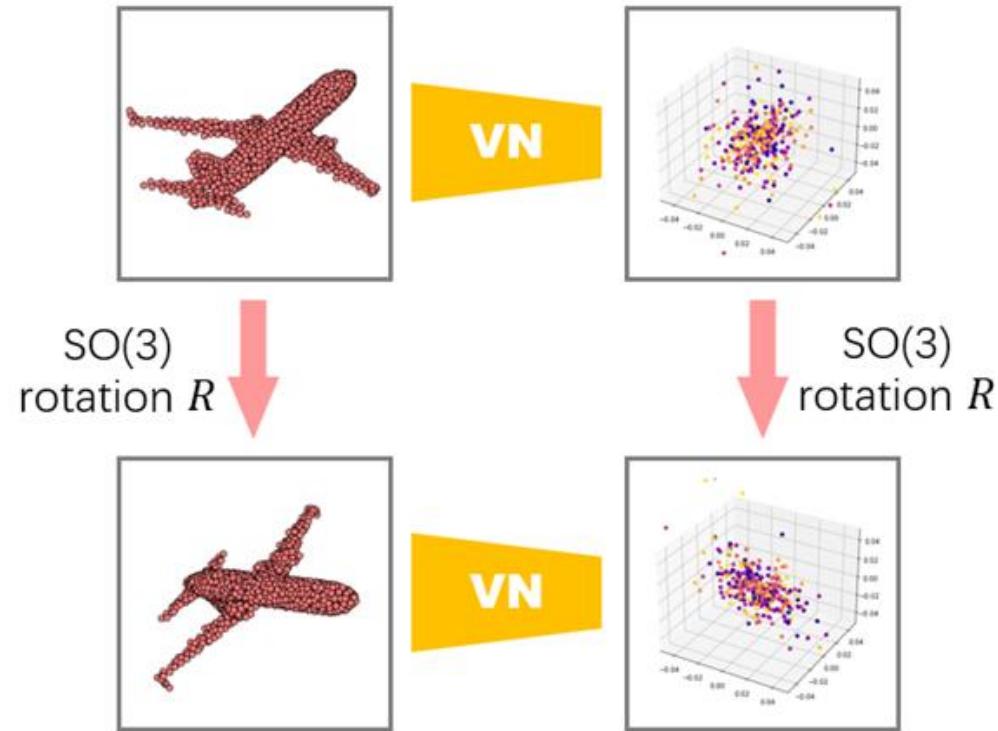


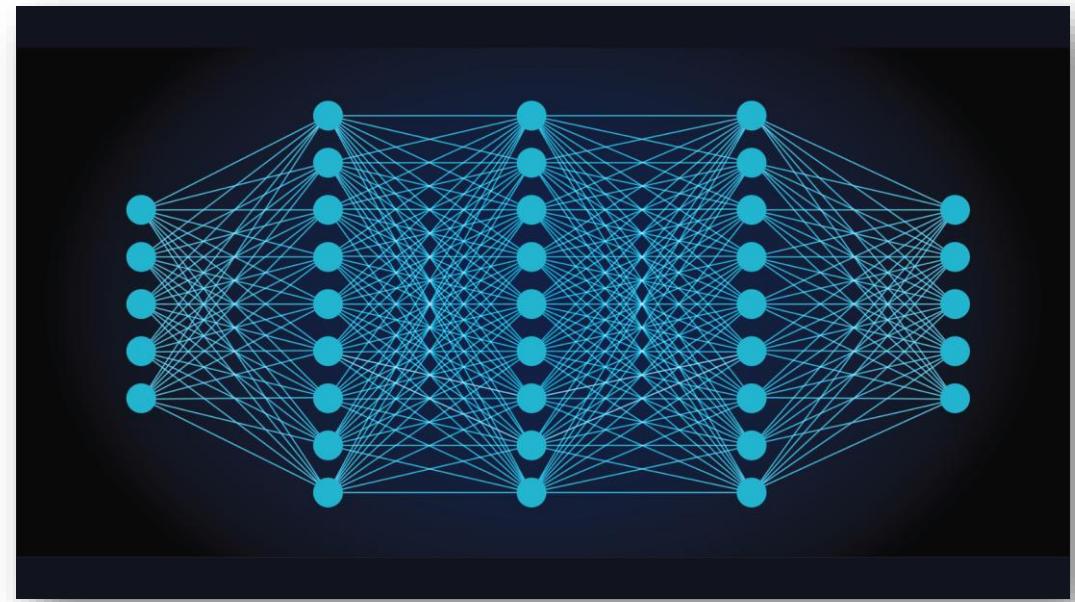
Injecting 3D Priors Into Neural Nets



Leonidas Guibas
Stanford University



Things Neural Nets Should Know



Things Neural Nets Should Know

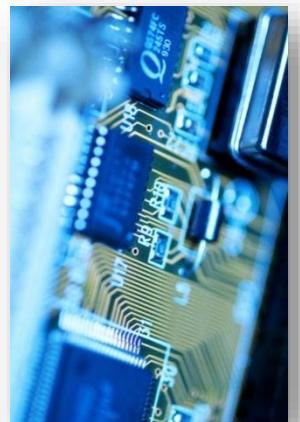
- Representation shaping: Distill knowledge into the neural representation

[Point Cloud 3D Rotation Equivariance](#)



- Network shaping: Distill knowledge into the neural architectures that use the representation

[Mutual Information Grouping in NeRFs and GSs](#)

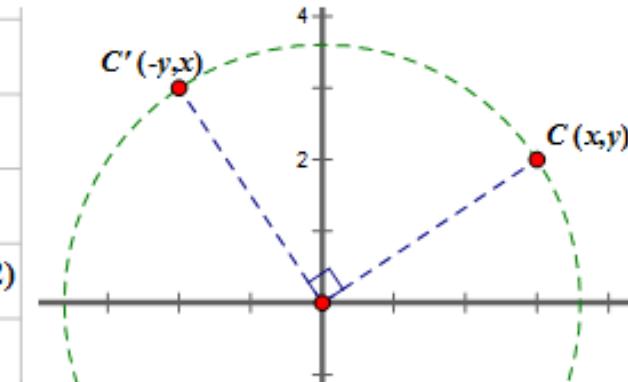
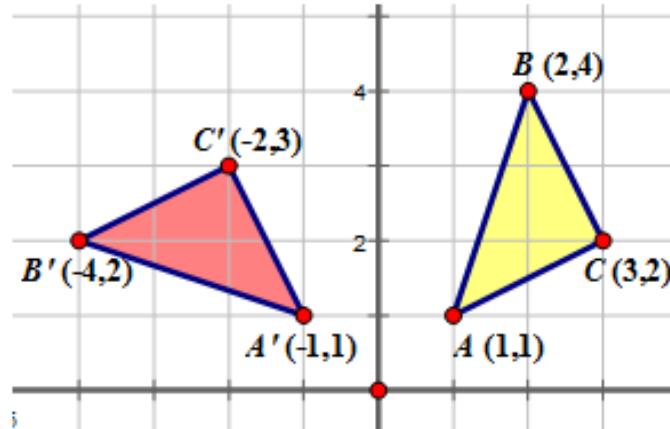


Enforcing Classic Invariance and Equivariance for 3D Objects and Scenes

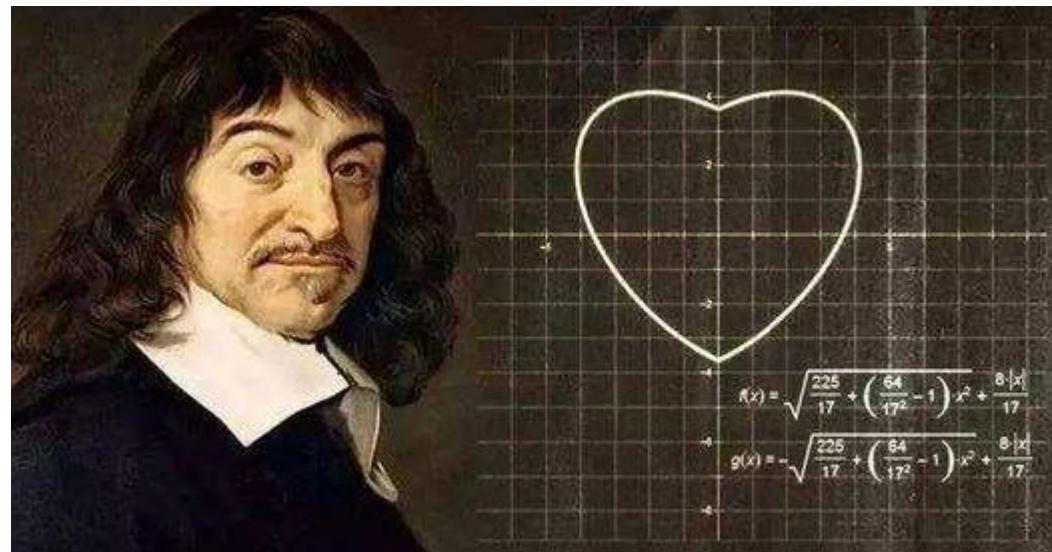
Geometry with Coordinates

Rotation of 90° about the Origin

When we rotate 90° about the origin, we see that the x and y coordinates are reversed and the new x coordinate is negated.

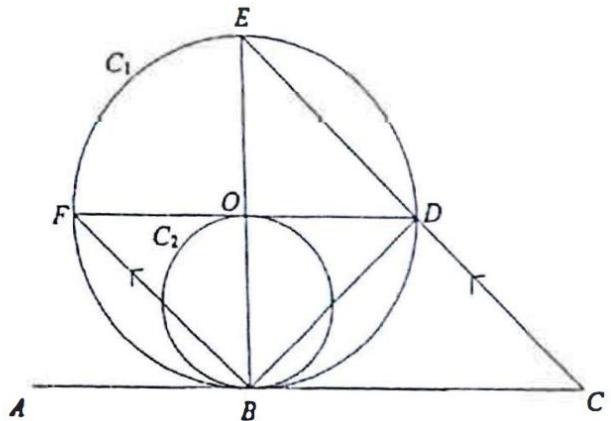


RULE FOR ROTATION BY 90° ABOUT THE ORIGIN $R_{o,90^\circ}(x, y) = (-y, x)$



René Descartes

Geometry without Coordinates



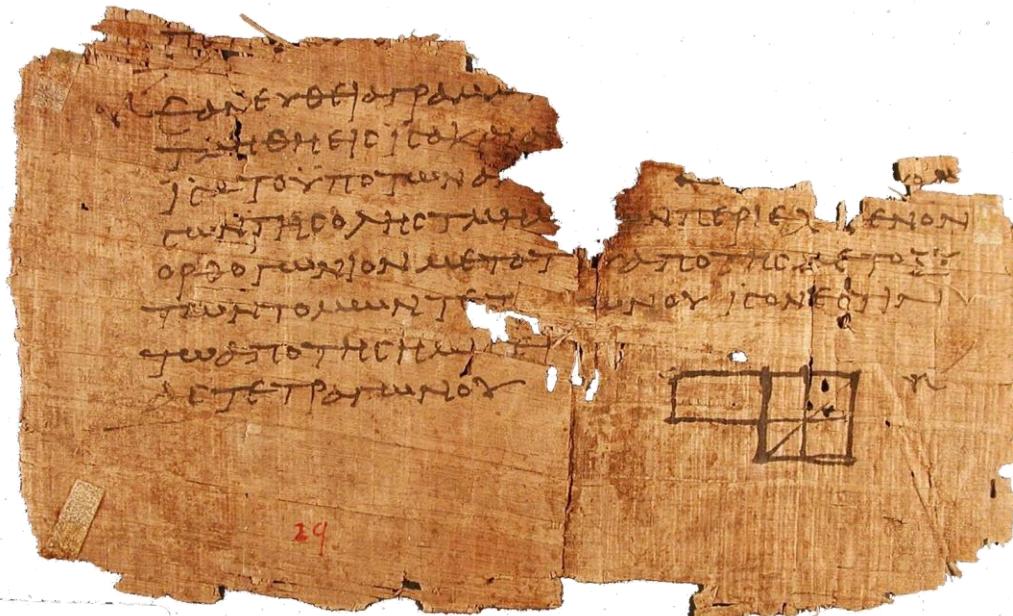
In the above diagram, C_1 is a circle with centre O , and C_2 is another circle passing through the point O and touching the circumference of C_1 at the point B . AC is the tangent to C_1 at B and meets the line ED produced at the point C . DF is the tangent to C_2 at O . BOE is a straight line.

Given that BF is parallel to CE , show that

(i) $\angle BDC$ is a right angle,

(ii) $\frac{OE}{BC} = \frac{1}{2}$, and

(iii) $CD \times CE = BC^2$



Euclid

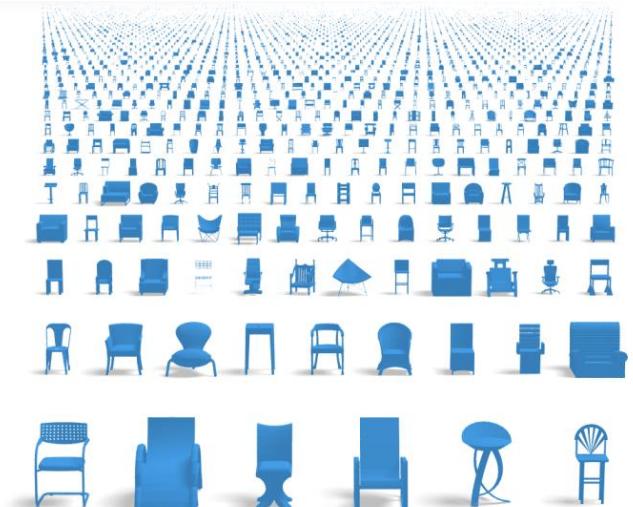
Factors of 3D Variation

- Geometric data is almost always given to us in a particular coordinate system.
- In many settings this frame or pose aspect of the presentation needs to be disentangled from the intrinsic data geometry, as it may be a nuisance factor.
- Depending on the application, this variation is captured by a transformation group that may include translation, rotation, scale, etc.

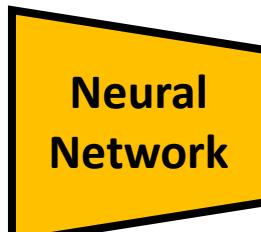


Learning on Geometric Shapes

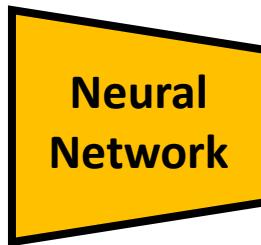
- Typical neural networks are trained on co-aligned collections of shapes



- Such networks do not generalize to objects in arbitrary poses



“chair”

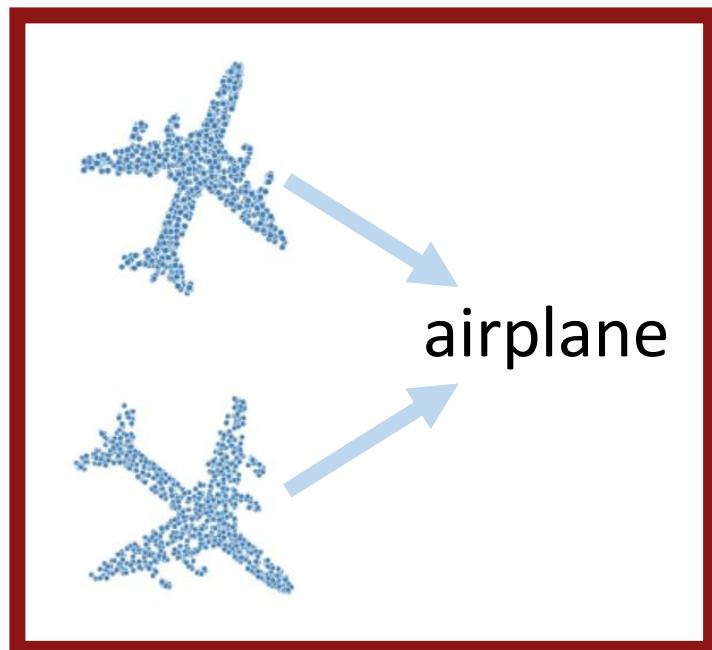


“???”



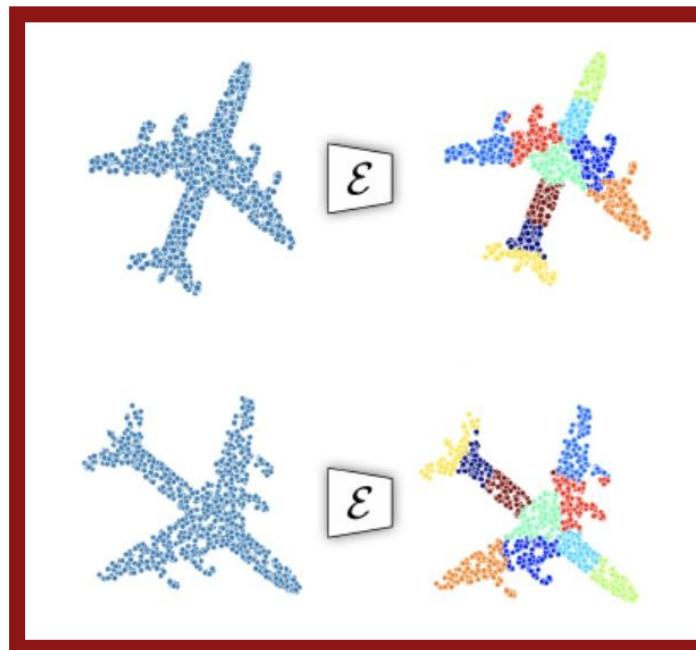
Key Requirement: Invariance and Equivariance

classification



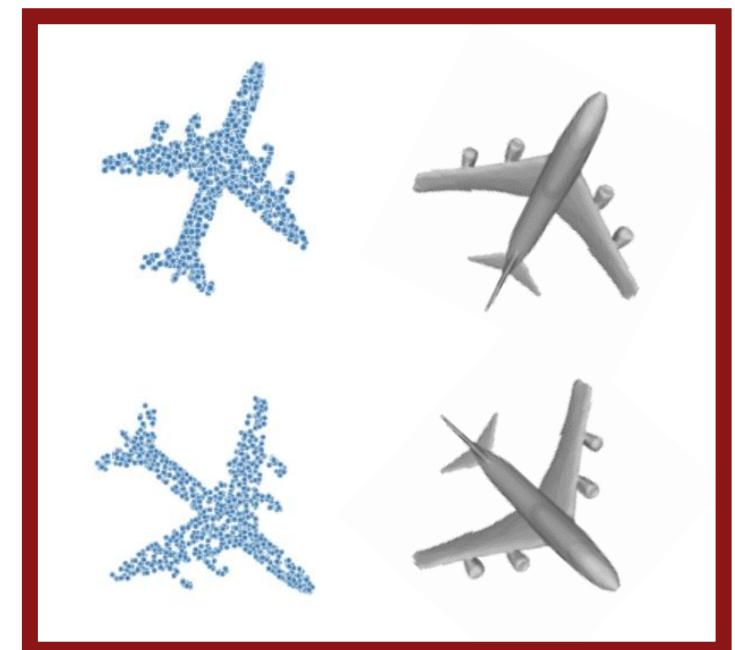
invariant

segmentation



equivariant
(or invariant, depending
on data representation)

reconstruction



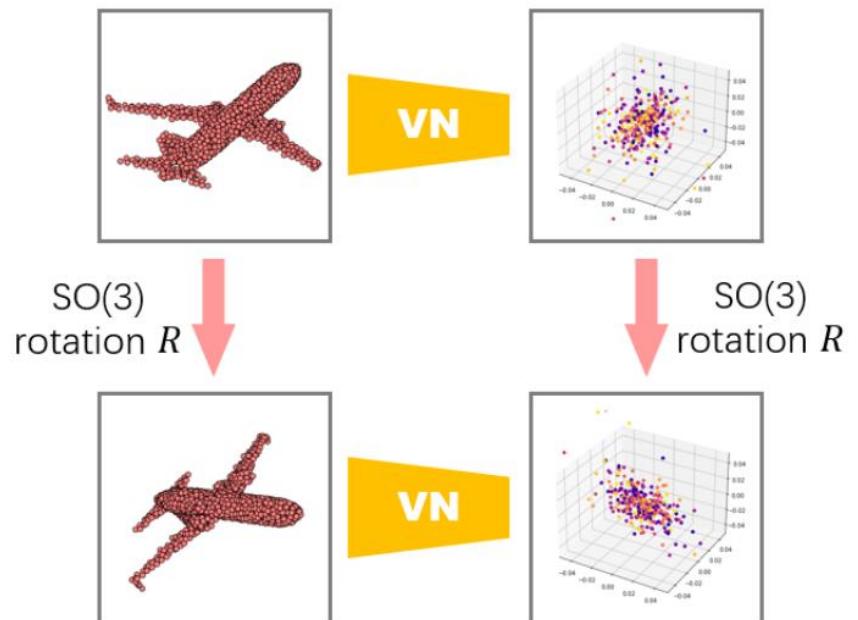
equivariant encoder
invariant decoder

In/Equivariance

We say a neural network $f(\cdot; \theta)$ is rotation equivariant, if for any 3D rotation $R \in \text{SO}(3)$ 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 identity mapping, it is invariance



A Naïve Solution: Data Augmentation

Apply random rotations 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



Vector Neurons for $SO(3)$ Equivariance

Congyue Deng, Or Litany, Yueqi Duan, Adrien Poulenard, Andrea Tagliasacchi, Leonidas Guibas,
ICCV'21

Scalar vs Vector Neurons

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

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

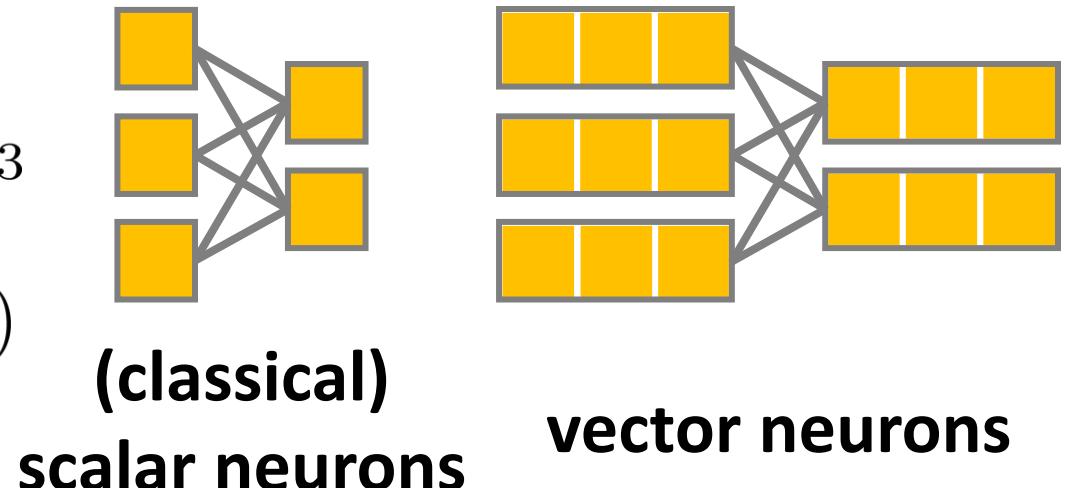
- For pointcloud with N points $\mathcal{V} = \{\mathbf{V}_1, \mathbf{V}_2, \dots, \mathbf{V}_N\} \in \mathbb{R}^{N \times C \times 3}$

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 $R \in \text{SO}(3)$

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



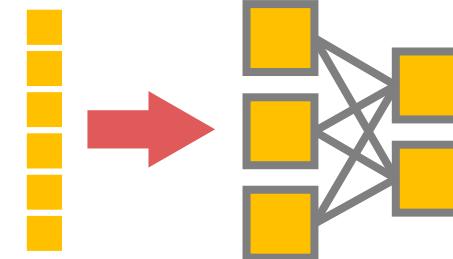
Expressing Transformations in the Latent Space

A network whose latent space understands rigid transformations

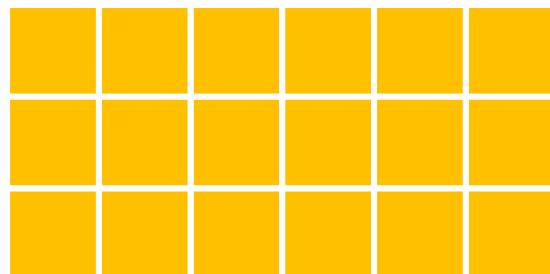
Classical Neuron: scalar channels



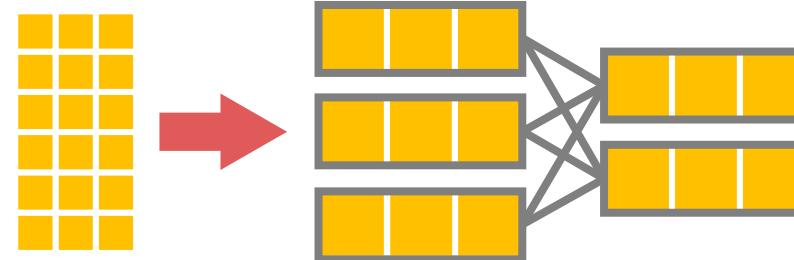
$C \times 1$ feature



Vector Neuron: 3D vector channels

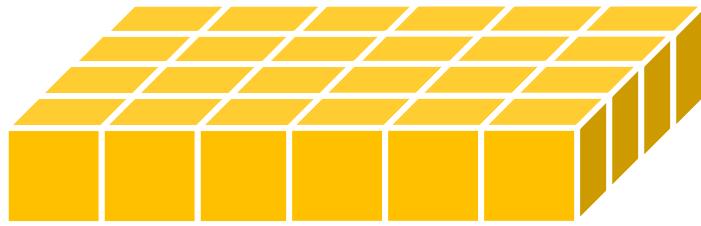


$C \times 3$ feature



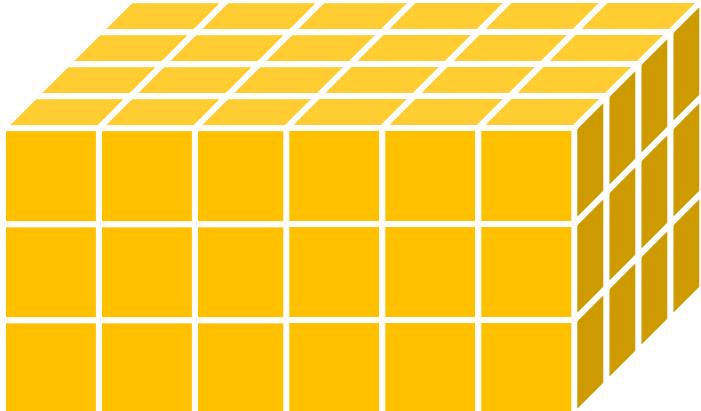
Vector Neuron Features for Point Cloud

Classical:

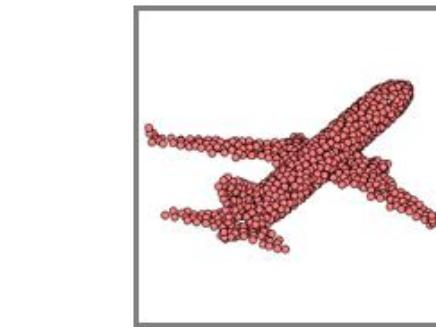


$N \times C \times 1$ feature

VN:



$N \times C \times 3$ feature

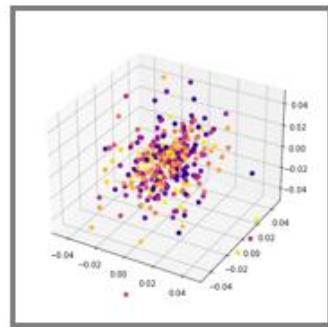


VN

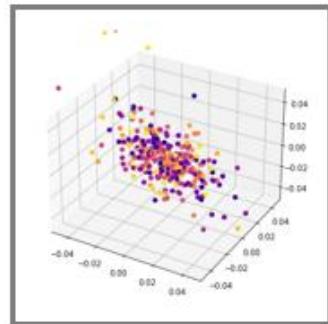
$\text{SO}(3)$
rotation R



VN



$\text{SO}(3)$
rotation R



Vector Neuron Linear Operations

Linear operator: left multiply by the learnable weight matrix

$$\begin{array}{ccc} \begin{matrix} \text{Red grid} \\ C' \times C \text{ weight} \end{matrix} & \times & \begin{matrix} \text{Yellow grid} \\ C \times 3 \text{ feature} \end{matrix} \\ & & = \\ & & \begin{matrix} \text{Gold grid} \\ C' \times 3 \text{ feature} \end{matrix} \end{array}$$

