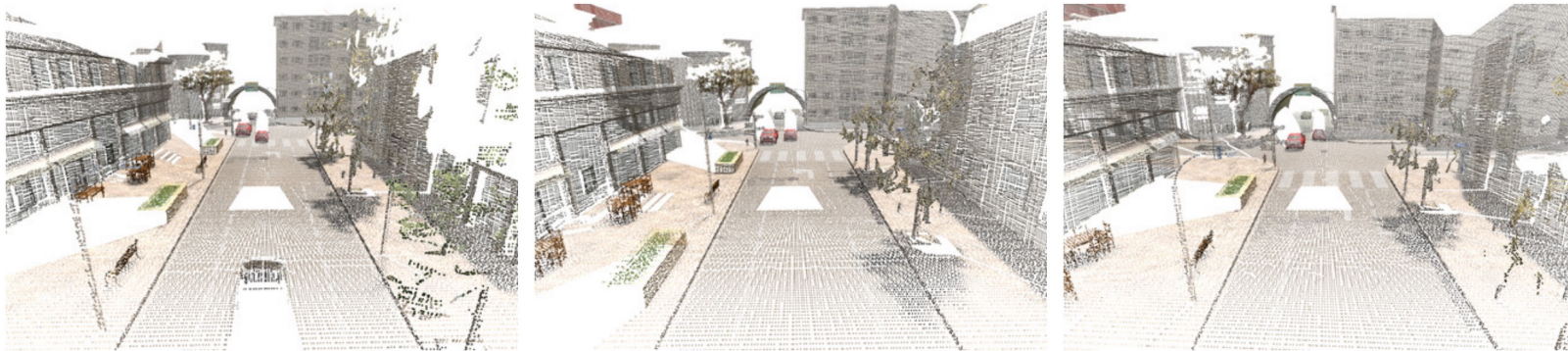


# 4D Spatio-Temporal ConvNets

Minkowski Convolutional Neural Networks

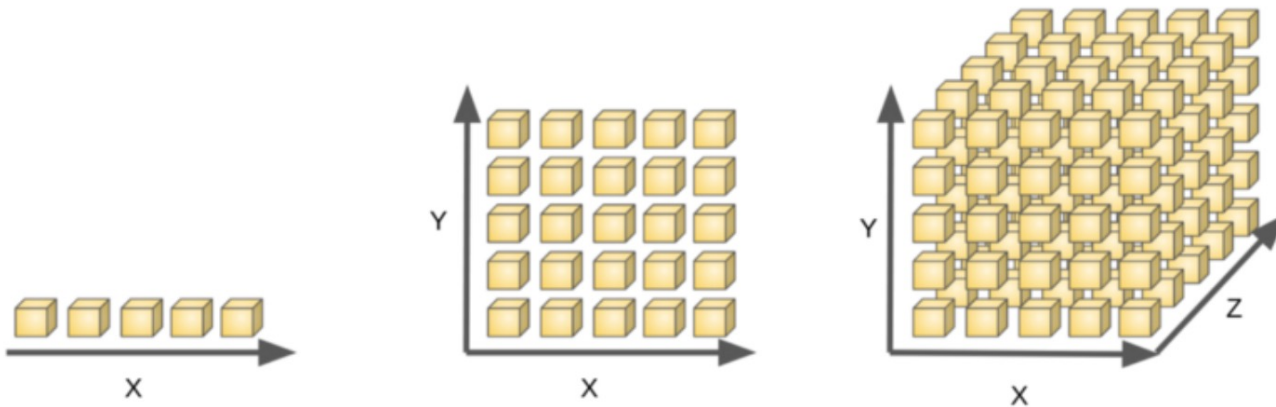


- 4D Space: 3 spatial dim., 1 time dim.
- 3D Pointcloud as a slice along the time axis
- Captured by depth sensing sensors over time
- Problem statement:
  - Input: 4D data
  - Goal: Segmentation of points

