

3D shape understanding

682: Neural Networks: A Modern Introduction

Subhransu Maji

April 21, 2026

College of
INFORMATION AND
COMPUTER SCIENCES



Agenda

3D representations

3D recognition architectures

- Multi-view methods
- Voxel-based methods
- Point-based methods

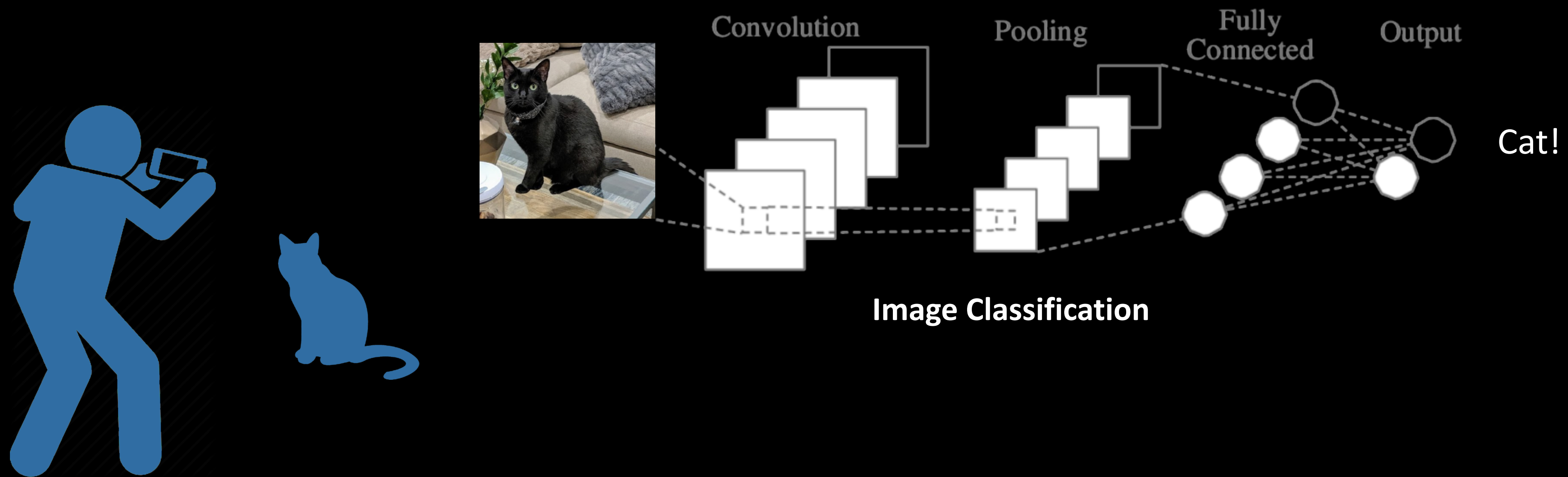
Recent trends

- Image to 3D

Slides credits: Hang Su & Hao Su

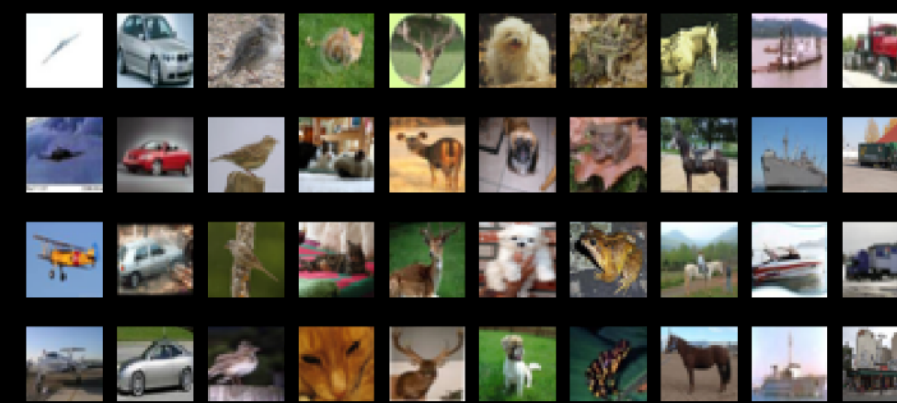
3D World Through 2D Lenses

Traditionally, visual recognition algorithms rely primarily on **2D input signals**, even though the targets are almost always **3D entities**.

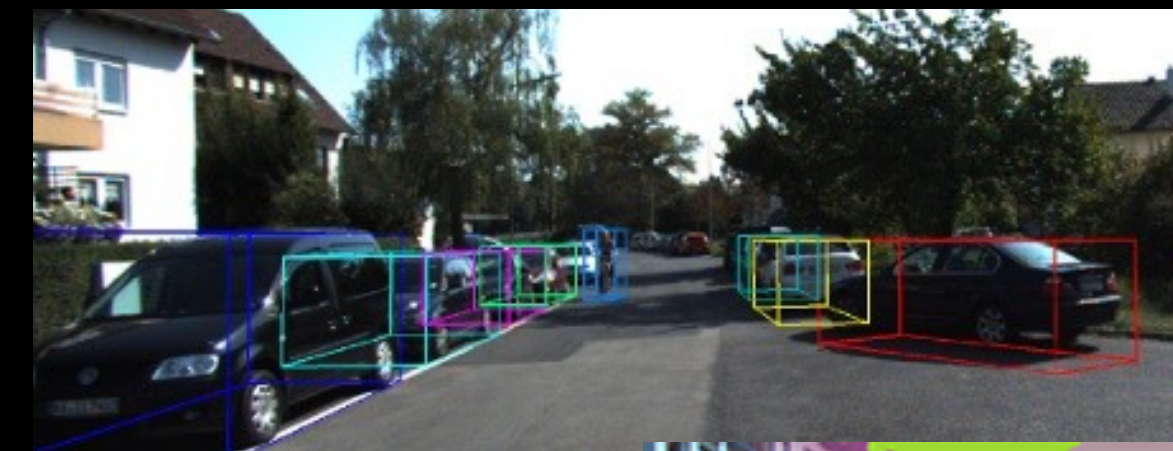


3D World Through 2D Lenses

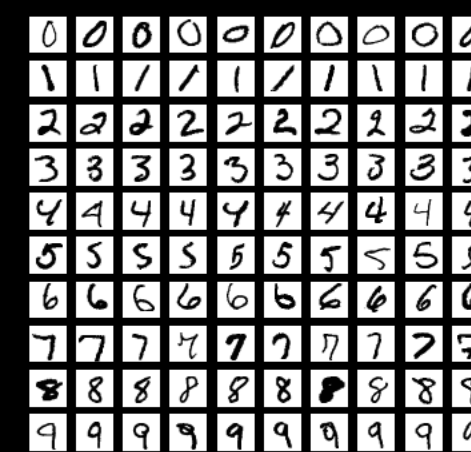
Traditionally, visual recognition algorithms rely primarily on **2D input signals**, even though the targets are almost always **3D entities**.



CIFAR-10



KITTI



MNIST

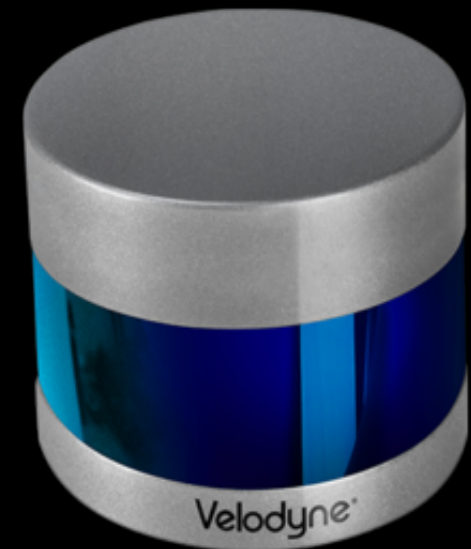


A Growing Interest in 3D Recognition Models

- 3D sensors have become better and more accessible.
- Applications in autonomous driving, robotics, consumer electronics, etc.



HDL-64E (\$85k, 2007)



VLP-16 (\$8k, 2014)



OS-1-16 (\$3.5k, 2017)



Livox Horizon (\$1k, 2020)

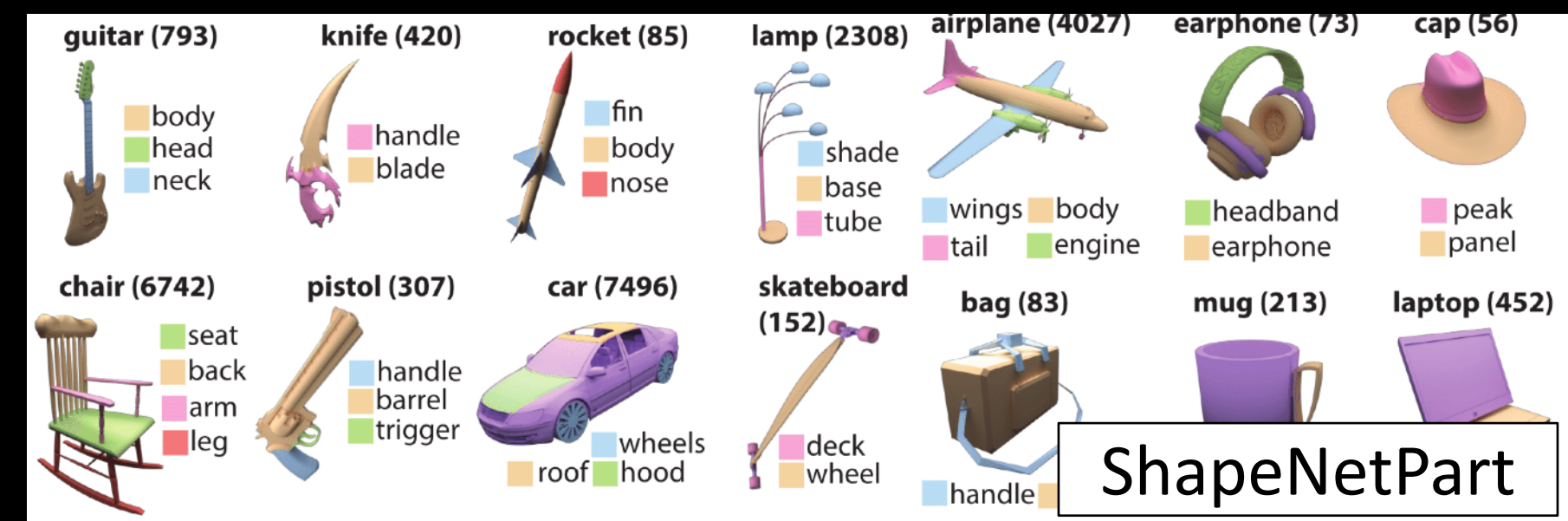
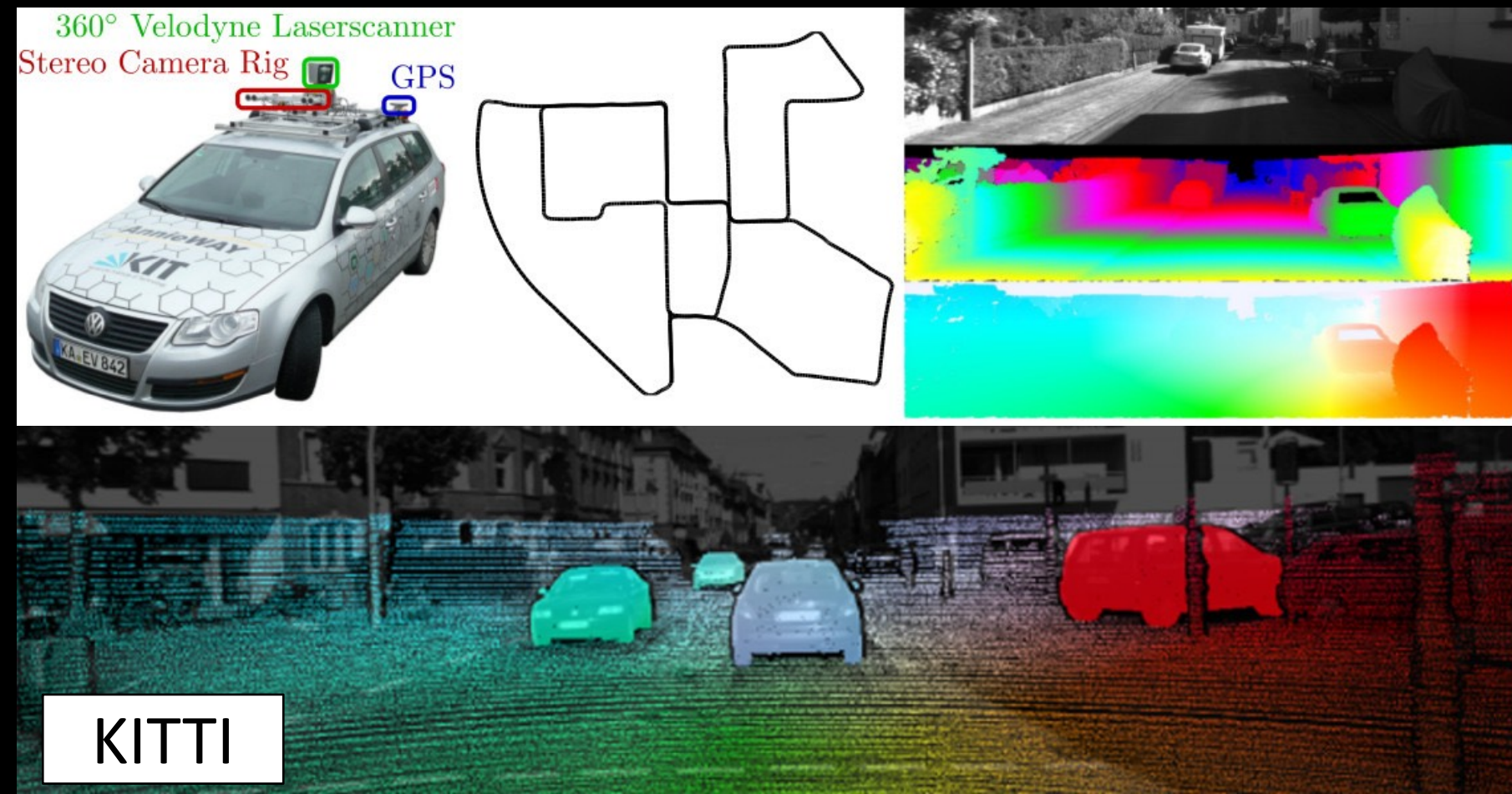


Velabit (\$100, 2020)



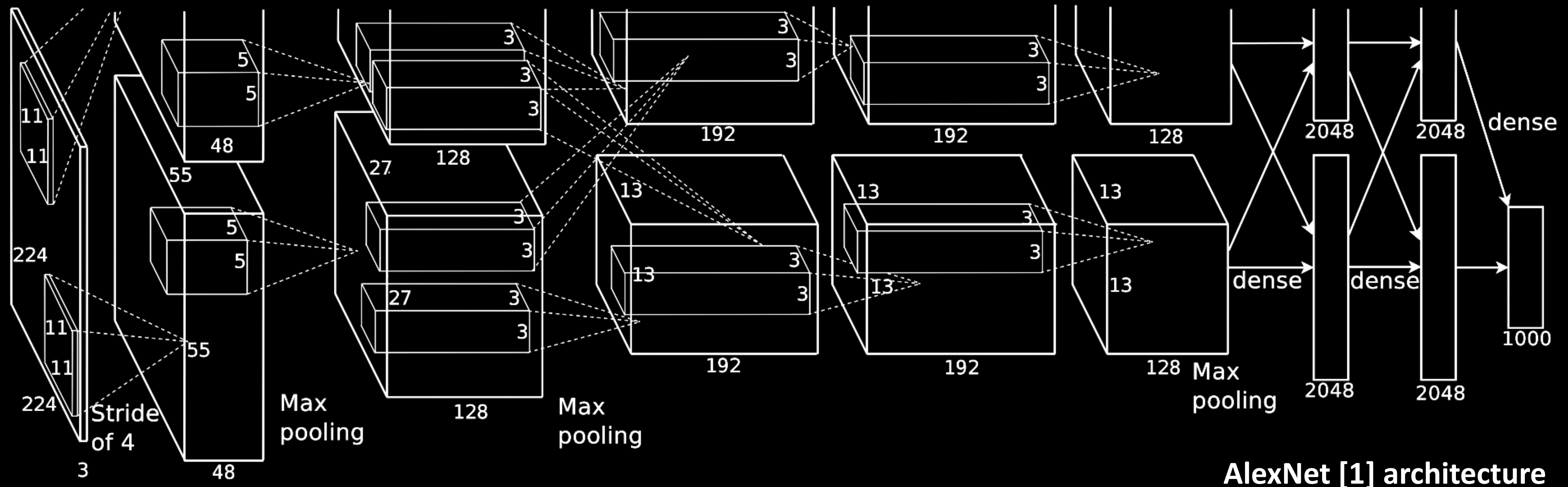
A Growing Interest in 3D Recognition Models

- Large-scale 3D datasets with annotations become available.



Challenges

- Convolutional Neural Networks (CNNs) achieve great success in 2D visual recognition tasks.
- Extending such models to 3D naively requires **dense regular grid structure**.

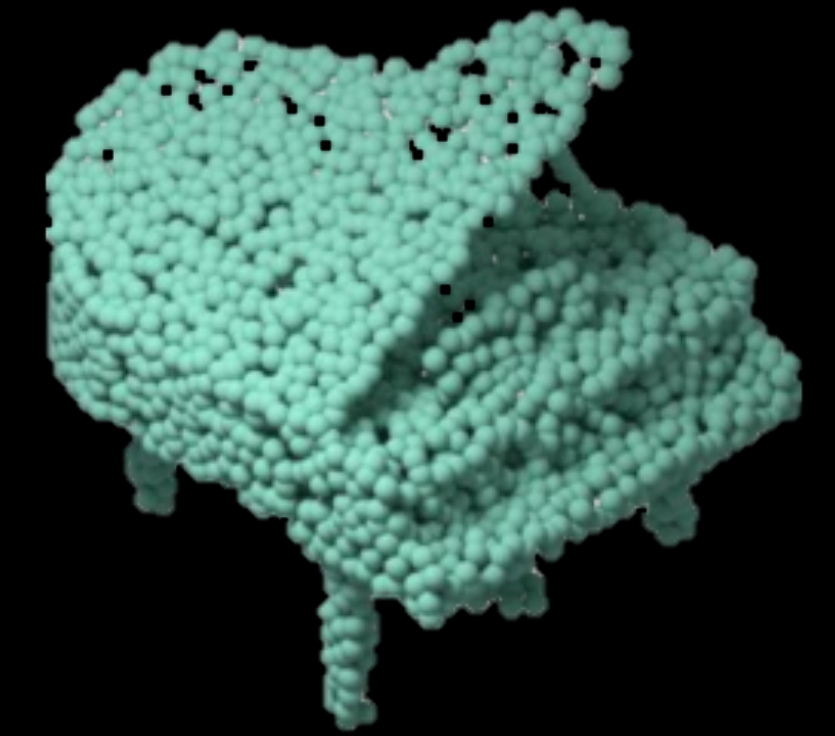


Challenges

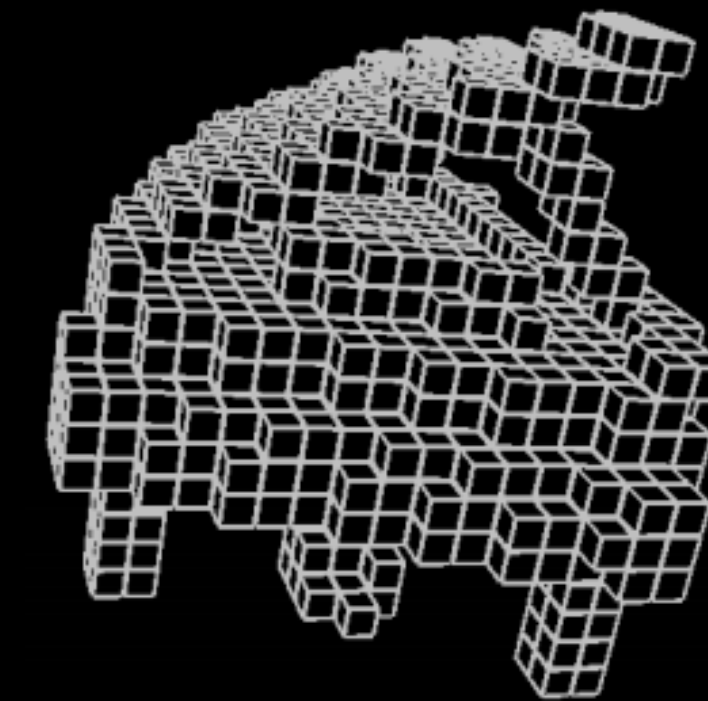
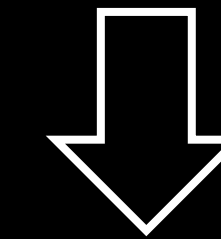
- Common representations of raw 3D data:
 - Polygon meshes (e.g. from CAD models)
 - Point clouds (e.g. from LiDAR sensors)
- Both lack regular grid structures to perform spatial convolution.
- A workaround: converting to voxels (volumetric formats)
 - ✓ 3D spatial convolutions readily apply
 - ✗ Losing details
 - ✗ Efficiency trade-offs



Polygon Meshes



Point Clouds



Voxels

3D shape classification



chair

Multi-view CNN

Image-Based 3D Representations

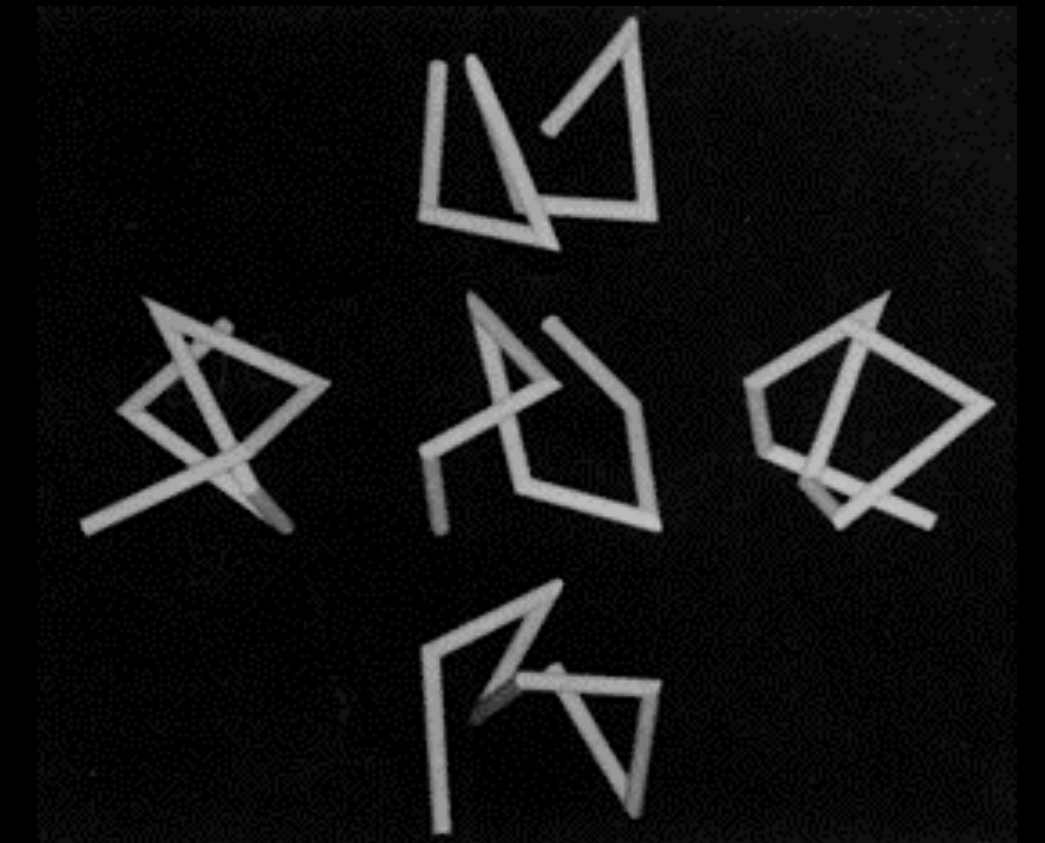
Can the 3D world be effectively represented with its 2D views?

Evidences human brain models objects as “collections of views”

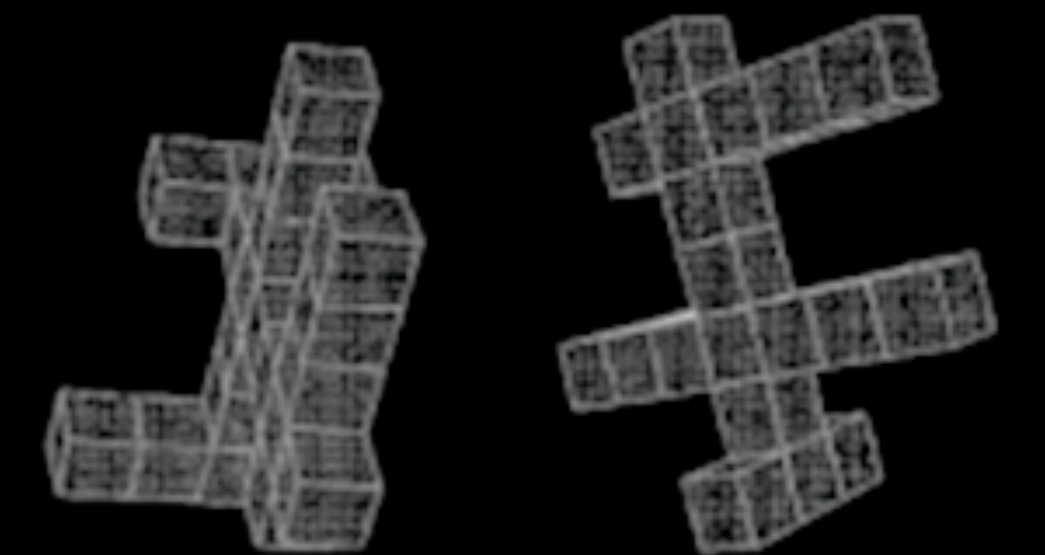
- *Faster and more accurate* to recognize when given some views of objects
- Some neurons were found to respond selectively to particular views

Some existing 3D descriptors are indeed image-based

- E.g. Light Field Descriptor
- Also work based on Fisher vectors, Gabor filters, SIFT features, etc.



[Bülthoff and Edelman, '91]



[Tarr, '95]

Image-Based 3D Representations

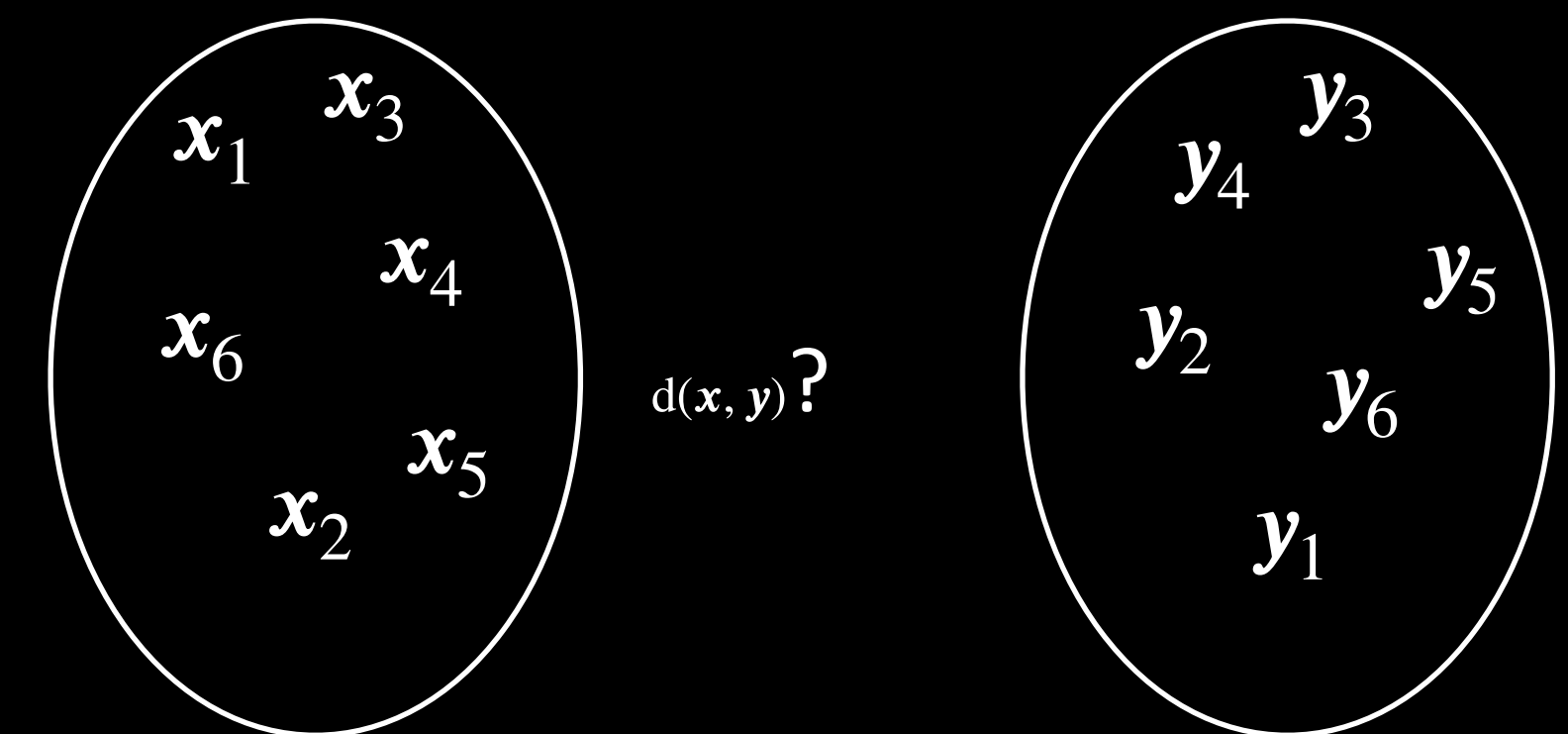
Can we use 2D CNNs for 3D recognition?