Equivariance: right multiply by the $\text{SO}(3)$ rotation matrix

$$\begin{array}{ccc} \left[\begin{matrix} \text{Red grid} \\ C' \times C \text{ weight} \end{matrix} \right] & \times & \left[\begin{matrix} \text{Yellow grid} \\ C \times 3 \text{ feature} \end{matrix} \right] \times \left[\begin{matrix} \text{Teal grid} \\ C' \times 3 \text{ feature} \end{matrix} \right] \\ & & = \\ & & \begin{matrix} \text{Gold grid} \\ C' \times 3 \text{ feature} \end{matrix} \end{array}$$

VN Linear Layer

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

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

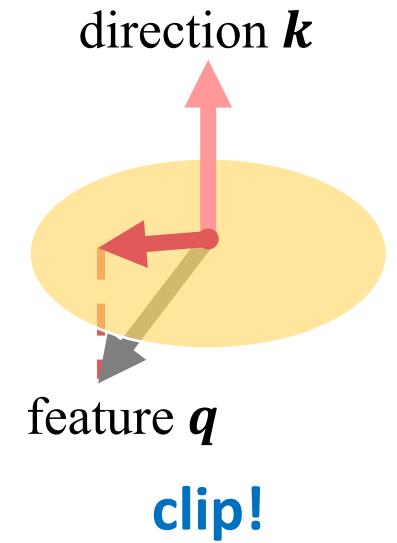
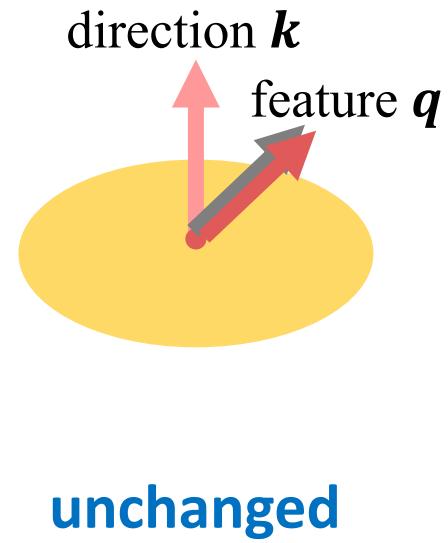
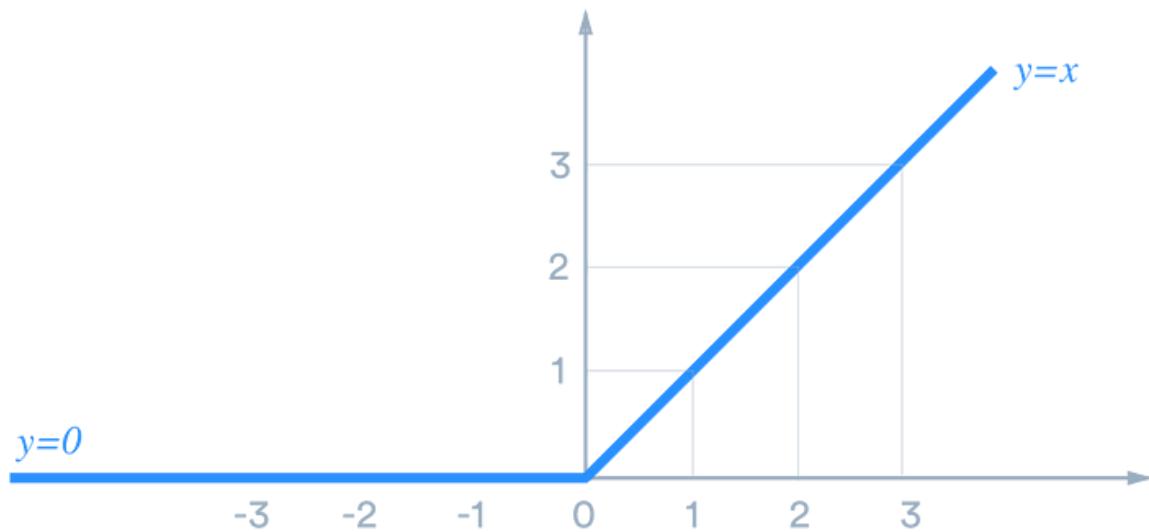
$$\mathbf{V}' = f_{\text{lin}}(\mathbf{V}; \mathbf{W}) = \mathbf{W}\mathbf{V} \in \mathbb{R}^{C' \times 3}$$

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

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

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

Vector Neuron ReLU Non-Linearity



Rectified Linear Unit: VN Non-Linearity

ReLU Non-Linearity

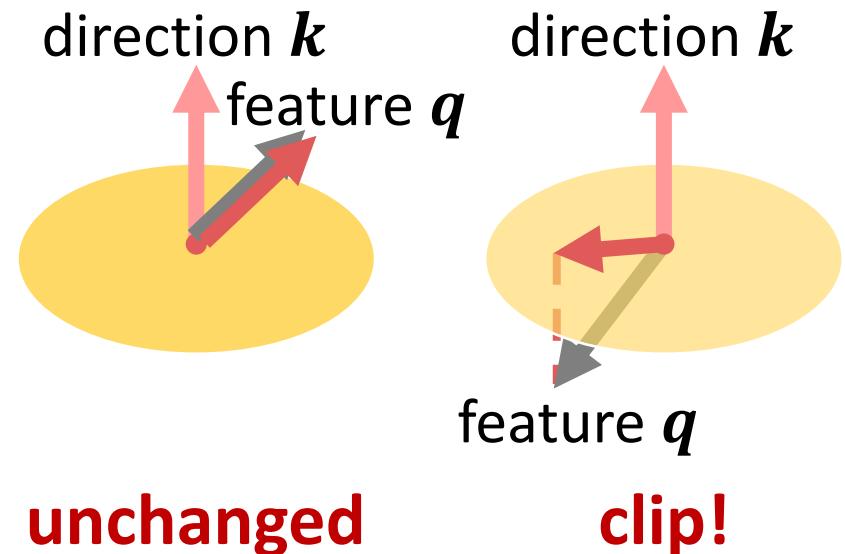
Weights $\mathbf{W} \in \mathbb{R}^{1 \times C}$ and $\mathbf{U} \in \mathbb{R}^{1 \times C}$

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

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

For each output vector neuron $v' \in V'$

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

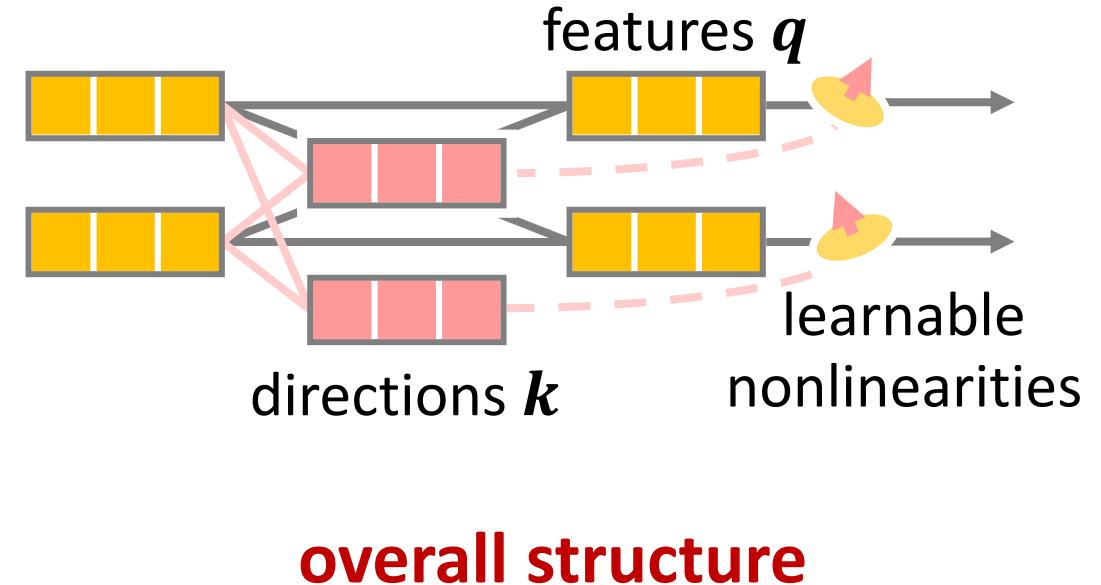


[Also works for uniform scaling]

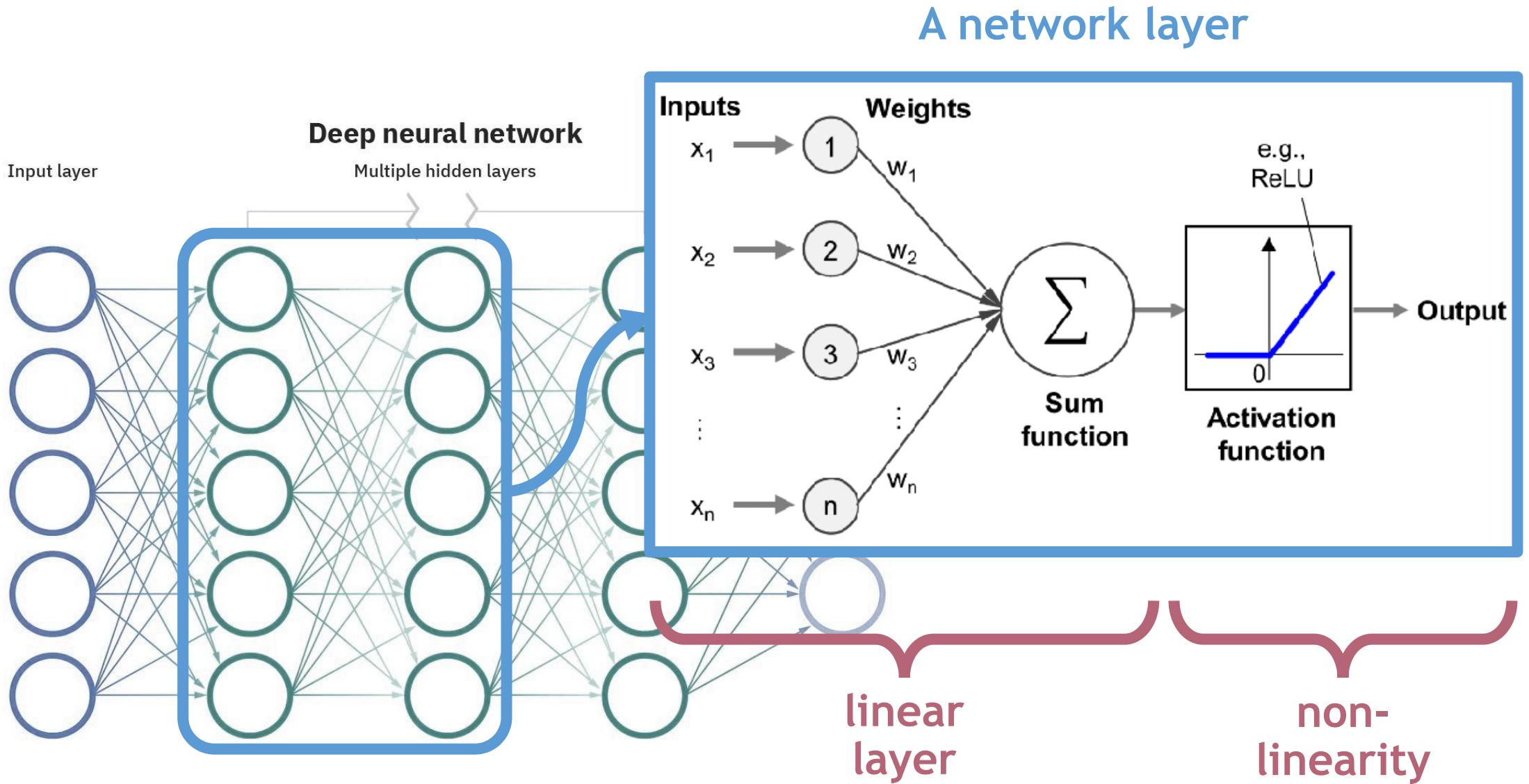
VN Non-Linearity: A High-d ReLU

Learnable ReLU Non-Linearity

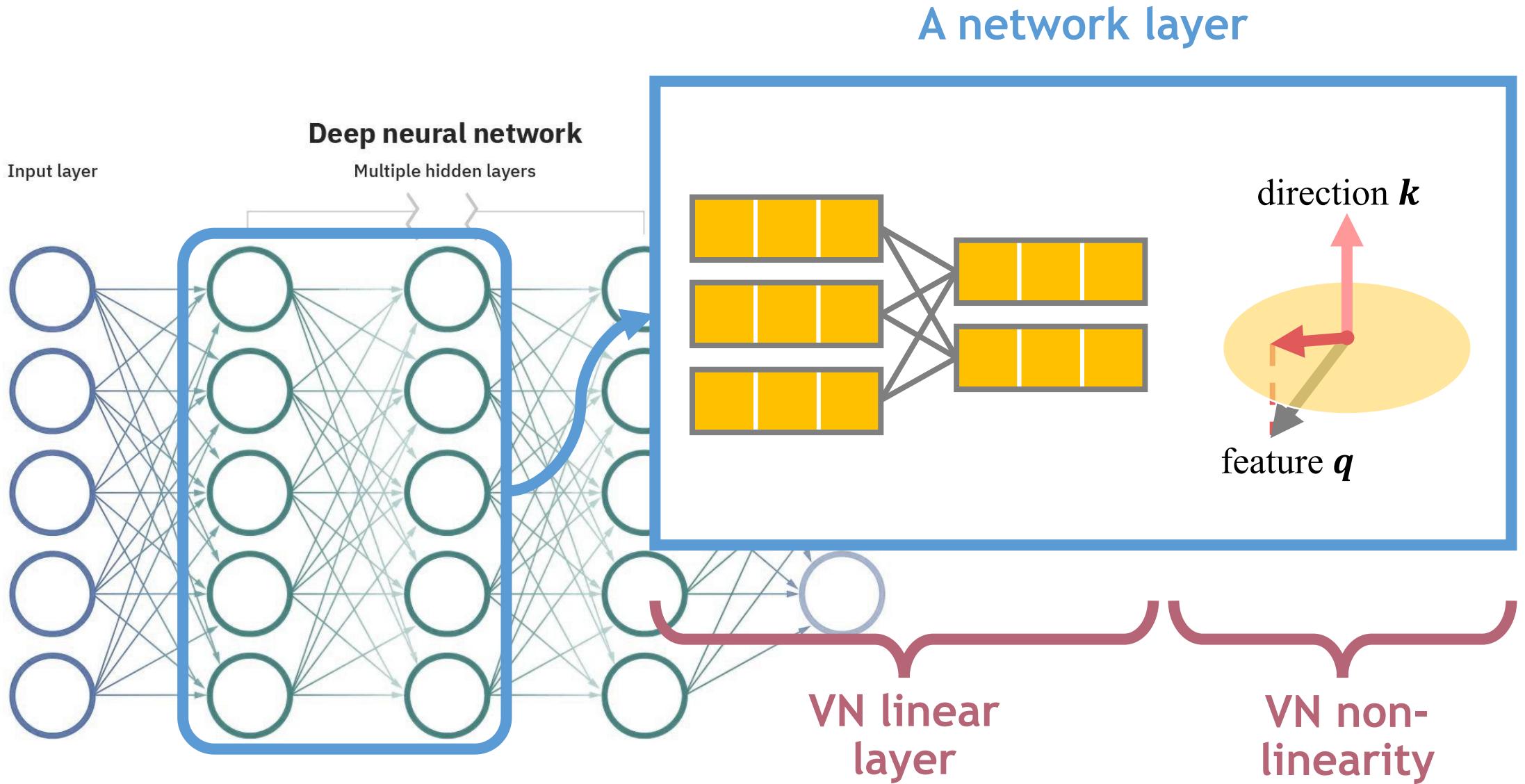
- **Non-linear layer** (with built-in linear layer)
= input linear transformation q + non-linearity k
- Other non-linearities



Network Layer: Scalar Network

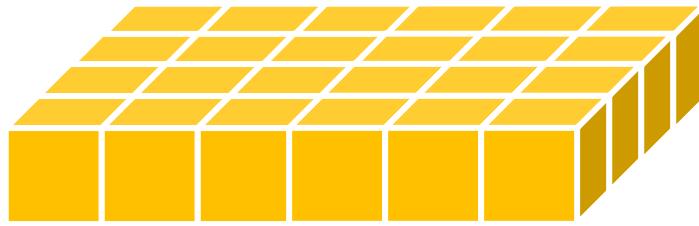


Network Layer: Vector Network



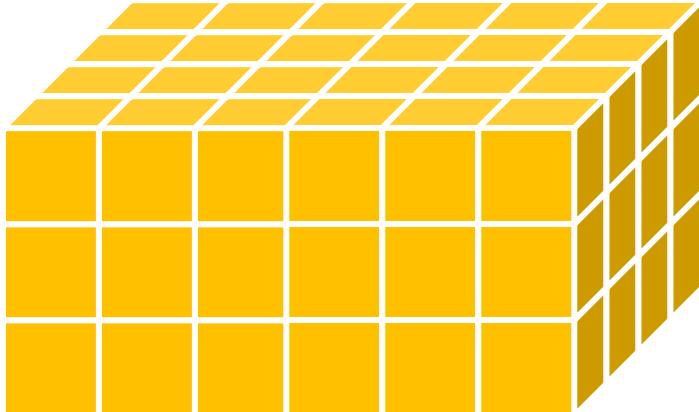
Vector Neuron Features for Point Cloud

Classical:

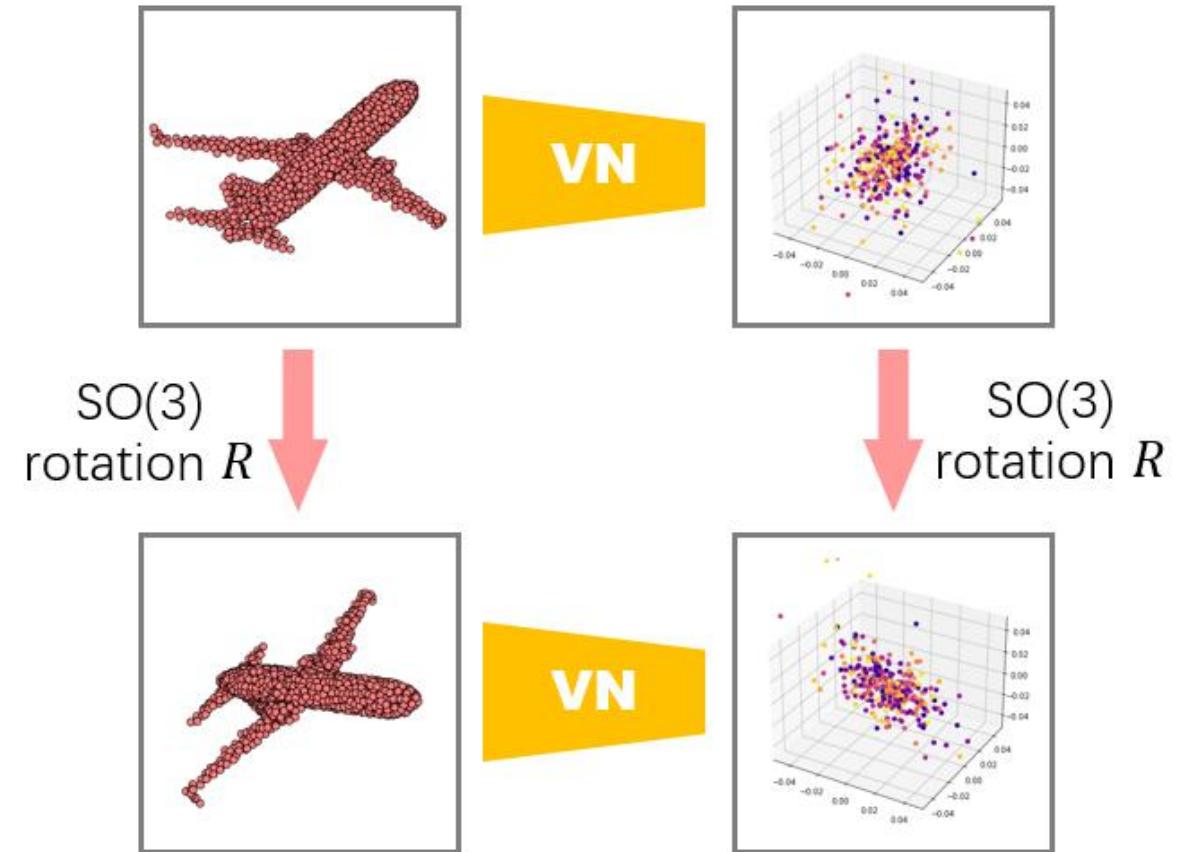


$N \times C \times 1$ feature

VN:



$N \times C \times 3$ feature

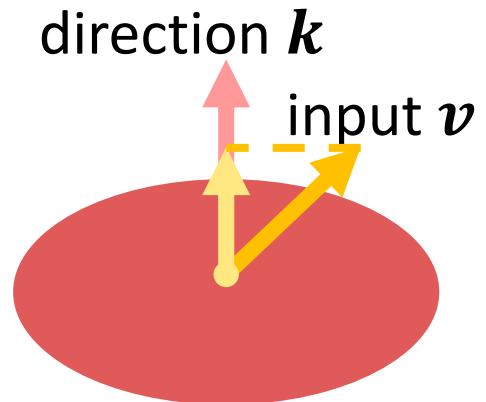
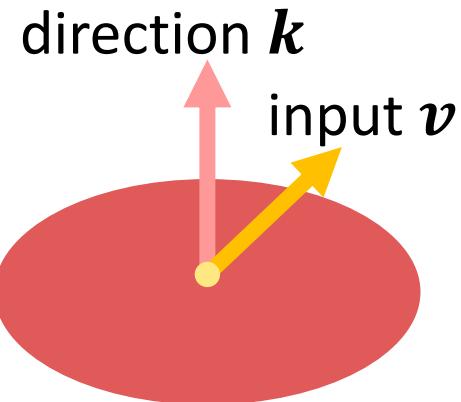


Vector Neuron Pooling

✓ Mean pooling

? Max pooling

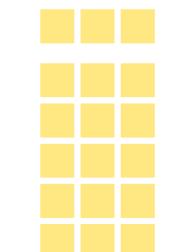
- (Similar to non-linearity)
- argmax alone learned directions



$C \times 3$
input

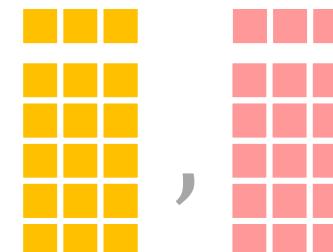


$C \times 3$
directions



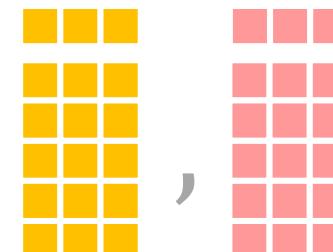
$C \times 3$
output

= argmax(



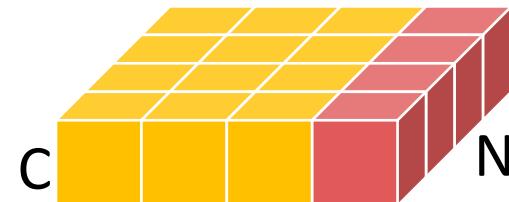
$C \times 3$
input

,

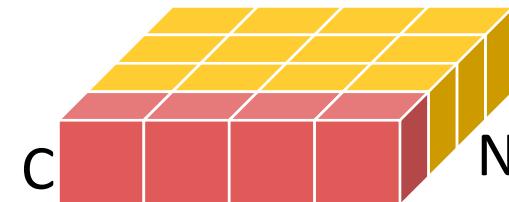


Vector Neuron Normalizations

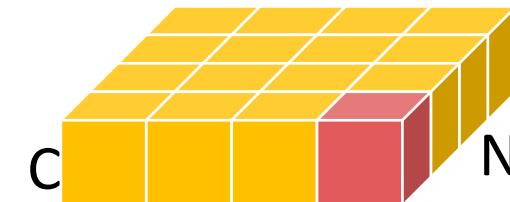
- ✓ LayerNorm
- ✓ InstanceNorm
- ✓ Dropout



LayerNorm



BatchNorm

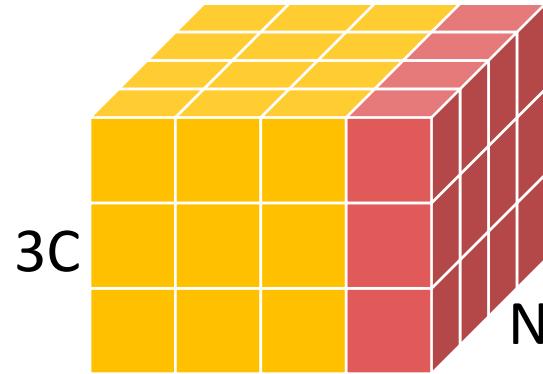


InstanceNorm

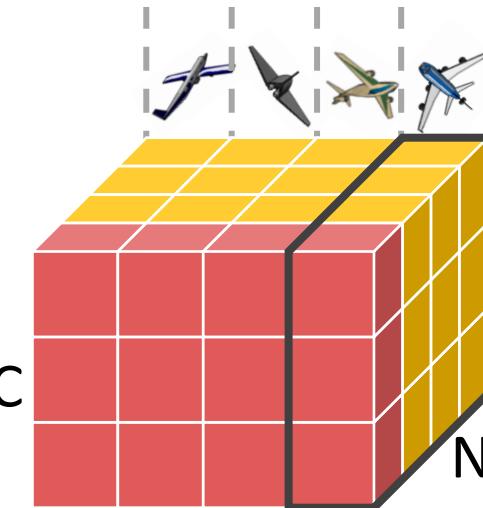
(classical) scalar neurons

? BatchNorm

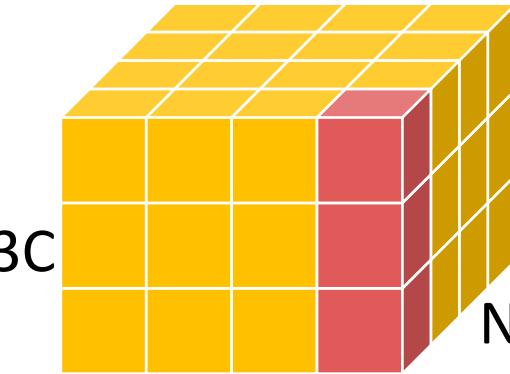
averaging across arbitrarily rotated inputs would not necessarily be meaningful