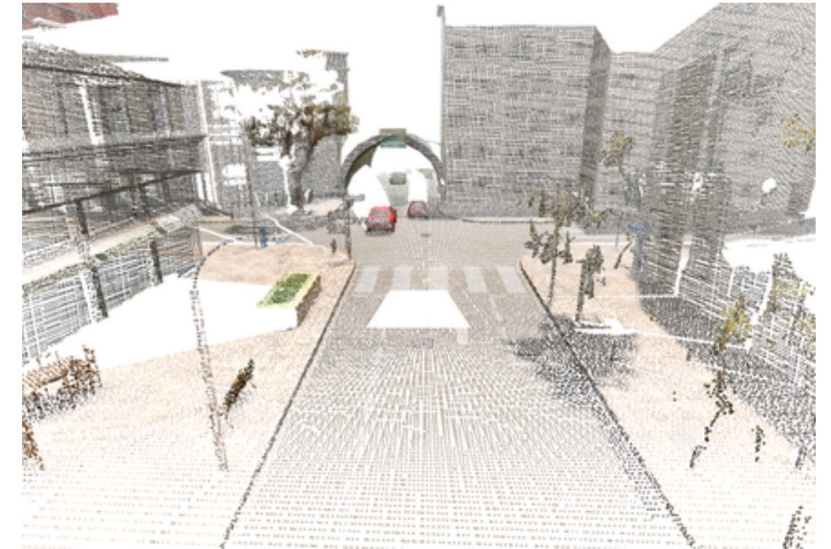


Source: 4D Spatio-Temporal ConvNets: Minkowski Convolutional Neural Networks

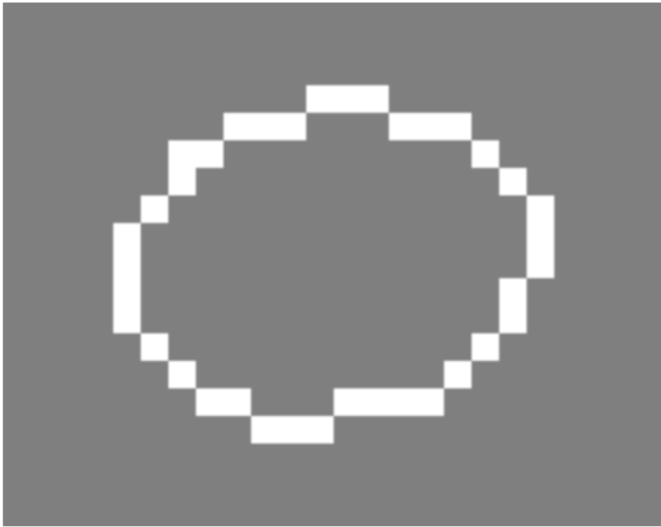
- More and more 4D data, e.g. in robotics
- But: Curse of dimensionality
  - Exponential computational effort
  - Sparse data



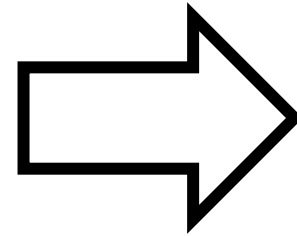
Source: <https://www.freecodecamp.org/news/the-curse-of-dimensionality-how-we-can-save-big-data-from-itself-d9fa0f872335/>



Source: 4D Spatio-Temporal ConvNets: Minkowski Convolutional Neural Networks



Source: <https://github.com/facebookresearch/SparseConvNet>



[Coordinate],  
[Value]