Benefits:

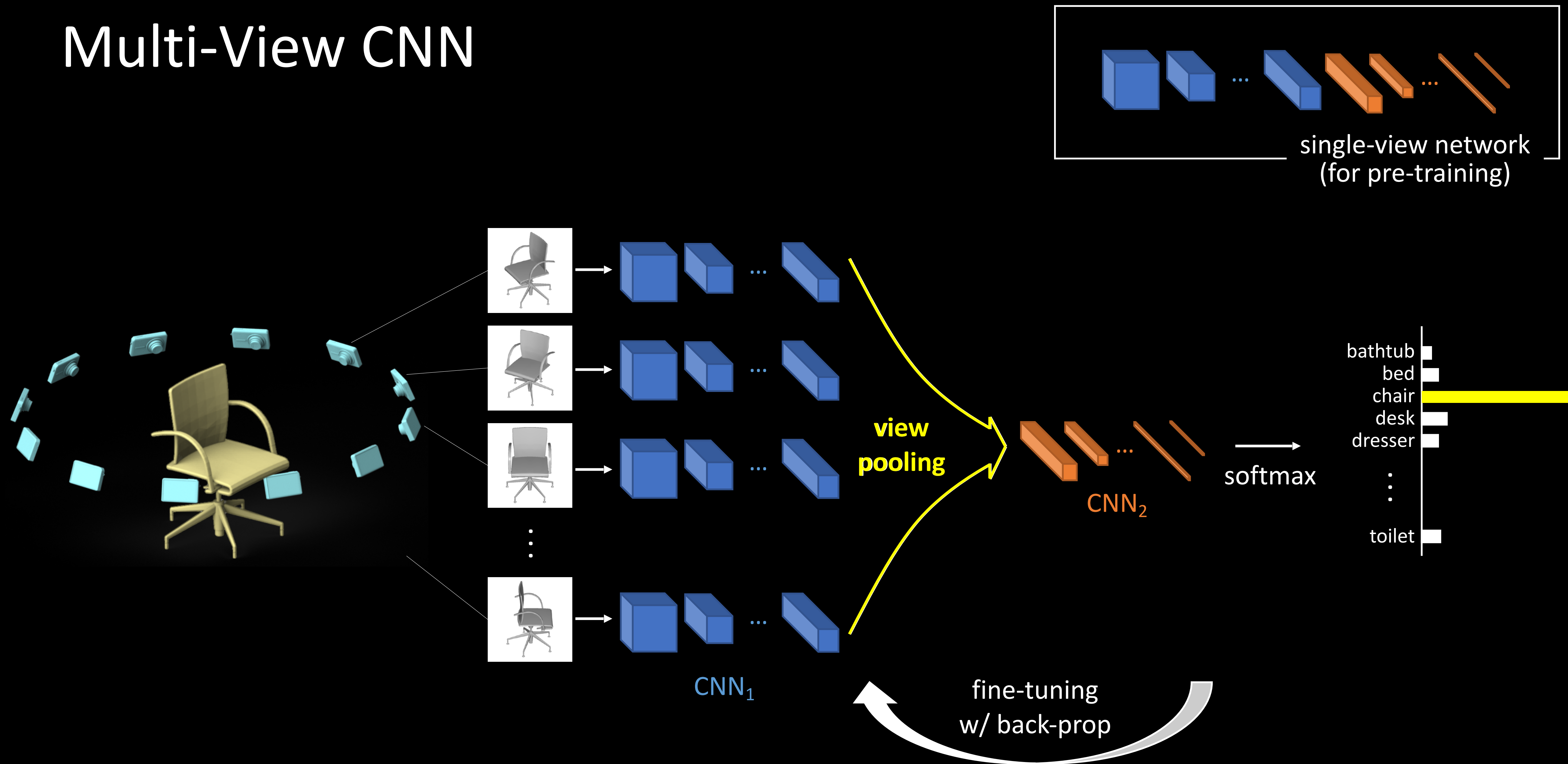
- Leverage powerful 2D architectures and pre-trained models
- Avoid expensive 3D operations
- Relate 3D shapes to images \rightarrow enables flexible applications

Challenge:

- A list/set of descriptors is not convenient to use for learning
- In need of *a single compact descriptor* for each shape



Multi-View CNN



Multi-View CNN

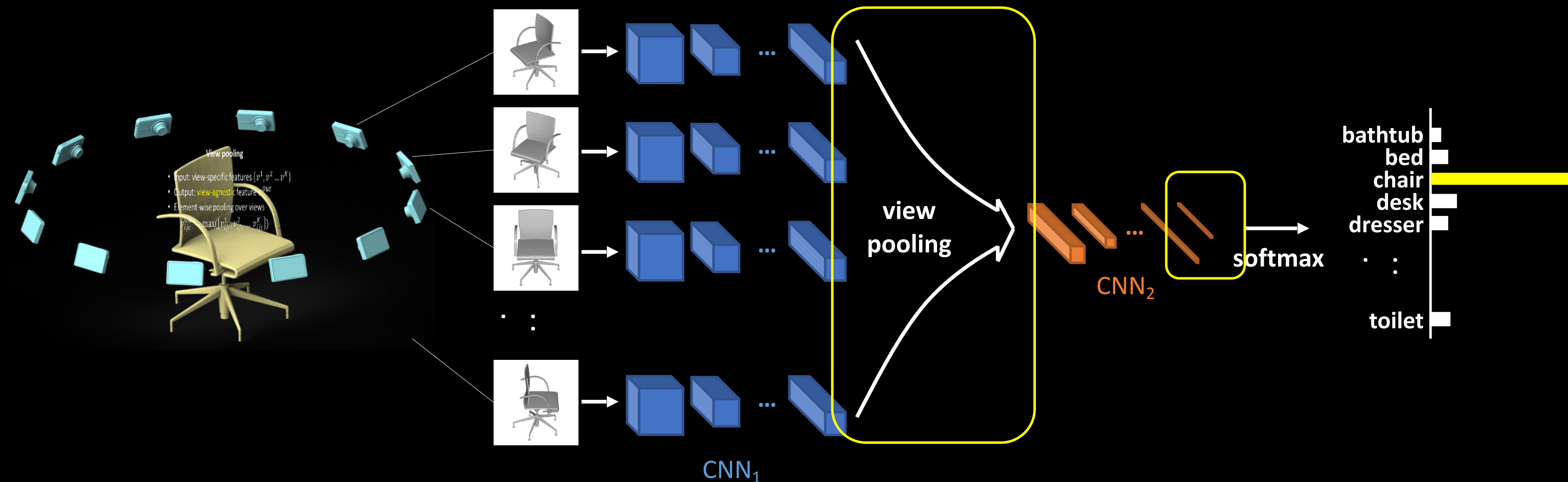
View pooling

- Input: view-specific features $\{v^1, v^2 \dots v^K\}$
- Output: **view-agnostic** feature v^{Out}
- Element-wise pooling over views

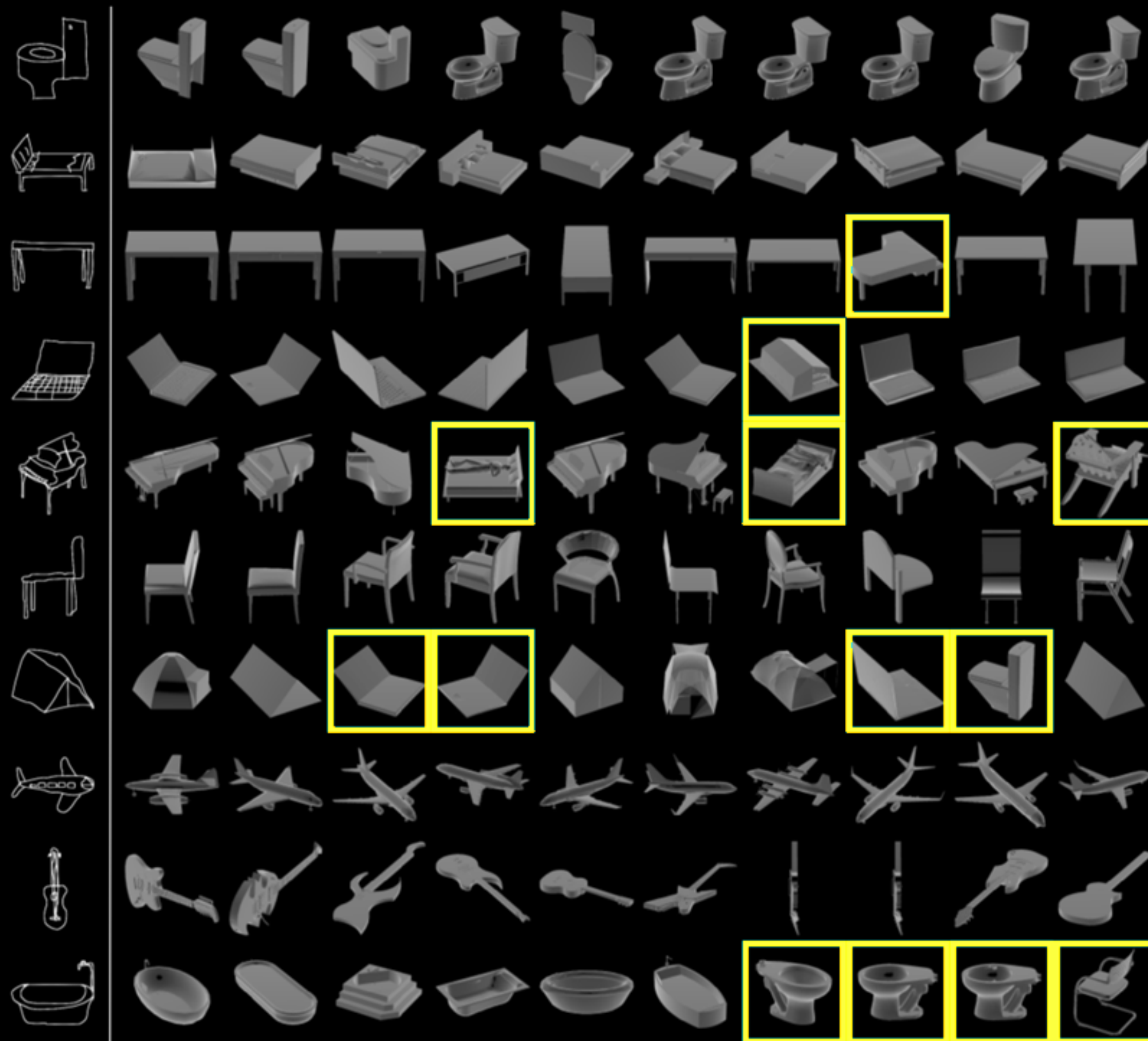
$$v_{ijc}^{Out} = \max(\{v_{ijc}^1, v_{ijc}^2 \dots v_{ijc}^K\})$$

More compact shape descriptor

- Learn a low-rank metric $W \in \mathbb{R}^{128 \times 4096}$
- 4096-dim (fc₇ layer) \rightarrow 128-dim
- Directly applicable to **retrieval**



Relating shapes to images



3D mesh



depth
buffer



Canny edge
detection



human
sketch

Multi-View CNN: Performance

State-of-the-art approach at the time of publishing:

- 75.5 (LFD) → 84.8 (FV) → 90.1 (MVCNN)
- **Ranked #1** in SHREC '16 ShapeNet Retrieval Contest

Still a very competitive approach with some recent updates

- Multi-Resolution [1]
- Improved shading [2]
- Stronger backbone [2]

	Input	Per Class Acc. (%)	Per Inst. Acc. (%)
Our MVCNN	Images	92.4	95.0
Rotation Net		-	94.8
Dominant Set Clustering		-	93.8
MVCNN-MultiRes		91.4	93.8
MVCNN		90.1	90.1
DynamicGraph	Point Clouds	90.2	92.2
Kd-Networks		-	91.8
PointNet++		-	90.7
PointNet		86.2	89.2
VRN Single		-	91.3
O-CNN	Voxels	-	90.6
VoxNet		-	83.0
3DShapeNets		77.3	84.7
PointNet++		-	91.9
FusionNet	Voxels+Images	-	90.8

Table from [2]

[1] C. Qi et al. Volumetric and Multi-View CNNs for Object Classification on 3D Data CVPR '16.

[2] J. Su et al. A Deeper Look at 3D Shape Classifiers. arXiv:1809.02560.

Multi-View CNN: Summary

- A collection of 2D views is highly informative for recognizing 3D objects.
- Using 2D CNNs allows leveraging powerful image architectures and available supervision.
- MVCNNs relate 3D shapes to 2D images, enabling cross-domain applications.



Volumetric CNN

Can we use CNNs without 2D-3D projection?

Yes, just use 3D convolutions!

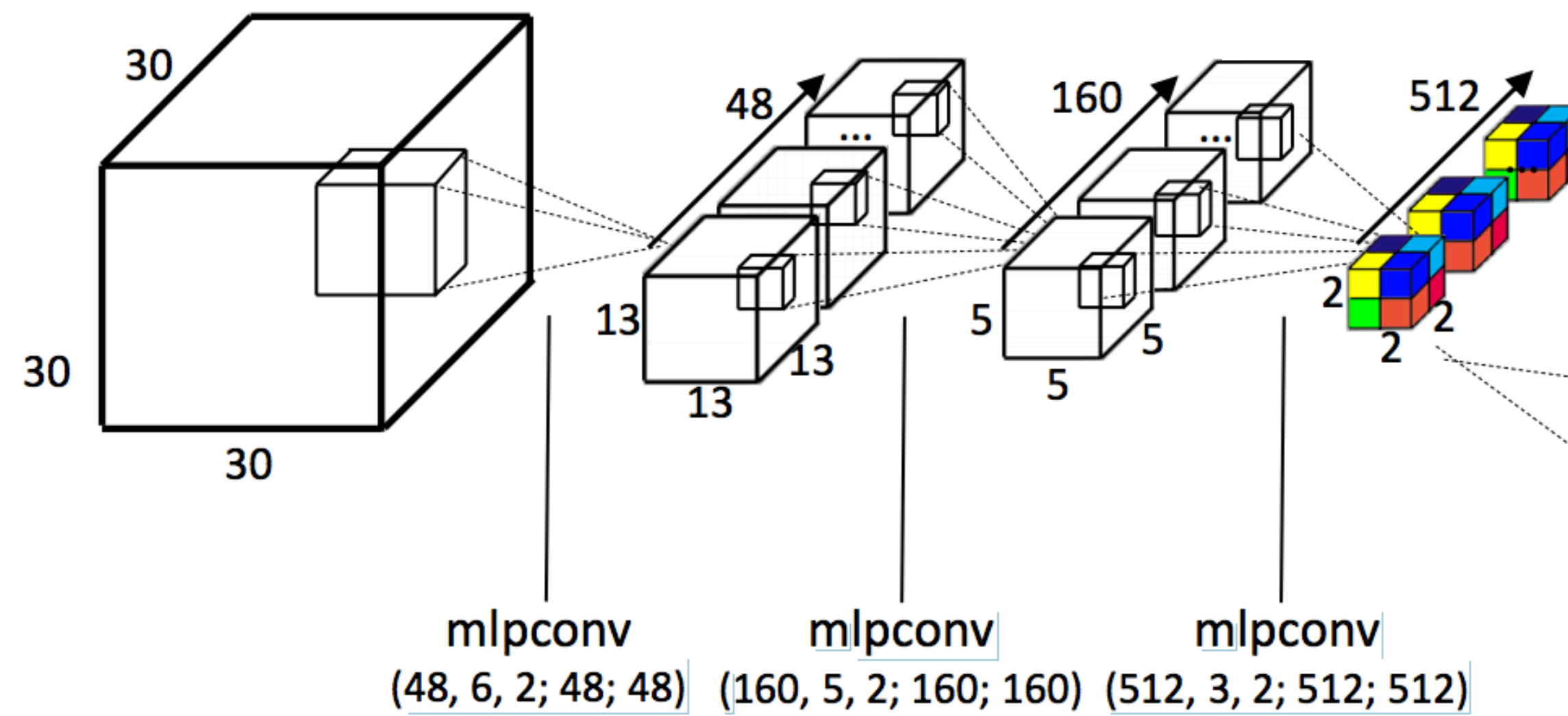
Voxelization

Represent the occupancy of regular 3D grids

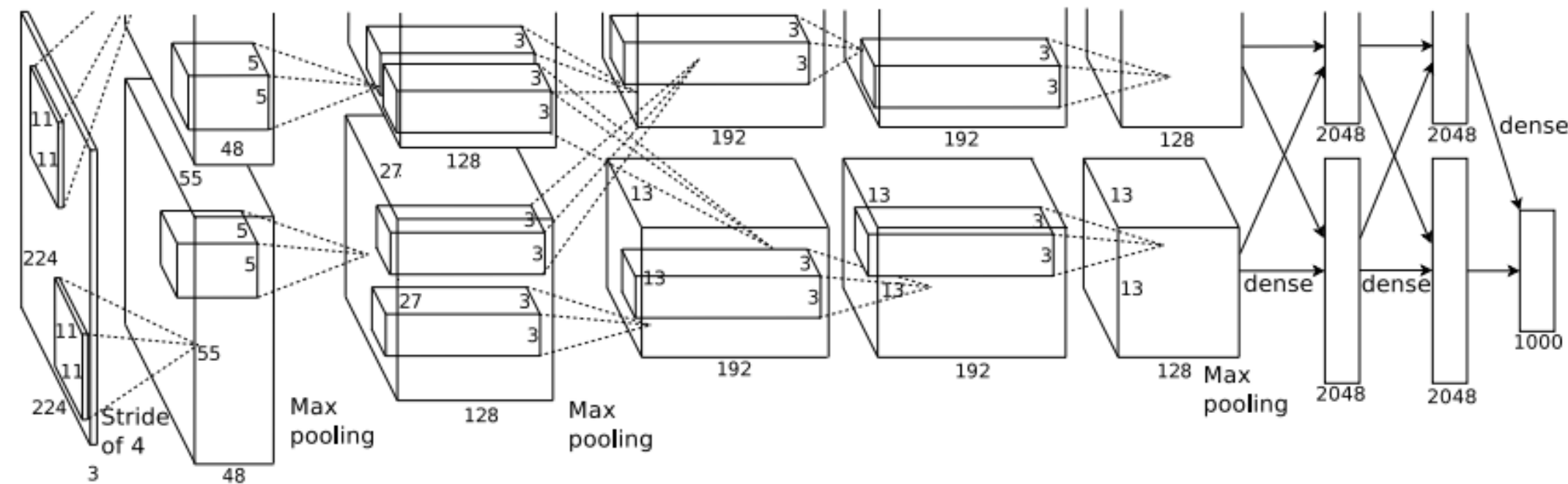


3D CNN on Volumetric Data

3D convolution uses 4D kernels



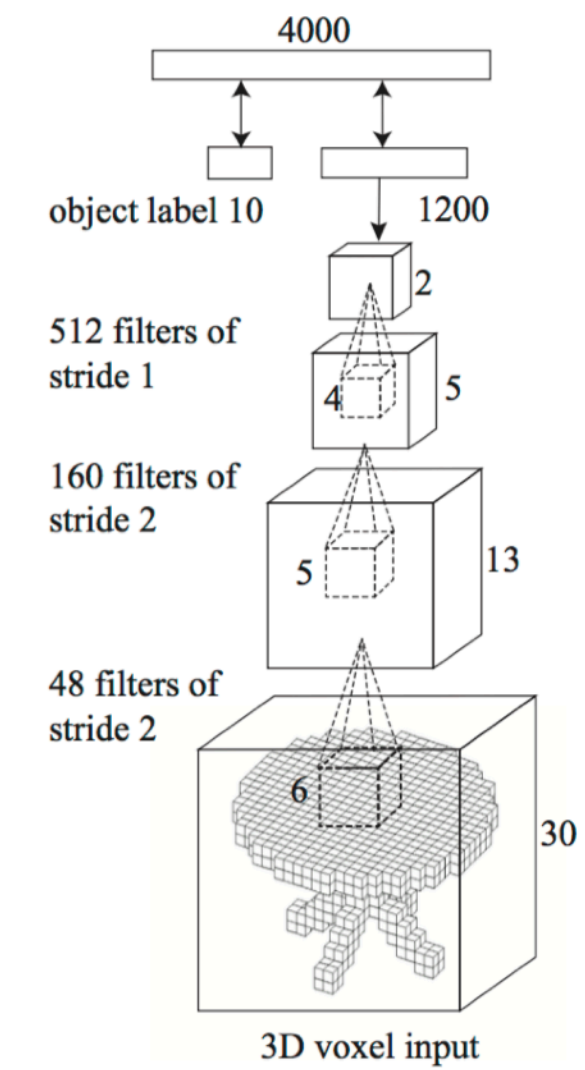
Complexity Issue



AlexNet, 2012

Input resolution: 224x224

$$224 \times 224 = 50176$$

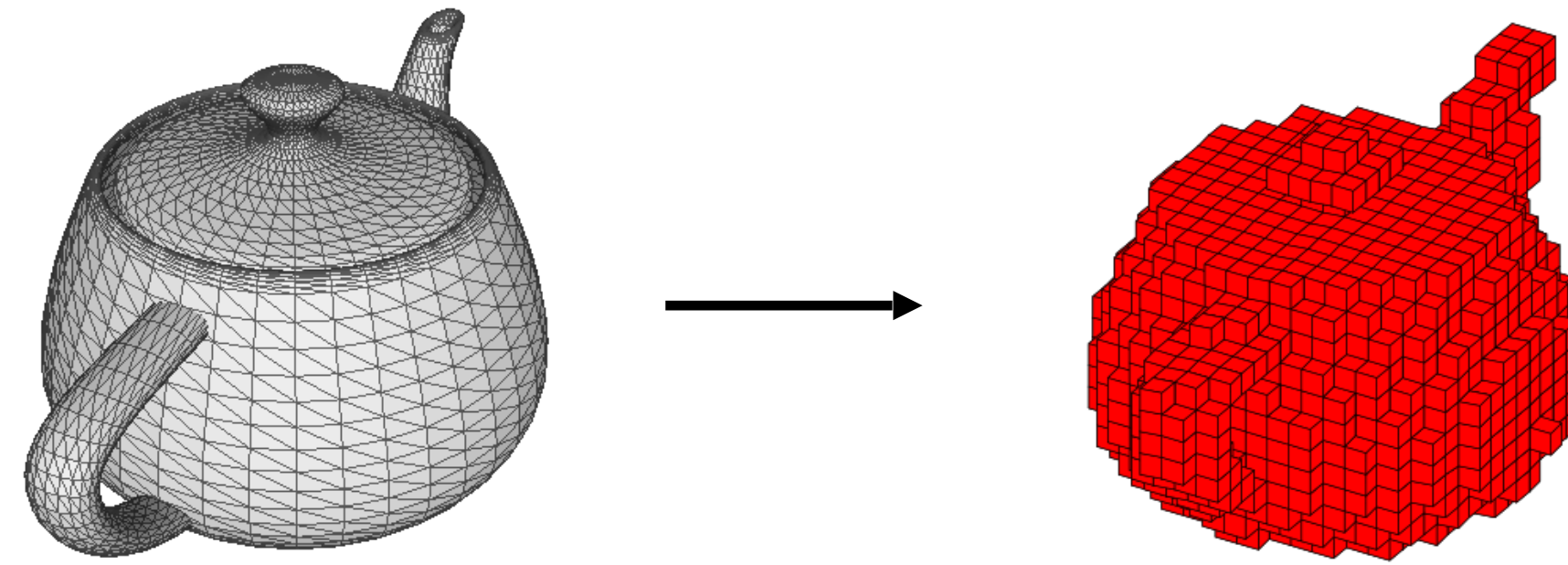


3DShapeNets,
2015

Input resolution: 30x30x30

$$224 \times 224 = 27000$$

Complexity Issue

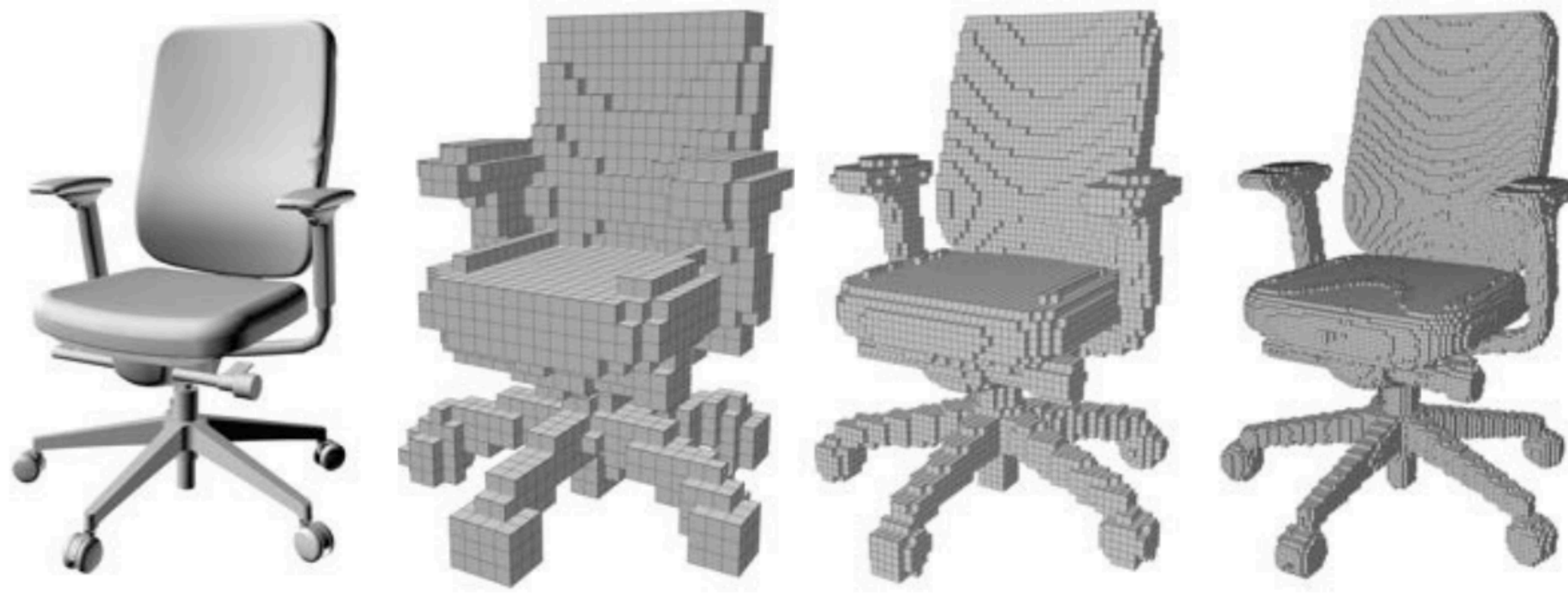


Polygon Mesh

Occupancy Grid
30x30x30

Information loss in voxelization

More Principled: Leverage Sparsity of 3D Surface Data



$$\frac{\#occupied\ grid}{\#total\ grid}$$

Occupancy: 10.41%

5.09%

2.41%

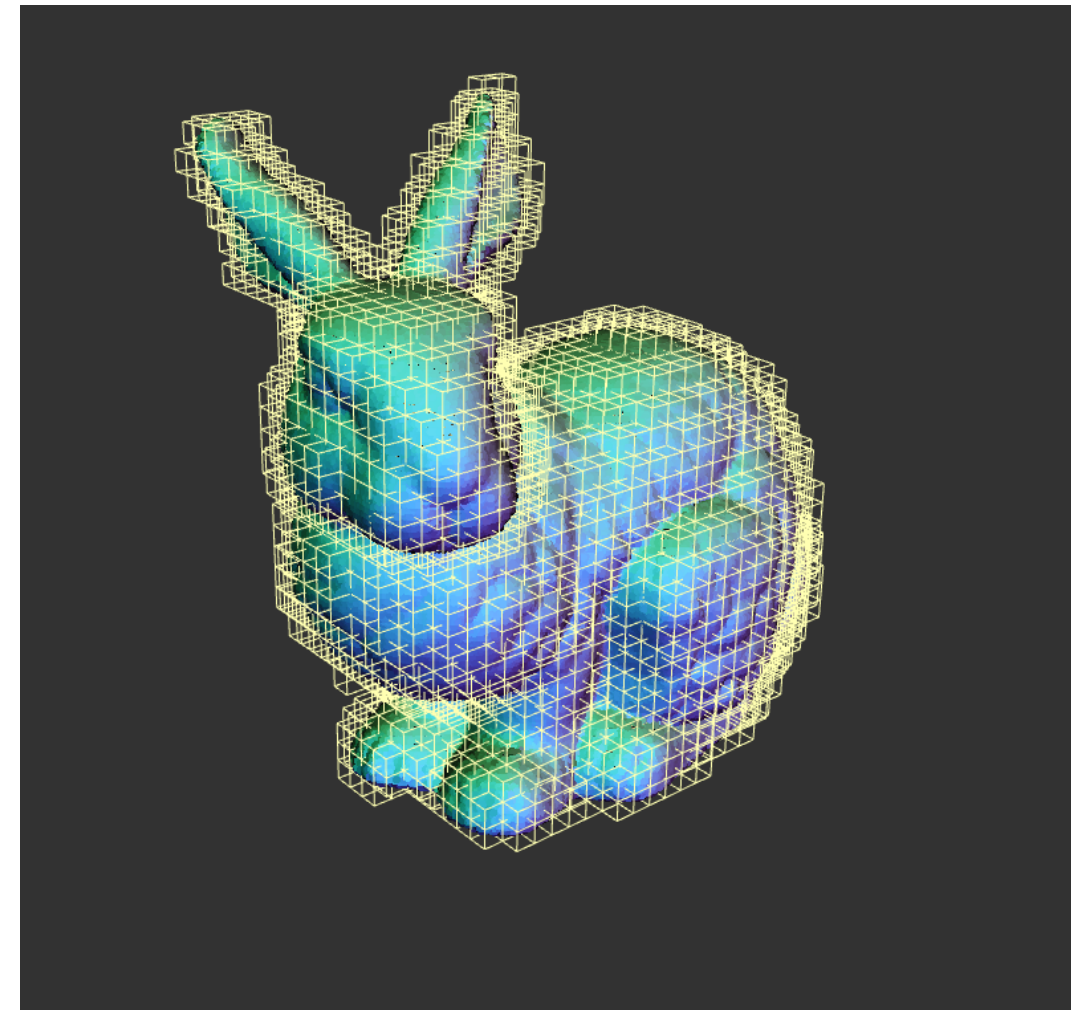
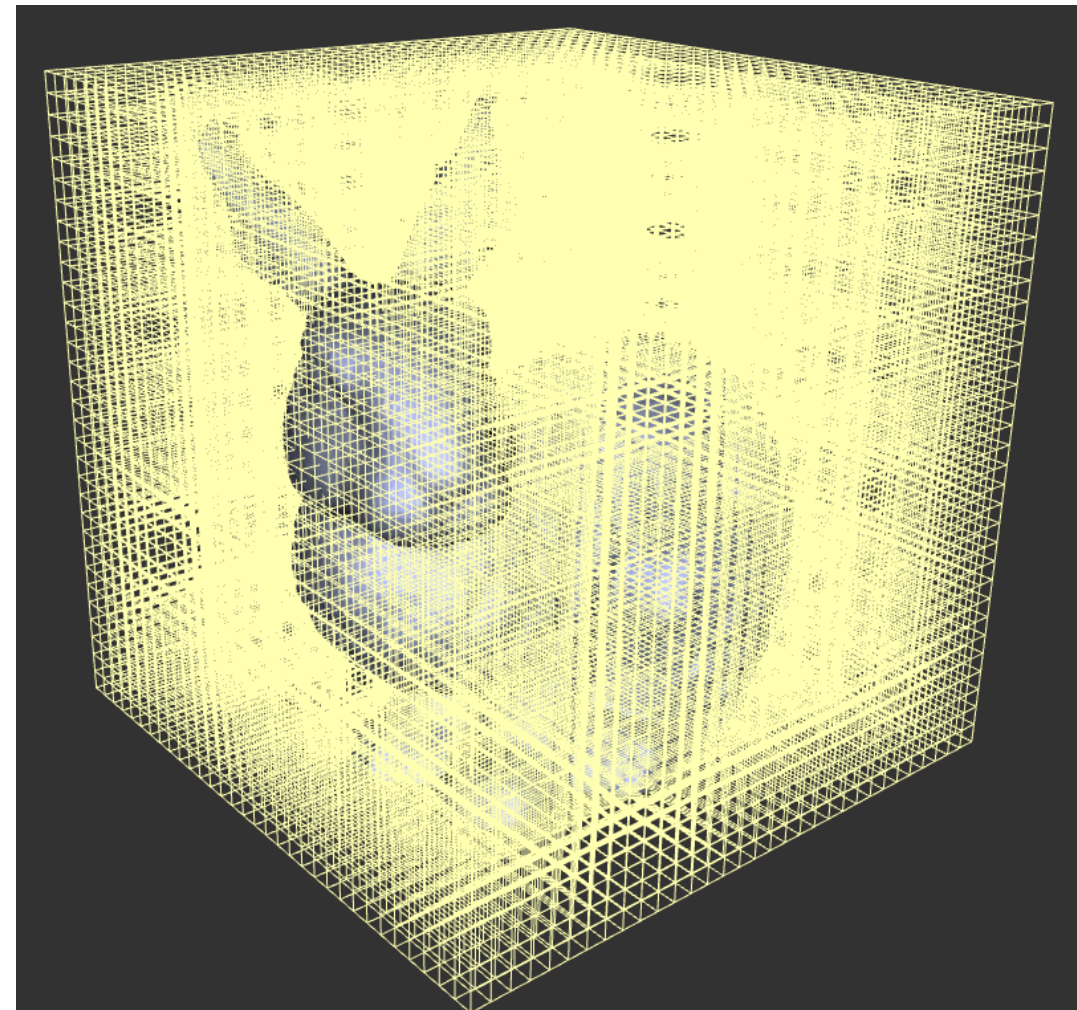
Resolution: 32

