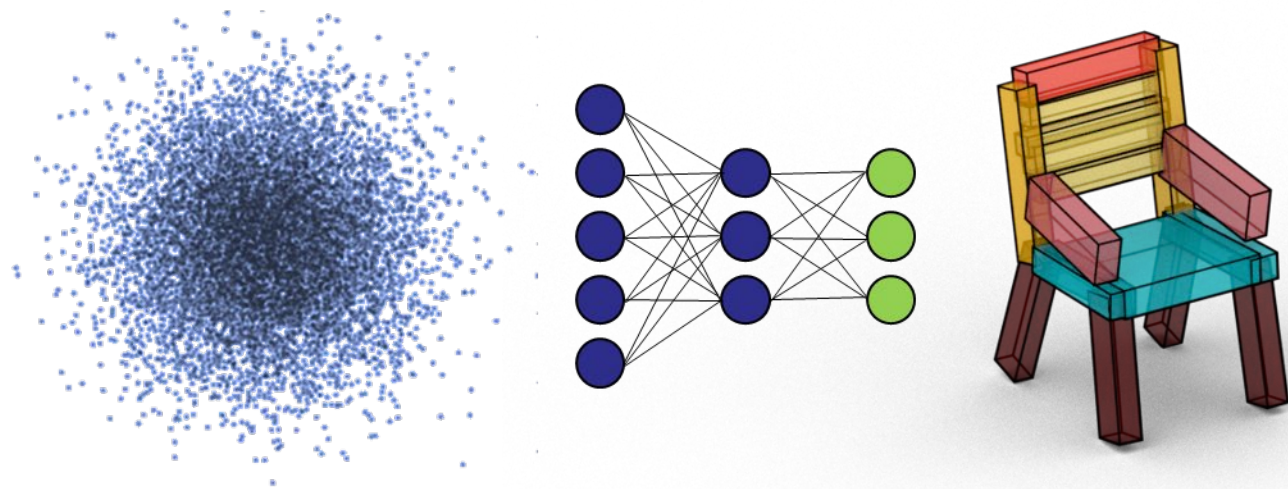


CS348n: Neural Representations and Generative Models for 3D Geometry

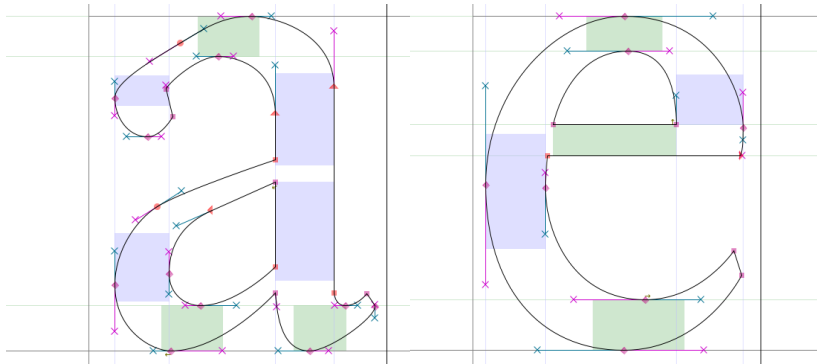


Leonidas Guibas
Computer Science Department
Stanford University

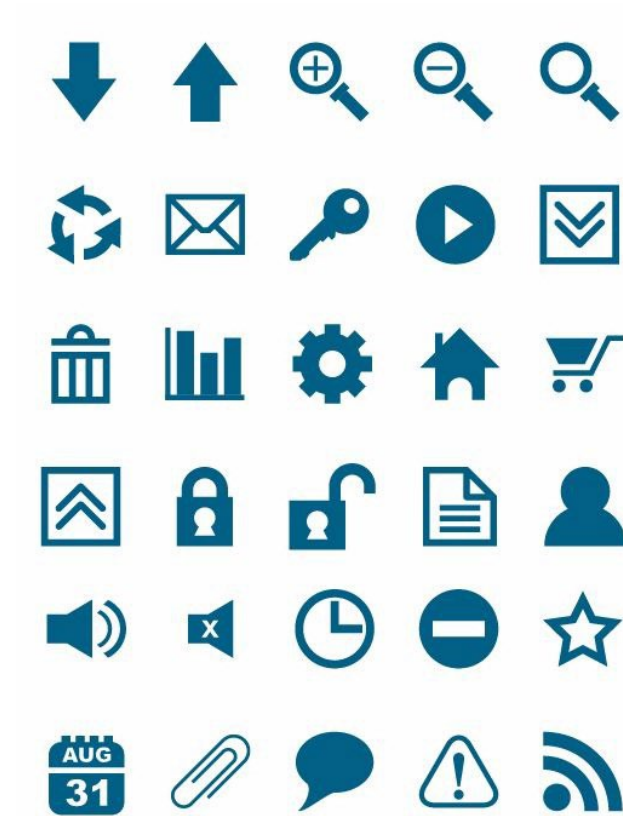


Last Time: Learning Vector and Mesh/CAD Representations

Vector (2D) and Mesh/CAD (3) Shape Models



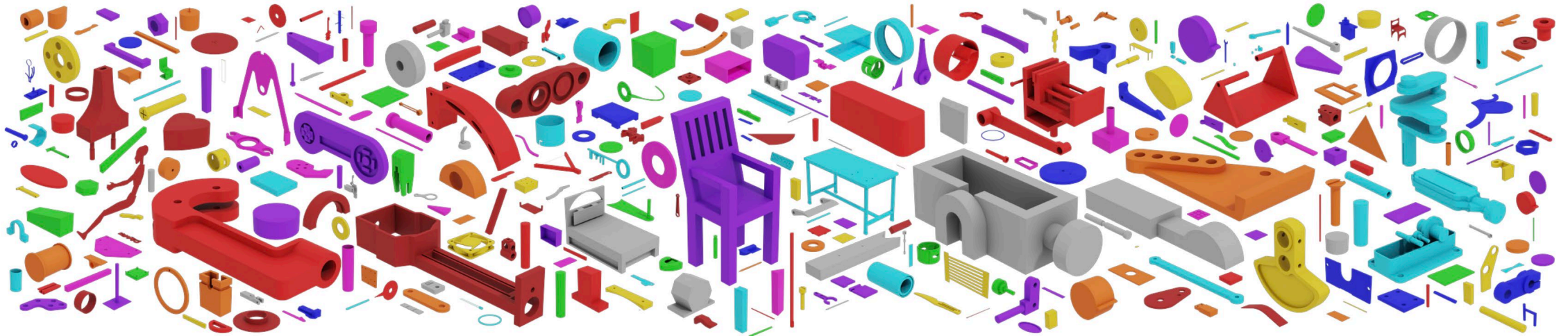
Fonts



Clip Art

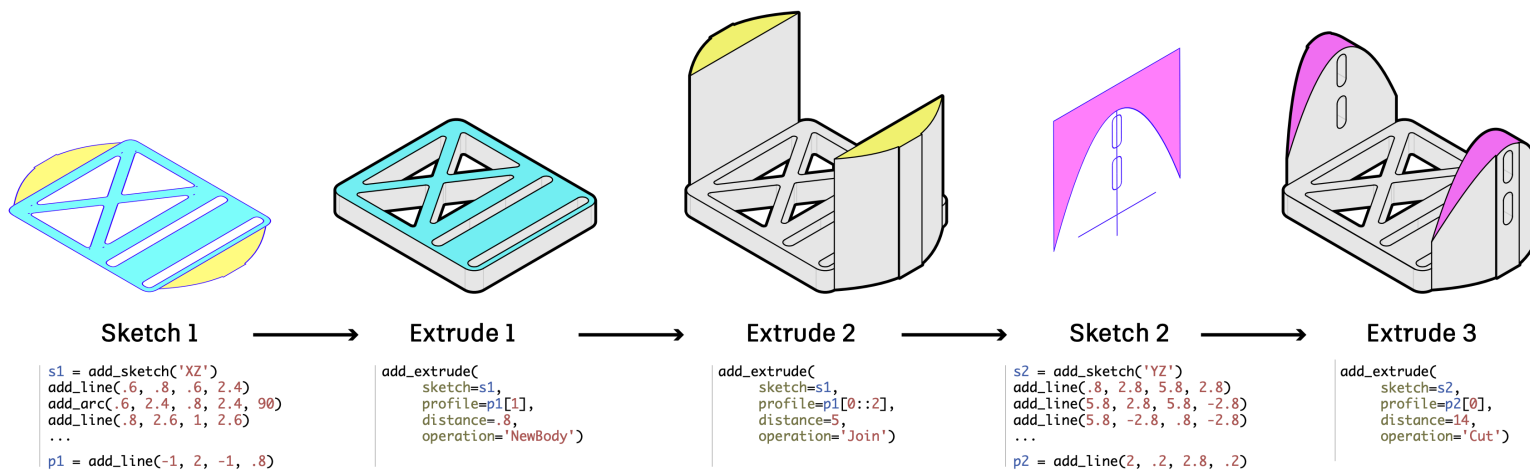
2D

Vector (2D) and Mesh/CAD (3D) Shape Models



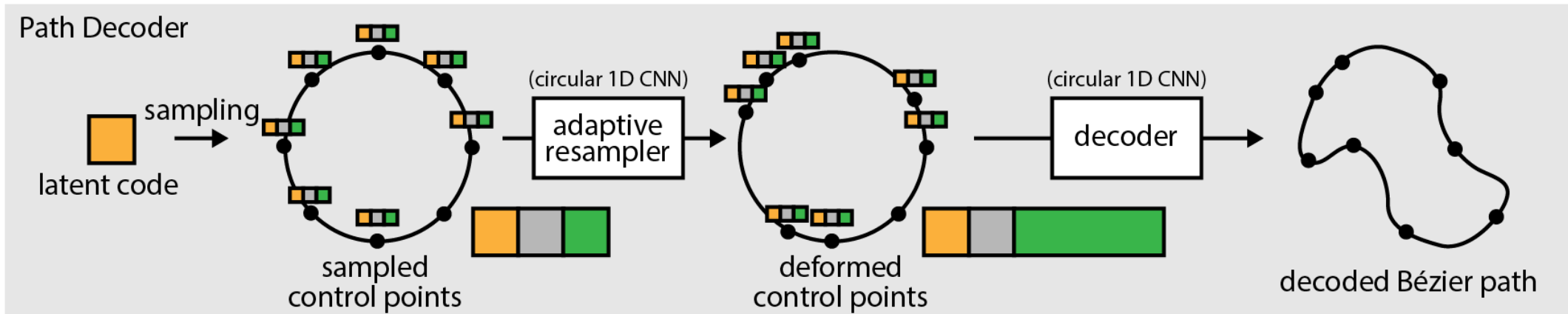
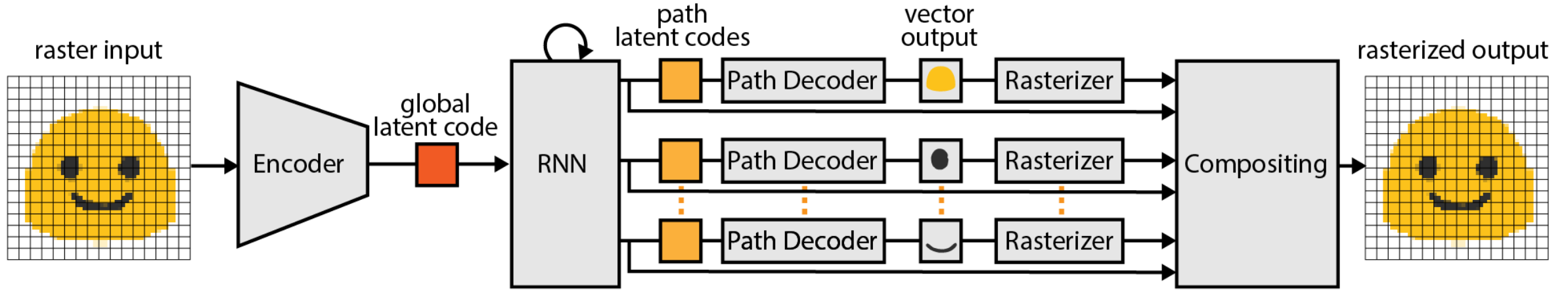
Autodesk Fusion 360 Dataset

3D



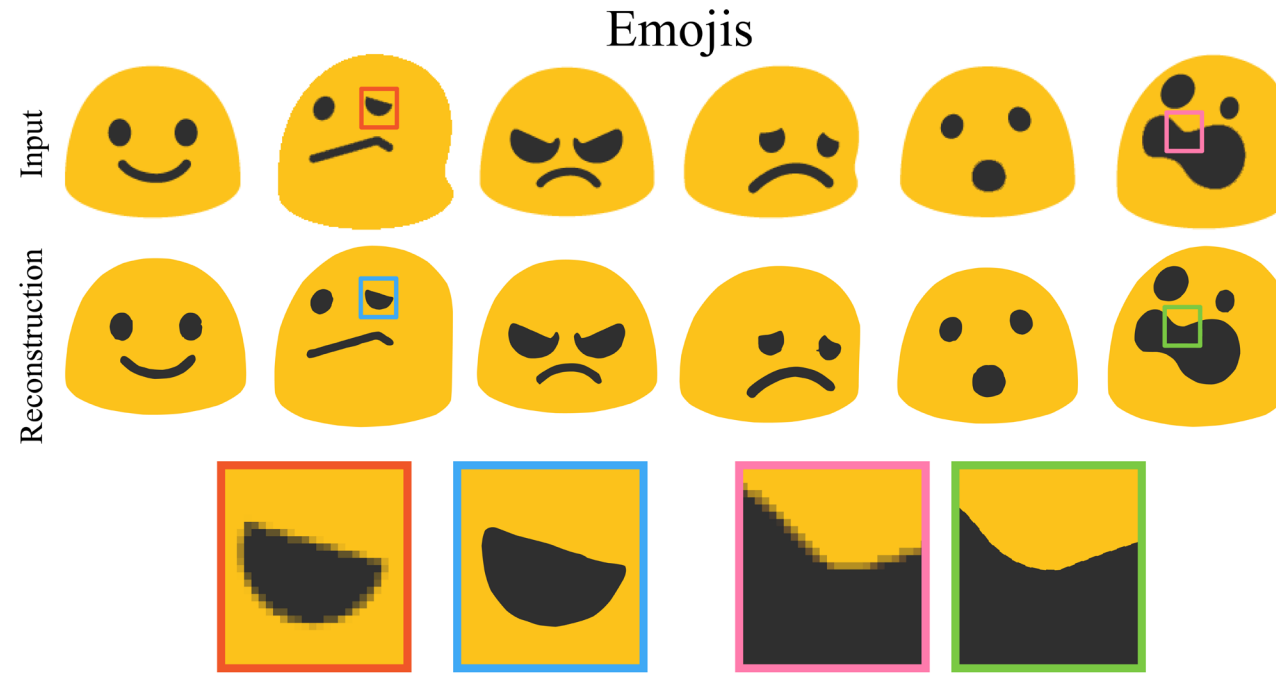
Im2Vec Inference Architecture

4 step process.



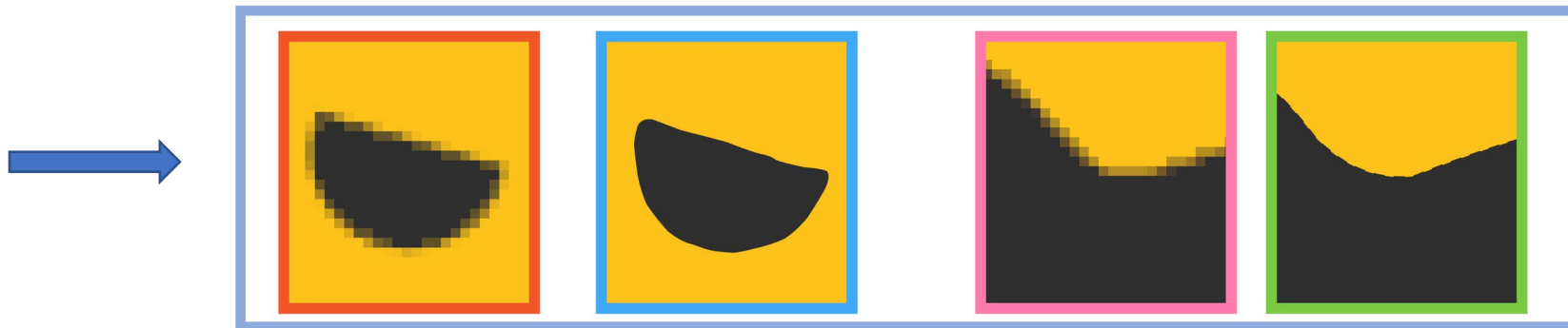
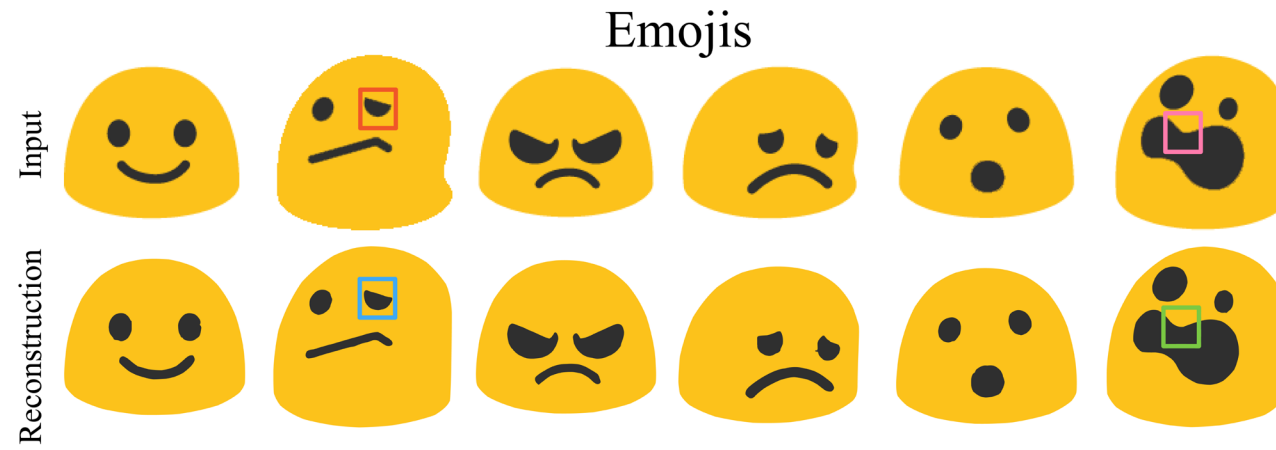
Reconstruction

Input images are 128x128



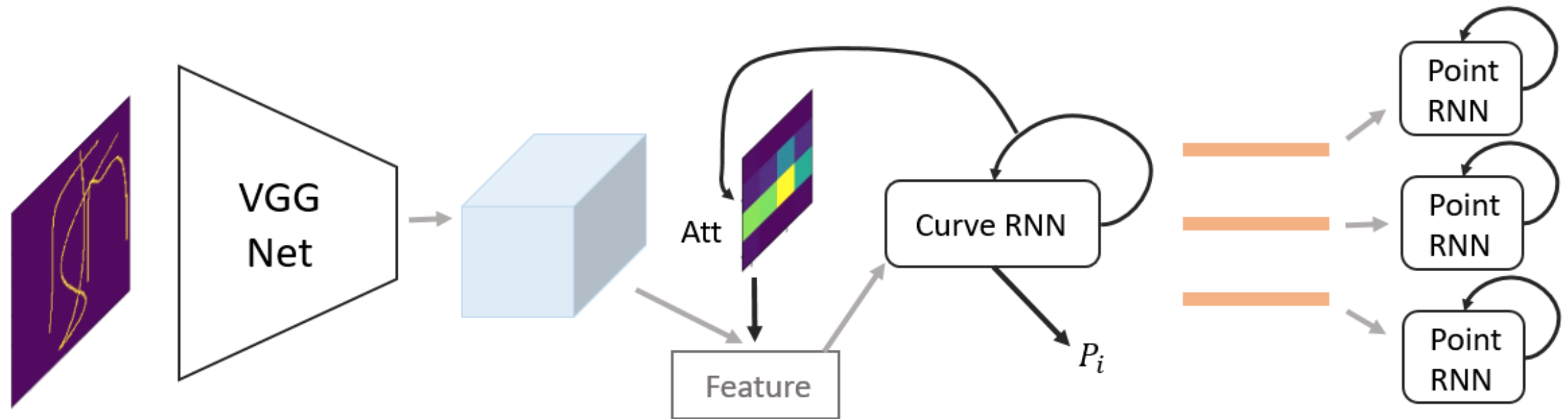
Reconstruction

Input images are 128x128

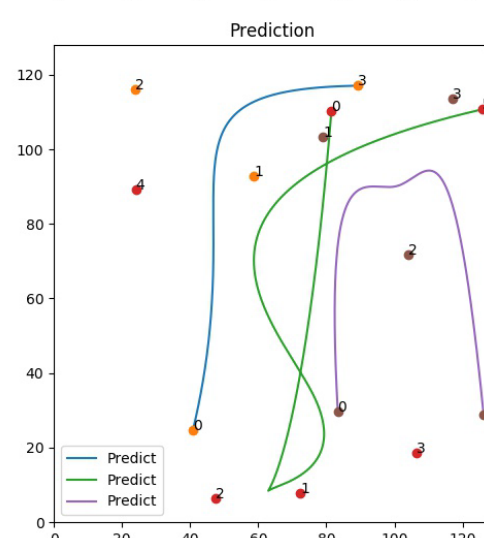
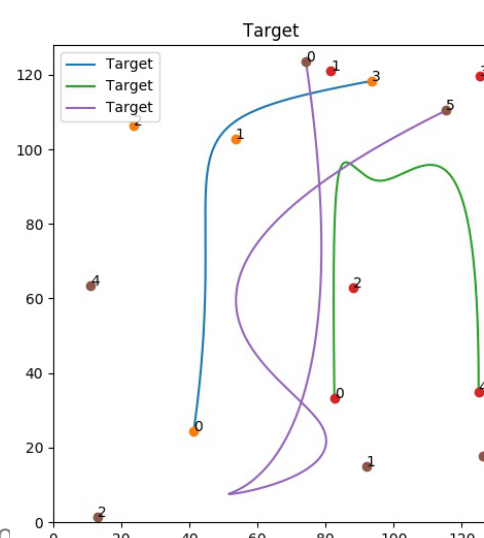
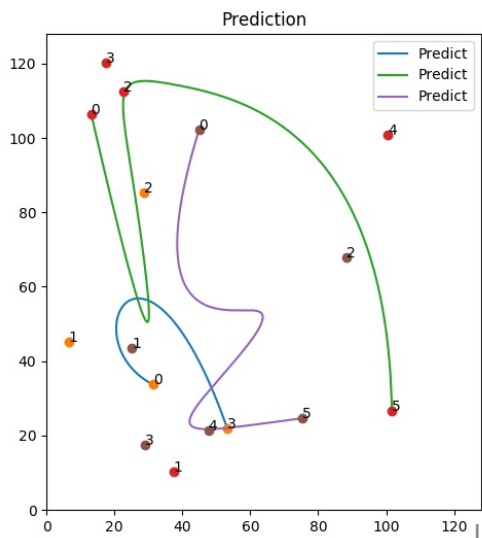
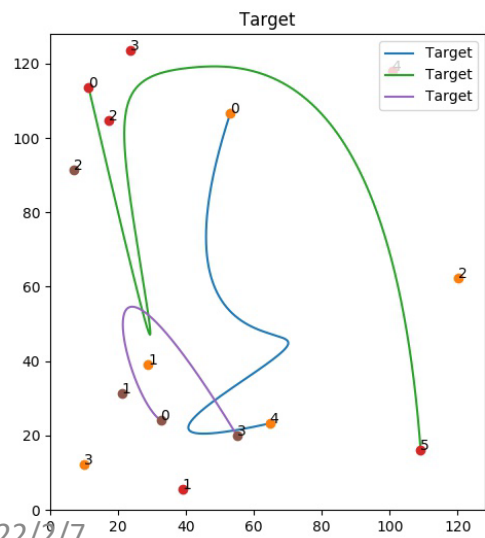
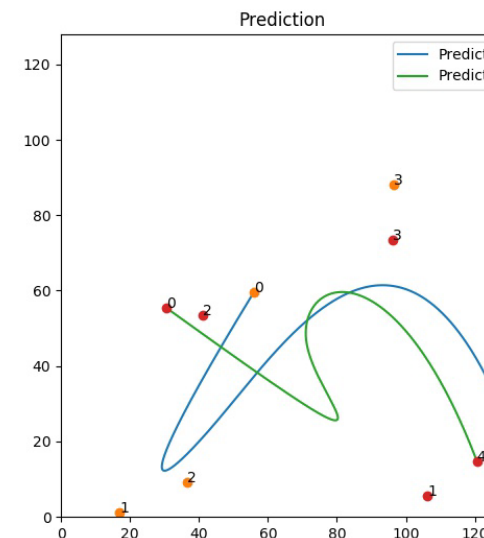
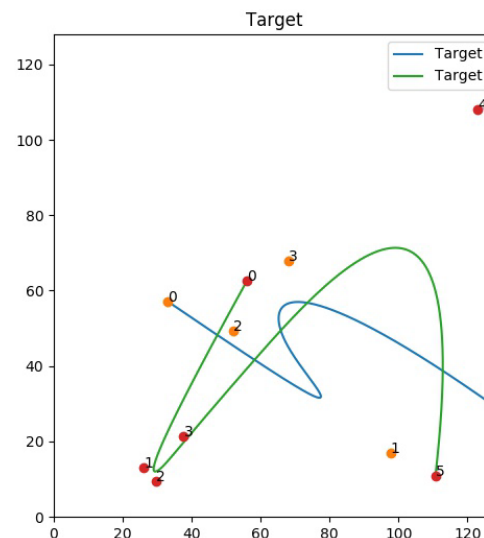
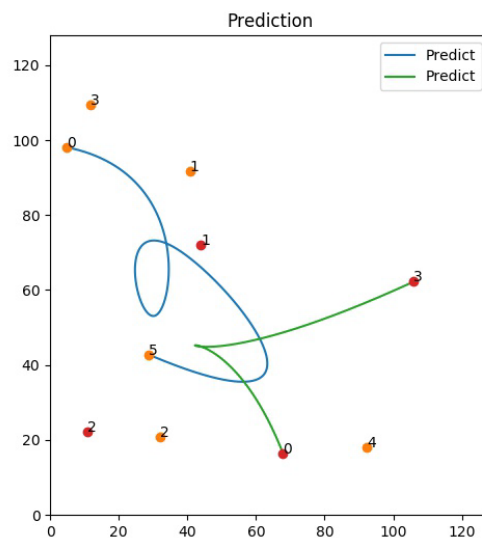
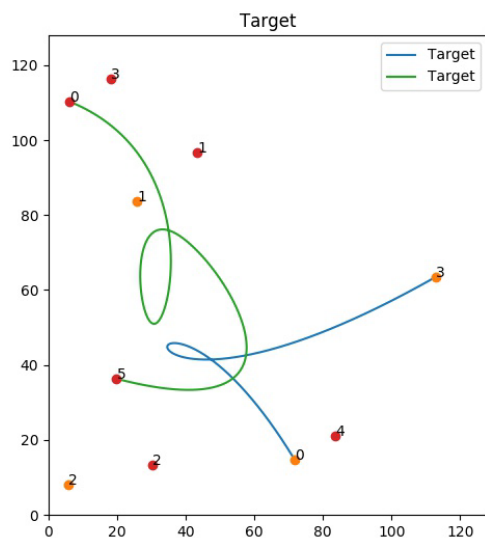