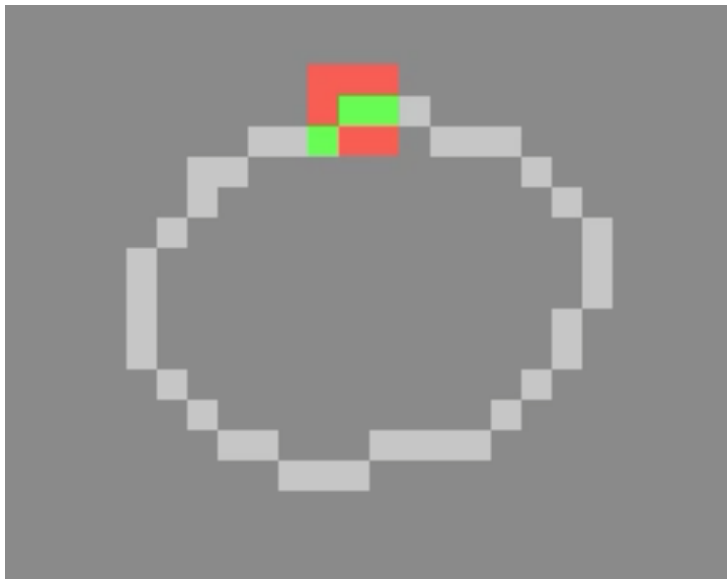
From:  $28 * 28 = 784$  *values*

To:  $50 * 2 + 50 * 1 = 150$  *values*

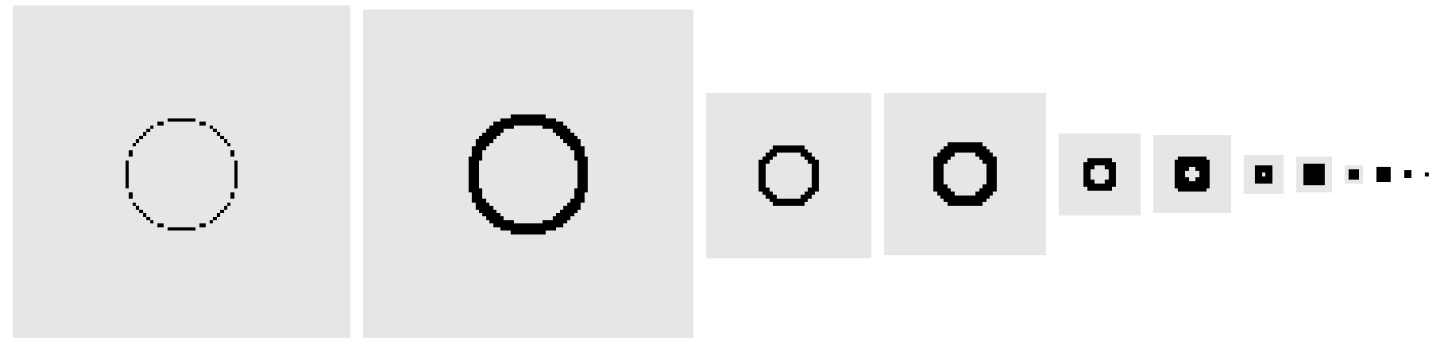
# Concepts explained – Sparse Convolutions

6

- Non-trivial convolutions – mapping from input to output
- Computation for just some outputs



