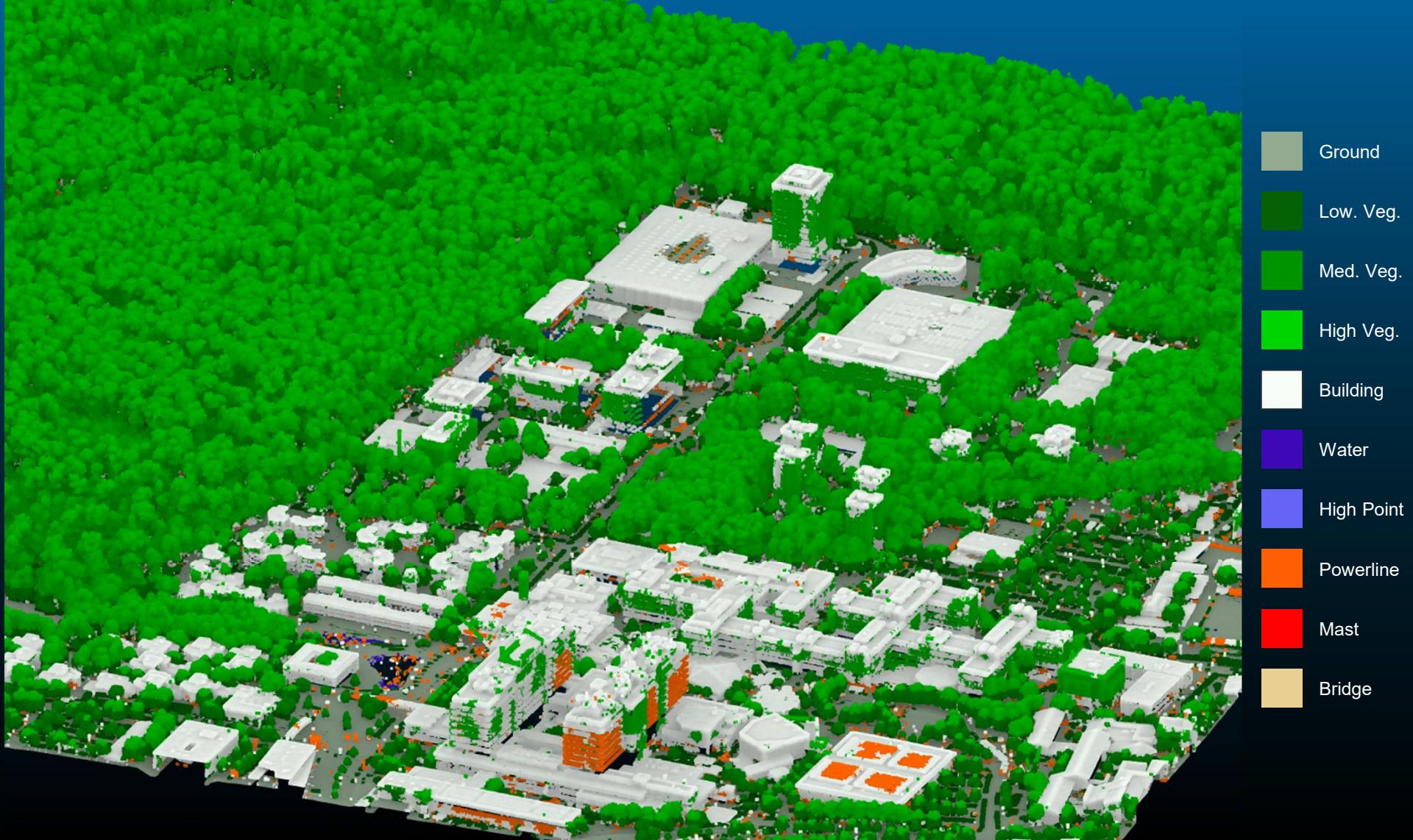
- intensity
- echo number
- # of echos
- scan angle