Multiple spline curves with variable numbers of control points

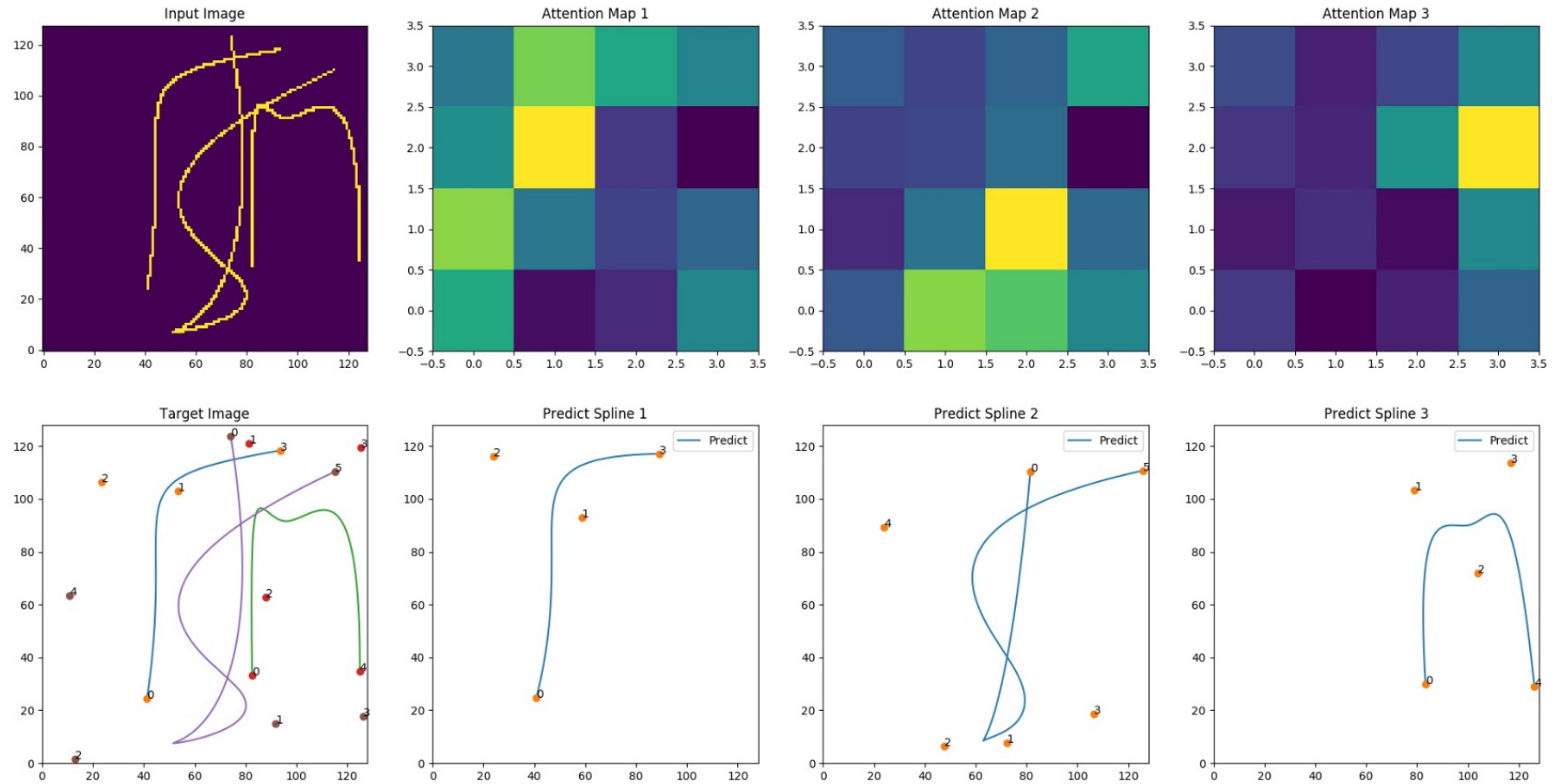
- Whole model

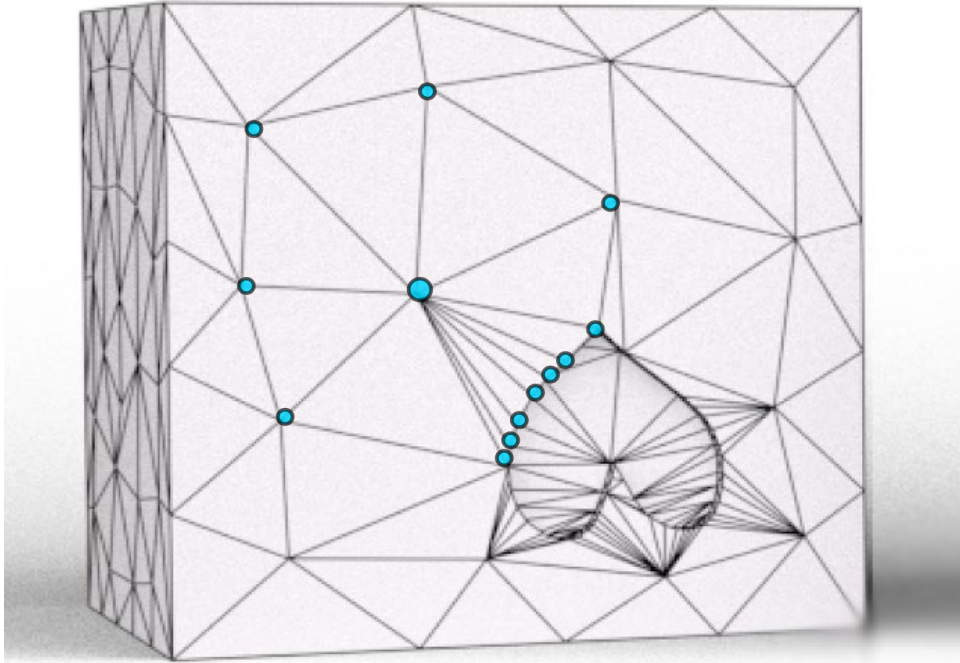


Results

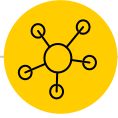


Attention Results





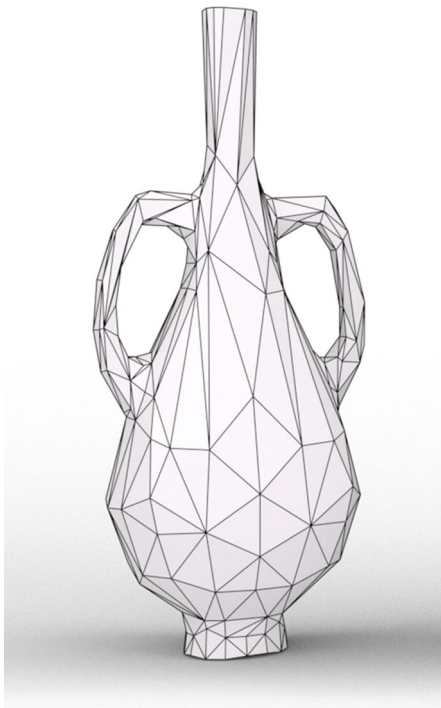
Irregular Unordered & Unoriented



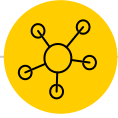
Goal: CNN directly on the irregular mesh elements

Classification

Segmentation

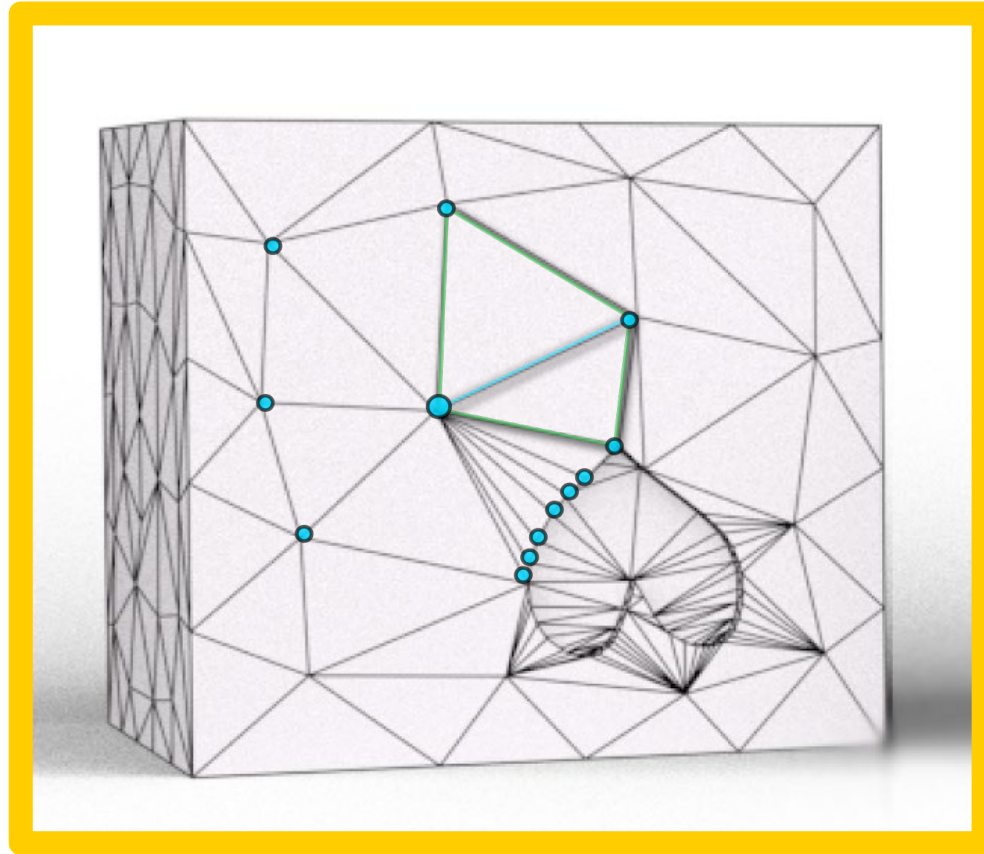


Vase



Fixed Size

Neighborhood



Vertices

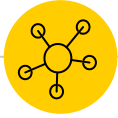
$\langle x, y, z \rangle$

Edges

$\langle v_i, v_j \rangle$

Faces

$\langle v_i, v_j, v_k \rangle$

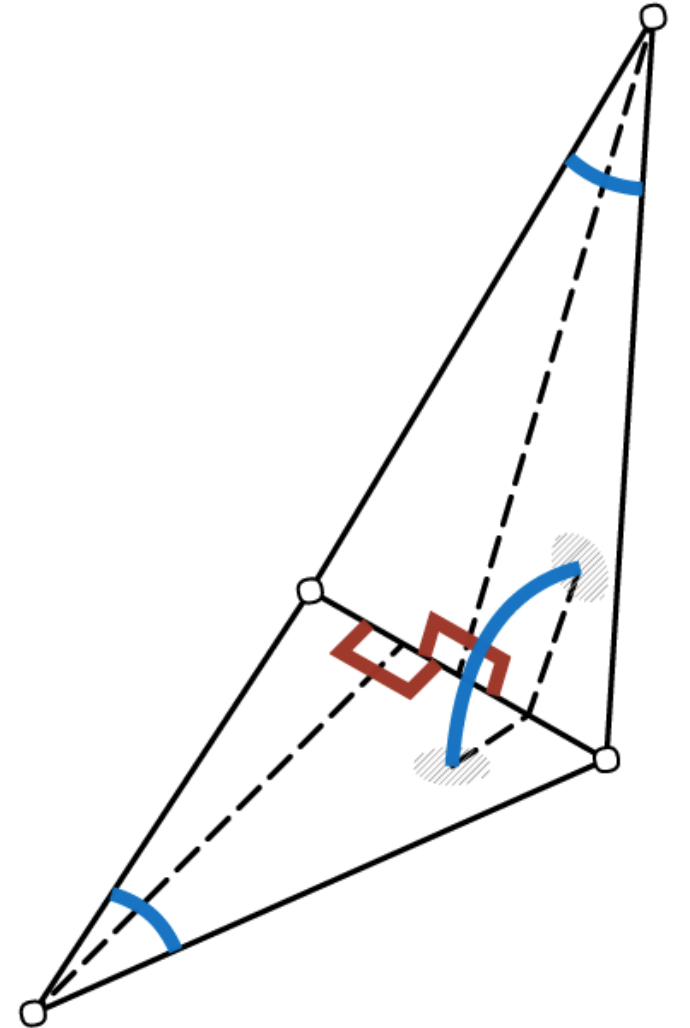


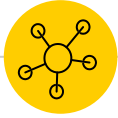
Input **Edge** Features

Relative Geometric Features

→ Invariant to *similarity* transformations

5-dimensional vector





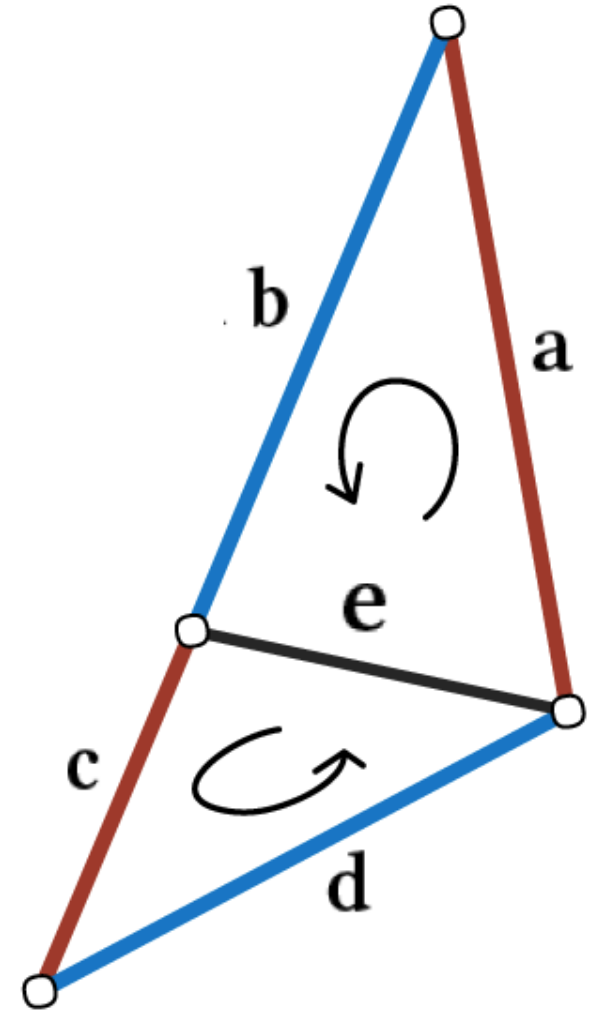
Mesh Convolution Order

Face normal

- Consistent ordering in each face
- Two *valid* orderings

Solution: build symmetric features

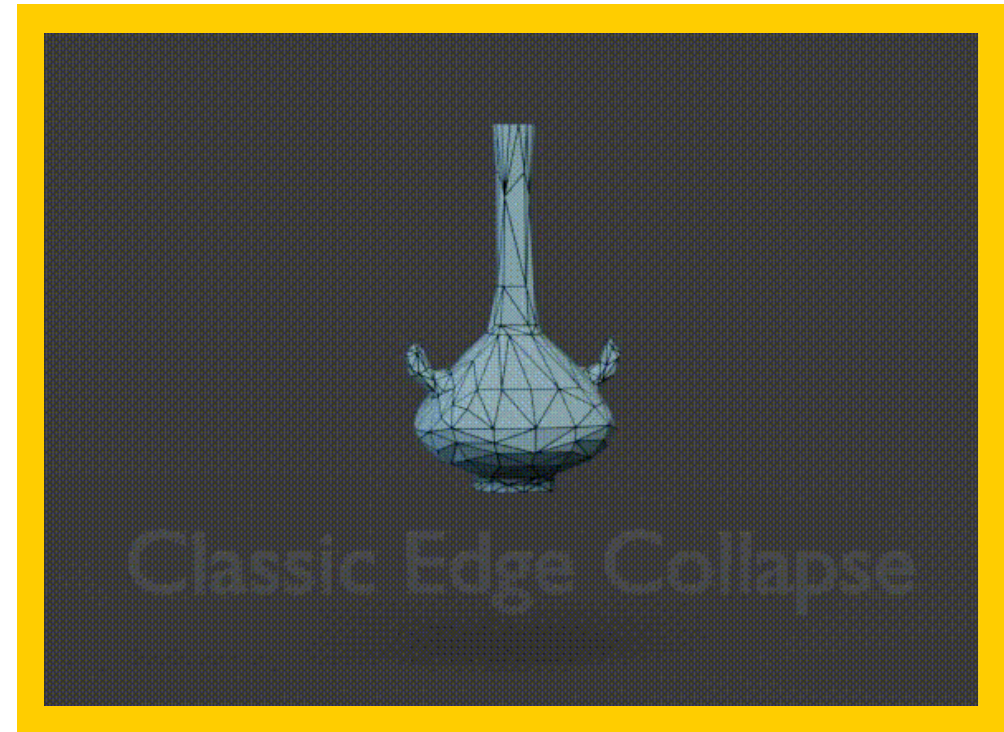
- $e \rightarrow (a+c, |a-c|, b+d, |b-d|)$

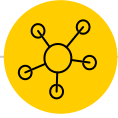




Learned **Edge** Collapse

- Network decides collapse
- Strengthens the learned representation
- Visual insights from network



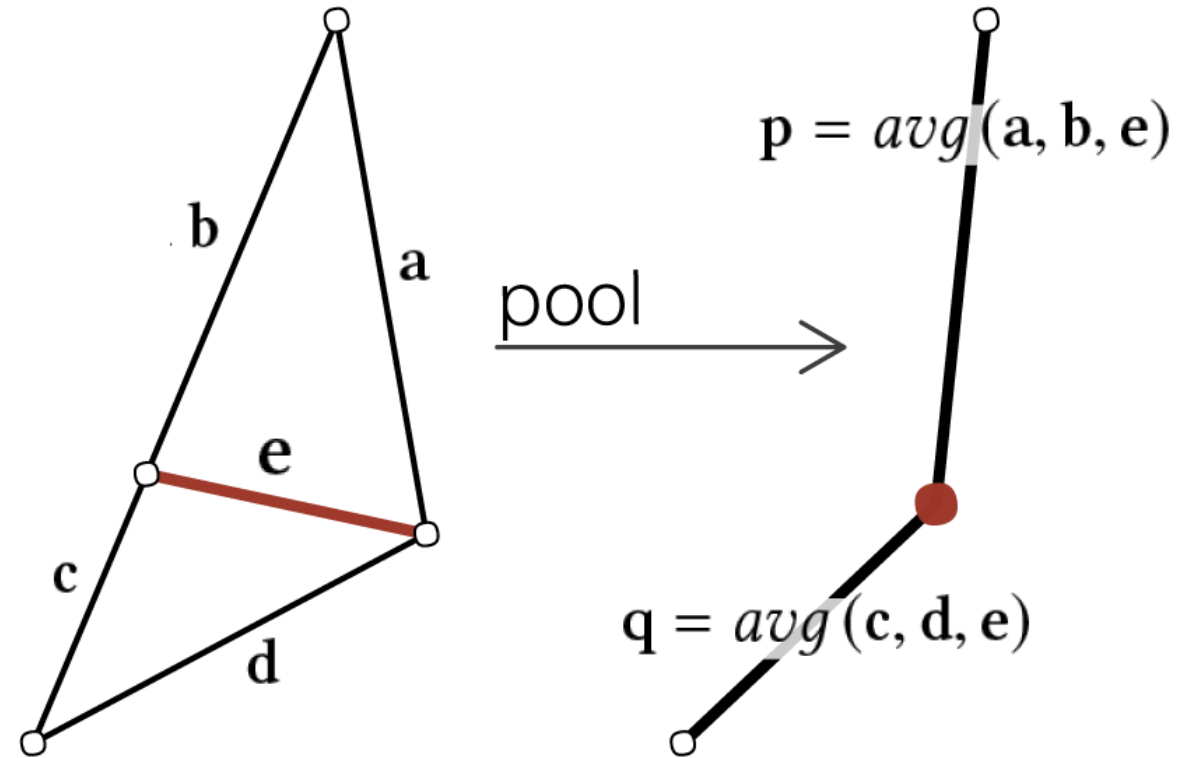


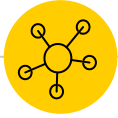
Mesh Pooling

Delete edge with smallest feature activations

→ Aggregate features

→ Update topology

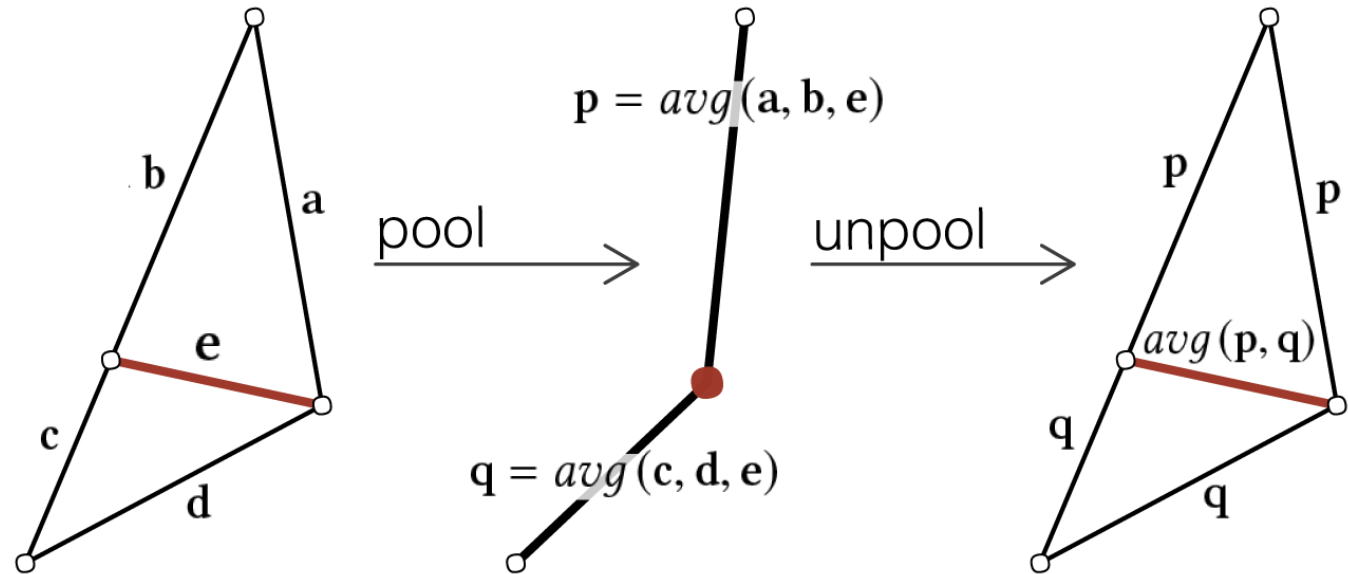




Mesh Unpooling

Partial Inverse of Pooling

- Restores upsampled topology (reversible)
- Unpooled features are a weighted combination of pooled features