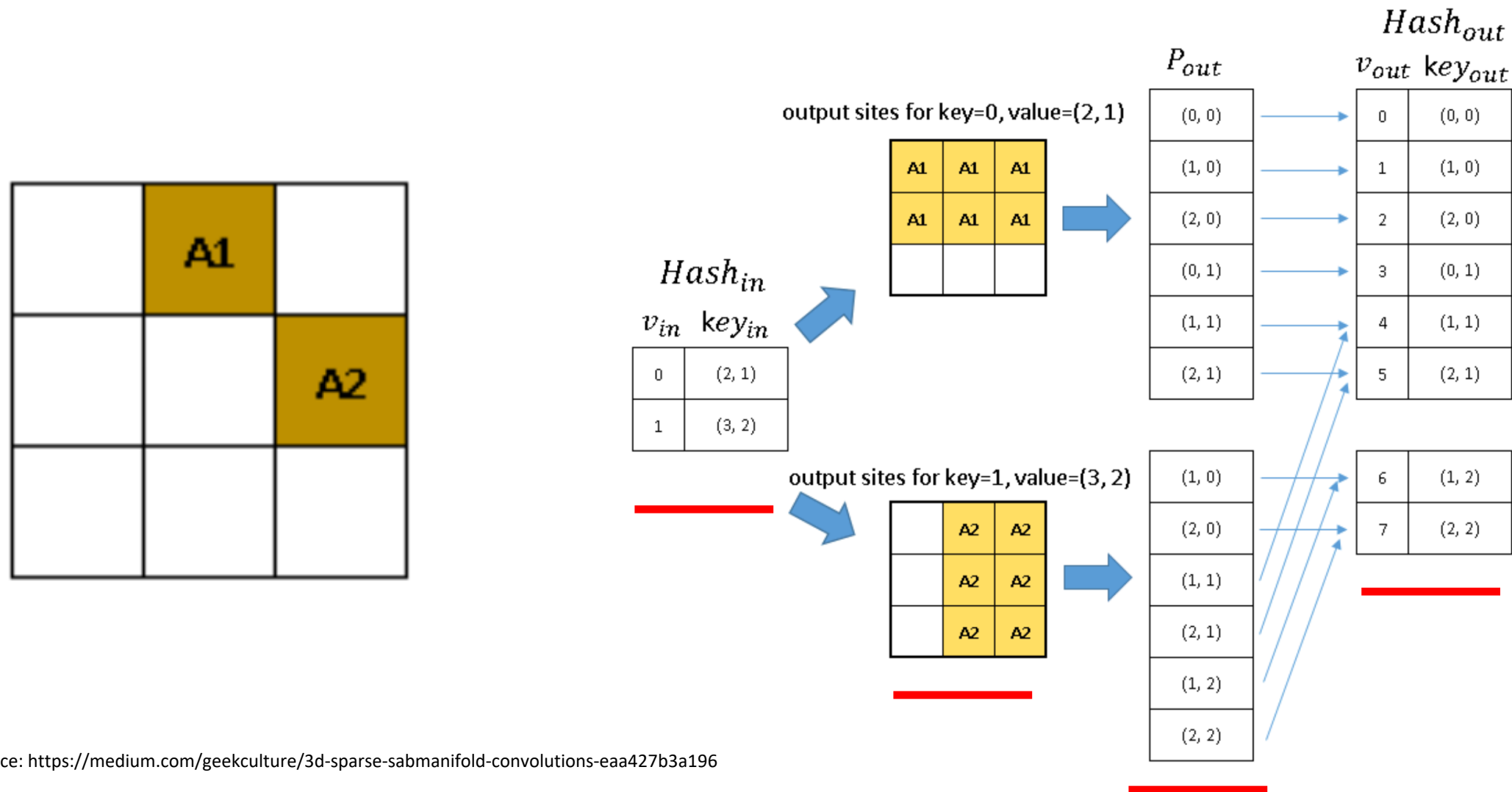
Source: <https://github.com/facebookresearch/SparseConvNet>



Source: Spatially-sparse convolutional neural networks

# Concepts explained – Sparse Convolutions

7



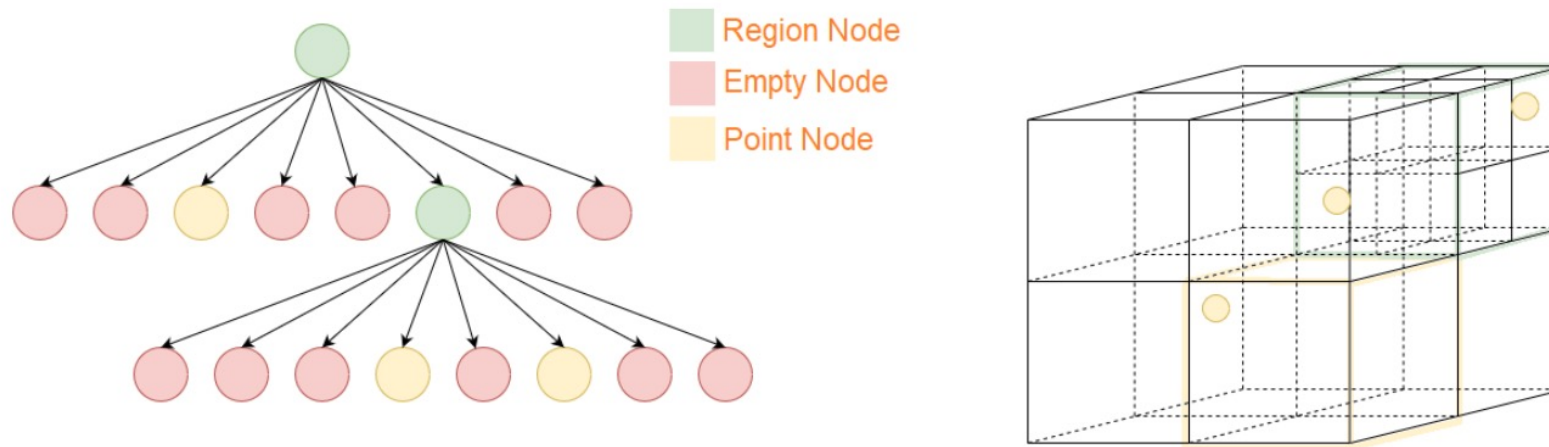
Source: <https://medium.com/geekculture/3d-sparse-sabmanifold-convolutions-eaa427b3a196>

