

**16.485: VNAV - Visual Navigation
for Autonomous Vehicles**

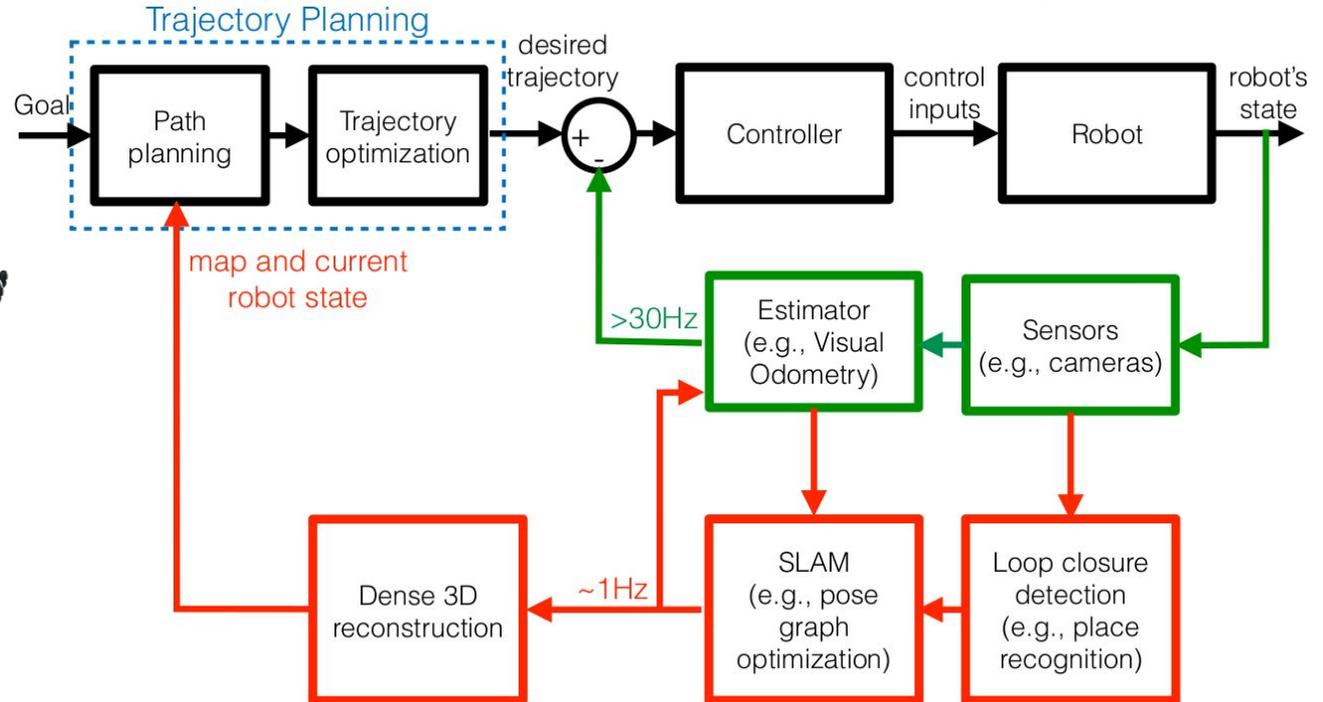
Rajat Talak