64

128

Store only the Occupied Grids

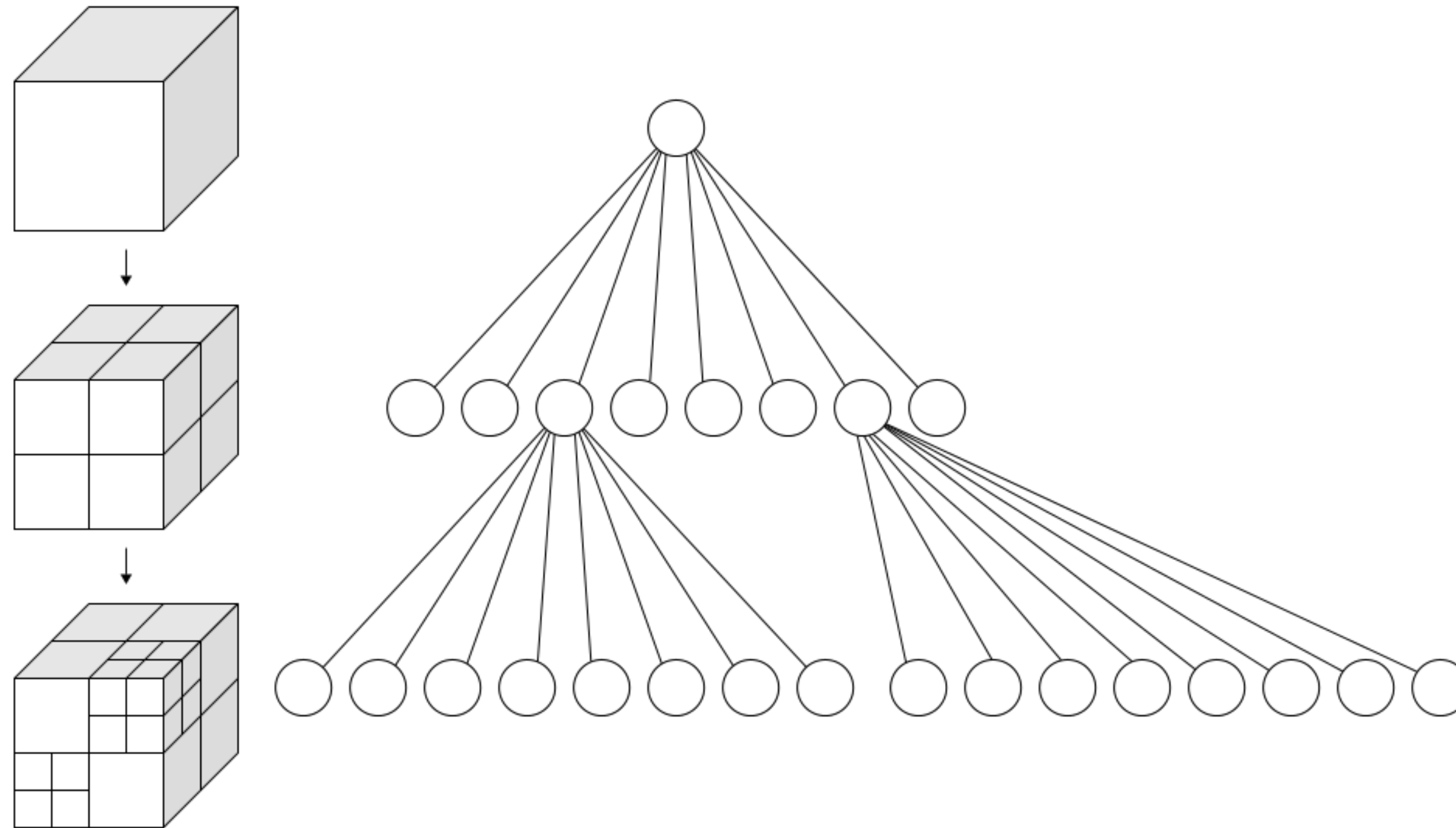
- Store the sparse surface signals
- Constrain the computation near the surface



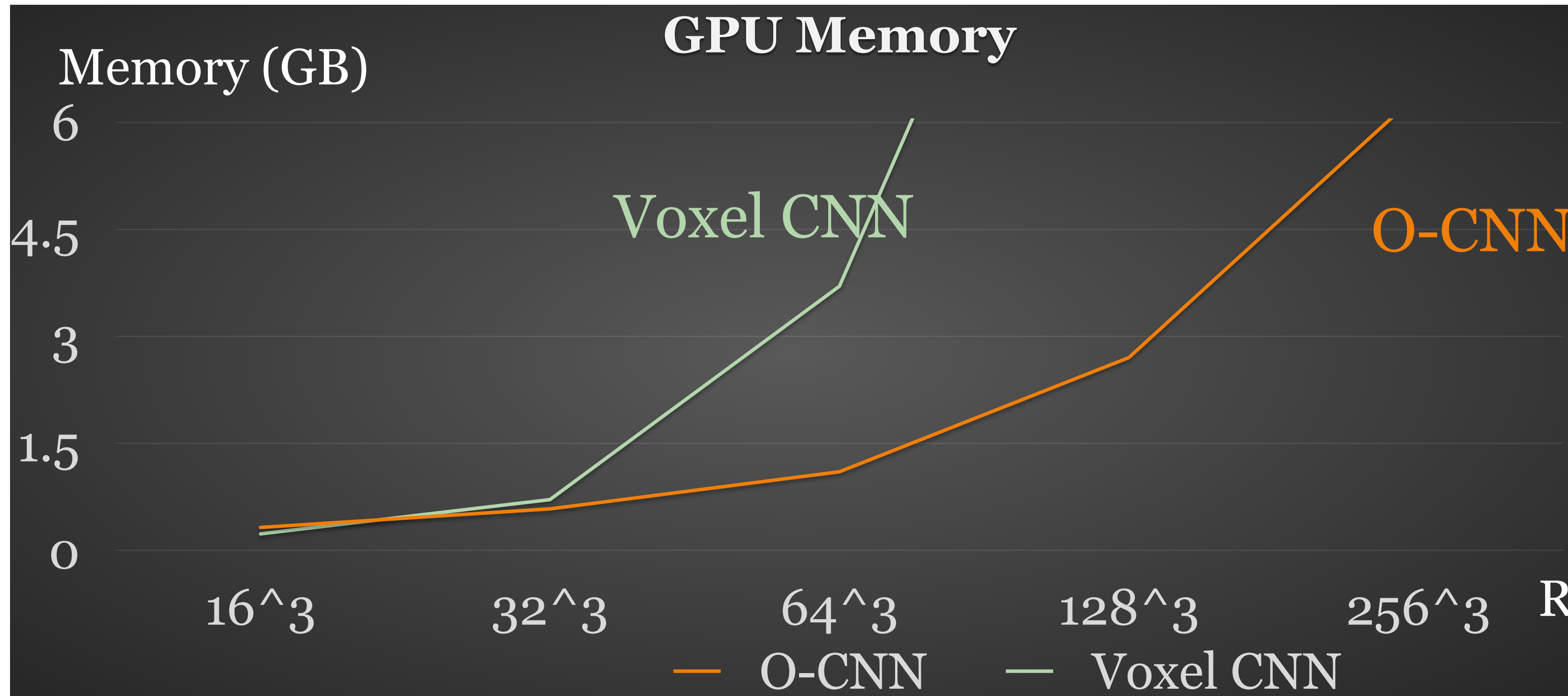
Octree: Recursively Partition the Space

Each **internal node** has exactly eight **children**

Neighborhood searching: Hash table

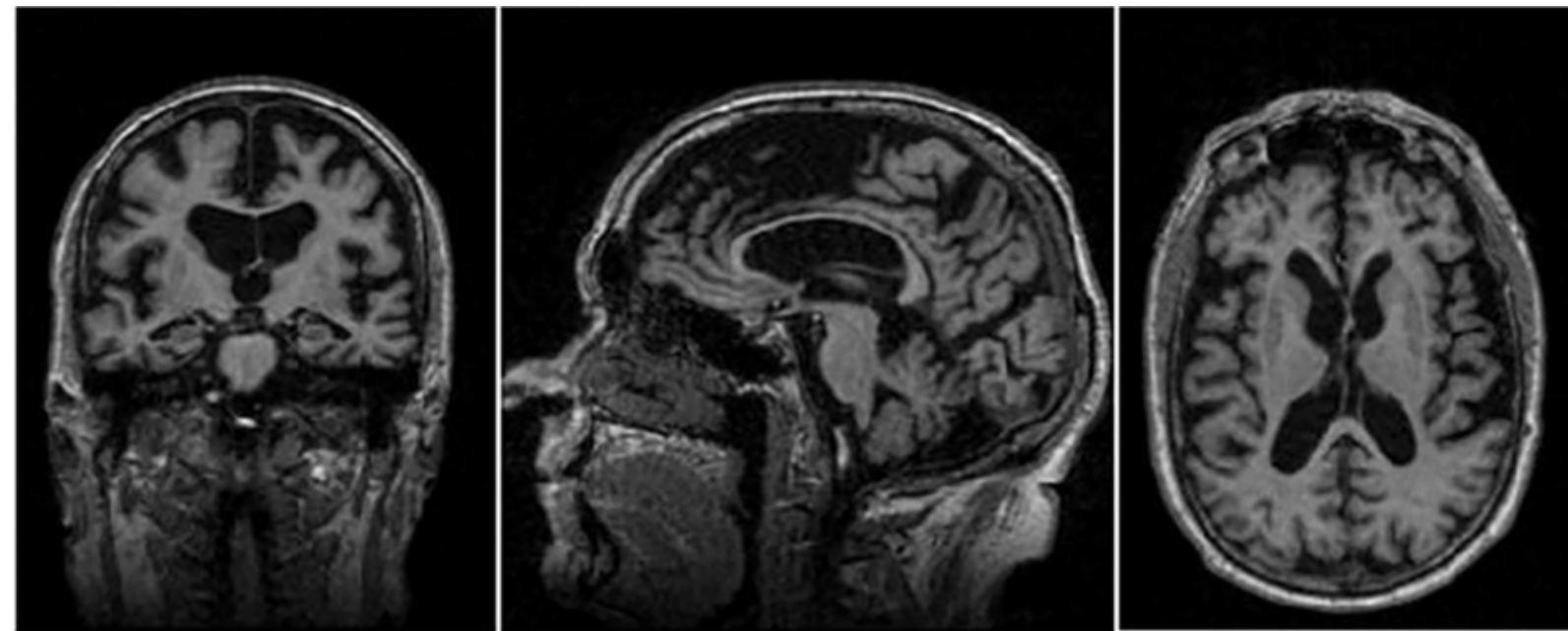


Memory Efficiency

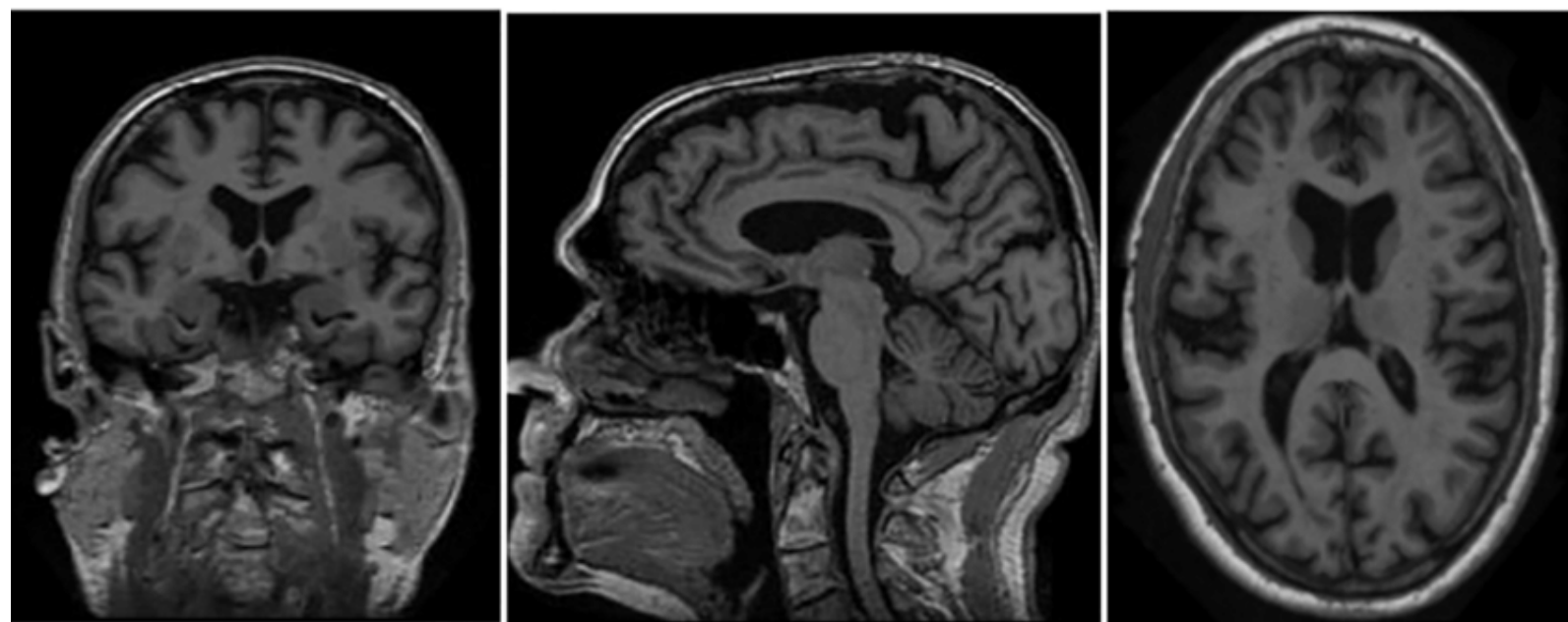


Sometimes the data is volumetric

CT Scan



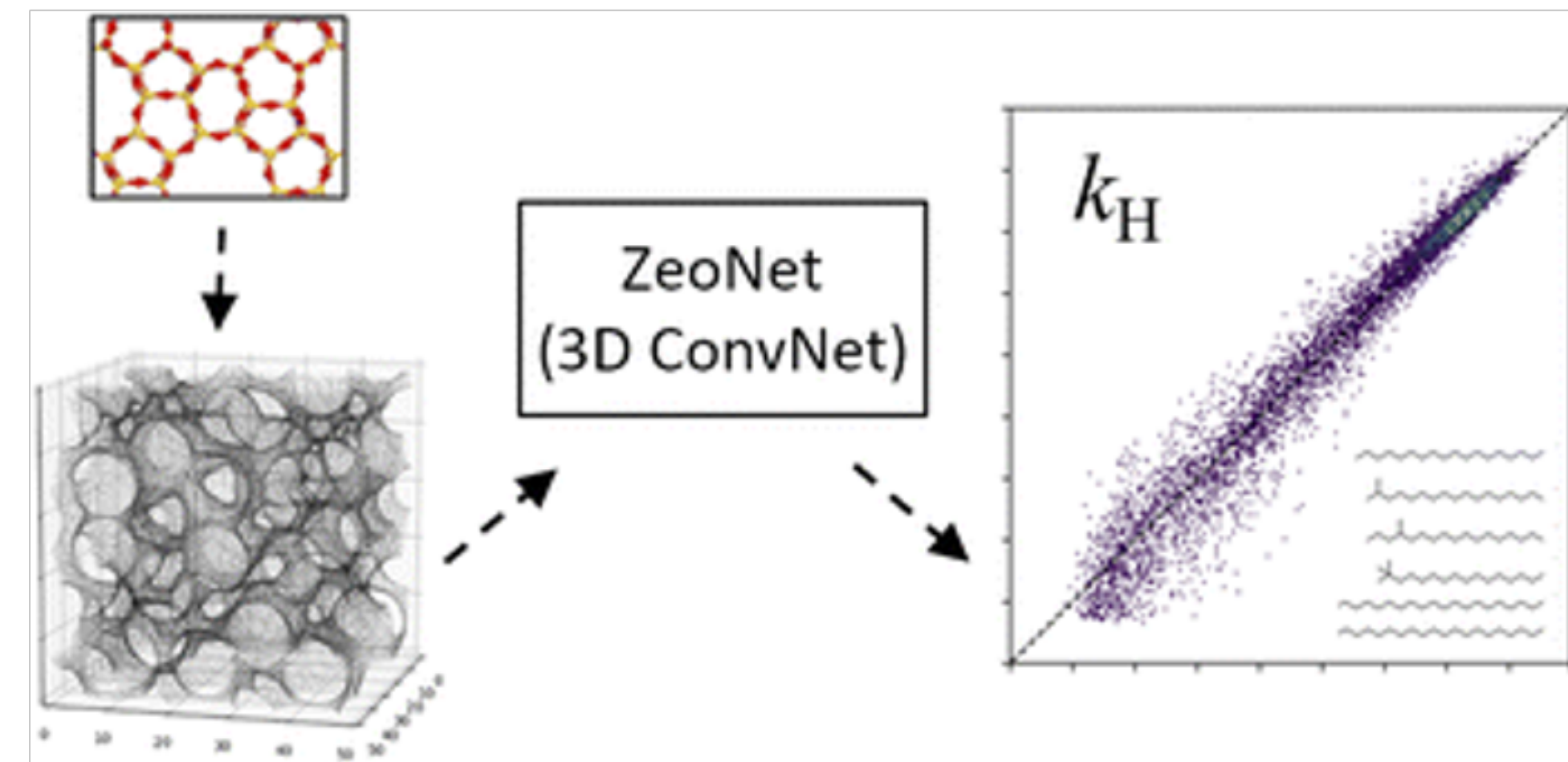
(a)



(b)

<https://ic3.center.ufl.edu/news-events/medical-imaging/medical-imaging-datasets/>

3D Materials (Zeolites)



Journal of
Materials Chemistry A



PAPER

[View Article Online](#)
[View Journal](#) | [View Issue](#)



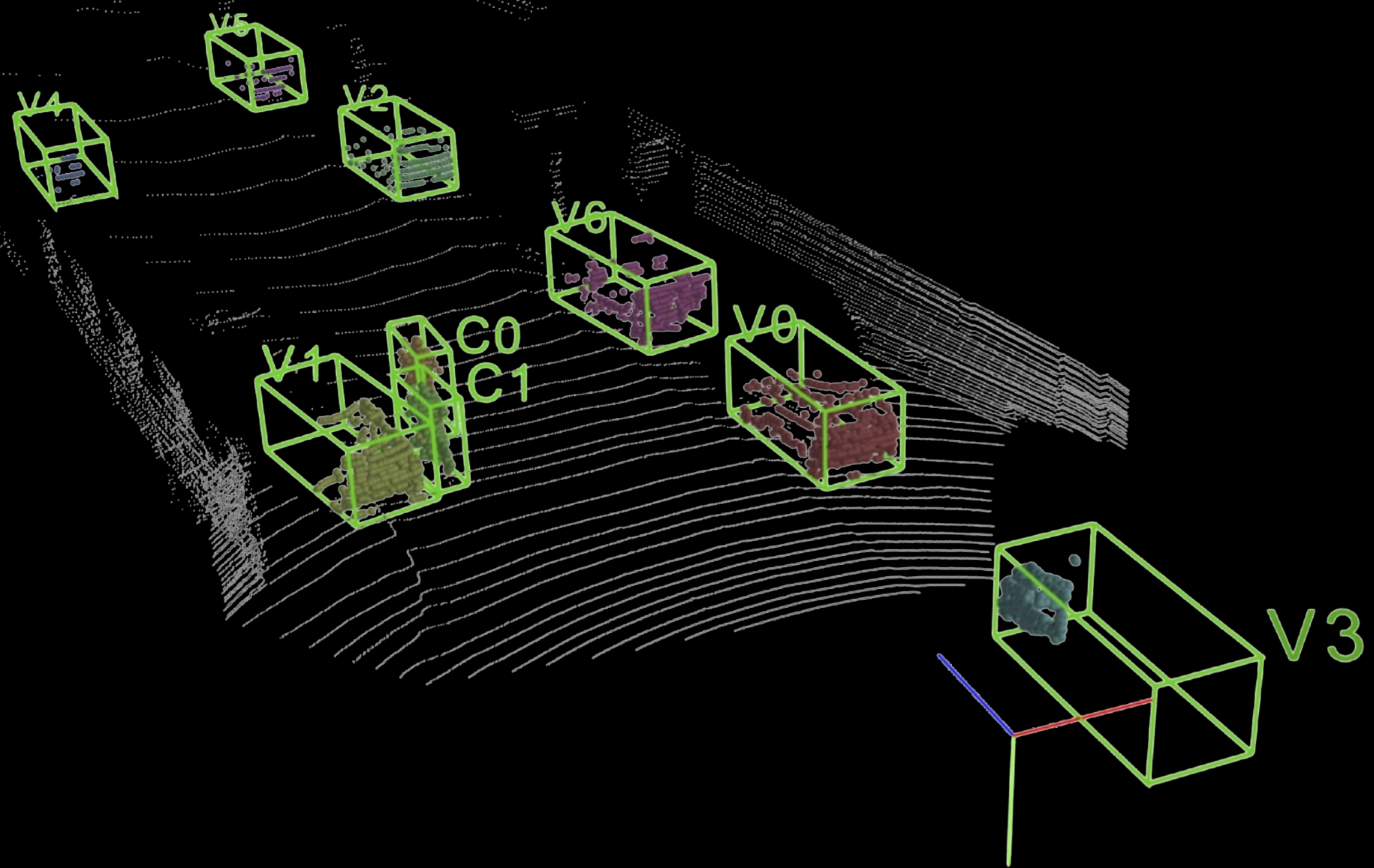
Cite this: *J. Mater. Chem. A*, 2023, 11, 17570

ZeoNet: 3D convolutional neural networks for predicting adsorption in nanoporous zeolites†

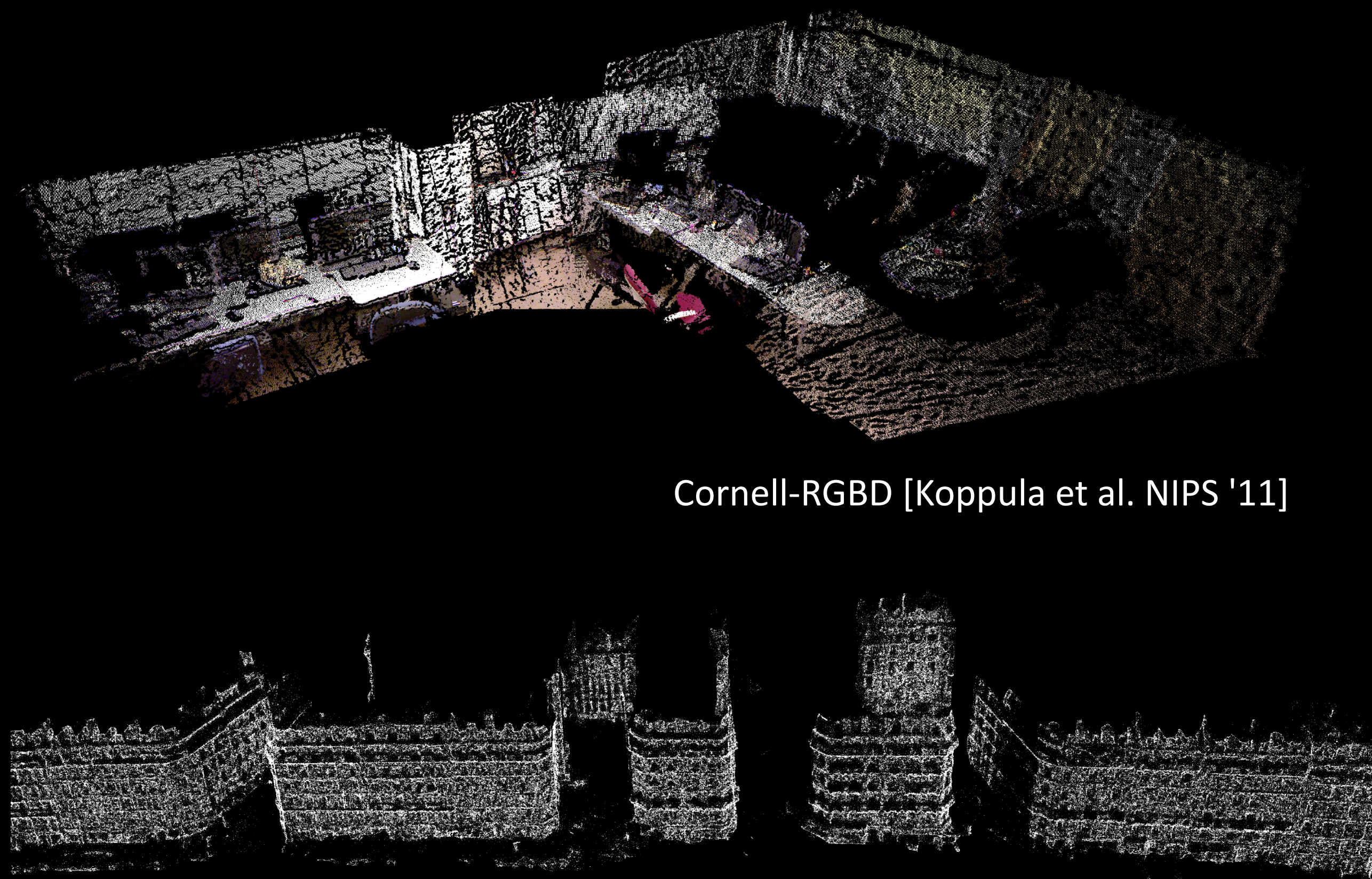
Yachan Liu,^{†a} Gustavo Perez,^{id†b} Zezhou Cheng,^b Aaron Sun,^b Samuel C. Hoover,^a Wei Fan,^{id^a} Subhansu Maji^{*b} and Peng Bai^{id^{*a}}

3D volumetric models are very well suited!

Point Networks



KITTI [Geiger CVPR '12, Qi et al. arXiv '17]



Cornell-RGBD [Koppula et al. NIPS '11]

Ruemonge2014 [Riemenschneider et al. ECCV '14]

Point clouds are highly **sparse**, and **lack of grid structure**.

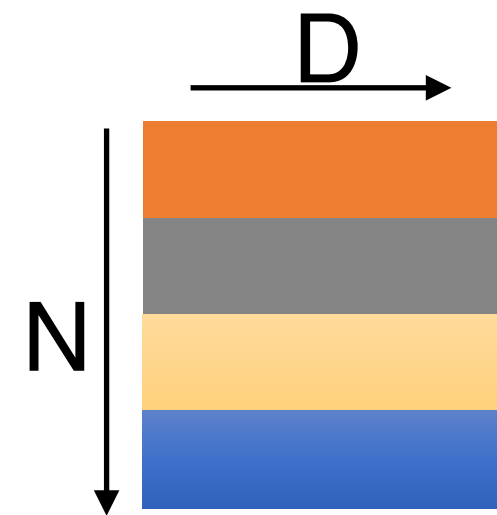
Directly Process Point Cloud Data

End-to-end learning for **unstructured,**
unordered point data



Permutation invariance

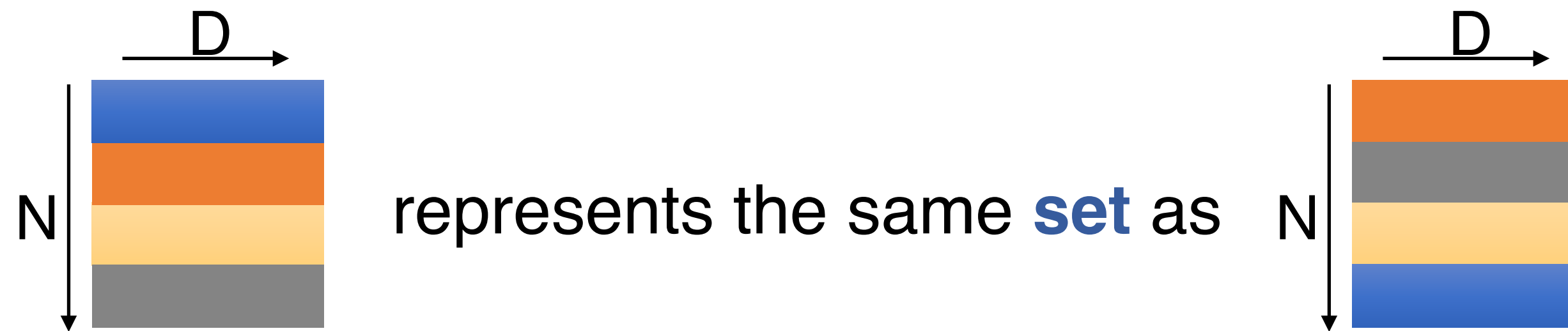
Point cloud: N **orderless** points, each represented by a D dim coordinate



2D array representation

Permutation invariance

Point cloud: N **orderless** points, each represented by a D dim coordinate



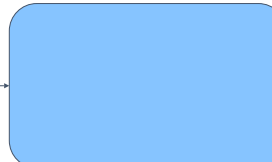
2D array representation


Construct a Symmetric Function


Observe:

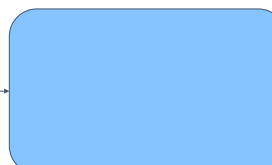
$f(x_1, x_2, \dots, x_n) = \gamma \circ g(h(x_1), \dots, h(x_n))$ is symmetric if g is symmetric

h

(1,2,3) — 

(1,1,1) — 

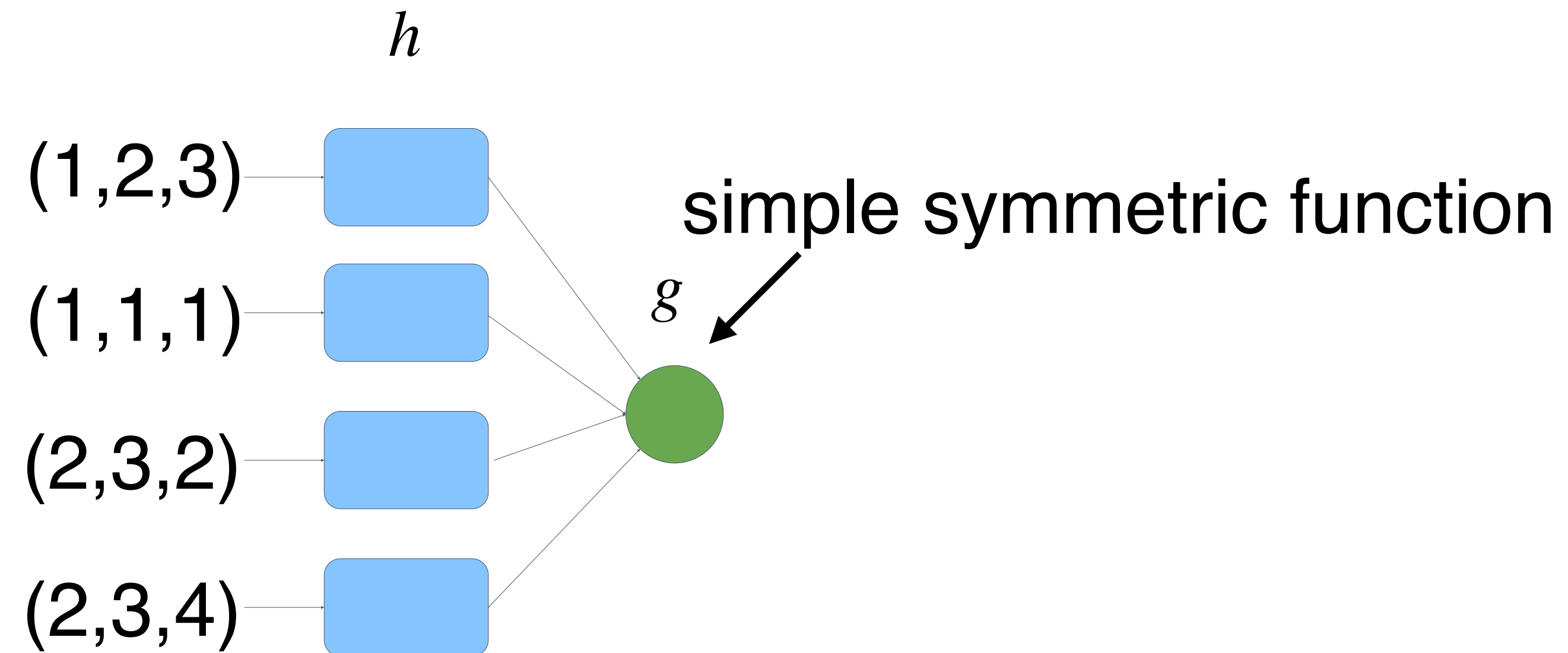
(2,3,2) — 

(2,3,4) — 

Construct a Symmetric Function

Observe:

