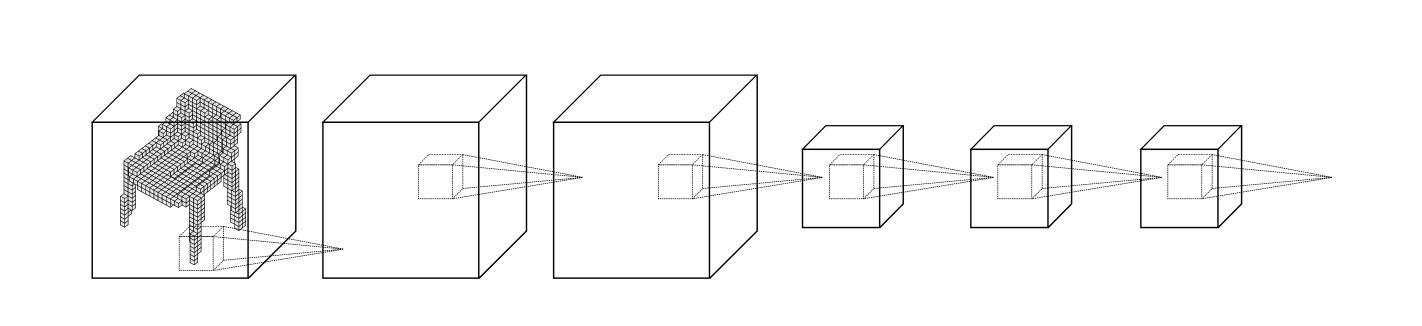


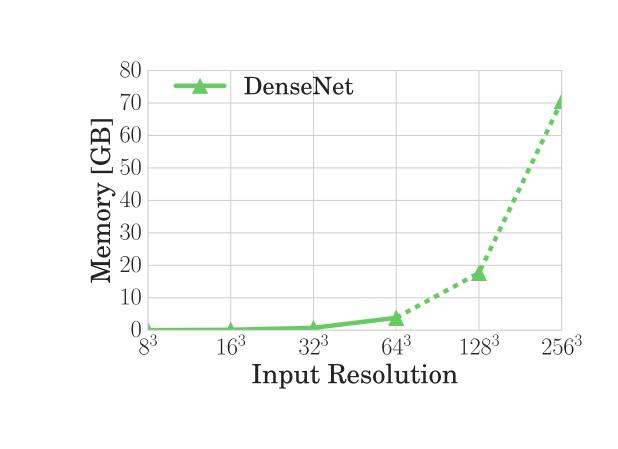
OctNet: Learning Deep 3D Representations at High Resolutions

Gernot Riegler¹, Ali Osman Ulusoy^{2,3}, Andreas Geiger^{2,4} 1 Graz University of Technology 2 MPI for Intelligent Systems Tübingen 3 Microsoft 4 ETH Zürich

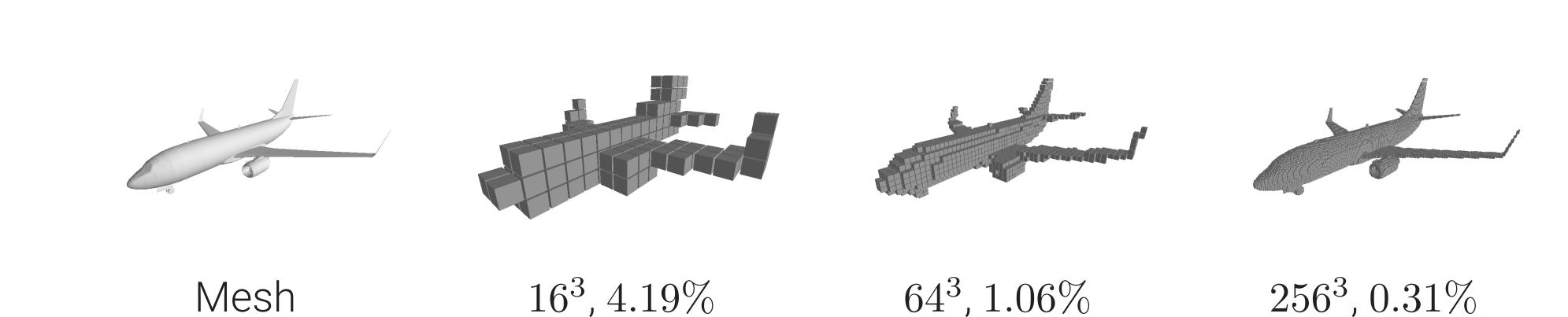
Motivation

In 3D convolutional networks memory requirements increase cubically wrt. input resolution