Lecture 32: Geometric Deep Learning

Recap: Motivation

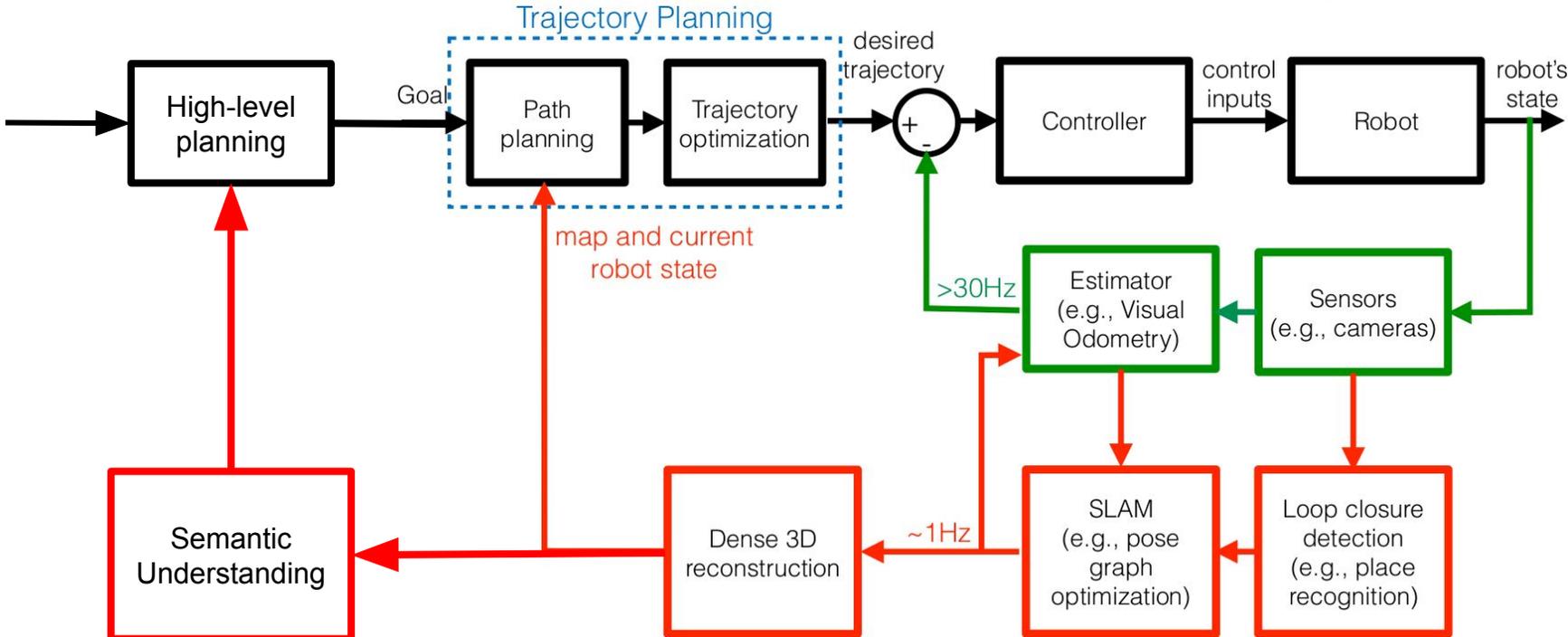
3D Geometric Reconstruction is not enough!