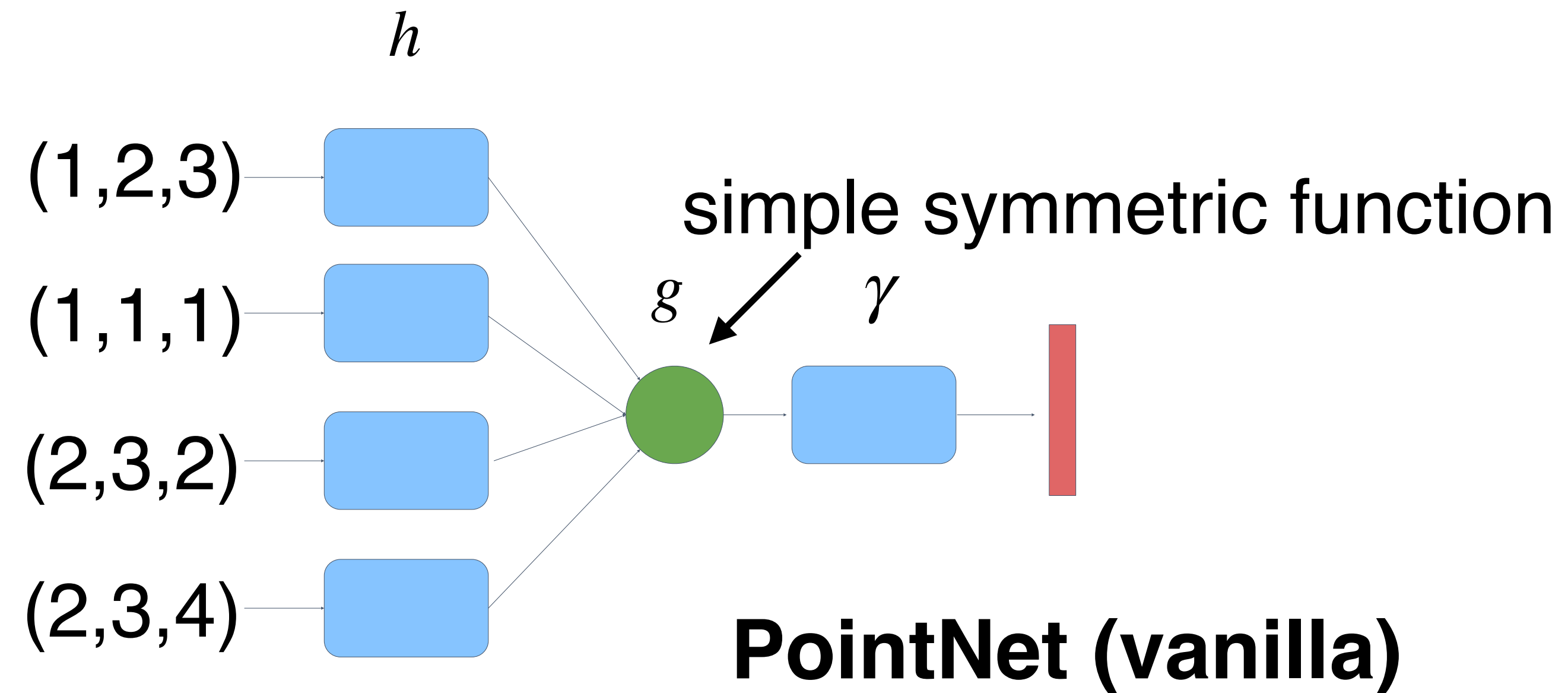
$f(x_1, x_2, \dots, x_n) = \gamma \circ g(h(x_1), \dots, h(x_n))$ is symmetric if g is symmetric



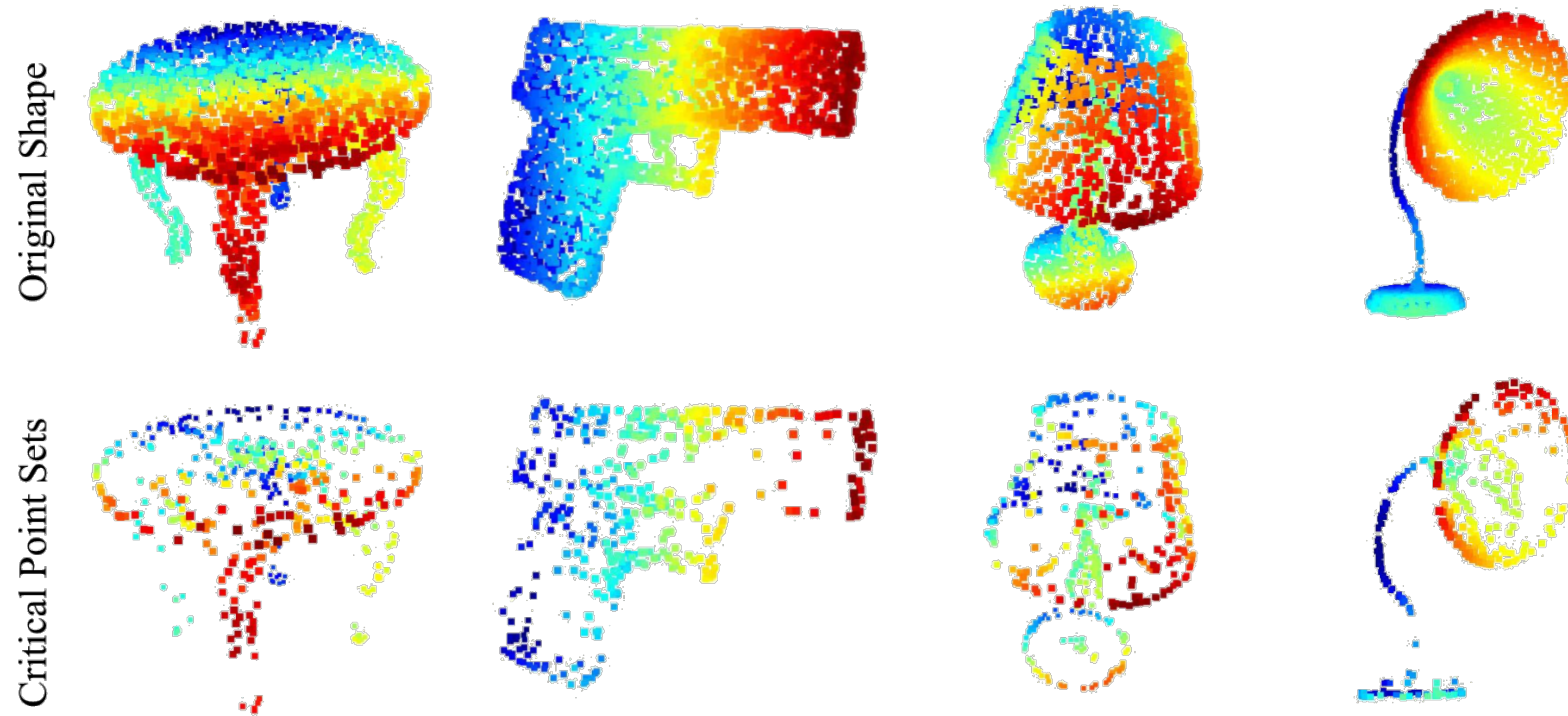
Construct a Symmetric Function

Observe:

$f(x_1, x_2, \dots, x_n) = \gamma \circ g(h(x_1), \dots, h(x_n))$ is symmetric if g is symmetric



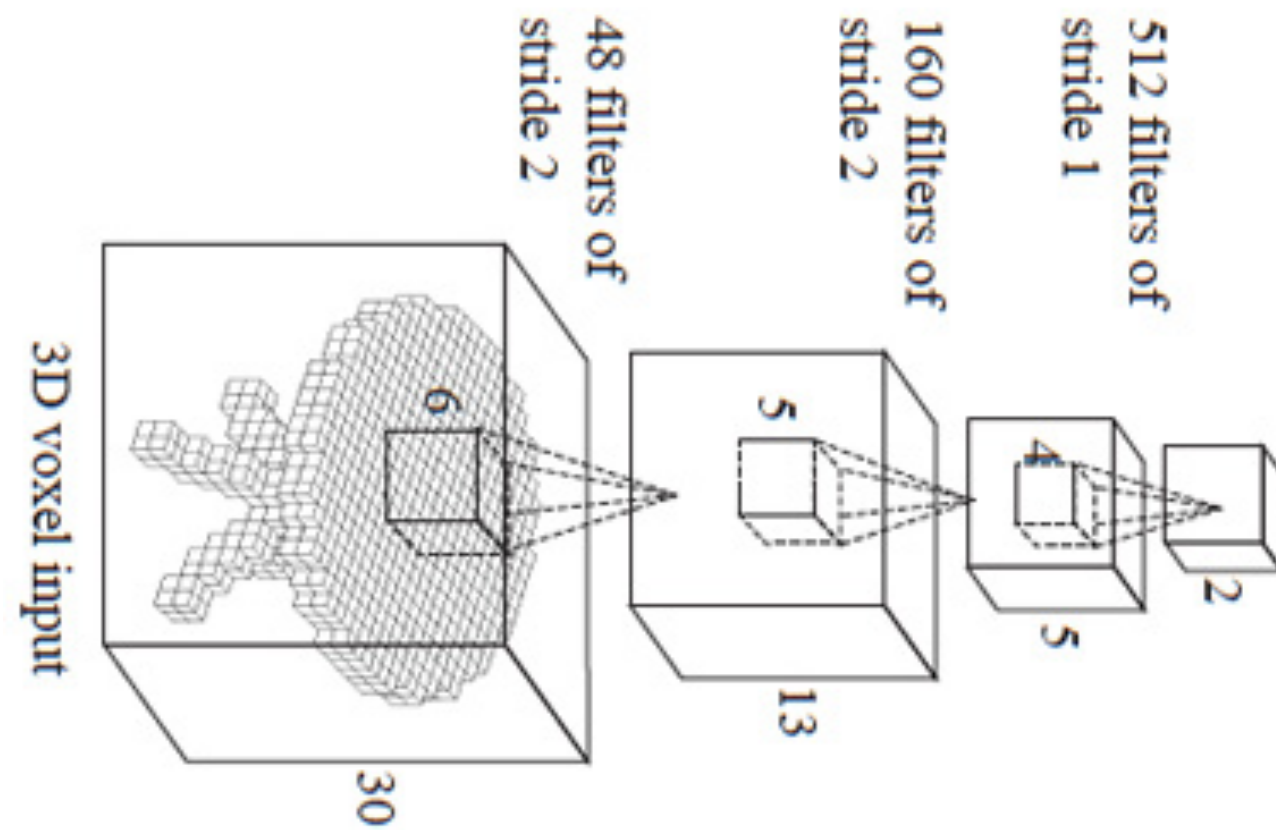
Visualize What is Learned by Reconstruction



Salient points are discovered!

Limitations of PointNet

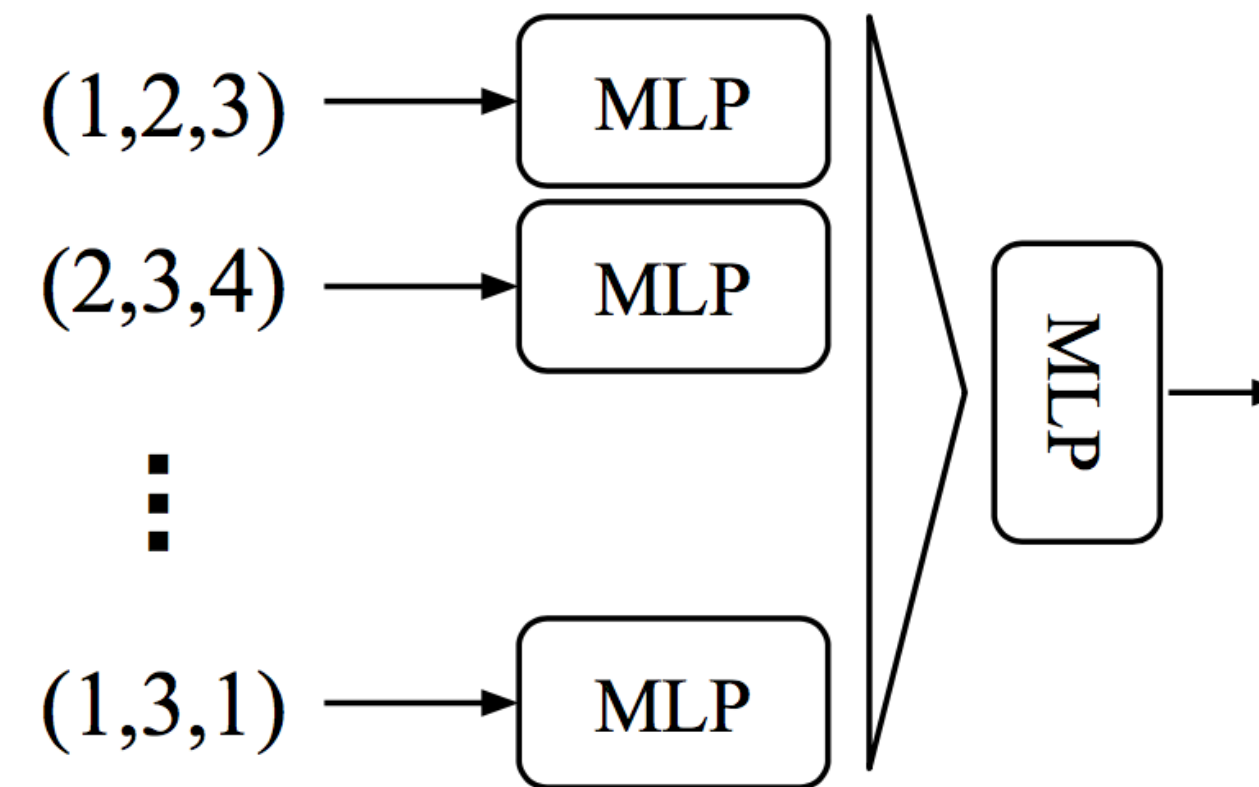
Hierarchical feature learning
Multiple levels of abstraction



3D CNN (Wu et al.)

v.s.

Global feature learning
Either one point or all points



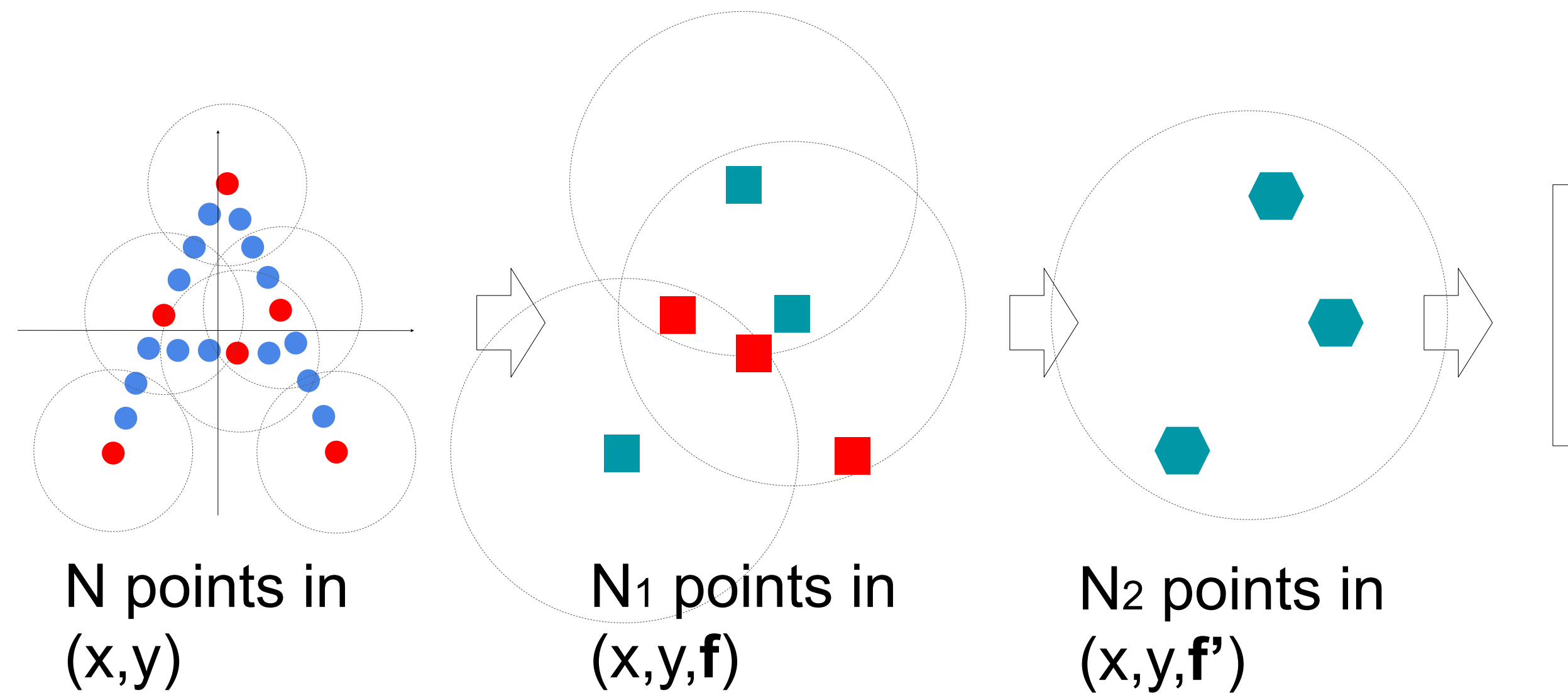
PointNet (vanilla) (Qi et al.)

- No local context for each point!
- Global feature depends on absolute coordinate. Hard to generalize to scene configurations!

Points in Metric Space

- Learn “kernels” in 3D space and conduct convolution
- Kernels have compact spatial support
- For convolution, we need to find neighboring points
- Possible strategies for range query
 - Ball query (results in more stable features)
 - k-NN query (faster)

PointNet v2.0: Multi-Scale PointNet

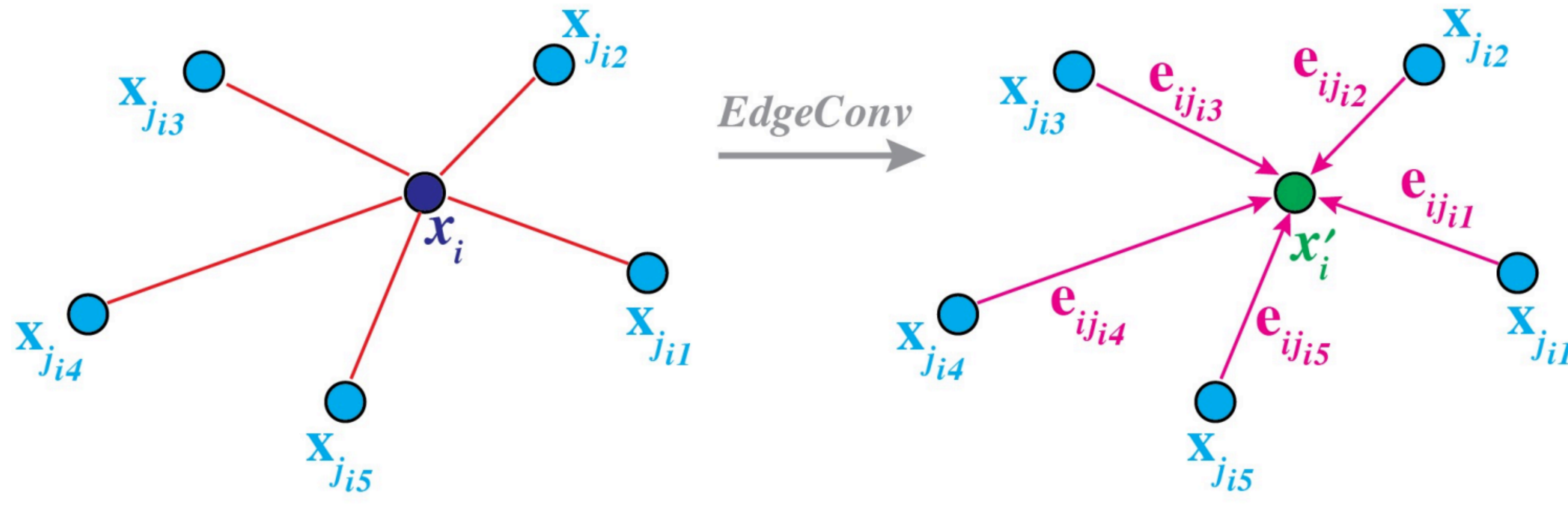


Repeat

- Sample anchor points
- Find neighborhood of anchor points
- Apply PointNet in each neighborhood to mimic convolution

Point Convolution As Graph Convolution

- Points \rightarrow Nodes
- Neighborhood \rightarrow Edges
- Graph CNN for point cloud processing



Wang et al., "Dynamic Graph CNN for Learning on Point Clouds",
Transactions on Graphics, 2019

Liu et al., "Relation-Shape Convolutional Neural Network for Point
Cloud Analysis", *CVPR* 2019

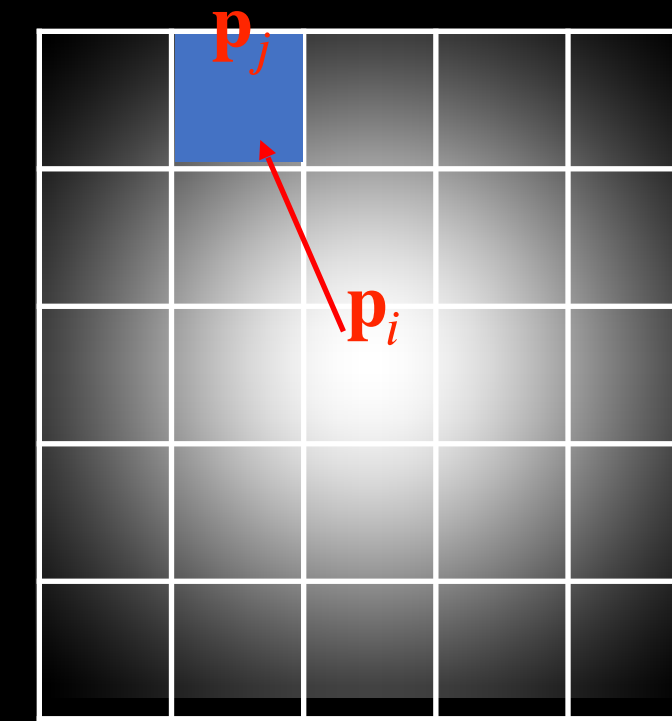
Spatial Convolution

Filter values based on position offsets w.r.t. center

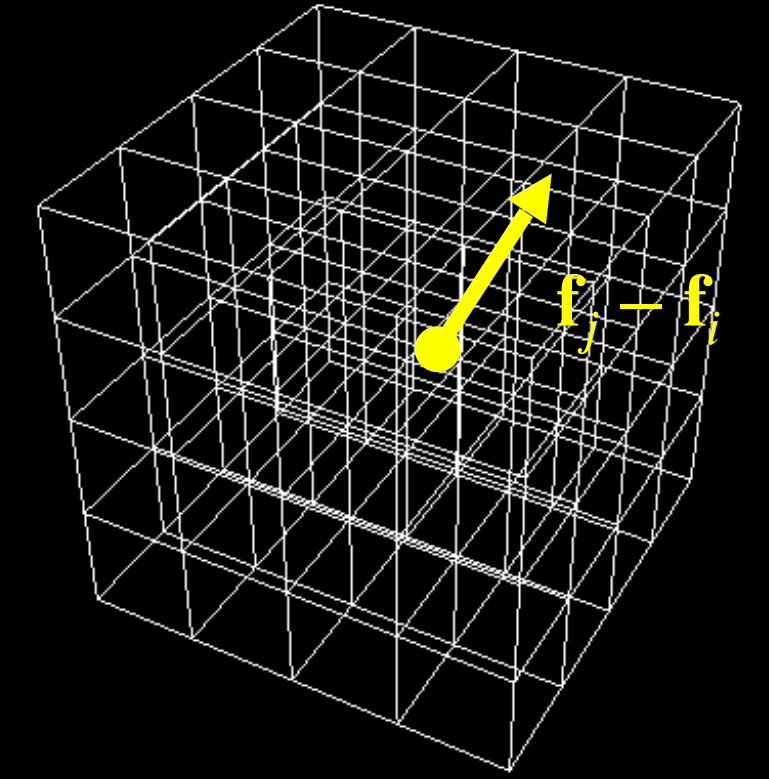
$$v'_i = \sum_{j \in \Omega(i)} w[p_j - p_i] v_j + b$$

↑ output value
 ↑ filter weights
 ↑ position offset
 ↑ input value

$$W[p_j - p_i]$$



$$W[f_j - f_i]$$



$$* \quad \begin{matrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{matrix} \quad =$$

$$W : e^{-\frac{1}{2\sigma^2} \|p_j - p_i\|^2}$$



After Gaussian filtering
(smooths image)

CNNs: w/ learned filters

High-Dimensional Filtering

$$\mathbf{v}'_i = \sum_{j \in \Omega(i)} \mathbf{w}[\mathbf{p}_j - \mathbf{p}_i] \mathbf{v}_j \quad \longrightarrow \quad \mathbf{v}'_i = \sum_{j \in \Omega(i)} \mathbf{w}[\mathbf{f}_j - \mathbf{f}_i] \mathbf{v}_j$$

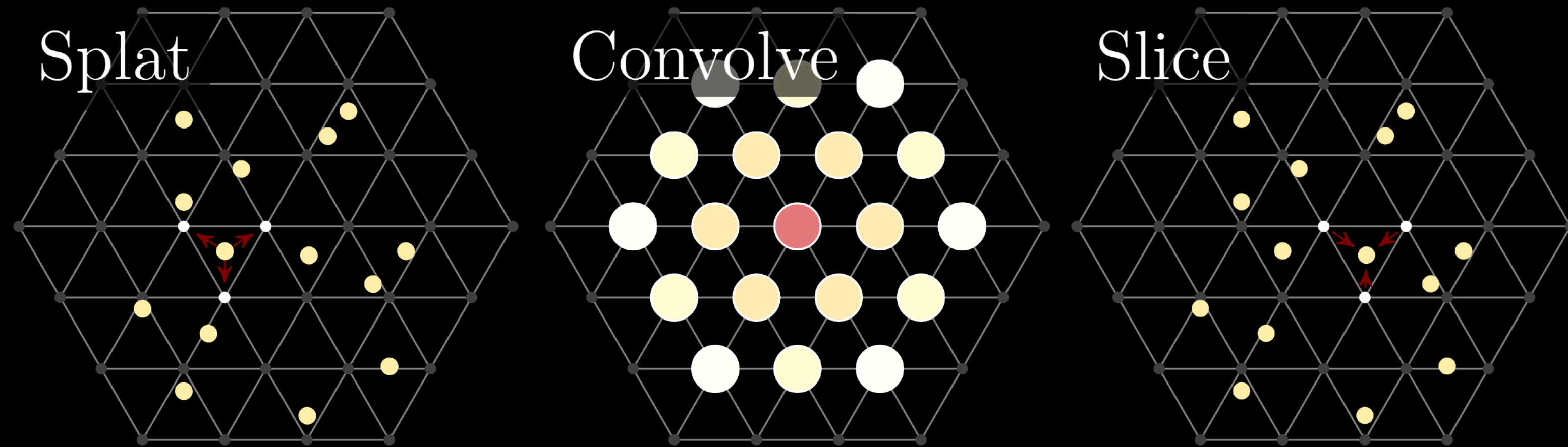
spatial convolution
high-dim. filtering

positions
 ↓ ↓
 input value
 ↙

		Bilateral filter	Conv2d	
Image	pixel value \mathbf{v}	(r, g, b)	(r, g, b)	$\mathbf{v}(\dots)$
	pixel position \mathbf{f}	(x, y)	(x, y, r, g, b)	$(x, y, depth)$
3D Point cloud	point value \mathbf{v}	1	(x, y, z)	(r, g, b)
	point position \mathbf{f}	(x, y, z)	(x, y, z)	(x, y, z, n_x, n_y, n_z)

Efficient Sparse High-Dimensional Filtering

Bilateral Convolution Layer (BCL) [1,2,3]

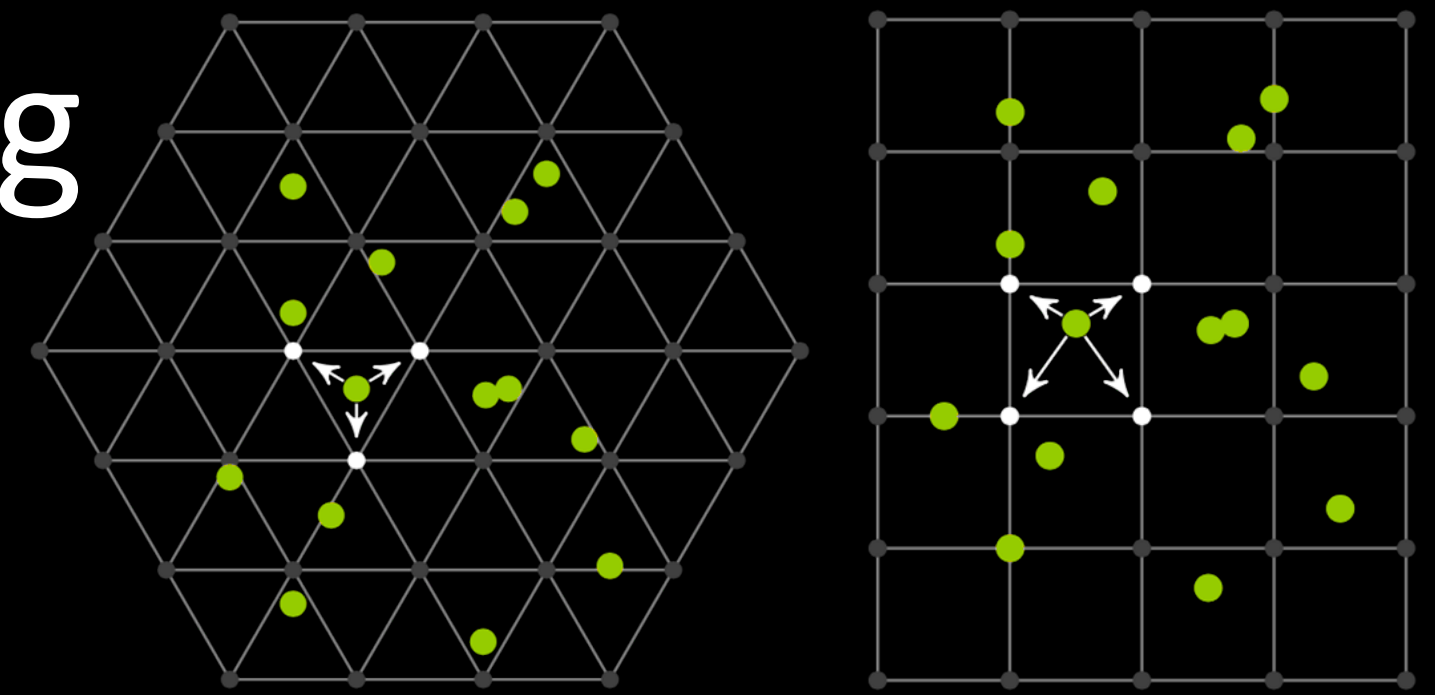


[1] A. Adams, J. Baek and M. A. Davis. Fast High-Dimensional Filtering Using the Permutohedral Lattice. Computer Graphics Forum '10

