

# Exploring explainable models for point cloud learning

Weiwei Sun

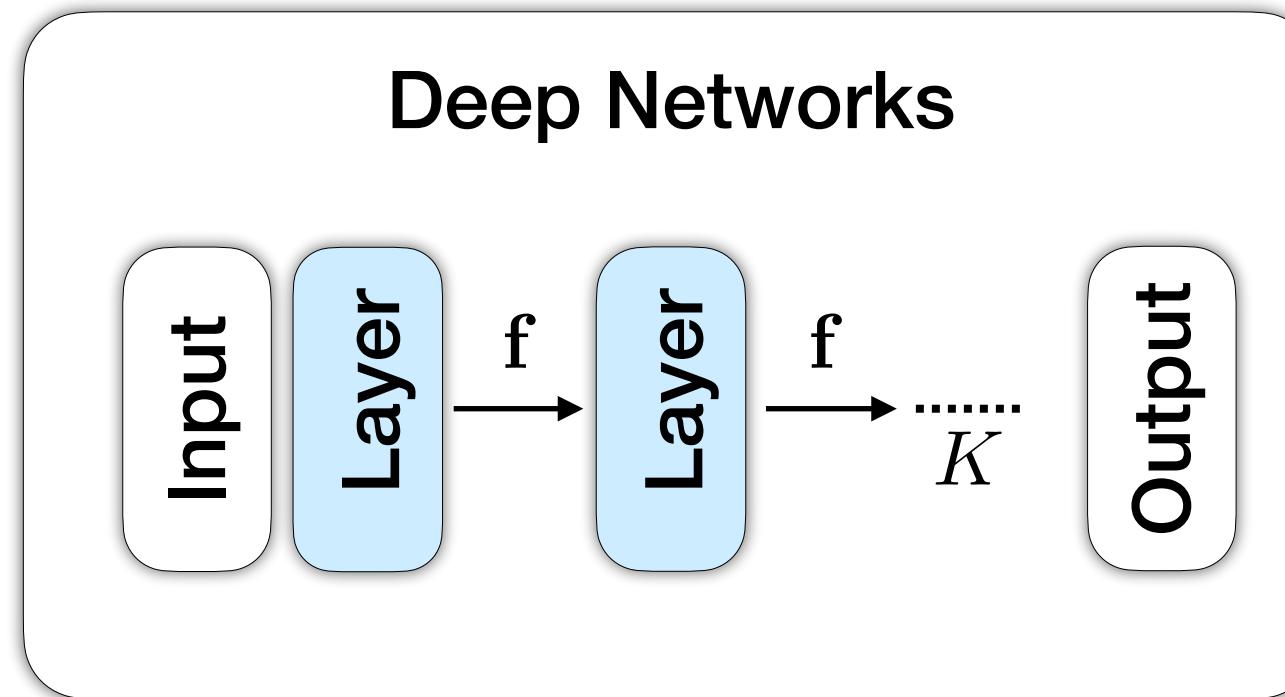
5th Ph.D. student @ UBC

under Dr. Andrea Tagliasacchi & Dr. Kwang Moo Yi



THE UNIVERSITY  
OF BRITISH COLUMBIA

# Explainability is essential



A stack of *K layers* and a set of intermediate *features f*.

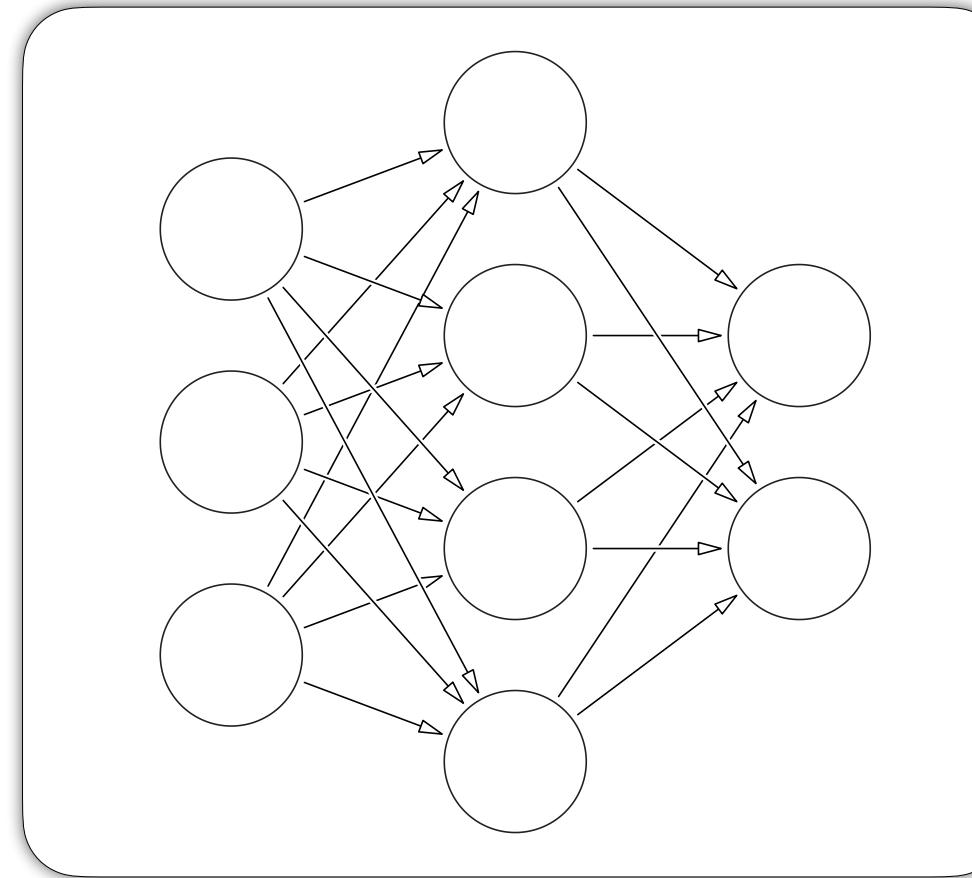
As a **black box**, deep networks work pretty good  
but lacks transparency, controllability, generalization ability.

How does deep networks **behave as human understand**? And, to what extent, we can **benefit from explainable models**?

# Point cloud learning

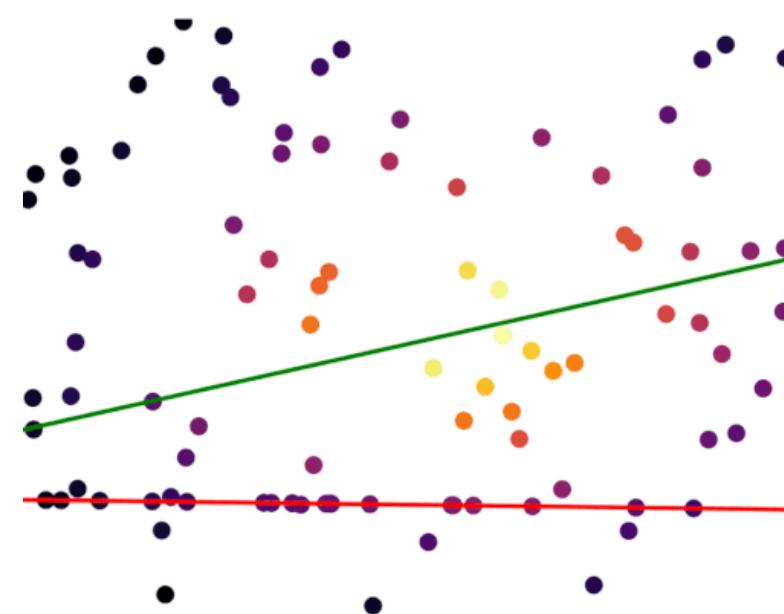
Permutation-equivariant point network

$$\mathbf{P} = \left\{ \begin{array}{l} \mathbf{p}_1 \\ \mathbf{p}_2 \\ \vdots \\ \mathbf{p}_N \end{array} \right\}$$



$$\mathbf{F} = \left\{ \begin{array}{l} \mathbf{f}_1 \\ \mathbf{f}_2 \\ \vdots \\ \mathbf{f}_N \end{array} \right\}$$

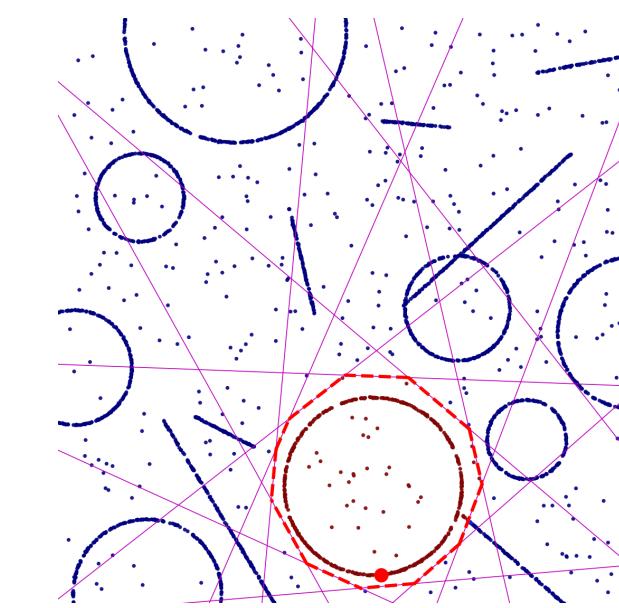
We present models explainable in either layer or features.