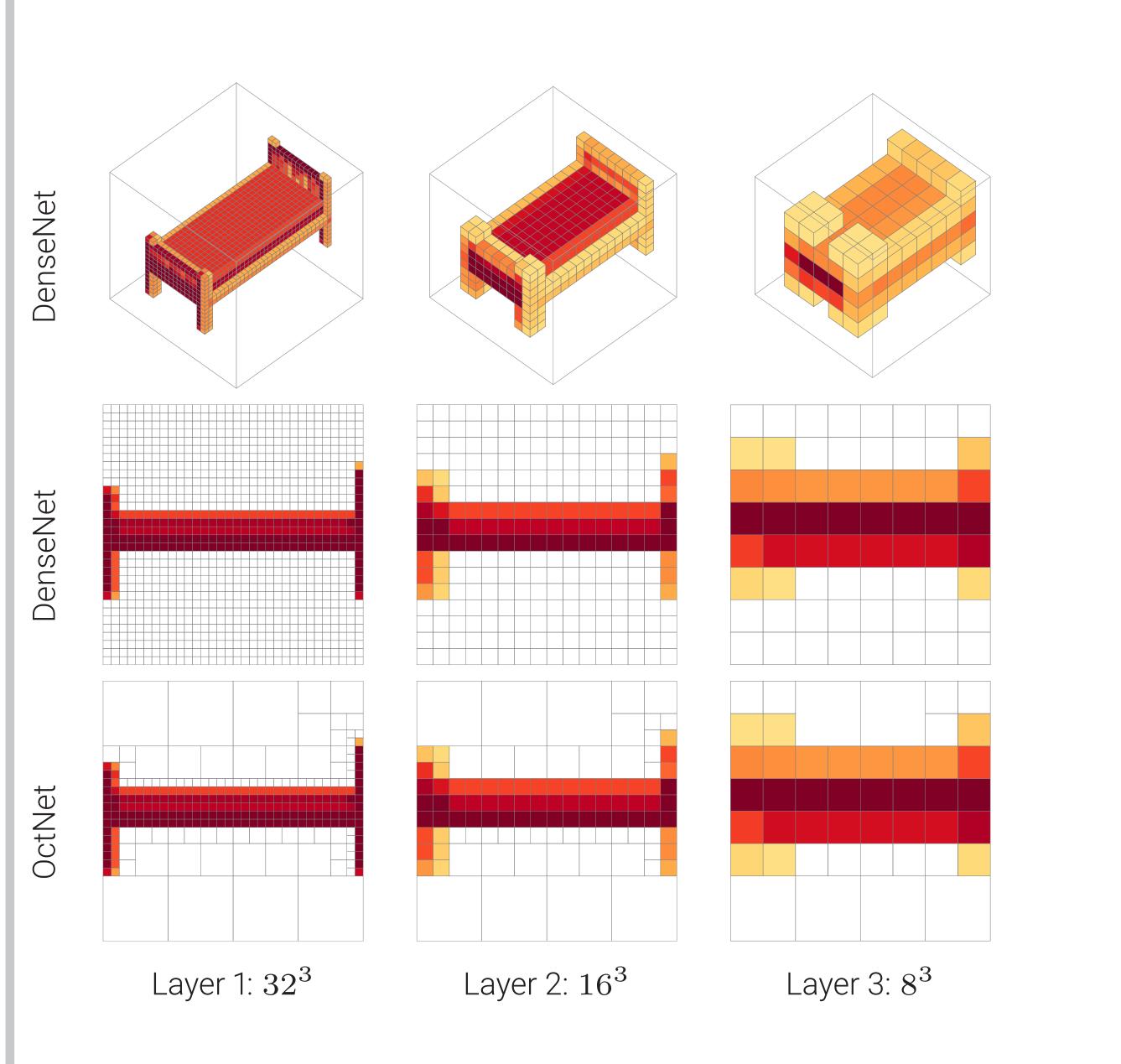


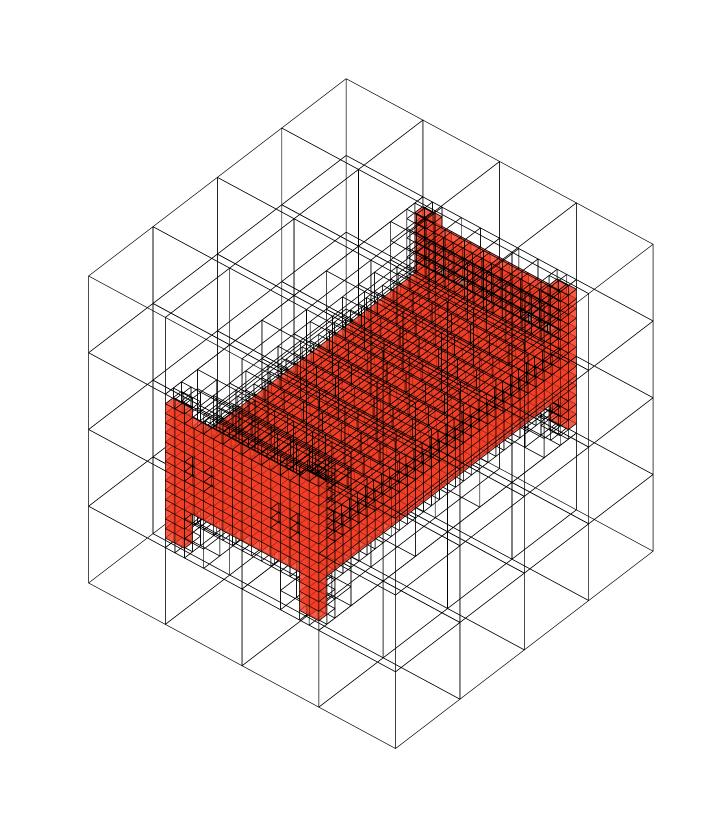
3D data is usually very sparse



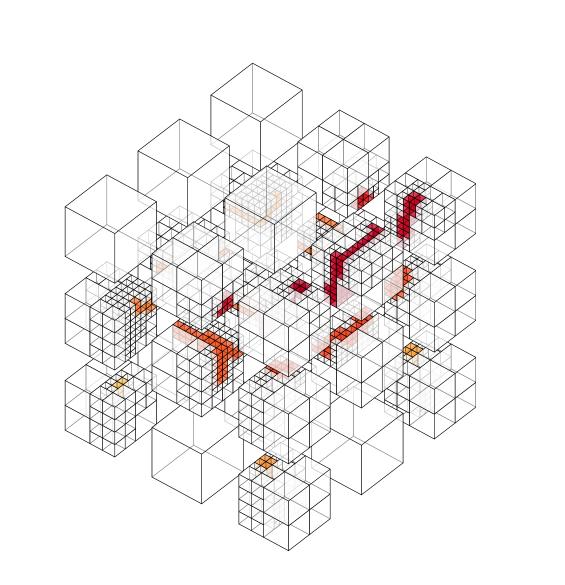
Idea

Use space partitioning