Need for Semantic Understanding



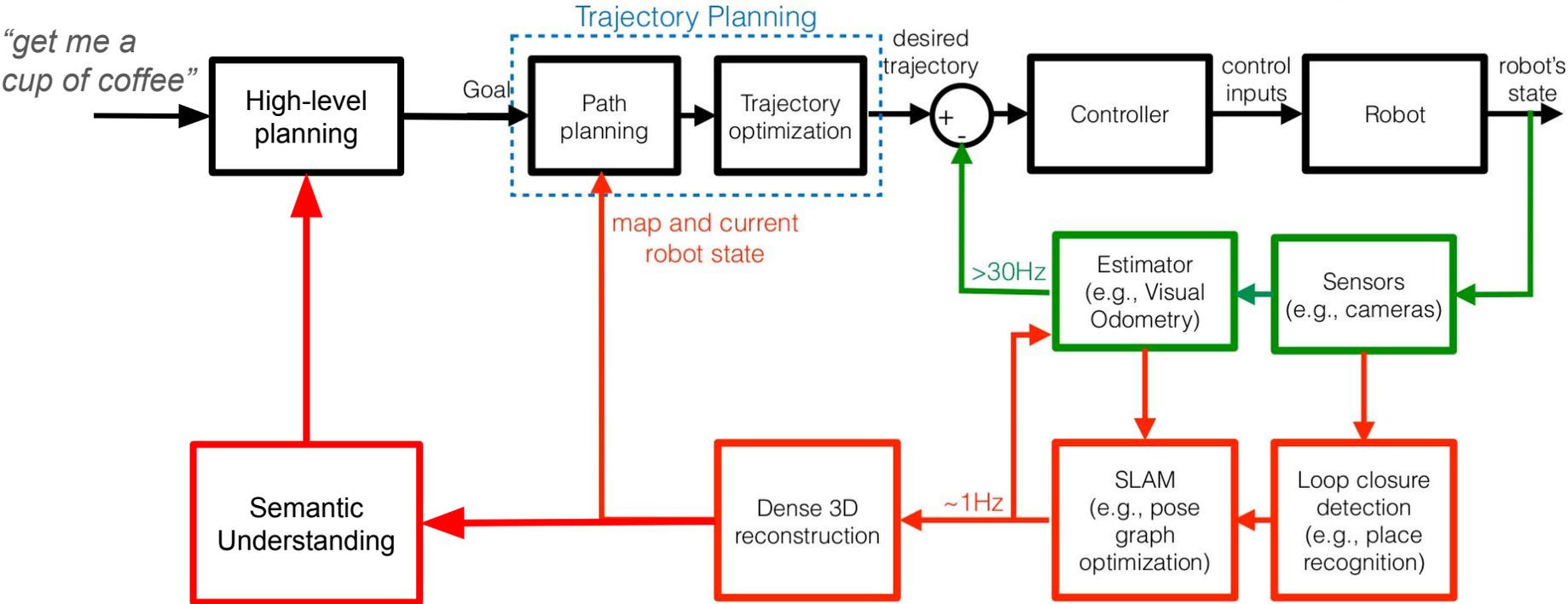
Recap: Motivation

3D Geometric Reconstruction is not enough!



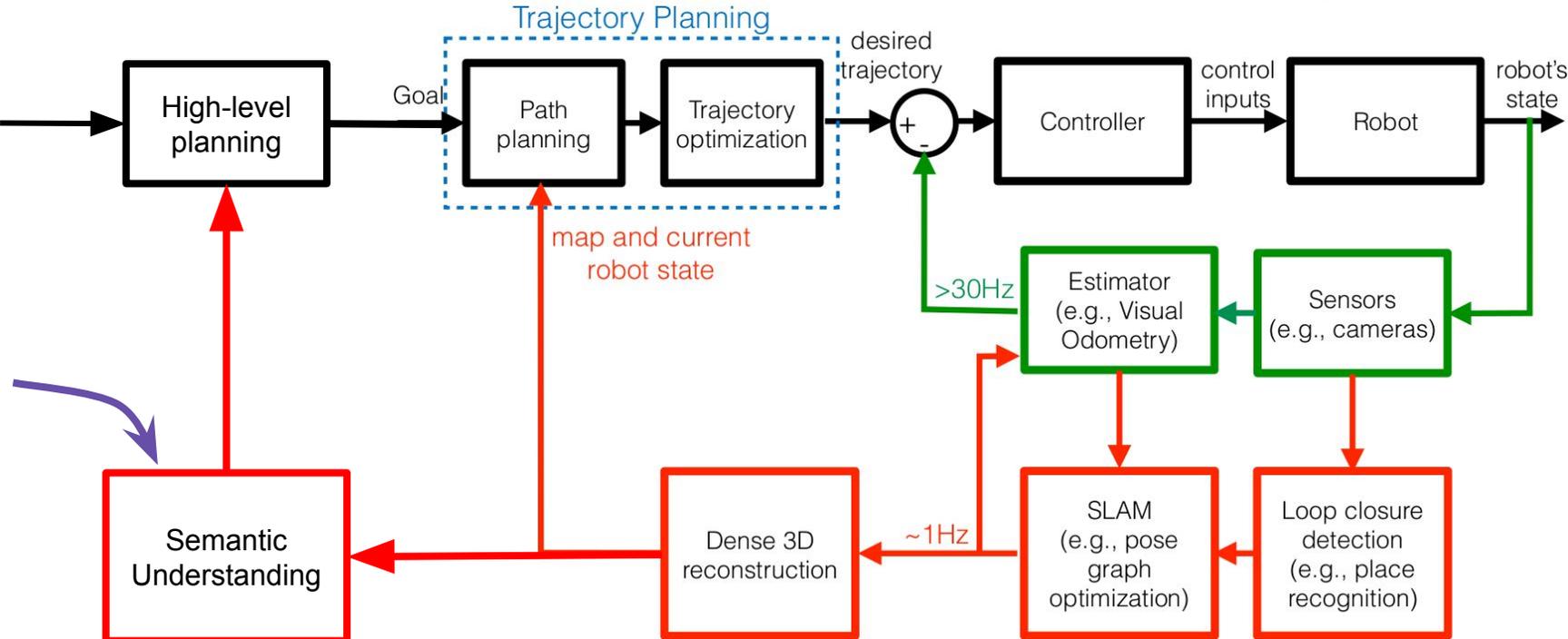
Recap: Motivation

3D Geometric Reconstruction is not enough!