LayerNorm



BatchNorm



InstanceNorm

Vector neurons

Vector Neuron Normalizations

BatchNorm

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

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

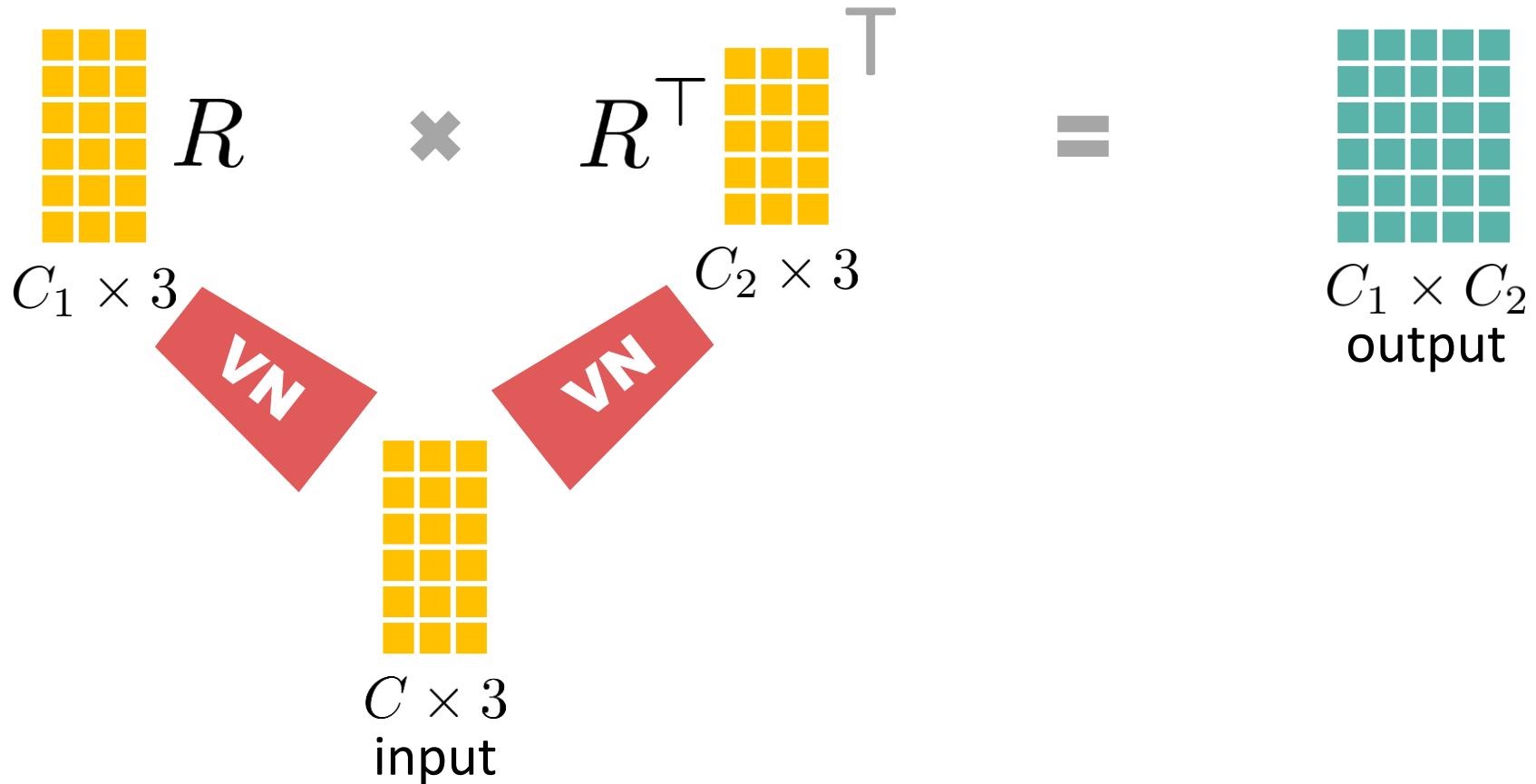
$$\{\mathbf{N}'_b\}_{b=1}^B = \text{BatchNorm}(\{\mathbf{N}_b\}_{b=1}^B)$$

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

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

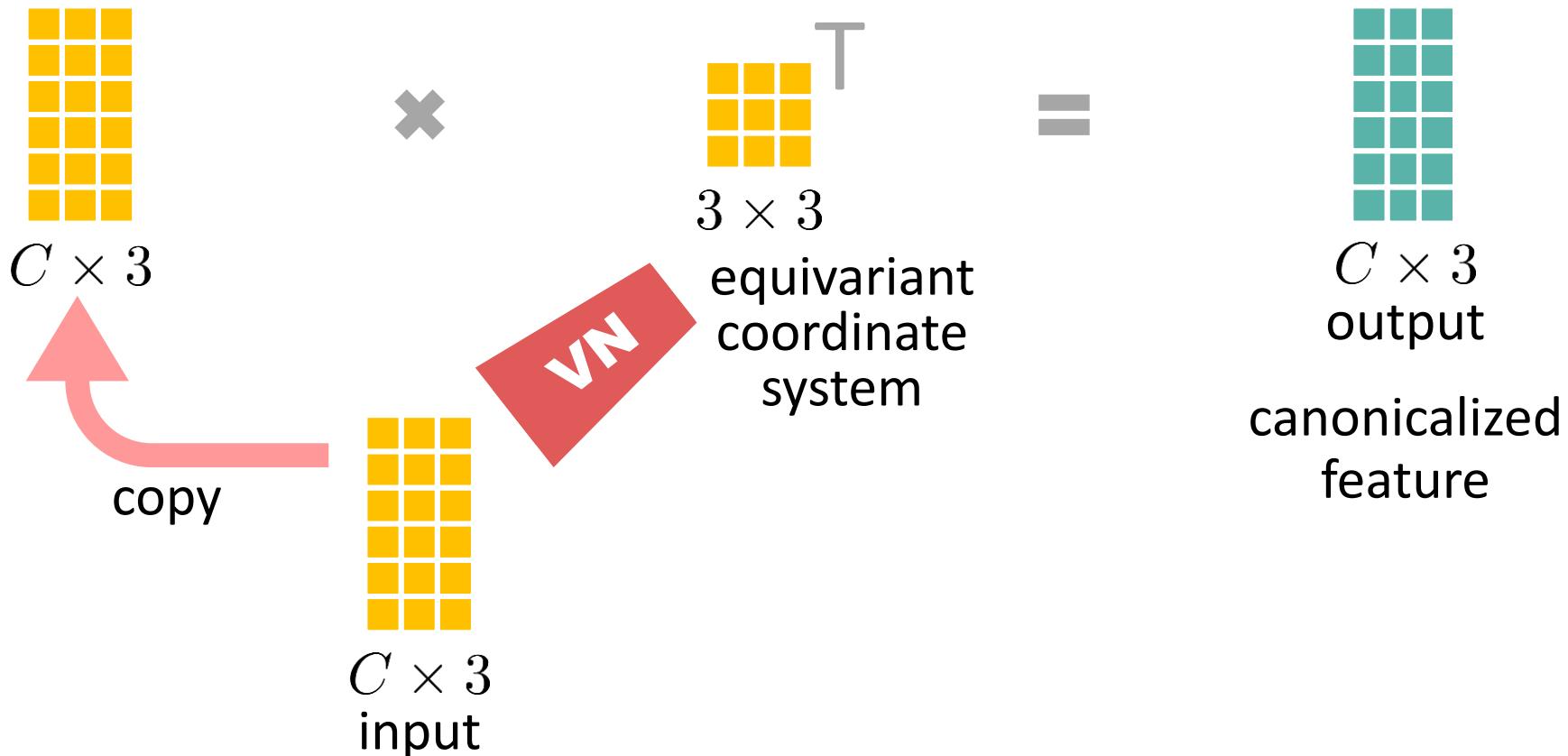
From Equivariance to Invariance

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



Vector Neuron Invariant Layer

Specifically...



Vector Neuron Invariant Layer

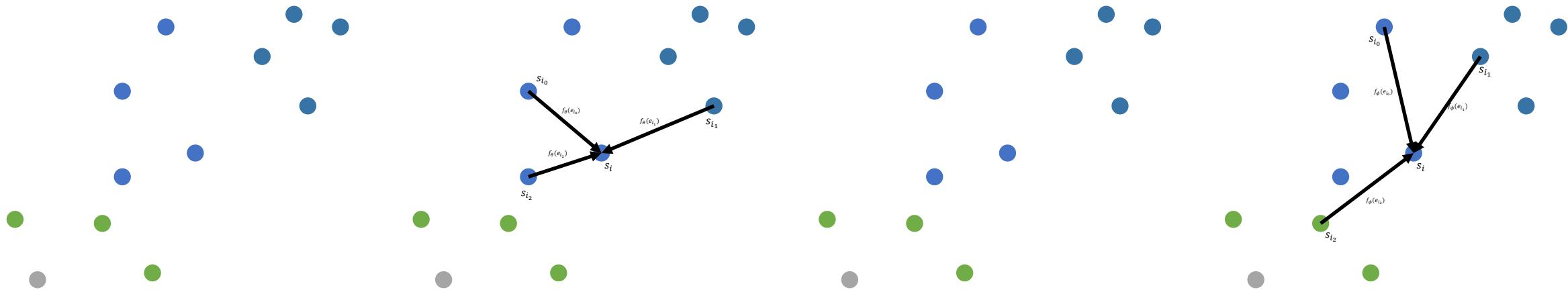
- Product of an equivariant signal $\mathbf{V} \in \mathbb{R}^{C \times 3}$ by the transpose of another equivariant signal $\mathbf{T} \in \mathbb{R}^{C' \times 3} \rightarrow$ invariant signal
- **Special case:** $\mathbf{T} \in \mathbb{R}^{3 \times 3}$ - an equivariant coordinate system
- For pointcloud,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

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

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

Vectorize Classical 3D Networks: DGCNN, PointNet

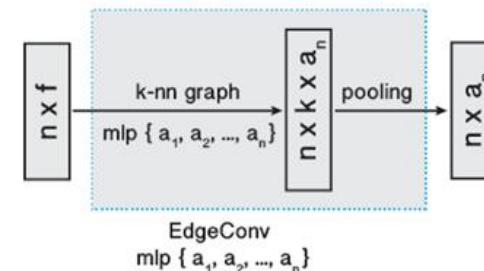
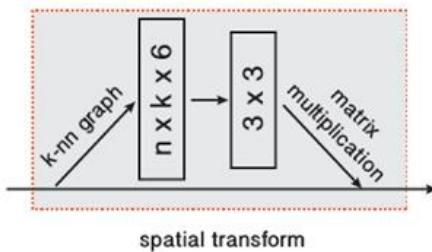
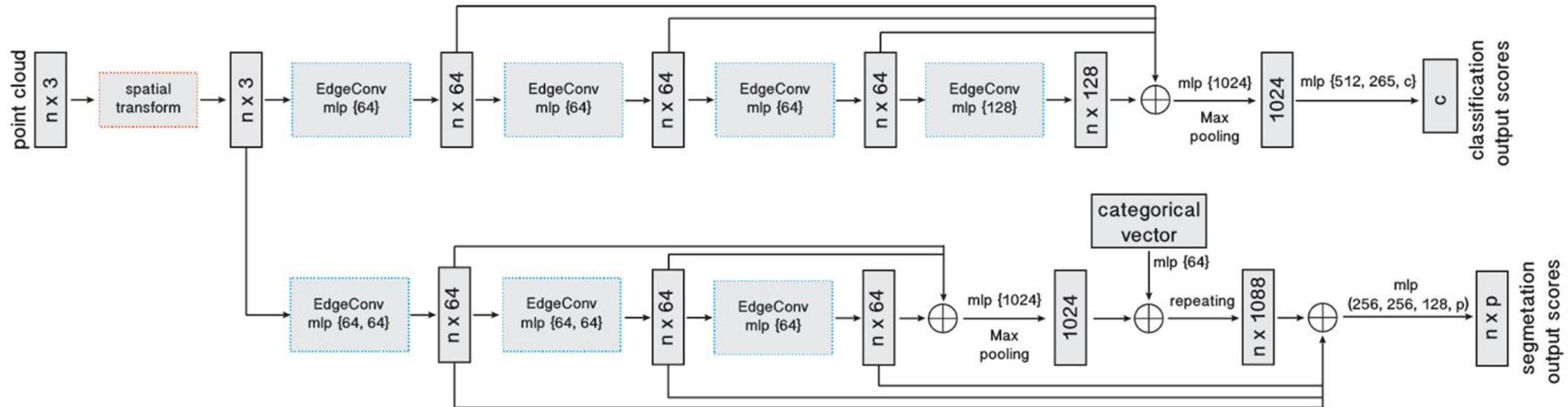
Dynamic Graph CNN (DGCNN)



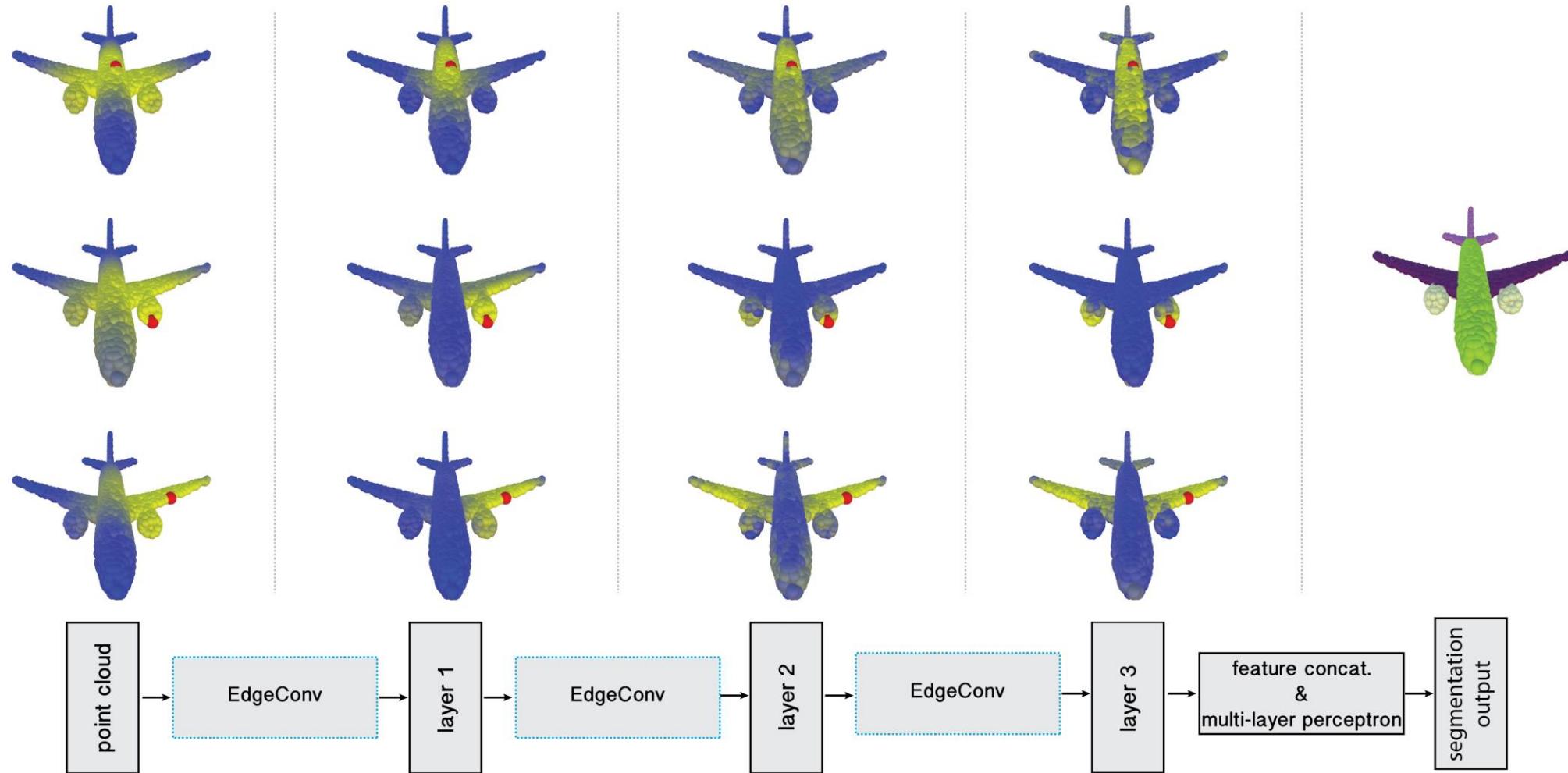
DGCNN alternates feature learning (EdgeConvs) and graph NN reconstruction

DGCNN Architecture: Alternating Processing

EdgeConv: Edge Convolutions



Dynamic Graph CNN (DGCNN)



Build VN Networks: VN-DGCNN

DGCNN

Edge feature:  $'_{nm} = \text{ReLU}(\Theta(\text{---}_m - \text{---}_n) + \Phi\text{---}_n)$

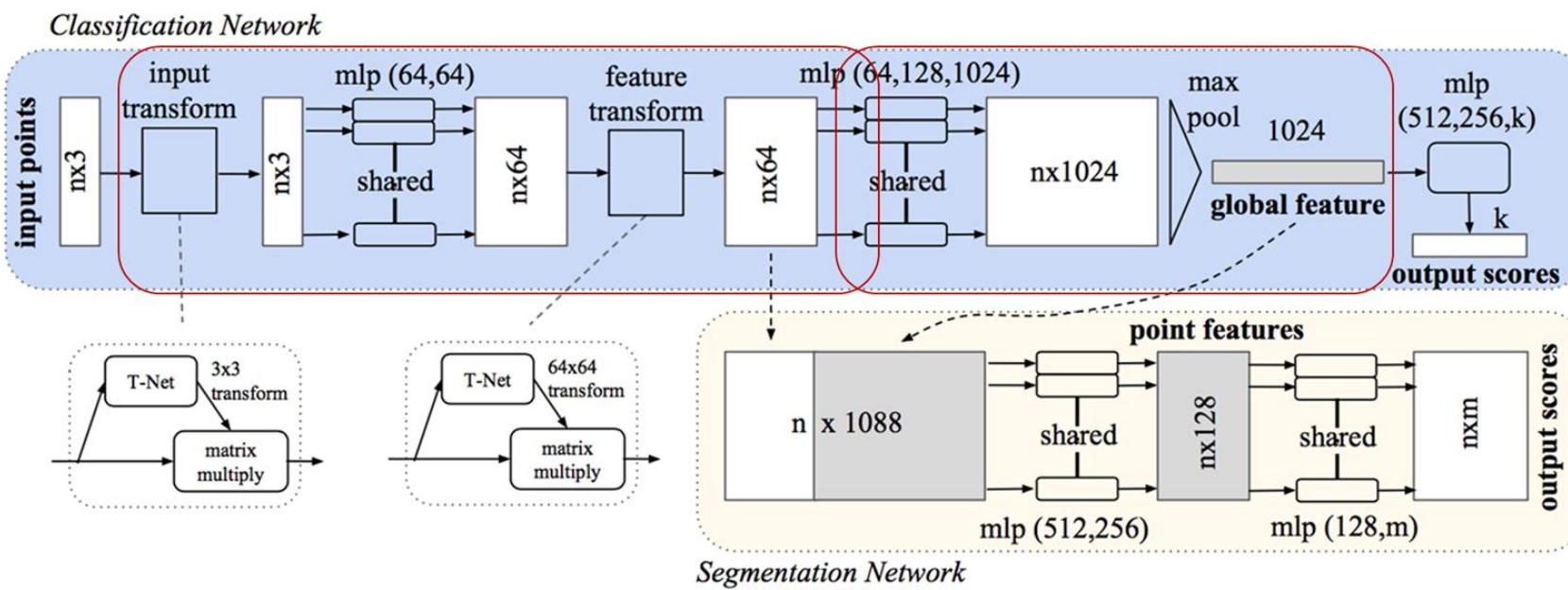
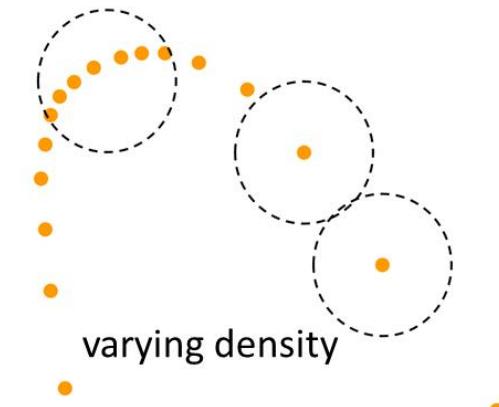
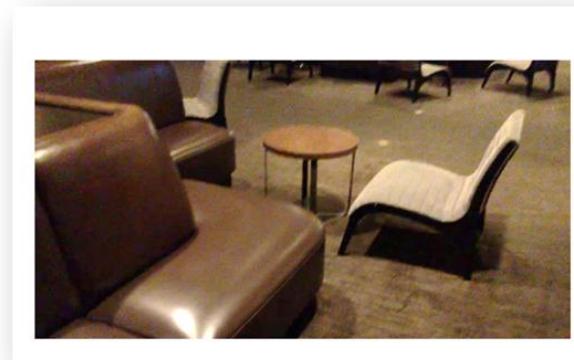
Aggregation:  $'_n = \text{Pool}_{m:(n,m) \in \mathcal{E}}(\text{---}'_{nm})$

VN-DGCNN

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

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

Deep Architectures: PointNet and PointNet++



Build VN Networks: VN-PointNet

PointNet

$$\text{█████}' = \text{Pool}_{x_n \in \mathcal{X}}(h(\text{█████}_1), h(\text{█████}_2) \dots, h(\text{█████}_N))$$

VN-PointNet

$$\text{████████}' = \text{VN-Pool}_{V_n \in \mathcal{V}}(f(\text{████████}_1), f(\text{████████}_2), \dots, f(\text{████████}_N))$$

Experiments of VN Use

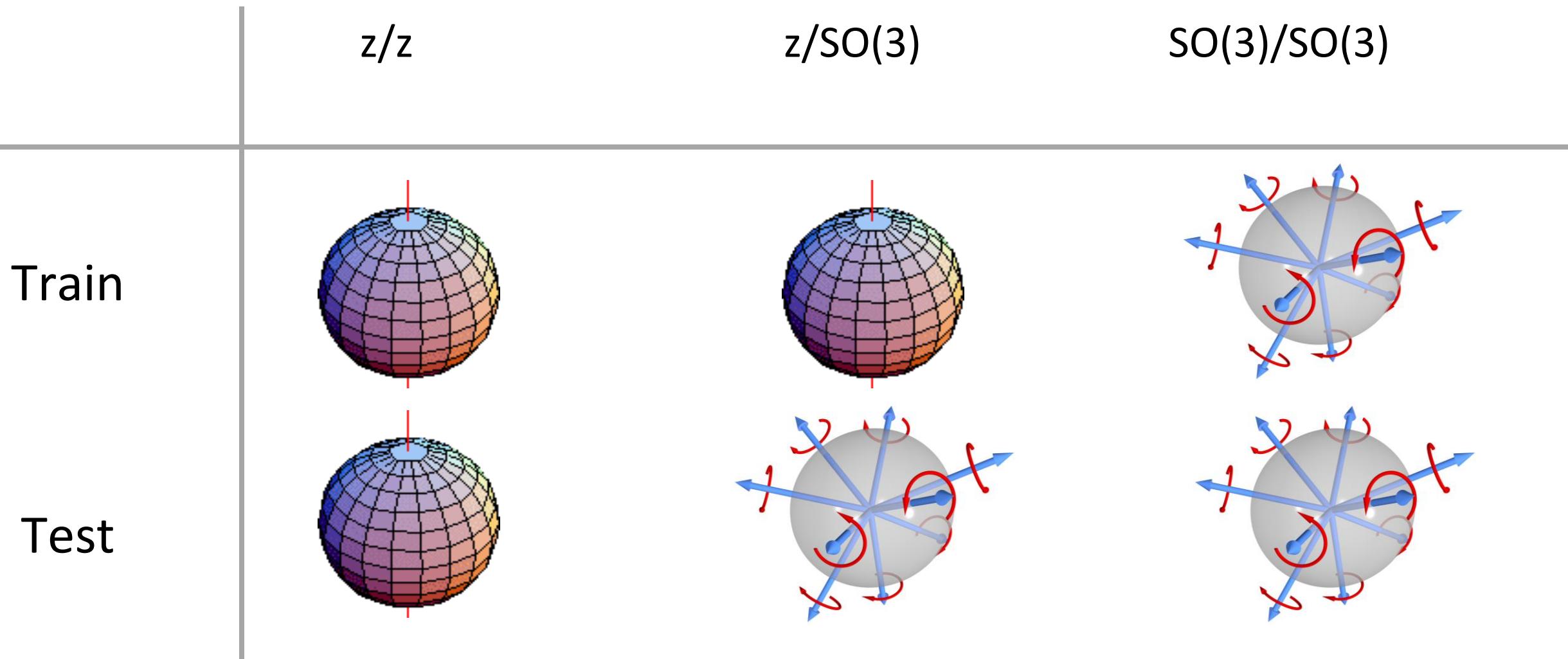


Figure: By BorisFromStockdale, CC BY-SA 3.0, <https://commons.wikimedia.org/w/index.php?curid=646939>

Figure: <https://quantum-journal.org/papers/q-2020-06-22-285/>

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^3S -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

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	$z/\text{SO}(3)$	$\text{SO}(3)/\text{SO}(3)$
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^3S -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

Part Segmentation

Results on ShapeNet (mIoU)

- Similarly...

VN Networks

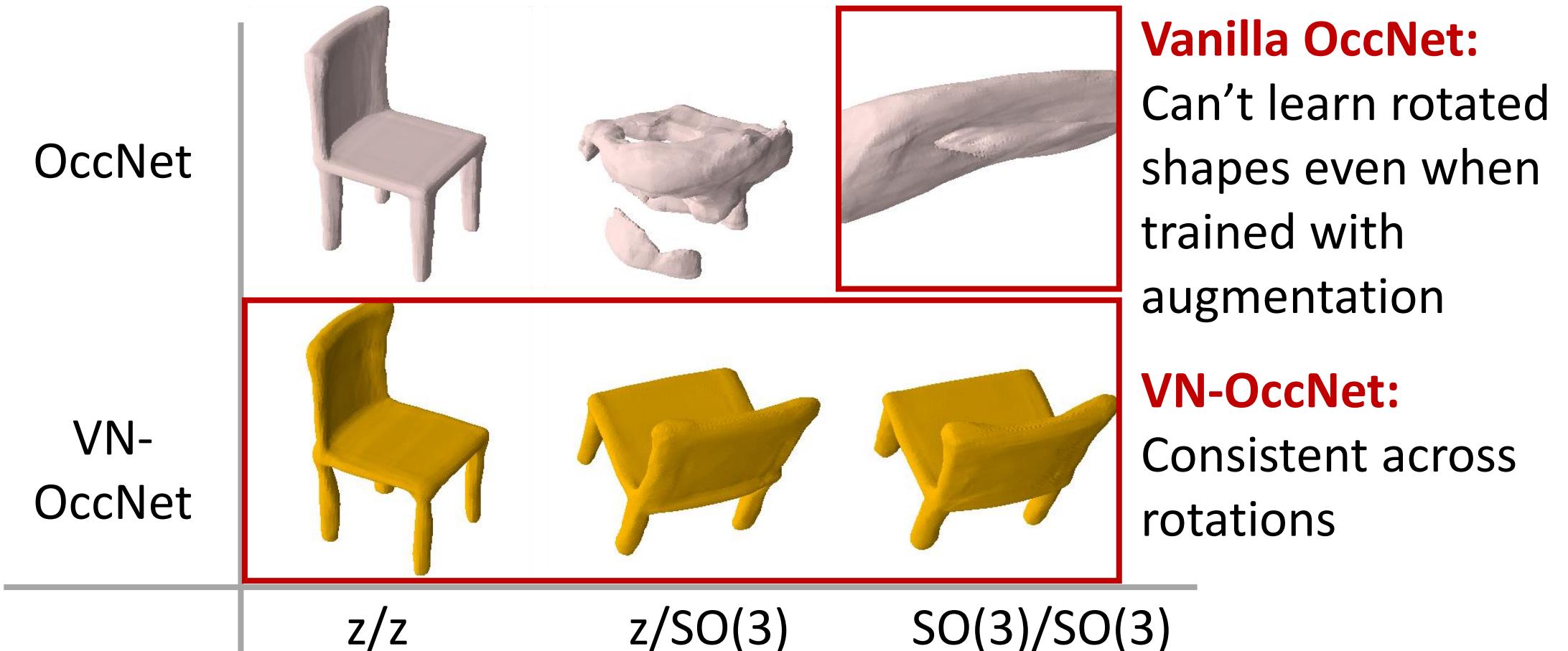
Rotation sensitive methods

Rotation robust methods

Methods	$z/\text{SO}(3)$	$\text{SO}(3)/\text{SO}(3)$
Point / mesh inputs		
PointNet [24]	38.0	62.3
DGCNN [34]	49.3	78.6
VN-PointNet	72.4	72.8
VN-DGCNN	81.4	81.4
PointCNN [19]	34.7	71.4
PointNet++ [25]	48.3	76.7
ShellNet [39]	47.2	77.1
RI-Conv [38]	75.3	75.3
TFN [31]	76.8	76.2
GC-Conv [40]	77.2	77.3
RI-Framework [17]	79.2	79.4