data-structure within convolutional networks





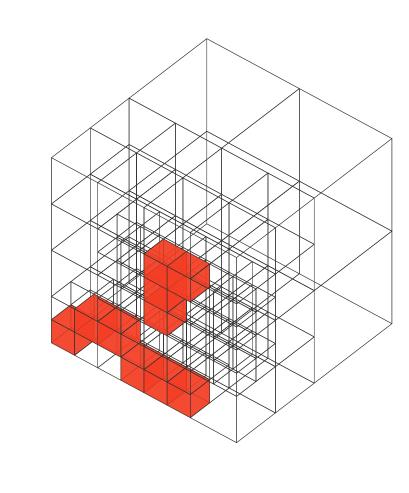
OctNet Representation

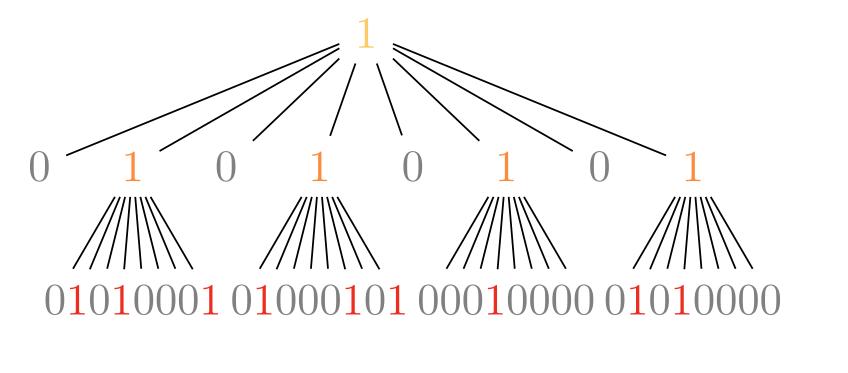


Data structure: grid of shallow octrees [1]