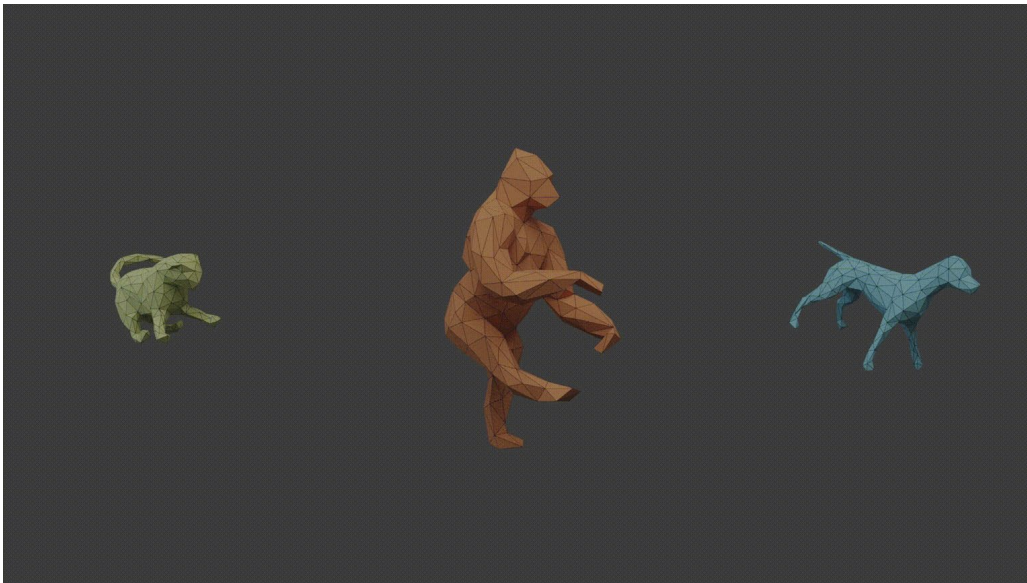




Applications of MeshCNN

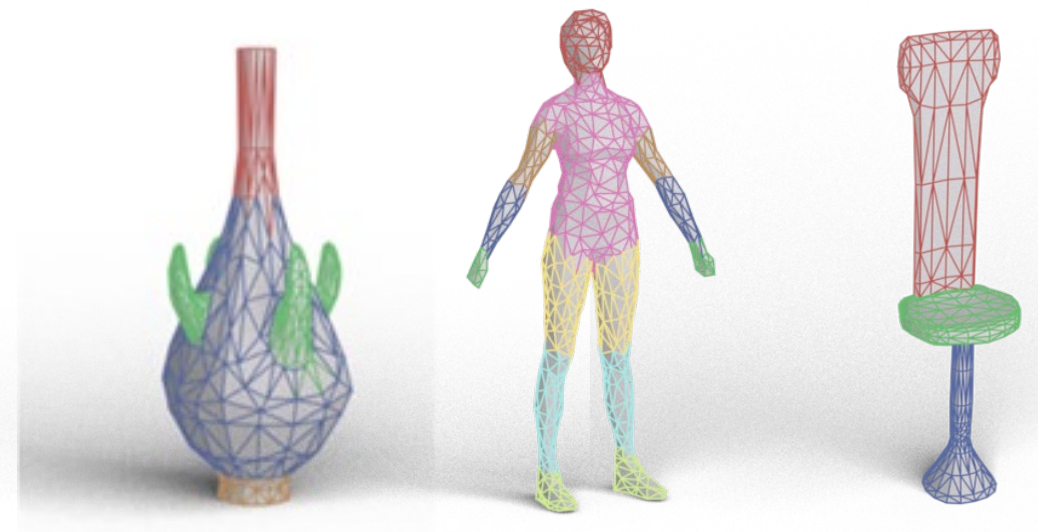
Classification

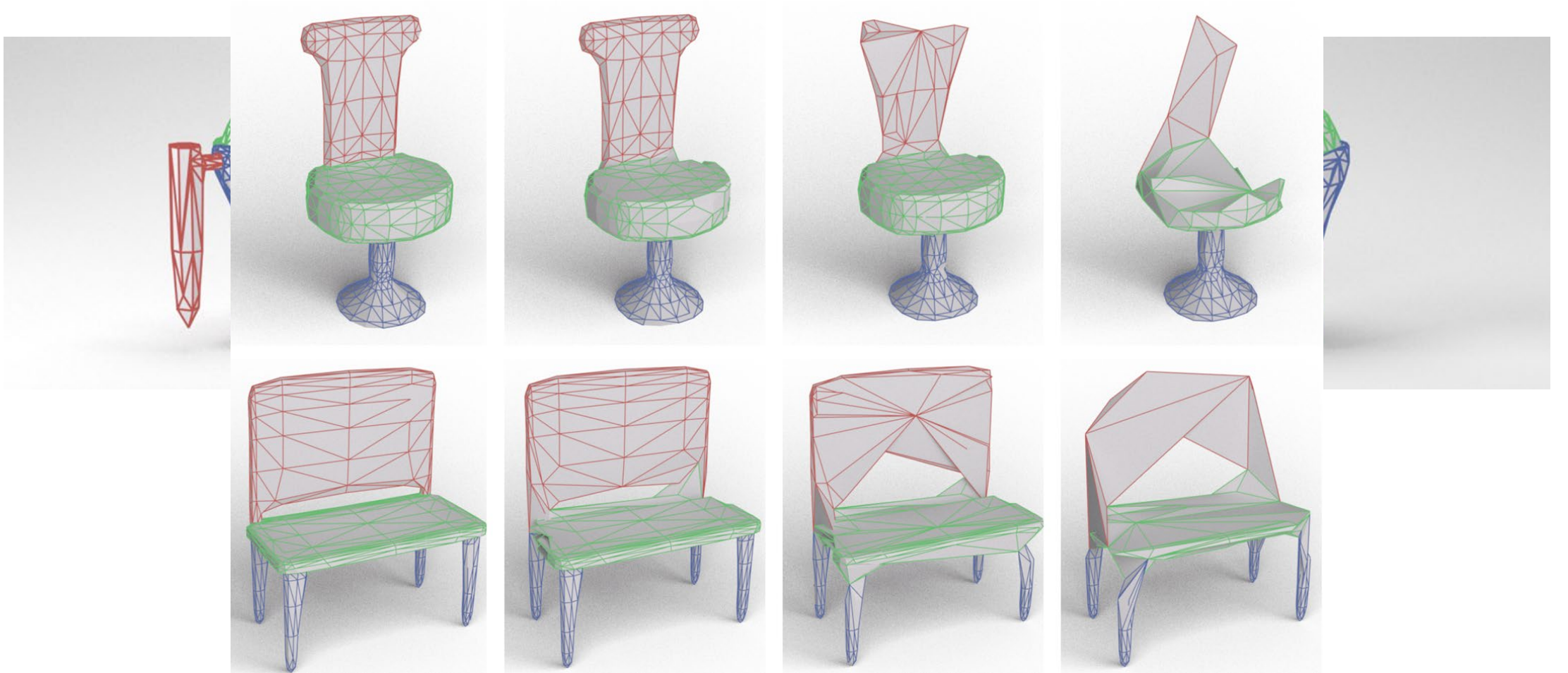
- Conv & Pooling Layers
- Fully-Connected Layers



Segmentation

- Fully convolutional
- Conv & Pooling & Unpooling





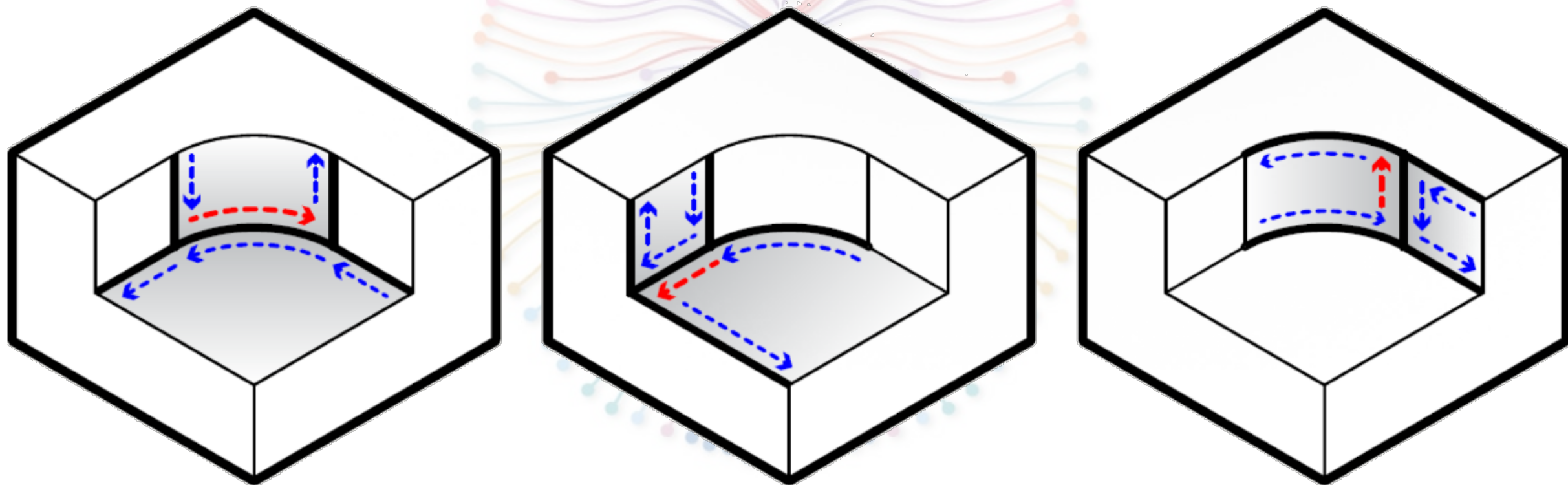
Segmentation



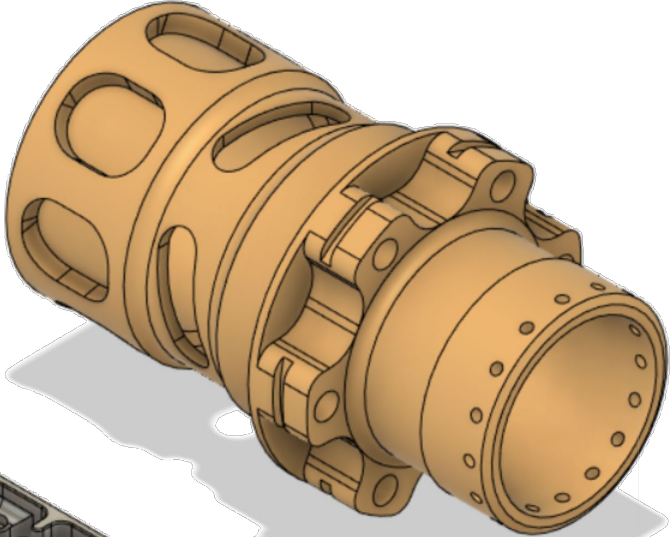
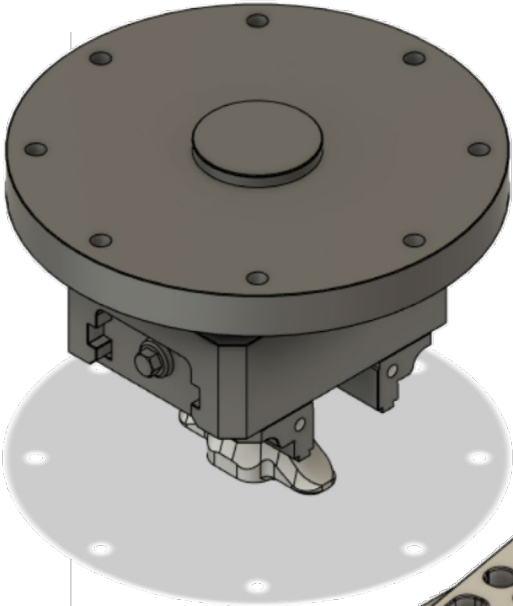
BRepNet: A topological message passing system for solid models

Joseph G. Lambourne¹, Karl D.D. Willis¹, Pradeep Kumar Jayaraman¹, Aditya Sanghi¹,
Peter Meltzer², Hooman Shayani¹

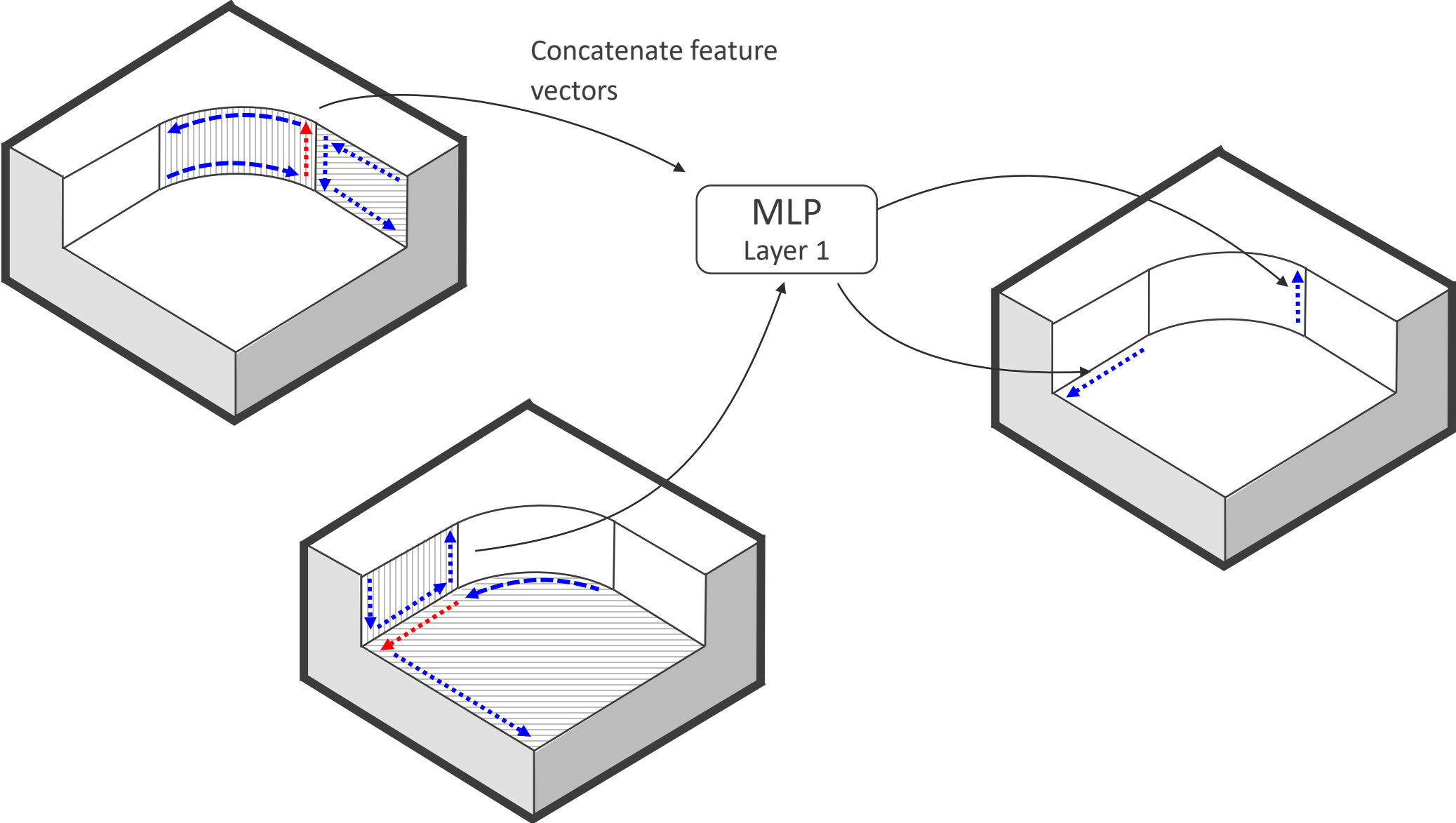
Autodesk Research¹, UCL, Computer Science²



Boundary representation models

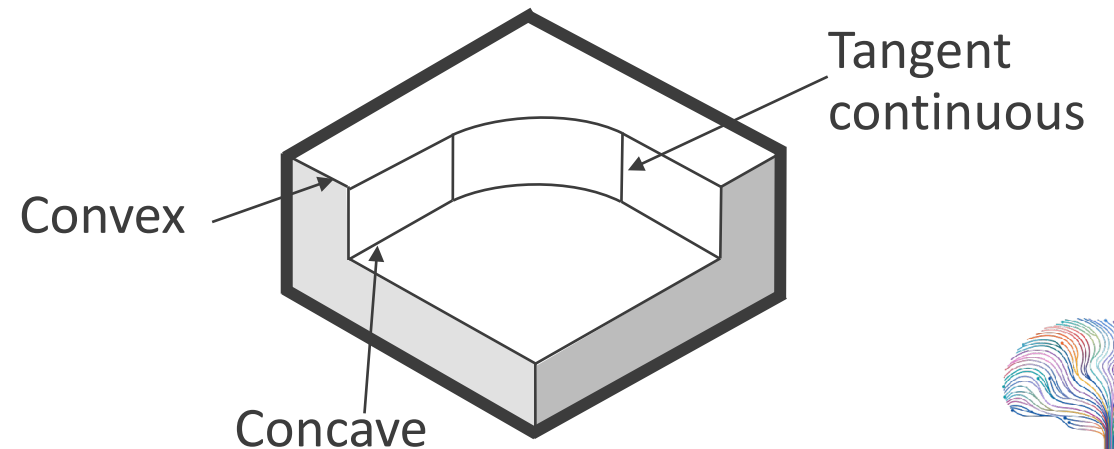
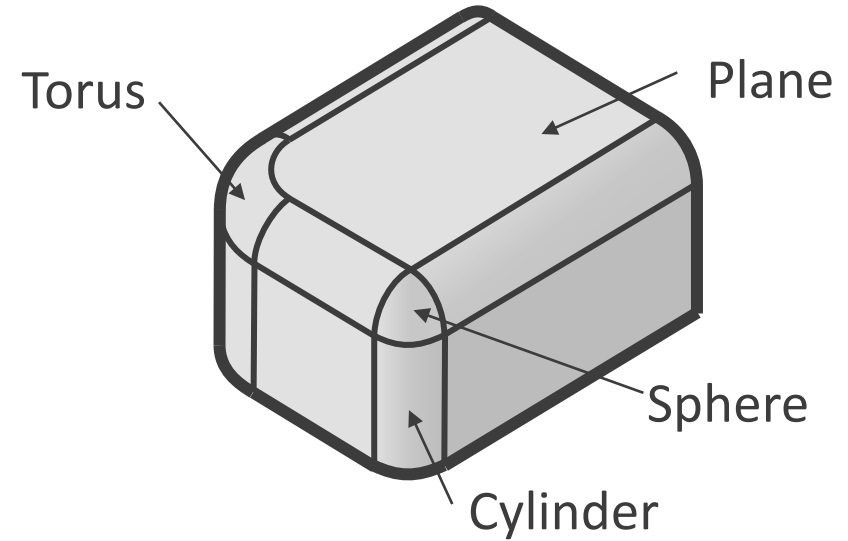


BRepNet convolution



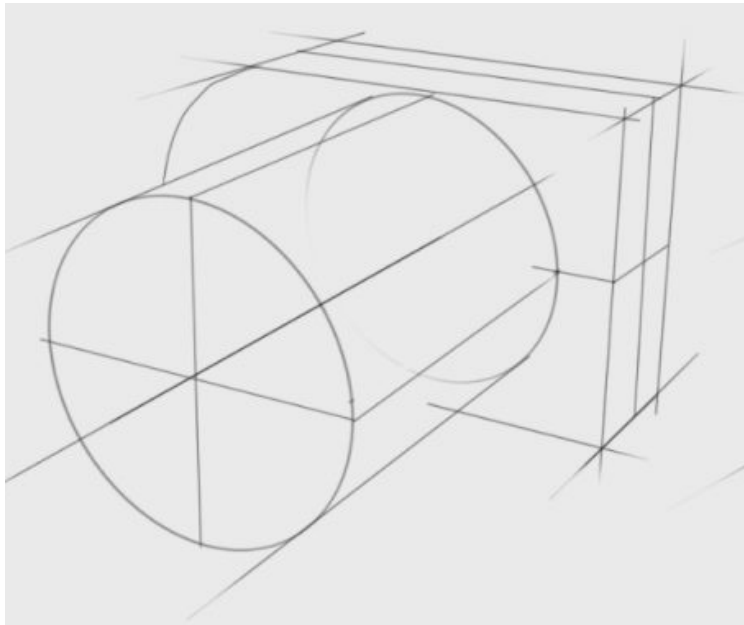
Input features

- Face features
 - Surface type (plane, cylinder, cone, sphere, torus, spline)
 - Face area
- Edge features
 - Curve type (line, circle, ellipse, intersection curve, spline)
 - Edge convexity (concave, convex, tangent plane continuous)
 - Edge length

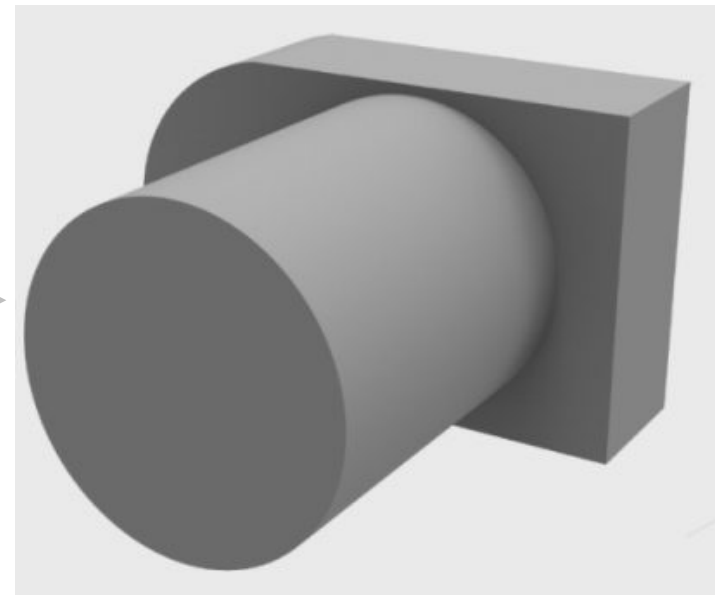
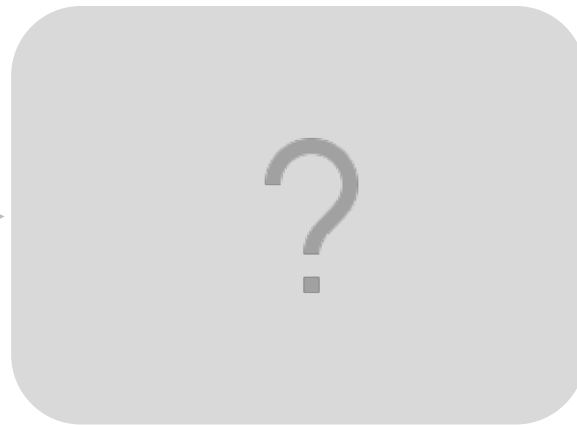


Sketch2CAD Goal

- Gap: converting sketches to a CAD model?