## Voxel size:

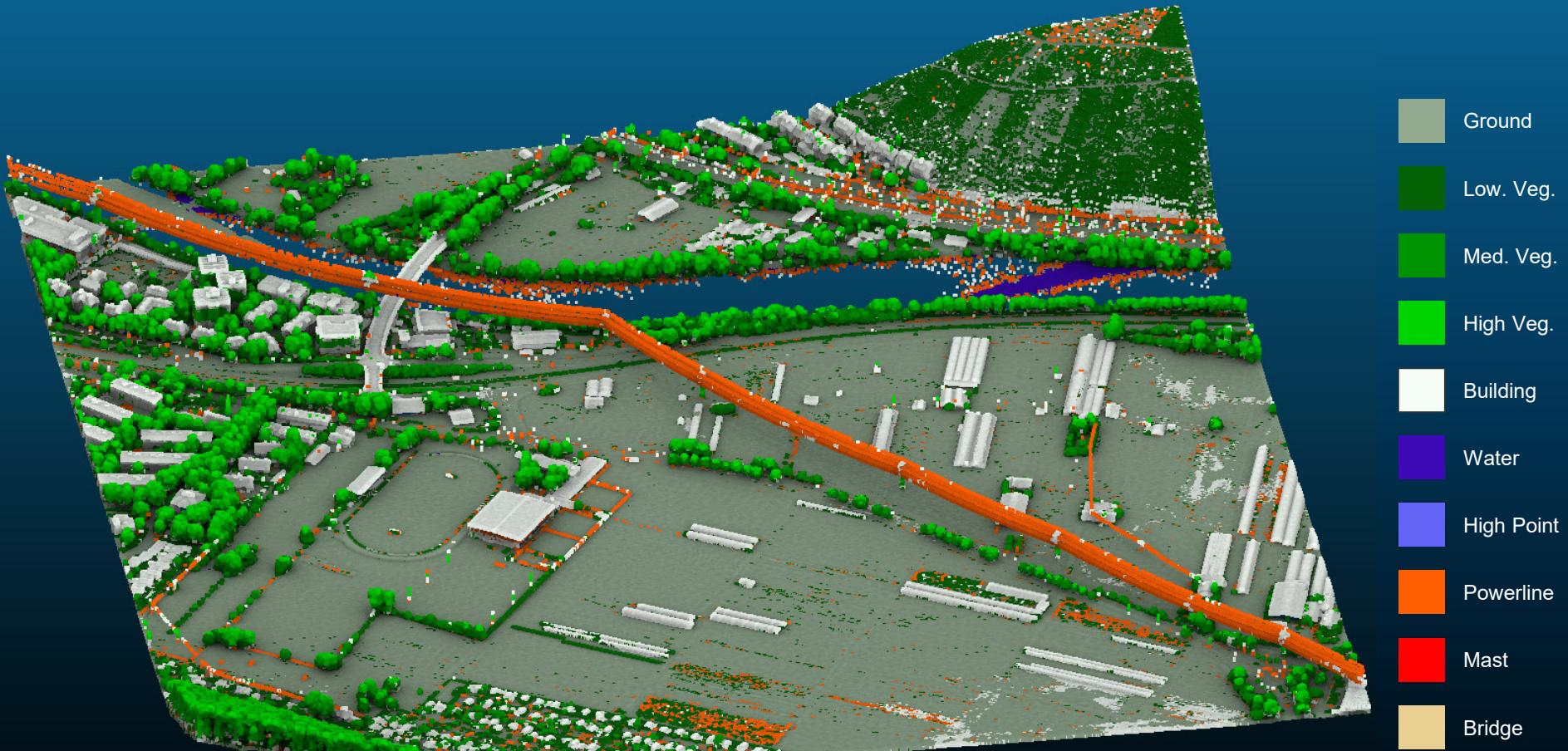
0.25 m

27 M points

⇒ Time: 173 s



# Vorarlberg → Stuttgart: First results



# Vorarlberg → Stuttgart: First results

## Features:

- intensity
- echo number
- # of echos
- scan angle

## Voxel size:

0.25 m



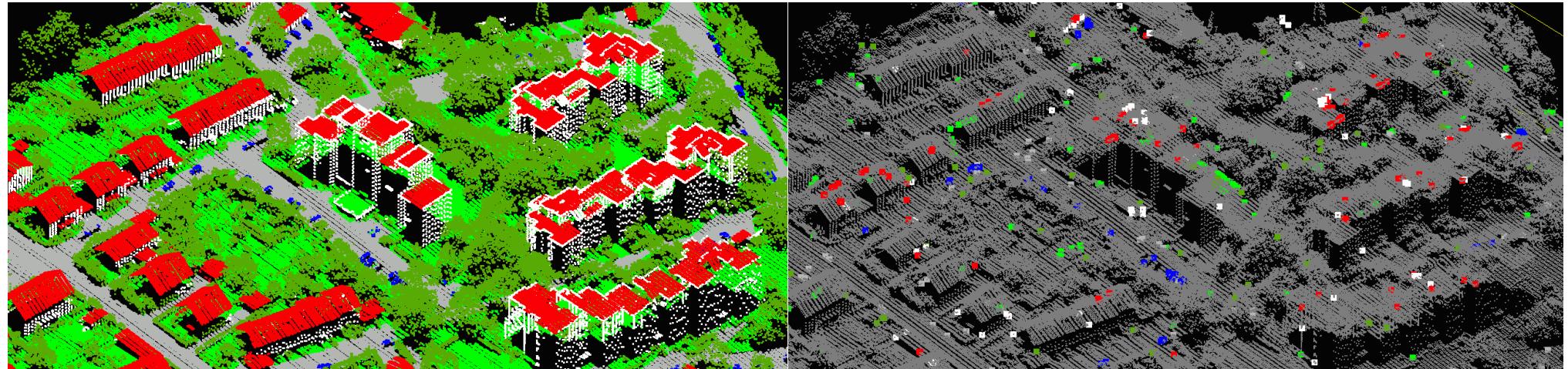
# Dense vs. sparse training dataset

Now with a sparse 3D CNN!

dense training set from expert

VS.

sparse training set from crowd



- 100 % of points labeled
- OA = **88.39 %**
- Costs = ???
- Expert(s) actively involved in labeling
- 0.4 % of points labeled
- OA = **85.43 %**
- Costs < 130 \$
- Only providing data and wages

# Conclusion

---

- No need for extra feature extraction or point-to-image conversion.
  - Implicit geometry is the most important feature.
  - State-of-the-art accuracy on V3D.
  - Very fast inference.
  - Good results also on sparsely labeled training data.
- 
- Not quite ready for production
  - “smart” training strategy needed for huge training data
  - Some “problematic classes”

# Thank you for listening!

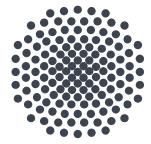
---

Many Thanks to:

- International Society for Photogrammetry and Remote Sensing (ISPRS)
- Publieke Dienstverlening Op de Kaart (PDOK)
- Landesamt für Geoinformation und Landentwicklung Baden-Württemberg (LGL)
- Landesamt für Vermessung und Geoinformation Vorarlberg

for proving the data that made this work possible!