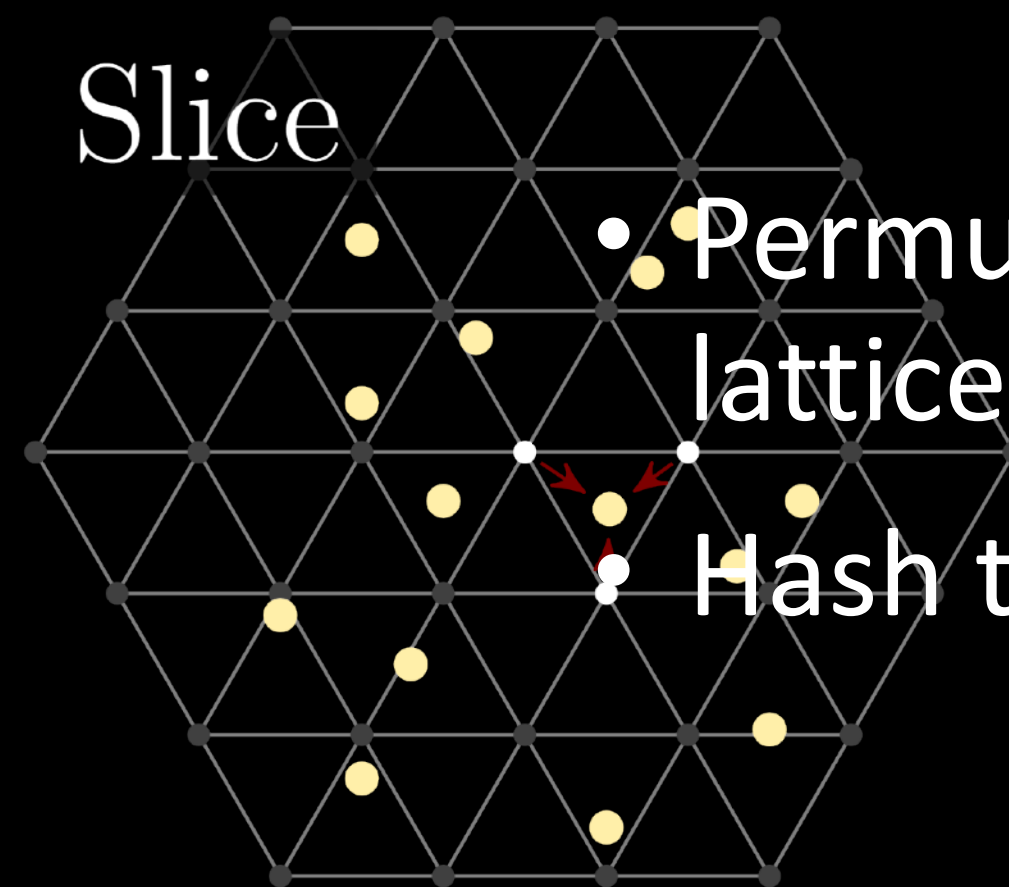
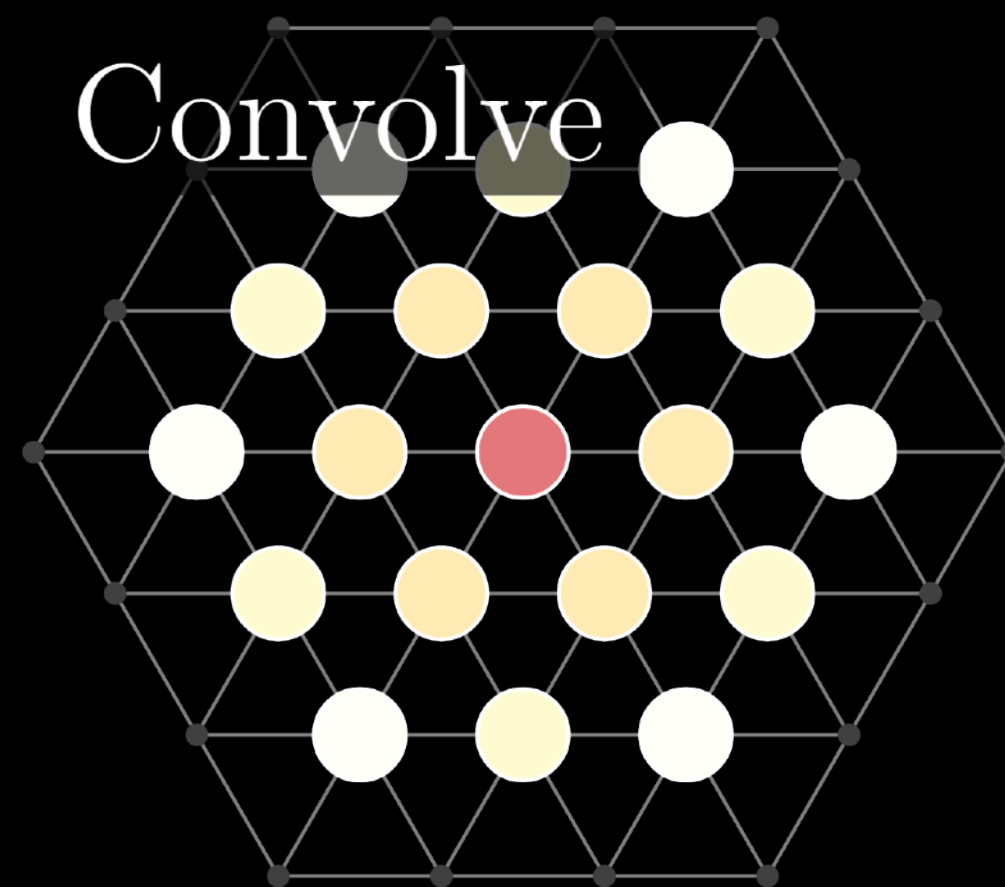
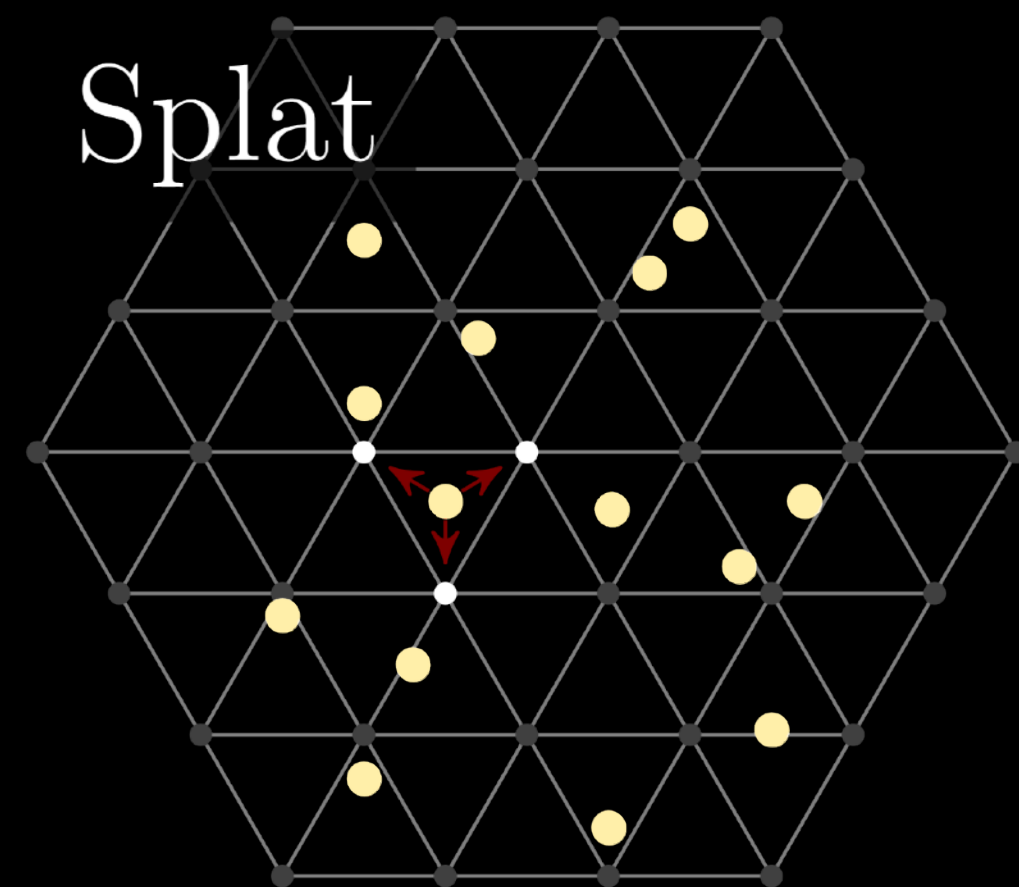
[2] V. Jampani, M. Kiefel and P. V. Gehler. Learning Sparse High Dimensional Filters: Image Filtering, Dense CRFs and Bilateral Neural Networks. CVPR '16 44

[3] M. Kiefel, V. Jampani and P. V. Gehler. Permutohedral Lattice CNNs. ICLR '15 workshops

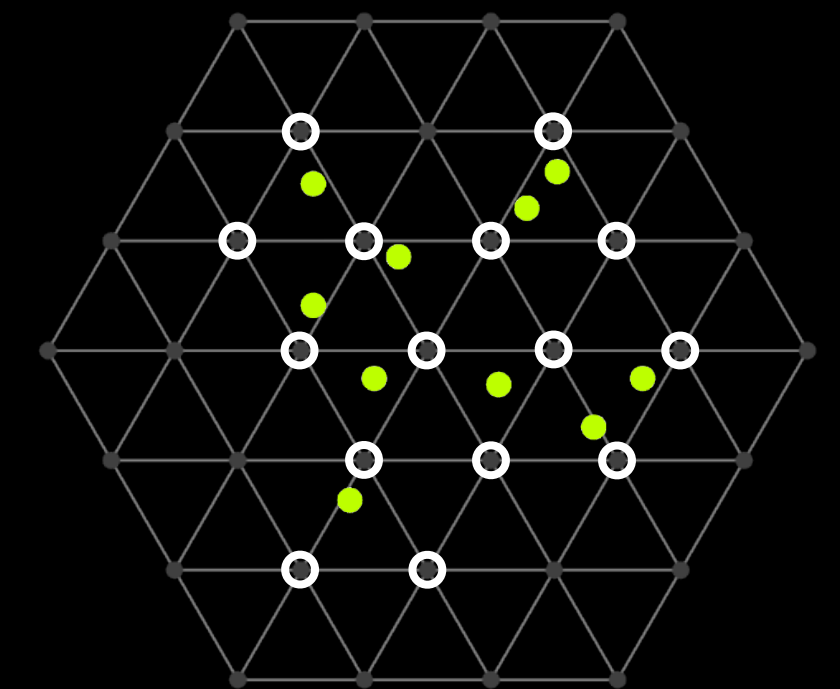
Efficient Sparse High-Dimensional Filtering



Data structures for efficient computation

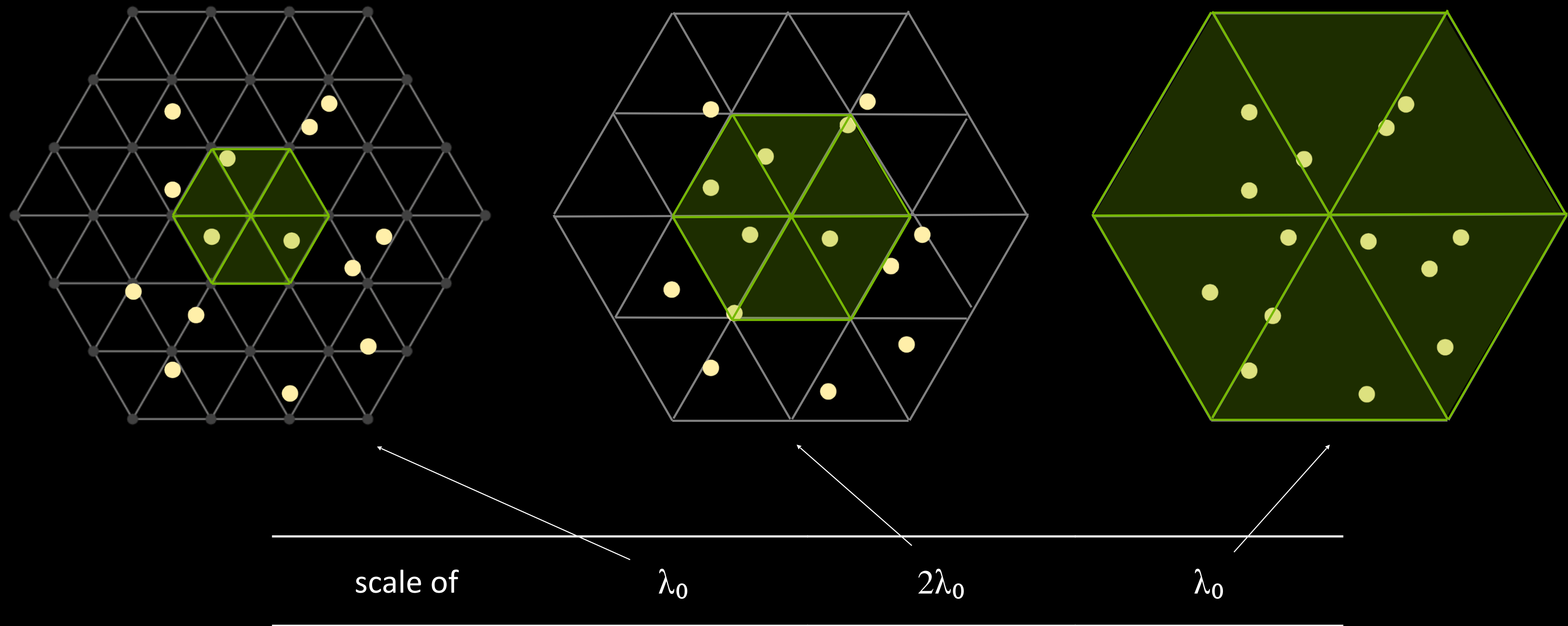


• Permutohedral lattice [1]
• Hash table



Controlling receptive fields

$$\mathbf{v}'_i = \sum_{j \in \Omega(i)} \mathbf{w} [\mathbf{f}_j - \mathbf{f}_i] \mathbf{v}_j$$





lattice feature: (x, y, z)

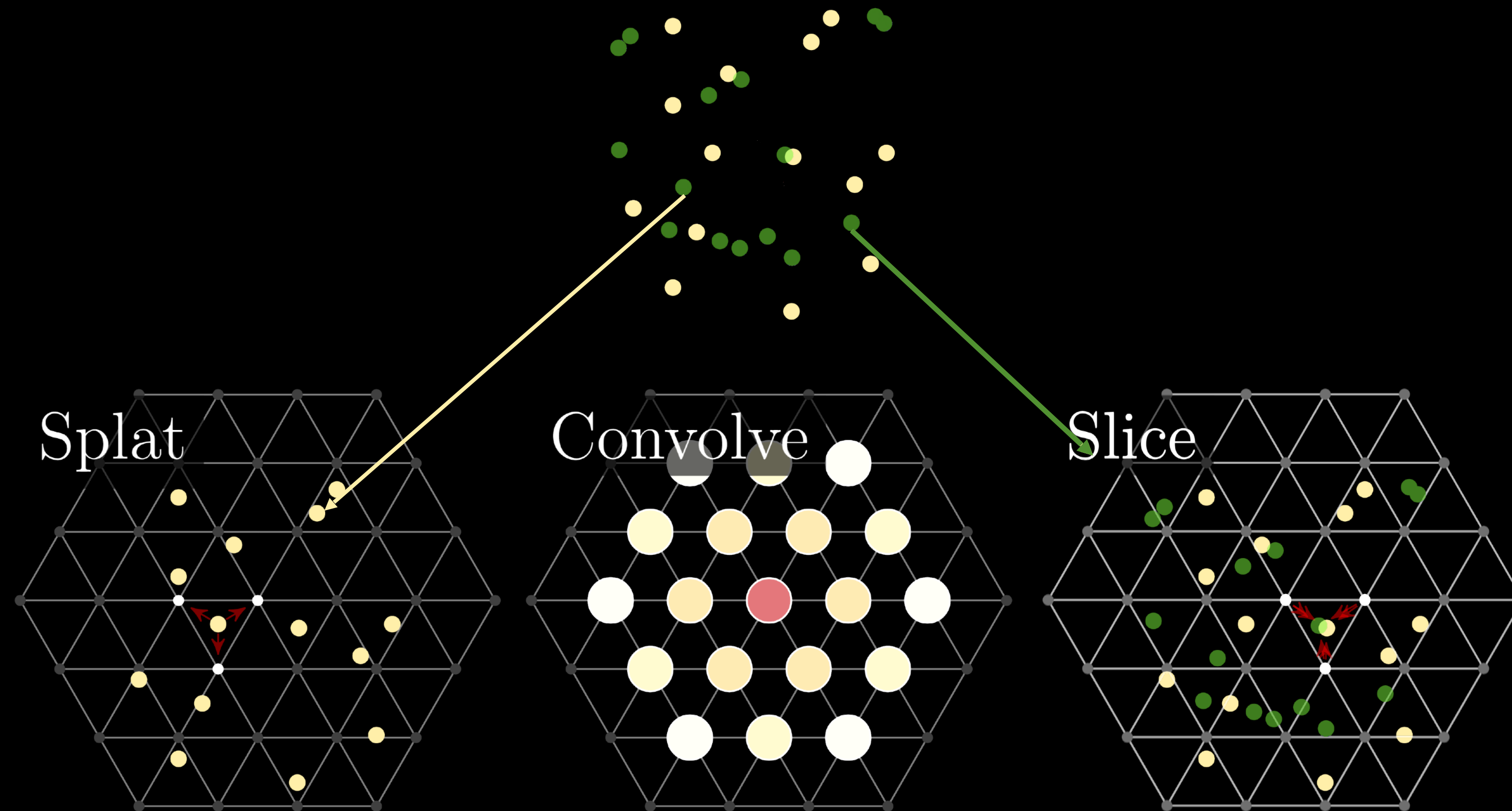
lattice scale: $8 \cdot \lambda_0$



lattice feature: (x, y, z)

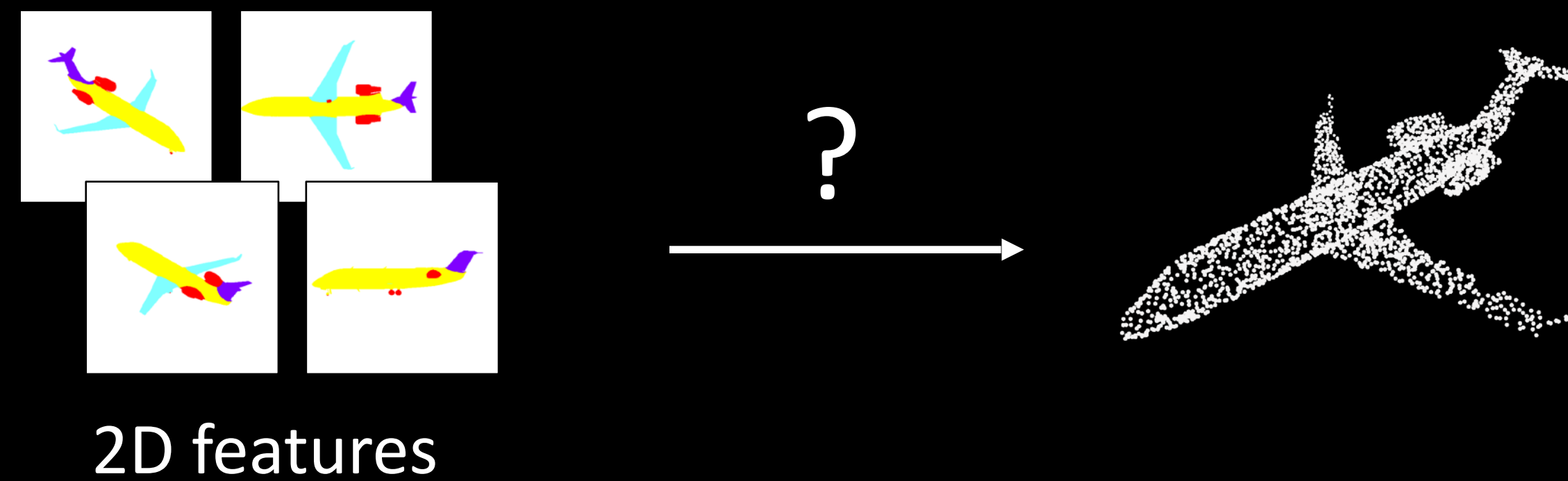
lattice scale: λ_0

Input and output can be at different points



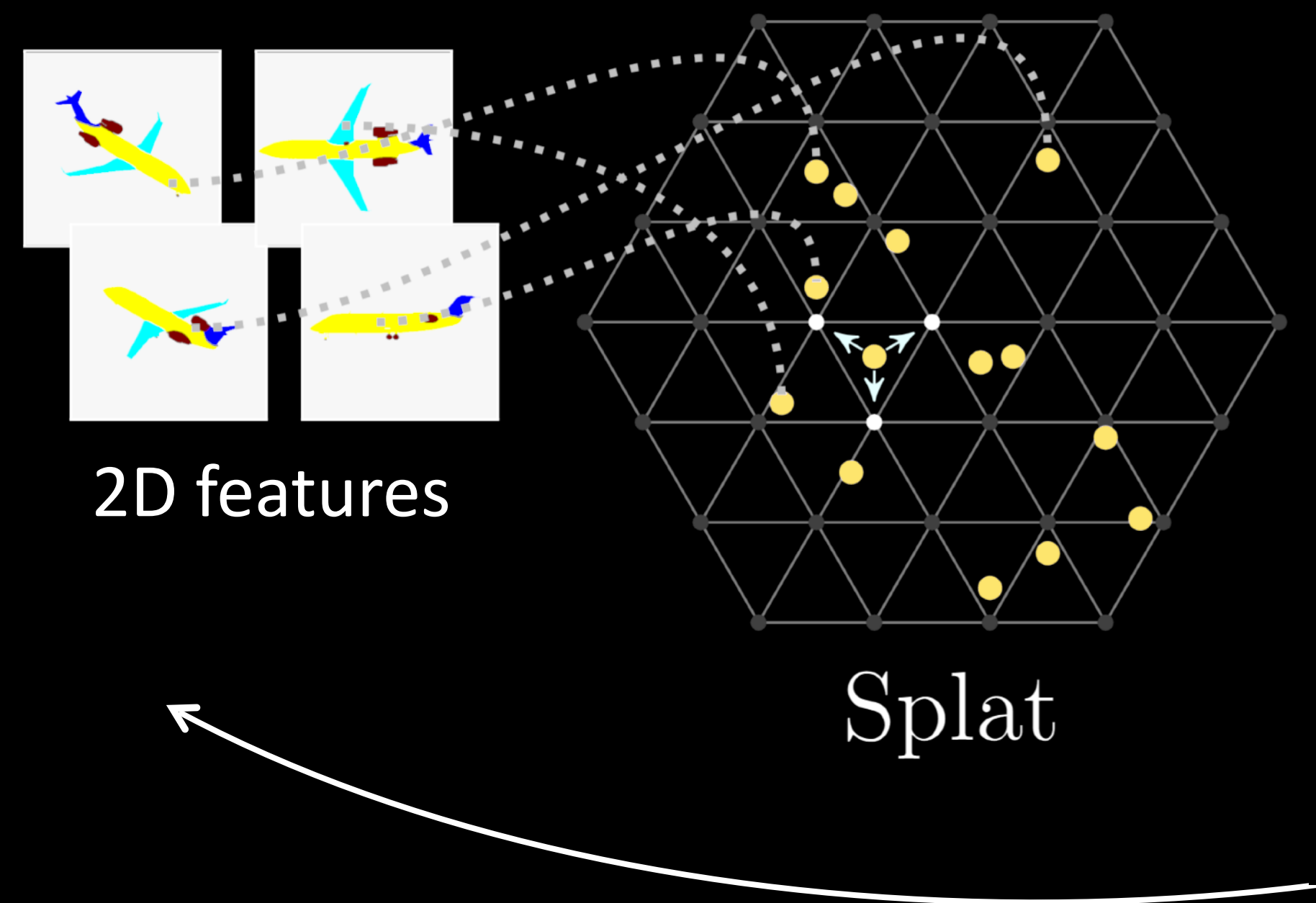
Smooth projection between 2D and 3D

- Each 2D pixel p : $f(p), (x, y, z)$

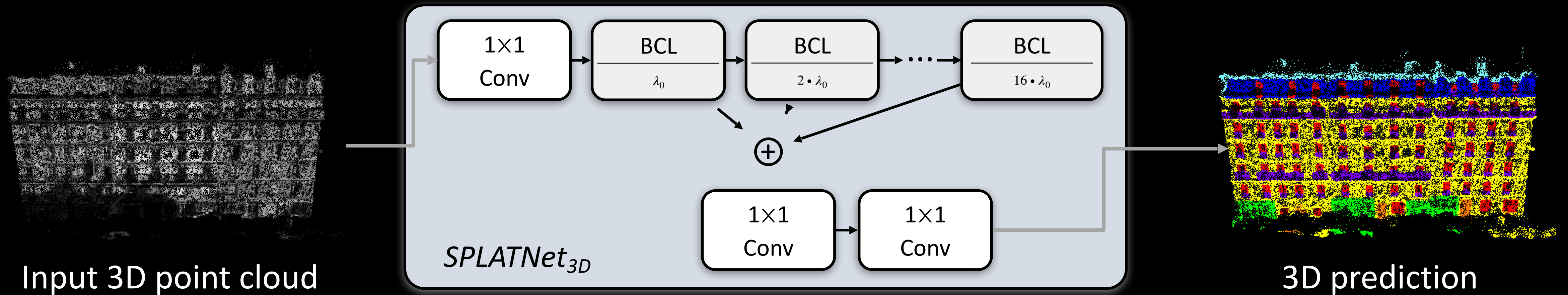


Smooth projection between 2D and 3D

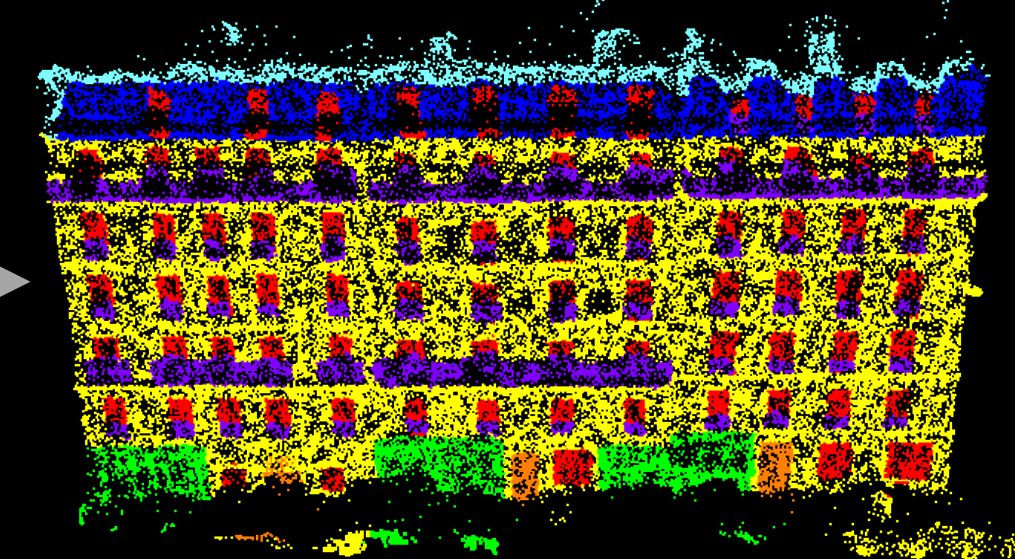
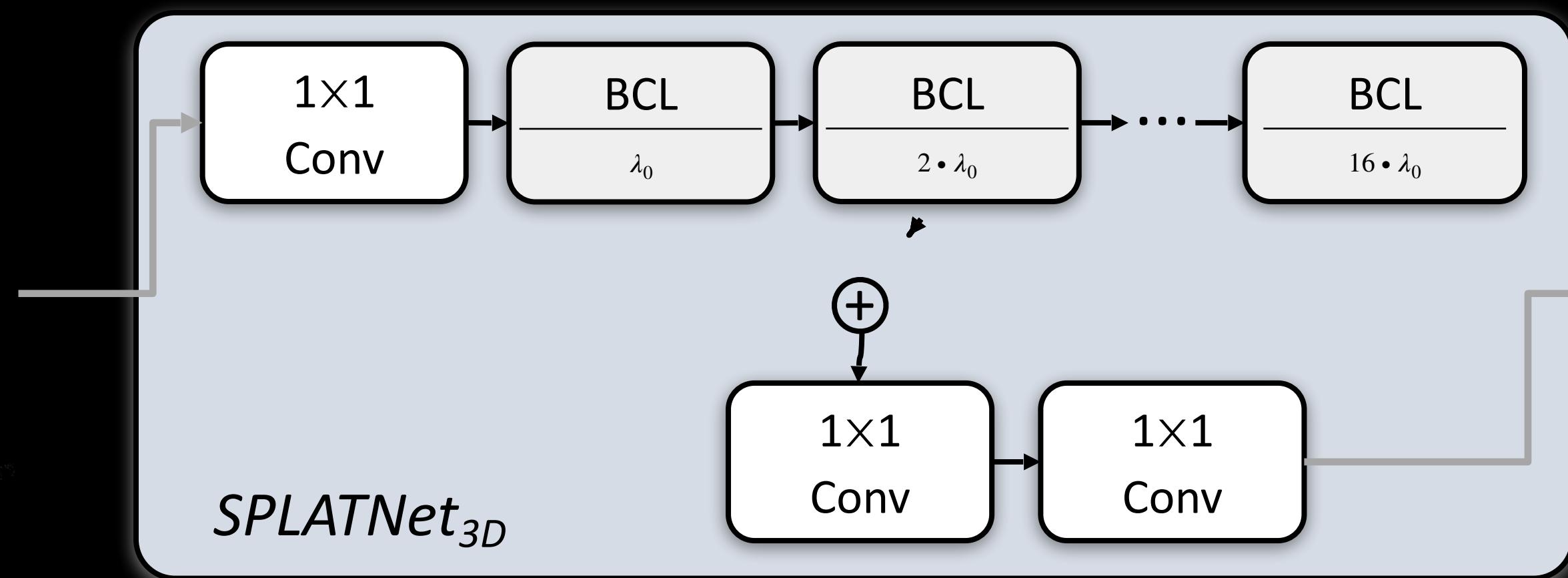
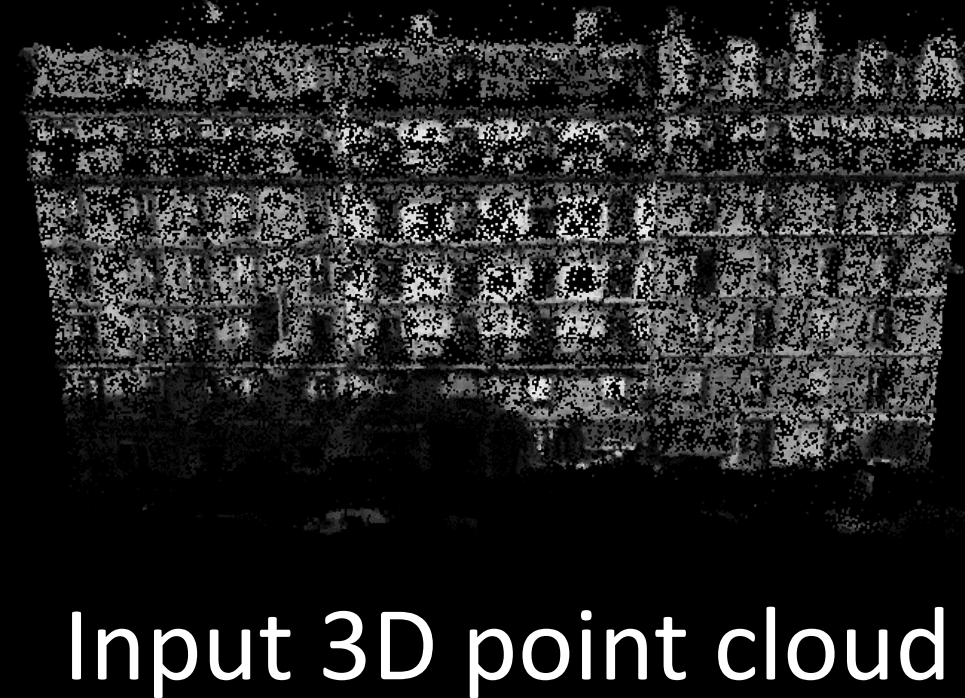
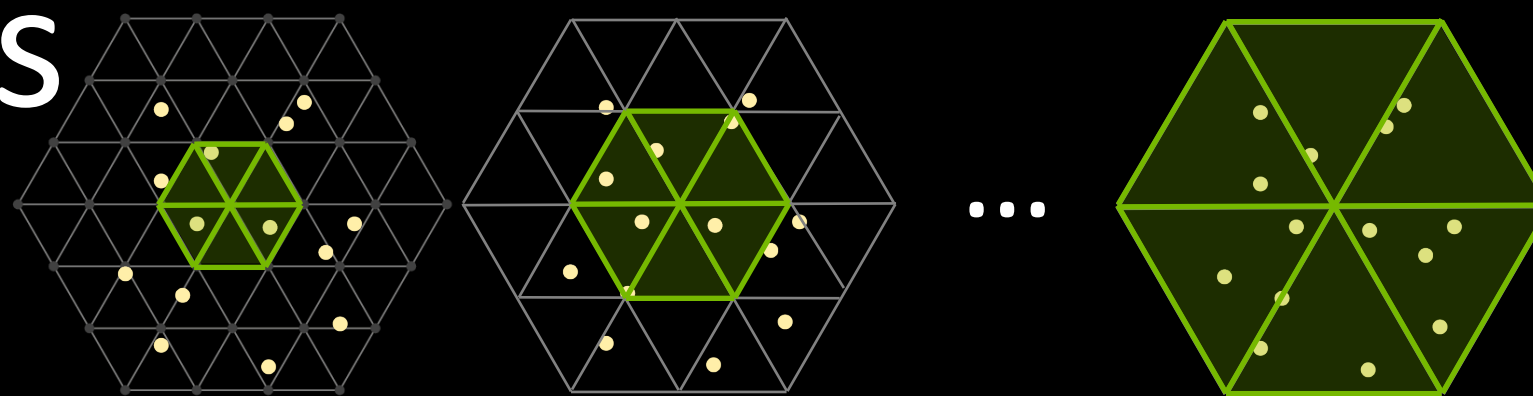
- Each 2D pixel p : $f(p), (x, y, z)$
 - Point value: $f(p)$
 - Point position: (x, y, z)