- Deals with the curse of dimensionality
- Sparse representation – more memory efficient
- Sparse convolutions – less computational complexity

- 4D perception without NN
  - 4D Markov random field
  - Spatio-temporal CNN
- 2D video perception
- Related Work is in 3D Perception
  - Perception utilizing 3D convolutions
  - Perception without 3D convolutions

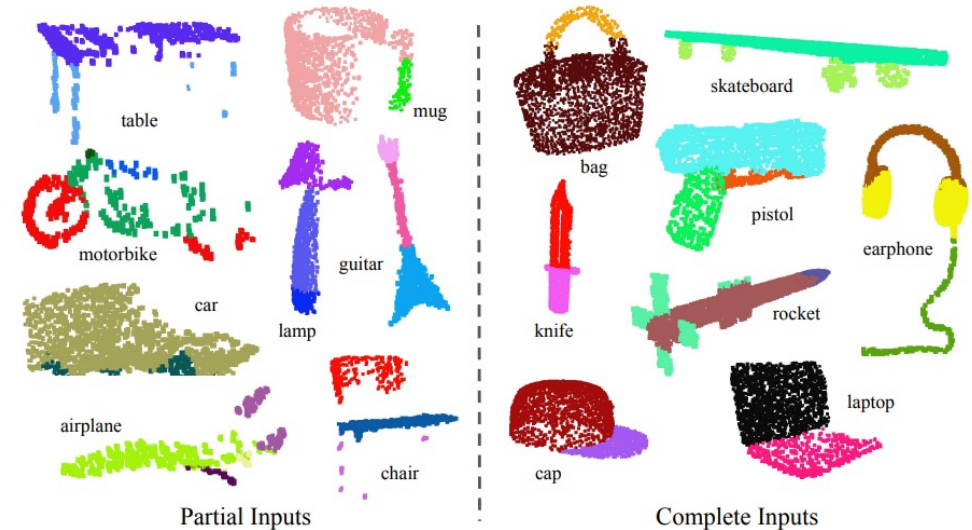


- Dense 3D convolutions
  - Suffer from the curse of dim. → OctNet
- Sparse 3D convolutions
- Continuous kernels in a continuous space



Source: <https://iq.opengenus.org/octree/>

- 2D convolutions for 3D perception
- Perception on pointclouds
  - PointNet/ PointNet++



PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation

- Coordinates with respective values
- 3 spatial dimensions and 1 temporal dimension
- Data needs to be quantized

$$C = \begin{bmatrix} x_1 & y_1 & z_1 & t_1 & b_1 \\ & & \vdots & & \\ x_N & y_N & z_N & t_N & b_N \end{bmatrix}, F = \begin{bmatrix} \mathbf{f}_1^T \\ \vdots \\ \mathbf{f}_N^T \end{bmatrix}$$

$$\begin{aligned} 1. \quad \mathbf{x}_{\mathbf{u}}^{\text{out}} &= \sum_{\mathbf{i} \in \mathcal{V}^D(K)} W_{\mathbf{i}} \mathbf{x}_{\mathbf{u}+\mathbf{i}}^{\text{in}} \quad \text{for } \mathbf{u} \in \mathbb{Z}^D \\ 2. \quad \mathbf{x}_{\mathbf{u}}^{\text{out}} &= \sum_{\mathbf{i} \in \mathcal{N}^D(\mathbf{u}, \mathcal{C}^{\text{in}})} W_{\mathbf{i}} \mathbf{x}_{\mathbf{u}+\mathbf{i}}^{\text{in}} \quad \text{for } \mathbf{u} \in \mathcal{C}^{\text{out}} \end{aligned}$$