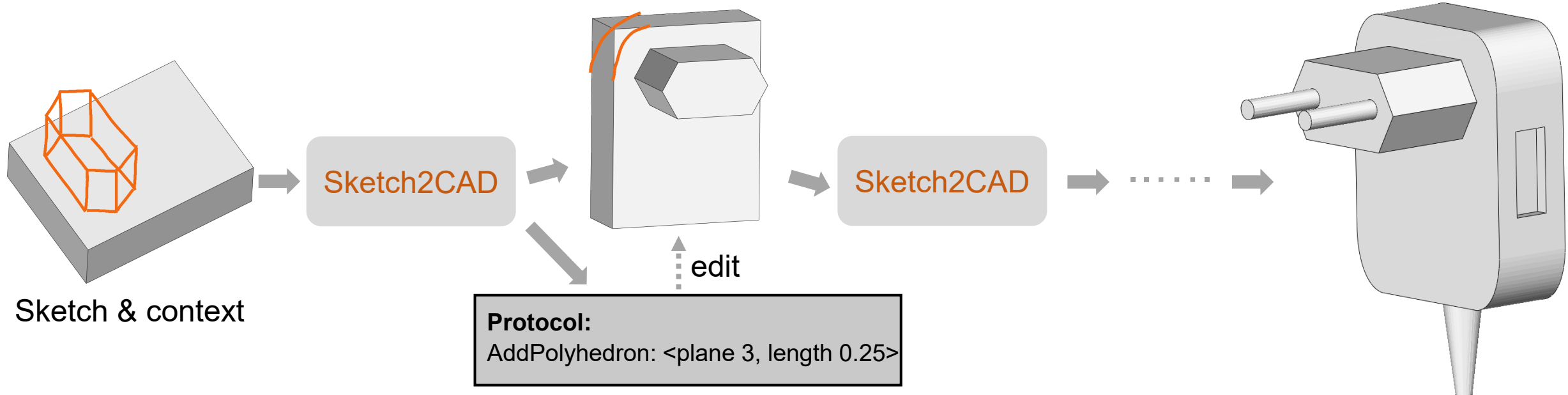
User sketches



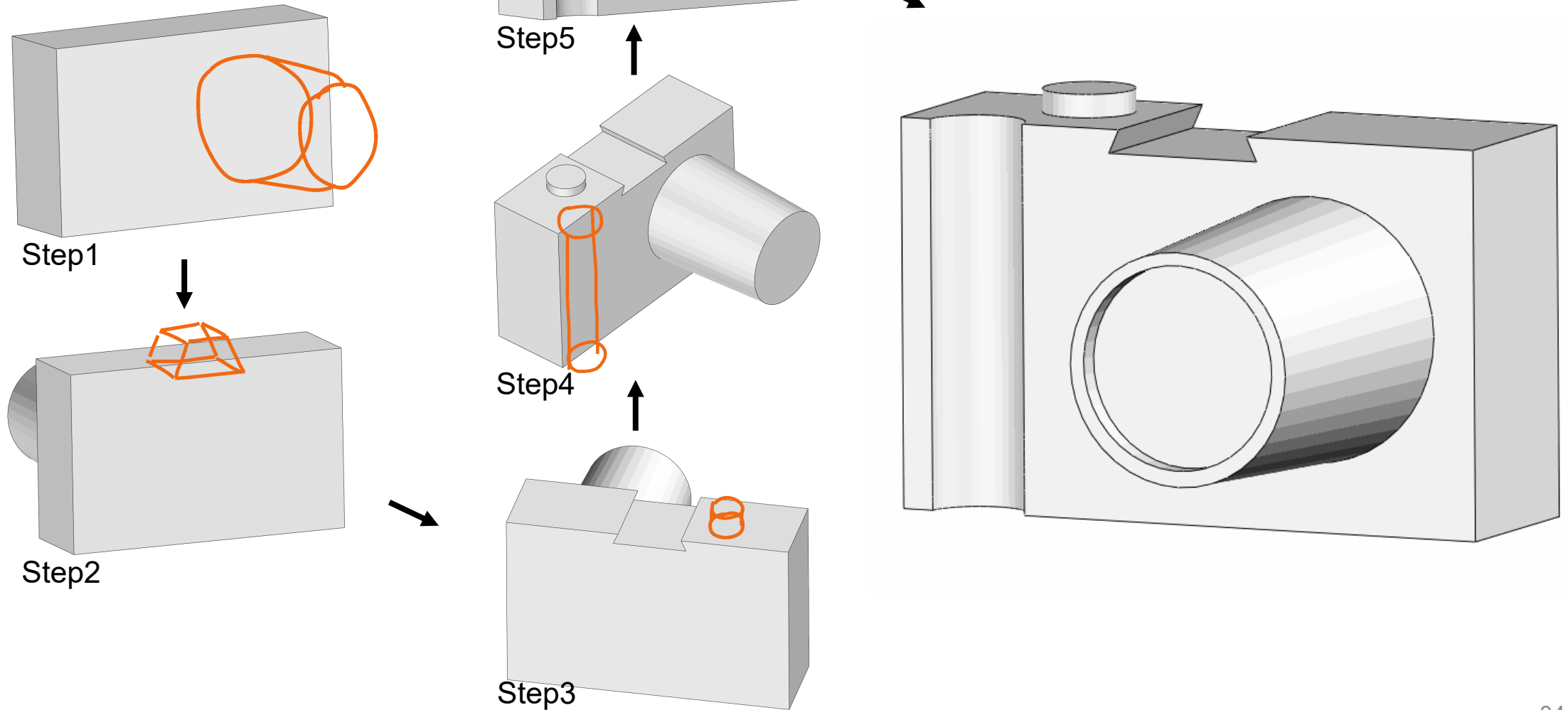
3D model

Sketch2CAD

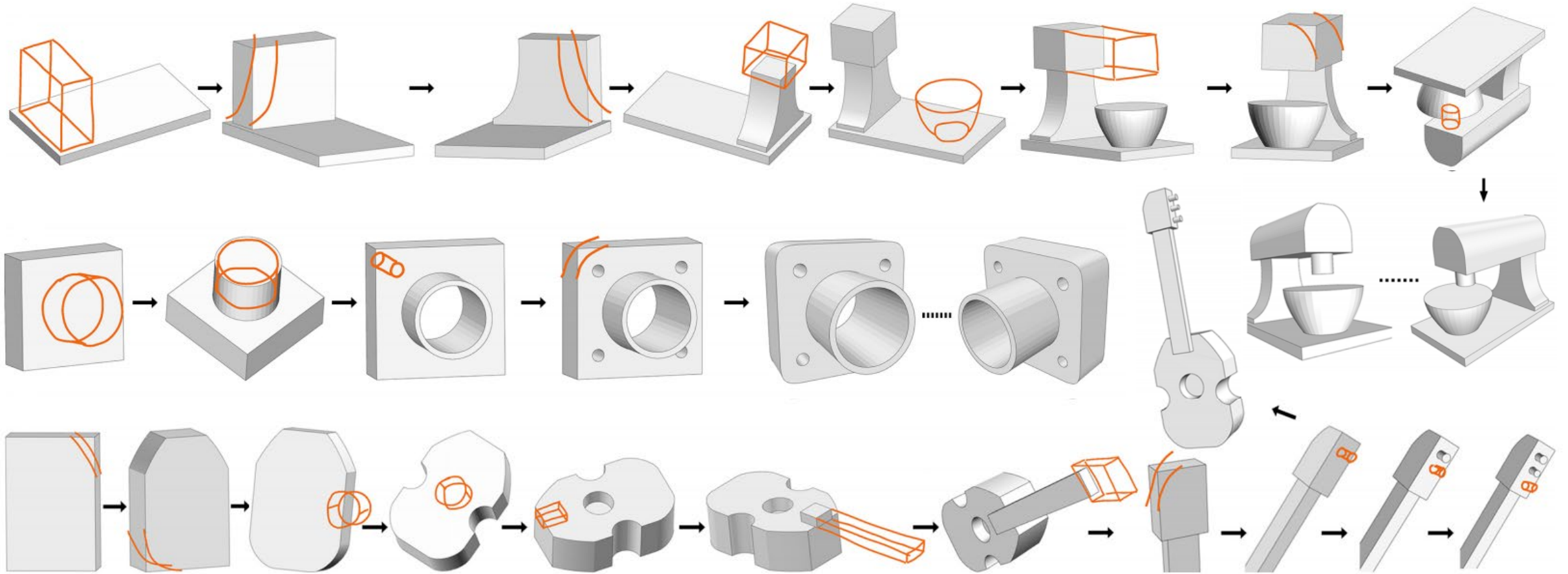
- A learning-based modeling system that translates **sketching operations** to their corresponding **CAD operations**, along with the associated **parameters**.



Results



Results



3D Equivariance and Invariance

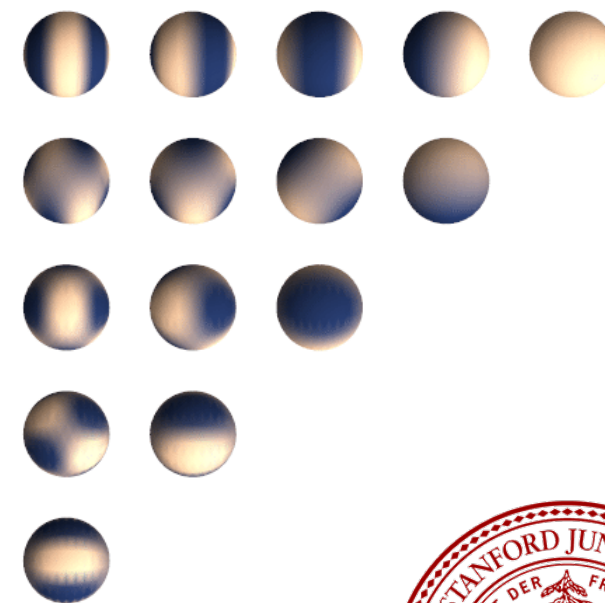
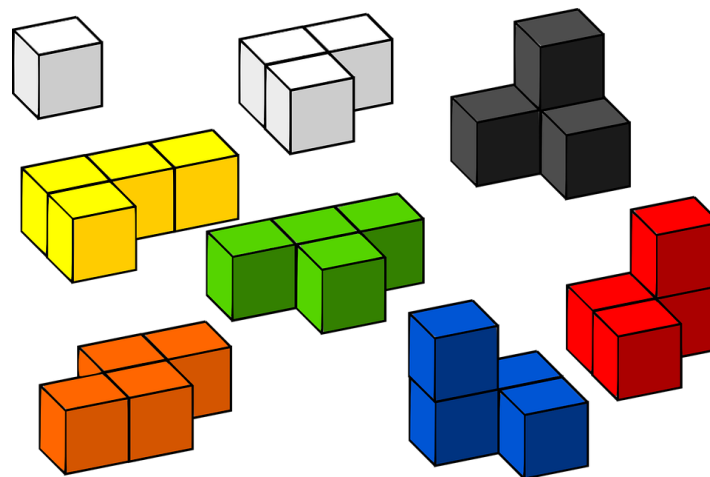
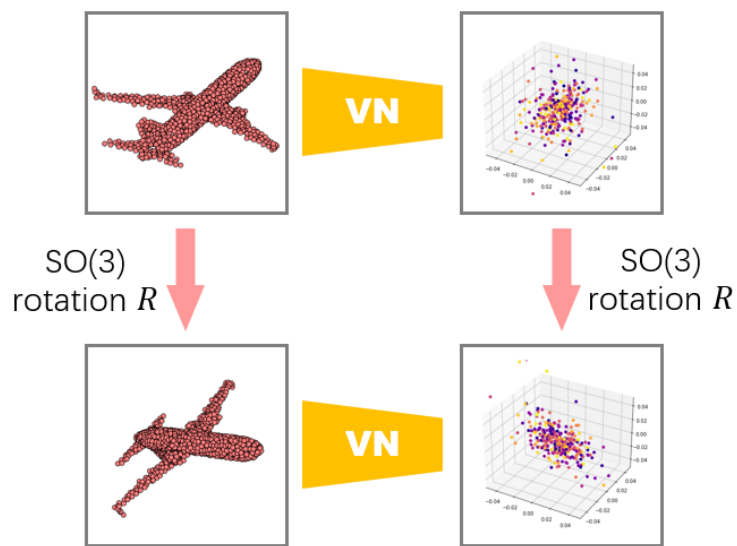
Class Questionnaire

Class Questionnaire

- Please take the questionnaire below to provide us with feedback on the class:
- <https://forms.gle/igFFpmnWaWL11Tfw9>

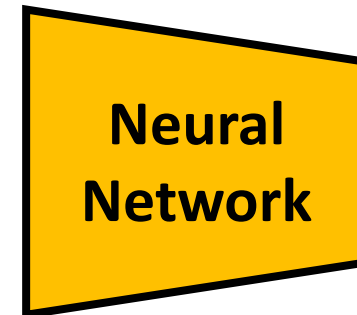
CS348N Lecture 10:

3D Rotation Equivariance and Invariance



Why Equivariance / Invariance?

- What neural networks can do?



“chair”

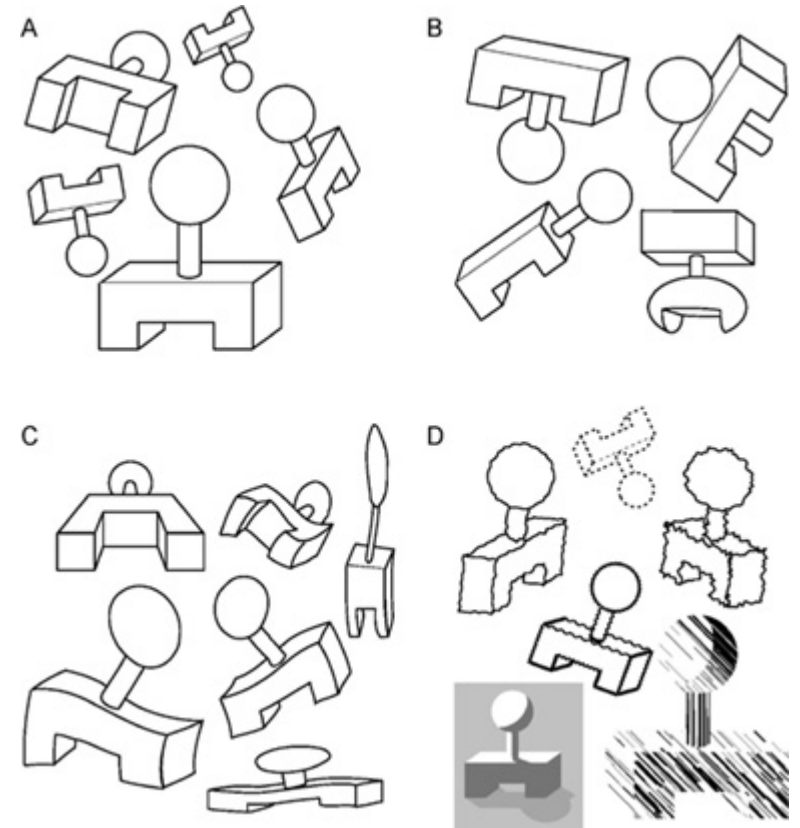
- But what is this?



Not too Hard for a Human...

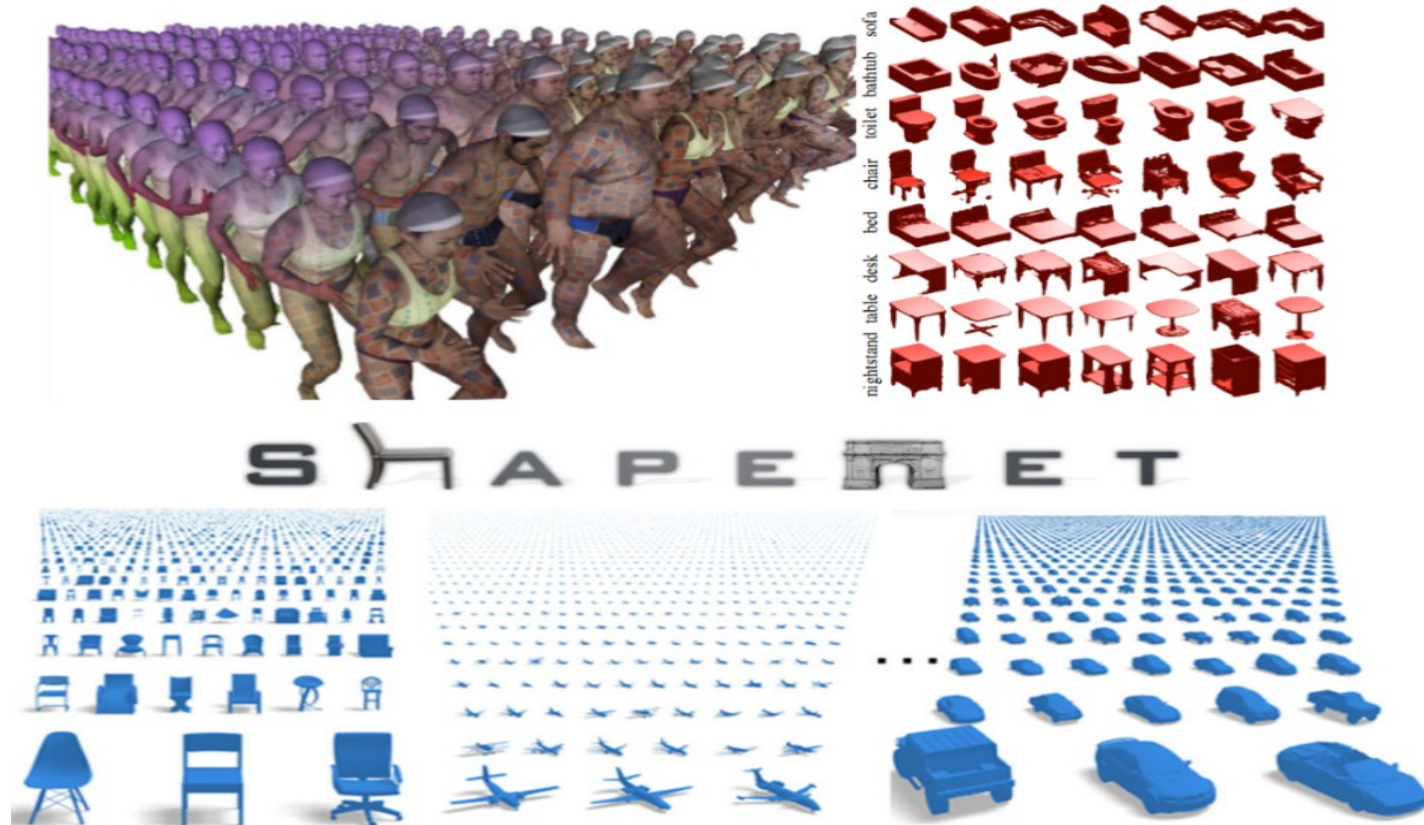
Gestalt theory – how human perceive objects

- Emergence
- Reification
- Multistability
- Invariance



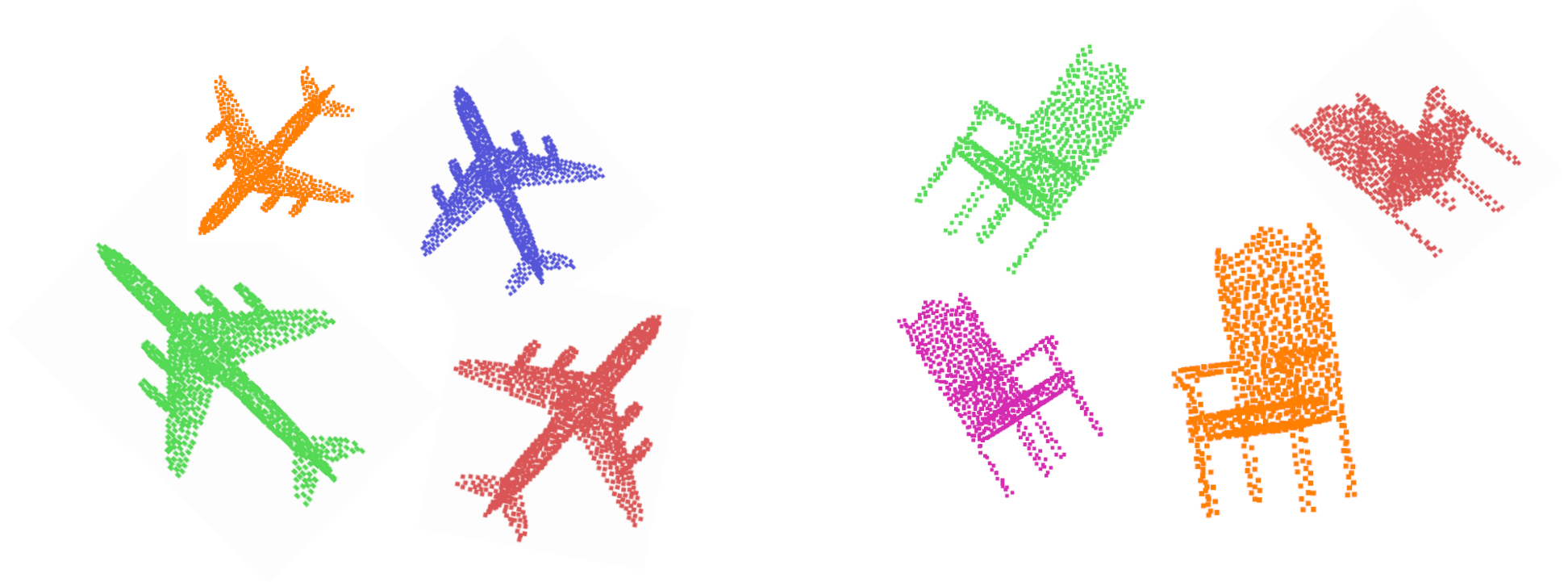
But Neural Networks Struggle...

- Many datasets are **aligned** (shapes under canonical poses)
- Networks trained aligned data **cannot generalize** to arbitrary poses



But Neural Networks Struggle...

Identical objects coming in different poses, scales, ratios, colors...
– viewed as totally irrelevant entities by classical neural networks

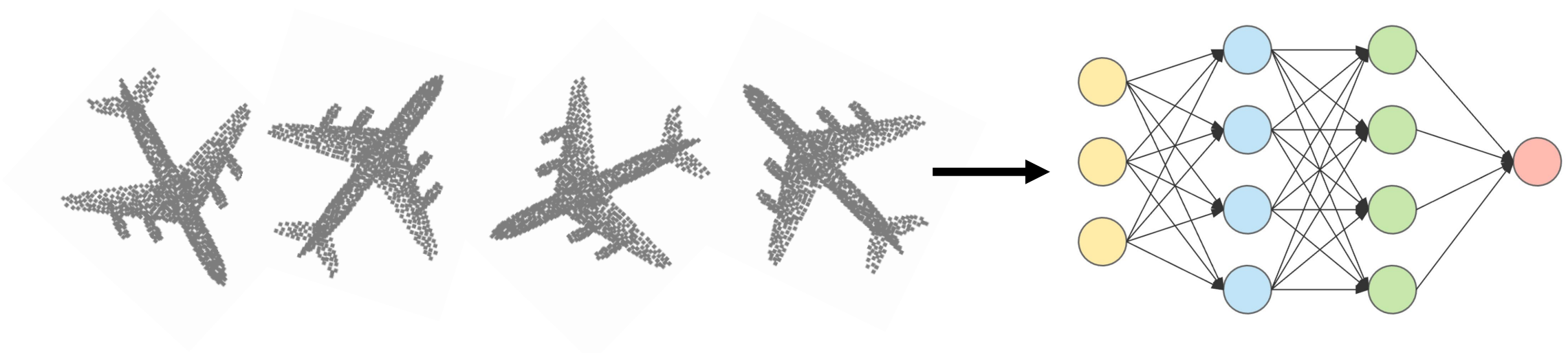


A Naïve Solution: Data Augmentation

Apply random rotation to the training data

So we let the network “see” and learn from all possible poses

- Reducing the generalization gap – but not eliminating it
- Sacrificing data-efficiency – longer training time
- Statistically equivariant/invariant – not guaranteed

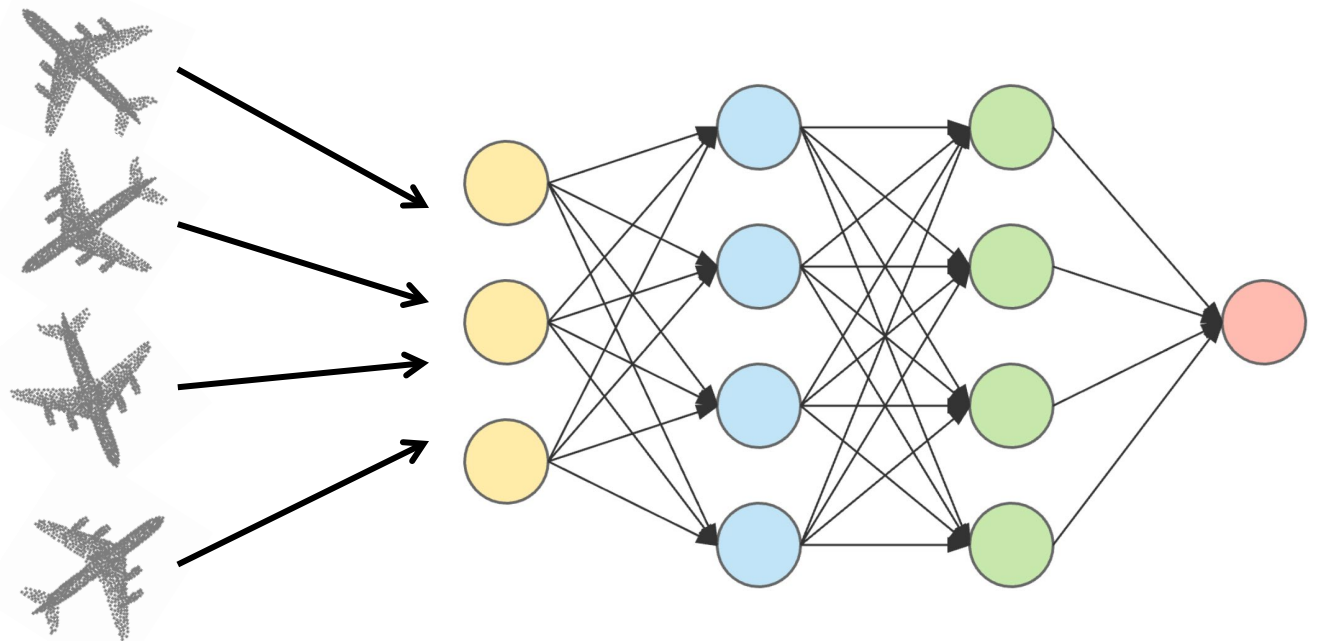


A **Less** Naïve Solution: Multi-view Approach

Feed multiple poses of the same object to the network **at once**

“If we don’t know what pose to look at, why not just look at **all** poses!”