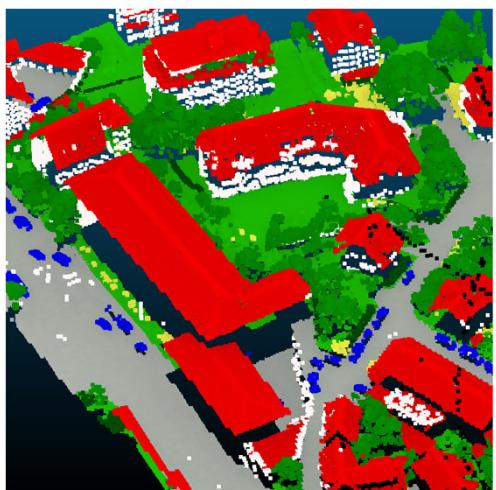
Universität Stuttgart

# Supplementary

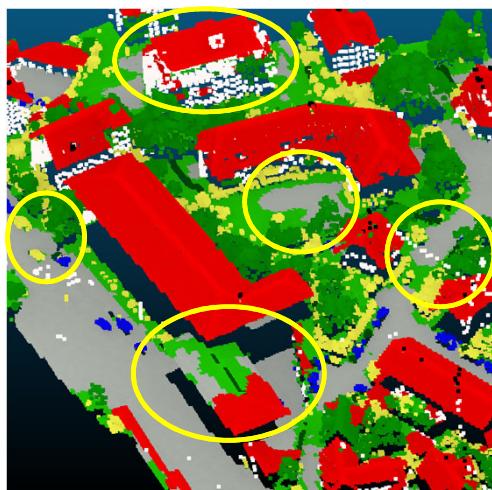


# Results on V3D

ground truth

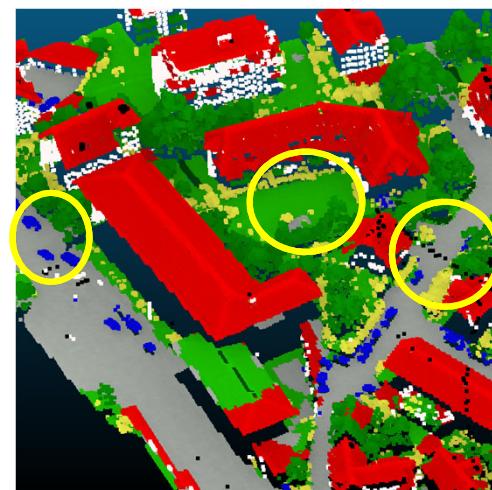


geometrie only  
(voxel value = 1)