ACNe@CVPR2020

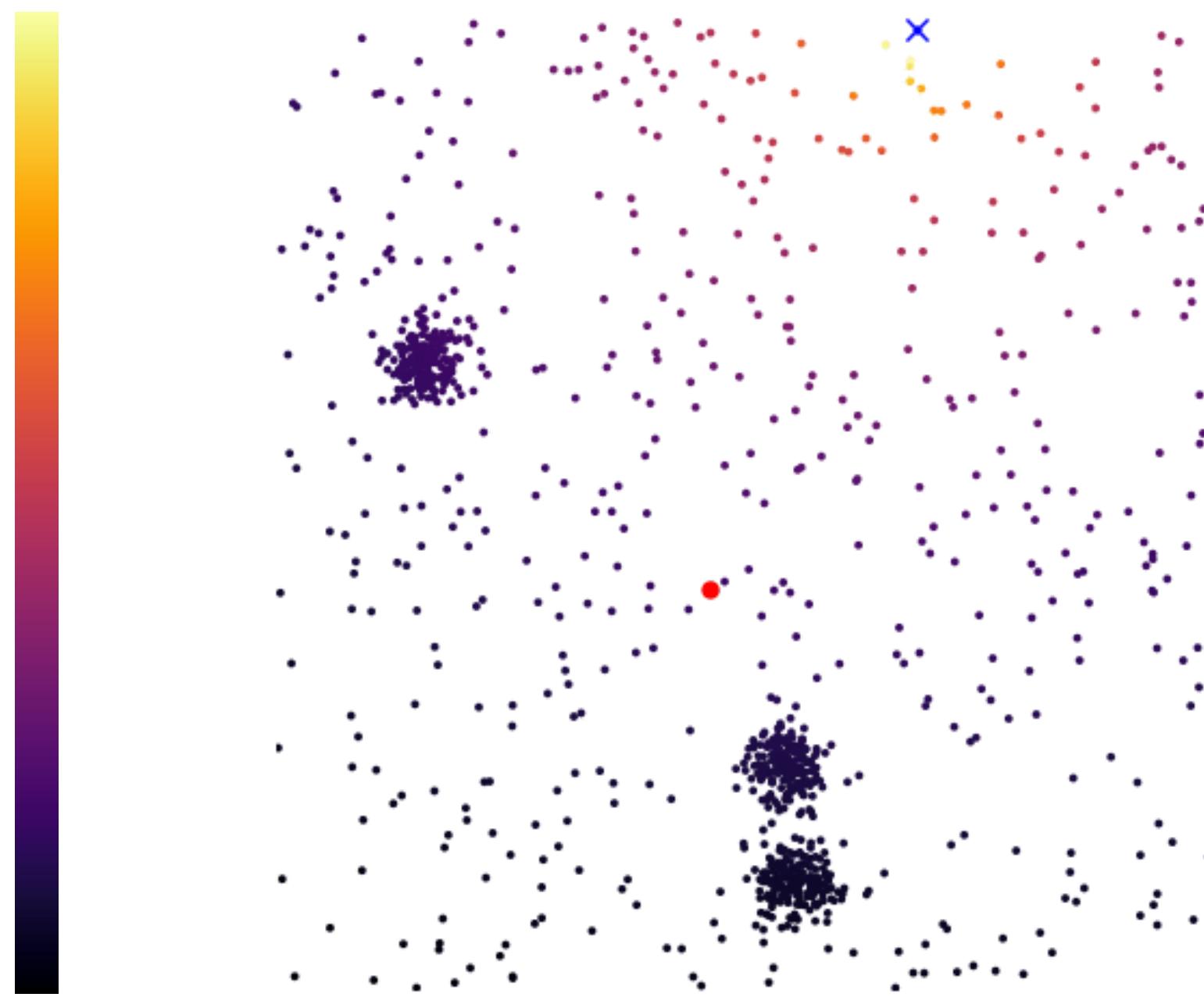


Canonical Capsules@NeurIPS2021



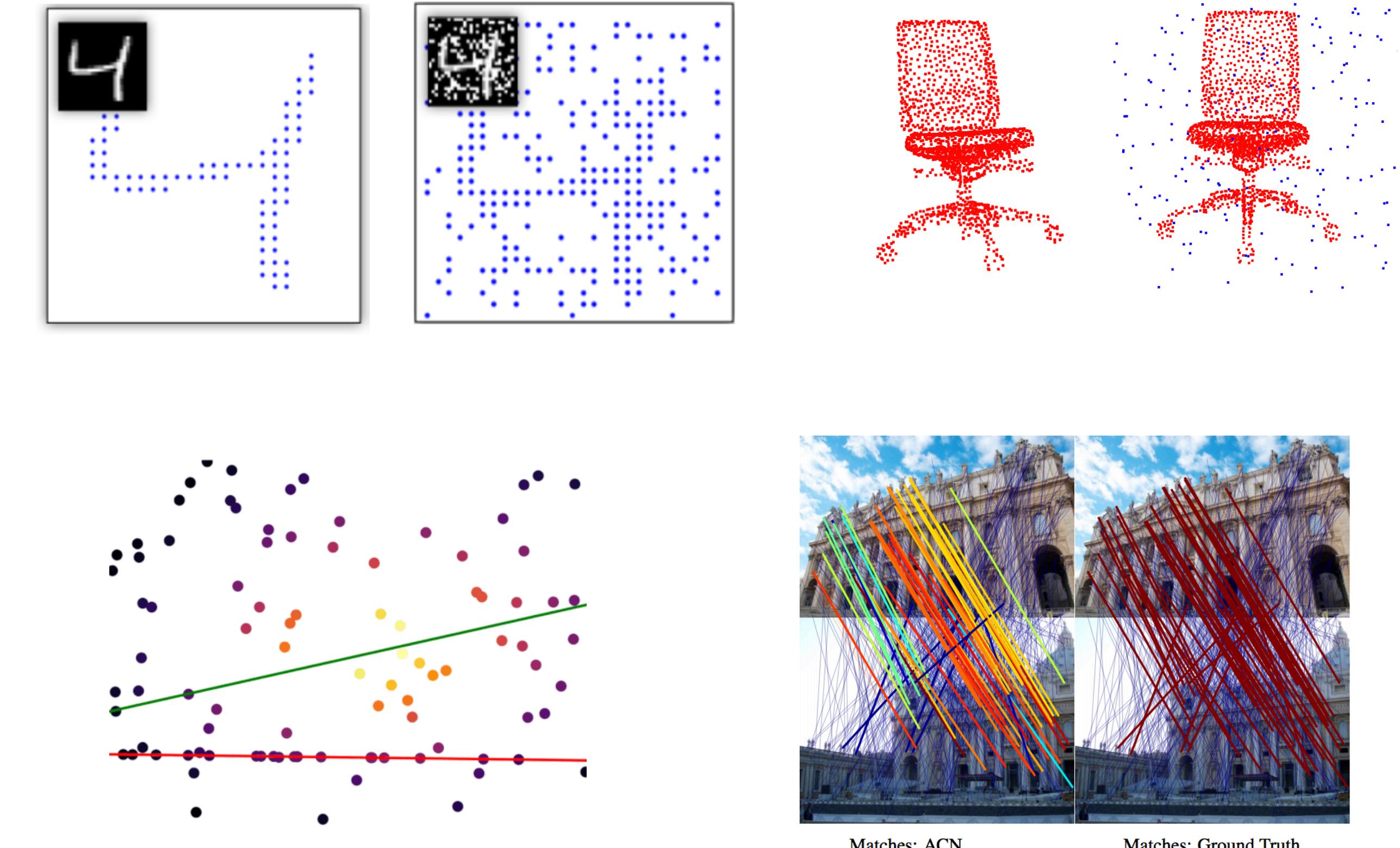
NeuralBF@WACV2023

# ACNe: unrolling optimization as deep network



- Mean
- ✖ Weighted mean

Inspired by IRLS — treating each step as a layer.



SOTA performance in robust point cloud learning tasks.

# Robust Learning of Point Clouds

Feature normalization relates points

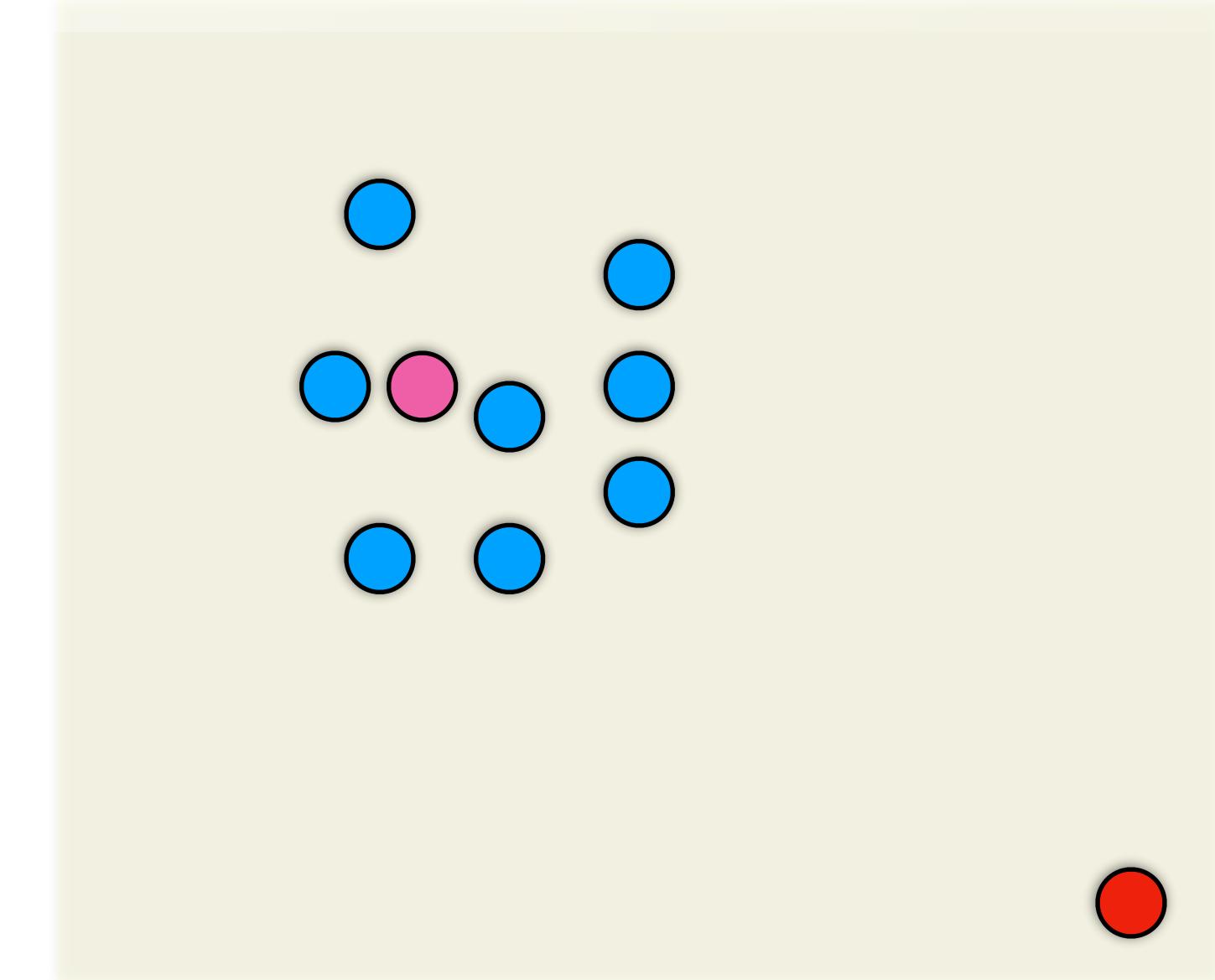
- Context normalization

$$\mathcal{N}_{CN}(\mathbf{f}) = (\mathbf{f} - \mu(\mathbf{f})) \oslash \sigma(\mathbf{f})$$

- Where mean and average can be rewritten as

$$\mu(\mathbf{f}) = \mathbb{E}[\mathbf{f}]$$

$$\sigma(\mathbf{f}) = \sqrt{\mathbb{E}[(\mathbf{f} - \mathbb{E}[\mathbf{f}])^{\circ 2}]}$$



- Mean point
- Input points
- Outlier

# Robust Learning of Point Clouds

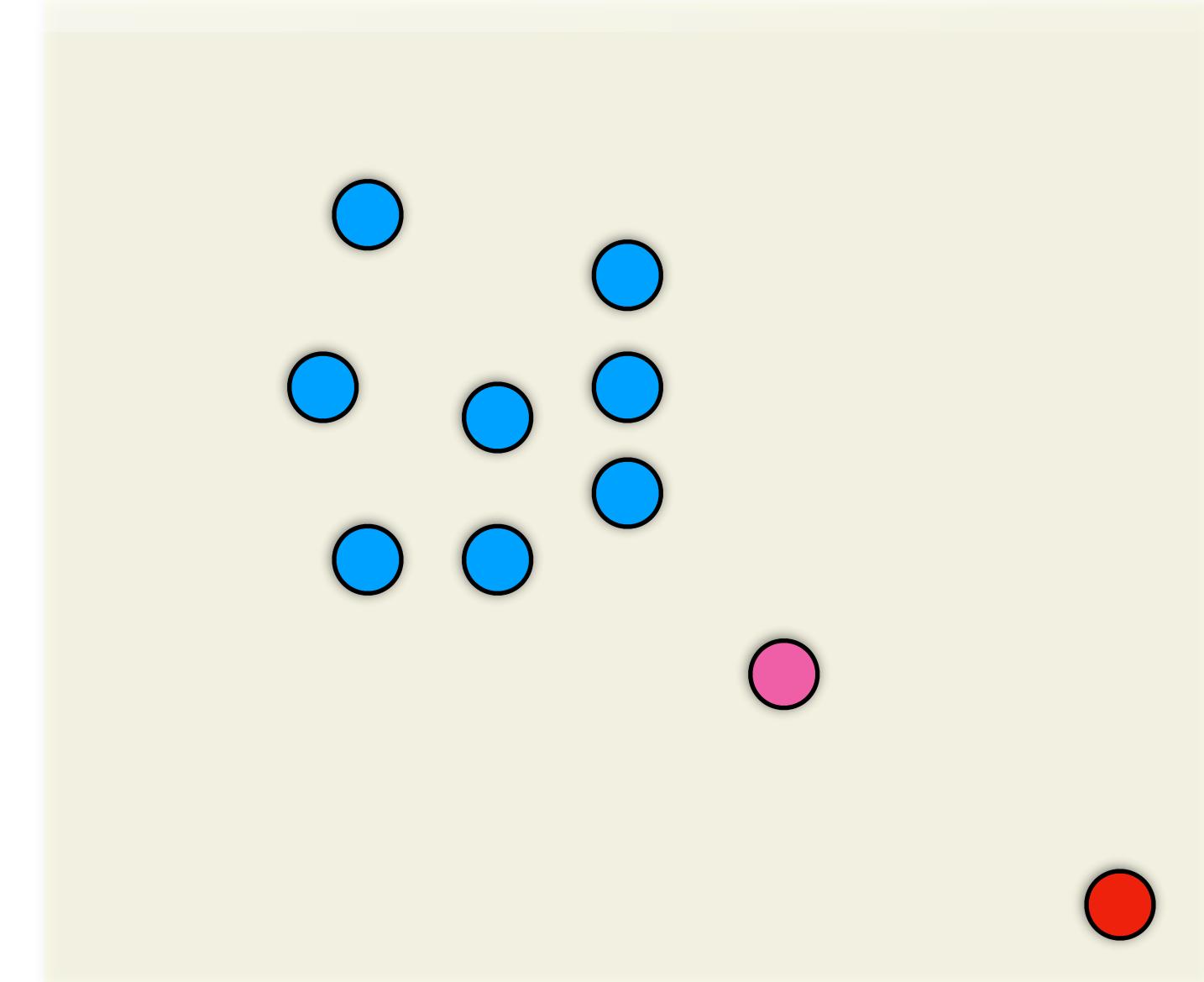
## Attentive Context Normalization (ACN)