SPLATNet Architectures

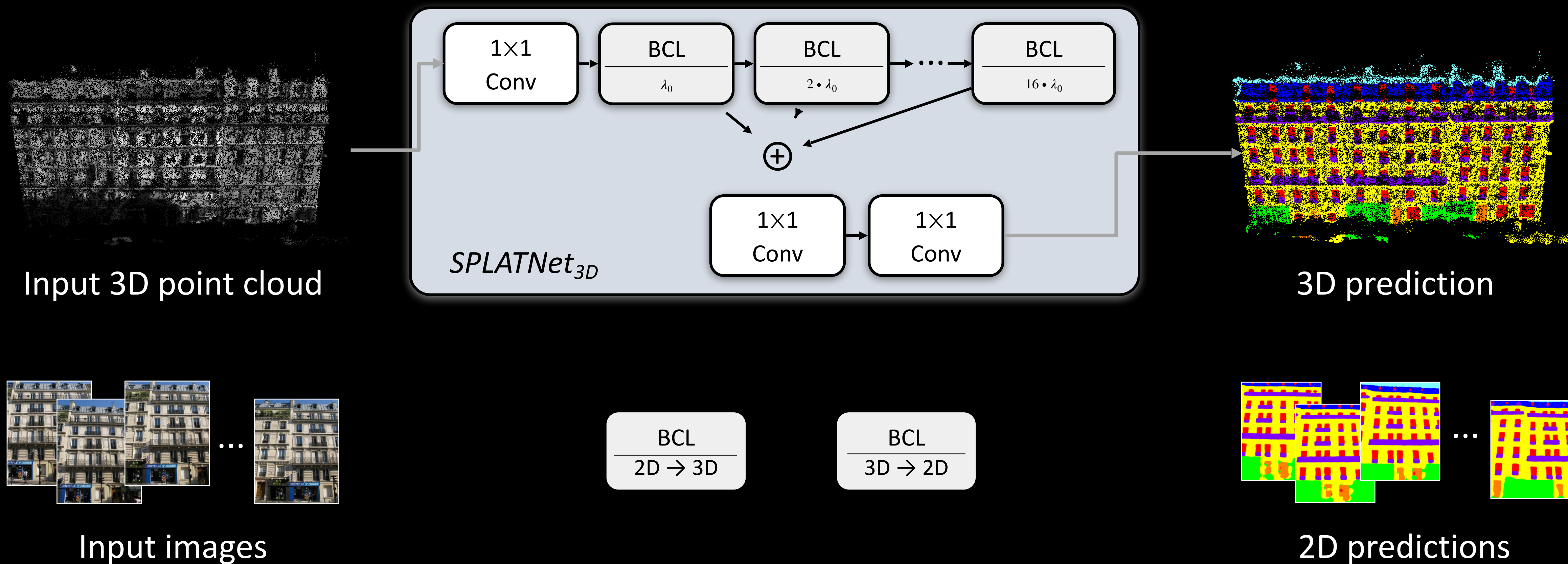


SPLATNet Architectures

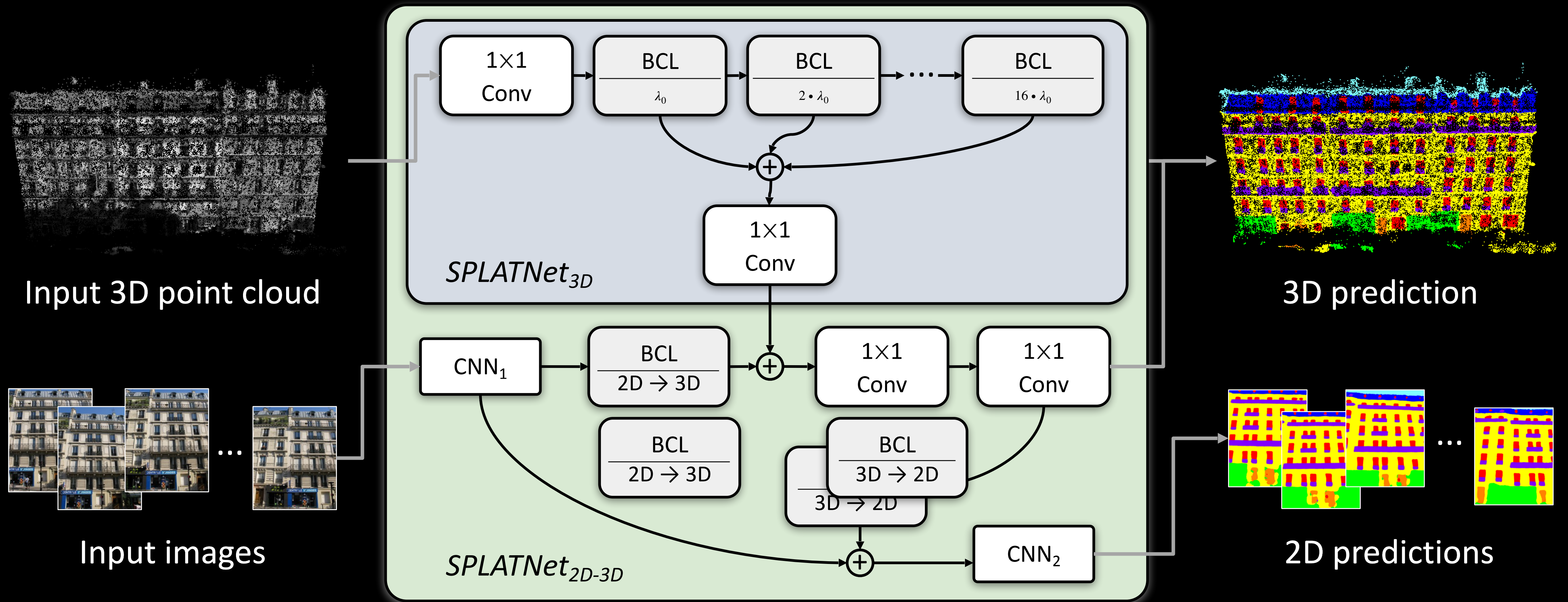


\oplus concatenation

SPLATNet Architectures



SPLATNet Architectures



Facade Segmentation (Ruemonge2014 [1])

with only 3D data

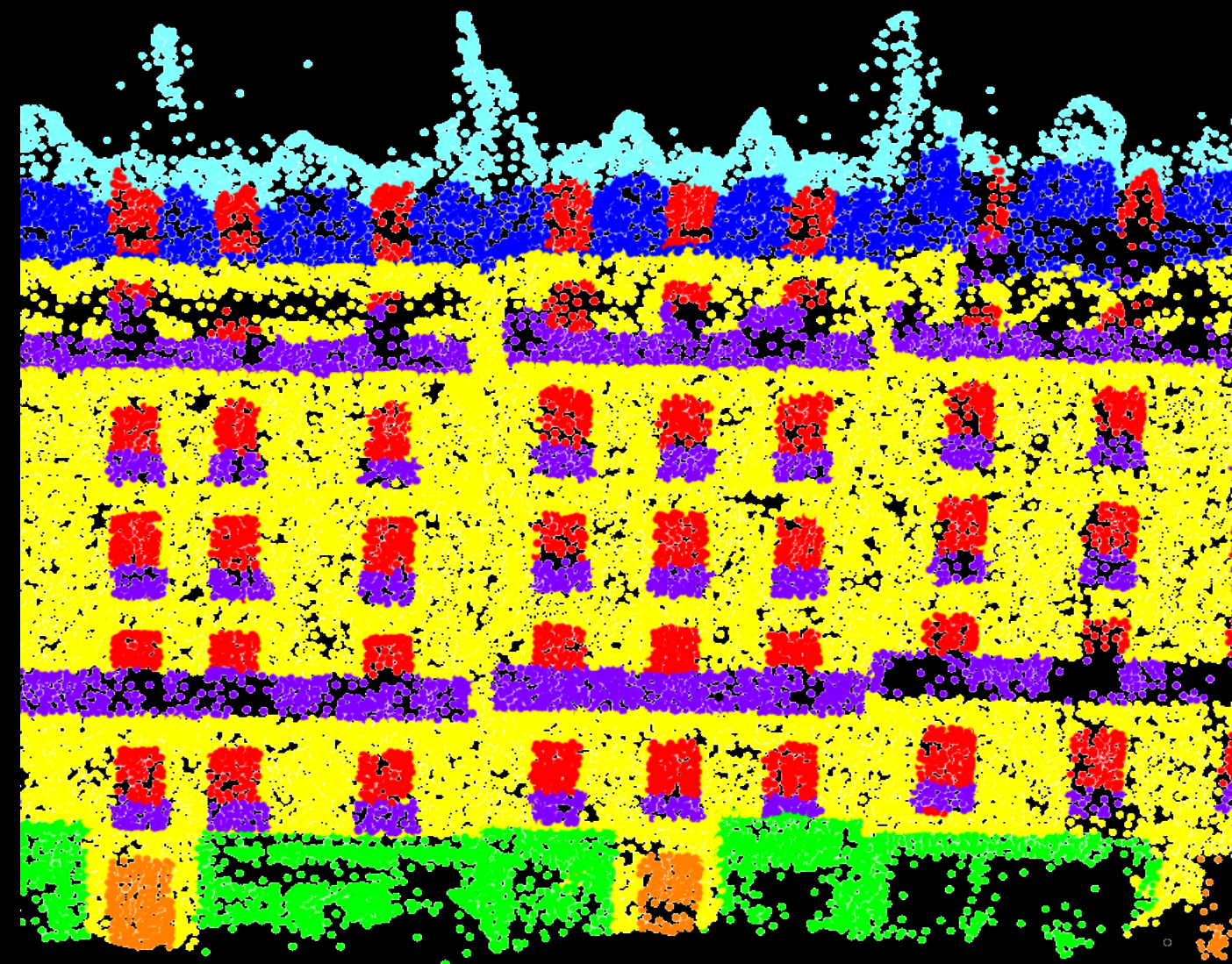
	IoU	runtime (min)
OctNet [2]	59.2	-
Autocontext _{3D} [3]	54.4	16
SPLATNet_{3D}	65.4	0.06

with both 2D & 3D data

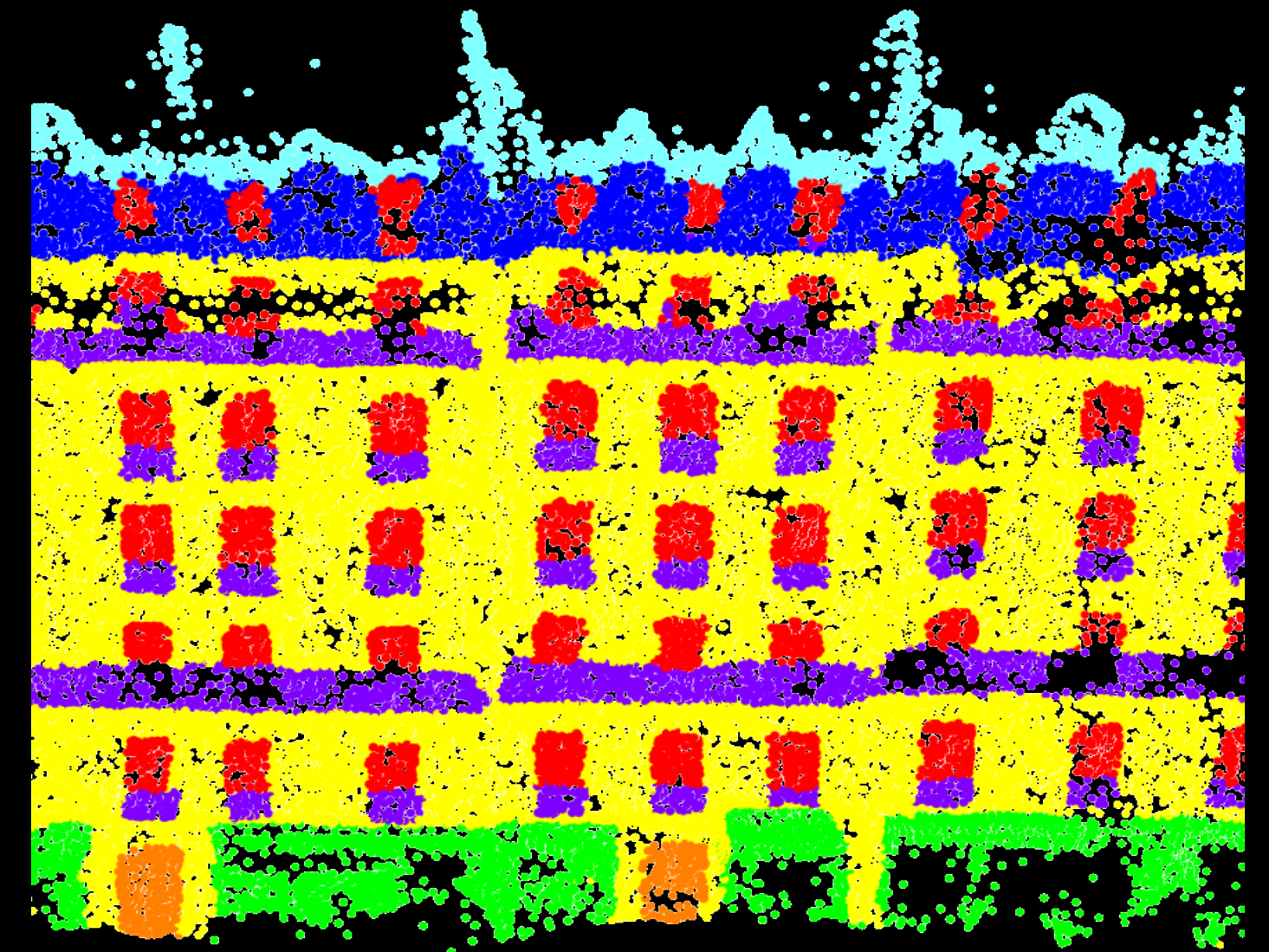
	IoU	runtime (min)
Autocontext _{2D-3D} [3]	62.9	87
SPLATNet_{2D-3D}	69.8	1.2



Input point cloud



Ground-truth



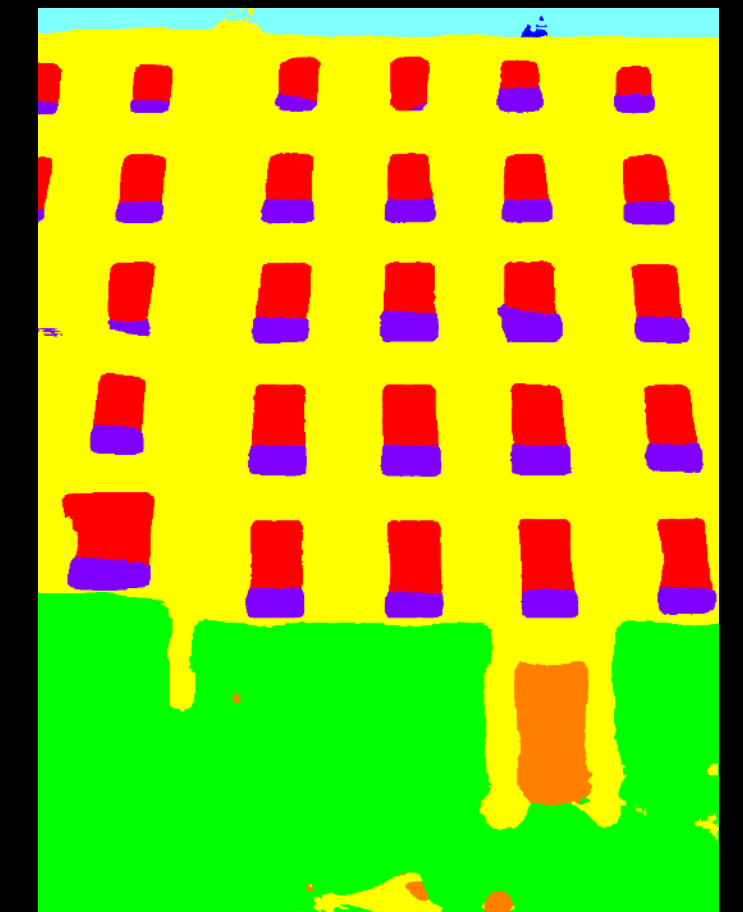
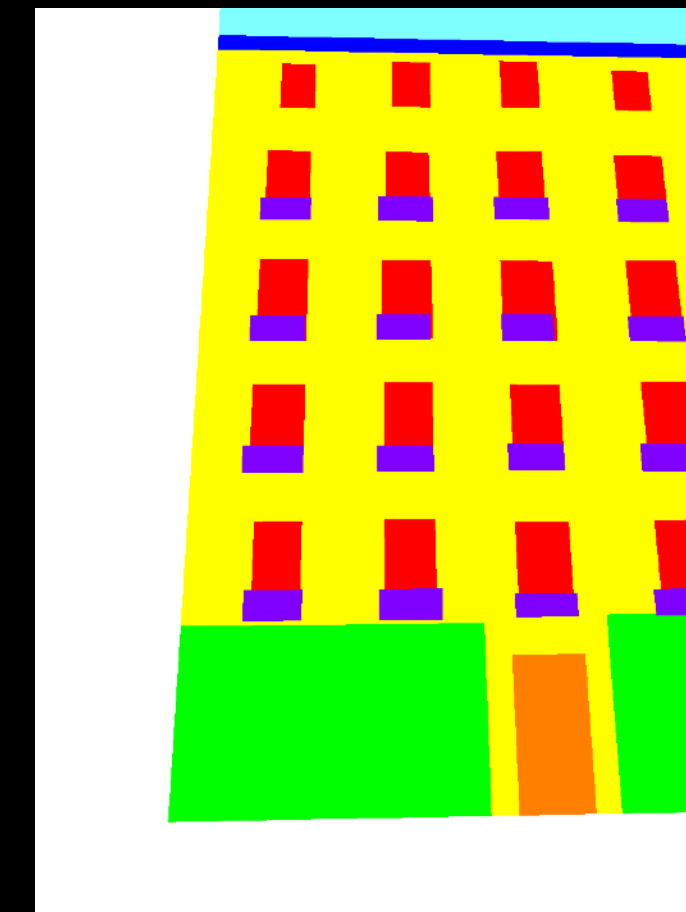
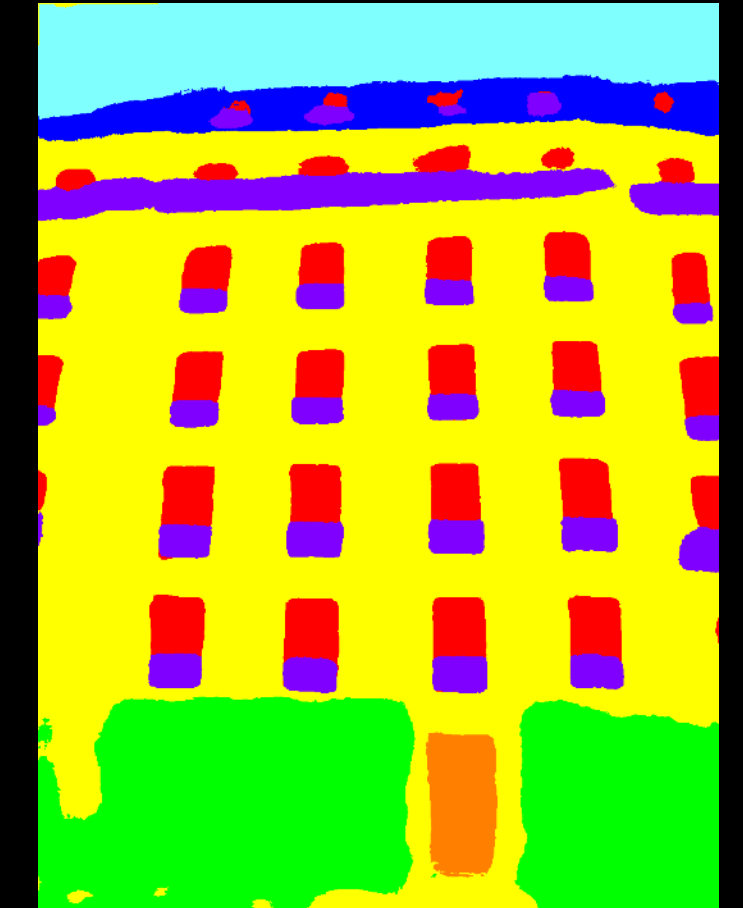
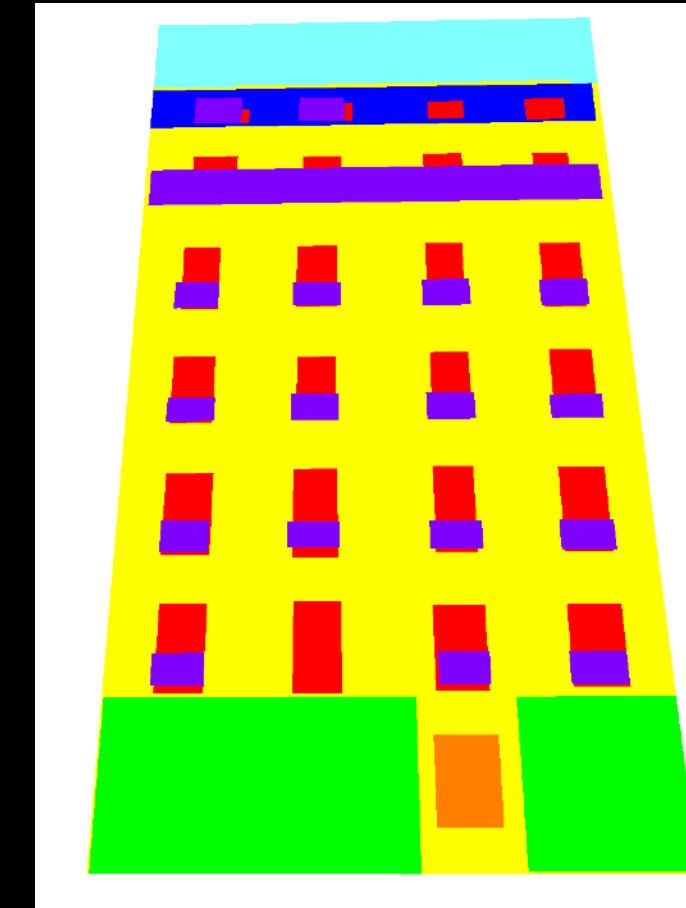
Ours (SPLATNet_{2D-3D})

2D Predictions

	IoU	runtime (min)
Autocontext _{2D} [1]	6055	117.7
Autocontext _{2D-3D} [1]	6277	146
2D CNN [2] only	6933	0.84
SPLATNet _{2D-3D}	7066	4.34

[1] Gadde *et al.* PAMI '17

[2] Chen *et al.* ICLR '15



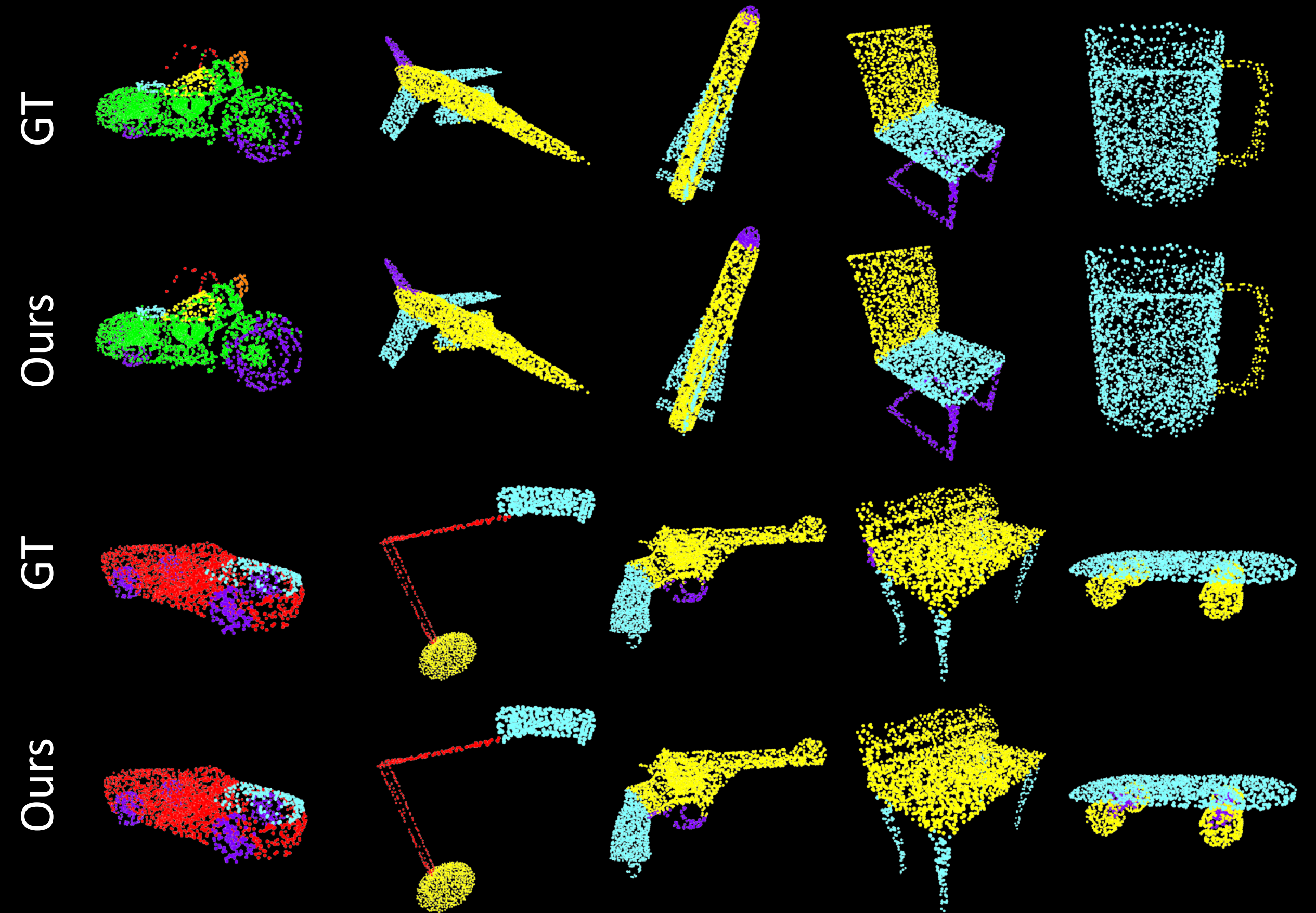
Input image

Ground-truth

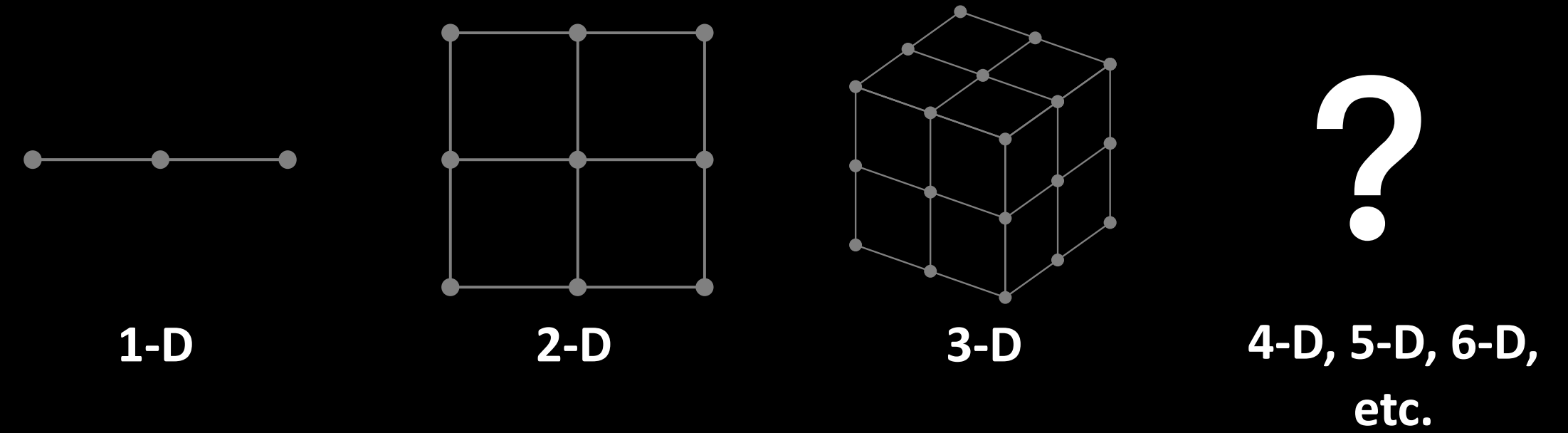
Our prediction

3D Object Part Labeling (ShapeNet [1])

	class avg. IoU	instance avg. IoU
Yi <i>et al.</i> [1]	79.0	81.4
3DCNN [2]	74.9	79.4
Kd-network [3]	77.4	82.3
PointNet [2]	80.4	83.7
PointNet++ [4]	81.9	85.1
SyncSpecCNN [5]	82.0	84.7
SPLATNet_{3D}	82.0	84.6
SPLATNet_{2D-3D}	83.7	85.4