- Can reproduce fairly good results
- Sacrifice data-efficiency – more memory consumption
- Theoretically equivariant – up to precision errors caused by discretization



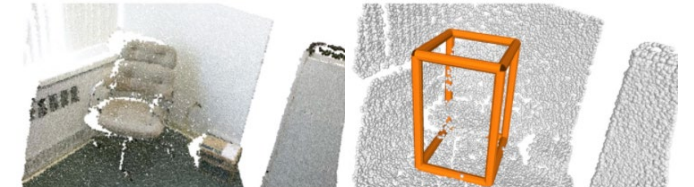
Summary: Why Equivariance / Invariance?

✓ **“Intelligence”**: More aligned with human perception



✓ **Generalization**: Eliminate the prior that all shapes in a dataset (e.g., ShapeNet) are aligned

✓ **Real-world applications**: Shapes (3D scans) in the wild may not have or come with canonical poses

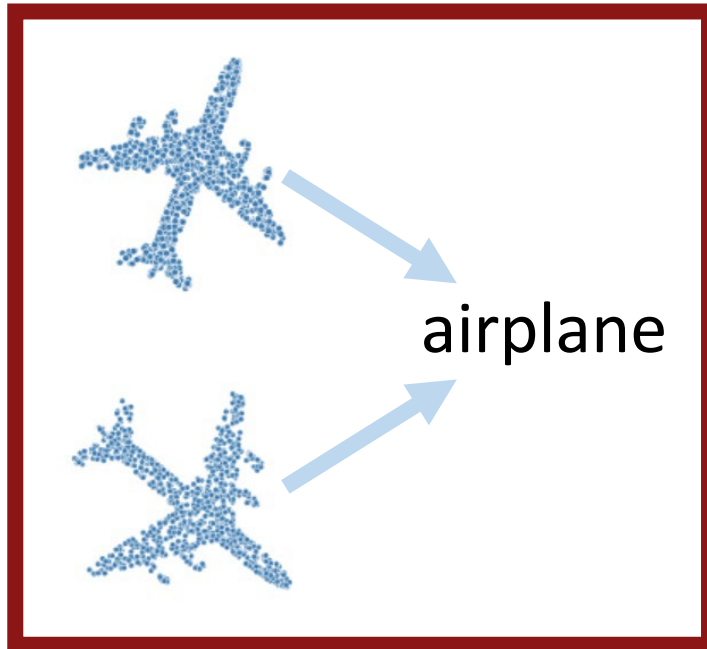


✓ **Data-efficiency**: Avoid exhaustive data augmentation



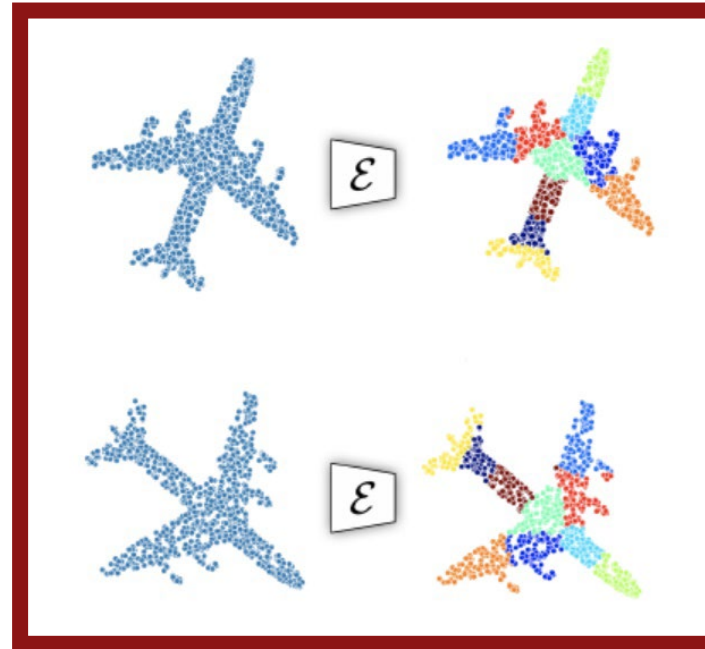
So We're Looking for...

classification



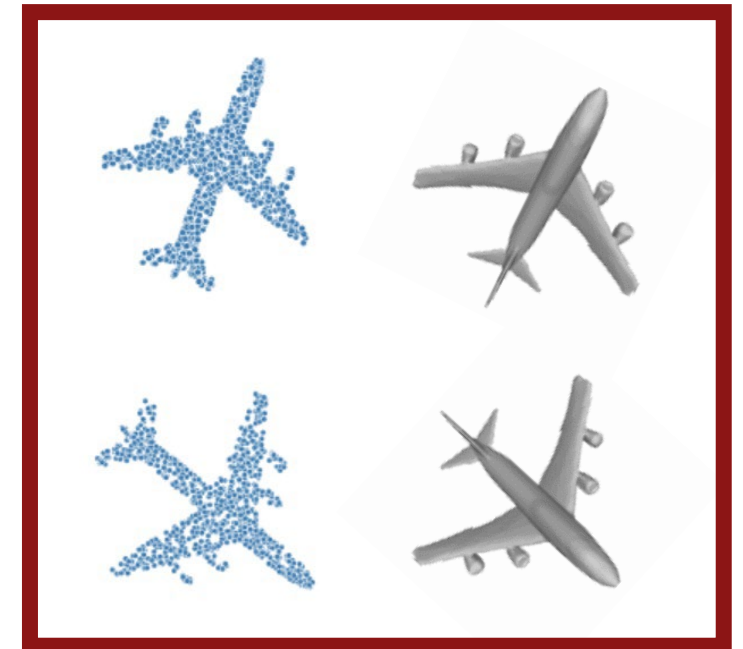
invariant

segmentation



equivariant
(or **invariant**, depending
on data representation)

reconstruction



equivariant encoder
invariant decoder

[W. Sun, A. Tagliasacchi, B. Deng, S. Sabour, S. Yazdani, G. Hinton, K. M. Yi, arXiv:2012.04718 (2020)]

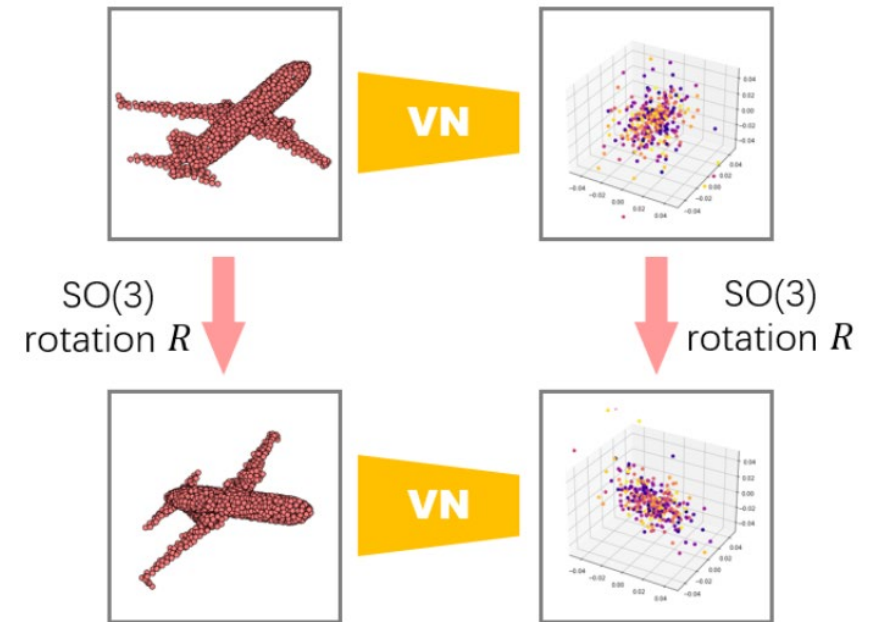
[J. J. Park, P. Florence, J. Straub, R. Newcombe, S. Lovegrove, CVPR 2019]

Equivariance!

We say a neural network $f(\cdot; \theta)$ is rotation equivariant, if for any 3D rotation $R \in SO(3)$ applied to its input \mathbf{x} , it is explicitly related to a transformation $D(R)$ on the network output satisfying

$$f(\mathbf{x}R; \theta) = f(\mathbf{x}; \theta)D(R)$$

- $D(R)$ should be independent of \mathbf{x}
- **Special case:** when $D(R) = R$ is the identity mapping, it is the common-sense “equivariance”
- **Special case:** when $D(R) = \mathbf{I}$ is the constant mapping, it is invariance



Vector Neurons (VN)

Vector Neurons

Classical (scalar) feature $\mathbf{z} = [z_1, z_2, \dots, z_C]^\top \in \mathbb{R}^C$, with $z_i \in \mathbb{R}$

Vector-list feature $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_C]^\top \in \mathbb{R}^{C \times 3}$, with $\mathbf{v}_i \in \mathbb{R}^3$

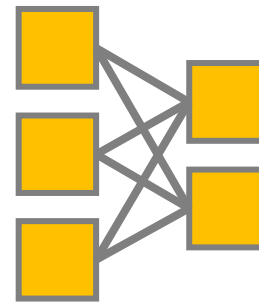
- For pointcloud with N points $\mathcal{V} = \{\mathbf{V}_1, \mathbf{V}_2, \dots, \mathbf{V}_N\} \in \mathbb{R}^{N \times C \times 3}$

Mapping between network layers:

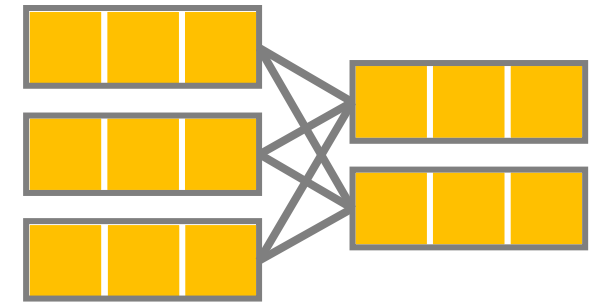
$$f(\cdot; \theta) : \mathbb{R}^{N \times C^{(d)} \times 3} \rightarrow \mathbb{R}^{N \times C^{(d+1)} \times 3}$$

? Equivariance to rotation $R \in \text{SO}(3)$:

$$f(\mathcal{V}R; \theta) = f(\mathcal{V}; \theta)R$$



**(classical)
scalar neurons**



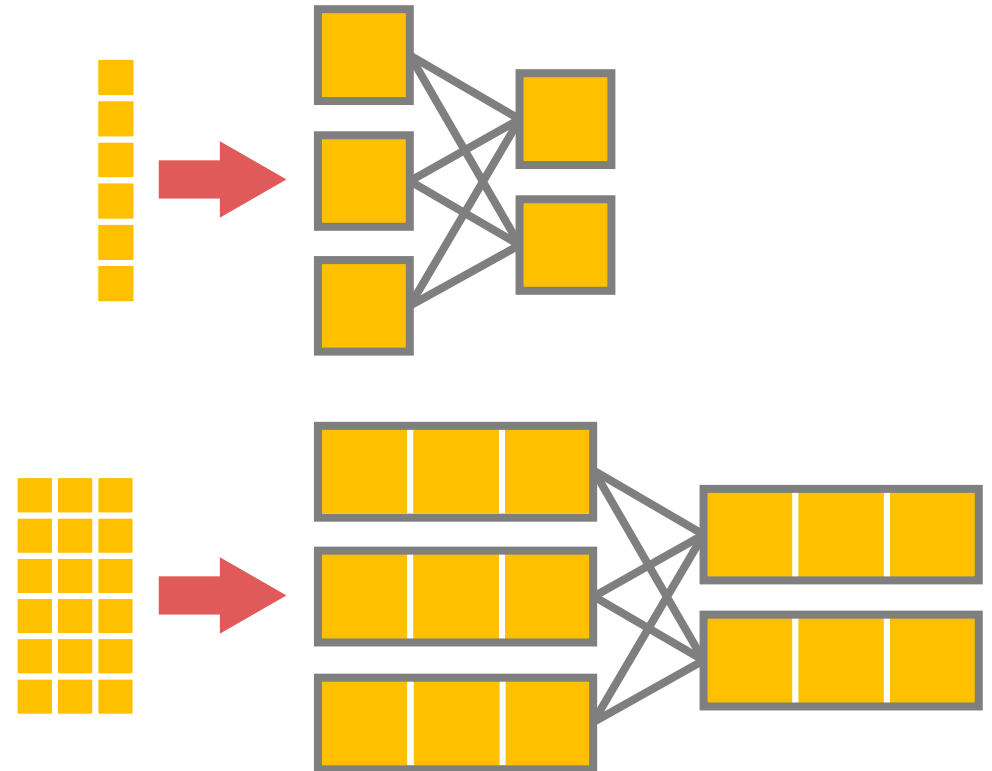
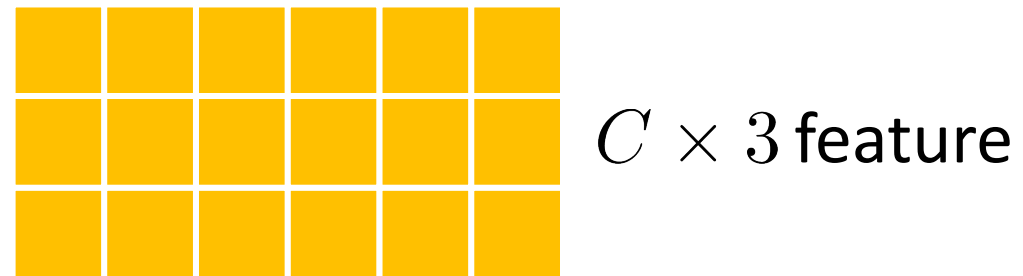
vector neurons

VN Features

Classical: scalar channels

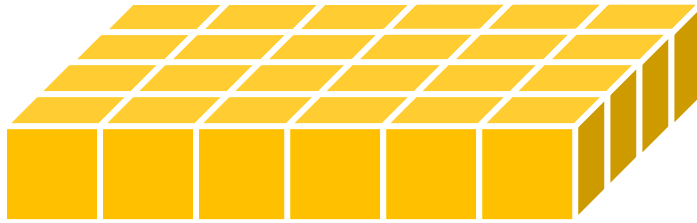


VN: 3D vector channels



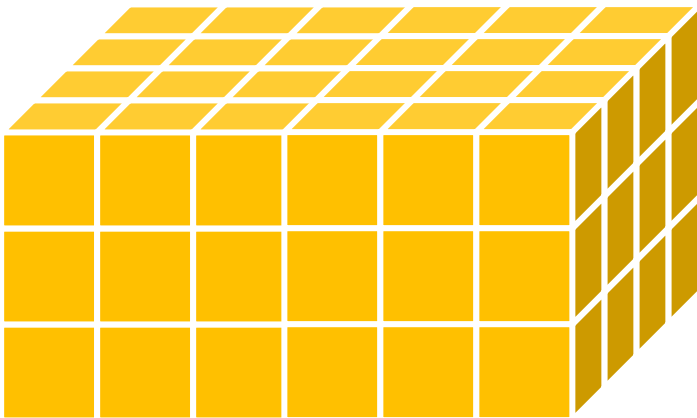
VN Features (for Point Cloud)

Classical:

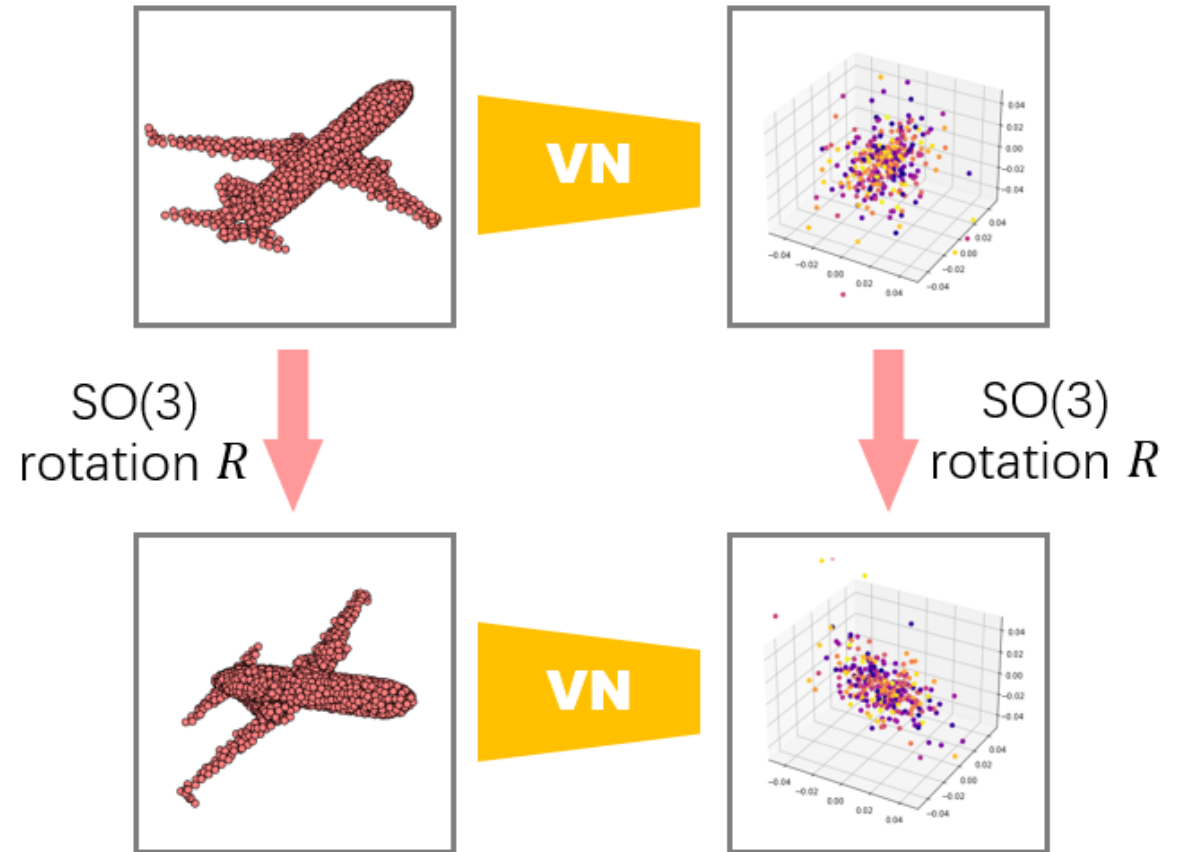


$N \times C \times 1$ feature

VN:



$N \times C \times 3$ feature




VN Linear Layer

Linear operator: left multiply by the learnable weight matrix



$C' \times C$ weight \times $C \times 3$ feature $=$ $C' \times 3$ feature

Equivariance: right multiply by the SO(3) rotation matrix



$C' \times C$ weight \times $\left[C \times 3 \text{ feature} \times \begin{matrix} \text{SO(3) rotation matrix} \end{matrix} \right] = \left[C' \times 3 \text{ feature} \times \begin{matrix} \text{SO(3) rotation matrix} \end{matrix} \right]$

VN Linear Layer

Vector-list feature $V \in \mathbb{R}^{C \times 3}$

Linear operator $f_{\text{lin}}(\cdot; \mathbf{W})$ with learnable weights $\mathbf{W} \in \mathbb{R}^{C' \times C}$:

$$V' = f_{\text{lin}}(V; \mathbf{W}) = \mathbf{W}V \in \mathbb{R}^{C' \times 3}$$

Equivariance to rotation $R \in \text{SO}(3)$:

$$f_{\text{lin}}(VR; \mathbf{W}) = \mathbf{W}VR = f_{\text{lin}}(V; \mathbf{W})R = V'R$$

- \mathbf{W} - left multiplication, R - right multiplication
- Note the absence of a bias term

VN Non-Linearity

ReLU Non-Linearity

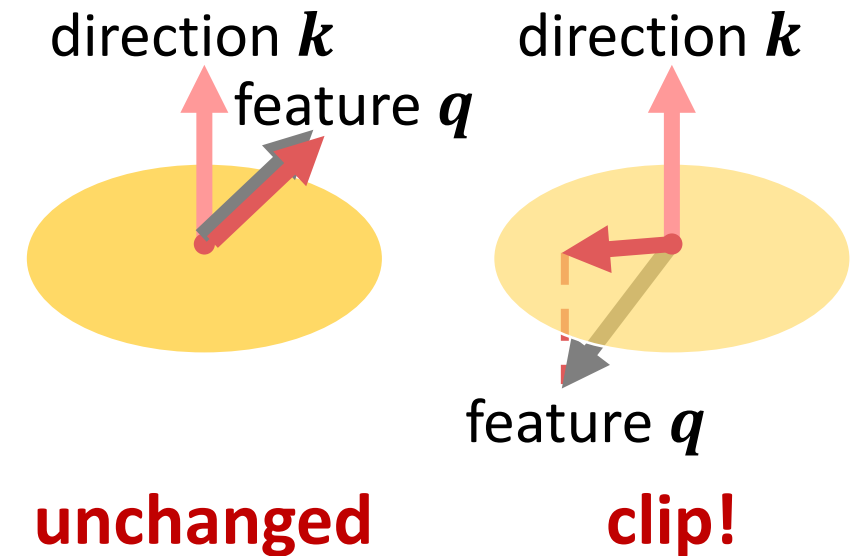
Weights $\mathbf{W} \in \mathbb{R}^{1 \times C}$ and $\mathbf{U} \in \mathbb{R}^{1 \times C}$

Learn a feature $\mathbf{q} = \mathbf{W}\mathbf{V} \in \mathbb{R}^{1 \times 3}$

Learn a direction $\mathbf{k} = \mathbf{U}\mathbf{V} \in \mathbb{R}^{1 \times 3}$

For each output vector neuron $\mathbf{v}' \in \mathbf{V}'$

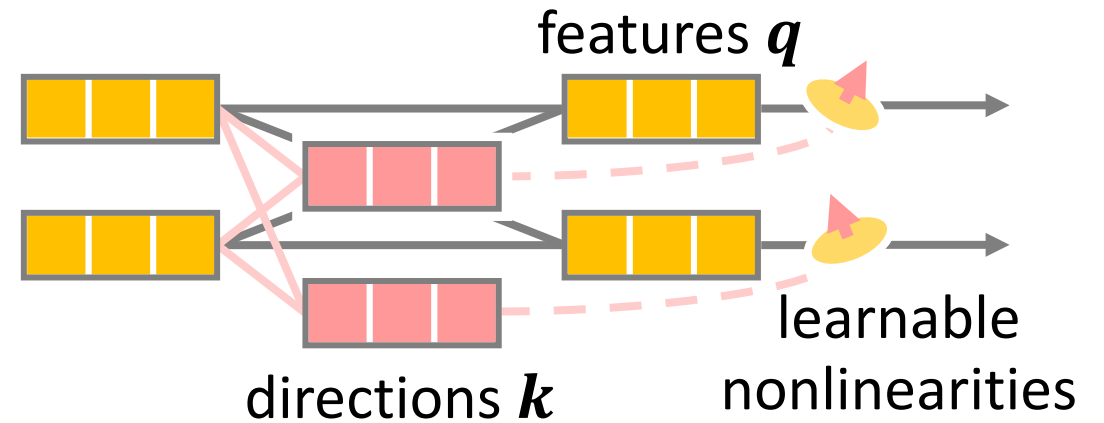
$$\mathbf{v}' = \begin{cases} \mathbf{q} & \text{if } \langle \mathbf{q}, \mathbf{k} \rangle \geq 0 \\ \mathbf{q} - \langle \mathbf{q}, \frac{\mathbf{k}}{\|\mathbf{k}\|} \rangle \frac{\mathbf{k}}{\|\mathbf{k}\|} & \text{otherwise} \end{cases}$$