- So we define our attention based on features

$$\mathbf{w} = \eta(\mathcal{W}_\omega(\mathbf{f})), \quad \eta(\mathbf{x}) = \mathbf{x}/\|\mathbf{x}\|_1.$$

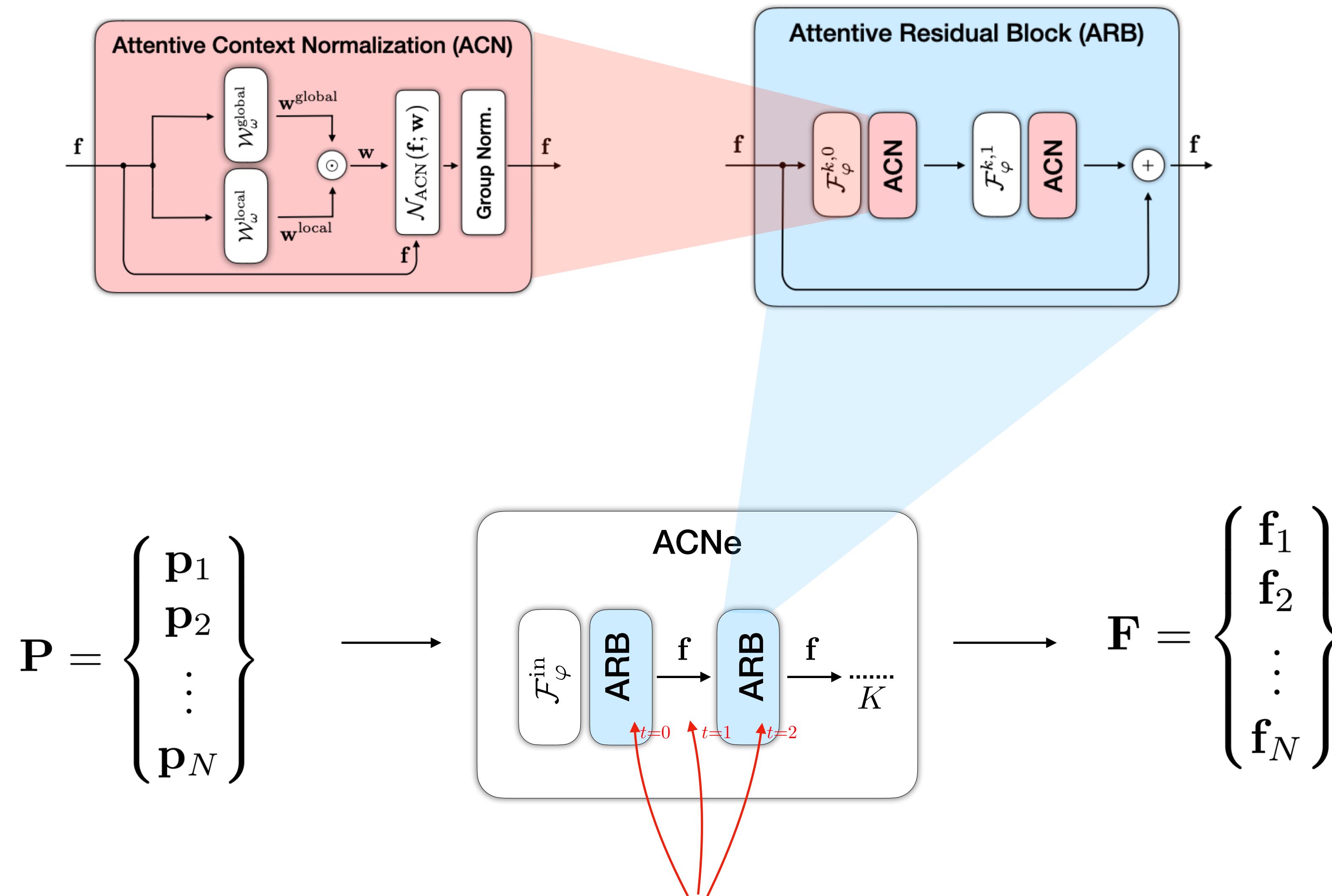
- And then define an attentive normalization

$$\mathcal{N}_{ACN}(\mathbf{f}; \mathbf{w}) = (\mathbf{f} - \mu_{\mathbf{w}}(\mathbf{f})) \oslash \sigma_{\mathbf{w}}(\mathbf{f}).$$



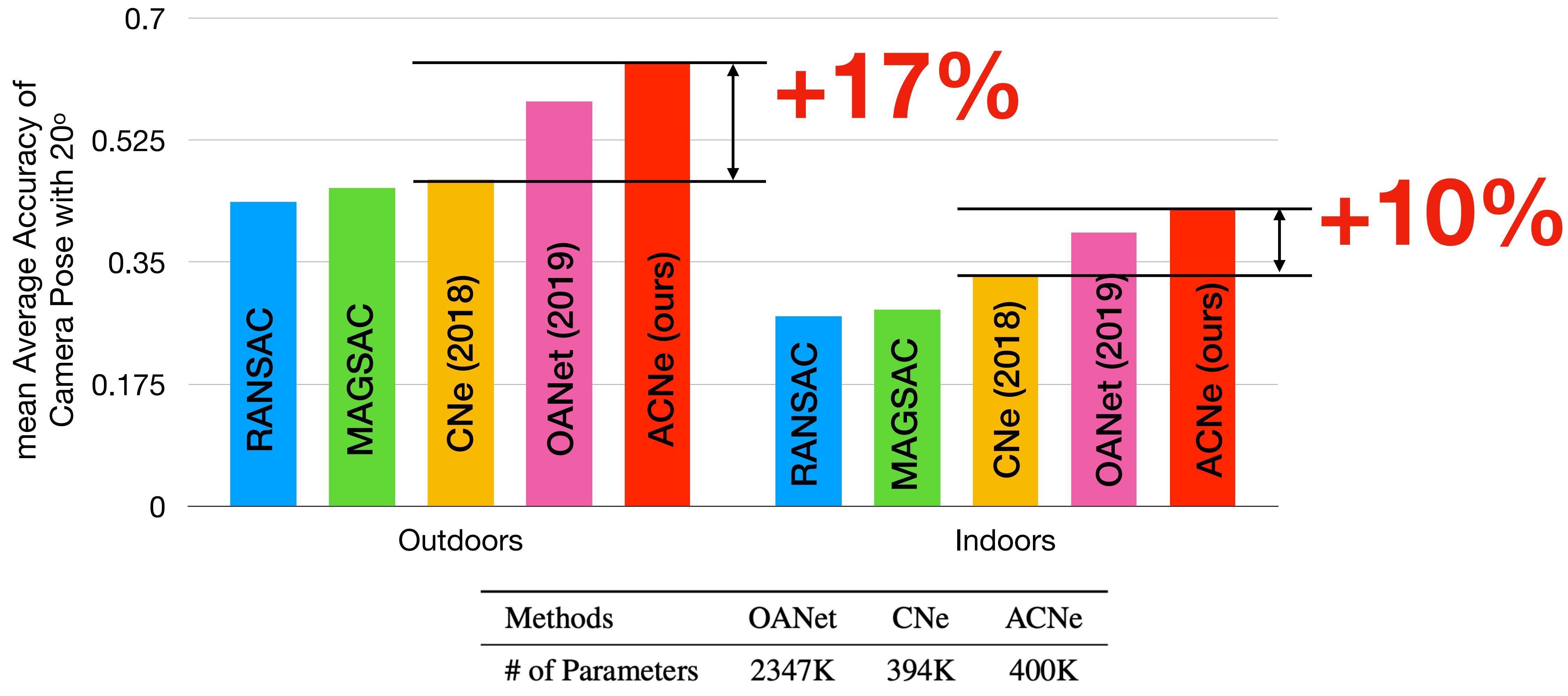
- Mean point
- Input points
- Outlier

# ACNe: an iterative network



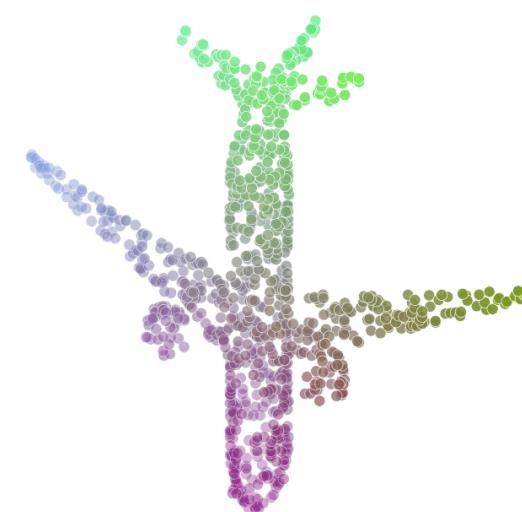
# ACNe: explainable model won with less parameters

In real-world wide-baseline stereo matching task, explainable ACNe achieves the best performance yet with less parameters.

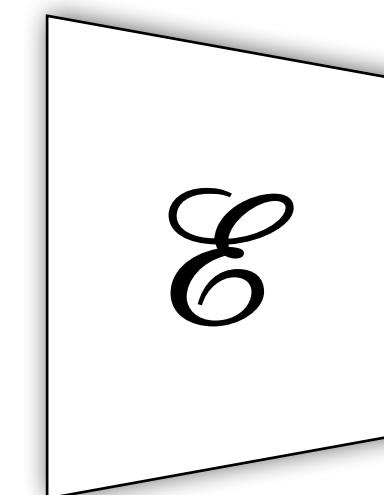


# Canonical Capsules: explainable features

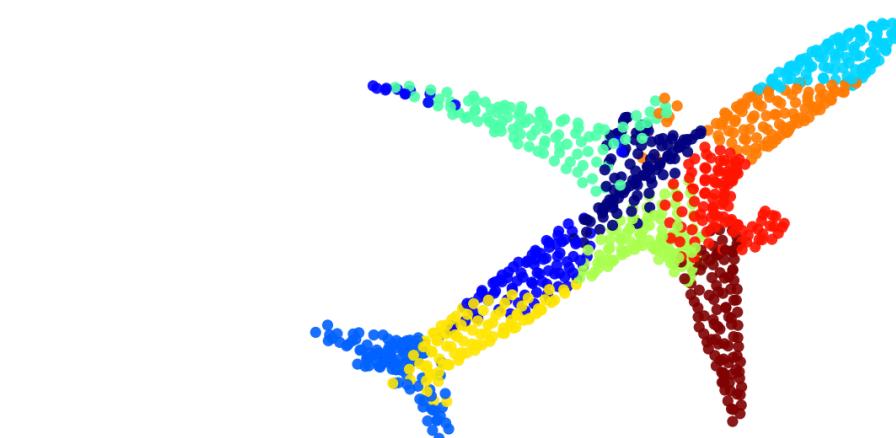
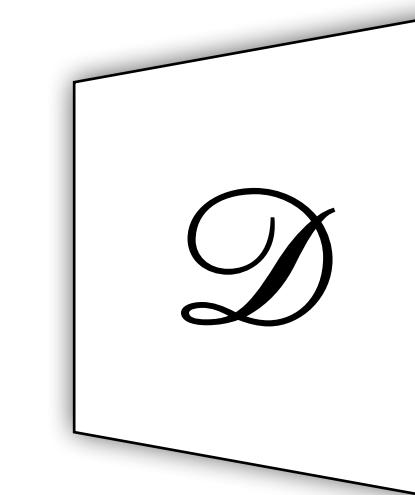
- A **self-supervised** training architecture for **point clouds** with:
  - SOTA performance in **auto-encoding**
  - SOTA performance in unsupervised **classification**
  - SOTA performance in learnt **canonicalization**



input point cloud  
(randomly transformed)