79,8%

intensity,  
echo number,  
number of echos



84,2%

+ CIR



85,0%



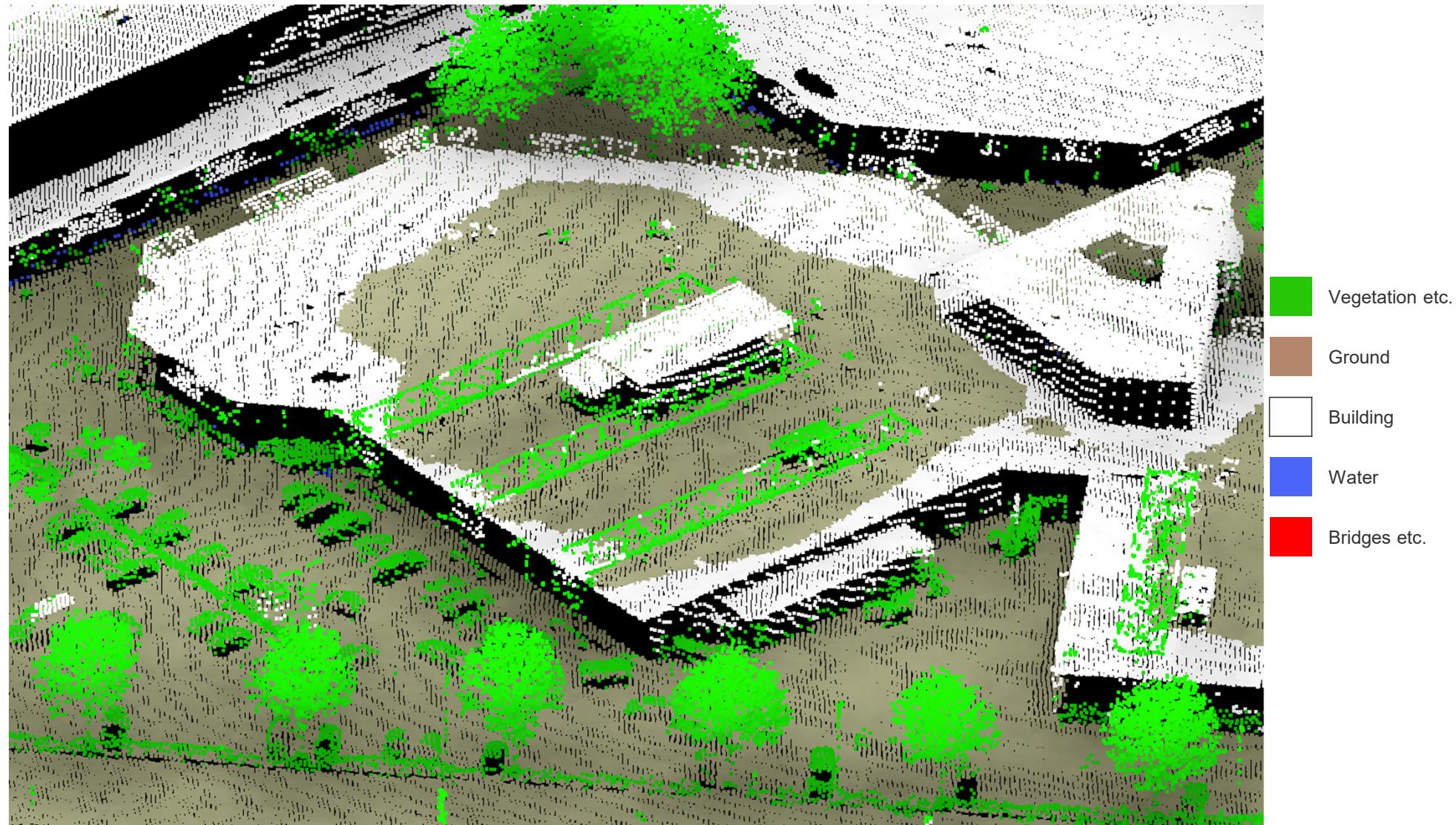
# Runtime and Memory on V3D

	Voxel size				
	<u>2.0 m</u>	<u>1.0 m</u>	<b>0.5 m</b>	<u>0.25 m</u>	<u>0.125 m</u>
<b>Graphics memory during training* (dense)</b>	1.5 GB	7.7 GB	-	-	-
<b>Graphics memory during training*</b>	0.9 GB	1.5 GB	<b>2.2 GB</b>	4.9 GB	7.9 GB