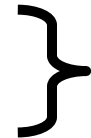
➤ We need an algorithm for this

- Max Pooling
- Global Pooling
- Average Pooling
- Sum Pooling



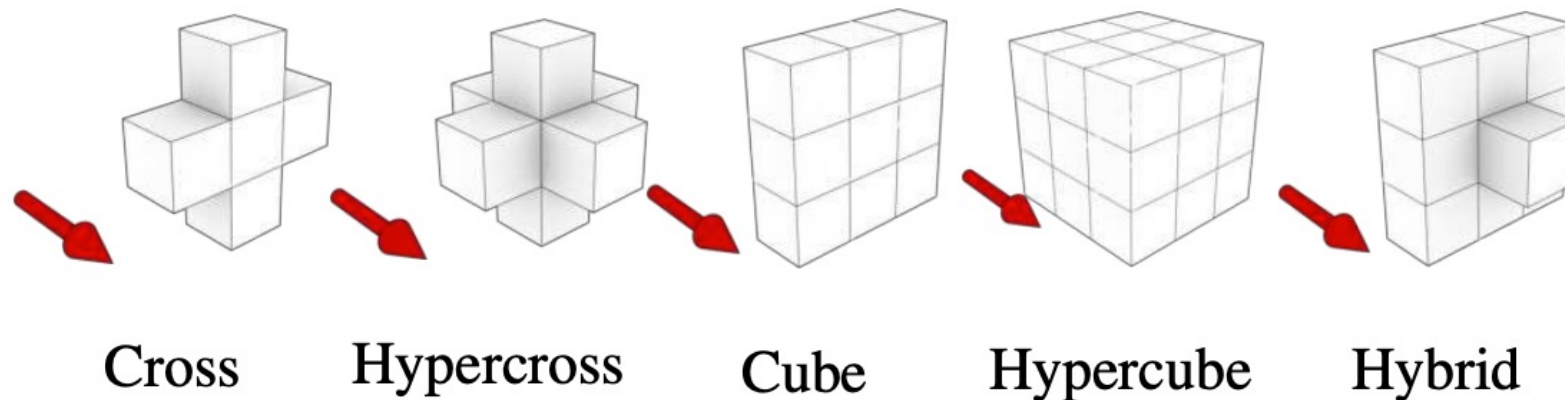
Non-trivial algorithm

- Non-linearity



Trivial

- Surfaces of objects in 3D are to be perceived
- Kernels are overparametrized
- Reduce params in kernels