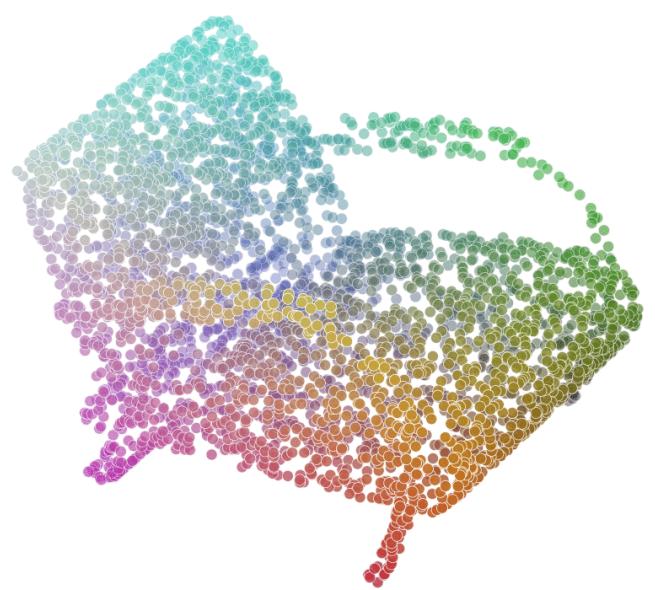
$\{ \text{red capsule, blue capsule, orange capsule, cyan capsule, green capsule} \}$   
set of capsules



auto-encoding  
(learnt canonical frame)

# Built-in biases in 3D datasets

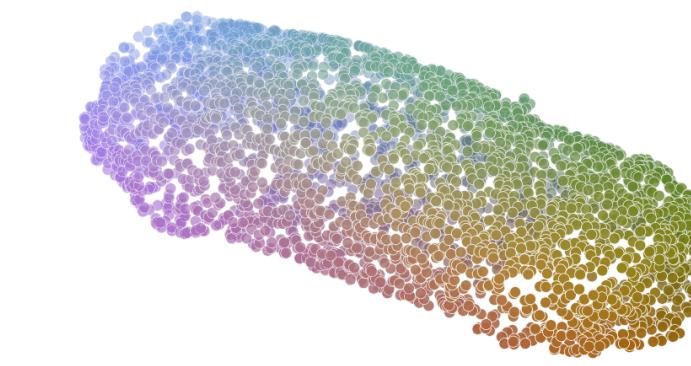
chairs



airplane



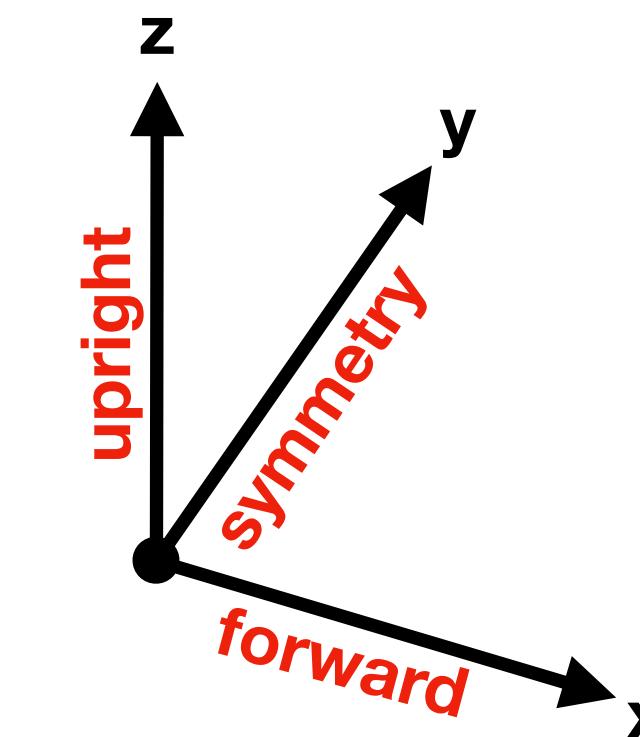
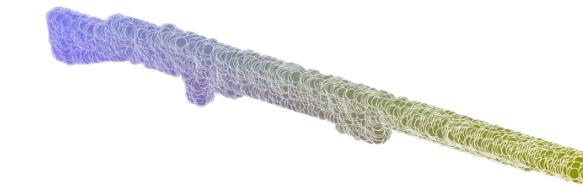
car



table

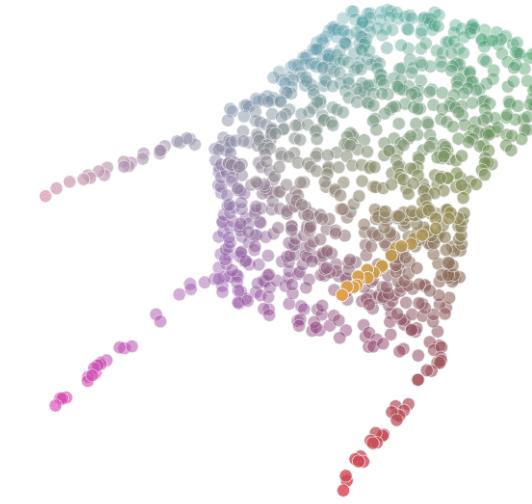


firearm

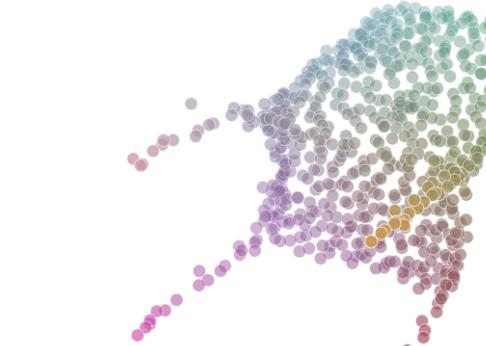
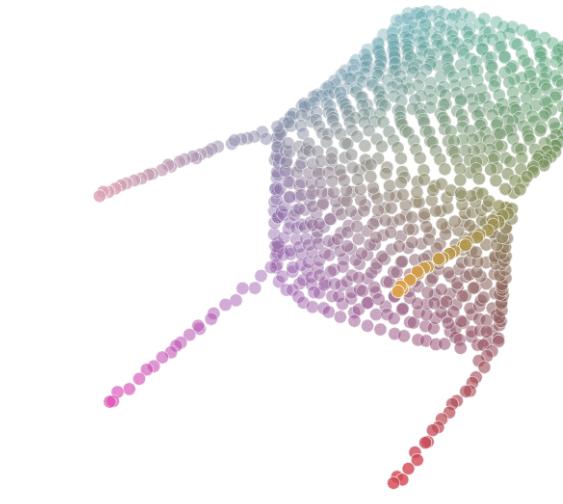
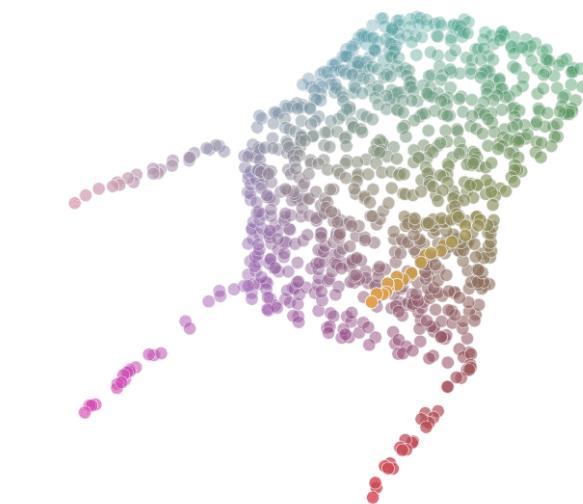


# Auto-encoding Point Clouds

ShapeNet  
aligned data



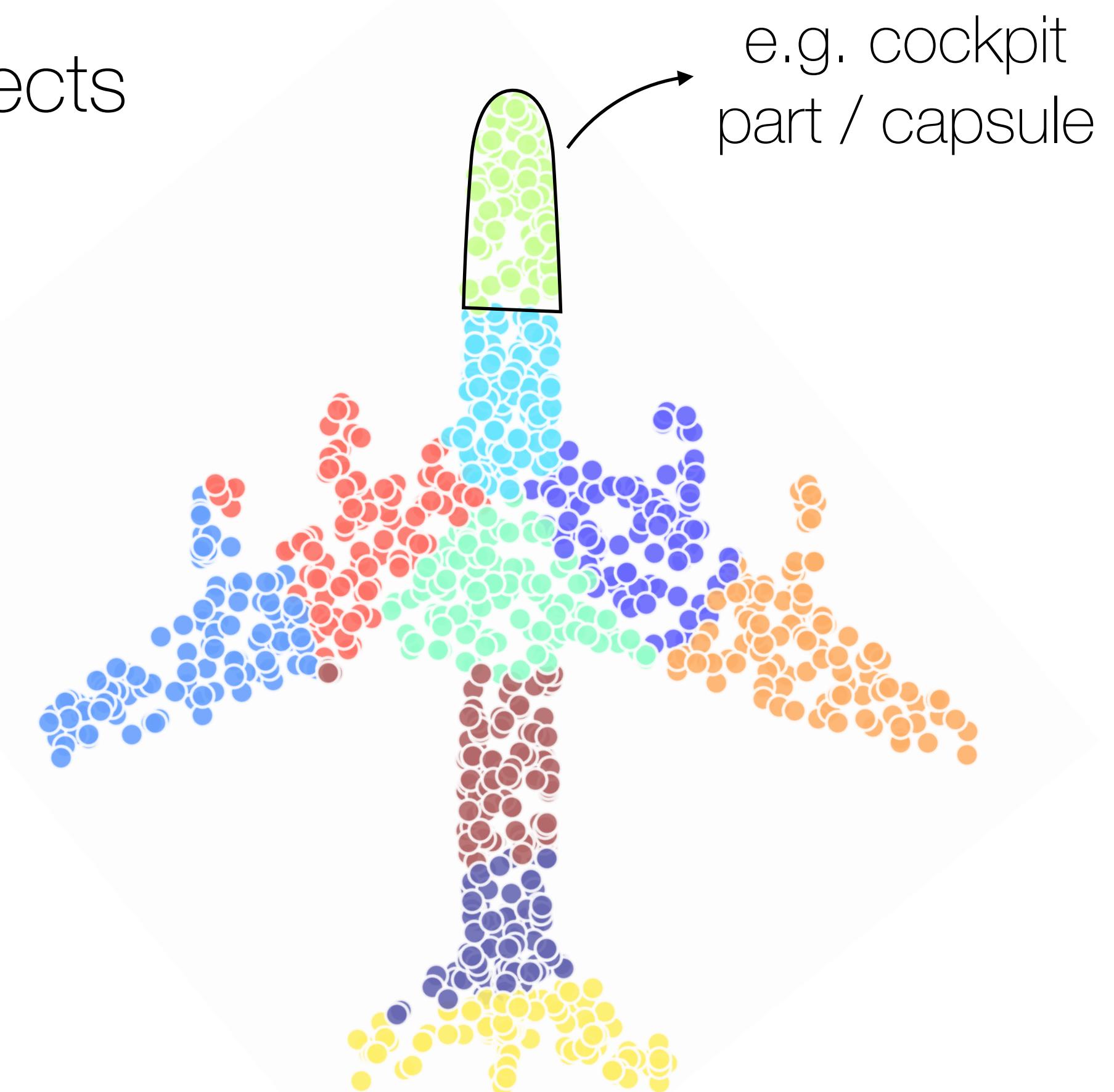
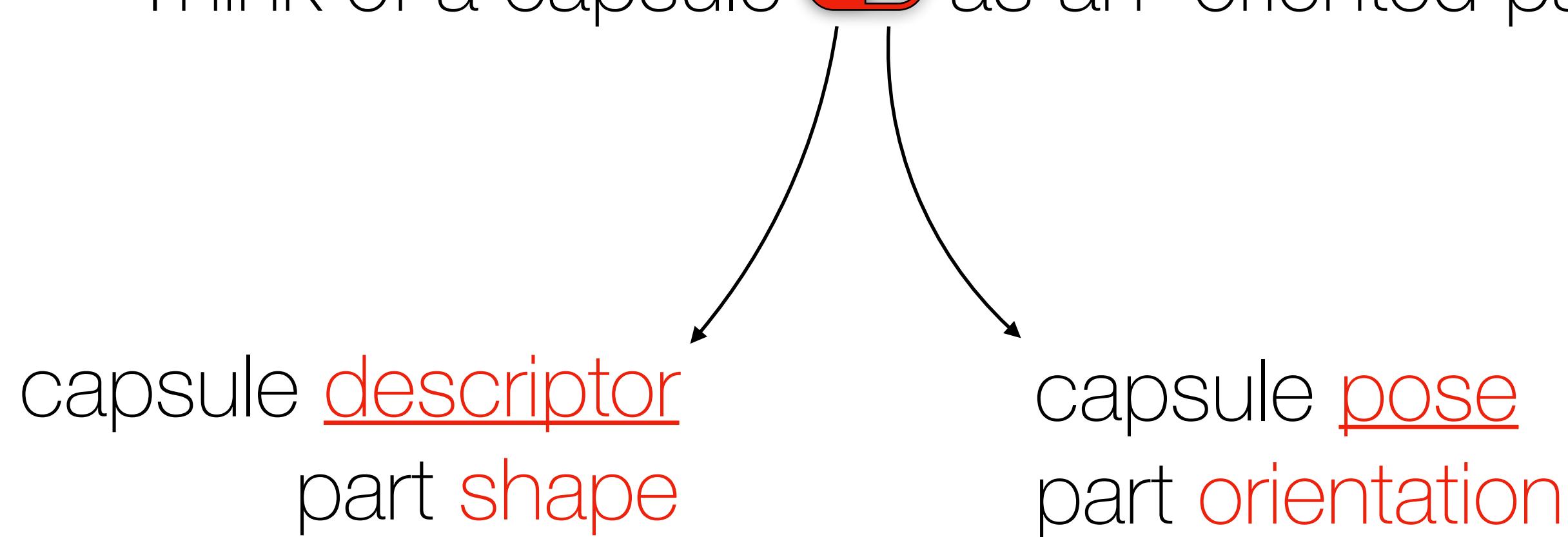
ShapeNet  
non-aligned data