<b>Training time*</b>	6 min	14 min	<b>30min</b>	63 min	107 min
<b>Inference time*</b>	0,3 s	0,4 s	<b>0,8 s</b>	1,4 s	2,0 s
<b>Accuracy**</b>	80,3 %	83,7 %	<b>84,2 %</b>	84,2 %	83,2 %

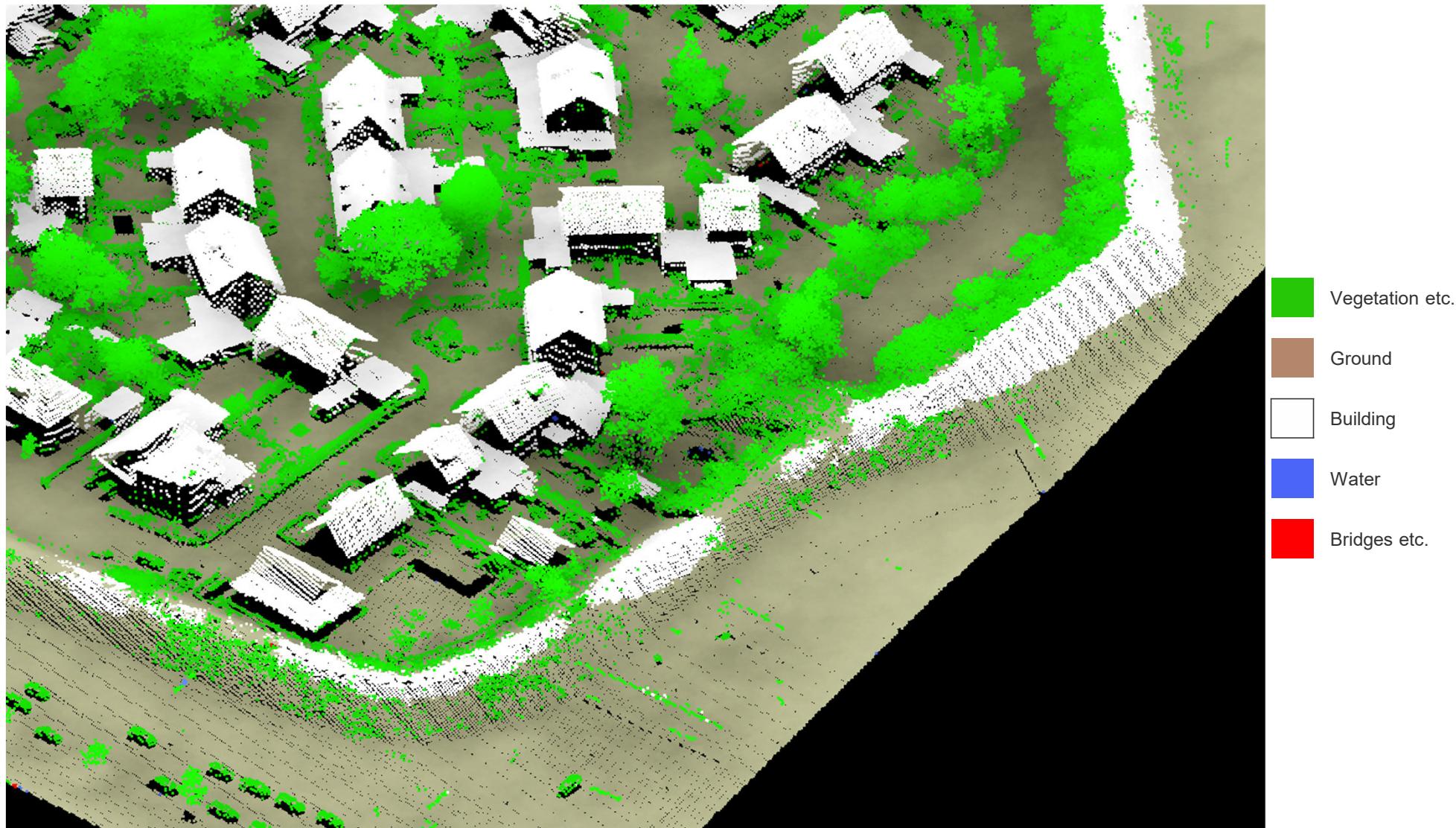
\* per network. For each configuration, results of 10 independent networks were used as an ensemble.