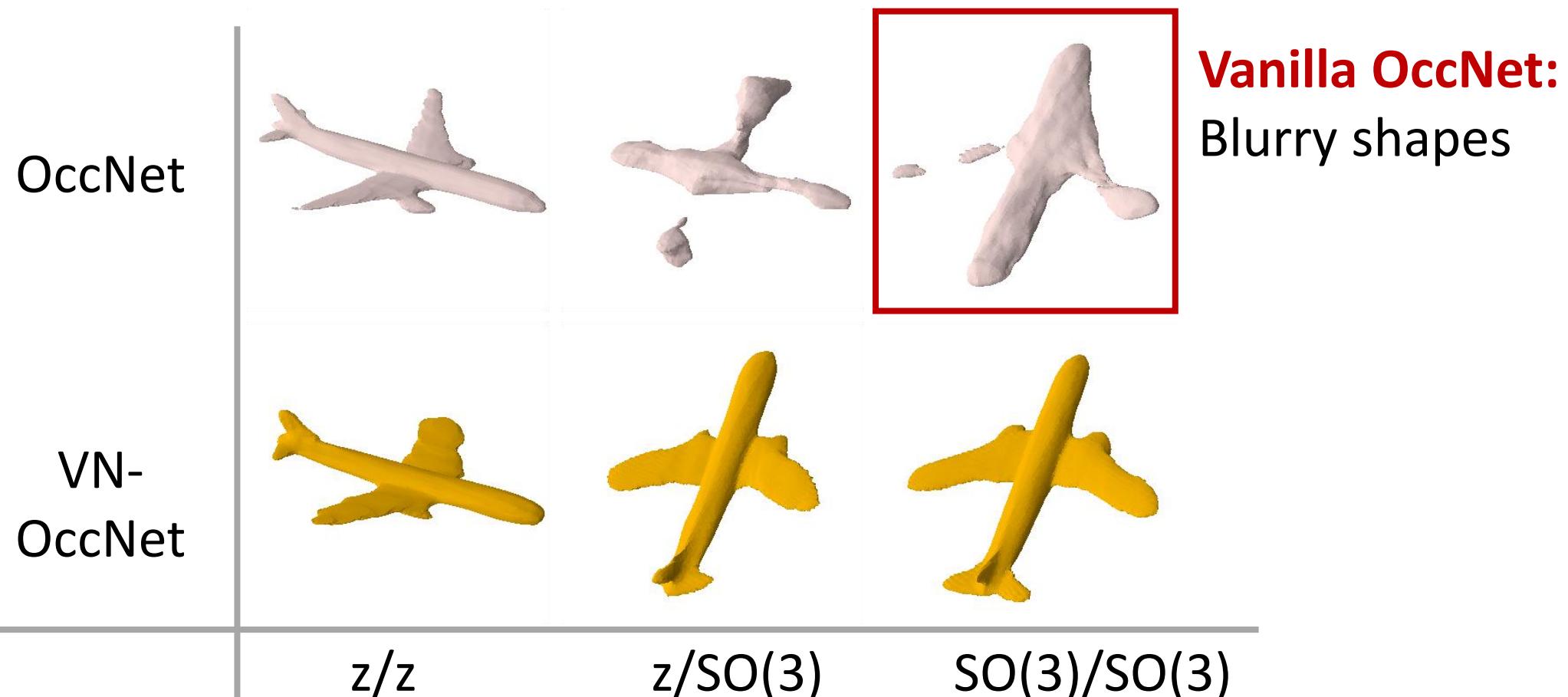
Neural Implicit Reconstruction

Results on ShapeNet (Examples)



Neural Implicit Reconstruction

Results on ShapeNet (Examples)



Applications of Vector Neurons

CVPR 2023

EFEM

Equivariant neural *Field Expectation Maximization* for 3D Object Segmentation **Without Scene Supervision**

Jiahui Lei¹ Congyue Deng² Karl Schmeckpeper¹ Leonidas Guibas² Kostas Daniilidis¹

¹ University of Pennsylvania

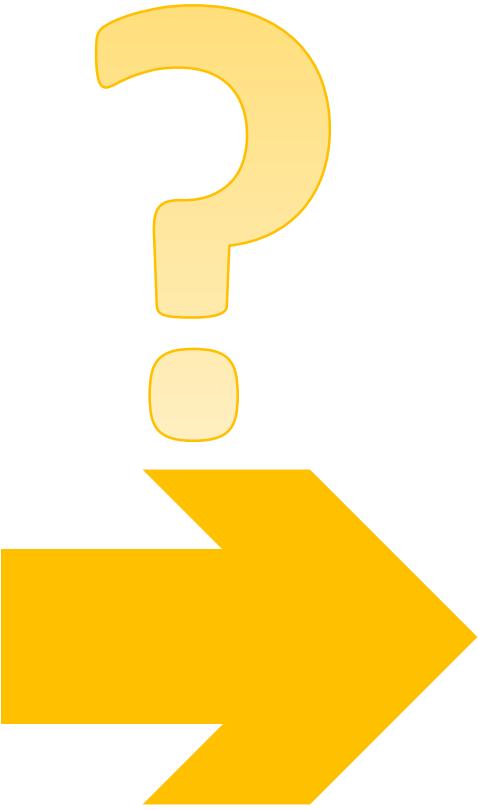
² Stanford University

{leijh, karls, kostas}@cis.upenn.edu, {congyue, guibas}@cs.stanford.edu

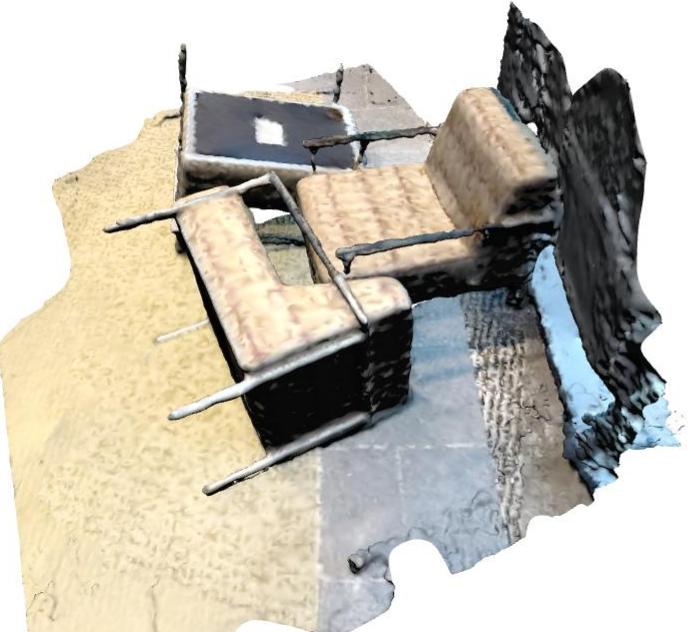


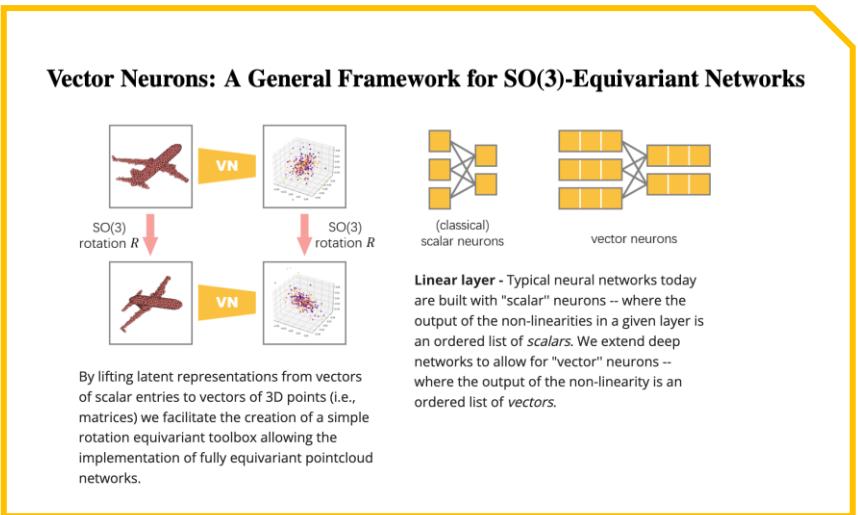
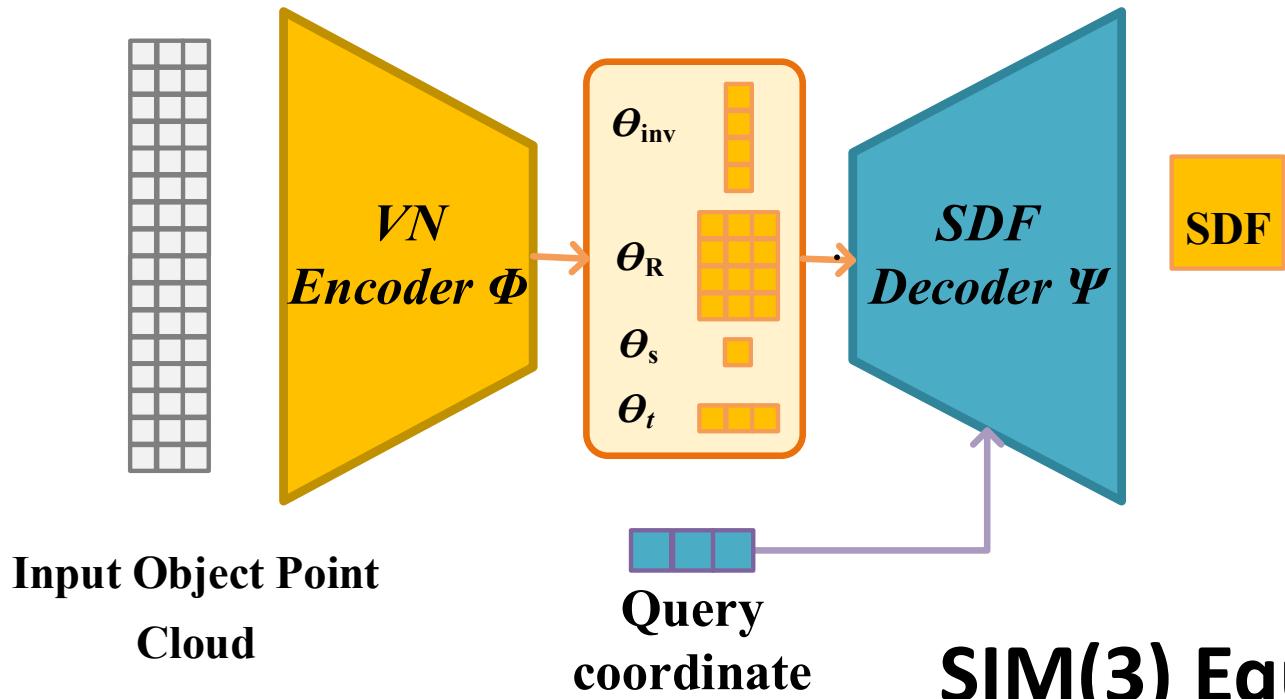


ShapeNet

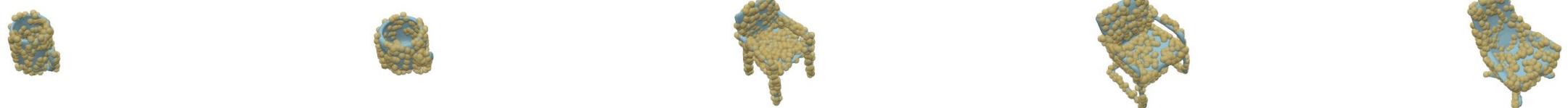


Real Scenes



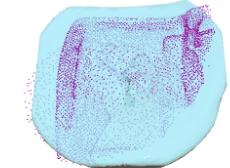
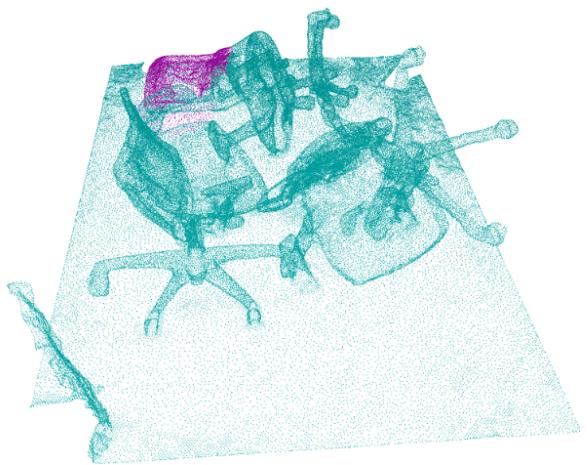


SIM(3) Equivariant Shape Prior



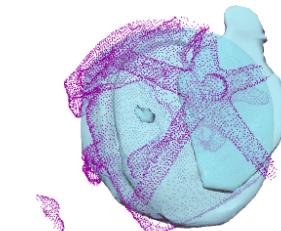
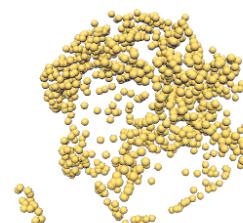
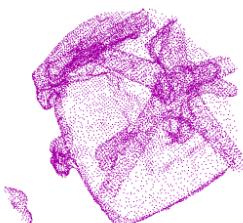
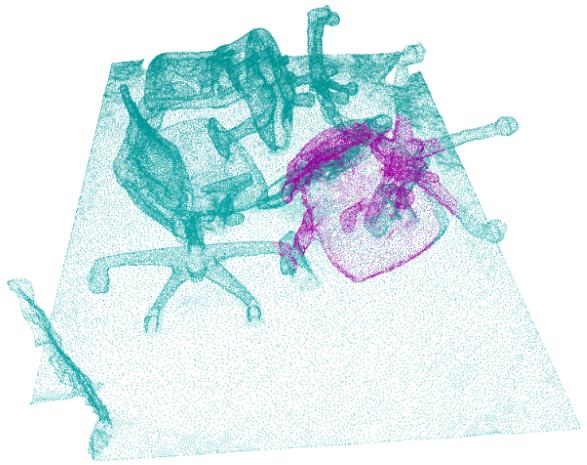
Output meshes follow input transformations

prop26-step0

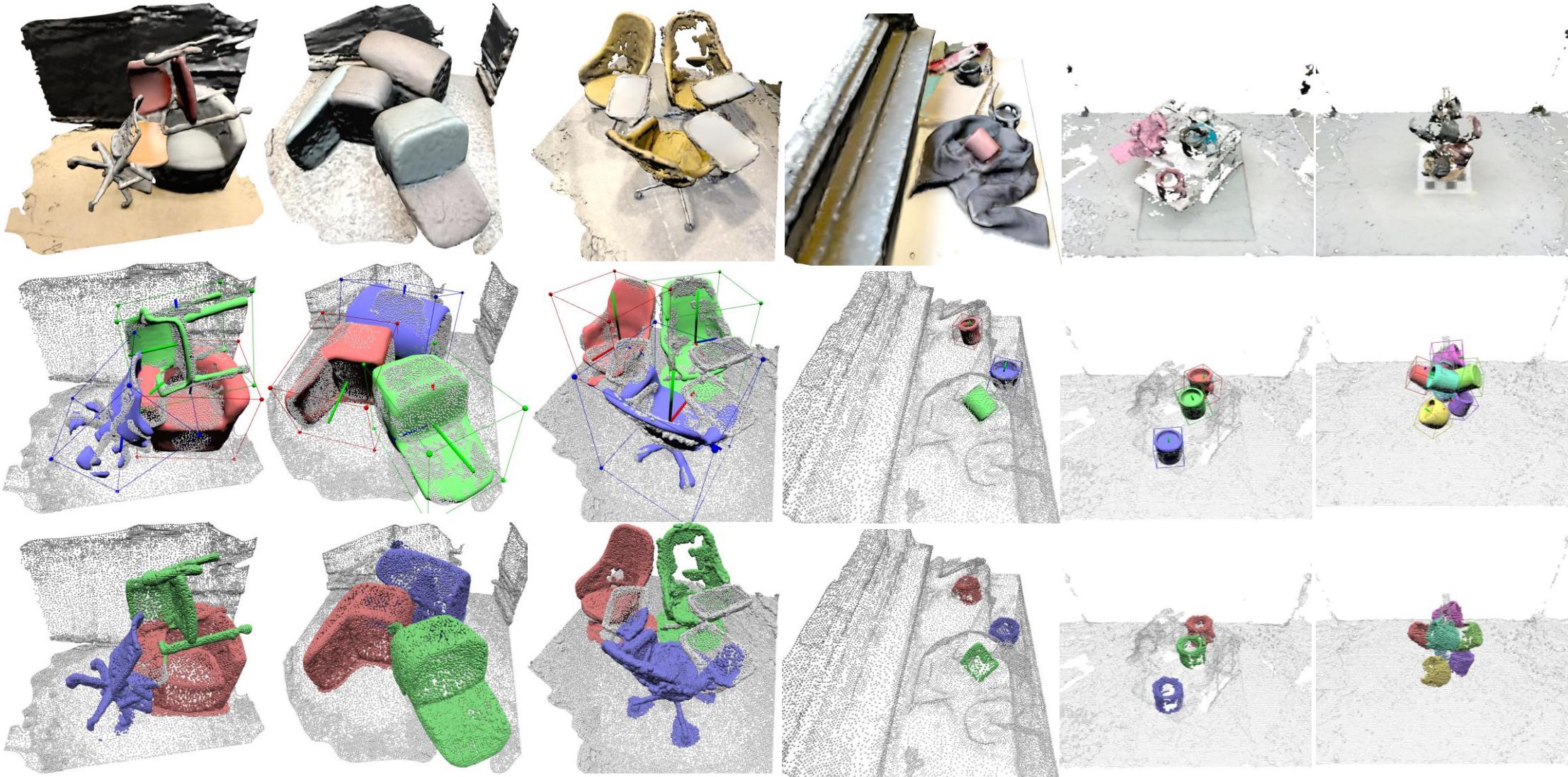


E-M iterative refinement

prop42-step0



Chairs and Mugs

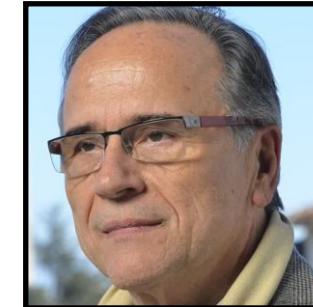
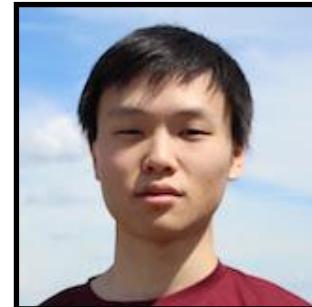
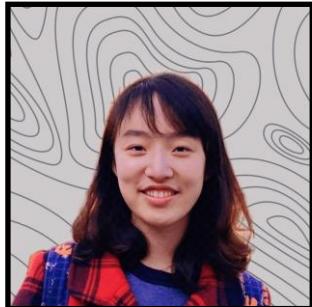
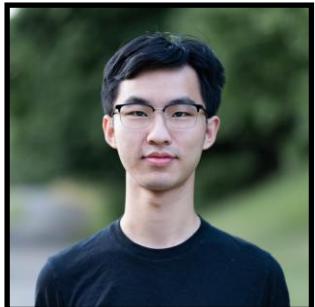


NeurIPS 2023

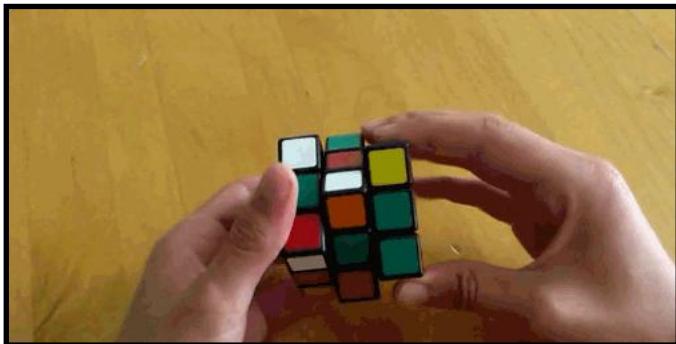
EquivAct: SIM(3)-Equivariant Visuomotor Policies beyond Rigid Object Manipulation

Jingyun Yang^{1*} Congyue Deng^{1*} Jimmy Wu² Rika Antonova¹ Leonidas Guibas¹ Jeannette Bohg¹

Stanford University¹ Princeton University²



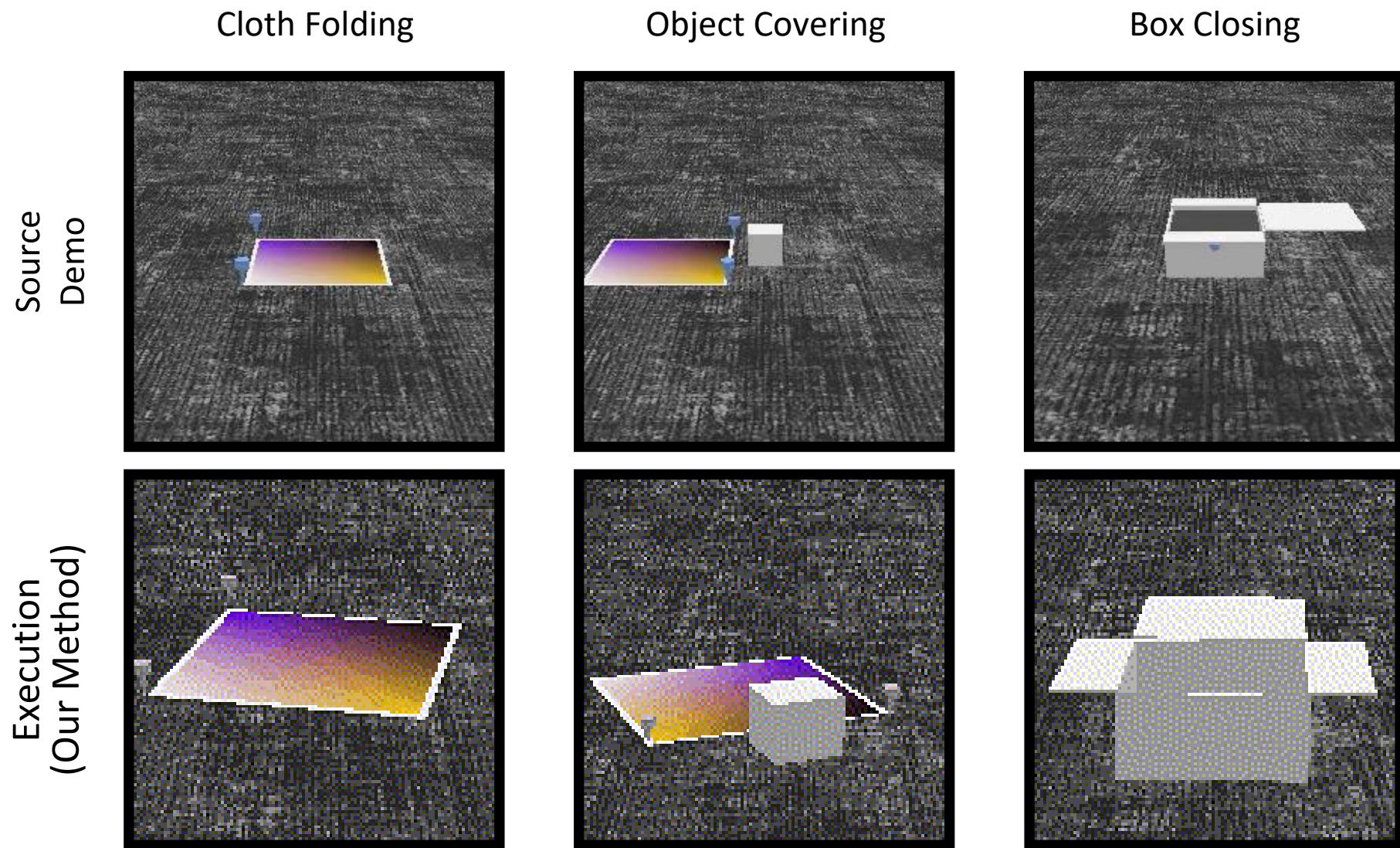
QUESTION ADDRESSED IN THIS WORK

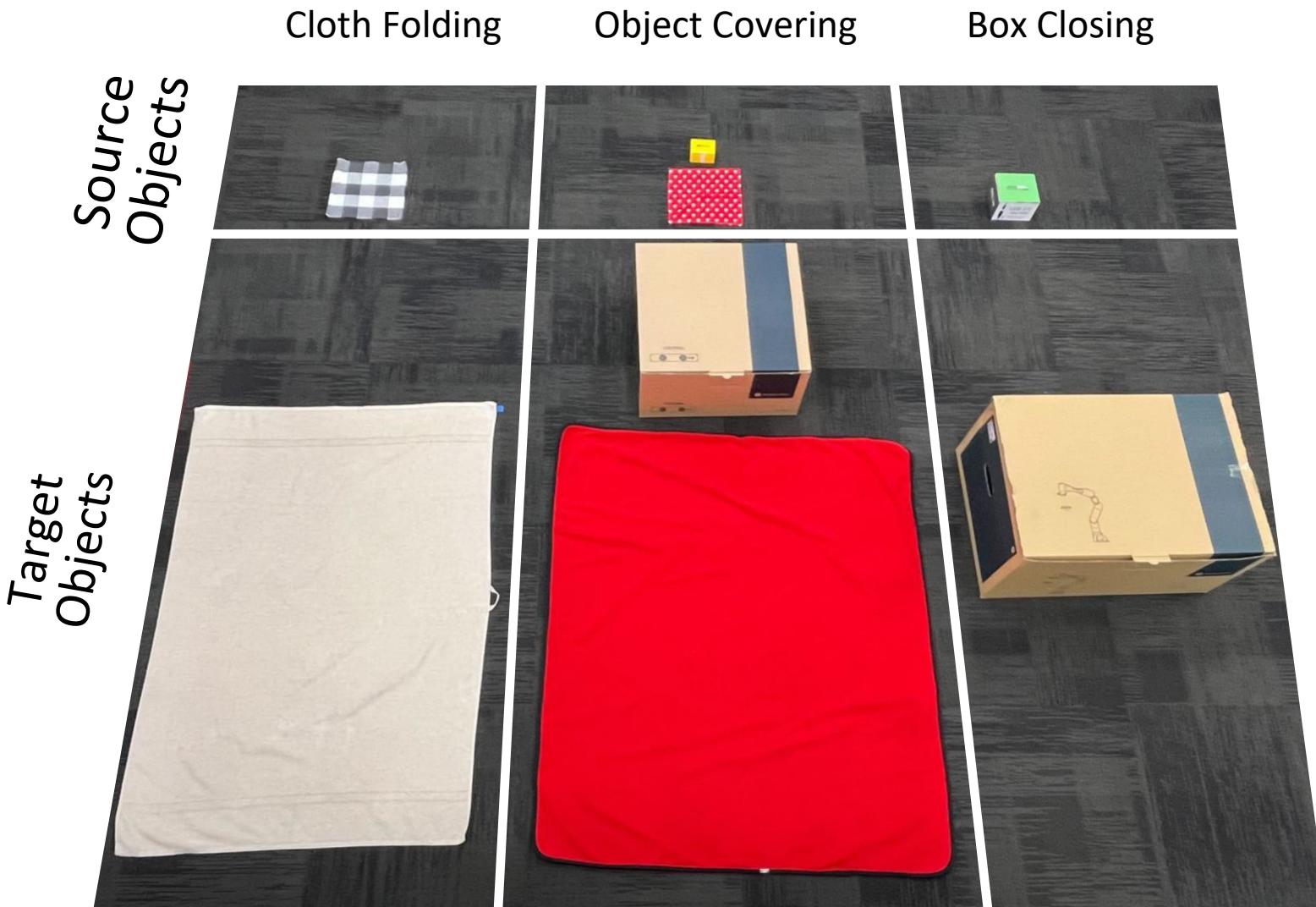


appearance
scales
poses

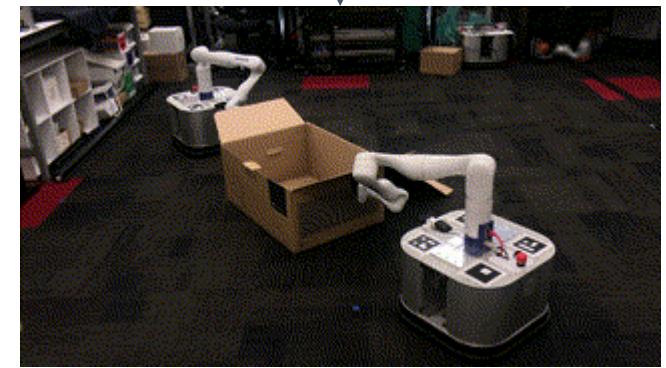
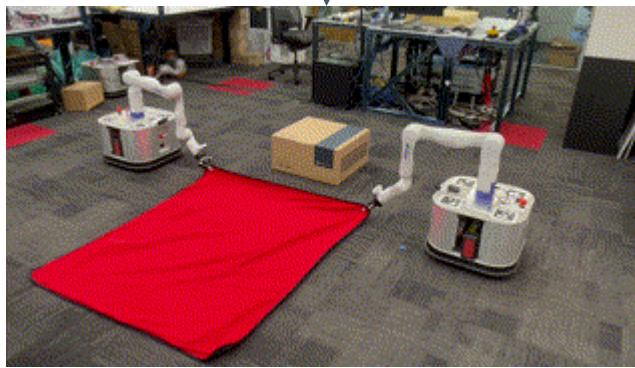
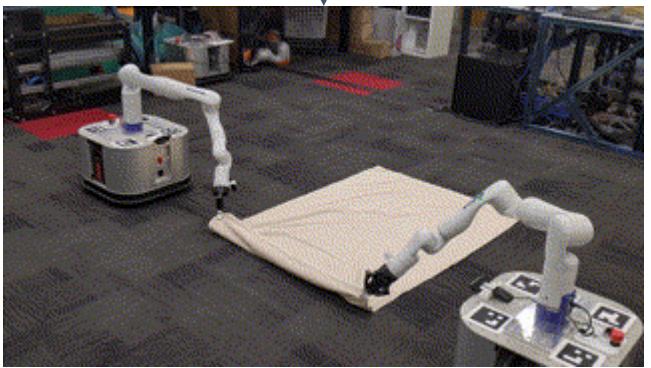
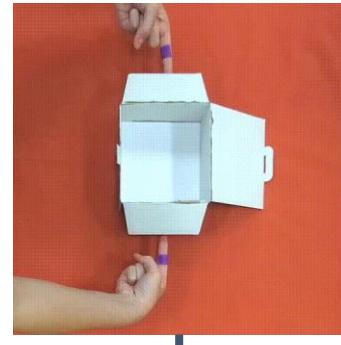


How can robots learn from a few example trajectories and generalize to scenarios with unseen visuals, scales, and poses?