performance  
drop

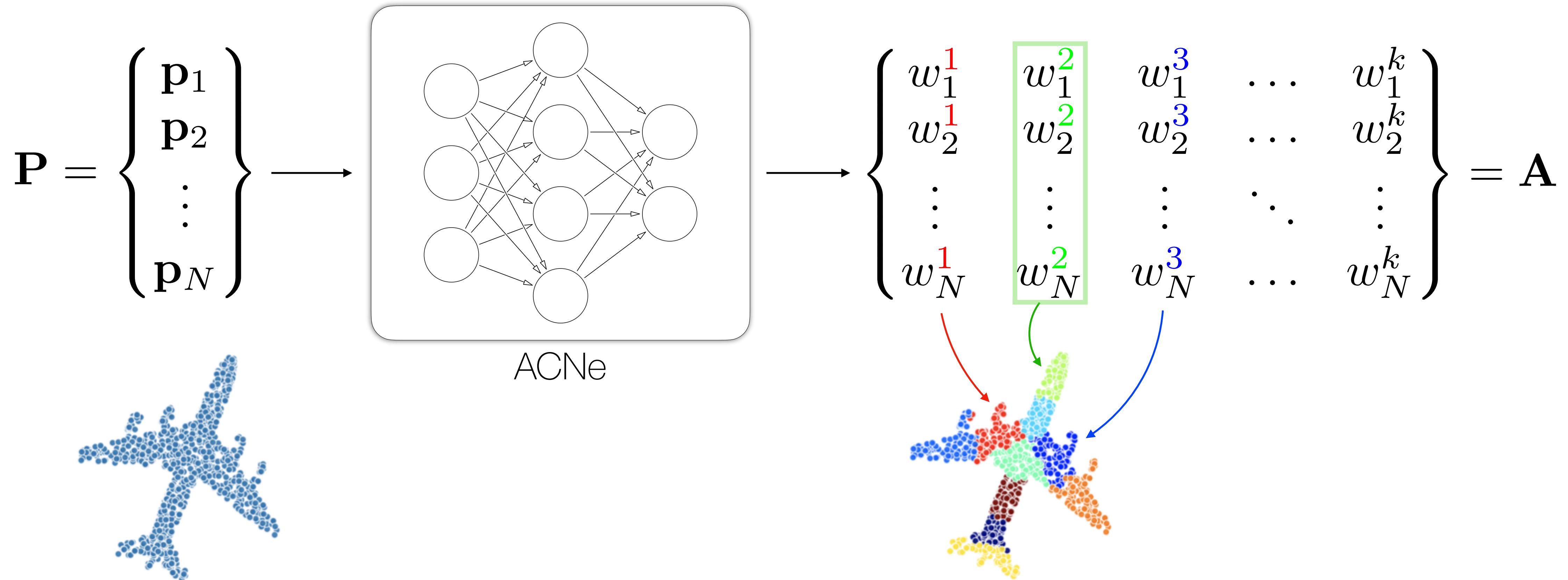
# But... what are capsules?

- An abstraction of **scene graphs** from computer graphics
- Capsules are **part-based** representations of objects
- Think of a capsule  as an “oriented part”



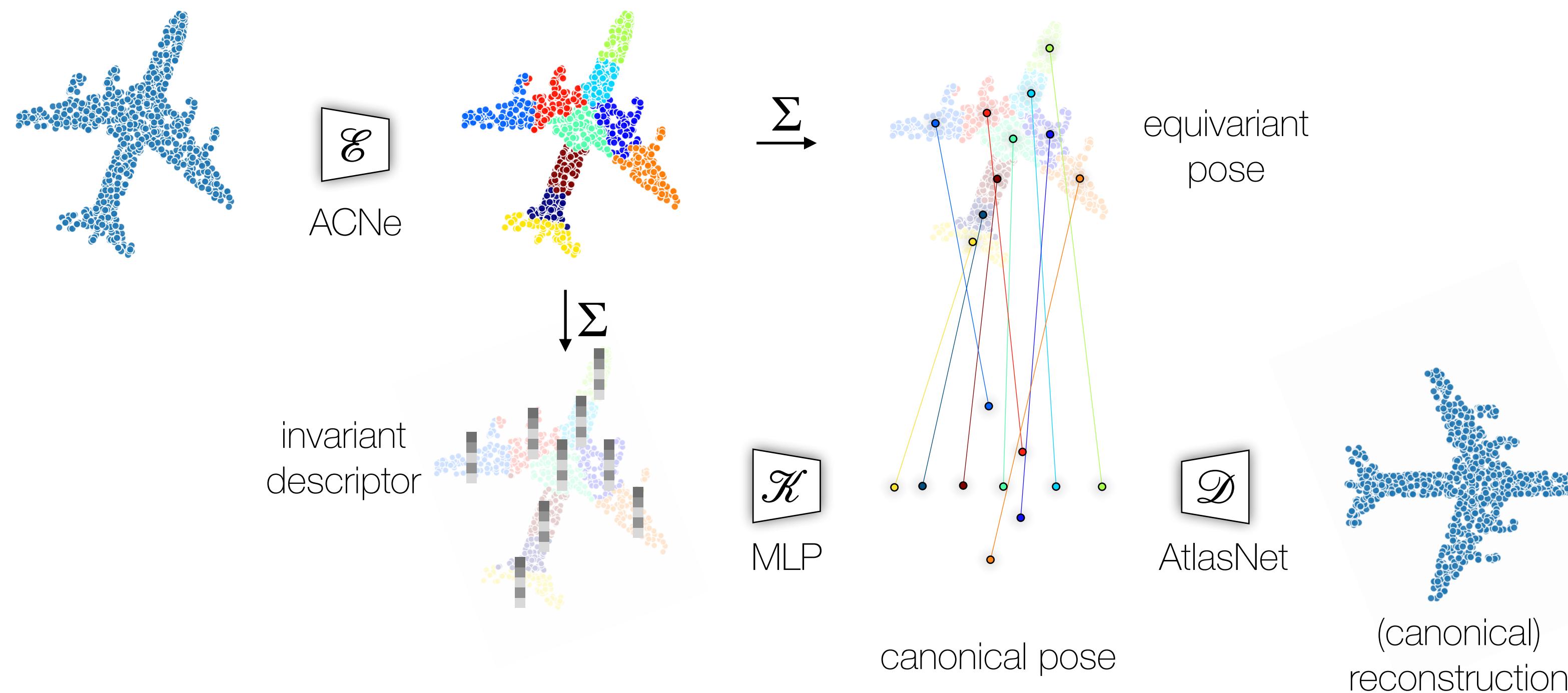
Credit to: Hinton et al. NeurIPS 2011 «Transforming Auto-encoders»

# Multi-head attention (ACNe)

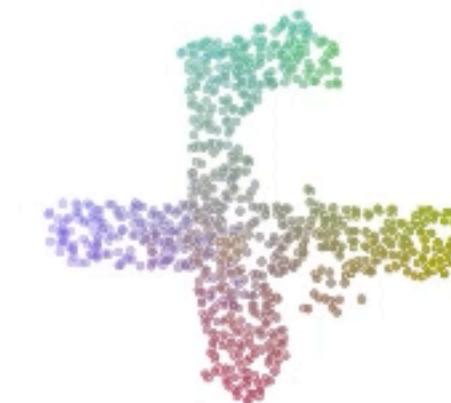
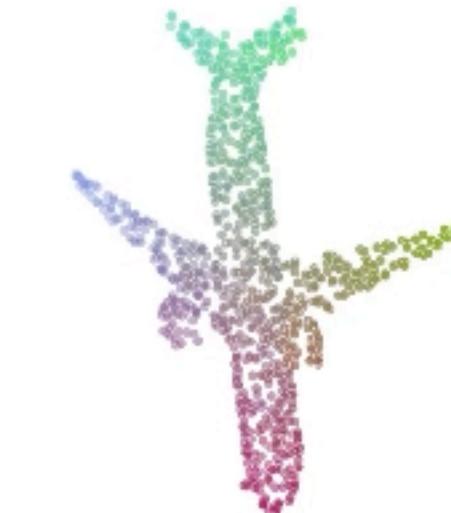
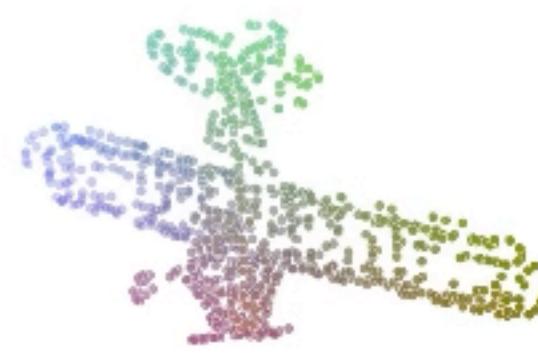


# Inference Procedure

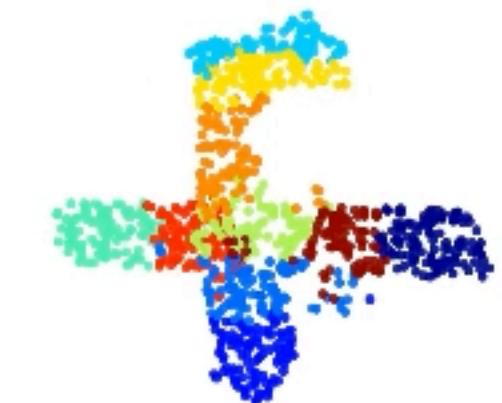
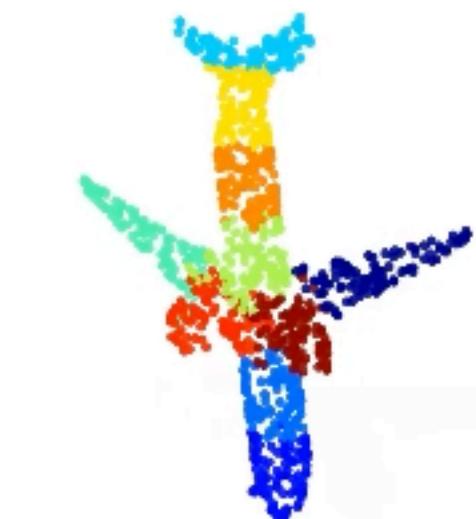
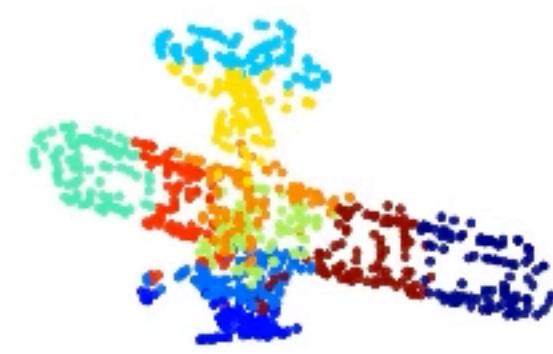
- Let us look at how the **inference** (forward pass) of our network architecture
  - Note the **input** is a point cloud in **any pose** (random rotation/translation)
  - While the **reconstruction** is executed in the **canonical frame**