\*\* without CIR

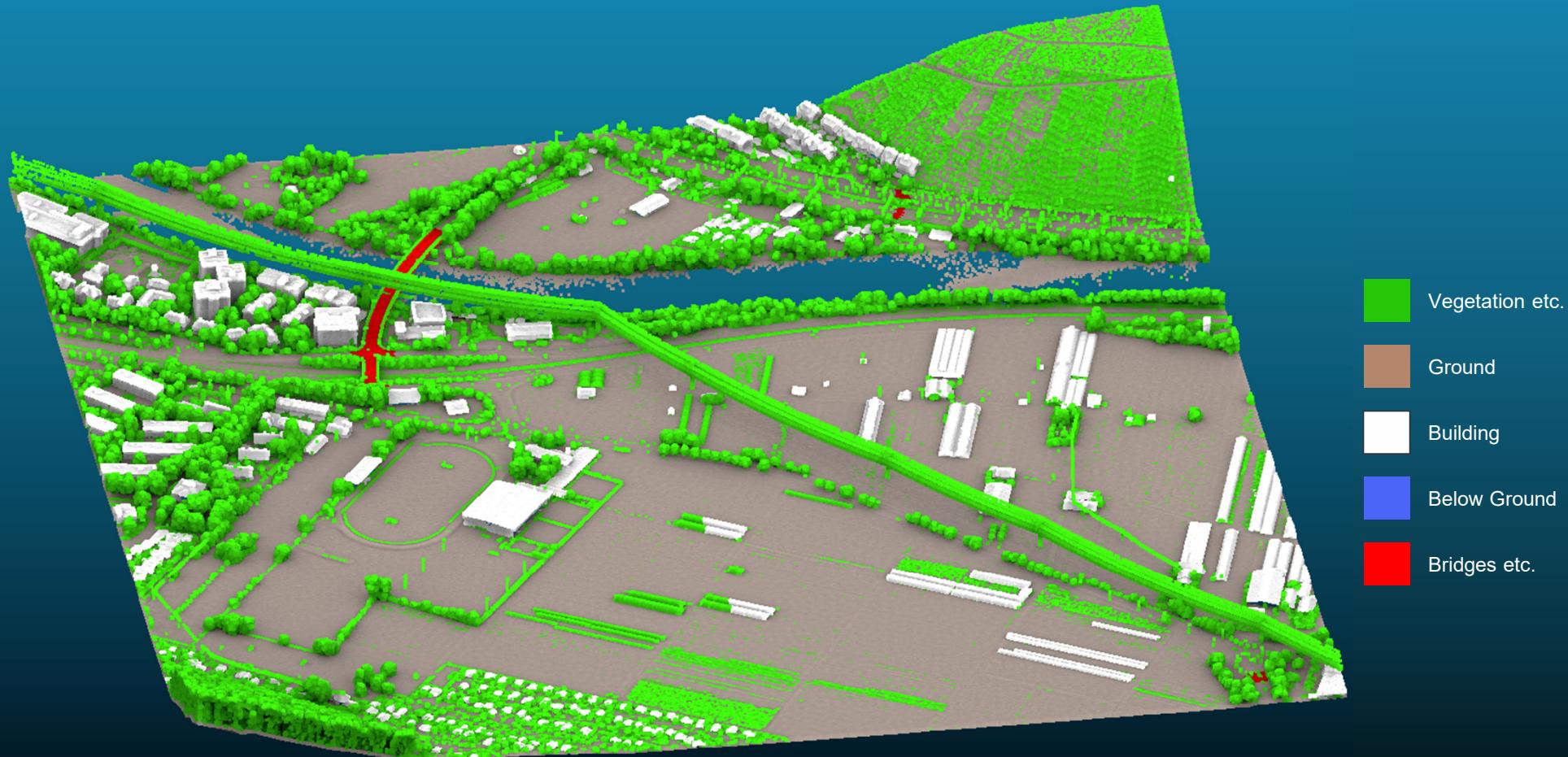
# Results on AHN3 – typical misclassifications