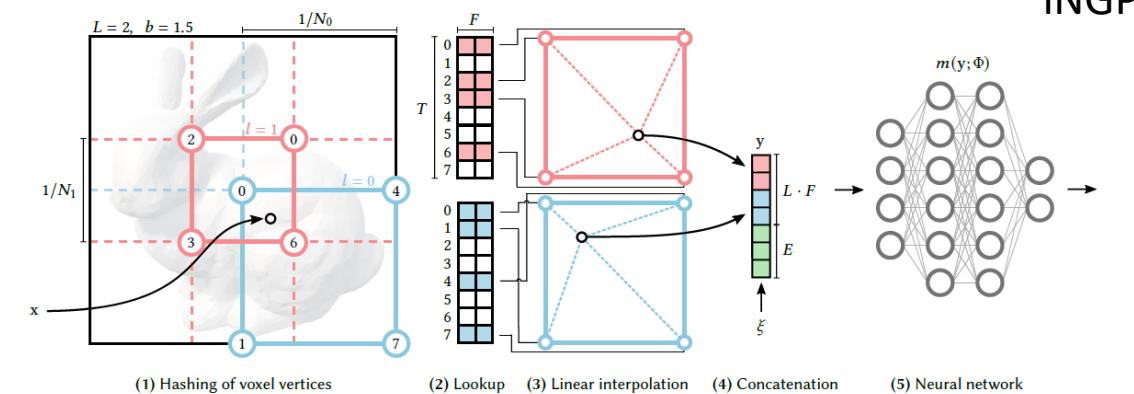
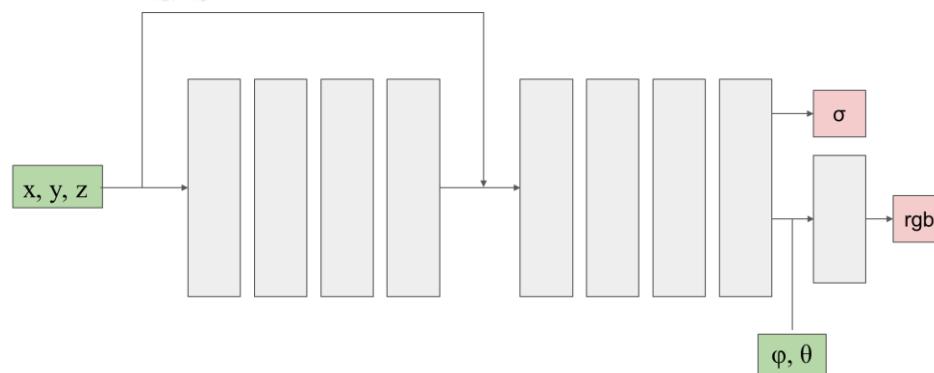
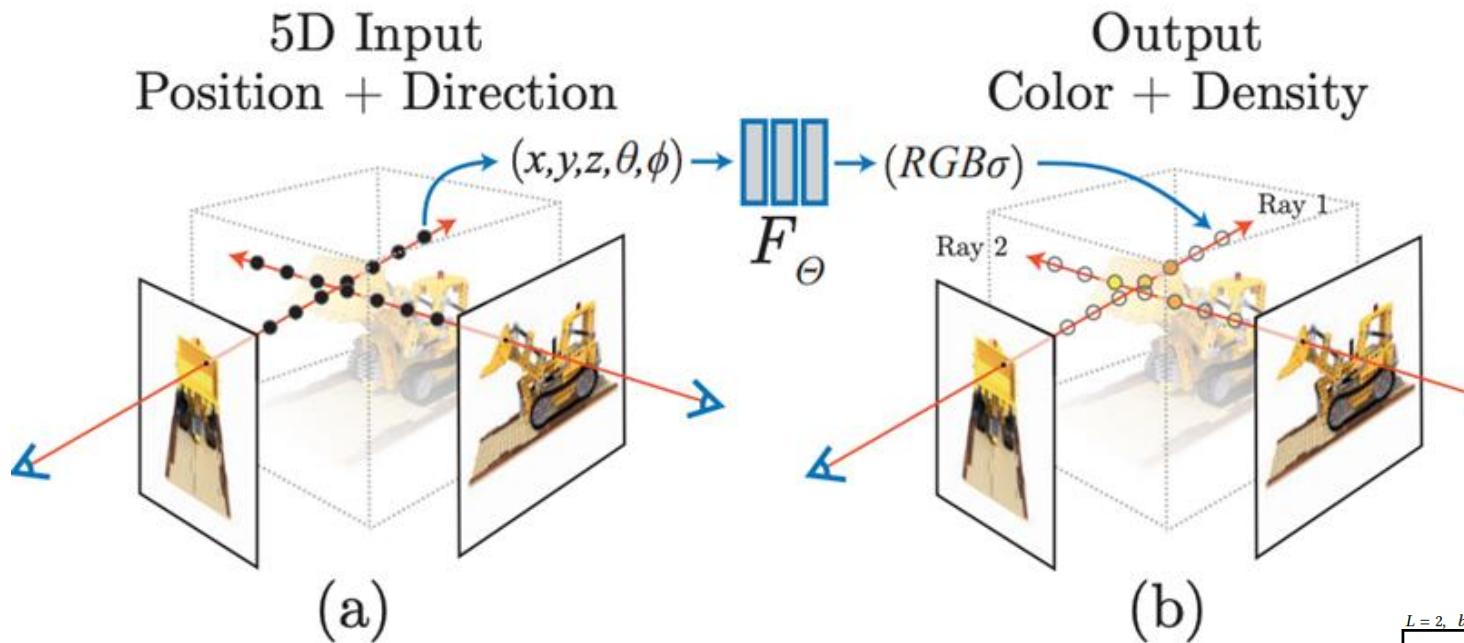


Applications in Robot Manipulation



Learned Priors: Semantic Structure and Compositionality in 3D Scenes

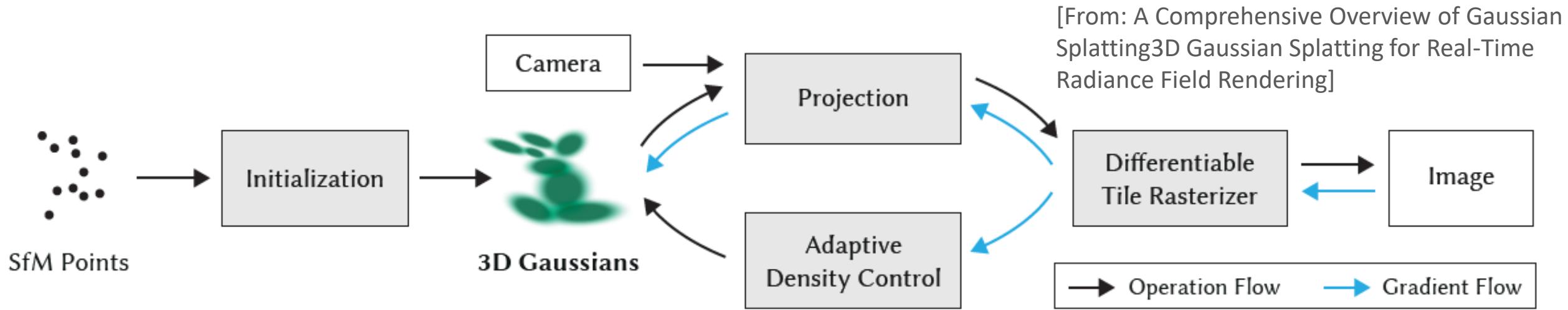
Neural Radiance Fields (NeRFs) for 3D Scenes



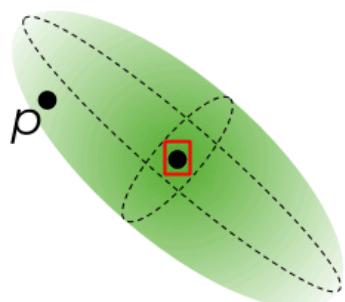
Mildenhall, B., Srinivasan, P. P., Tancik, M., Barron, J. T., Ramamoorthi, R., & Ng, R. (2021). NeRF: Representing scenes as neural radiance fields for view synthesis. Communications of the ACM, 65(1), 99-106.

Müller, T., Evans, A., Schied, C., & Keller, A. (2022). Instant neural graphics primitives with a multiresolution hash encoding. ACM Transactions on Graphics (ToG), 41(4), 1-15.

Gaussian Scene Representations: Gaussian Splatting



$$f_i(p) = \sigma(\alpha_i) \exp\left(-\frac{1}{2}(p - \boxed{\mu_i}) \Sigma_i^{-1} (p - \boxed{\mu_i})\right)$$



GS Reps (and NeRFs, too) are Unstructured



Just millions of individual Gaussians...

How can structure them?

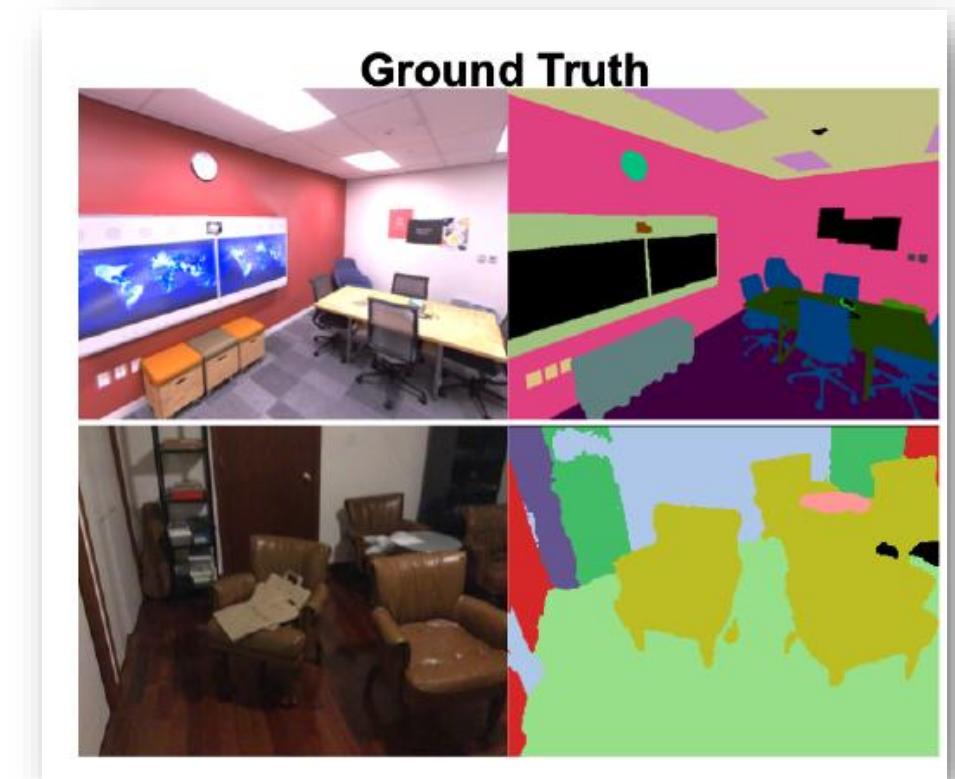
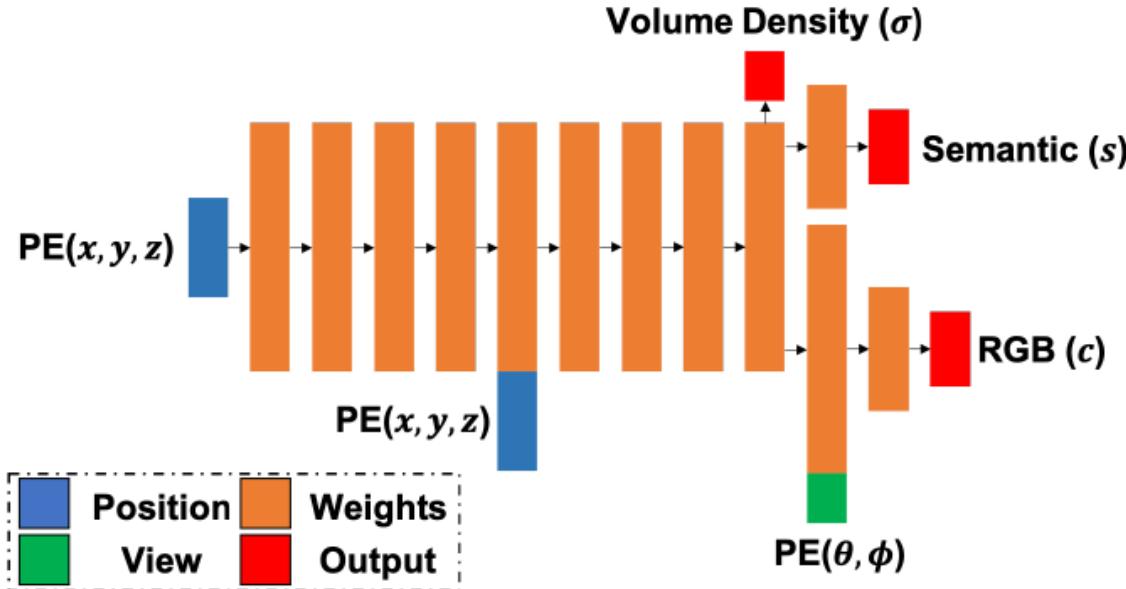
How can we manipulate the scene at the object / entity level?

NeRFs / GSs with Semantic Channels

Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, Andrew J. Davison.

In-Place Scene Labelling and Understanding with Implicit Scene Representation.
ICCV 2021.

Semantic NeRF



Sosuke Kobayashi, Eiichi Matsumoto, Vincent Sitzmann.
Decomposing NeRF for Editing via Feature Field Distillation.
NeurIPS 2022

NeRF Shaping with Sparse Semantic Supervision

Xiamen Xu, Yanchao Yang, Kaichun Mo, Boxiao Pan, Li Yi, Leonidas Guibas.
JacobiNeRF: NeRF Shaping with Mutual Information Gradients
(CVPR 2023).

Supervision for NeRFs: 1st Order vs 2nd Order

- In typical NeRF training, we supervise via pixel values in views (1st order info):

- this pixel's color is red ...
 - this pixel's semantics is "car" ...



- But we can also supervise with value relationships (2nd order info)

- these two pixels should have the same entity ID ...



- Semantics / composition can be implicitly encoded in the correlations, or mutual information, between pixels



- How can we directly encode into the neural representation semantic correlations?
 - the NeRF variation space

An Indoor Scene

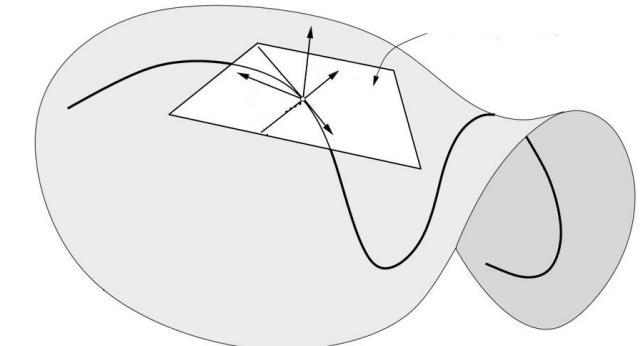


Understand how the scene is, *and how it could be ...*

The Compositional Structure of a Scene is Reflected in its Variations



The scene semantic variations:
tangent space



this chair has moved

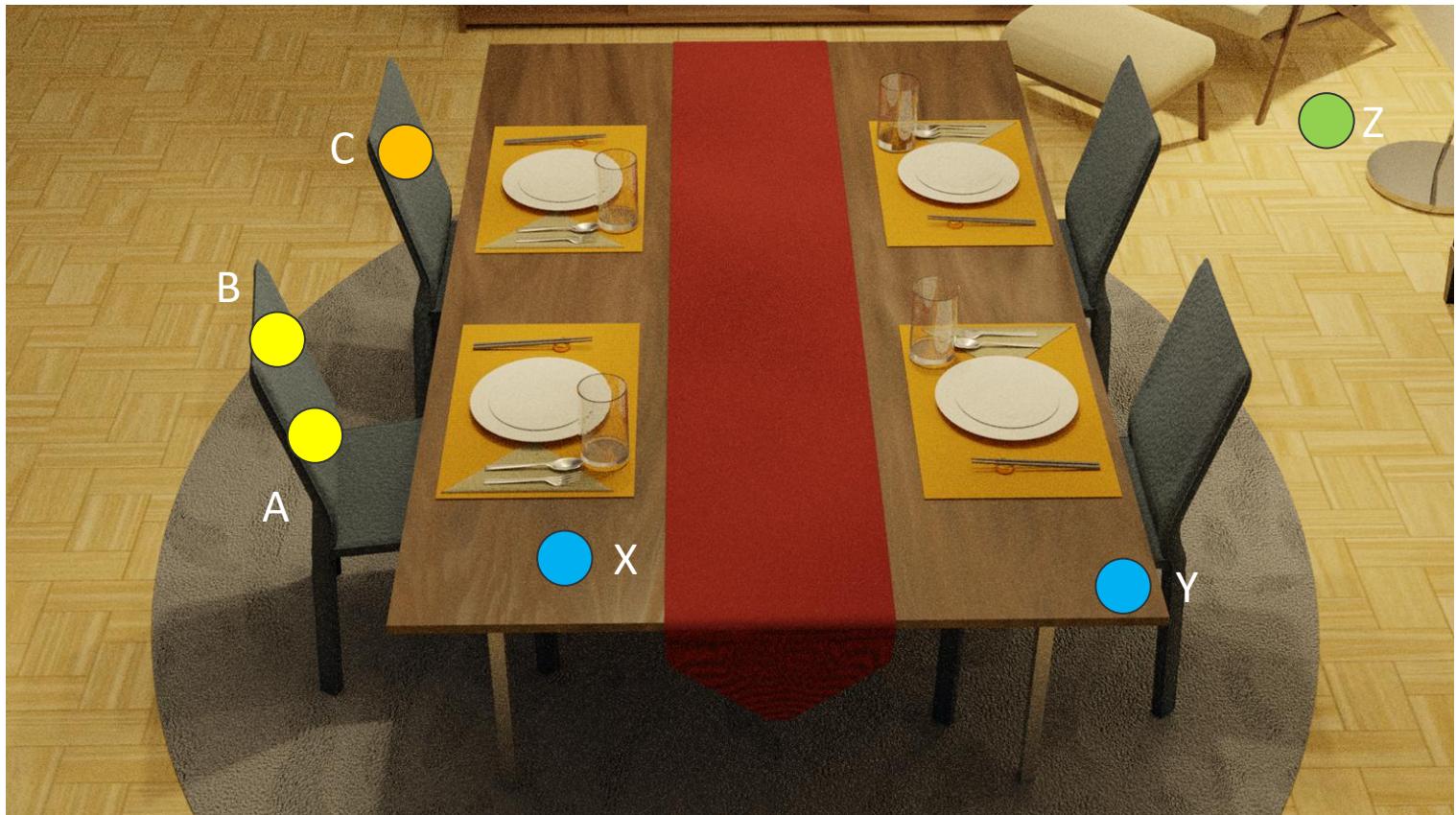


the table became longer



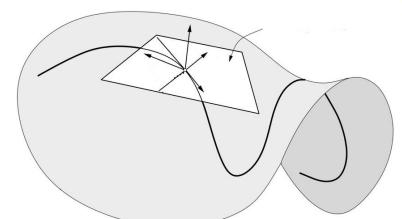
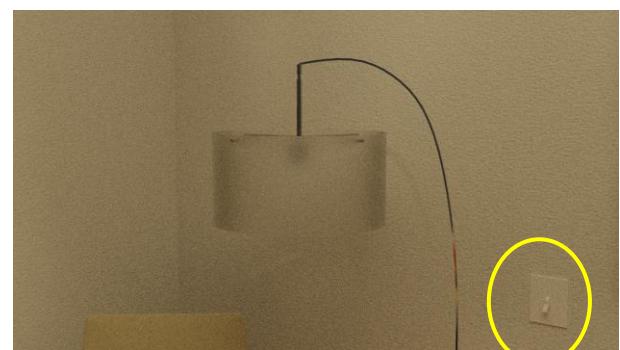
the table became darker

Mutual Information and 2nd Order Relationships

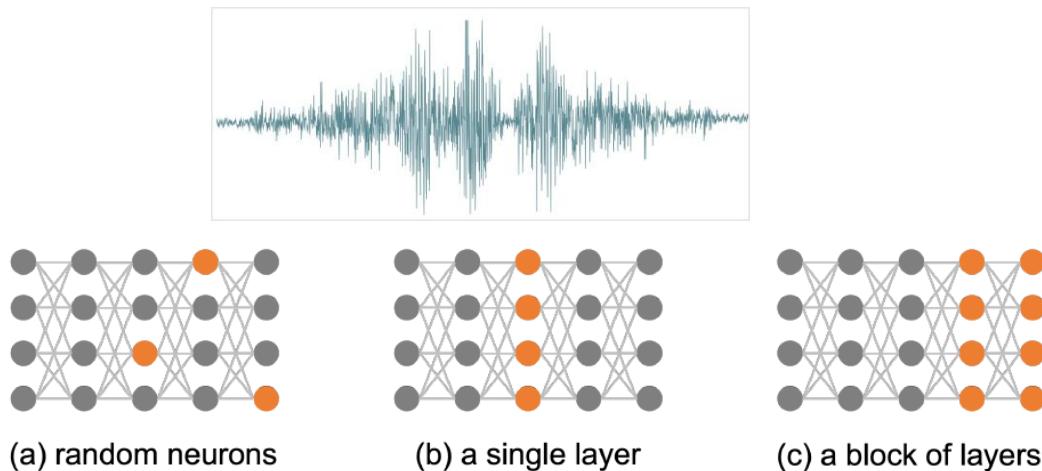


A is more correlated with B than with C

X is more correlated with Y than with Z

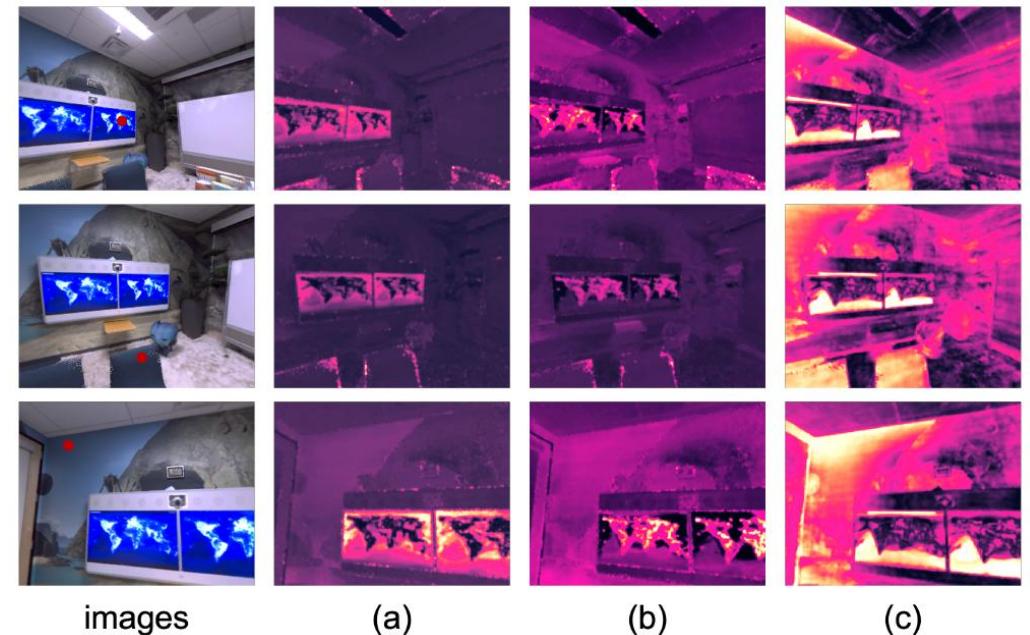


NeRF Variation Space Through its Parameters



Unfortunately, these variations are not semantically meaningful

“Jiggle” the neuron weights of the NeRF



Key Idea: Operate on the NeRF to modify its weights, so as to align the scene semantic variation space with the NeRF parametric variations

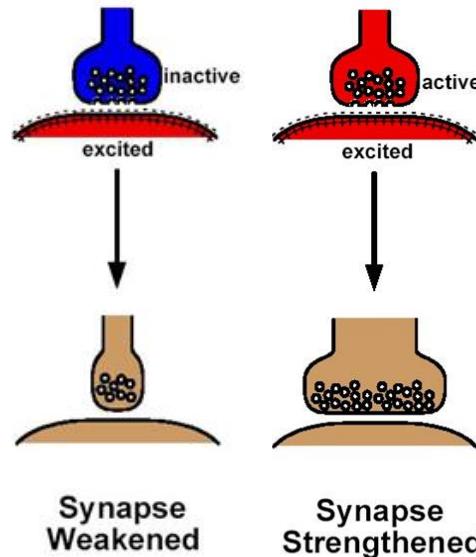


“NeRF Shaping”
NeRF Operators

Shaping Neural Representations

The Hebbian hypothesis:

Neurons that fire together, wire together



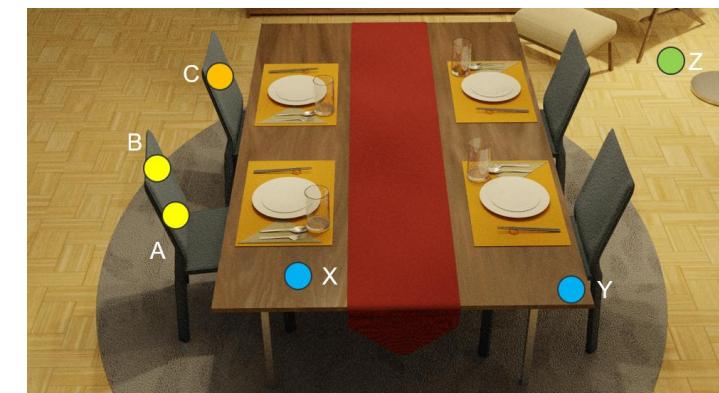
Donald O. Hebb
1904-1985

Can we build neural scene representations that better reflect mutual information correlations in the scene?