VN Non-Linearity

ReLU Non-Linearity

- **Non-linear layer** (with built-in linear layer)
 - = **input linear transformation q** + **non-linearity k**
- A single non-linear layer detached from linear layer: $q = v$
 - introduce additional network depth!
- Other non-linearities



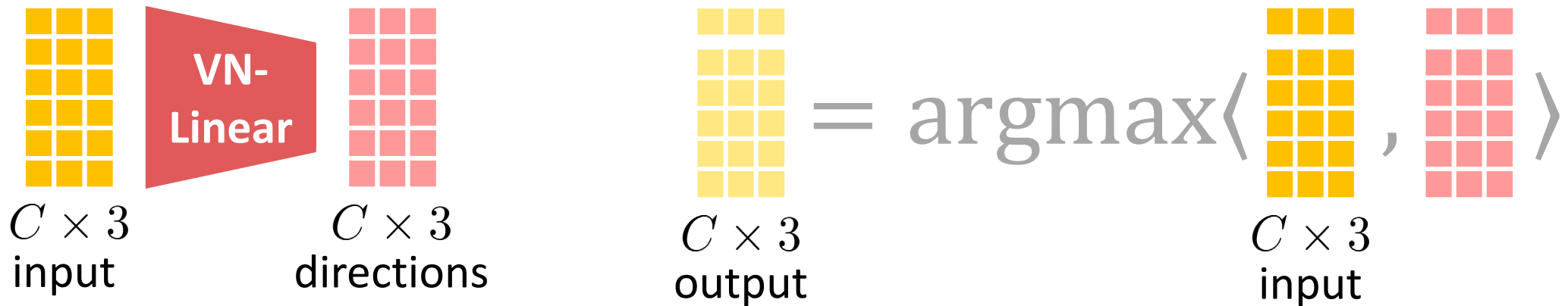
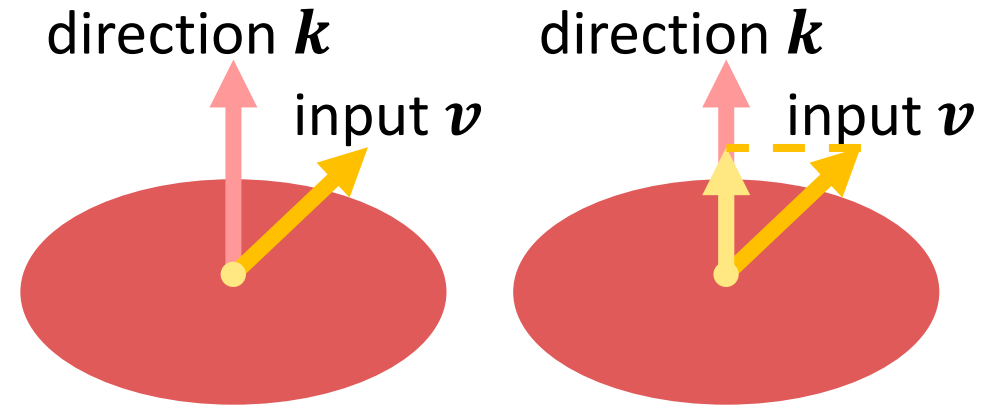
overall structure

VN Pooling

✓ Mean pooling

? Max pooling

- (Similar to non-linearity)
- argmax alone learned directions



VN Pooling

✓ Mean pooling

? Max pooling

- similar to non-linearities → use learnable directions

Given a set of vector-lists $\mathcal{V} = \{\mathbf{V}_1, \mathbf{V}_2, \dots, \mathbf{V}_N\} \in \mathbb{R}^{N \times C \times 3}$

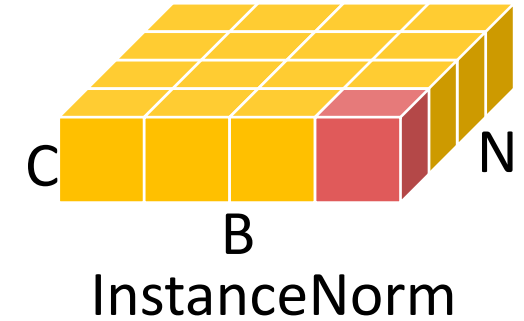
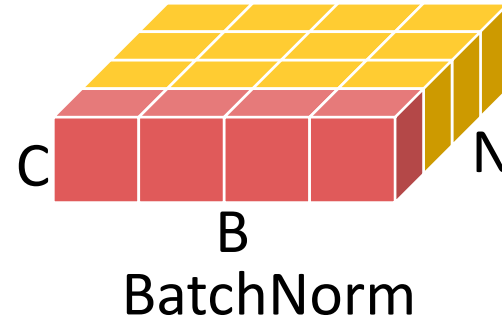
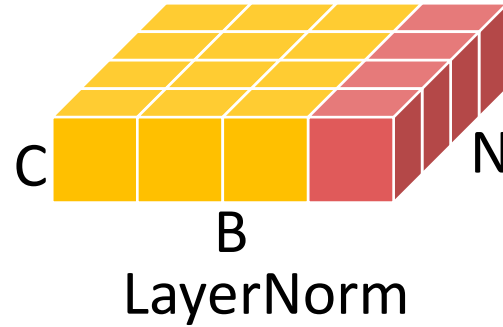
Learn data-dependent directions $\mathcal{K} = \{\mathbf{W}\mathbf{V}_n\}_{n=1}^N \in \mathbb{R}^{N \times C \times 3}$ with learnable weights $\mathbf{W} \in \mathbb{R}^{C \times C}$

Max pool along \mathcal{K} directions: for each channel $c \in [C]$

$$f_{\text{MAX}}(\mathcal{V})[c] = \mathbf{V}_{n^*}[c] \quad \text{where } n^* = \operatorname{argmax} \langle \mathbf{W}\mathbf{V}_n[c], \mathbf{V}_n[c] \rangle$$

VN Normalizations

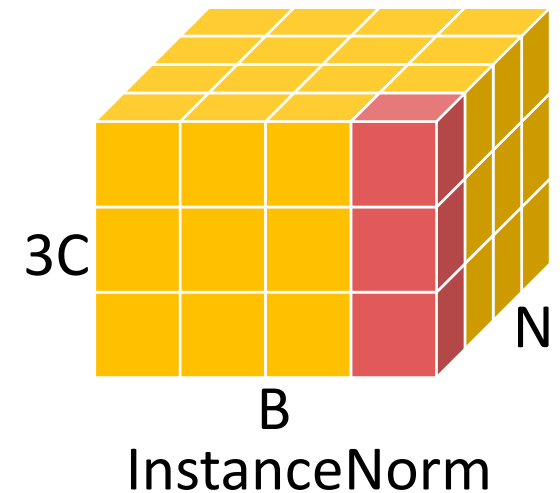
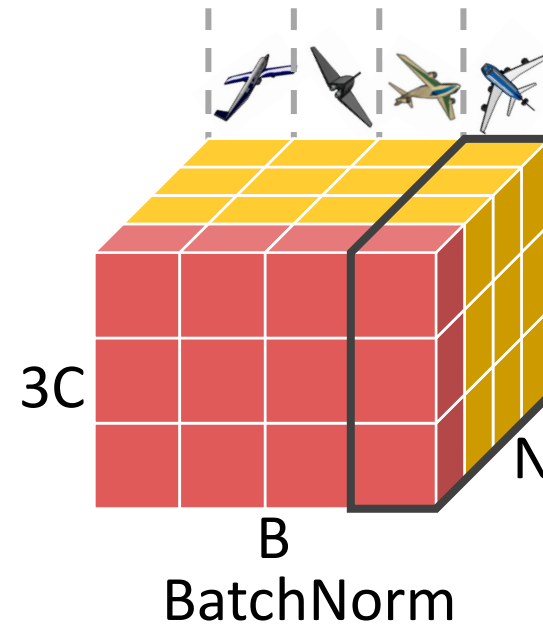
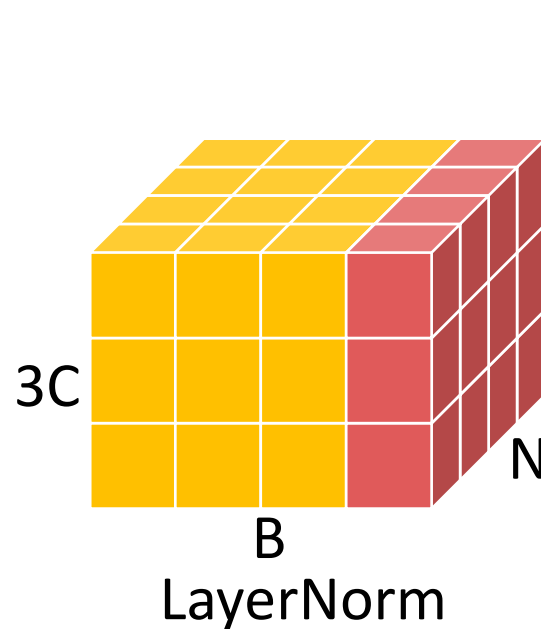
- ✓ LayerNorm
- ✓ InstanceNorm
- ✓ Dropout



(classical) scalar neurons

? BatchNorm

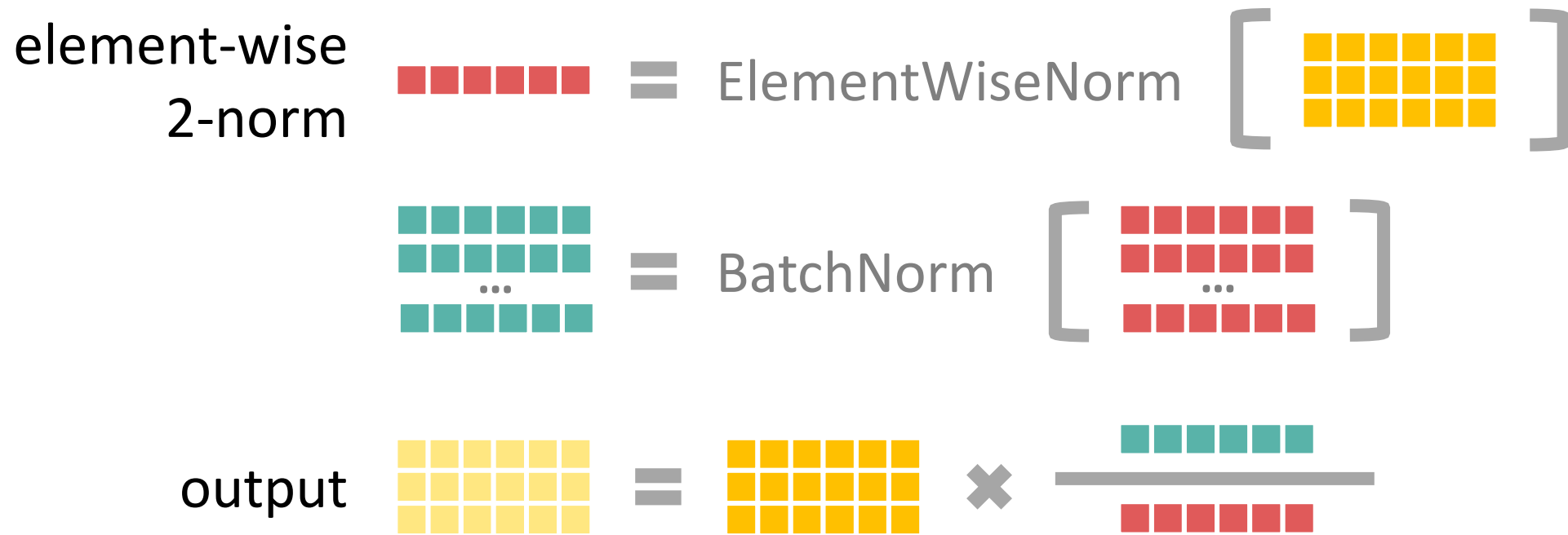
averaging across arbitrarily rotated inputs would not necessarily be meaningful



Vector neurons

VN Normalizations

BatchNorm



VN Normalizations

BatchNorm

- Normalize the **2-norm (invariant component)** of the vector-list feature

$$\mathbf{N}_b = \text{ElementWiseNorm}(\mathbf{V}_b) \in \mathbb{R}^{N \times 1}$$

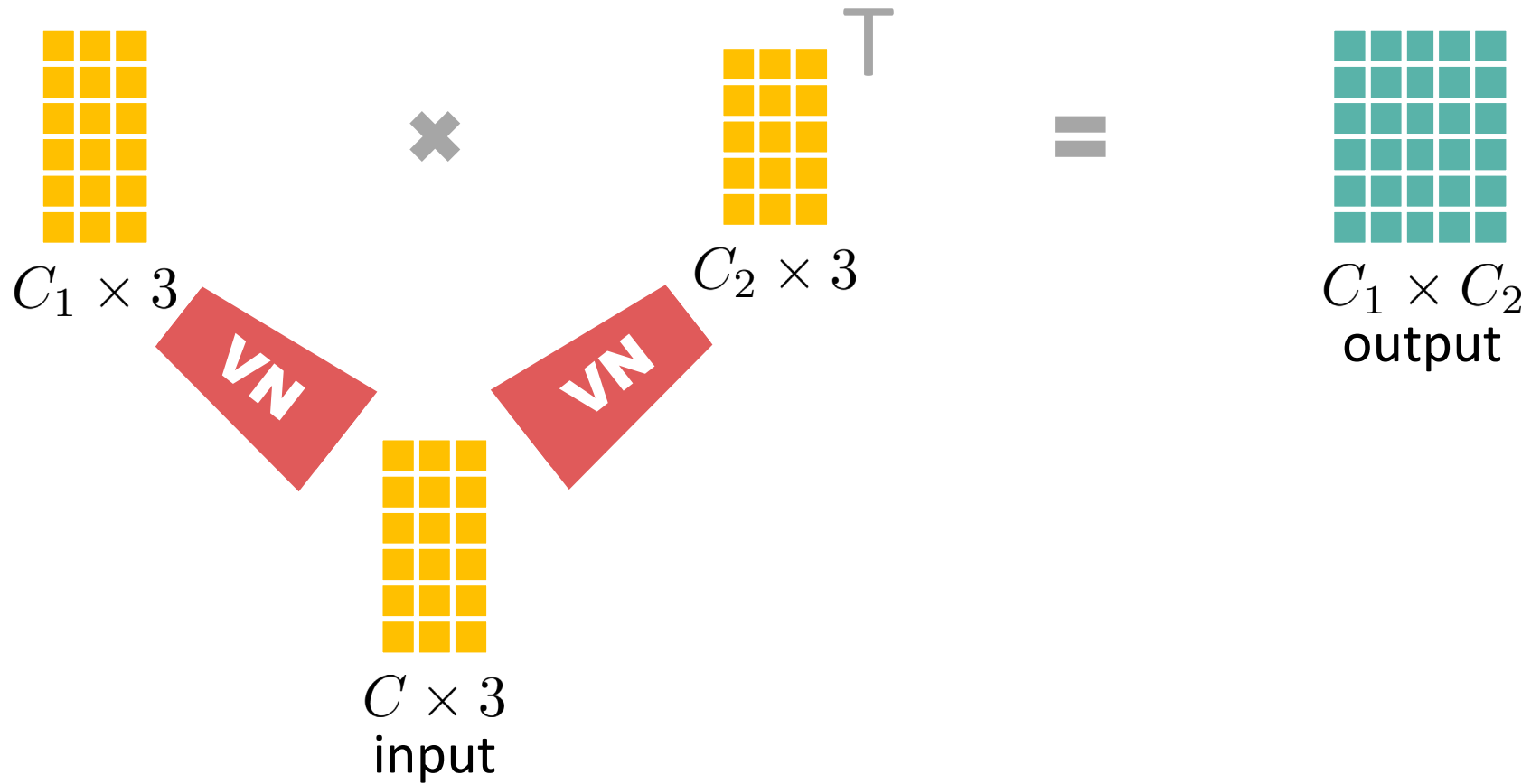
$$\{\mathbf{N}'_b\}_{b=1}^B = \text{BatchNorm}(\{\mathbf{N}_b\}_{b=1}^B)$$

$$\mathbf{V}'_b = \mathbf{V}_b[c] \frac{\mathbf{N}'_b[c]}{\mathbf{N}_b[c]} \quad \forall c \in [C]$$

- Element-wise norm: 2-norm for each vector $\mathbf{v}_c \in \mathbf{V}_b$

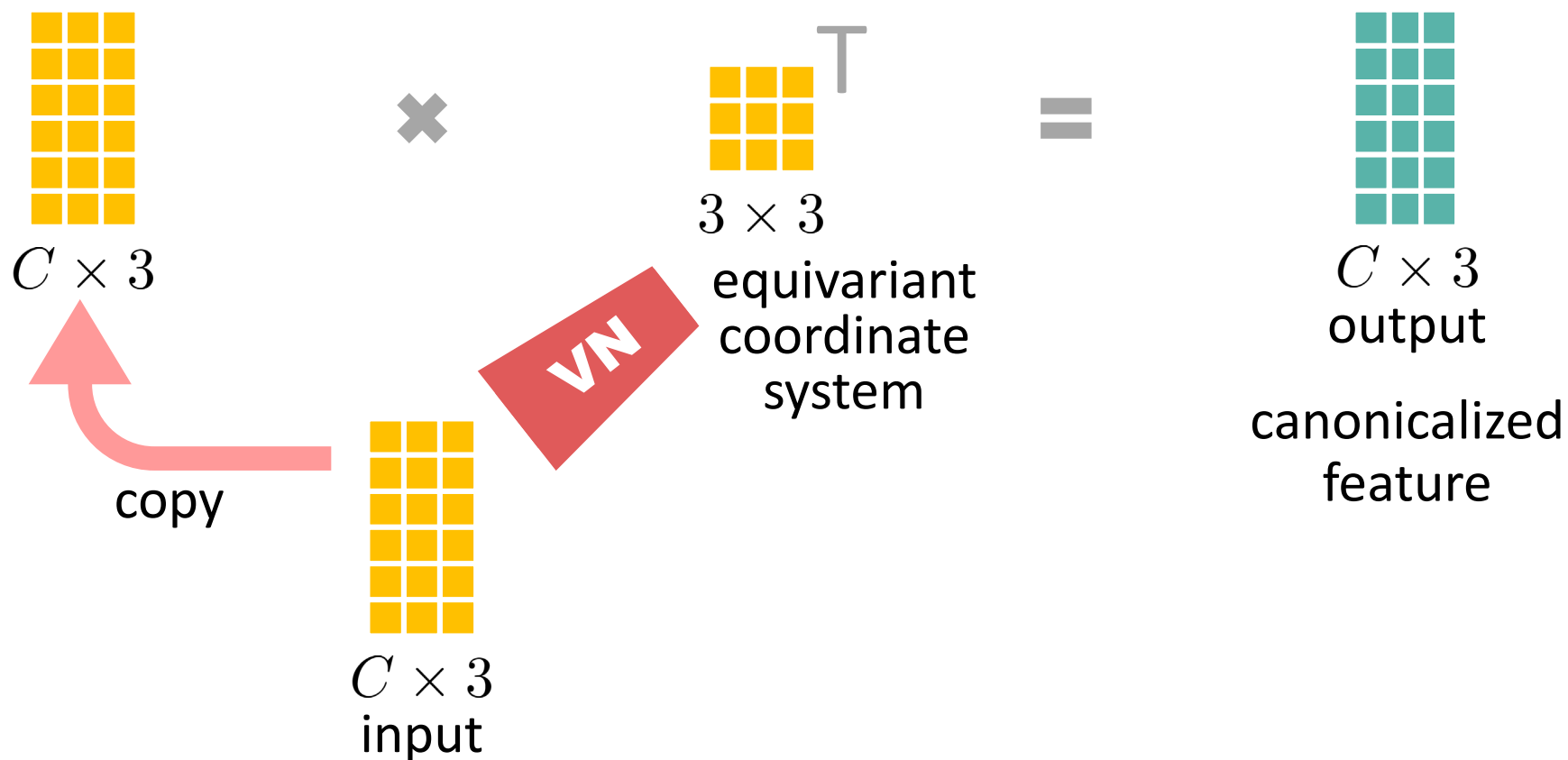
VN Invariant Layer

(**equivariant** feature) \times (**equivariant** feature)^T = (**invariant** feature)



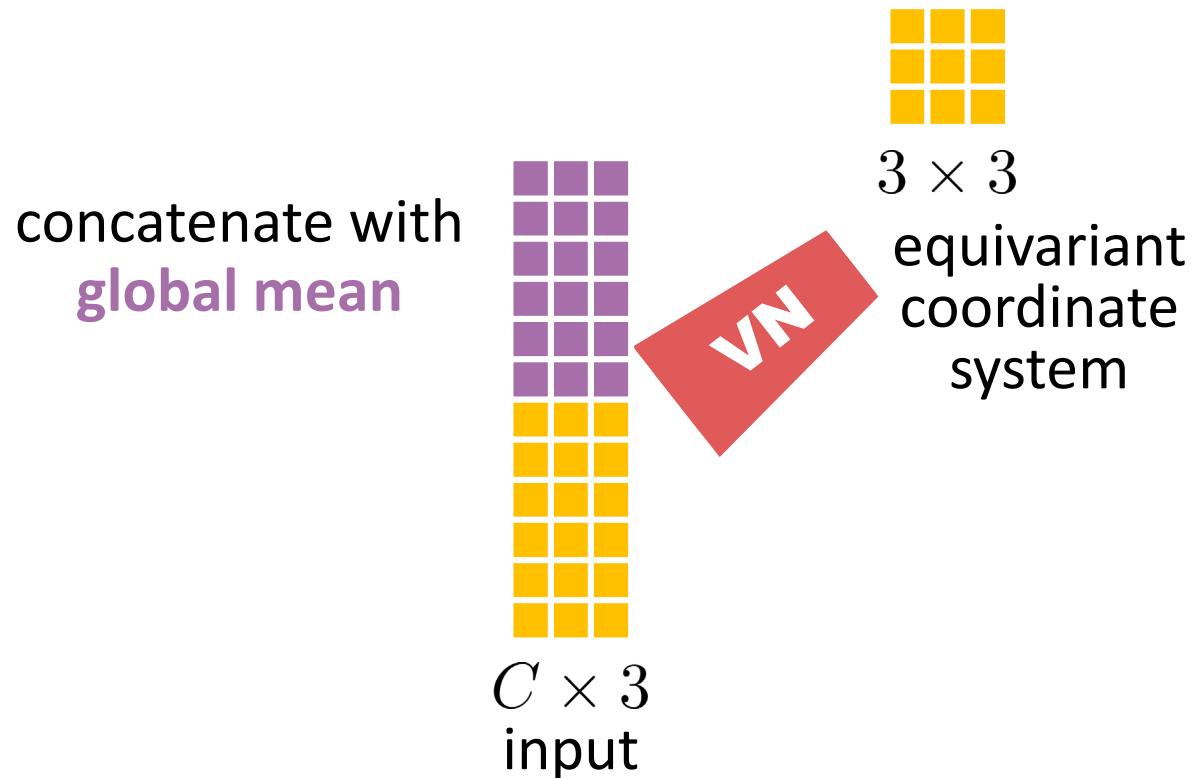
VN Invariant Layer

Specifically...