- Shallow octrees have fixed depth
- Efficiently encoded as bit-string
- Bit indicates if node is split, or not
- Fast neighbour address resolution via bit operations

01010000





$$\operatorname{pa}(i) = \left\lfloor \frac{i-1}{8} \right\rfloor, \quad \operatorname{ch}(i) = 8 \cdot i + 1,$$

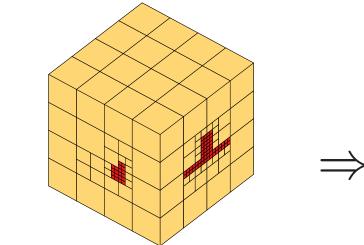
$$\operatorname{data_idx}(i) = 8 \sum_{j=0}^{\operatorname{pa}(i)-1} \operatorname{bit}(j) + 1 - \sum_{j=0}^{i-1} \operatorname{bit}(j) + \underbrace{\operatorname{mod}\,(i-1,8)}_{\text{offset}}$$
 #split nodes pre i

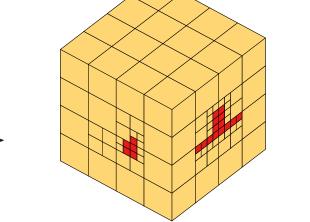
OctNet Pooling

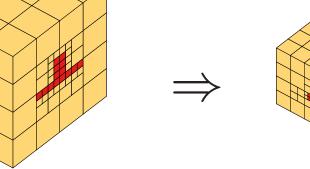
Combines shallow octrees Data on finest resolution is pooled

$$O^{\text{out}}[i, j, k] = \begin{cases} O^{\text{in}}[2i, 2j, 2k] & \text{if } \text{vxd}(2i, 2j, 2k) < 3 \\ P & \text{else} \end{cases}$$

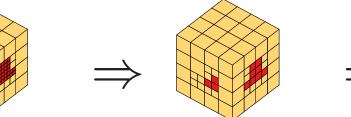
$$P = \max_{l, m, n \in [0, 1]} (O^{\text{in}}[2i + l, 2j + m, 2k + n])$$