## Results on AHN3 – typical misclassifications



# Stuttgart – Ground Truth



# Vorarlberg – Results (training set)

		Confusion Matrix																
		Confusion Matrix																
ground truth	precision	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(12)	(14)	(15)	(17)	recall				
Ground (2)	97.5%	1.4%	0.0%	0.0%	0.3%	0.0%	0.0%	0.8%	0.0%	0.0%	0.0%	0.0%	0.0%	97.5%	150,787,216	154,607,814 = 51.6%		
Low Veg. (3)	14.8%	77.3%	4.4%	0.4%	2.3%	0.0%	0.0%	0.4%	0.0%	0.2%	0.0%	0.0%	0.0%	77.3%	3854,201,130,900	26,030,460 = 8.7%		
Medium Veg. (4)	0.0%	2.2%	79.7%	16.8%	1.3%	0.0%	0.0%	0.0%	0.0%	0.1%	0.0%	0.0%	0.0%	79.7%	2,136,44,230,227,918	39,192,738 = 13.1%		
High Veg. (5)	0.0%	0.0%	0.1%	99.9%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	99.9%	0,37,743,798	37,770,384 = 12.6%		
Building (6)	0.4%	0.5%	0.1%	0.2%	98.6%	0.0%	0.0%	0.0%	0.0%	0.1%	0.0%	0.0%	0.0%	98.6%	0,28,135,834	28,523,874 = 9.5%		
Low Point (7)	39.7%	46.7%	0.2%	0.8%	11.8%	0.0%	0.0%	0.6%	0.0%	0.3%	0.0%	0.0%	0.0%	0.0%	0,0	141,906 = 0.0%		
High Point (8)	18.4%	38.2%	2.2%	0.3%	38.4%	0.0%	0.0%	0.0%	0.0%	2.5%	0.0%	0.0%	0.0%	0.0%	0,0	1,563,950,37,408,270,28,680	8,483,940 = 2.8%	
Water (9)	1.1%	0.0%	0.0%	0.0%	0.1%	0.0%	0.0%	98.6%	0.0%	0.1%	0.0%	0.0%	0.0%	98.6%	0,2175,210	2,205,726 = 0.7%		
Other Structures (12)	nan%	nan%	nan%	nan%	nan%	nan%	nan%	nan%	nan%	nan%	nan%	nan%	nan%	nan%	0,0	0 = 0.0%		
Power Lines (14)	0.1%	0.0%	0.0%	0.0%	0.4%	0.0%	0.0%	0.0%	0.0%	99.5%	0.0%	0.0%	0.0%	99.5%	0,186,782	589,986 = 0.2%		
Masts (15)	0.2%	4.4%	6.5%	20.1%	47.4%	0.0%	0.0%	0.0%	0.0%	21.3%	0.0%	0.0%	0.0%	0.0%	0,0	65,340 = 0.0%		
Bridges (17)	33.6%	0.4%	0.0%	0.0%	66.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0,0	2,147,070 = 0.7%		
precision	96.0%	75.7%	95.7%	84.7%	81.8%	3.2%	0.0%	62.0%	0.0%	62.4%	0.0%	0.0%	0.0%	90.3%	0,027,0,788,214			
F1	96.7%	76.5%	87.0%	91.7%	89.4%	0.0%	nan%	76.2%	nan%	76.7%	nan%	nan%	nan%	nan%				