Create “neuronal resonances” through contrastive supervision

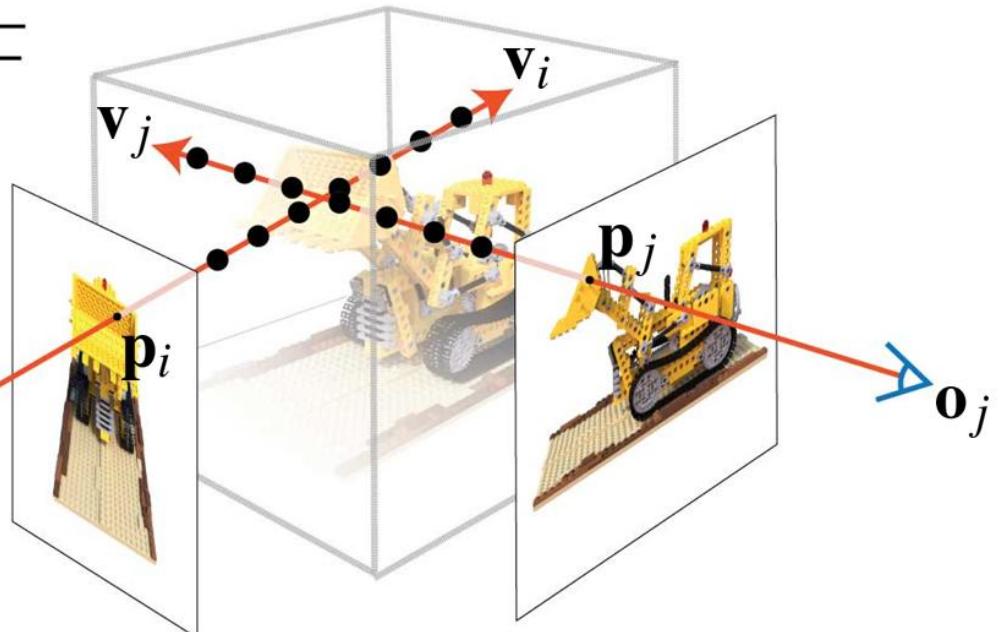
$$\mathbb{I}(A,B) > \mathbb{I}(A,C)$$
$$\mathbb{I}(X,Y) > \mathbb{I}(X,Z)$$

Mutual information \mathbb{I}



Mutual Information via NeRF Gradients

NeRF



$$\mathbf{p}(t) = \mathbf{o} + t\mathbf{v} \mid t \geq 0$$

$$I(\mathbf{p}) = \Phi(\mathbf{o}, \mathbf{v}; \theta) = \int_0^{+\infty} w(t; \theta) c(\mathbf{p}(t), \mathbf{v}; \theta) dt$$

$$I(p_i) = \Phi(\mathbf{o}_i, \mathbf{v}_i; \theta)$$

$$I(p_j) = \Phi(\mathbf{o}_j, \mathbf{v}_j; \theta)$$

Mutual information \mathbb{I}

$$\hat{I}(p_i) = \Phi(\mathbf{o}_i, \mathbf{v}_i; \theta^D + \mathbf{n})$$

$$\hat{I}(p_j) = \Phi(\mathbf{o}_j, \mathbf{v}_j; \theta^D + \mathbf{n})$$

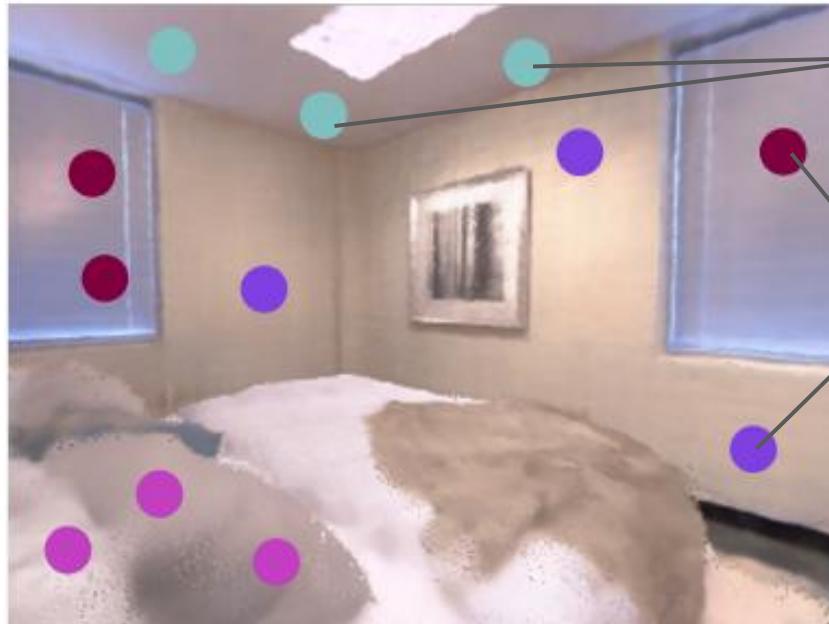
$$\mathbb{I}\left(\hat{I}(p_i), \hat{I}(p_j)\right) \approx \cos\left(\frac{\partial \Phi_i}{\partial \theta^D}, \frac{\partial \Phi_j}{\partial \theta^D}\right)$$

D = set of parameters selected for shaping
 \mathbf{n} = noise added

Inter-pixel correlations are captured by cosine similarity of the NeRF Jacobians

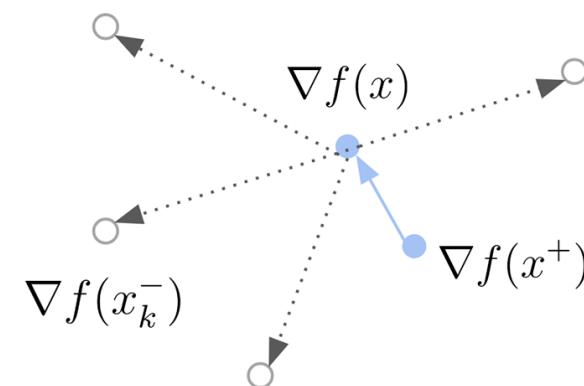
NeRF Mutual Information Shaping

Setting up Semantic “Neuronal Resonances” for Correlated Pixels through Contrastive Learning

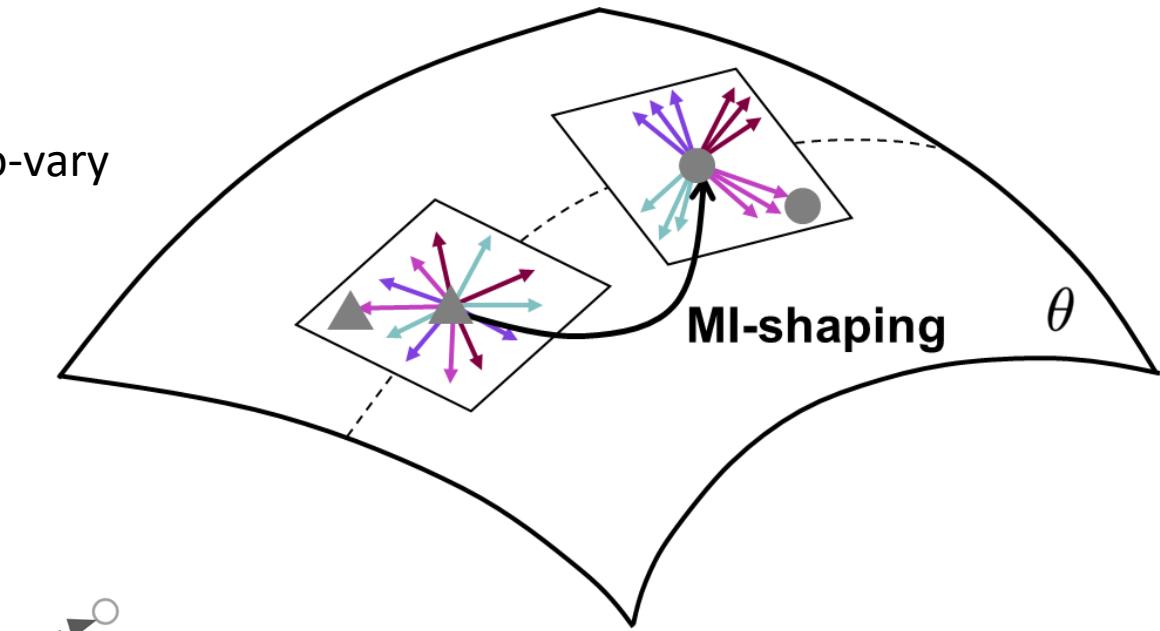


These pixels should co-vary

But these should not



Use DINO features
(alternatively, LSeg-CLIP)

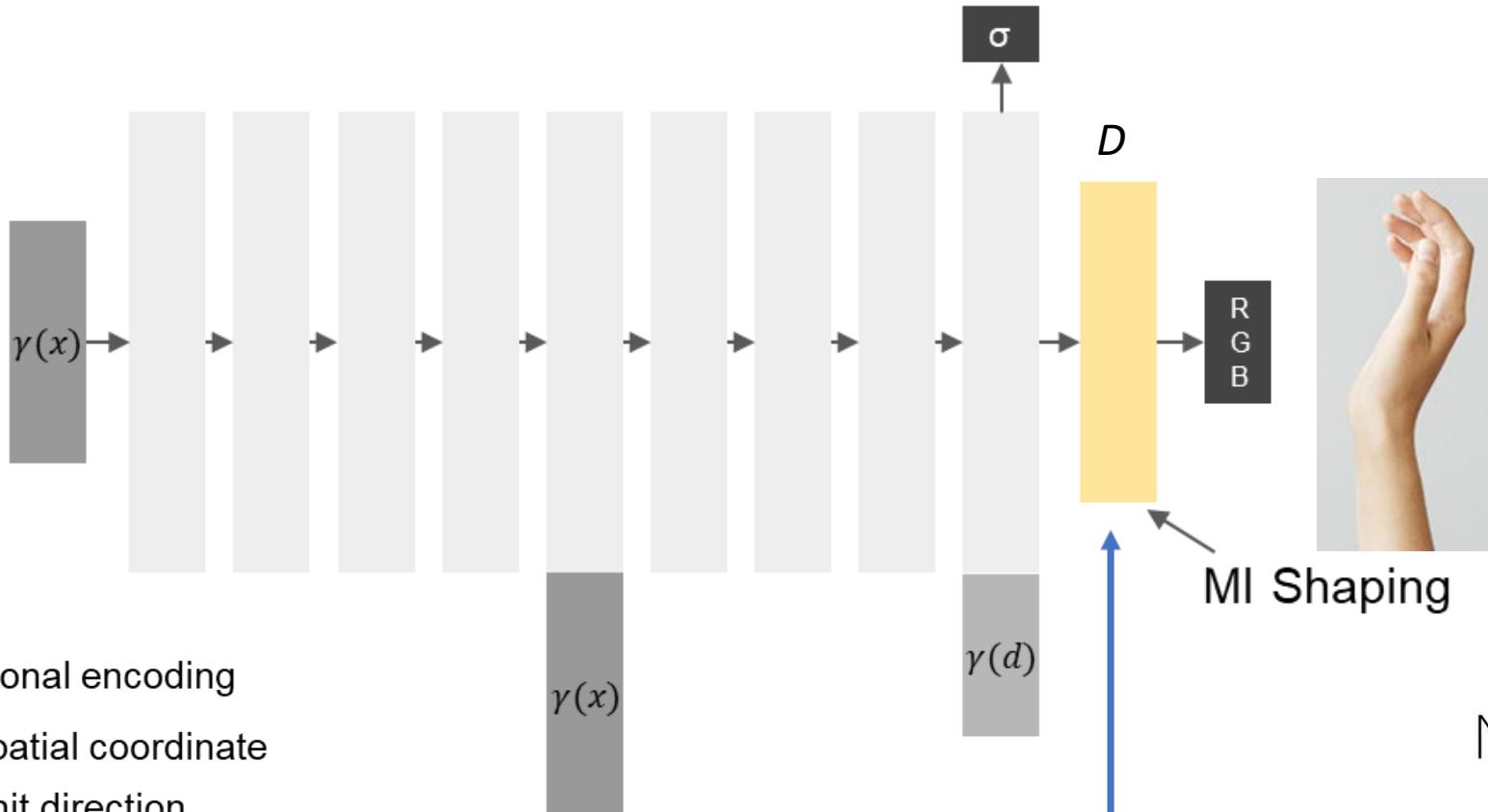


Shaping co-aligns gradients of correlated pixels (here points of the same semantic class)

JacobiNeRF (J-NeRF): Shaping via Mutual Information Jacobians

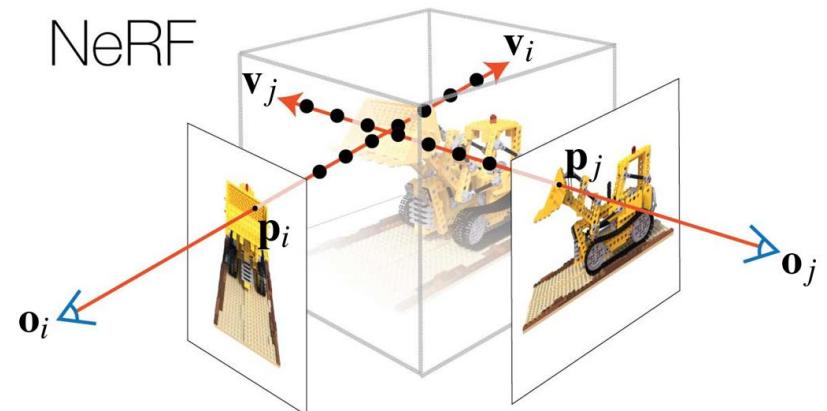


γ positional encoding
 x 3D spatial coordinate
 d 2D unit direction
 σ point density



$$\cos \left(\frac{\partial \Phi_i}{\partial \theta^D}, \frac{\partial \Phi_j}{\partial \theta^D} \right)$$

Apply gradient shaping on pixel gray scale values
(last three layers of RGB branch)



Carl Gustav Jacob Jacobi
1804-1851

Lightweight Contrastive Training Regimen (After Photometric)

Use DINO features for contrastive learning

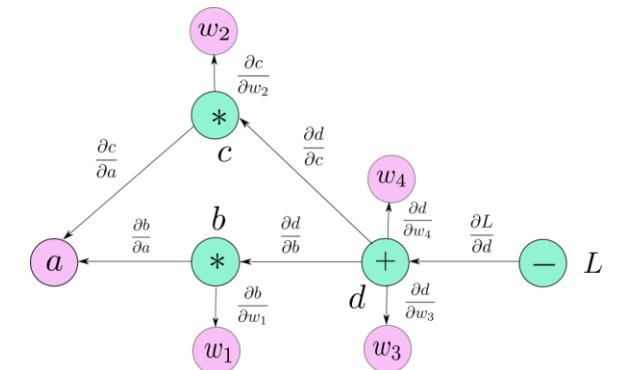
64 batches of 64 rays/pixels across all views, or

64 batches of 64 rays/pixels all in one view

Gradients are on pixel gray level

InfoNCE loss

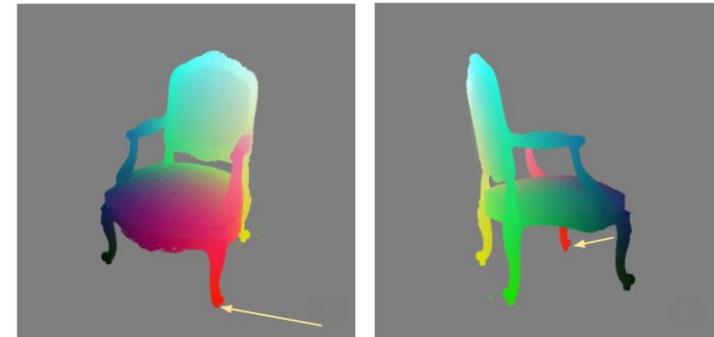
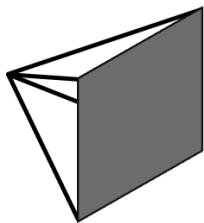
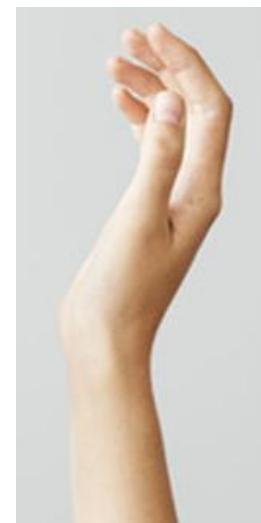
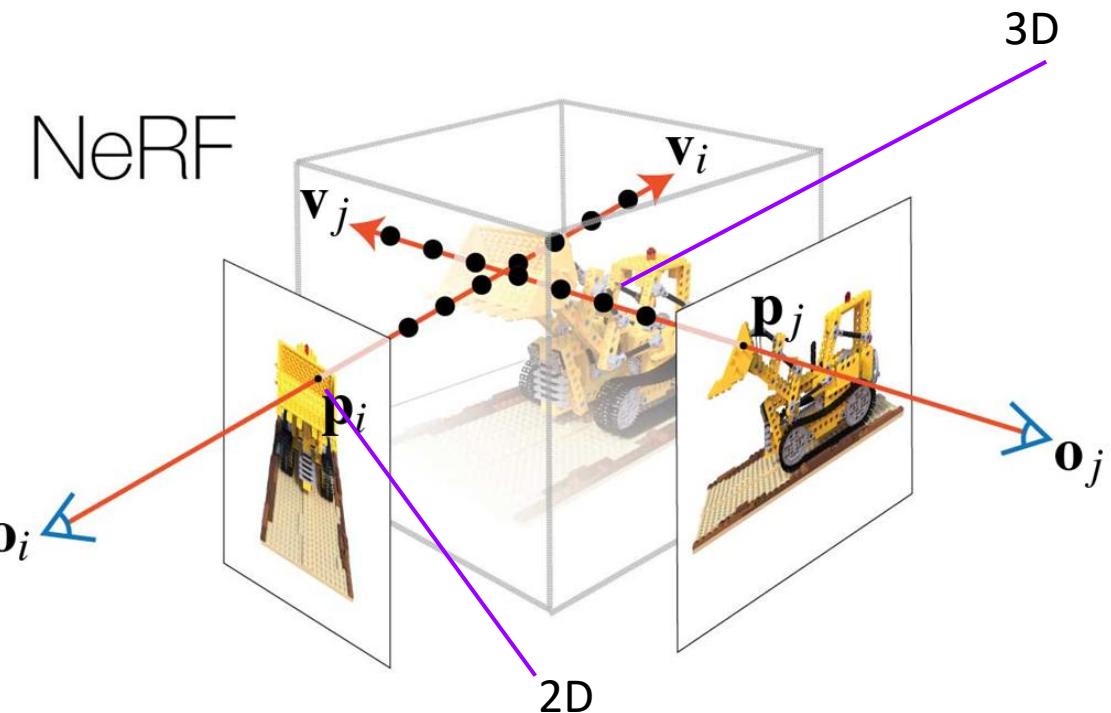
Exploit Autograd for gradients



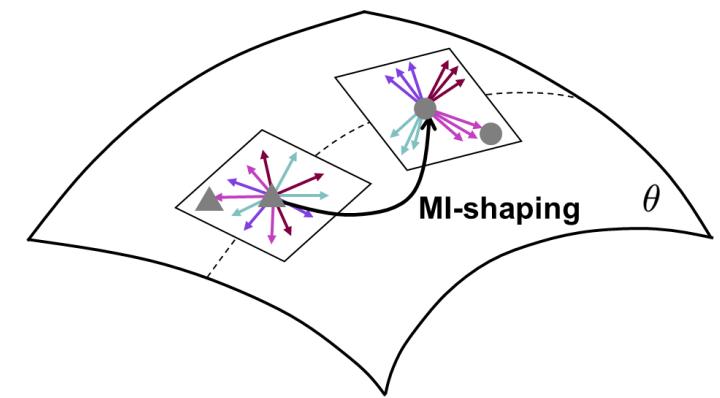
Resonances are transitive

2D vs 3D Shaping

Use NOCS to lift 2D pixels to 3D points.

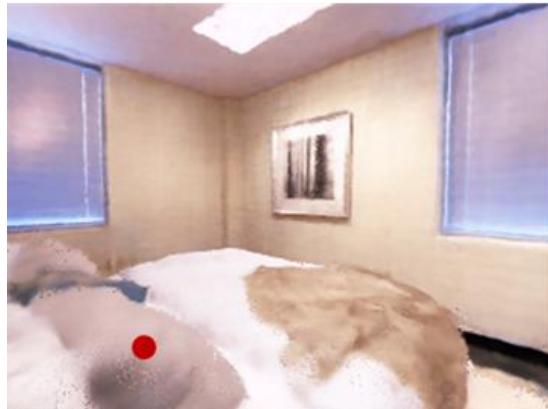
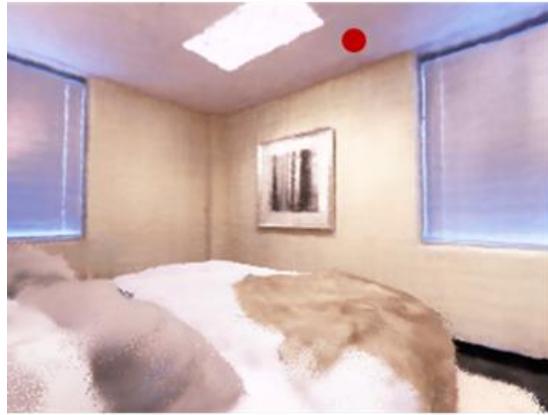


NOCS (x, y, z) coordinates rendered as RGB



Shaping can be applied to either 2D (J-NeRF 2D) or 3D (J-NeRF 3D)
network gradients

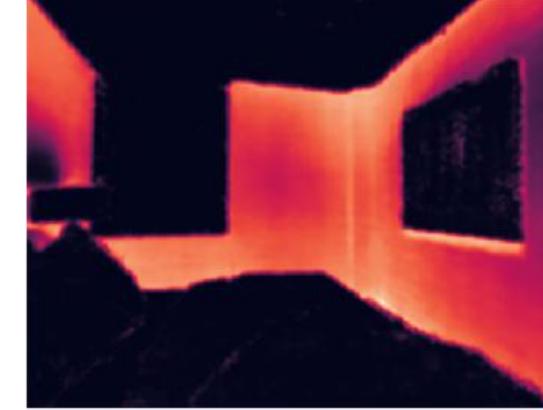
Image View



Unshaped NeRF



Shaped NeRF



After shaping:

From a single pixel
we can select an
entire semantic
entity.

Re-coloring via Resonances

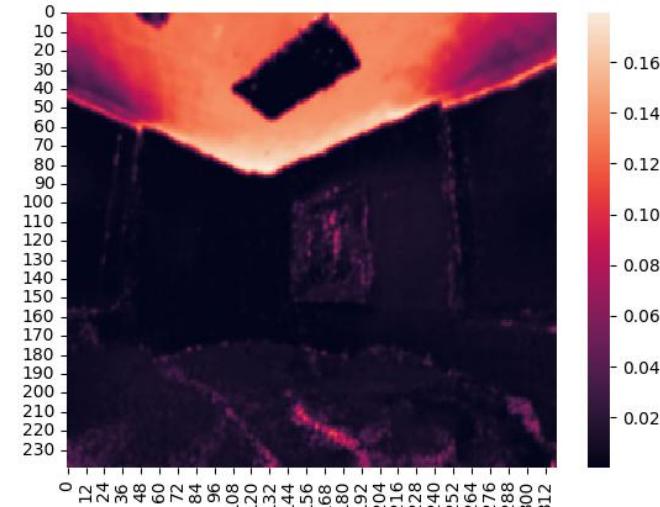
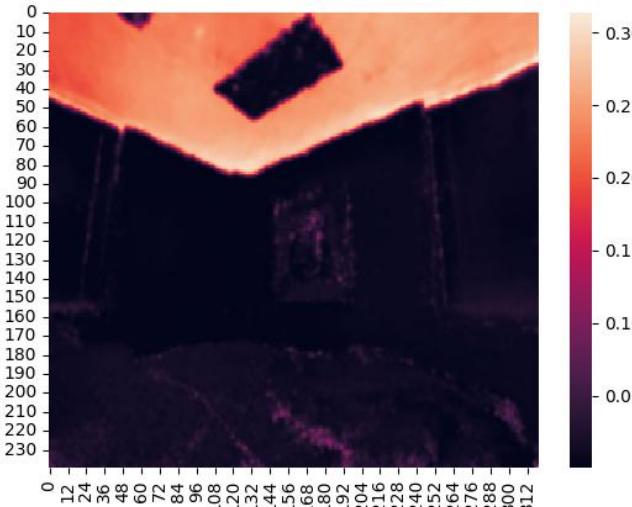
NeRF re-coloring after shaping

Shape NeRF by aligning grey scale gradients of correlated pixels/points [same as before].

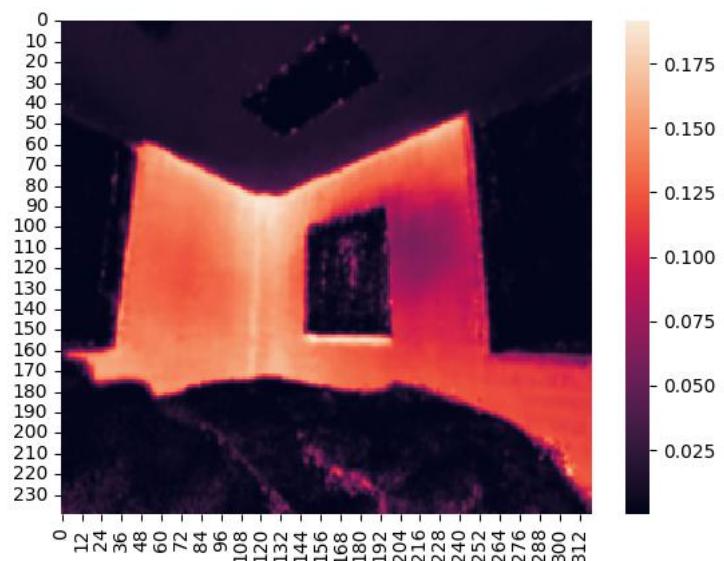
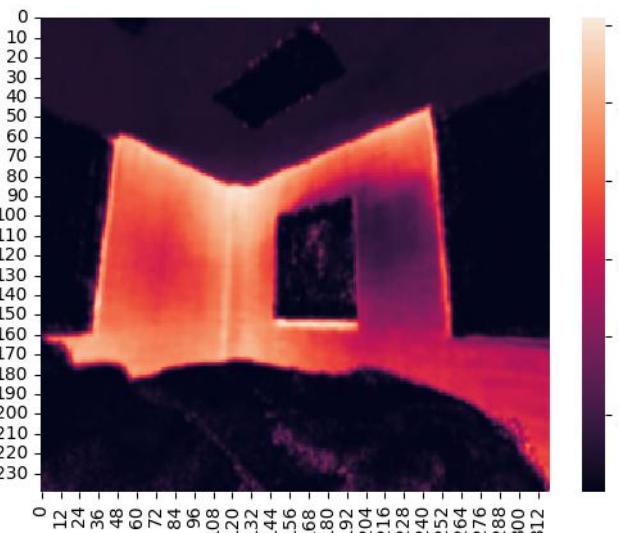
Calculate separate R, G, B gradients; select **one pixel in one view** and push the network parameters along these gradients to reach a desired color value at that pixel.

All “resonating” pixels in this and other NeRF views get also automatically recolored ...