Remarks



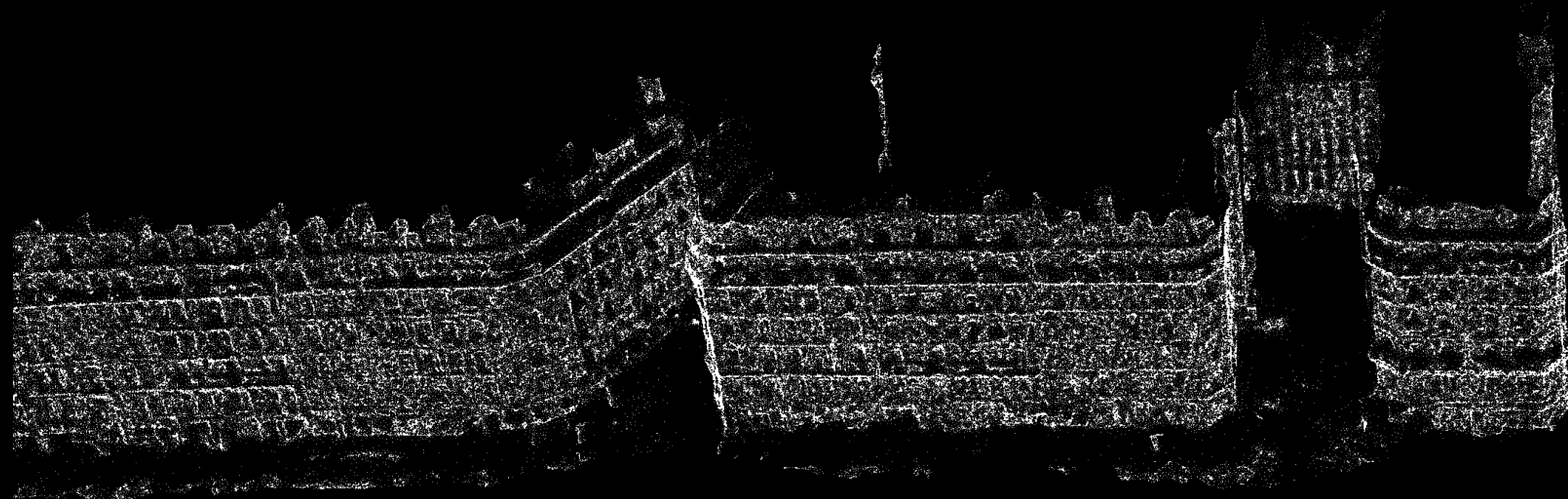
Equip neural networks with a richer class of **sparse high-dimensional operations**

1. Data that are inherently sparse and high-dimensional

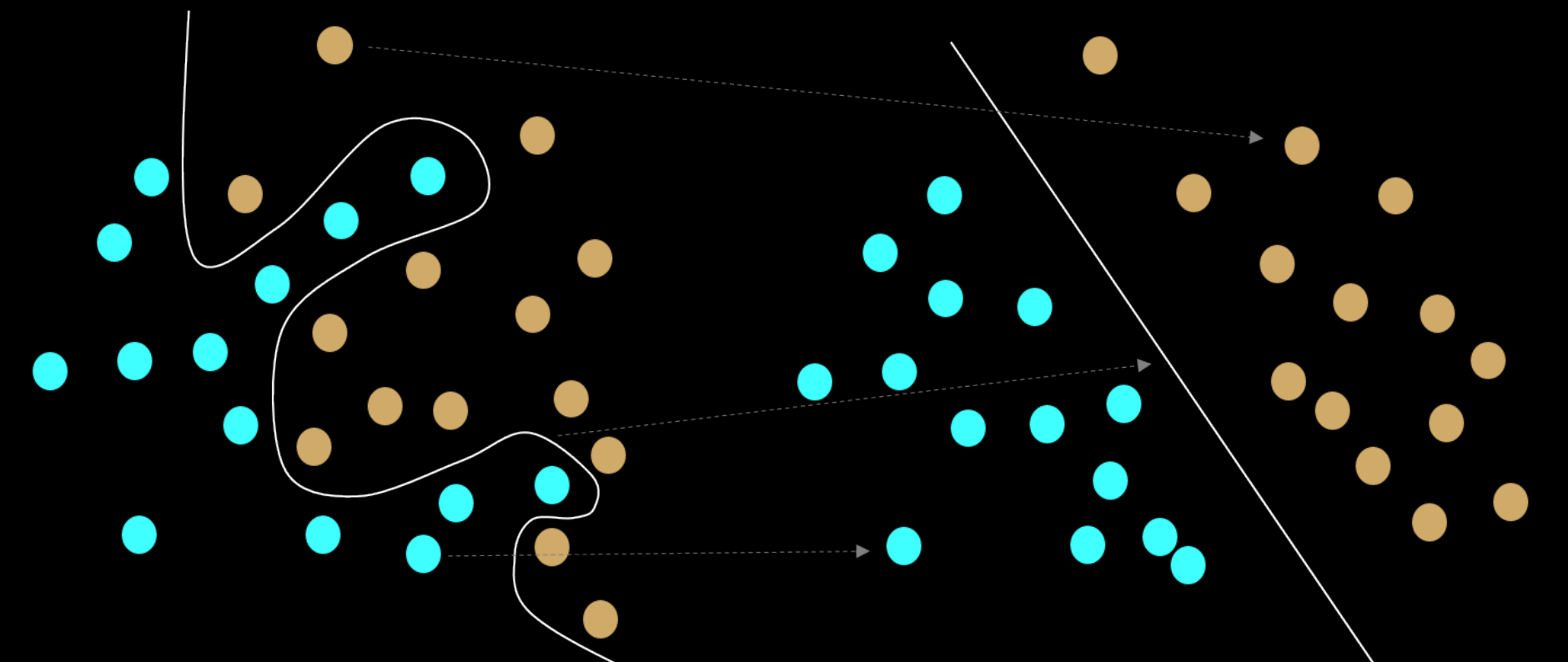
E.g. 3D data, BRDF material properties, video

2. Project data to higher dimensions

E.g. bilateral filter



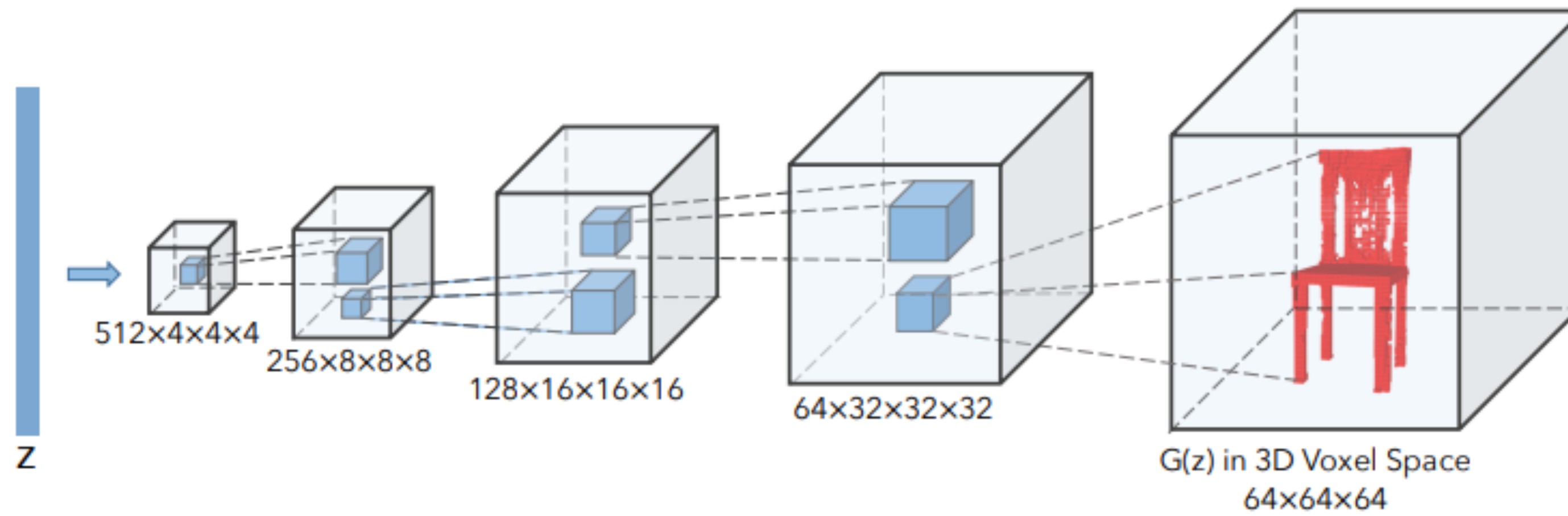
Ruemonge2014 [Riemenschneider et al. ECCV '14]



Input Space

High dimensional space

Volumetric Generation



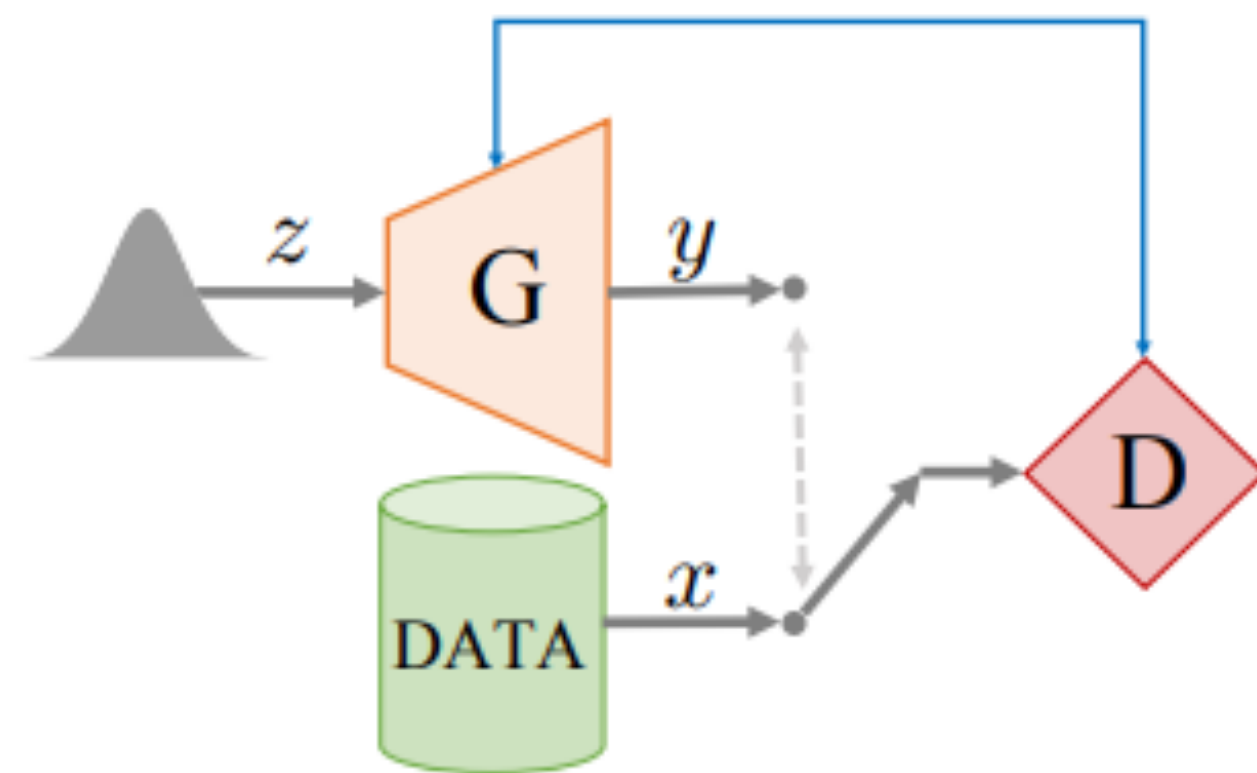
$$\log D(x) + \log(1 - D(G(z)))$$



Wu et al., "Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling", *NeurIPS 2016*

Point Cloud Generation

- FC layer as generator
- PointNet as discriminator
- WGAN



Many Issues

- Still cannot generate high quality local details
- Still hard to generate complex structures
- Partly limited by 3D training data, which is much smaller than images

Recent trends

Image (s) to 3D

- **Task:** Given input image(s), estimate the underlying 3D structure.
- **Scale ambiguity:** Outputs are often estimated only up to an unknown scale
- **Monocular depth estimation:** Single image to depth
 - Early models were ConvNet-based and trained on small datasets
 - Recent models are based on diffusion models and transformers, and are trained on massive datasets.
- **Multi-view 3D:** Multiple images as input; simultaneously estimate camera poses and the underlying geometry for a more complete 3D reconstruction
- Recent methods are often transformer-based and trained on massive datasets

Monocular depth estimation

Repurposing Diffusion-Based Image Generators for Monocular Depth Estimation

Bingxin Ke Anton Obukhov Shengyu Huang
Nando Metzger Rodrigo Caye Daudt Konrad Schindler
Photogrammetry and Remote Sensing, ETH Zürich



Monocular depth estimation

Fine-tuned “Stable Diffusion” to model $P(\text{depth} | \text{image})$

Training takes 2.5 days on a single NVIDIA RTX 4090 GPU

Trained on **Hypersim** (a photorealistic indoor scene dataset)

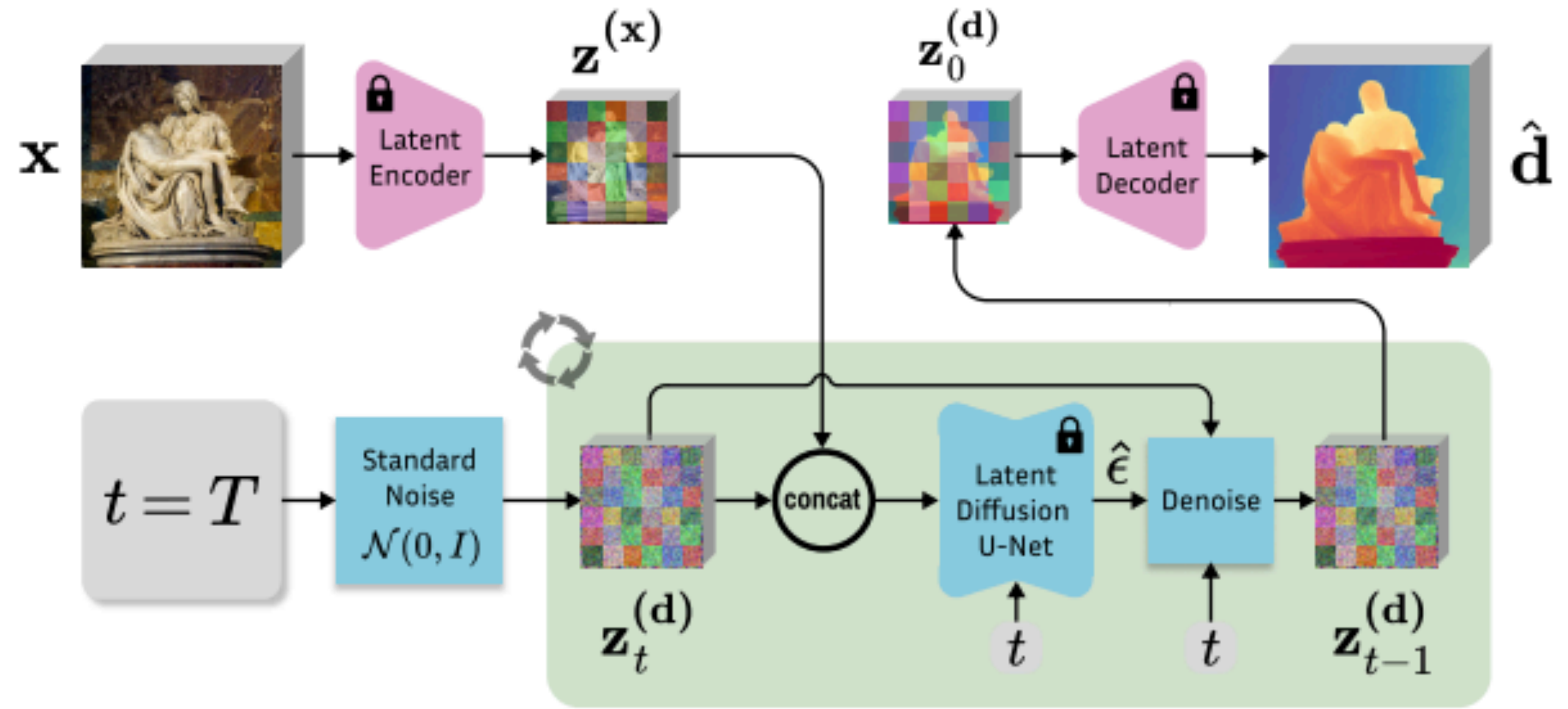


Figure 3. **Overview of the Marigold inference scheme.** Given an input image x , we encode it with the original Stable Diffusion VAE into the latent code $z^{(x)}$, and concatenate with the depth latent $z_t^{(d)}$ before giving it to the modified fine-tuned U-Net on every denoising iteration. After executing the schedule of T steps, the resulting depth latent $z_0^{(d)}$ is decoded into an image, whose 3 channels are averaged to get the final estimation \hat{d} . See Sec. 3.4 for details.

Monocular depth estimation

Input RGB Image

MiDaS

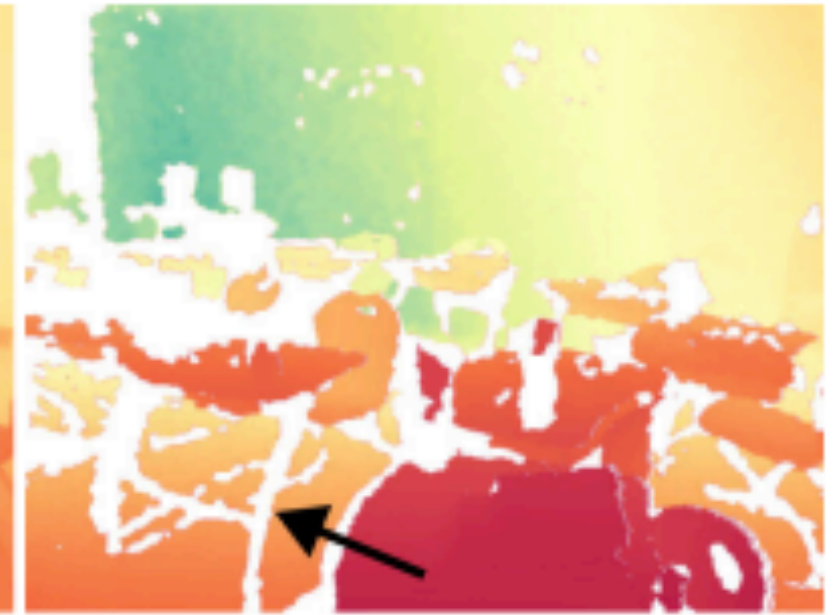
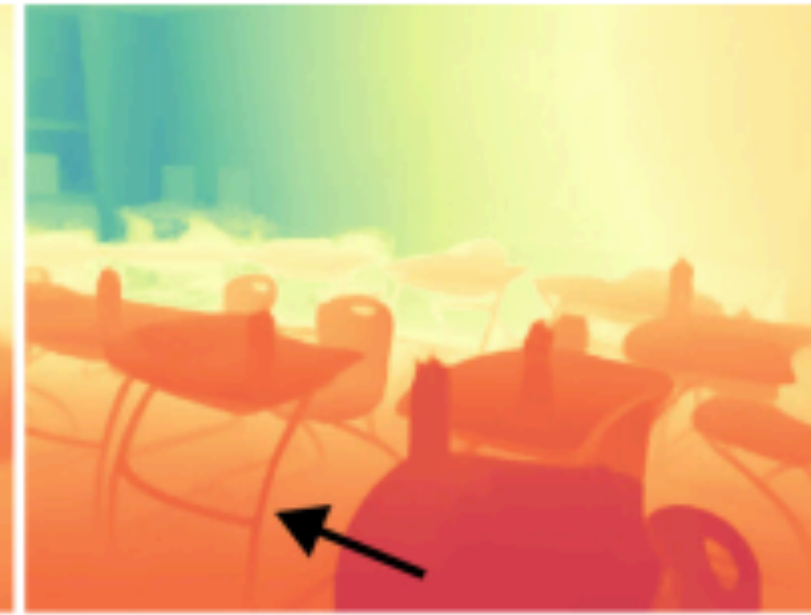
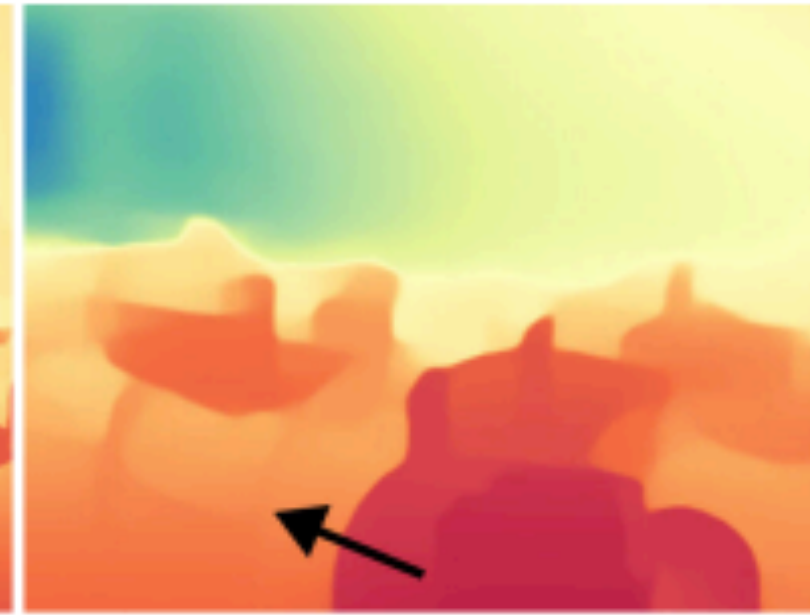
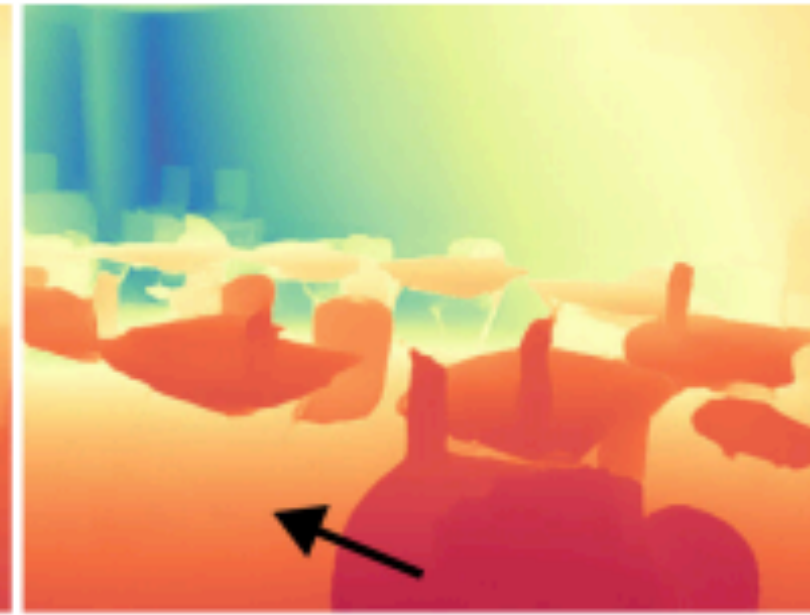
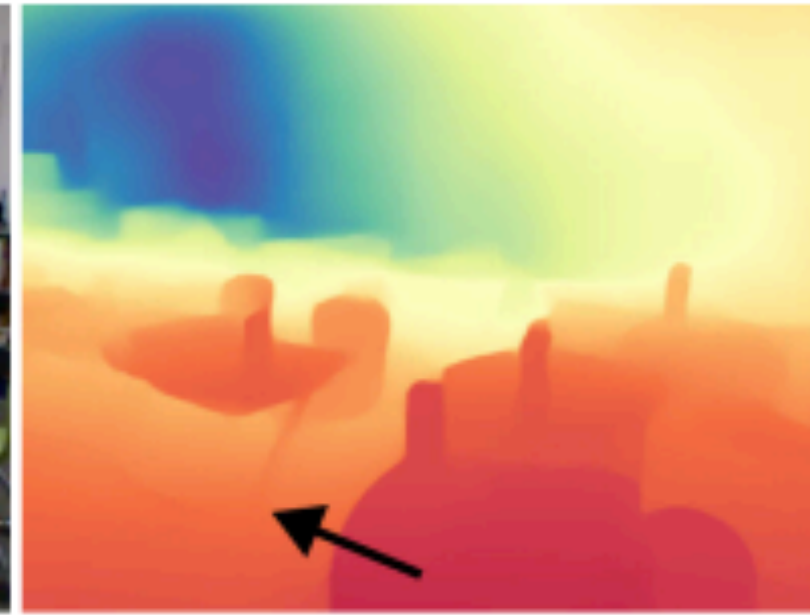
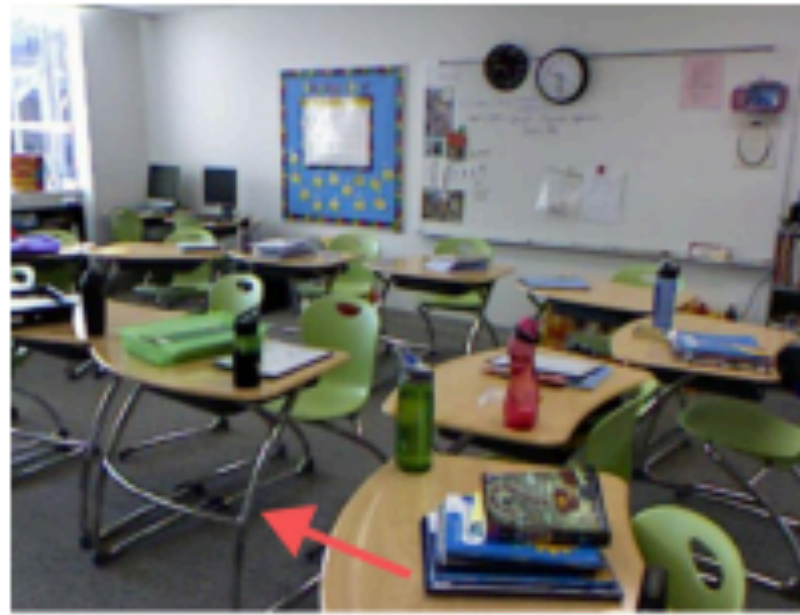
Omnidata

DPT

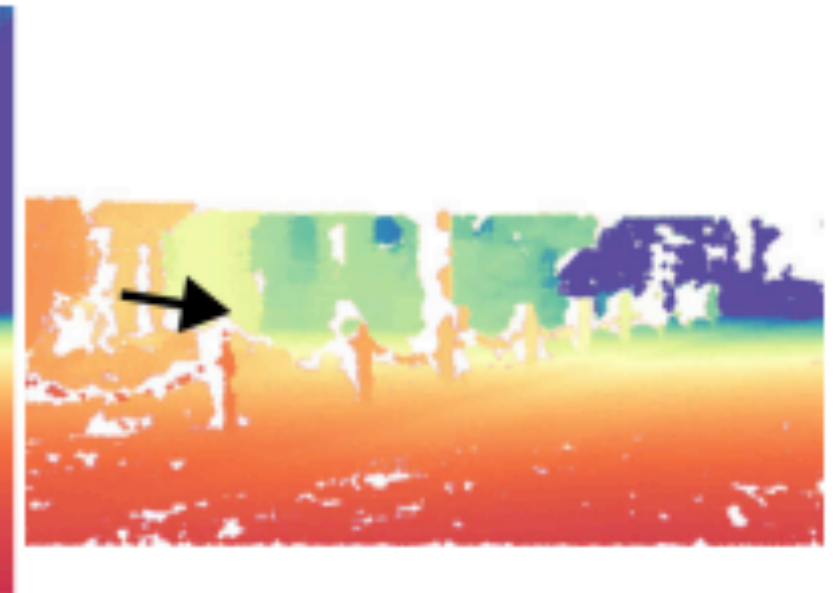
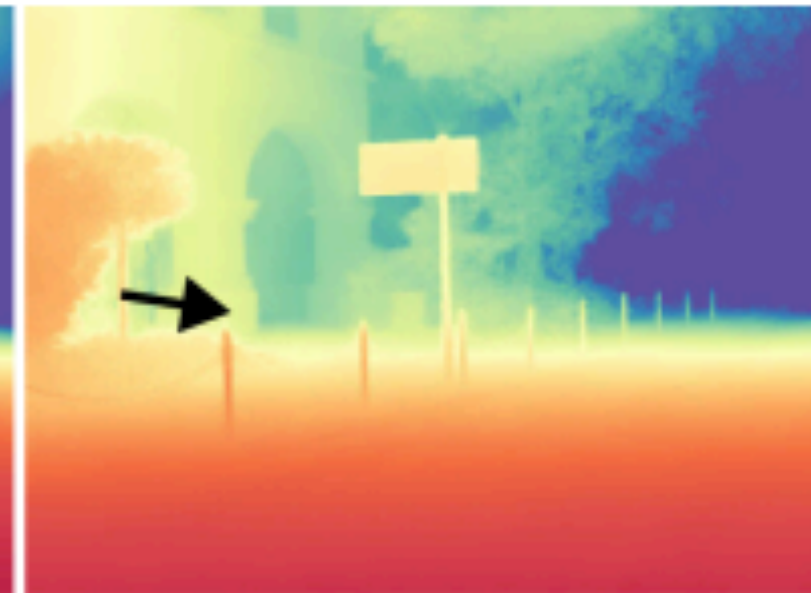
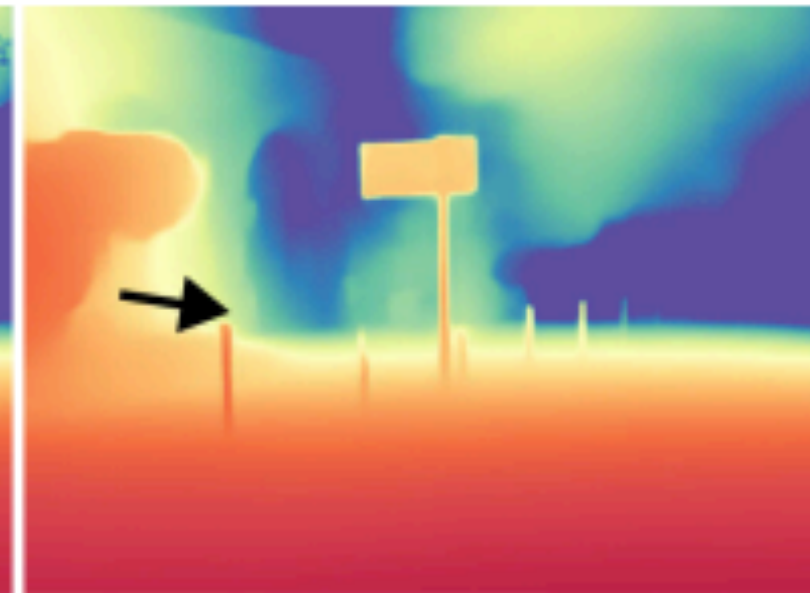
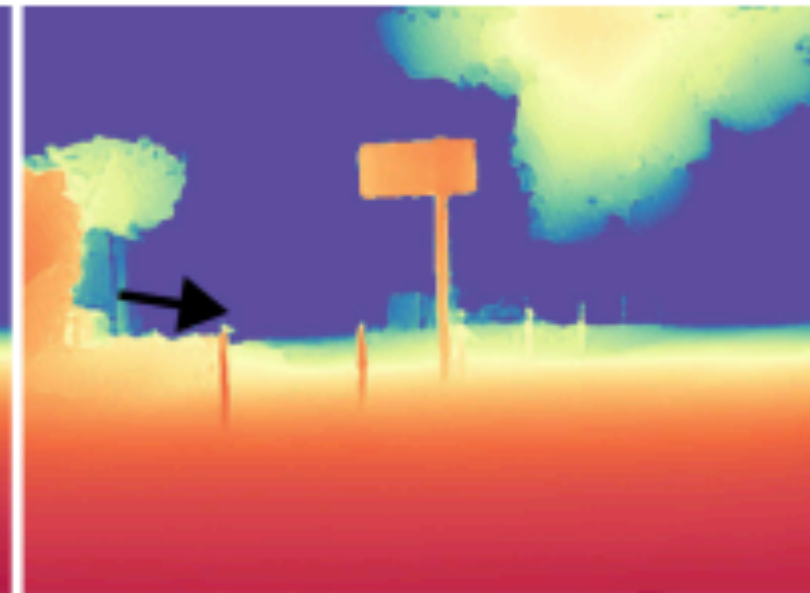
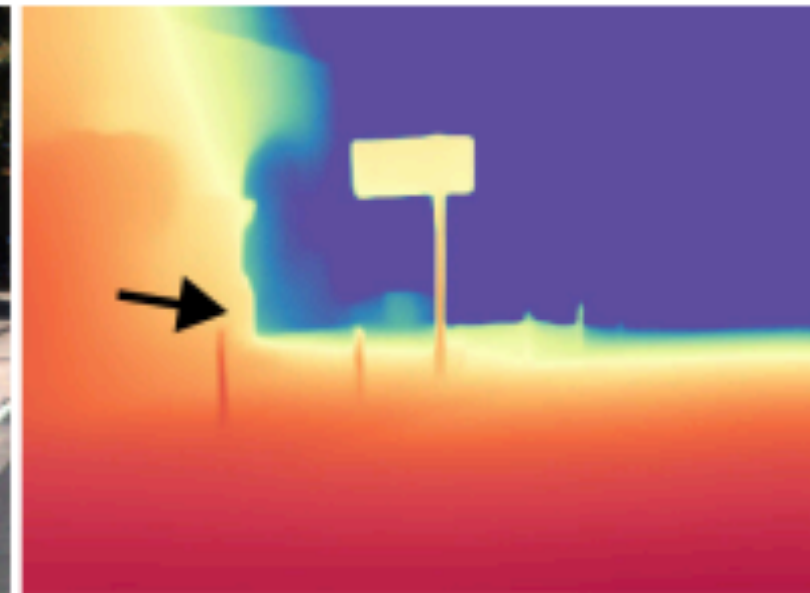
Marigold (ours)