# Vorarlberg – Results (validation set)

		Confusion Matrix														
		(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(12)	(14)	(15)	(17)	recall		
ground truth	Ground (2)	95.0%	3.1%	0.0%	0.1%	1.3%	0.0%	0.0%	0.5%	0.0%	0.0%	0.0%	0.0%	95.0%		
	45,285,680 = 47.9%	43,019,770	289,734	2,804	50,756	83,056	0	0	228,152	0	9,404	0	0	43,019,776		
	Low Veg. (3)	13.8%	76.6%	4.5%	1.0%	3.6%	0.0%	0.0%	0.2%	0.0%	0.3%	0.0%	0.0%	76.6%		
	7,480,116 = 7.9%	1,033,520	29,058	88,477	76,344	67,320	16	0	12,712	0	23,172	0	0	0,5729,056		
	Medium Veg. (4)	0.0%	3.6%	71.4%	21.6%	3.3%	0.0%	0.0%	0.0%	0.0%	0.1%	0.0%	0.0%	71.4%		
	10,205,468 = 10.8%	632	367,452	285,421	201,286	86,968	0	0	0	0	13,716	0	0	0,7285,432		
	High Veg. (5)	0.0%	0.0%	0.0%	99.9%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	99.9%		
	23,719,804 = 25.1%	28	1,244	11,526	6,699,050	684	0	0	0	0	252	0	0	0,23,699,050		
	Building (6)	1.2%	1.1%	0.2%	0.2%	97.2%	0.0%	0.0%	0.0%	0.0%	0.1%	0.0%	0.0%	97.2%		
	6,426,636 = 6.8%	76,868	68,892	15,736	10,016	249,360	20	0	56	0	5,692	0	0	0,6249,360		
	Low Point (7)	74.8%	19.4%	0.0%	0.7%	3.4%	0.0%	0.0%	1.7%	0.0%	0.0%	0.0%	0.0%	0.0%		
	24,444 = 0.0%	18,272	4,744	0	164	840	0	0	416	0	8	0	0	0.0%	0	0
	High Point (8)	17.1%	59.2%	2.4%	0.4%	17.9%	0.0%	0.0%	0.0%	0.0%	3.0%	0.0%	0.0%	0.0%	0.0%	
	1,042,768 = 1.1%	178,296	17,442	25,156	3,836	86,468	68	0	268	0	31,252	0	0	0.0%	0	0
	Water (9)	17.7%	0.2%	0.0%	0.0%	0.0%	0.0%	0.0%	82.0%	0.0%	0.1%	0.0%	0.0%	82.0%		
	14,864 = 0.0%	2,632	28	0	0	0	0	0	12,184	0	20	0	0	0	12,184	
	Other Structures (12)	nan%	nan%	nan%	nan%	nan%	nan%	nan%	nan%	nan%	nan%	nan%	nan%	nan%	0.0%	
	0 = 0.0%	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	Power Lines (14)	0.9%	0.1%	0.0%	0.1%	2.4%	0.0%	0.0%	0.9%	0.0%	95.5%	0.0%	0.0%	95.5%		
	163,204 = 0.2%	1,536	176	0	244	3,948	0	0	1,520	0	155,780	0	0	0	155,780	
	Masts (15)	0.0%	2.0%	3.0%	34.9%	49.9%	0.0%	0.0%	0.0%	0.0%	10.2%	0.0%	0.0%	0.0%	0.0%	
	36,500 = 0.0%	4	736	1,112	12,736	18,196	0	0	0	0	3,716	0	0	0	0	
	Bridges (17)	33.5%	1.1%	0.0%	0.8%	64.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	
	81,784 = 0.1%	27,424	932	20	648	52,760	0	0	0	0	0	0	0	0	0	
	precision	97.0%	70.0%	94.8%	91.0%	81.1%	0.0%	0.0%	4.8%	0.0%	64.1%	0.0%	0.0%	91.2%		
		43,019,570	289,052	285,421	6,699,050	249,360	0	0	12,184	0	155,780	0	0	0,86150,648		
	F1	96.0%	73.2%	81.5%	95.2%	88.4%	nan%	nan%	9.0%	nan%	76.7%	nan%	nan%	nan%		