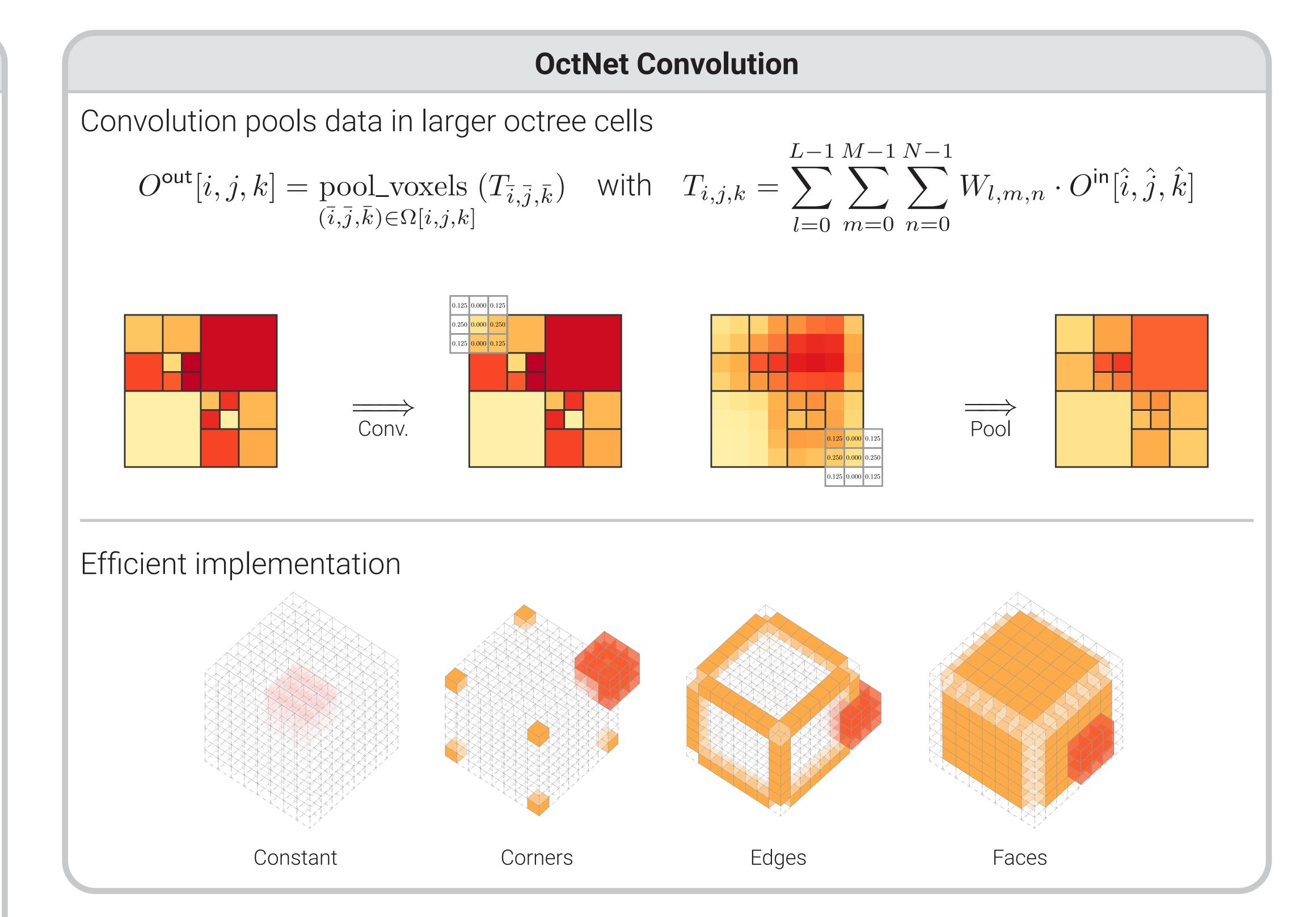








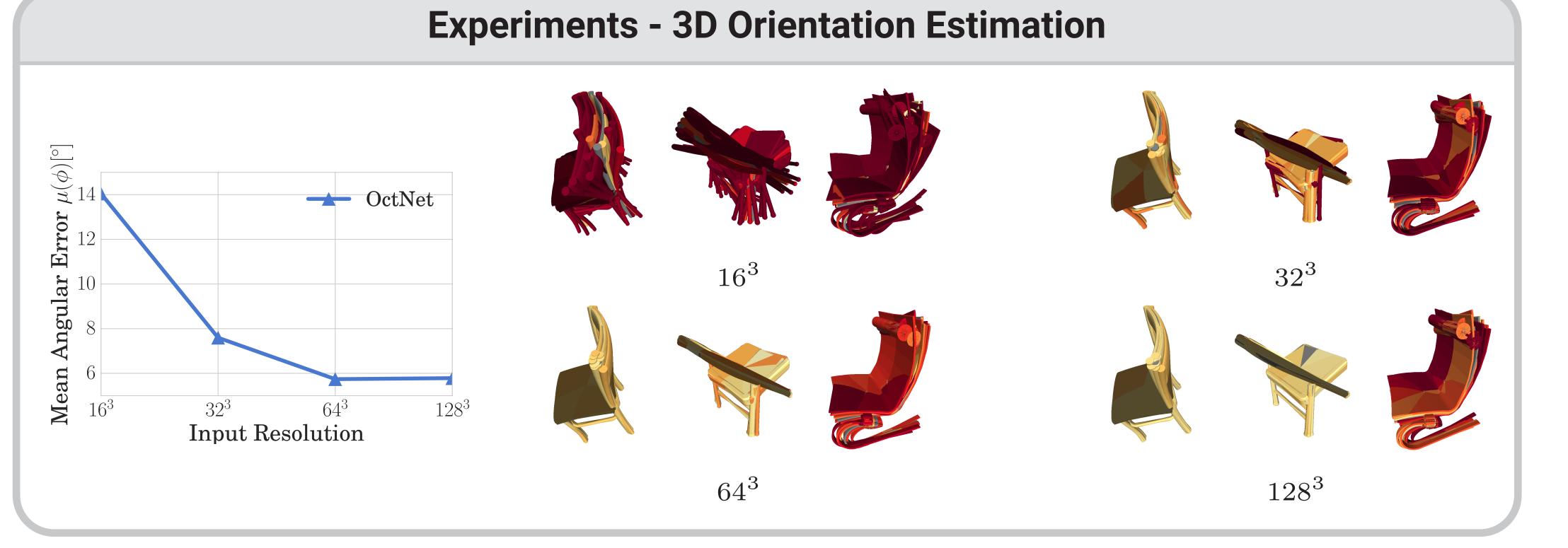


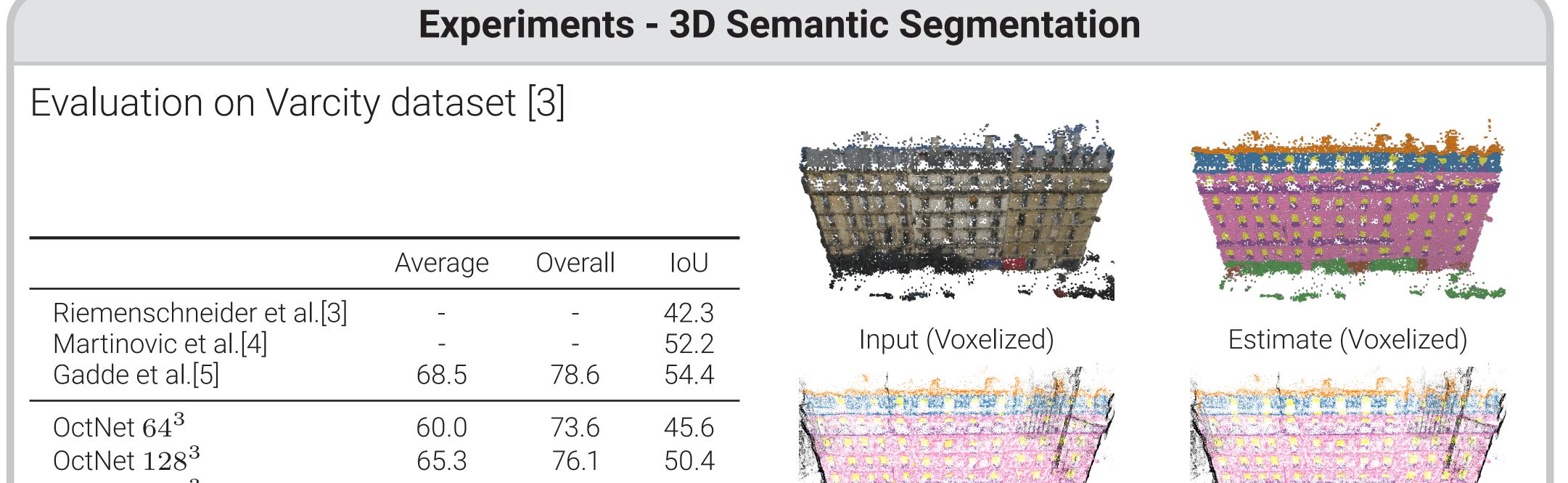


Experiments - 3D Classification Evaluation on ModelNet dataset [2] Memory & Runtime — Dense 0.8^{3} 16^{3} 32^{3} 64^{3} 128^{3} 256^{3} Input Resolution 0.00 0.01 0.02 0.03 0.04 0.05 0.06 256^{3} Memory Runtime Accuracy 2 VoxNet 0.70^{3} 16^{3} 32^{3} 64^{3} 128^{3} 256^{3} 8^3 16^3 32^3 64^3 128^3 256^3 Input Resolution ModelNet40 ModelNet10









Conclusion

OctNet makes high resolution 3D convolutional networks tractable Can be extended to generate high resolution 3D output [6] Code is online: https://github.com/griegler/octnet



References

- [1] Andrew Miller et al. "Real-time rendering and dynamic updating of 3-d volumetric data". In: *GPGPU*. 2011.
- [2] Zhirong Wu et al. "3D ShapeNets: A deep representation for volumetric shapes". In: *CVPR*. 2015.
- [3] Hayko Riemenschneider et al. "Learning Where to Classify in Multi-view Semantic Segmentation". In: ECCV. 2014.
- [4] Andelo Martinovic et al. "3D All The Way: Semantic Segmentation of Urban Scenes from Start to End in 3D". In: CVPR. 2015
- [5] Raghudeep Gadde et al. "Efficient 2D and 3D Facade Segmentation using Auto-Context". In: arXiv.org 1606.06437 (2016).
- [6] Gernot Riegler et al. "OctNetFusion: Learning Depth Fusion from Data". In: arXiv preprint arXiv:1704.01047 (2017).