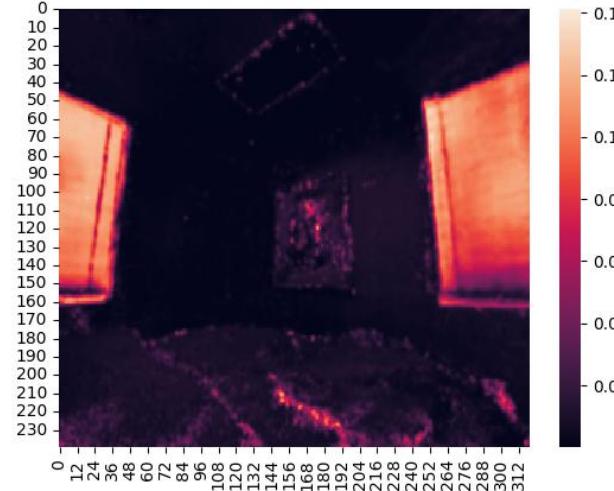
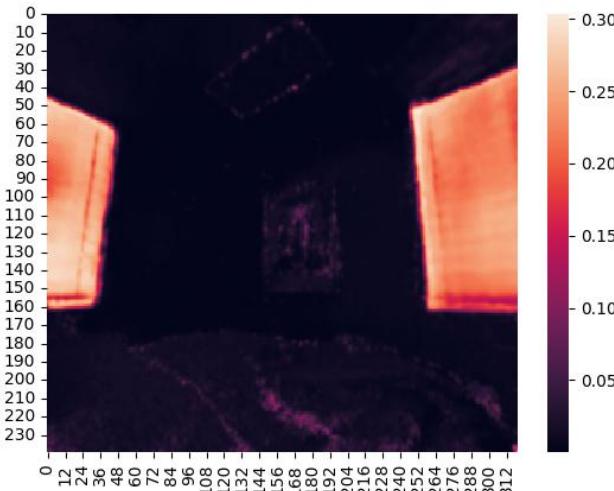
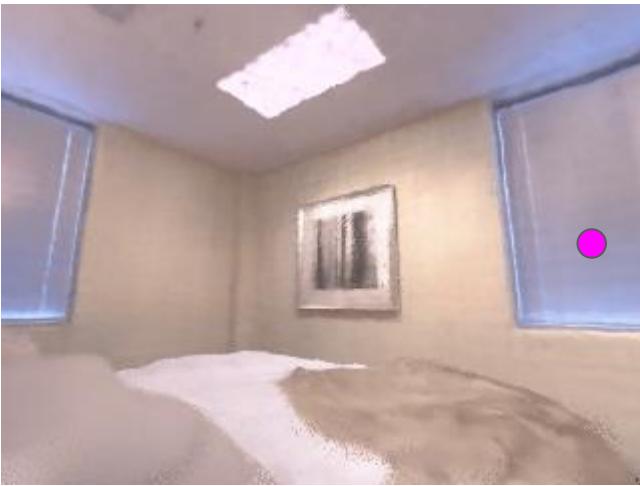
Ceiling re-colored
(yellow, blue)



Walls re-colored
(yellow, blue)



Windows re-colored
(yellow, blue)



Info Propagation Through Resonances in Views

2D version (JacobiNeRF-2D):

for each labeled pixel

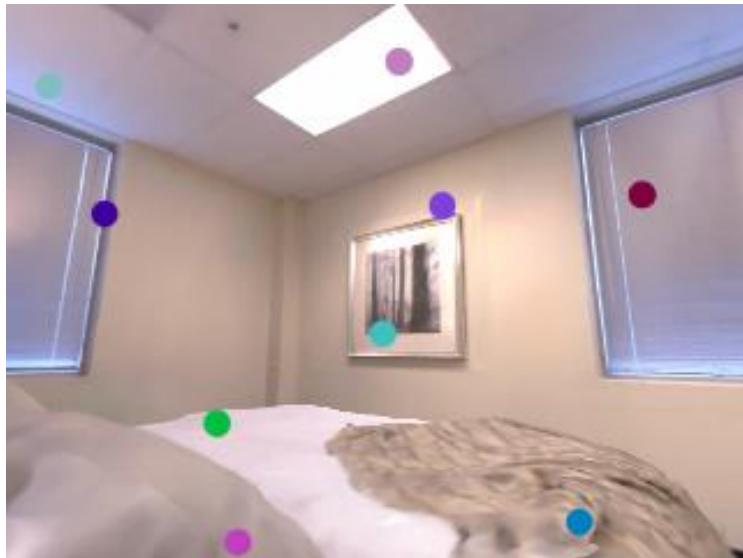
- perturb the NeRF along the gradient of the gray value of that pixel (e.g., change the network parameters)
- synthesize the target view from the perturbed NeRF
- calculate the perturbation response at each pixel for every source
- assign in target view pixels to the class generating the maximal response (argmax)

Make a move



See who follows

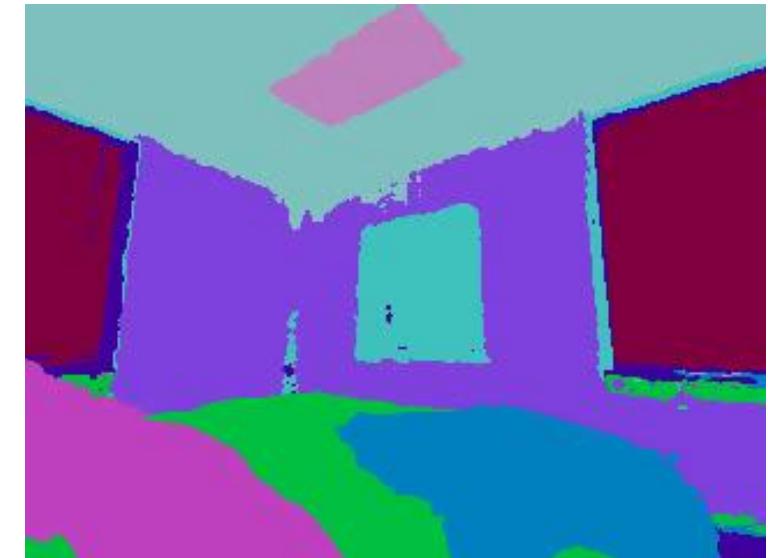
Semantic Segmentation (sparse, Replica)



Given label

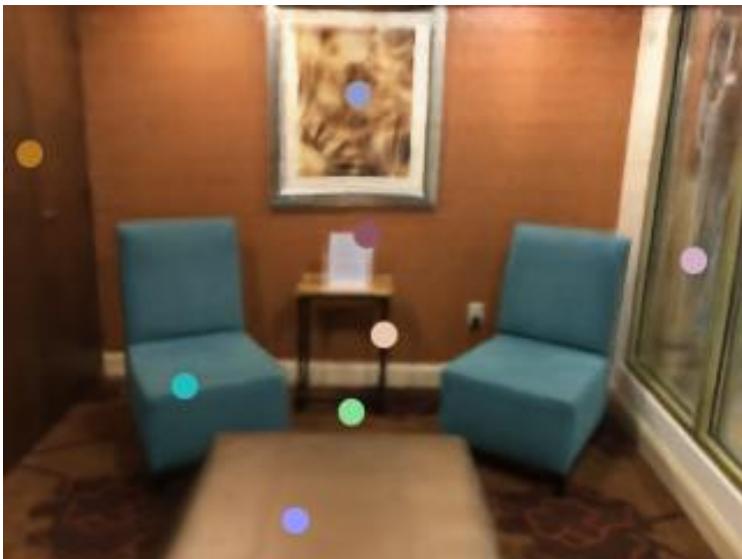


Regenerated view



J-NeRF 3D Semantics
Propagation

Semantic Segmentation (**sparse**, ScanNet)



Given label



J-NeRF 3D

Semantic Segmentation (dense, Replica)



Given label

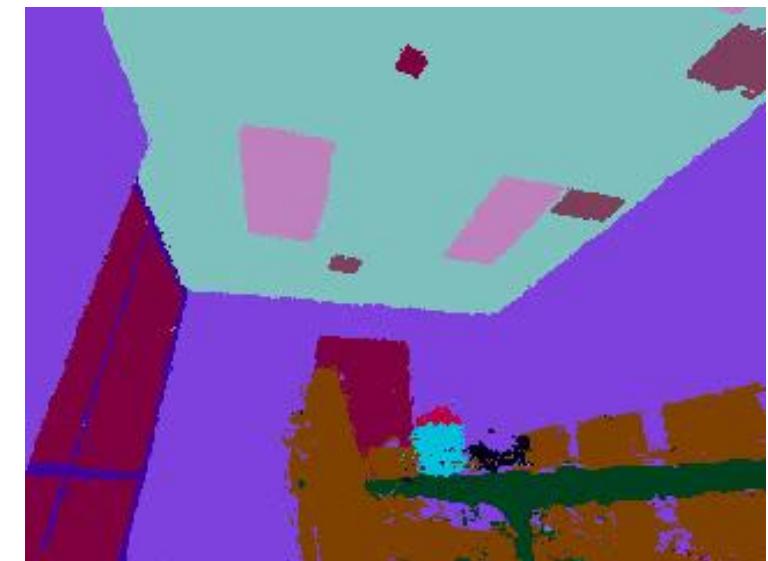


J-NeRF 3D

Semantic Segmentation (**dense**, Replica)



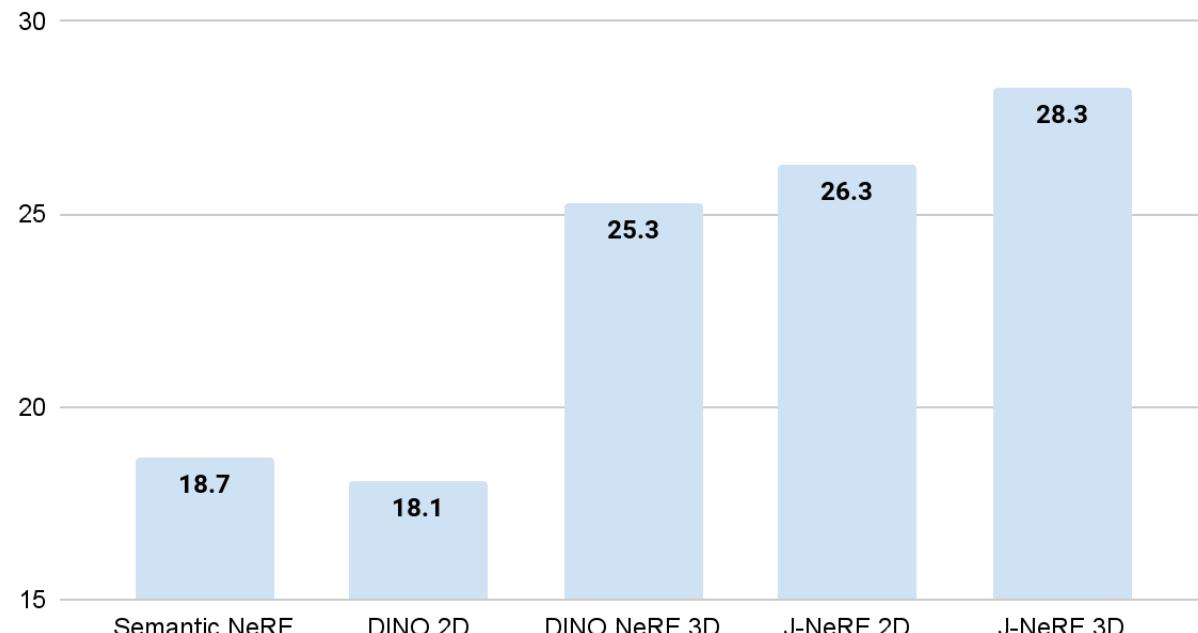
Given label



J-NeRF 3D

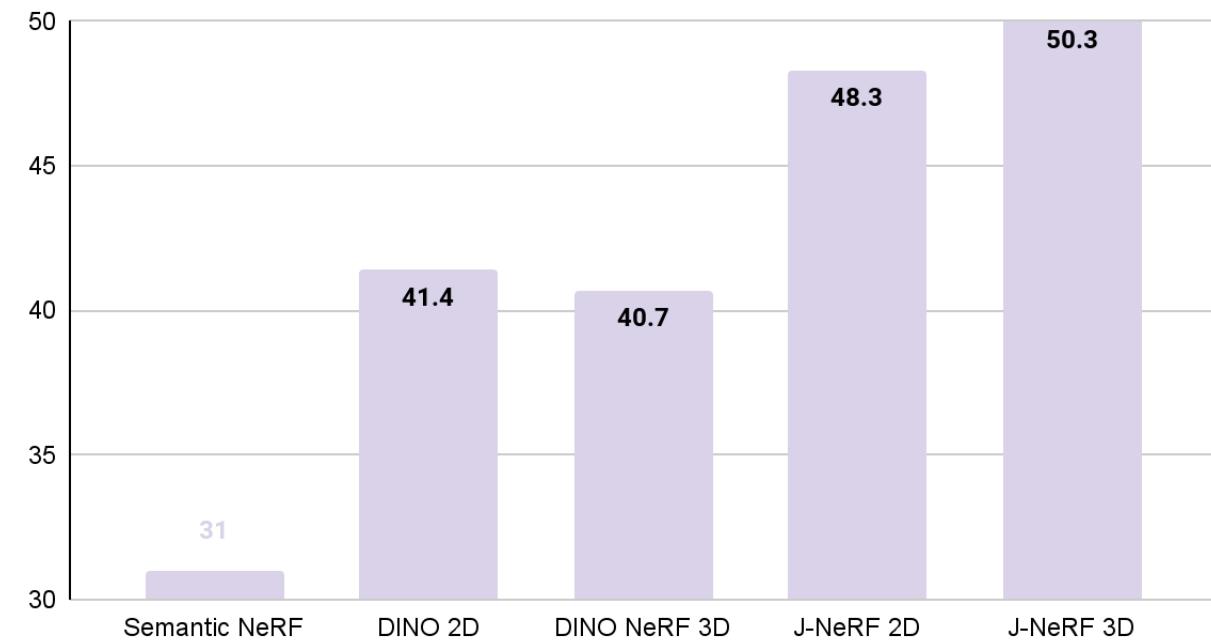
Semantic Segmentation (sparse 1pix/class, Replica)

1 pix/class 1 view



mIoU

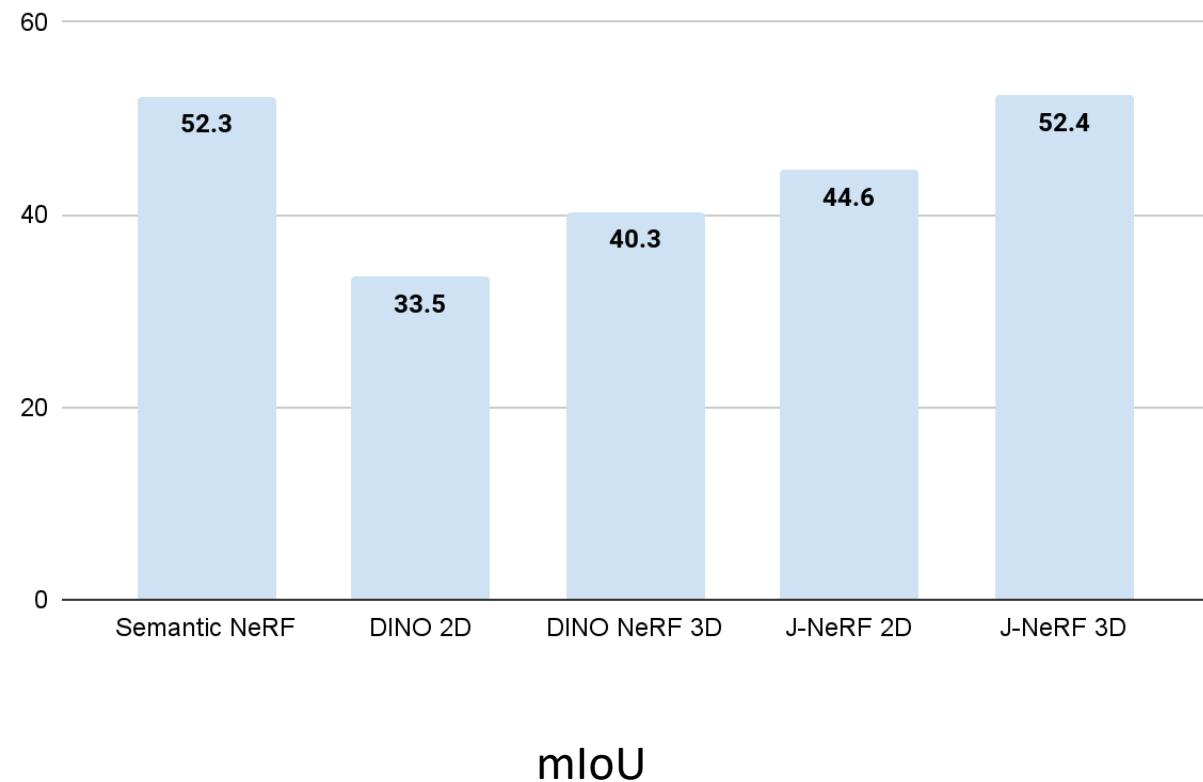
1 pix/class 1 view



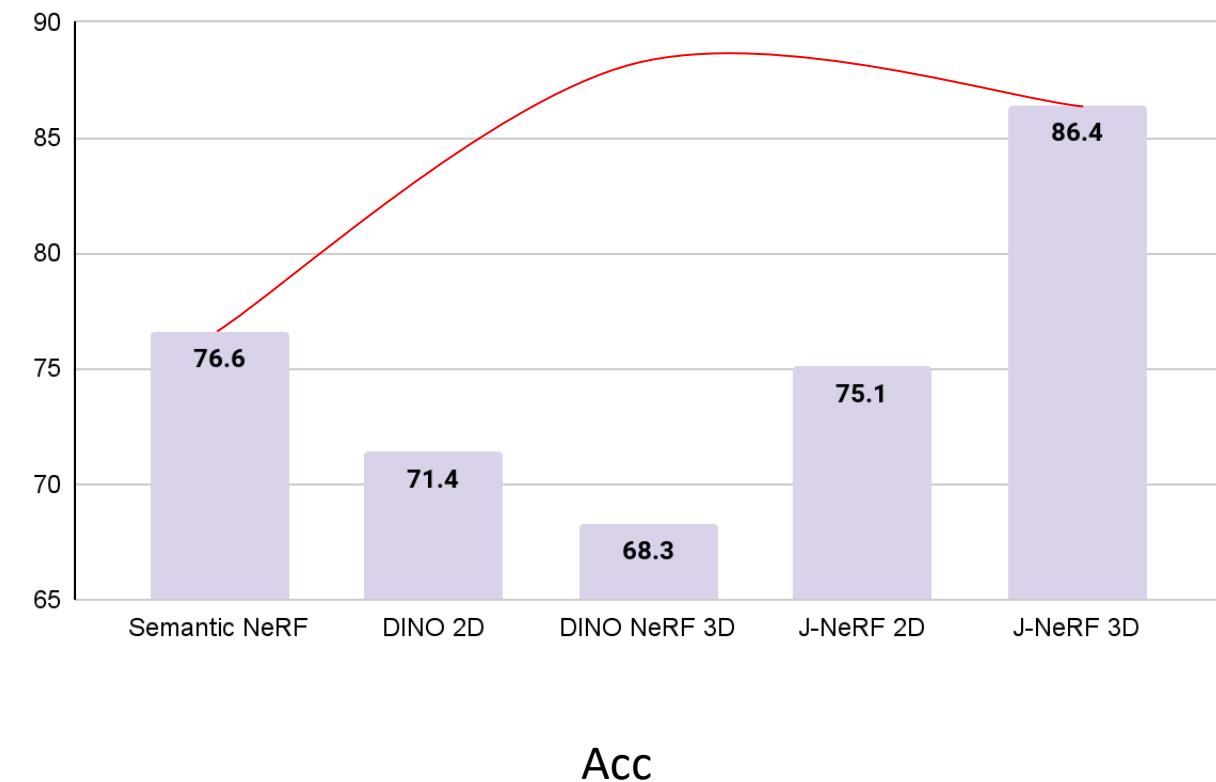
Acc

Semantic Segmentation (1 view, **dense** labels, Replica)

Dense label 1 view



Dense label 1 view



Light Supervision for Structure Emergence in Gaussian Fields

InfoGaussian: Structure-Aware Dynamic Gaussians through Lightweight Information Shaping. Yunchao Zhang, Guandao Yang, Leonidas Guibas, Yanchao Yang. ICRL 2025.

Static Scene Reconstruction with Gaussian Splatting

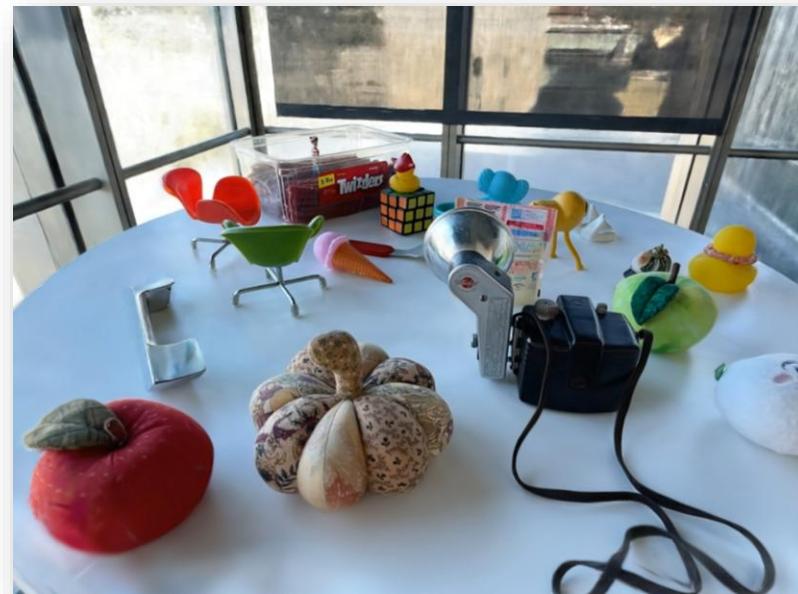


SAM “Segment Anything” 2D Dense Instance Supervision



Gaussian Grouping: Segment and Edit Anything in 3D Scenes. Mingqiao Ye, Martin Danelljan, Fisher Yu, Lei Ke.

GS Network Shaping for Object Motion



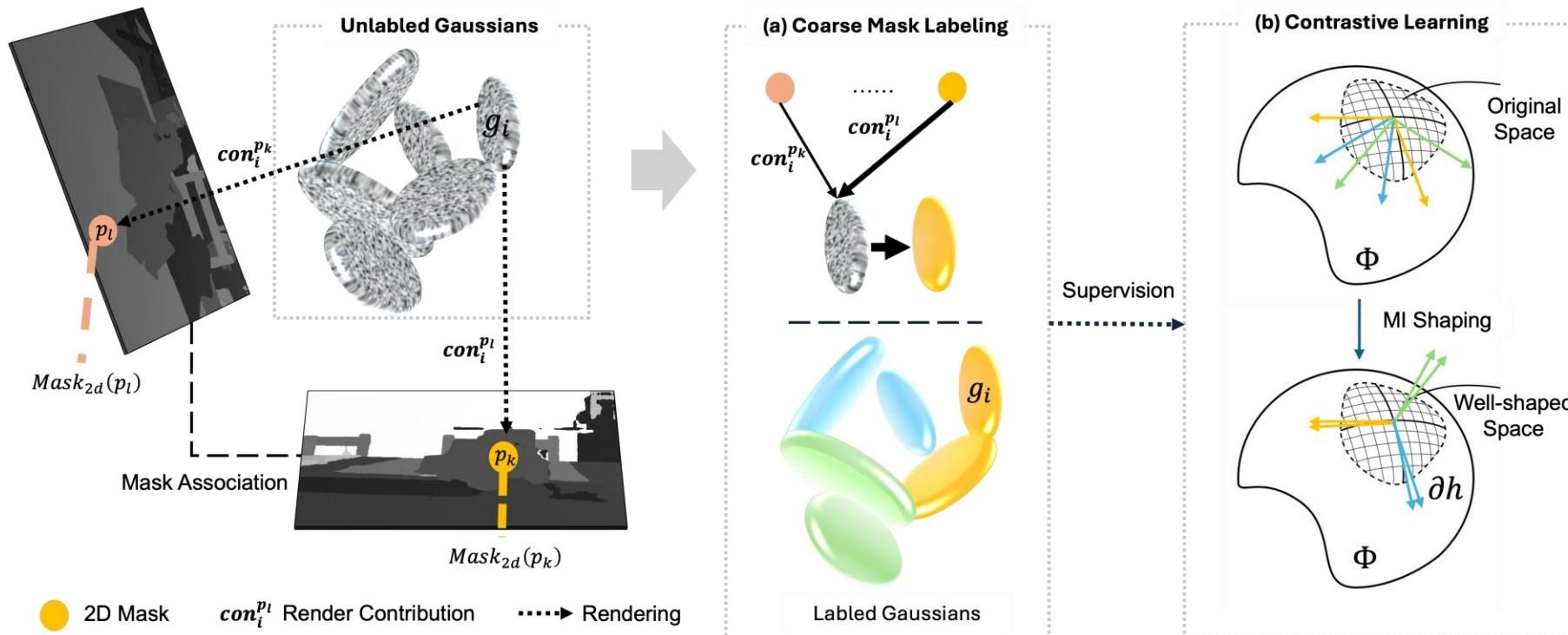
Pixels in the Same Mask Have High Mutual Information (MI)



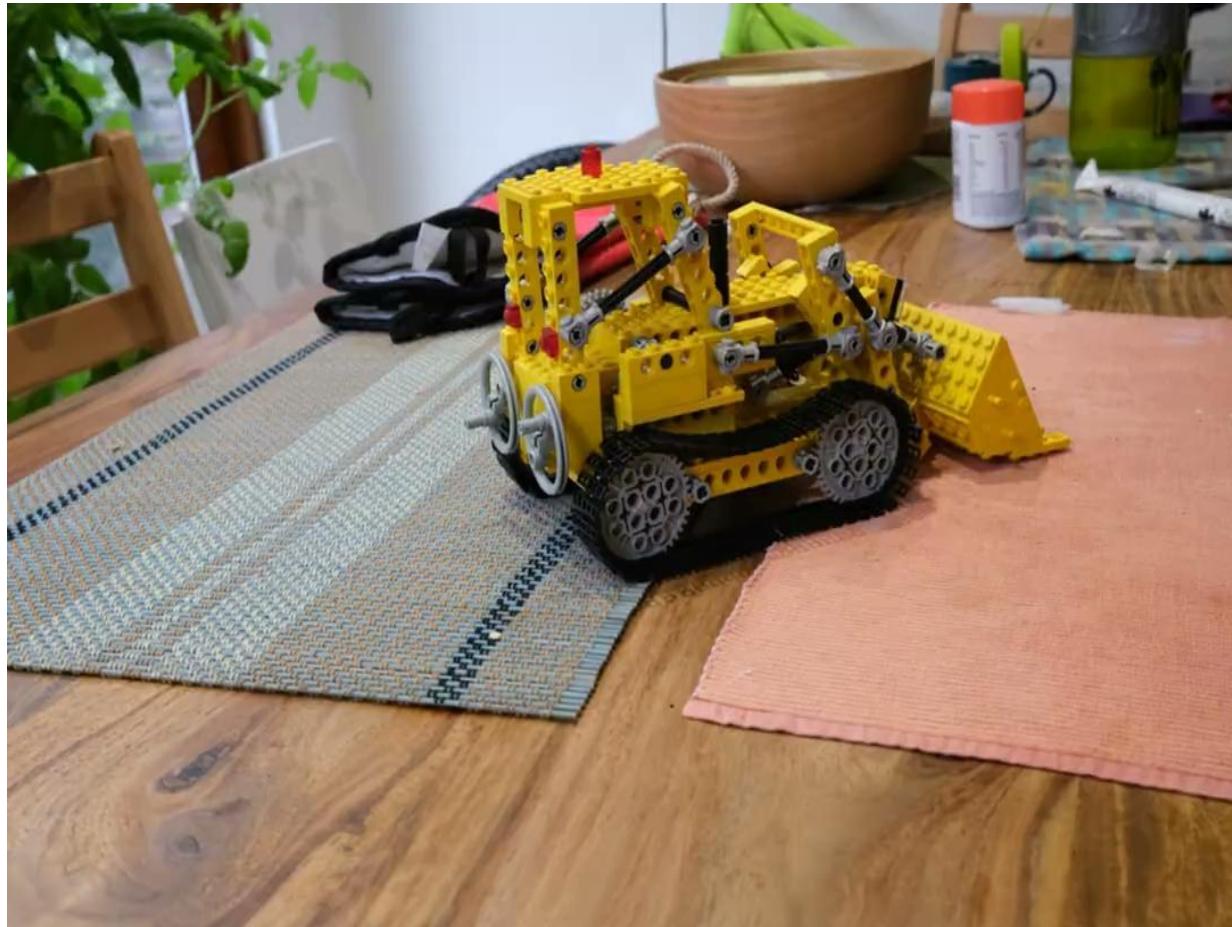
They are likely to change in correlated ways, as objects move.