Source: 4D Spatio-Temporal ConvNets: Minkowski Convolutional Neural Networks

# Methods - Full Architecture

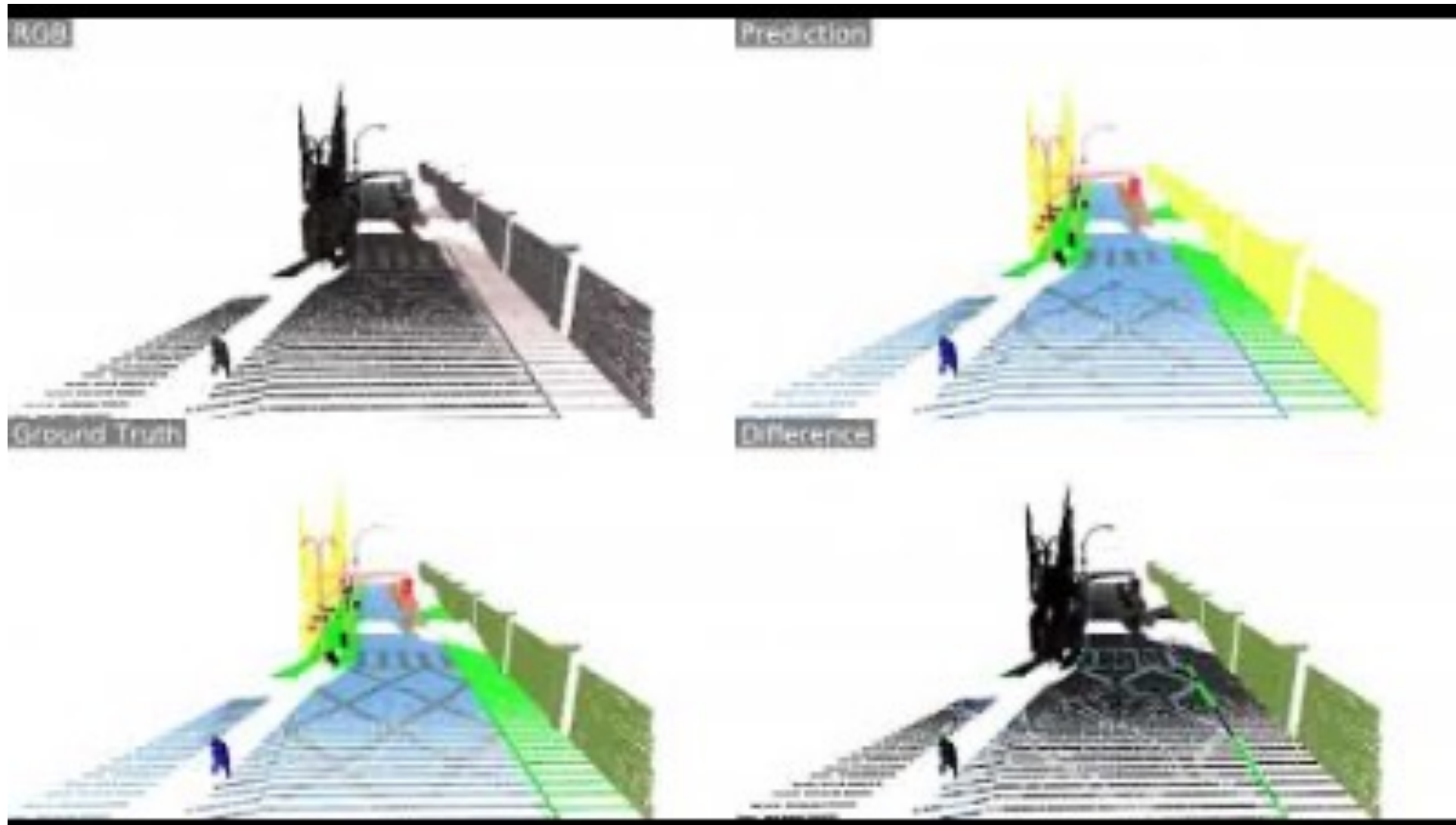
20

4D input



2D input









Source: <http://www.scan-net.org/>

Table 1: 3D Semantic Label Benchmark on ScanNet<sup>†</sup> [5]

| Method                            | mIOU        |
|-----------------------------------|-------------|
| ScanNet [5]                       | 30.6        |
| SSC-UNet [10]                     | 30.8        |
| PointNet++ [24]                   | 33.9        |
| ScanNet-FTSDF                     | 38.3        |
| SPLATNet [29]                     | 39.3        |
| TargetConv [30]                   | 43.8        |
| SurfaceConv [21]                  | 44.2        |
| 3DMV <sup>‡</sup> [6]             | 48.4        |
| 3DMV-FTSDF <sup>‡</sup>           | 50.1        |
| PointNet++SW                      | 52.3        |
| MinkowskiNet42 (5cm)              | <b>67.9</b> |
| MinkowskiNet42 (2cm) <sup>†</sup> | 72.1        |
| SparseConvNet [10] <sup>†</sup>   | <b>72.5</b> |

<sup>†</sup>: post-CVPR submissions. <sup>‡</sup>: uses 2D images additionally. Per class IoU in the supplementary material. The parenthesis next to our methods indicate the voxel size.