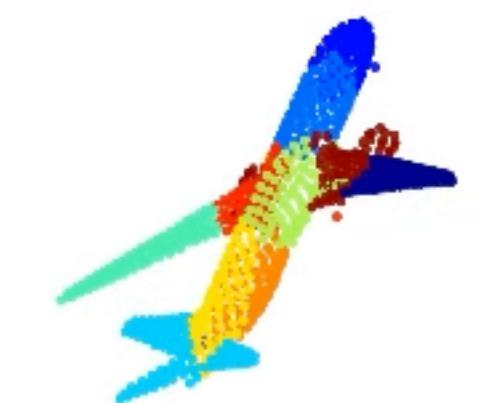
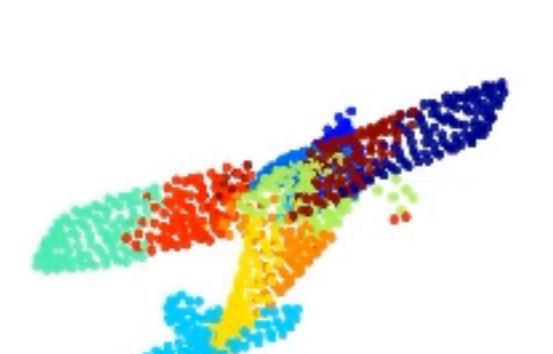
# Auto Encoding – Controllability and better performance



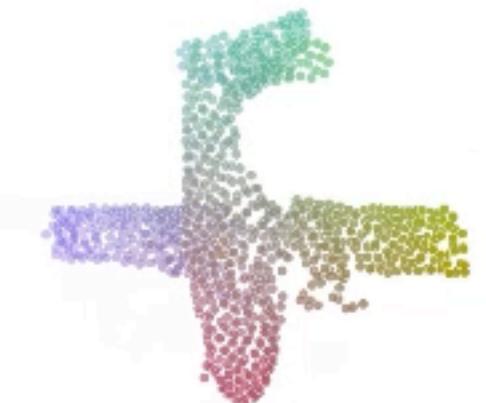
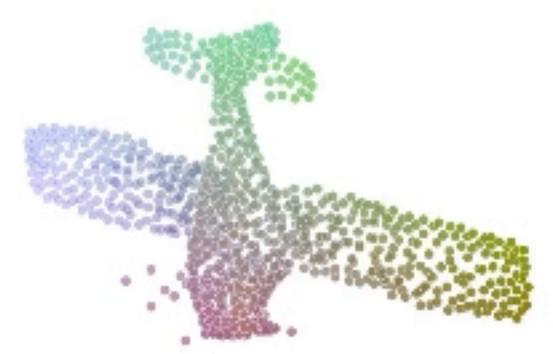
Input  
(Point Cloud)



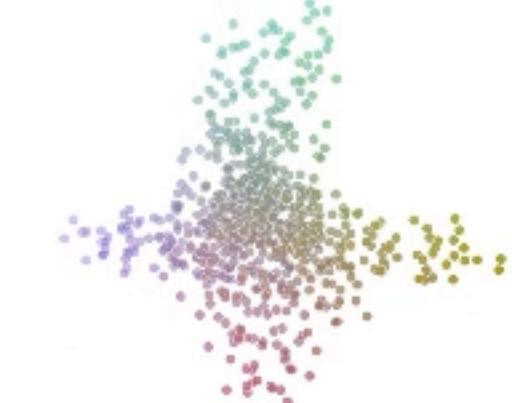
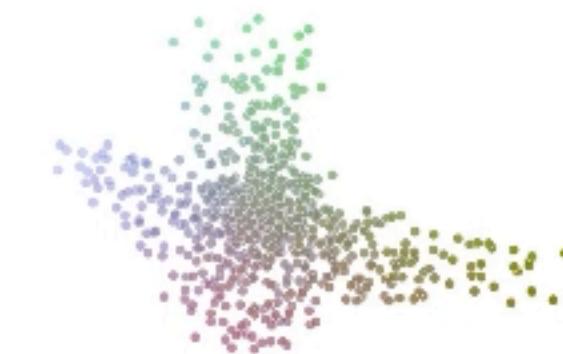
Decomposition  
(original frame)



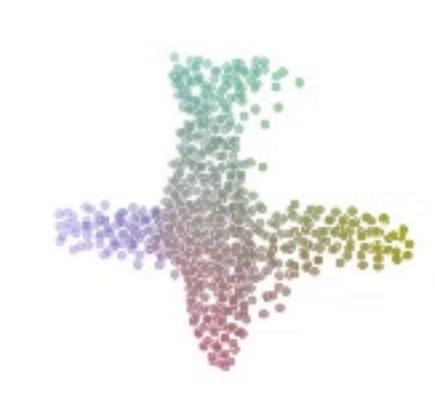
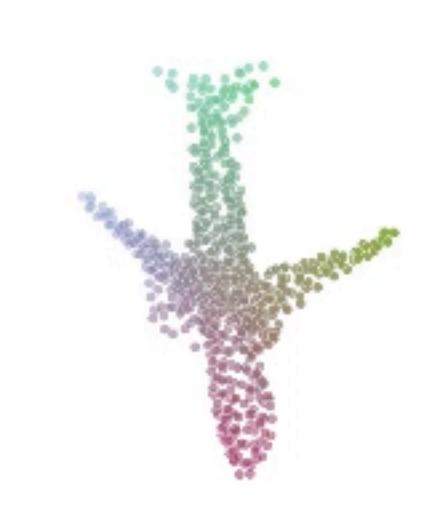
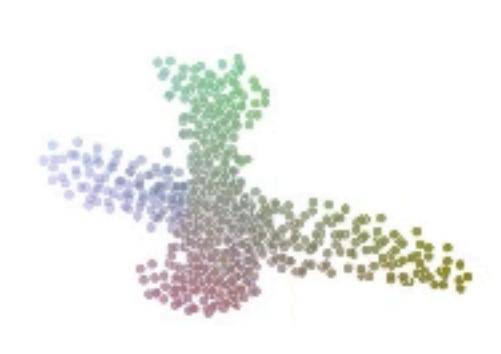
Reconstruction  
(canonical frame)



Reconstruction  
(original frame)



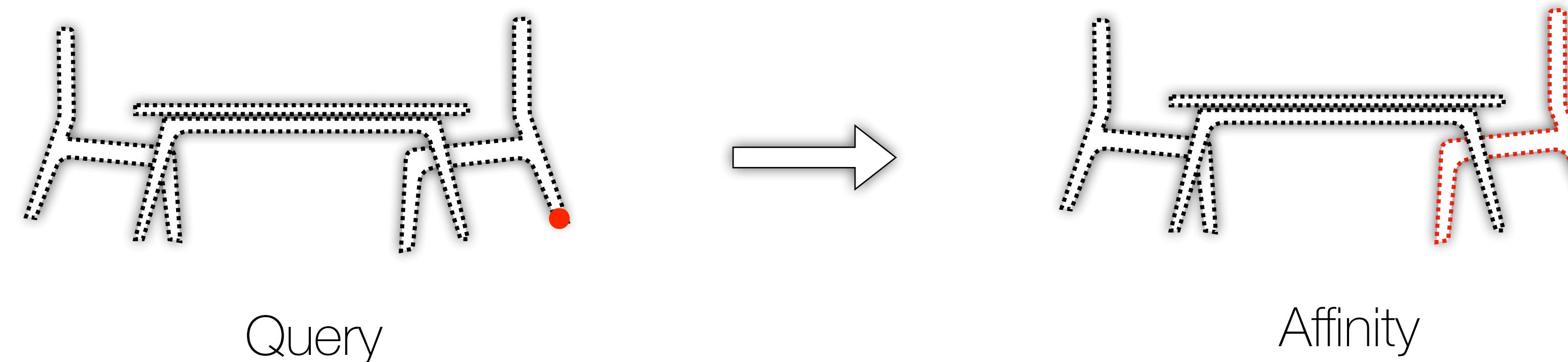
PointCapsNet  
(CVPR 2019)



AtlasNetV2  
(NeurIPS 2019)

# NeuralBF: Explainable feature and layer

- An **iterative** layer that
  - is inspired by traditional **bilateral filtering**
  - generates **query-conditioned** instance proposals for 3D point clouds
  - gives state-of-the-art **top-down** 3D instance segmentation.



# Direct bounding box regression

- A commonly used method for instance proposal generation
- Widely adopted for top-down 2D instance segmentation
- However, bottleneck the performance of top-down 3D instance segmentation