Ground Truth

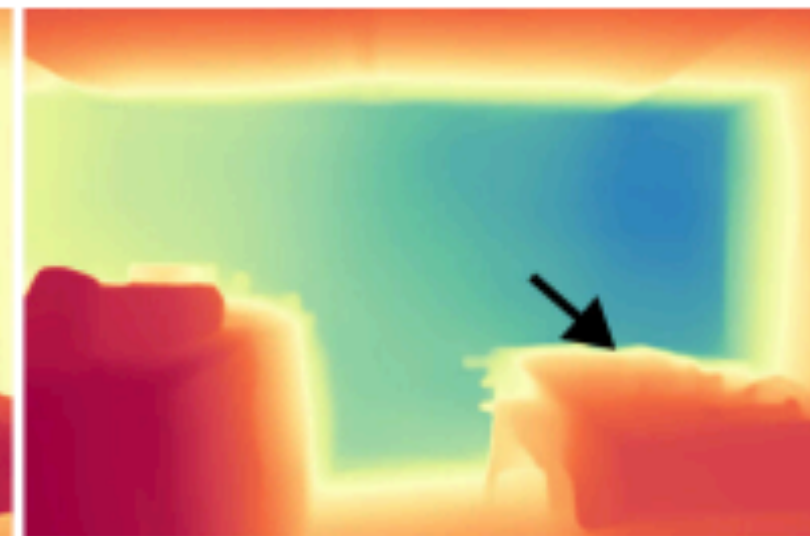
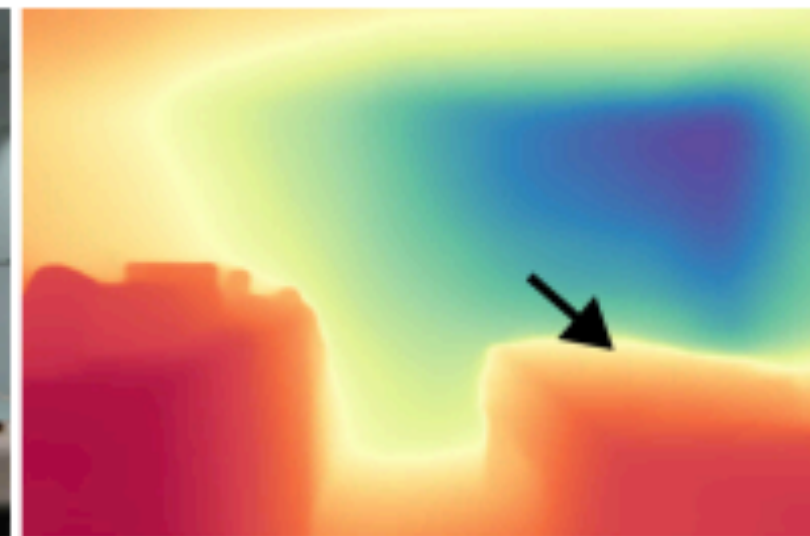
NYUv2



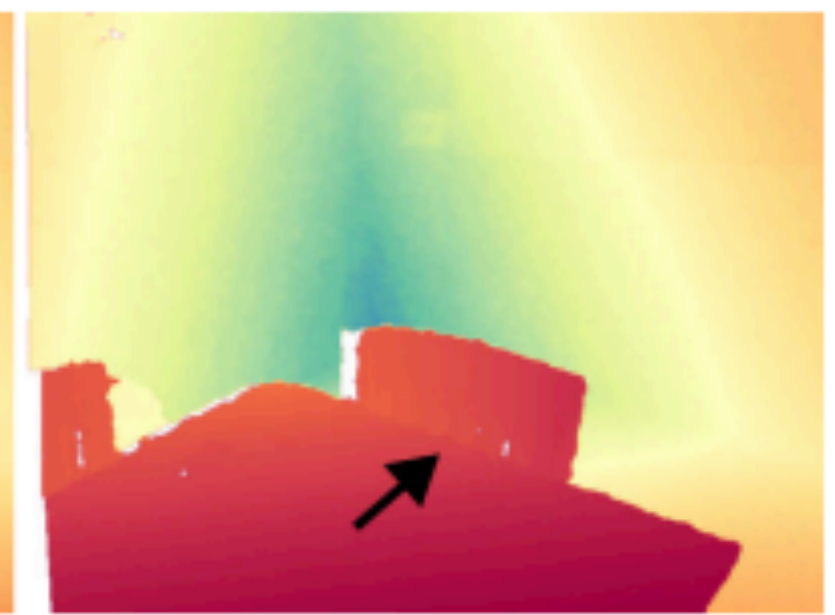
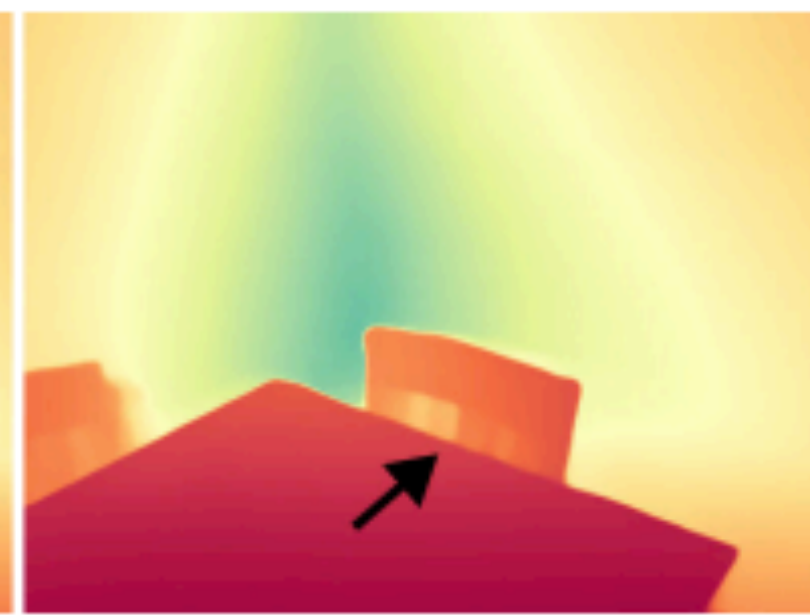
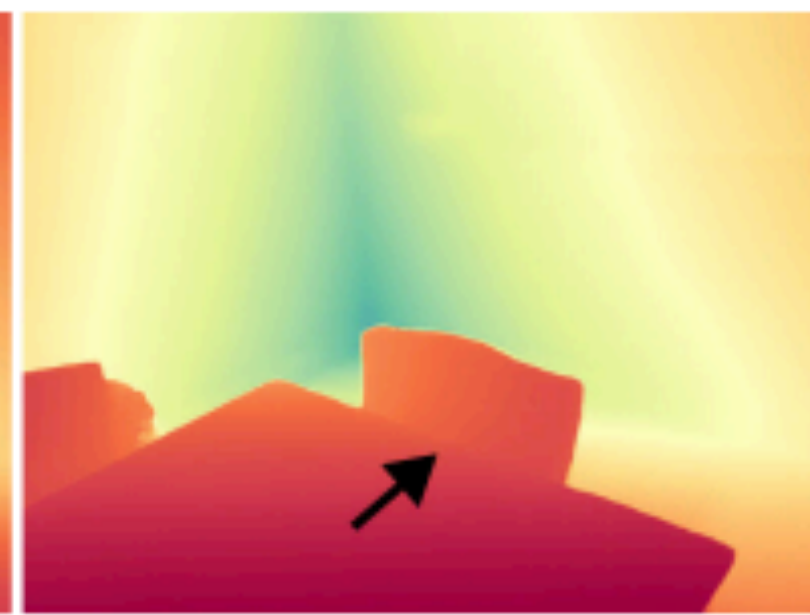
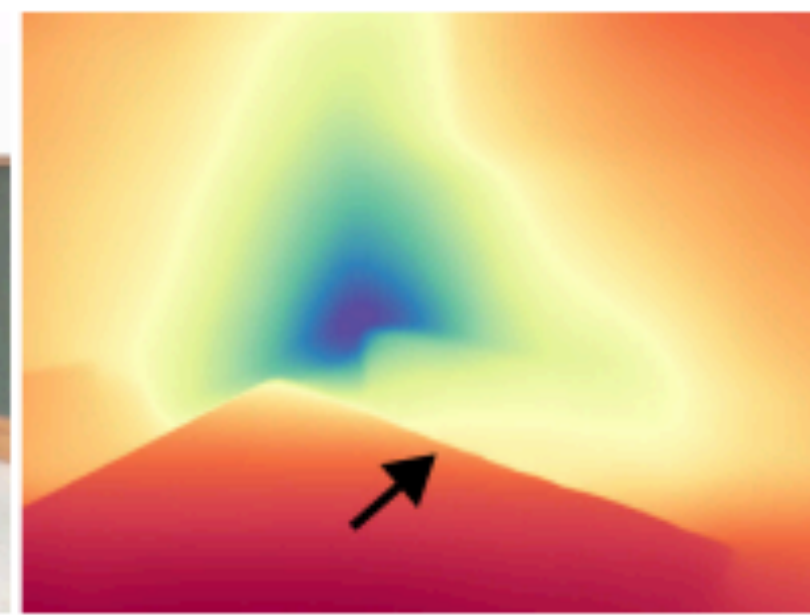
KITTI



ETH3D



Scannet



DUST3R: Geometric 3D Vision Made Easy

Shuzhe Wang^{*}, Vincent Leroy[†], Yann Cabon[†], Boris Chidlovskii[†] and Jerome Revaud[†]

^{*}Aalto University

[†]Naver Labs Europe

shuzhe.wang@aalto.fi

firstname.lastname@naverlabs.com

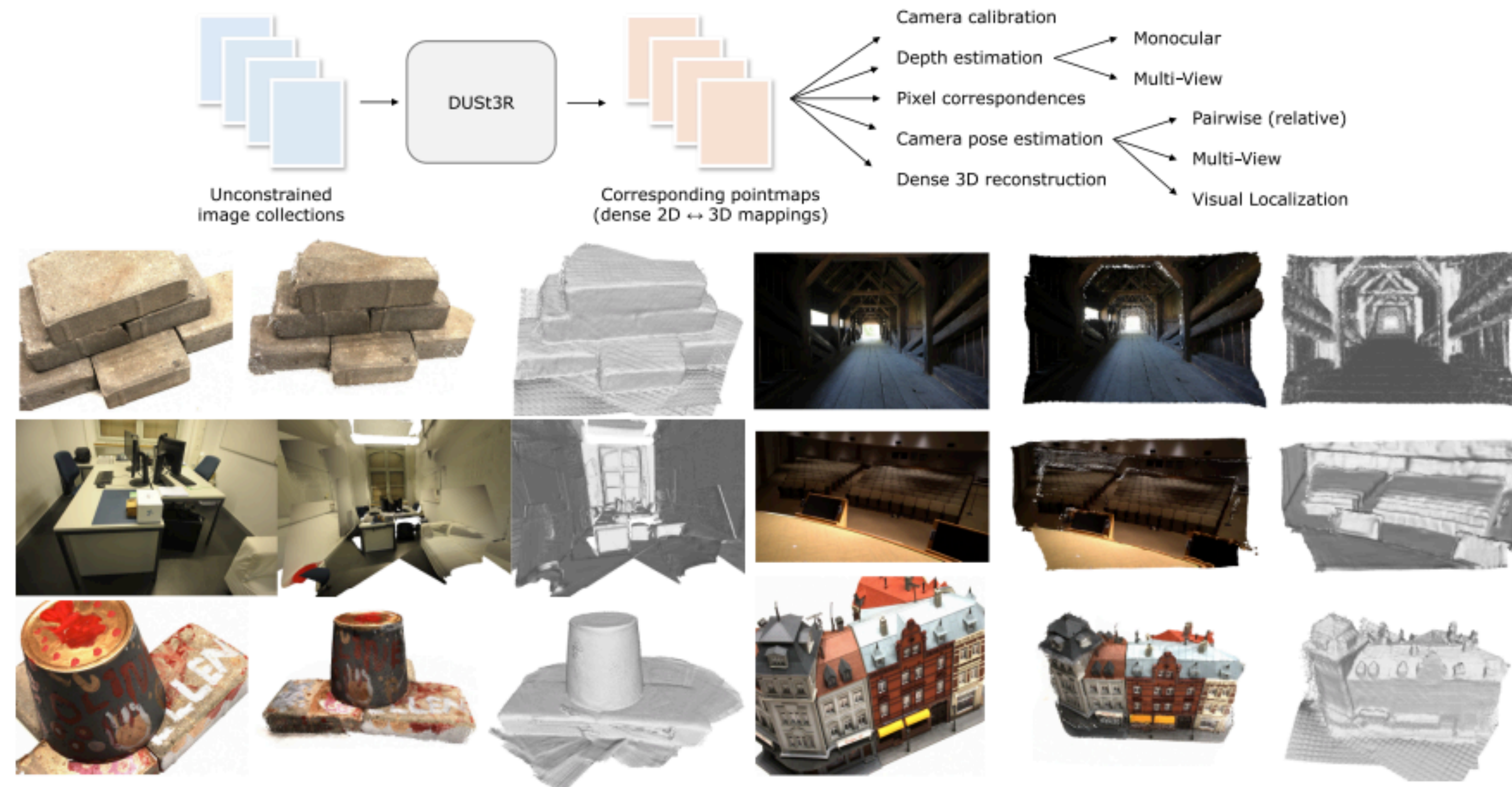


Figure 1. **Overview:** Given an unconstrained image collection, *i.e.* a set of photographs with unknown camera poses and intrinsics, our proposed method **DUST3R** outputs a set of corresponding *pointmaps*, from which we can straightforwardly recover a variety of geometric quantities normally difficult to estimate all at once, such as the camera parameters, pixel correspondences, depthmaps, and fully-consistent 3D reconstruction. Note that DUST3R also works for a single input image (*e.g.* achieving in this case monocular reconstruction). We also show **qualitative examples** on the DTU, Tanks and Temples and ETH-3D datasets [1, 51, 108] obtained **without** known camera parameters. For each sample, from *left to right*: input image, colored point cloud, and rendered with shading for a better view of the underlying geometry.

DUSt3R: Geometric 3D Vision Made Easy

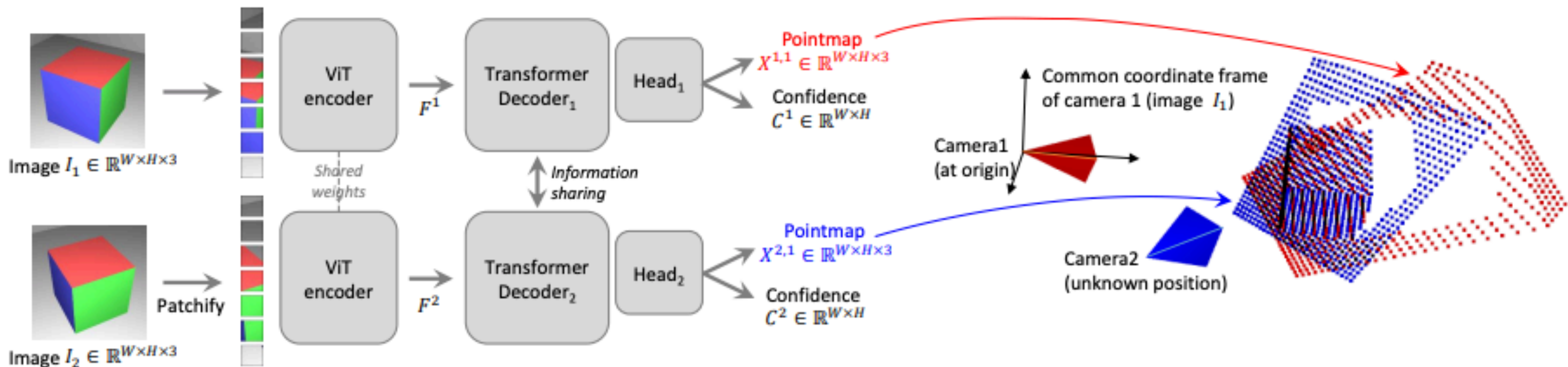
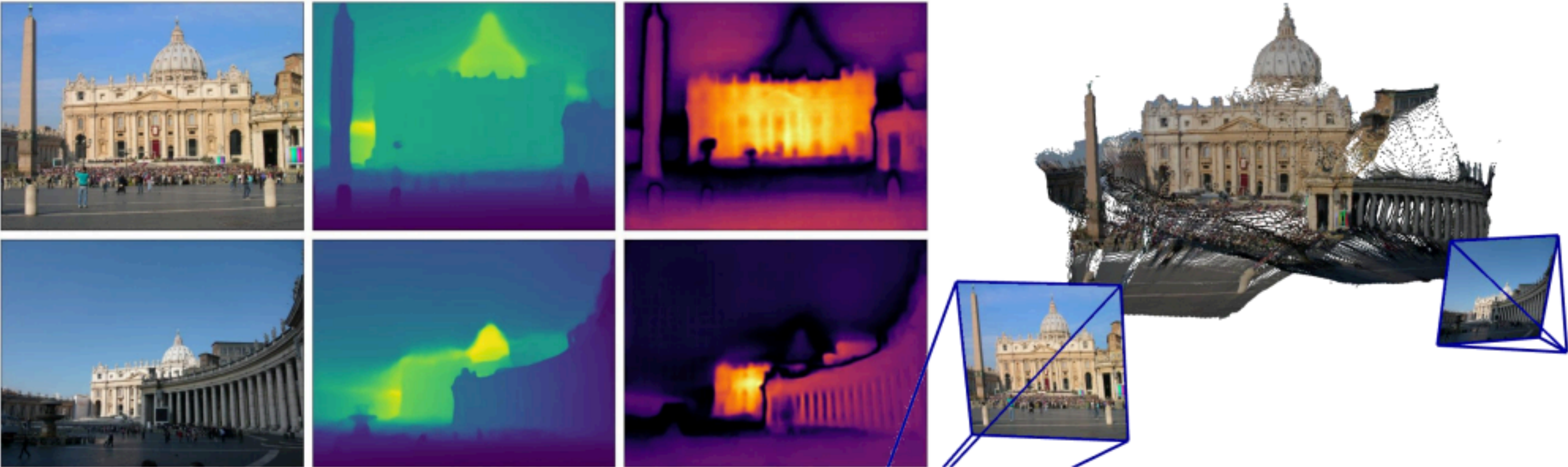
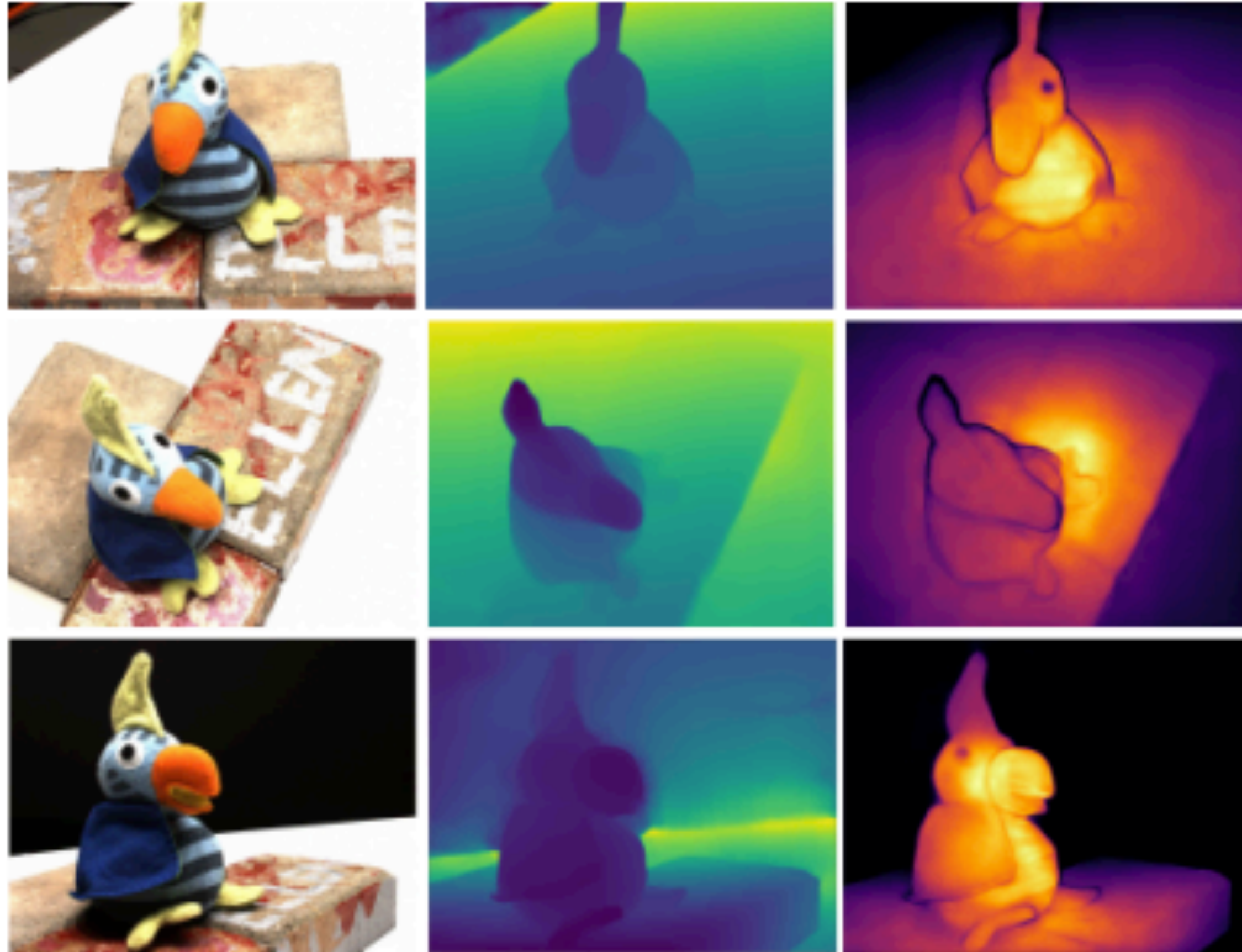


Figure 2. **Architecture of the network \mathcal{F} .** Two views of a scene (I^1, I^2) are first encoded in a Siamese manner with a shared ViT encoder. The resulting token representations F^1 and F^2 are then passed to two transformer decoders that constantly exchange information via cross-attention. Finally, two regression heads output the two corresponding pointmaps and associated confidence maps. Importantly, the two pointmaps are expressed in the same coordinate frame of the first image I^1 . The network \mathcal{F} is trained using a simple regression loss (Eq. (4))

DUSt3R: Geometric 3D Vision Made Easy



DUSt3R: Geometric 3D Vision Made Easy



VGGT: Visual Geometry Grounded Transformer

Jianyuan Wang^{1,2}

Minghao Chen^{1,2}

Nikita Karaev^{1,2}

Andrea Vedaldi^{1,2}

Christian Rupprecht¹

David Novotny²

¹Visual Geometry Group, University of Oxford

²Meta AI



Figure 1. **VGGT** is a large feed-forward transformer with minimal 3D-inductive biases trained on a trove of 3D-annotated data. It accepts up to hundreds of images and predicts cameras, point maps, depth maps, and point tracks for all images at once in less than a second, which often outperforms optimization-based alternatives without further processing.

VGGT: Visual Geometry Grounded Transformer

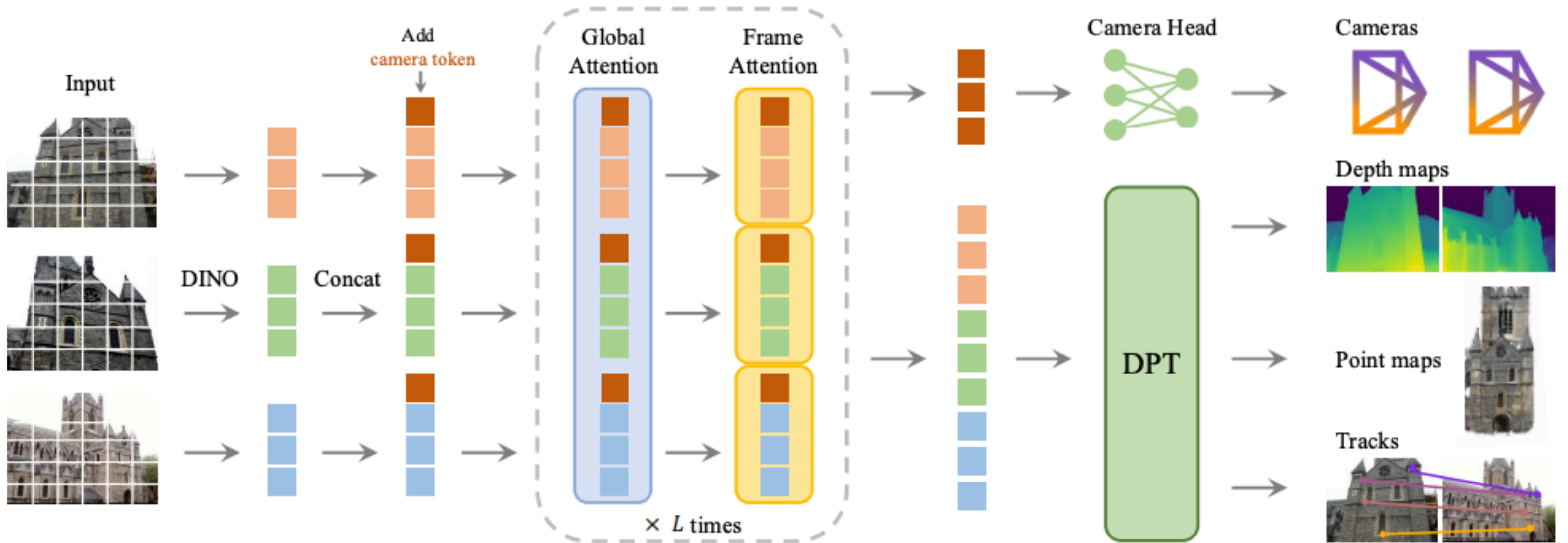


Figure 2. **Architecture Overview.** Our model first patchifies the input images into tokens by DINO, and appends camera tokens for camera prediction. It then alternates between frame-wise and global self attention layers. A camera head makes the final prediction for camera extrinsics and intrinsics, and a DPT [87] head for any dense output.

Depth Anything 3

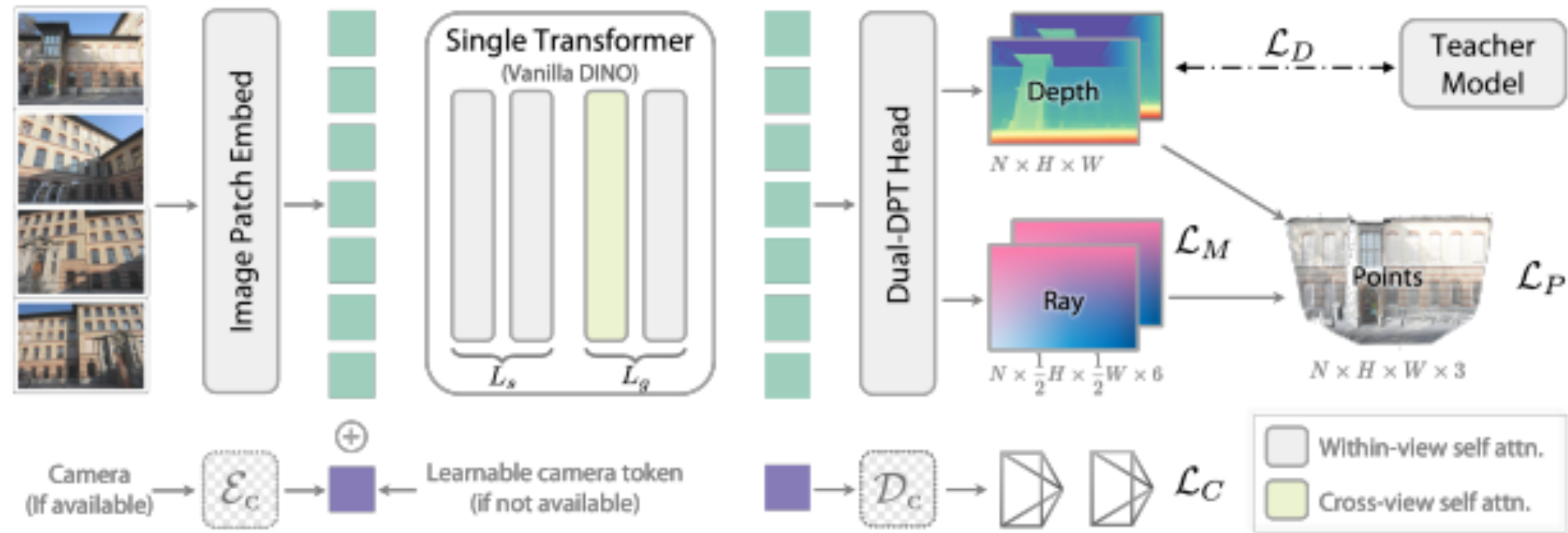


Figure 2 Pipeline of Depth Anything 3. Depth Anything 3 employs a single transformer (vanilla DINOv2 model) without any architectural modifications. To enable cross-view reasoning, an input-adaptive cross-view self-attention mechanism is introduced. A dual-DPT head is used to predict depth and ray maps from visual tokens. Camera parameters, if available, are encoded as camera tokens and concatenated with patch tokens, participating in all attention operations.

Depth Anything 3: Recovering the Visual Space from Any Views

Haotong Lin*, Sili Chen*, Jun Hao Liew*, Donny Y. Chen*, Zhenyu Li, Guang Shi,
Jiashi Feng, Bingyi Kang*,†

ByteDance Seed



Figure 1 Given any number of images and optional camera poses, **Depth Anything 3** reconstructs the visual space, producing consistent depth and ray maps that can be fused into accurate point clouds, resulting in high-fidelity 3D Gaussians and geometry. It significantly outperforms VGGT in multi-view geometry and pose accuracy; with monocular inputs, it also surpasses Depth Anything 2 while matching its detail and robustness.

Summary

- Interests in **3D recognition models** are rapidly growing as 3D sensors and datasets become more accessible.
- Extending CNNs to 3D is non-trivial because 3D data like **polygon meshes** or **point clouds** lack regular grid structure.
- Recent advances in image to 3D combining **classical geometric pipelines** for **multi-view reconstruction, synthetic data, transformer architectures**.
- Many techniques based on **view-based methods, sparse convolutions, graph neural networks, transformers**, etc.
- However, **3D data** is still lacking relative to **images**.