Correlation Shaping on Attribute Decoding

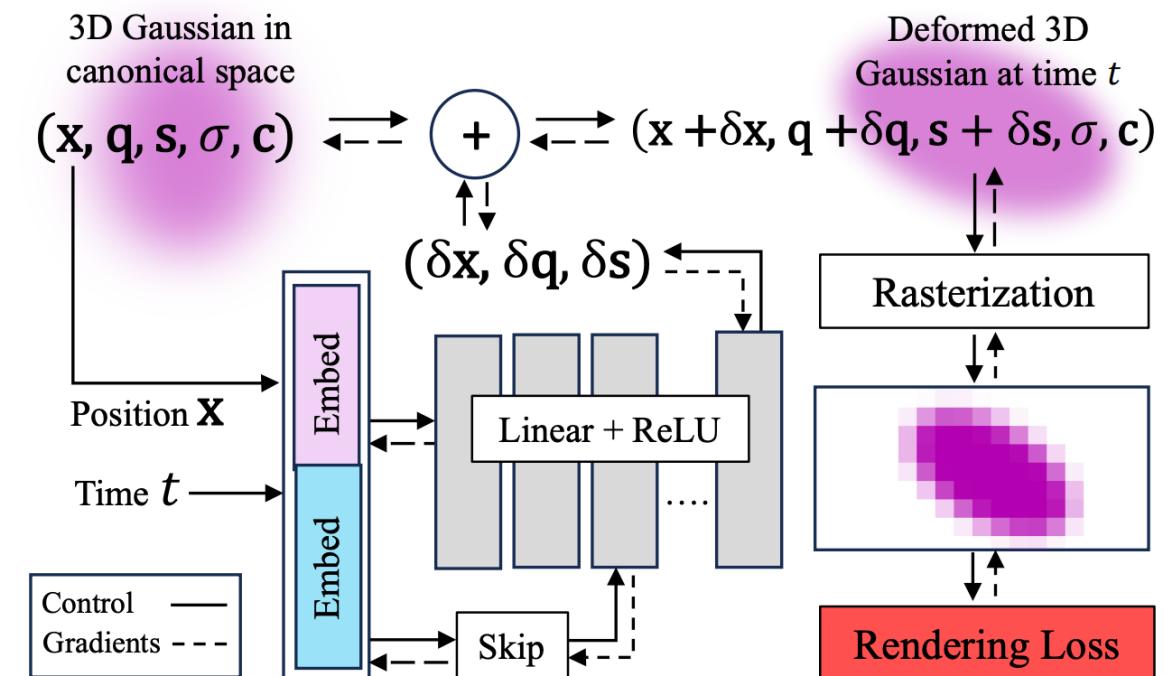
1. Use pretrained vision model (SAM) to generate 2D masks.
2. Label the 3D masks of Gaussians by 2D masks of the pixel.
3. Conduct contrastive learning for mutual information shaping.



Object Motion Without Explicit Grouping



MotionMLP



Motion MLP Gradient Alignments



Aligning MLP gradients via contrastive learning, to force high MI correlations

$$\text{III} \left(\hat{f}(g_i), \hat{f}(g_j) \right) = \log \left(\frac{1}{\sqrt{1 - \cos^2(\gamma)}} \right) + \text{const.}$$

$$\hat{f}(g_i) = \Phi(x_i; \theta_D + n)$$

$$\hat{f}(g_j) = \Phi(x_j; \theta_D + n)$$

$$\gamma = \cos \left(\frac{\partial \Phi(x_i; \theta)}{\partial \theta_D}, \frac{\partial \Phi(x_j; \theta)}{\partial \theta_D} \right)$$

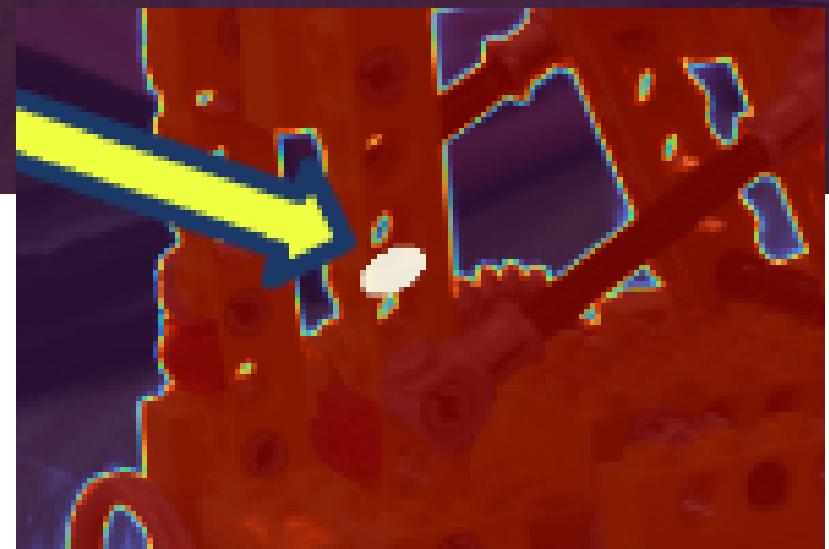
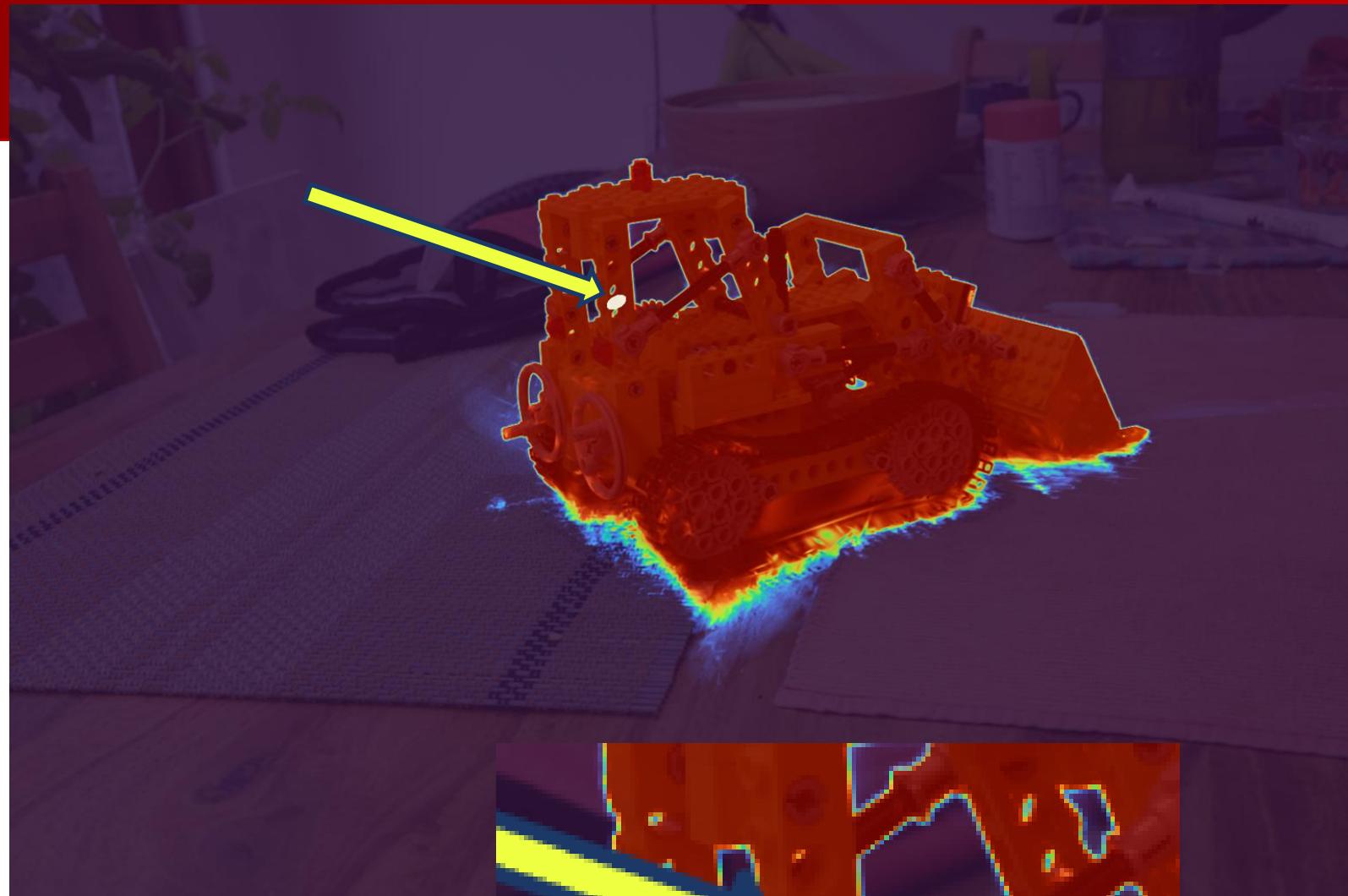
Experiment (Mips-NeRF 360 Scene)



Resonances

~7,900 Gaussians in the bulldozer

~460 participated in the contrastive training

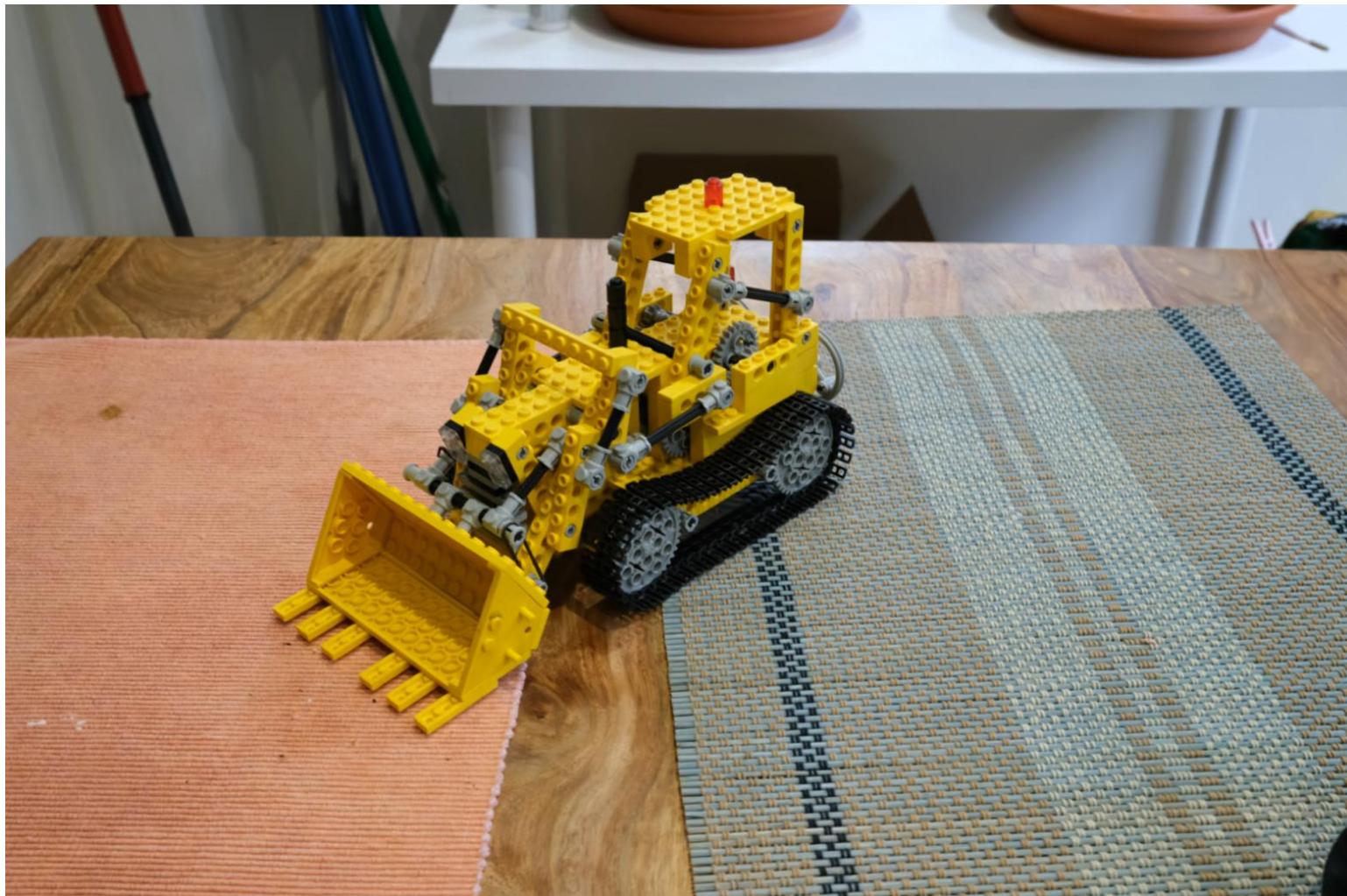


Experiment (Mips-NeRF 360 Scene)



Post-Shaping

Experiment (Mips-NeRF 360 Scene)

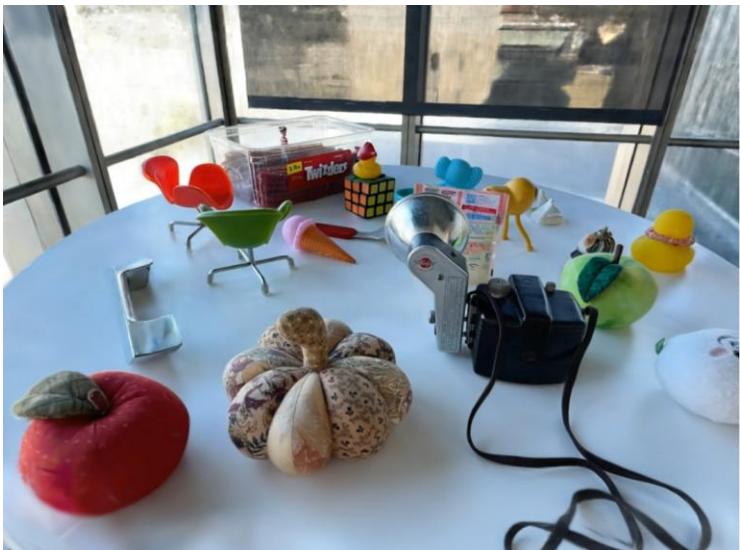


Post-Shaping

Experiment (LeRF Scene)



Resonances



Experiment (LeRF Scene)



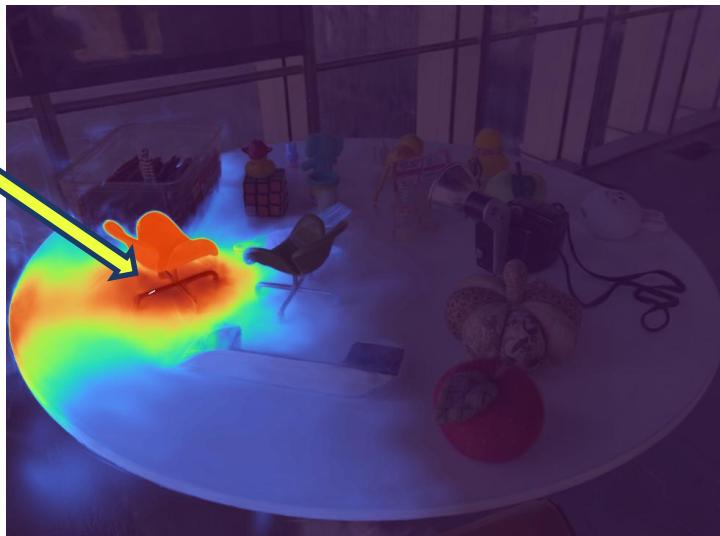
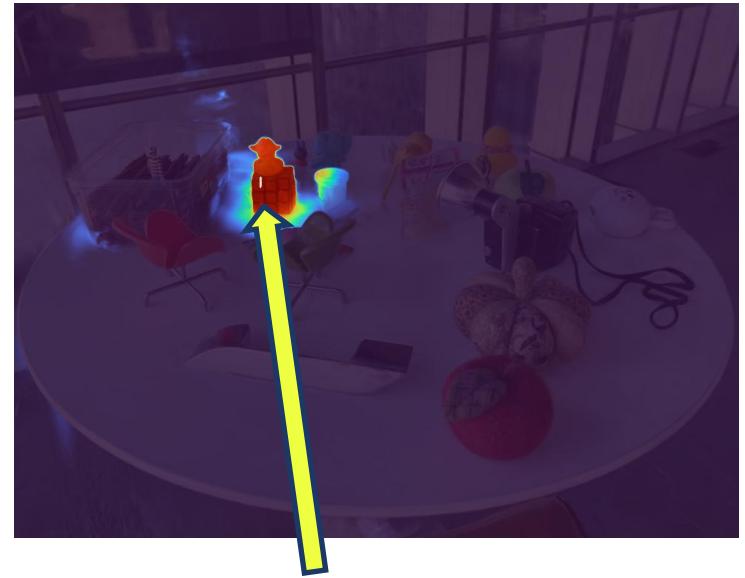
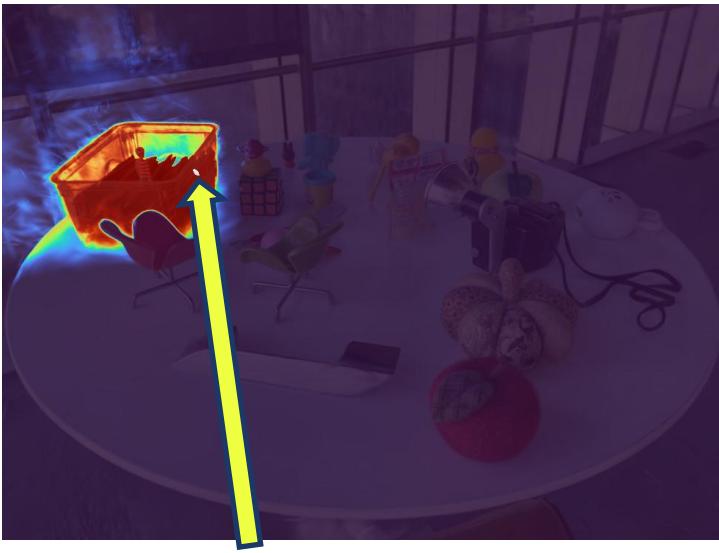
Post-Shaping

Experiment (LeRF Scene)



Post-Shaping

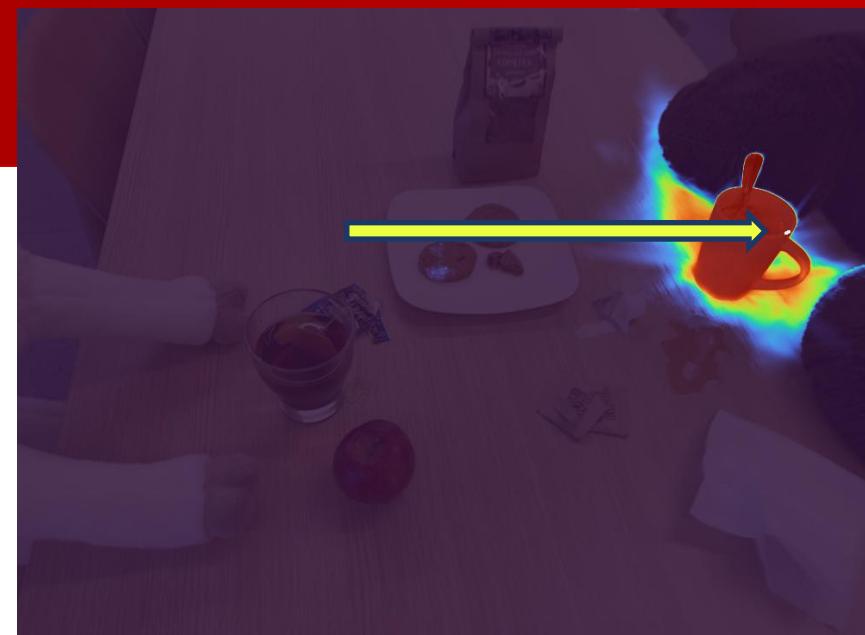
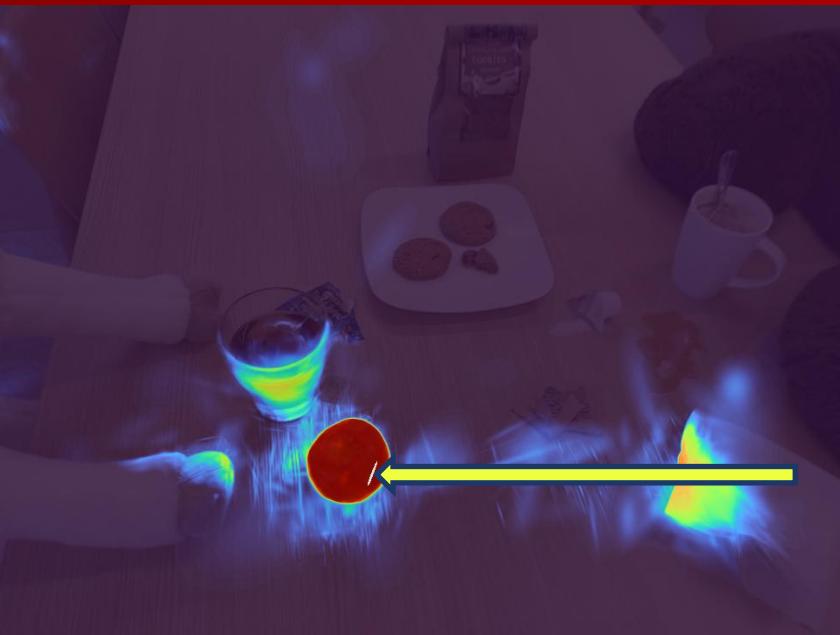
Resonances



Experiment (LeRF Scene)



Resonances



Experiment (LeRF Scene)



Post-Shaping

Summary: Enforcing a Prior (Equivariance)

- **Vector Neurons:**

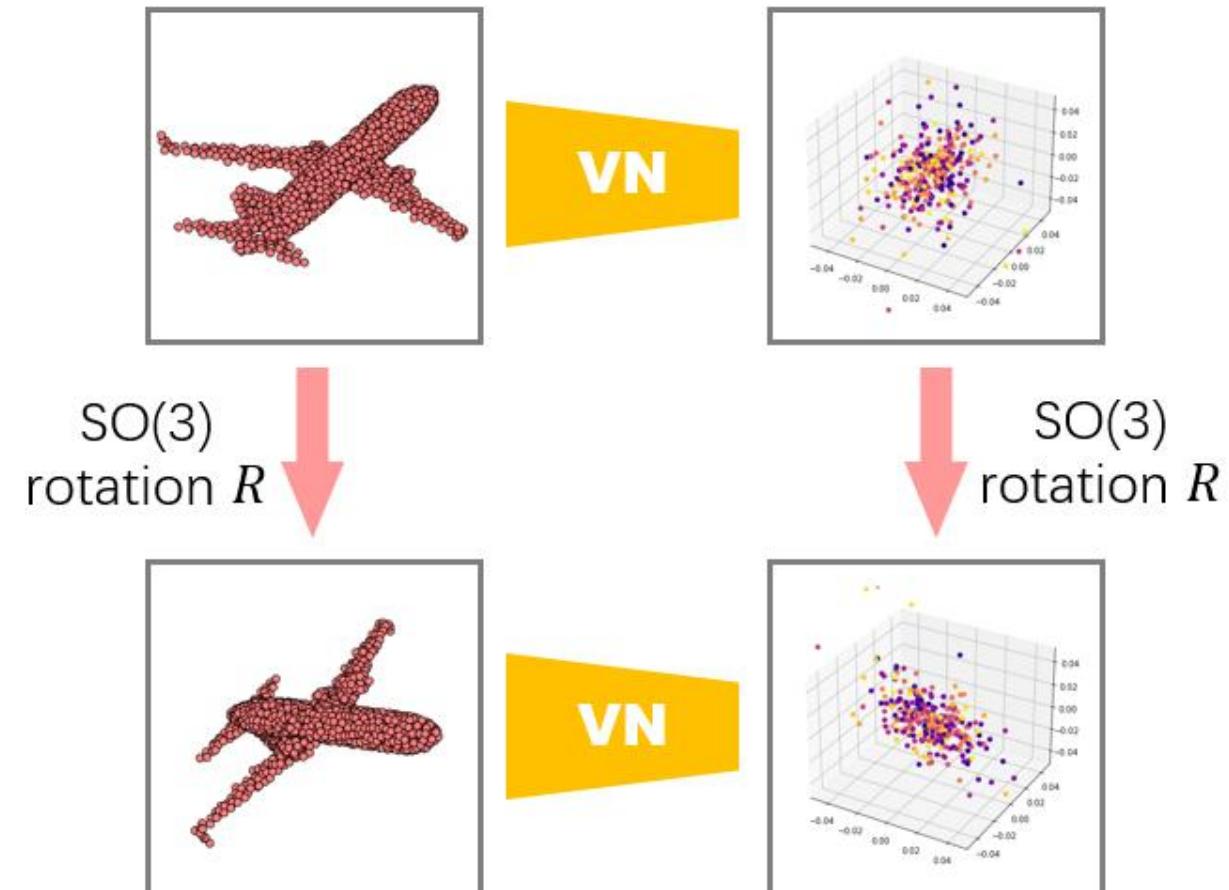
- Lift latent features to 3D vector lists

- **Building blocks:**

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

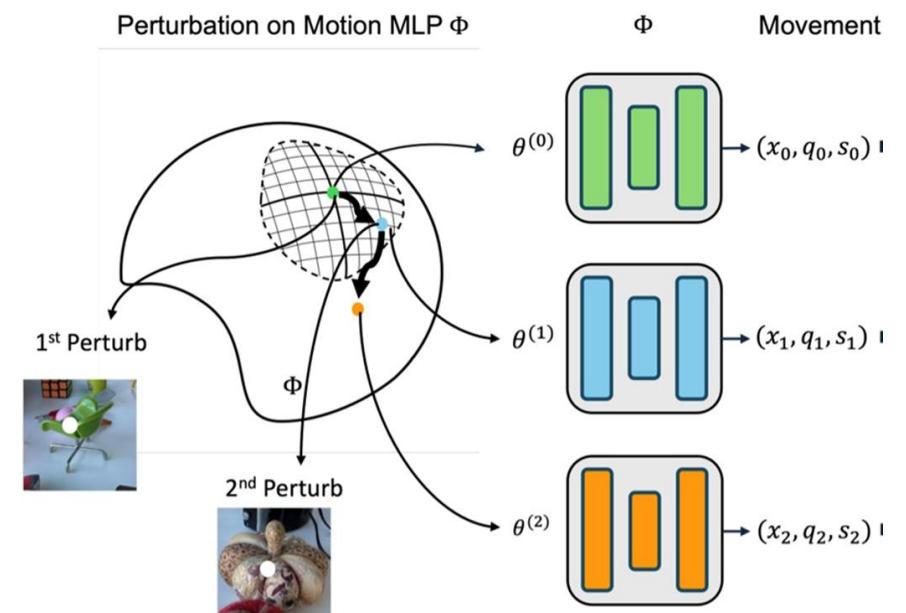
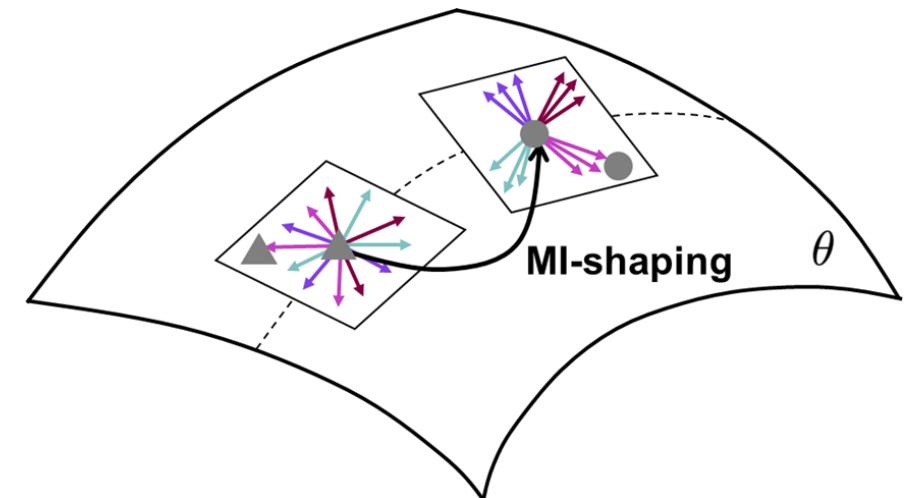
- **Network examples:**

- VN-DGCNN
- VN-PointNet



Summary: Learning a Prior (Grouping, Equivariance)

- **Scene Structure from Network Shaping:**
Gradient alignments according to mutual information
- **Semantic resonances:**
 - Learned from DINO features
 - Allow coherent edit propagation of
 - semantics
 - appearance
- **Motion resonances:**
 - Learned from SAM instance masks
 - Allow coherent entity motions



Thanks



That's All