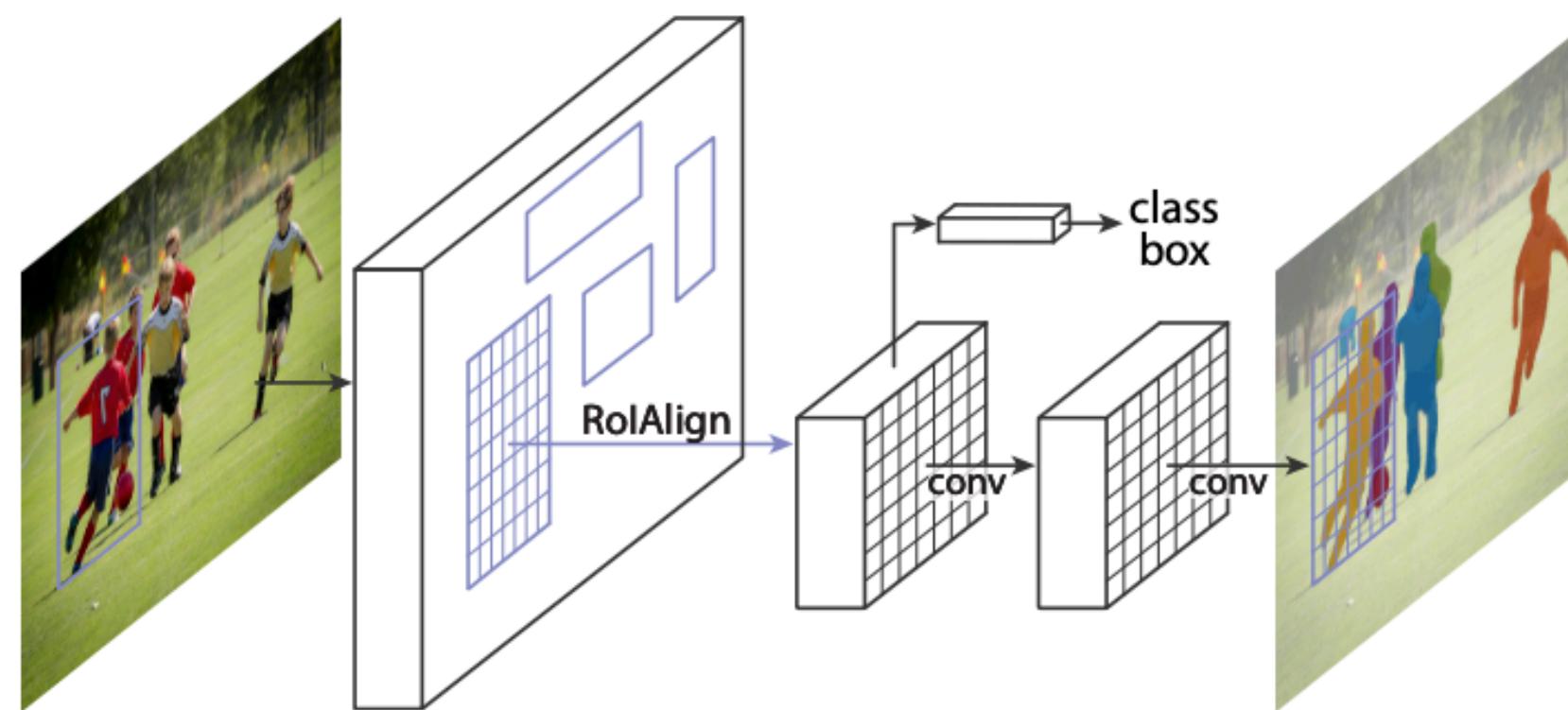
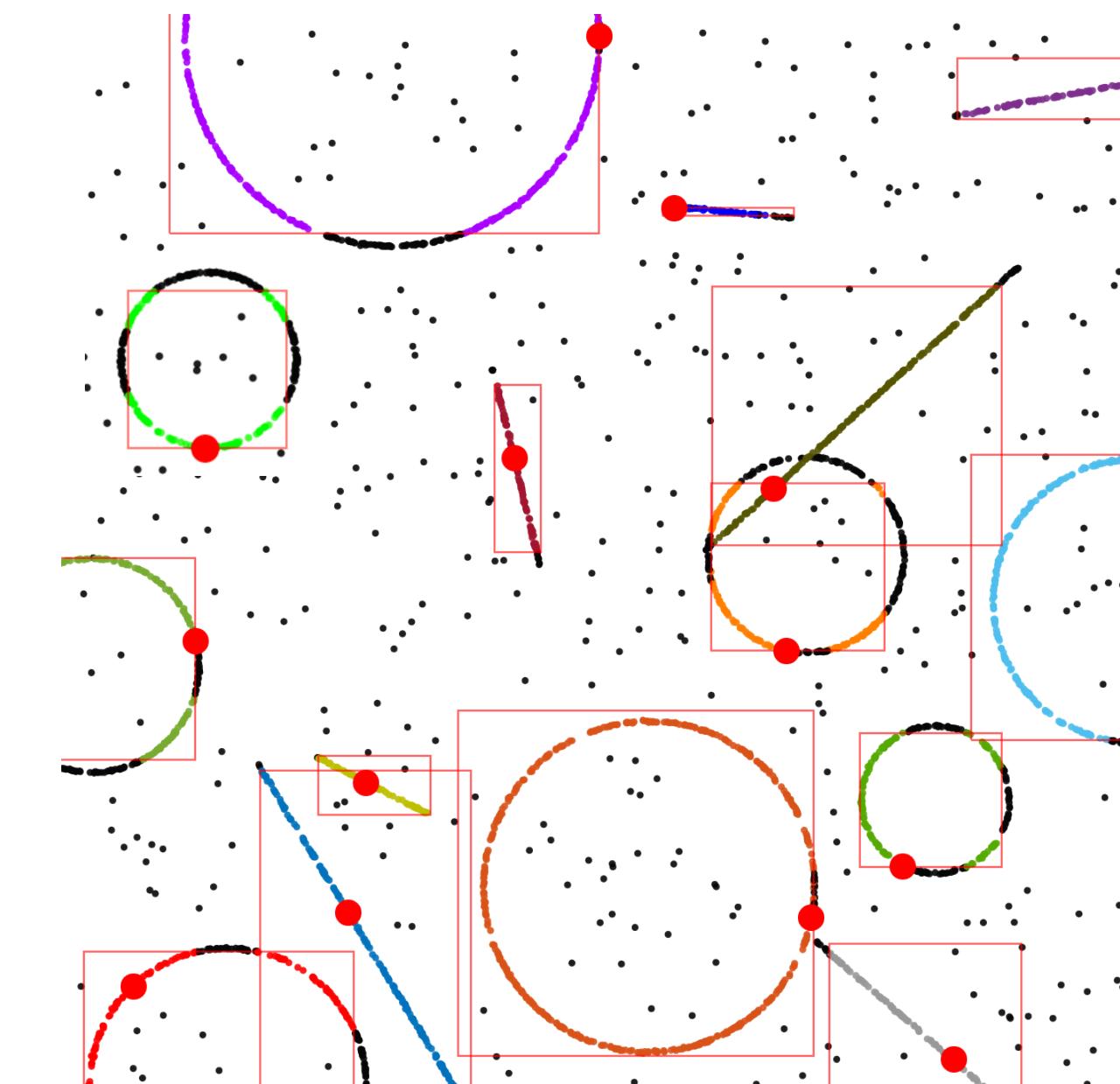


Figure 1. The **Mask R-CNN** framework for instance segmentation.

He et al. Mask R-CNN. ICCV 17.

Great success in 2D images

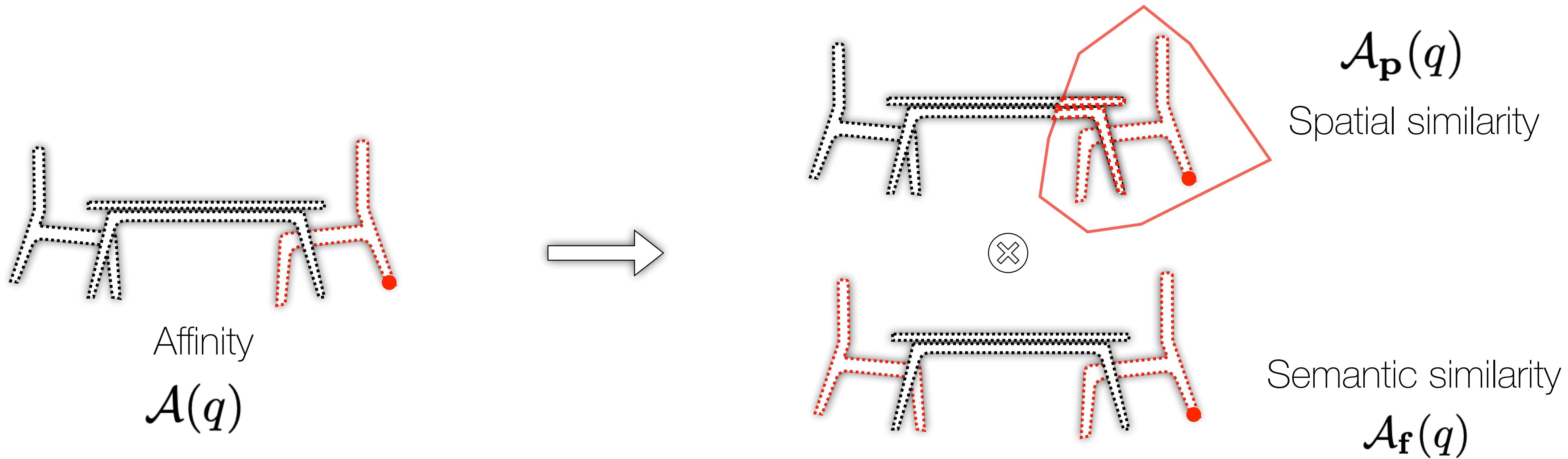
To 3D point cloud



mAP < 50%

# Affinity as proposal

- Input a given query ( $q$ ) on the input point cloud.
- Output affinity score map of range  $[0, 1]$ .
- Consists of *spatial* and *semantic* affinity.



# Semantic similarity

- A simple **exponential kernel**
  - L2 distance between semantic similarity features

$$\mathcal{A}_f(q)[n] = \exp(-\tau_f \cdot \mathcal{K}_f(q, n))$$

The diagram shows the formula  $\mathcal{A}_f(q)[n] = \exp(-\tau_f \cdot \mathcal{K}_f(q, n))$  with annotations pointing to its parts:

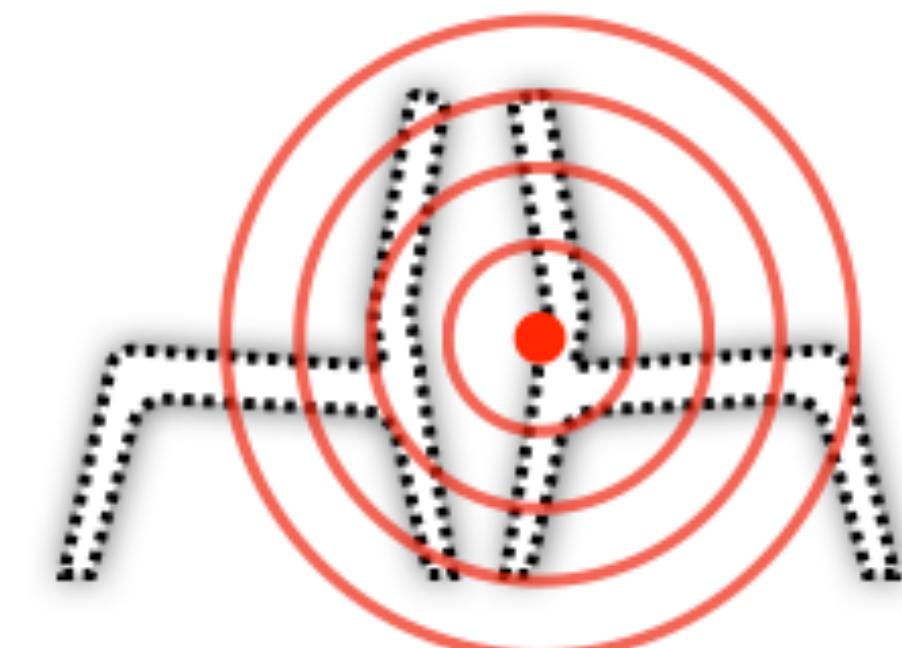
- Query**: Points to the variable  $q$  in  $\mathcal{A}_f(q)[n]$ .
- Indexing n-th element**: Points to the index  $[n]$  in  $\mathcal{A}_f(q)[n]$ .
- Backbone feature**: Points to the term  $\mathcal{K}_f(q, n)$ .
- Temperature**: Points to the parameter  $\tau_f$ .
- L2 distance kernel**: Points to the entire formula  $\mathcal{A}_f(q)[n] = \exp(-\tau_f \cdot \mathcal{K}_f(q, n))$ .

# Spatial similarity

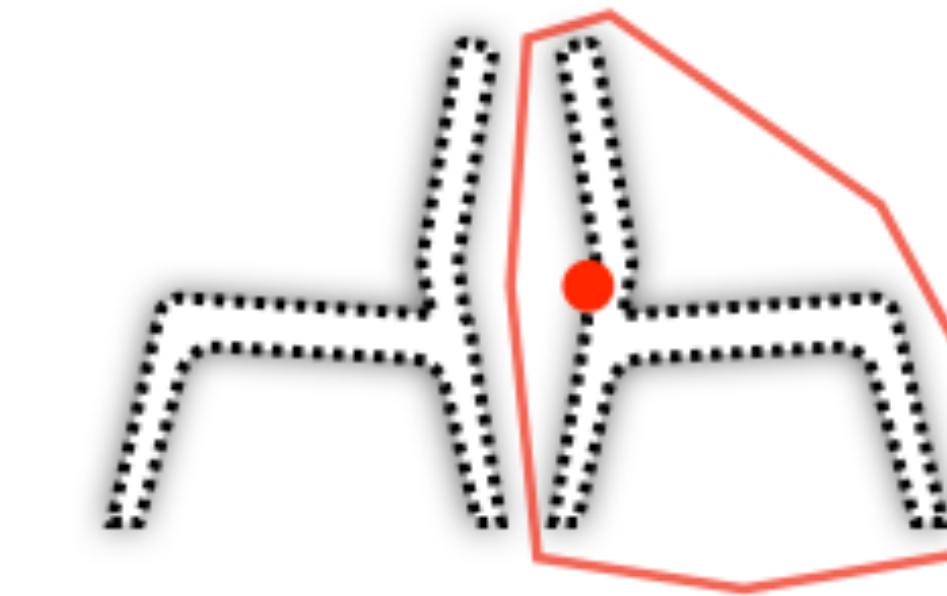
- A non-isotropic kernel based on local convex hulls
  - Is able to isolate the queried instance

$$\mathcal{A}_{\mathbf{p}}(q)[n] = \exp(-\boxed{\tau_{\mathbf{p}}} \cdot \mathcal{K}_{\mathbf{p}}(q, n))$$