VN Invariant Layer

For pointcloud: combine **global information** with local features



VN Invariant Layer

- Product of an equivariant signal $\mathbf{V} \in \mathbb{R}^{C \times 3}$ by the transpose of another equivariant signal $\mathbf{T} \in \mathbb{R}^{C' \times 3} \rightarrow$ invariant signal
- **Special case:** $\mathbf{T} \in \mathbb{R}^{3 \times 3}$ - an equivariant coordinate system
- For pointcloud, concatenate local feature $\mathbf{V} \in \mathbb{R}^{C \times 3}$ with global mean $\bar{\mathbf{V}} = \frac{1}{N} \sum_{n=1}^N \mathbf{V}_n \in \mathbb{R}^{C \times 3}$

Invariant layer: $\mathbf{T}_n = \text{VN-MLP}([\mathbf{V}_n, \bar{\mathbf{V}}])$

$$\text{VN-In}(\mathbf{V}_n) = \mathbf{V}_n \mathbf{T}_n^\top$$

Build VN Networks: VN-DGCNN

DGCNN

Edge feature: $\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare'_{nm} = \text{ReLU}(\Theta(\blacksquare\blacksquare\blacksquare_m - \blacksquare\blacksquare\blacksquare_n) + \Phi\blacksquare\blacksquare\blacksquare_n)$

Aggregation: $\blacksquare\blacksquare\blacksquare\blacksquare'_n = \text{Pool}_{m:(n,m) \in \mathcal{E}}(\blacksquare\blacksquare\blacksquare\blacksquare'_{nm})$

VN-DGCNN

Edge feature: $\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare'_{nm} = \text{VN-ReLU}(\Theta(\blacksquare\blacksquare\blacksquare_m - \blacksquare\blacksquare\blacksquare_n) + \Phi\blacksquare\blacksquare\blacksquare_n)$

Aggregation $\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare'_n = \text{VN-Pool}_{m:(n,m) \in \mathcal{E}}(\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare\blacksquare'_{nm})$

Build VN Networks: VN-DGCNN

Convolution in latent spaces, graph edges not embedded in \mathbb{R}^3

DGCNN: $\mathbf{x} \in \mathbb{R}^C$, $\mathbf{e}'_n, \mathbf{x}'_n \in \mathbb{R}^{C'}$ scalar features

Edge conv $\mathbf{e}'_{nm} = \text{ReLU}(\Theta(\mathbf{x}_m - \mathbf{x}_n) + \Phi \mathbf{x}_n)$

Aggregation $\mathbf{x}'_n = \text{Pool}_{m:(n,m) \in \mathcal{E}}(\mathbf{e}'_{nm})$

VN-DGCNN: $\mathbf{V} \in \mathbb{R}^{C \times 3}$, $\mathbf{E}'_n, \mathbf{V}'_n \in \mathbb{R}^{C' \times 3}$ vector-list features

Edge conv $\mathbf{E}'_{nm} = \text{VN-ReLU}(\Theta(\mathbf{V}_m - \mathbf{V}_n) + \Phi \mathbf{V}_n)$

Aggregation $\mathbf{V}'_n = \text{VN-Pool}_{m:(n,m) \in \mathcal{E}}(\mathbf{E}'_{nm})$

Build VN Networks: VN-PointNet

PointNet

$$\blacksquare\blacksquare\blacksquare\blacksquare' = \text{Pool}_{\mathbf{x}_n \in \mathcal{X}}(h(\blacksquare\blacksquare\blacksquare\blacksquare_1), h(\blacksquare\blacksquare\blacksquare\blacksquare_2), \dots, h(\blacksquare\blacksquare\blacksquare\blacksquare_N))$$

VN-PointNet

$$\blacksquare\blacksquare\blacksquare\blacksquare' = \text{VN-Pool}_{\mathbf{V}_n \in \mathcal{V}}(f(\blacksquare\blacksquare\blacksquare\blacksquare_1), f(\blacksquare\blacksquare\blacksquare\blacksquare_2), \dots, f(\blacksquare\blacksquare\blacksquare\blacksquare_N))$$

Build VN Networks: VN-PointNet

No convolutions, only point-wise feature transformations

PointNet: $\mathbf{x}' = \text{Pool}_{\mathbf{x}_n \in \mathcal{X}}(h(\mathbf{x}_1), \dots, h(\mathbf{x}_N))$

VN-PointNet: $\mathbf{V}' = \text{VN-Pool}_{\mathbf{V}_n \in \mathcal{V}}(f(\mathbf{V}_1), \dots, f(\mathbf{V}_N))$

- h MLP, f VN-MLP
- First layer single channel \rightarrow edge conv to lift dimensionality
 - “Cannot apply per-channel transformation to gray-scale images”

Classification

Classification results on ModelNet40

VN Networks

Rotation sensitive methods

Rotation robust methods

Methods	z/z	$z/\text{SO}(3)$	$\text{SO}(3)/\text{SO}(3)$
Point / mesh inputs			
PointNet [25]	85.9	19.6	74.7
DGCNN [35]	90.3	33.8	88.6
VN-PointNet	77.5	77.5	77.2
VN-DGCNN	89.5	89.5	90.2
PCNN [2]	92.3	11.9	85.1
ShellNet [40]	93.1	19.9	87.8
PointNet++ [26]	91.8	28.4	85.0
PointCNN [20]	92.5	41.2	84.5
Spherical-CNN [11]	88.9	76.7	86.9
$a^3\text{S-CNN}$ [21]	89.6	87.9	88.7
SFCNN [27]	91.4	84.8	90.1
TFN [32]	88.5	85.3	87.6
RI-Conv [39]	86.5	86.4	86.4
SPHNet [24]	87.7	86.6	87.6
ClusterNet [6]	87.1	87.1	87.1
GC-Conv [41]	89.0	89.1	89.2
RI-Framework [18]	89.4	89.4	89.3

Classification

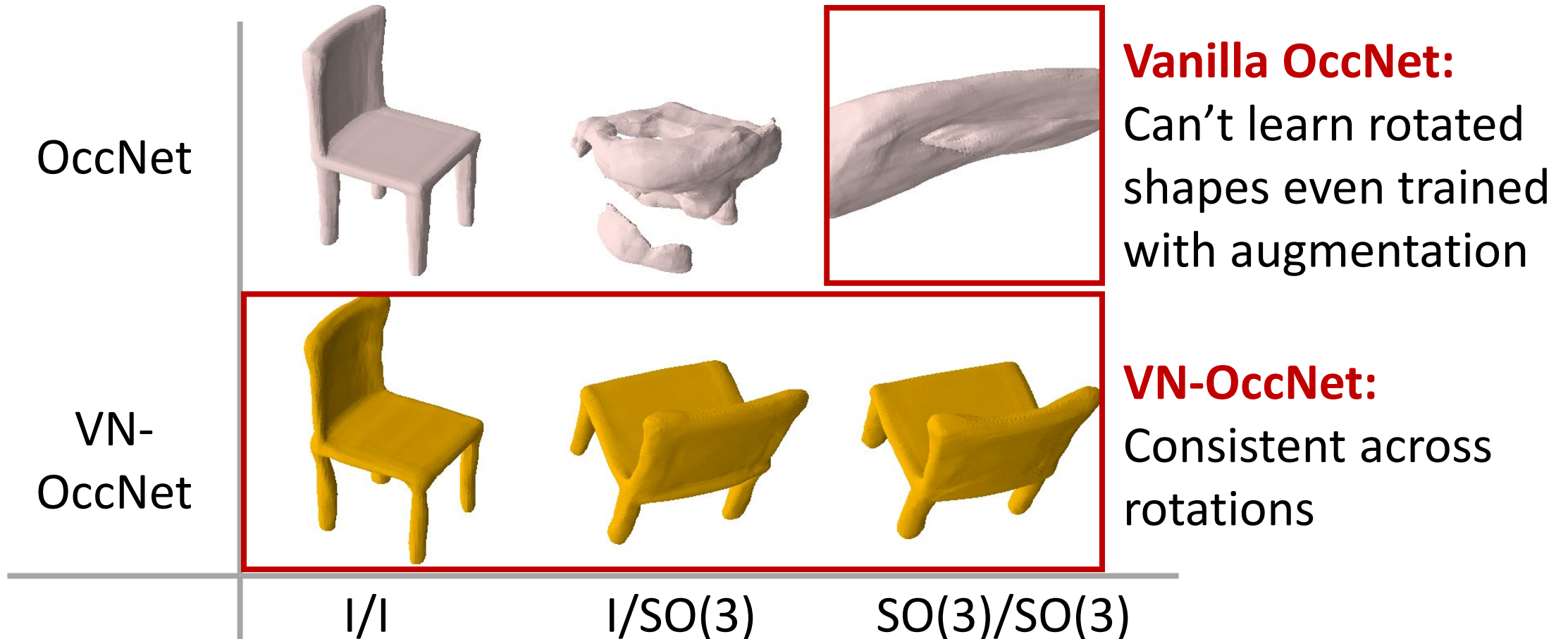
Classification results on ModelNet40

- VN networks are **robust to (seen & unseen) rotations**
- Excellent performance compared with other methods
- **SO(3)/SO(3)**: equivariance by construction is better than rotation augmentation

Methods	z/z	$z/\text{SO}(3)$	SO(3)/SO(3)
Point / mesh inputs			
PointNet [25]	85.9	19.6	74.7
DGCNN [35]	90.3	33.8	88.6
VN-PointNet	77.5	77.5	77.2
VN-DGCNN	89.5	89.5	90.2
PCNN [2]	92.3	11.9	85.1
ShellNet [40]	93.1	19.9	87.8
PointNet++ [26]	91.8	28.4	85.0
PointCNN [20]	92.5	41.2	84.5
Spherical-CNN [11]	88.9	76.7	86.9
$a^3\text{S-CNN}$ [21]	89.6	87.9	88.7
SFCNN [27]	91.4	84.8	90.1
TFN [32]	88.5	85.3	87.6
RI-Conv [39]	86.5	86.4	86.4
SPHNet [24]	87.7	86.6	87.6
ClusterNet [6]	87.1	87.1	87.1
GC-Conv [41]	89.0	89.1	89.2
RI-Framework [18]	89.4	89.4	89.3

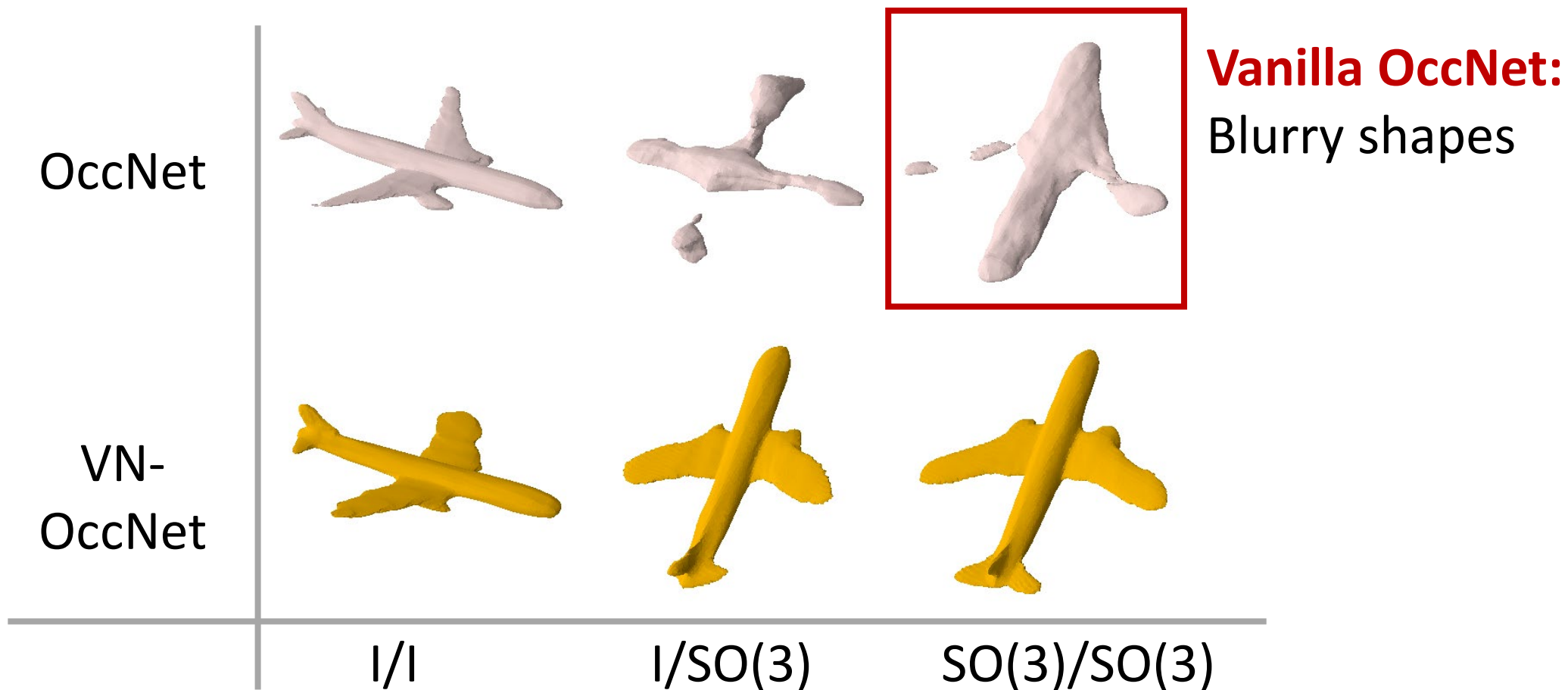
Neural Implicit Reconstruction

Results on ShapeNet (Examples)



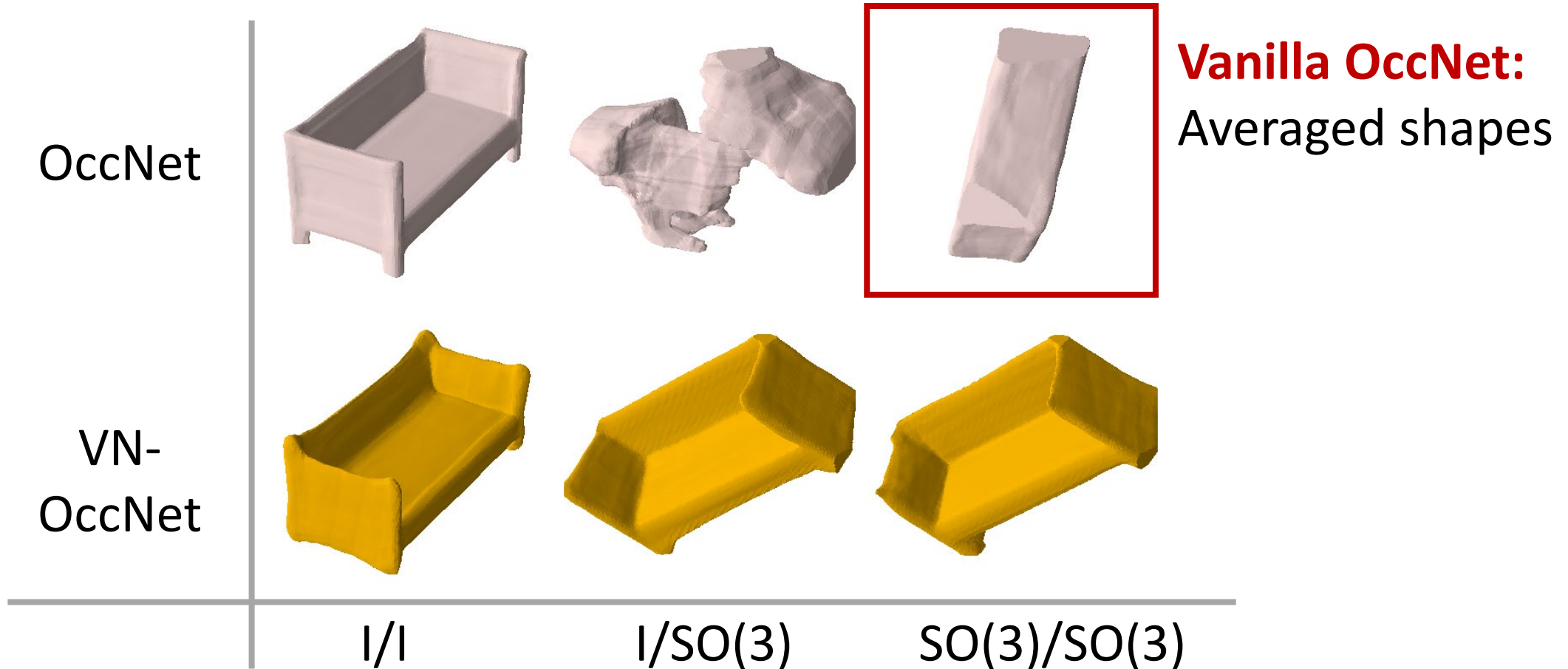
Neural Implicit Reconstruction

Results on ShapeNet (Examples)



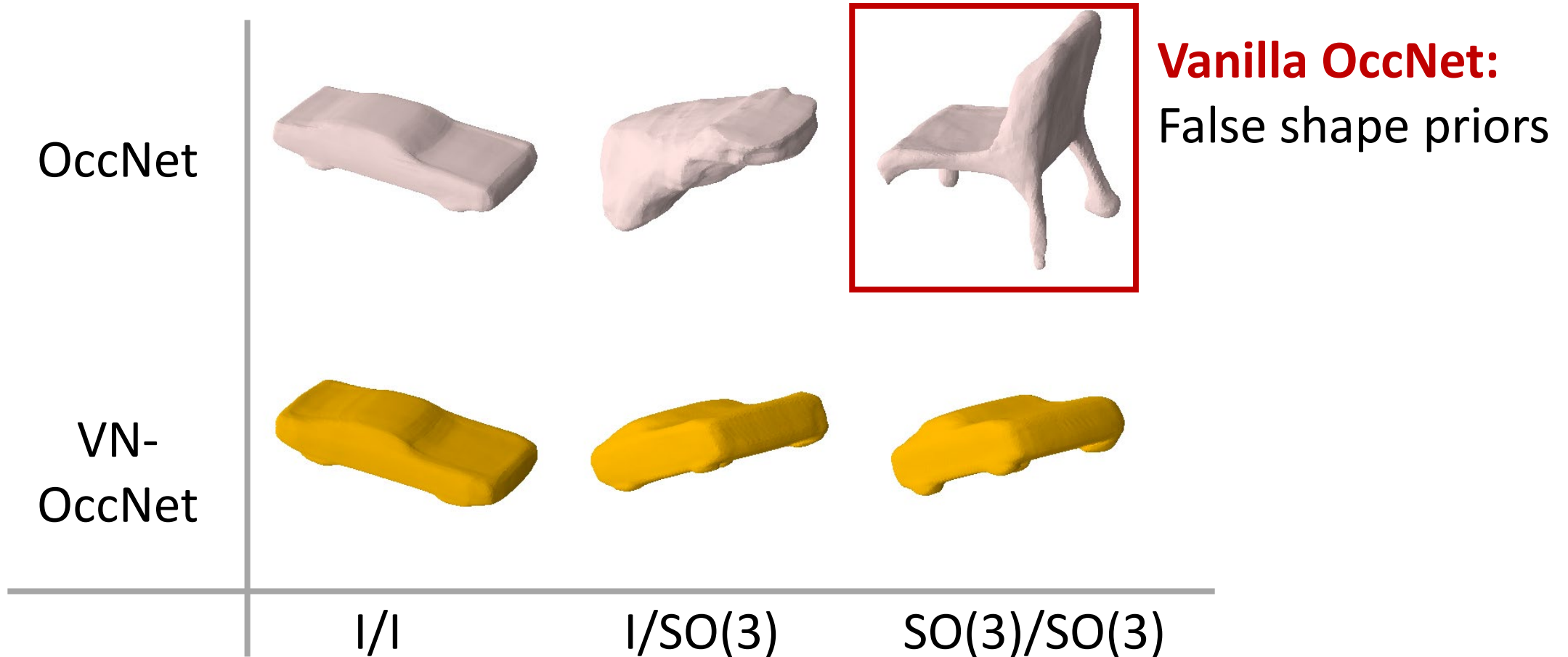
Neural Implicit Reconstruction

Results on ShapeNet (Examples)



Neural Implicit Reconstruction

Results on ShapeNet (Examples)



Summary: Vector Neurons

- **Vector Neurons:**

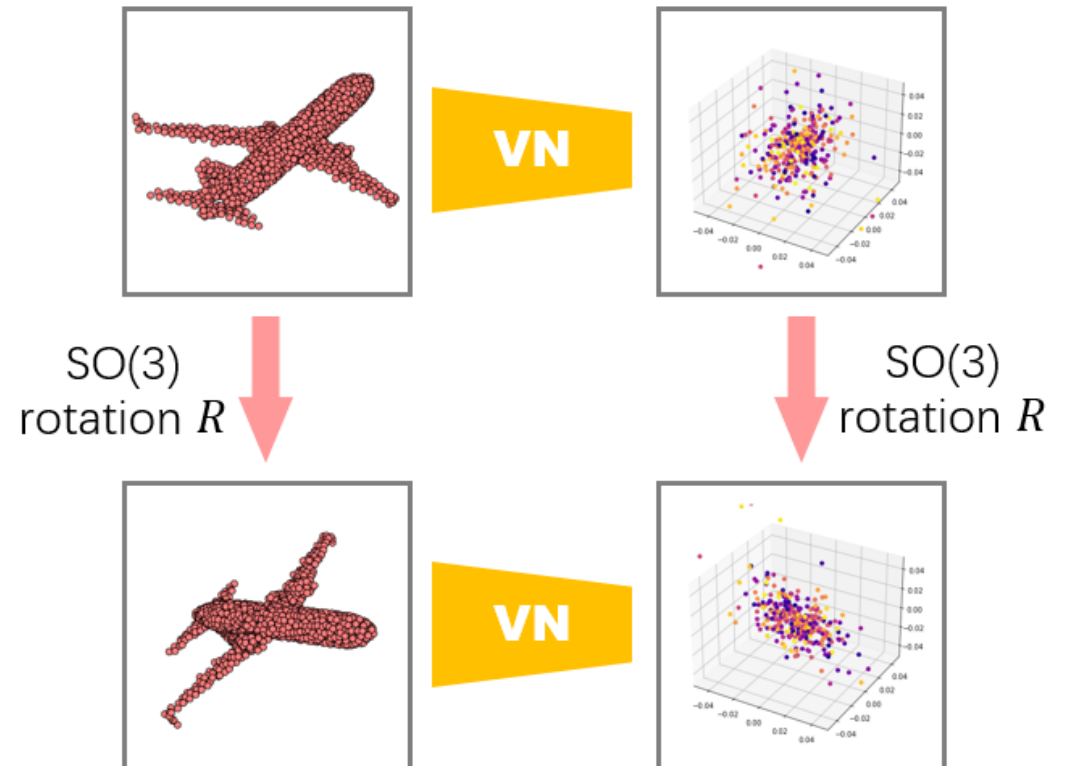
Lift latent features to 3D vector lists

- **Building blocks:**

- Linear layer
- Non-linearity (ReLU)
- Pooling (MaxPool)
- Normalizations (BatchNorm)
- Invariance

- **Network examples:**

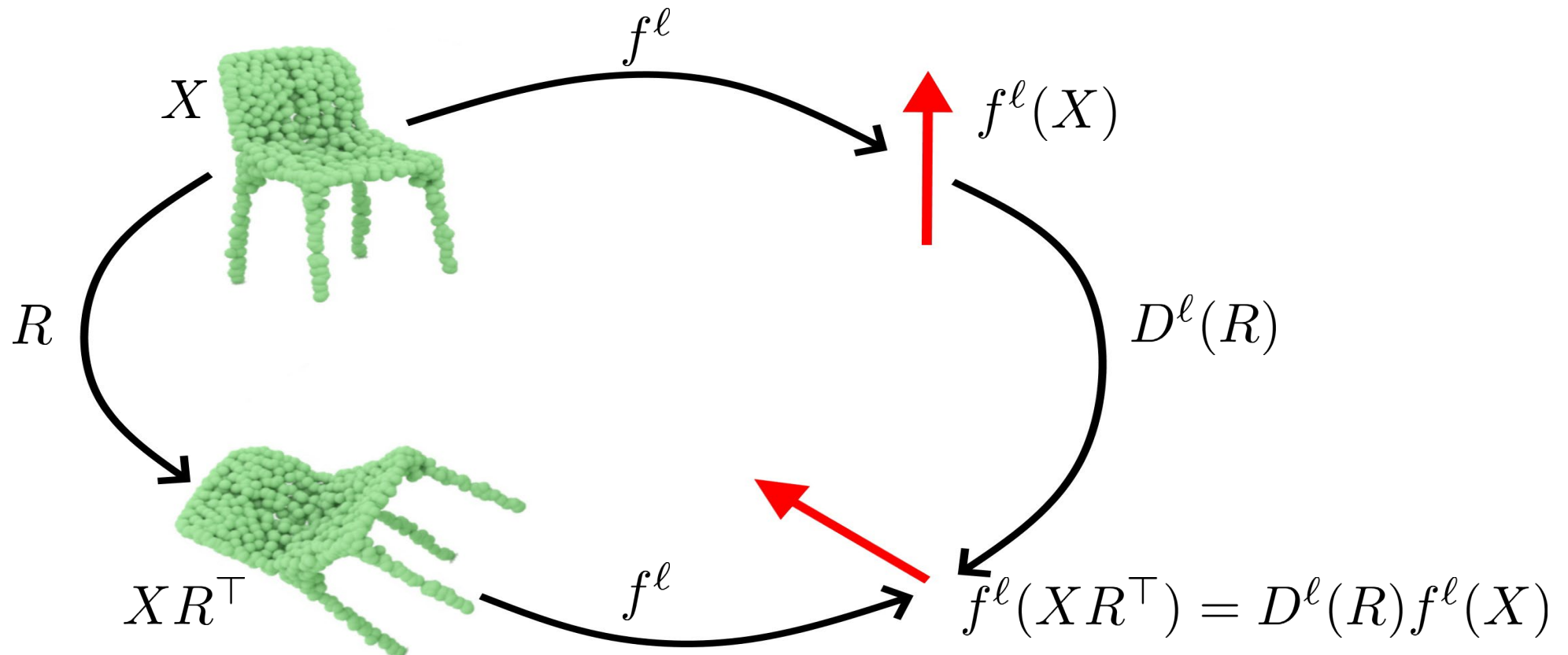
- VN-DGCNN
- VN-PointNet



Tensor Field Networks and Equivariant CNNs

SO(3) Equivariant Features

- A type $\ell \in \mathbb{N}$ equivariant feature map $f^\ell : \mathbb{R}^{n \times 3} \rightarrow \mathbb{R}^{2\ell+1}$ satisfying $f^\ell(XR^\top) = D^\ell(R)f^\ell(X)$ for $R \in \text{SO}(3)$, where $D^\ell(R) \in \text{SO}(2\ell+1)$ is the type ℓ wigner matrix.