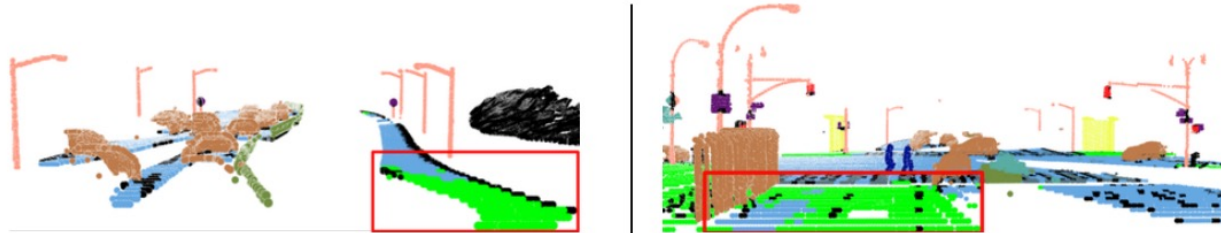
Source: 4D Spatio-Temporal ConvNets: Minkowski Convolutional Neural Networks

Table 3: Segmentation results on the noisy Synthia 4D dataset

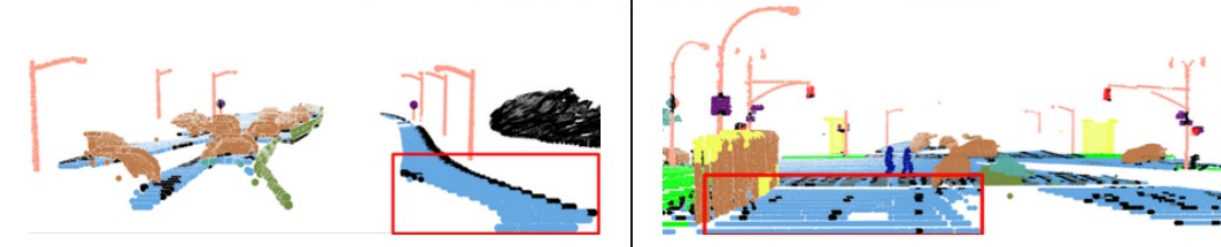
| IoU                    | Building      | Road          | Sidewalk      | Fence         | Vegetation    | Pole          | Car           | Traffic Sign  | Pedestrian    | Lanemarking   | Traffic Light | mIoU          |
|------------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| 3D MinkNet42           | 87.954        | 97.511        | 78.346        | 84.307        | 96.225        | 94.785        | 87.370        | 42.705        | <b>66.666</b> | 52.665        | 55.353        | 76.717        |
| 3D MinkNet42 + TA      | 87.796        | 97.068        | 78.500        | 83.938        | 96.290        | 94.764        | 85.248        | 43.723        | 62.048        | 50.319        | 54.825        | 75.865        |
| 4D Tesseract MinkNet42 | <b>89.957</b> | 96.917        | 81.755        | 82.841        | 96.556        | 96.042        | 91.196        | 52.149        | 51.824        | <b>70.388</b> | <b>57.960</b> | 78.871        |
| 4D MinkNet42           | 88.890        | <b>97.720</b> | <b>85.206</b> | <b>84.855</b> | <b>97.325</b> | <b>96.147</b> | <b>92.209</b> | <b>61.794</b> | 61.647        | 55.673        | 56.735        | <b>79.836</b> |

TA denotes temporal averaging. As the input pointcloud coordinates are noisy, averaging along the temporal dimension introduces noise.

3D:



4D:



- Generalizes (sparse) convolutions to any dimension
  - Introduces different kernel shapes
- Coming up with concrete algorithms for these ideas

- Deep Learning on Dynamic 3D Point Cloud Sequences:
  - “Quantization introduces error”
  - Their method looks at the neighbours of each point instead of the space itself