Temperature



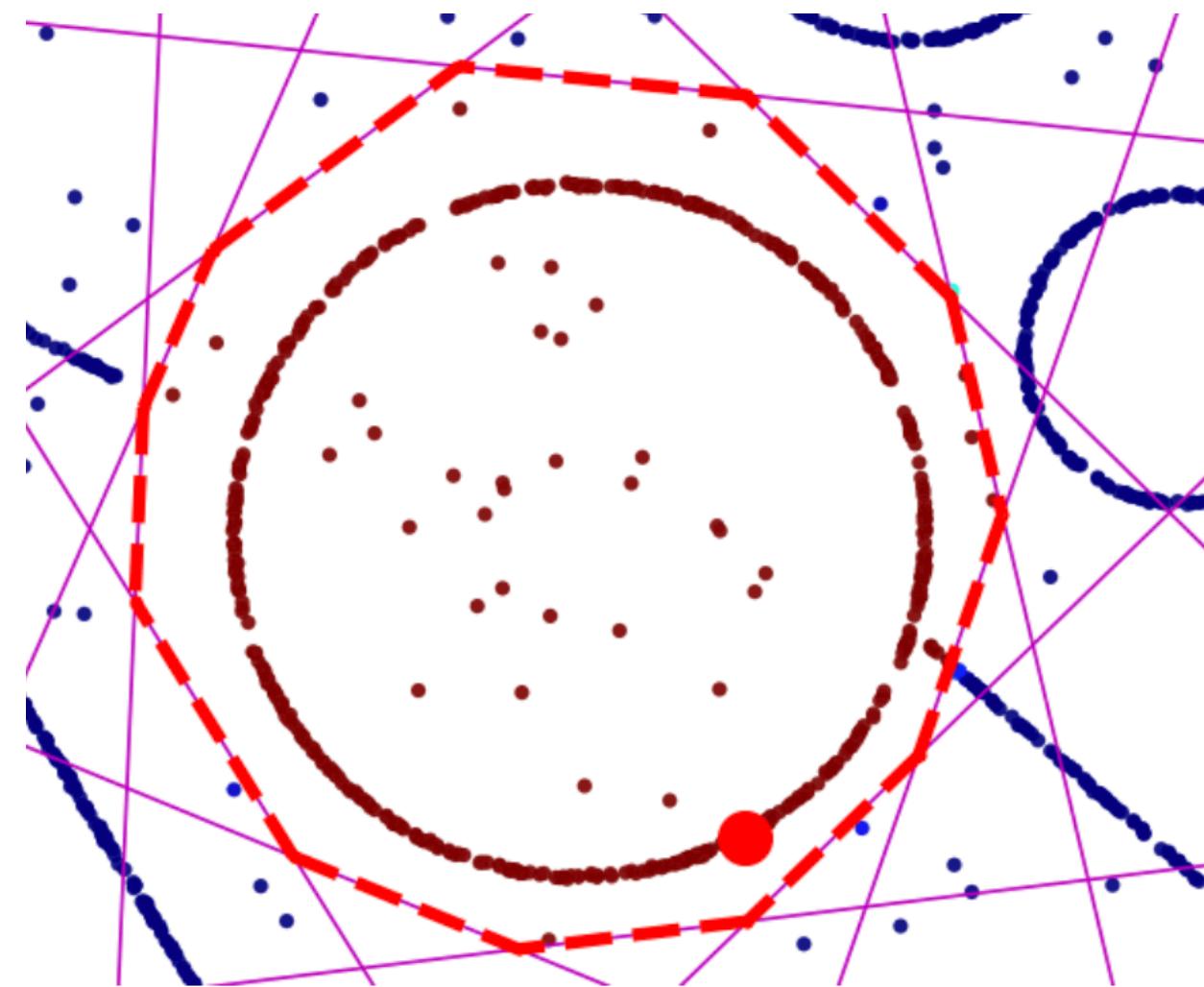
isotropic affinity



non-isotropic affinity

# Spatial similarity — in more details

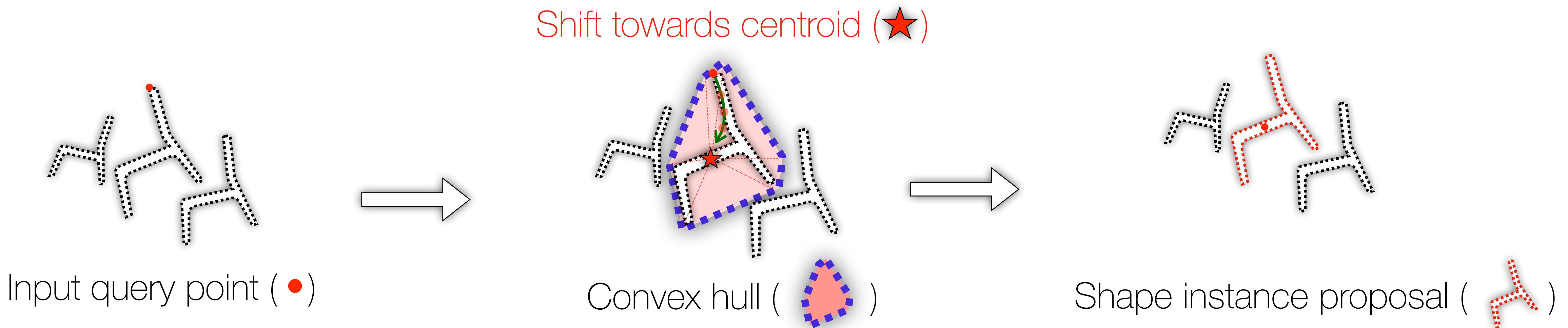
- A hyper coordinate network -- convex hull
  - Consists of a set of hyperplanes conditional on query
  - Is an explainable neural field—comparing with recent MLP-based neural field.
  - Is very small and memory efficient



$$C(x; f) \begin{cases} = 0 & \text{if } x \text{ inside convex defined by } f, \\ > 0 & \text{otherwise (}\approx \text{boundary distance).} \end{cases}$$

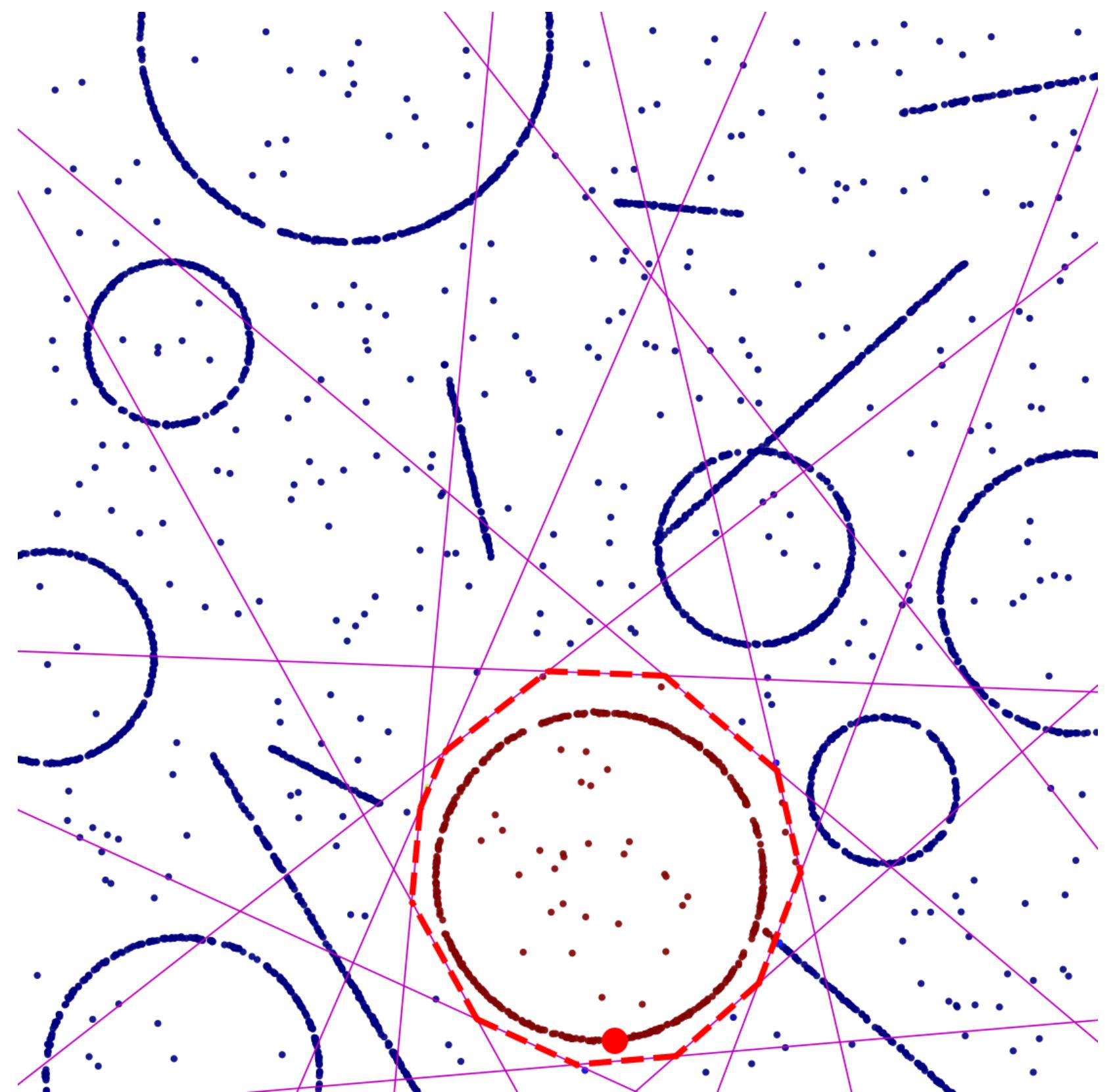
# Bilateral filtering

- Iterative inference
  - Update query as the weighted center of input
  - Shift query into the centroid of instance



# Results – Learned convex hull of high quality

Convex hulls in 2D



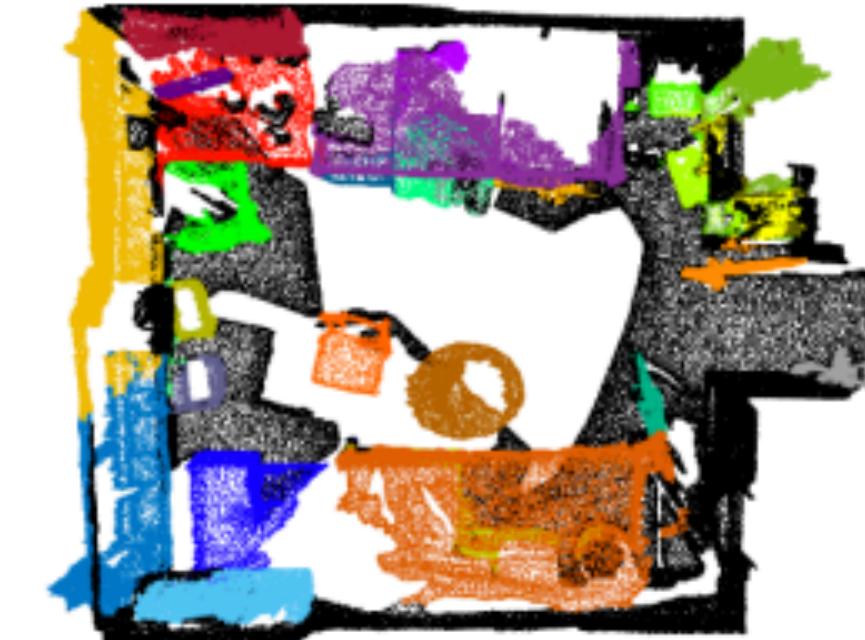
Convex hulls in 3D



# Results – Instance segmentation in ScanNet.

- Our method is the **best among purely top-down** methods
- Top-down methods in general lags behind the bottom-up/Mix strategy

Qualitative



Quantitative

	Methods	Validation			Test		
		mAP	AP <sub>50</sub>	AP <sub>25</sub>	mAP	AP <sub>50</sub>	AP <sub>25</sub>
Bottom-up	PointGroup [22]	34.8	56.7	71.3	40.7	63.6	77.8
	SSTNet [26]	49.4	64.3	74.0	50.6	69.8	78.9
	HAIS [2]	43.5	64.4	75.6	45.7	69.9	80.3
Mix	Dyco3D [18]	35.4	57.6	72.9	39.5	64.1	76.1
	SoftGroup [44]	–	67.6	78.9	50.4	76.1	86.5
Top-down	3D-SIS [20]	–	18.7	35.7	16.1	38.2	55.8
	GSPN [53]	19.3	37.8	53.4	–	30.6	–
	3D-Bonet [52]	–	–	–	25.3	48.8	68.7
	<b>Ours</b>	<b>36.0</b>	<b>55.5</b>	<b>71.1</b>	<b>35.3</b>	<b>55.5</b>	<b>71.8</b>

# Summary



Explainable models have lots of potential. When carefully designed, these models can also have good performance!

[wsunid.github.io](https://wsunid.github.io)