Tensor Field Networks (TFN)

- Given a pointcloud X a TFN can produce pointwise or global equivariant features of different types ($\ell \in \mathbb{N}$):

Pointwise features

$$f^\ell(X)_{i,c,:}$$

point channel

global features

$$g^\ell(X)_{c,:}$$

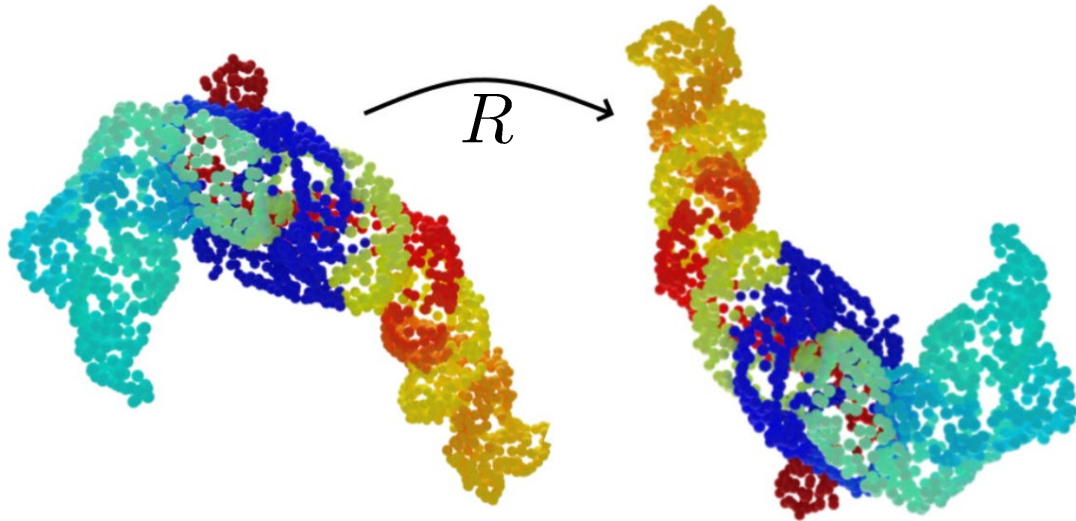
channel

Examples of Type 0 Features

- Type 0 features are rotation invariant as $D^0(R) = 1$:

Segmentation

$$f^0(X)_{i,c,:} = f^0(XR^\top)$$



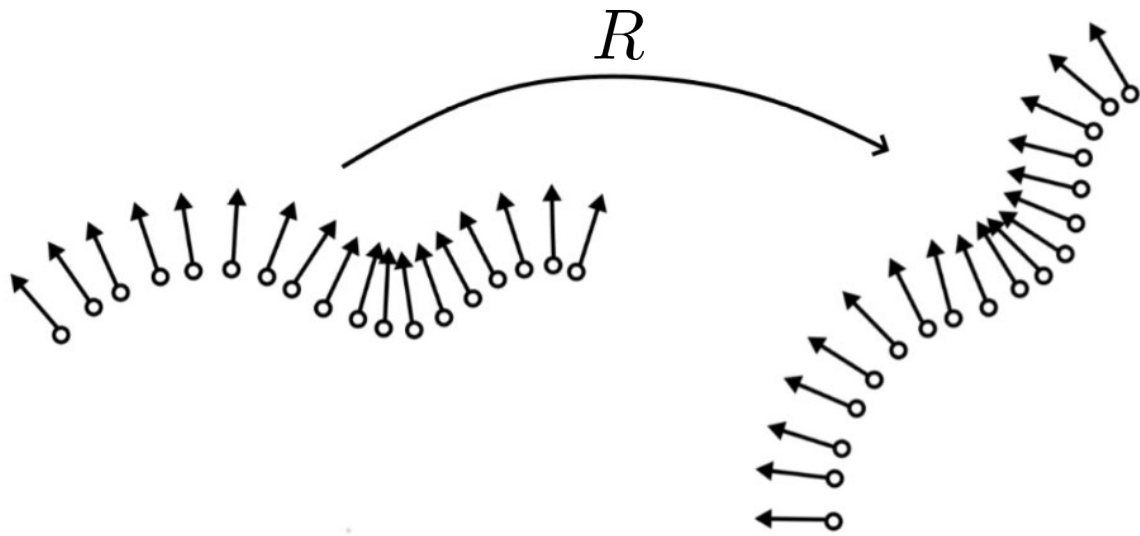
Classification

$$g^0 \left(\begin{array}{c} \text{Chair} \end{array} \right) = \text{Chair}$$

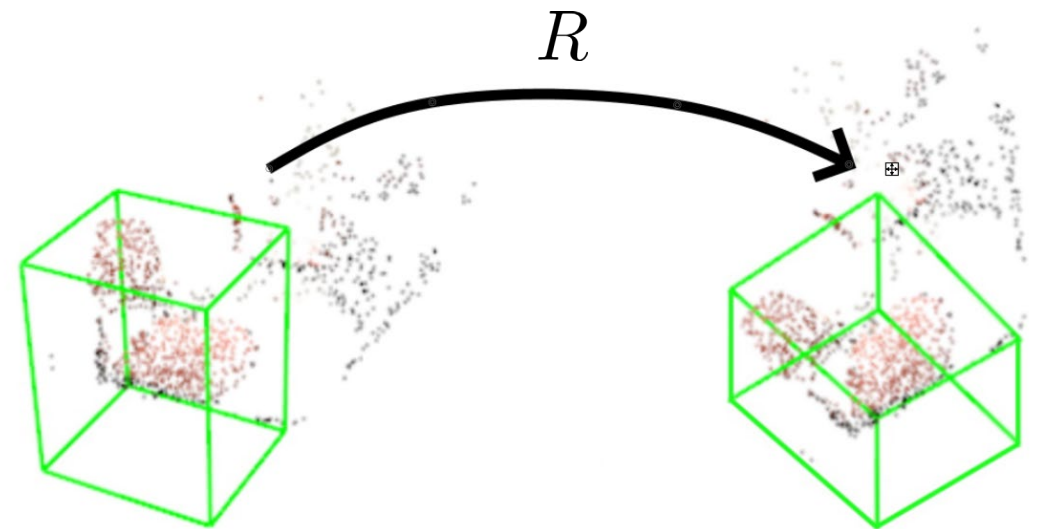
$$g^0 \left(\begin{array}{c} \text{Chair} \end{array} \right) = \text{Chair}$$

Examples of type 1 Features

- We have $D^1(R) = R$, therefore type 1 features are 3D vectors rotating with the pointcloud X .



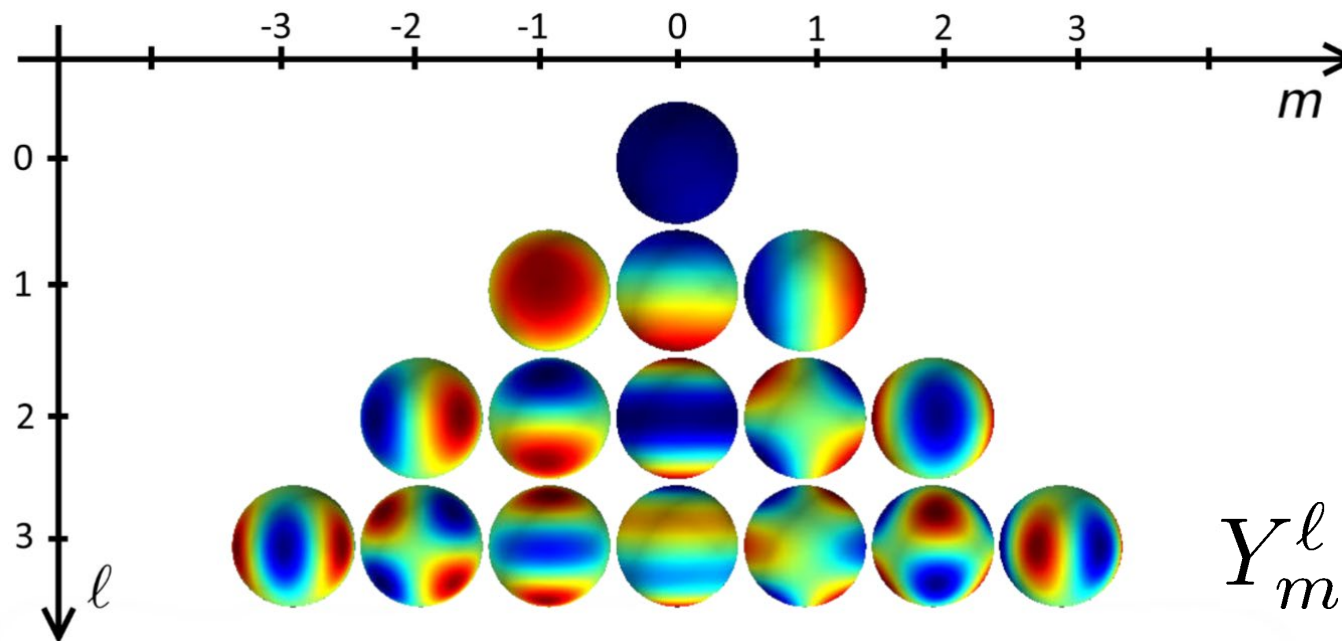
Pointcloud normals are type 1 features.



Bounding box center and principal directions are type 1 features, lengths are type 0 features

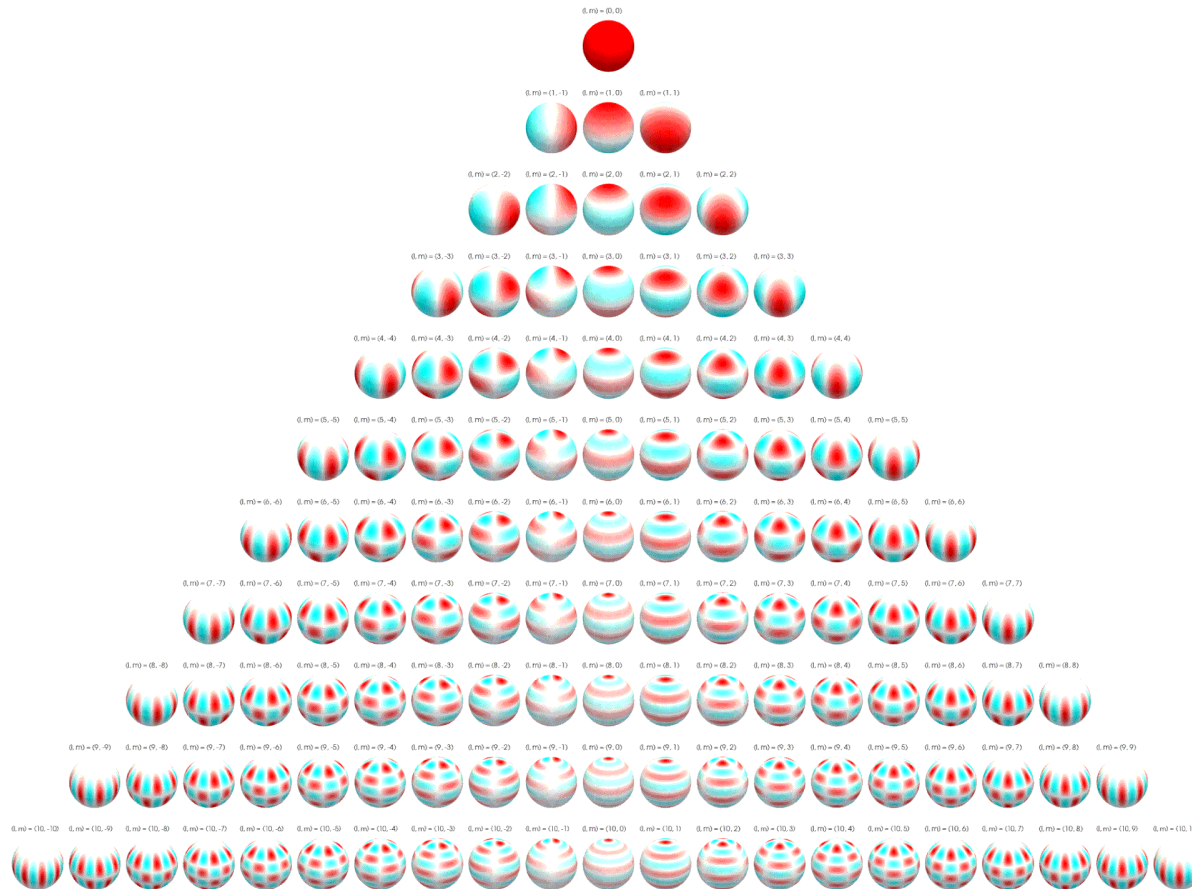
Spherical Harmonics & Higher Degree Features

- Spherical harmonics are homogeneous polynomials on \mathbb{R}^3 , their restriction to \mathcal{S}_2 form an orthonormal basis of $L^2(\mathcal{S}_2)$.
- Just like type ℓ equivariant features the vector of degree ℓ spherical harmonics $Y^\ell(x) \in \mathbb{R}^{2\ell+1}$, satisfies $Y^\ell(Rx) = D^\ell(R)Y^\ell(x)$.



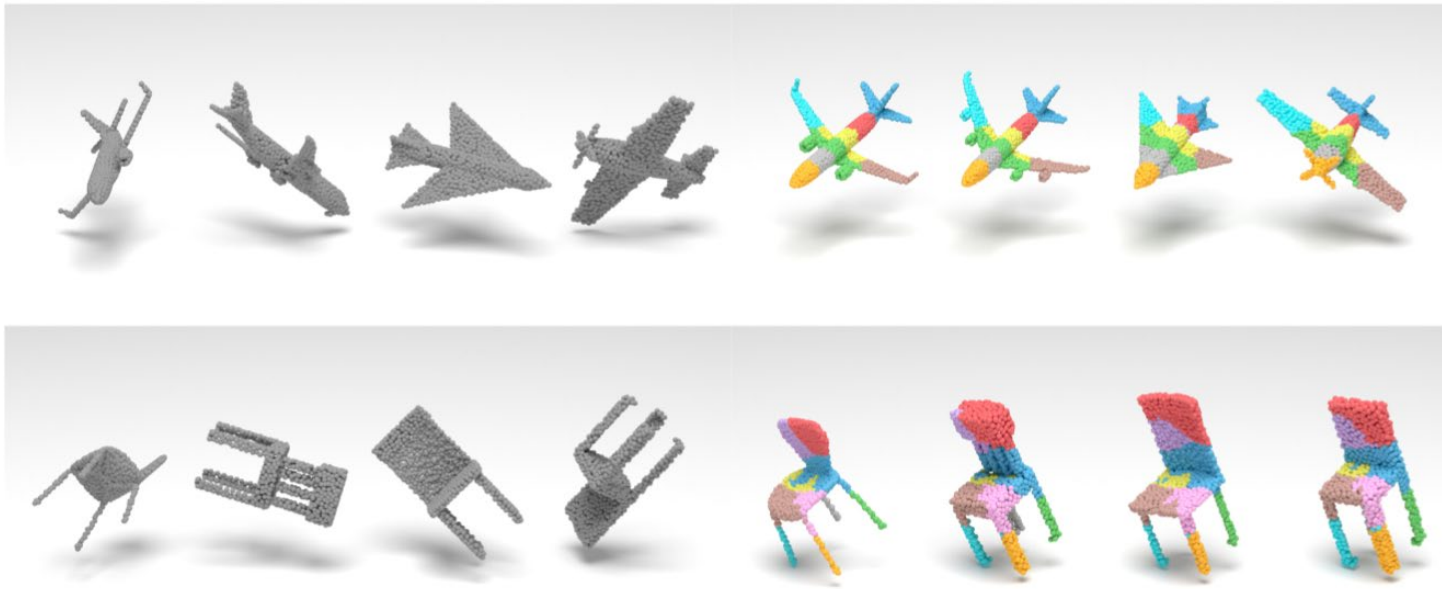
Spherical Harmonics & Higher Degree Features

- Spherical harmonics rotating around the z axis:



SH & TFN Feature Applications (ConDor)

- Dot products of type ℓ features are invariant.
- We can view TFN features as coefficients of functions in the SH basis.
- The dot prods $\langle f^\ell(X)_{ic}, Y^\ell(X_i) \rangle$ give an invariant embedding of X .

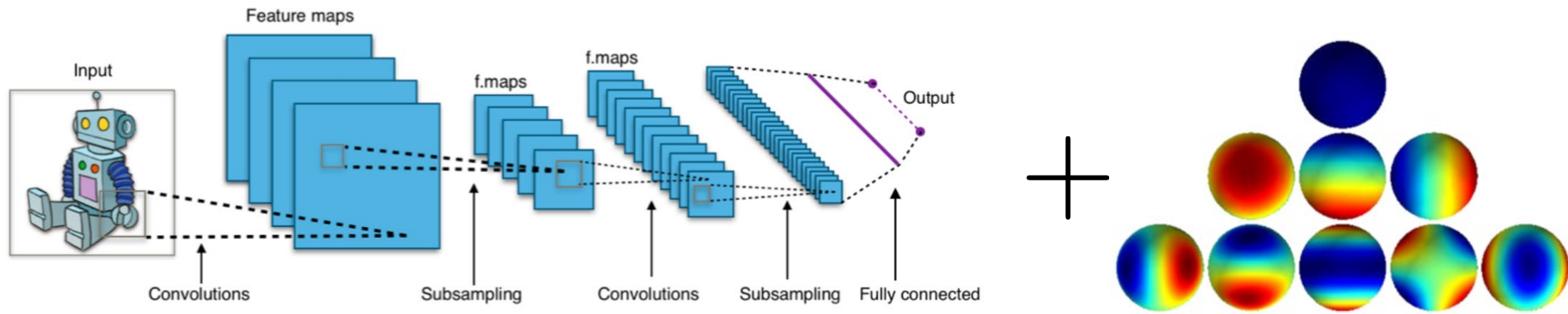


- Invariant embeddings can be used to:
 1. canonicalize shape pose.
 2. segment the shape(see ConDor).

How does TFN work ?

- TFN is a convolutional architecture.
- It inherits its equivariance properties from SH kernels.

TFN =



Convolutions on Different Domains

- The convolution between a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and a kernel $\kappa : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined by:

$$f * \kappa(x) := \int_{\mathbb{R}^d} f(t) \kappa(t - x) dt.$$

- Similarly the convolution a function f over a pointcloud X and a kernel $\kappa : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined by:

$$f *_X \kappa(X_i) := \sum_j f_j \kappa(X_j - X_i).$$

- We can view images as functions over the pointcloud \mathbb{Z}^d ($d = 2, 3$).

Convolutional Layers in Neural Nets

- Typically, we are given a collection of functions $(f_i)_i$ on a pointcloud X and a collection of kernel filters $(\kappa_j)_j$ and a weight tensor W_{kij} .
- A convolutional layer outputs a collection of functions by taking linear combinations of convolutions:

$$\text{Conv}_X(f, \kappa, W)_k := \sum_{ij} W_{kij} f_i *_{X} \kappa_j.$$

Convolutions and Rotations

- For any rotation R we have:

$$f *_{XR^\top} \kappa_j(X_i) = \sum_k f_k \kappa_j(R(X_k - X_i)).$$

- Problem what is the relation between $(f *_{X} \kappa_j)_j$ and $(f *_{XR^\top} \kappa_j)_j$?

Steerable Kernel Basis

- A steerable kernel basis, is a kernel basis $(\kappa_k)_k$ such that, the rotation of any kernel κ_j linearly decomposes onto the basis by a rotation matrix $D(R)$:

$$\kappa_j(Rx) = \sum_k D(R)_{jk} \kappa_k(x).$$

- Using a steerable basis we can relate convolutions on X and on XR^\top by a simple linear relation:

$$f *_{XR^\top} \kappa_j(x) = \sum_j D(R)_{jk} f *_X \kappa_k(x)$$

3D Steerable Basis

- For each ℓ we have a steerable basis of the form:

$$\kappa_{rm}^{\ell}(x) := \varphi_r(\|x\|_2) Y_m^{\ell} \left(\frac{x}{\|x\|_2} \right)$$

- Where φ_r is any radial function e.g. a gaussian shell:

$$\varphi_r(y) := \exp \left(\frac{-(y - \rho_r)^2}{2\sigma^2} \right)$$

Example of Steerable Basis (2D)

- For each ℓ in \mathbb{Z} we can build a steerable kernel basis on \mathbb{R}^2 using polar coordinates:

$$\kappa_r^\ell(\theta, \rho) := \varphi_r(\rho)e^{i\ell\theta}$$

here we have $D^\ell(t) = e^{i\ell t}$.

$r \setminus \ell$	0 Re	1 Re Im		2 Re Im		3 Re Im		4 Re Im		
3										...
2										...
1										
0										

Example of steerable kernels,
 ... (source Learning Steerable Filters
 with Rotation Equivariant CNNs)

Steerable Convolution

- We consider a collection of steerable kernel bases $(\kappa_r^\ell)_{lr}$ over \mathbb{R}^d indexed by an equivariance type ℓ such that for any rotation $R \in \text{SO}(d)$ we have:

$$\kappa^\ell(Rx) = D^\ell(R)\kappa^\ell(x)$$

- For each ℓ we define the associated steerable convolution operator:

$$\text{SConv}_X^\ell(f, \kappa_r^\ell, W)_{km} := \sum_{ir} W_{kir} f_i *_X \kappa_{rm}^\ell.$$

Comparing Regular and Steerable Convolution

- Regular convolution takes linear combination over all indices of the kernel basis:

$$\text{Conv}_X(f, \kappa, W)_k := \sum_{ij} W_{kij} f_i *_X \kappa_j.$$

- Steerable convolution does not take linear combination over the "steerable" index m in order to preserve equivariance:

$$\text{SConv}_X^\ell(f, \kappa_r^\ell, W)_{km} := \sum_{ir} W_{kir} f_i *_X \kappa_{rm}^\ell.$$

Steerable Convolution is Equivariant

- A rotation R of the pointcloud X induces a linear transform $D^\ell(R)$ which is independent of X :

$$\begin{aligned}\text{SConv}_{X R^\top}^\ell(f, \kappa_r^\ell, W)_k &= \sum_{ir} W_{kir} f_i *_{X R^\top} \kappa_r^\ell \\ &= \sum_{ir} W_{kir} f_i *_X D^\ell(R) \kappa_r^\ell \\ &= D^\ell(R) \text{SConv}_X^\ell(f, \kappa_r^\ell, W)_k\end{aligned}$$

Why not Mix Different Types?

- In 2D a type ℓ features makes ℓ turns when the input makes one turn:

$$v(e^{i\theta} X) = e^{i\ell\theta} v(X).$$

- There is no single rotation that rotates different types, or a linear combination of different types.

Stacking Steerable Convolutions

- We have described steerable convolution of rotation invariant pointwise features with a steerable basis.
- TFN extends steerable convolution to equivariant inputs, by linearly decomposing the resulting convolutional features into equivariant features of different types.

Equivariant non-Linearities

- Dot products and norms of type ℓ features are rotation invariant:

$$\langle f(R.X), g(R.X) \rangle = \langle D^\ell(R)f(X), D^\ell(R)g(X) \rangle = \langle f(X), g(X) \rangle.$$

- Non linearities can be applied to the norms of features as they are invariant:

$$\xi(f^\ell(X), b) = \xi(\|f^\ell(X)\|_2 + b)f^\ell(X).$$

J. STOLFI
1-89