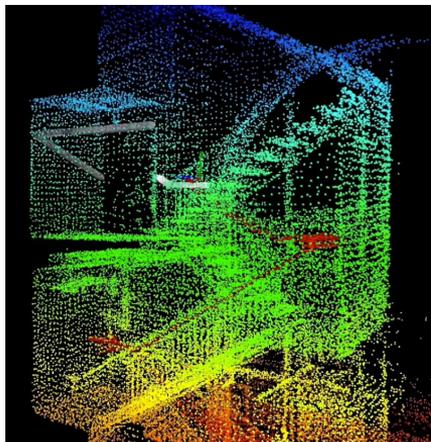
Recap: Motivation

3D Geometric Reconstruction is not enough!



Semantic Understanding

Need for Deep Learning Architectures on Point Clouds, Voxels, Meshes, Graphs

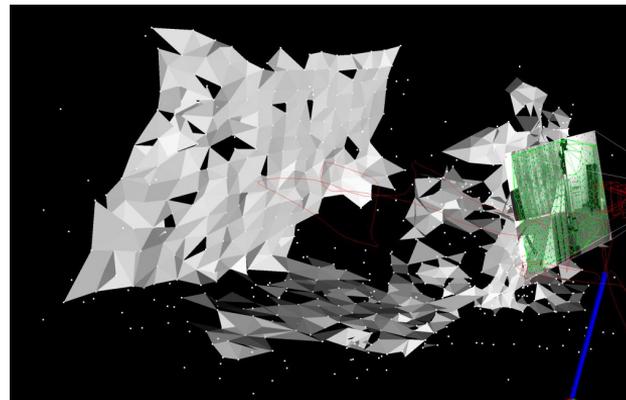


Point Cloud

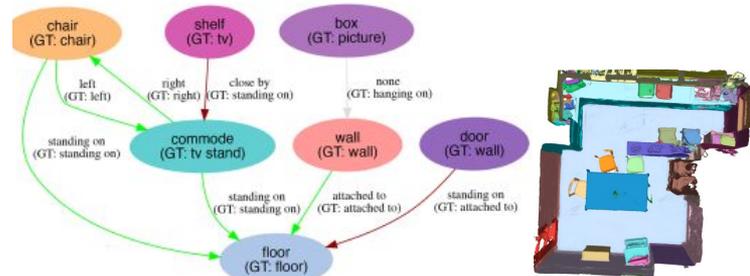


Voxel

Vespa et al. "Efficient Octree-based Volumetric SLAM Supporting Signed-Distance and Occupancy Mapping" RAL 2017



Mesh



Graph

Wald et al. "Learning 3D Semantic Scene Graphs from 3D Indoor Reconstruction" 2020

Recap: Architectures Discussed

Voxel

VoxNet

OctNet

Point Cloud

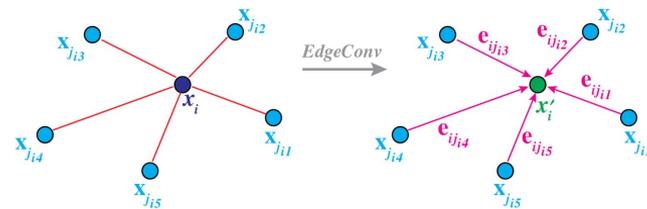
PointNet

PointNet++

EdgeConv

KPConv

Point Transformer



Today! For the first part ...

Voxel

VoxNet

OctNet

Point Cloud

PointNet

PointNet++

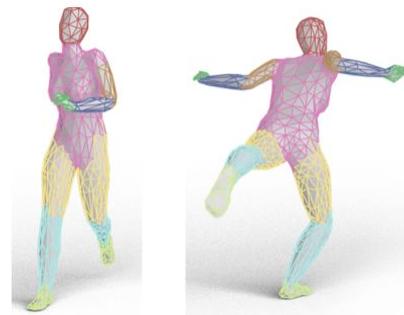
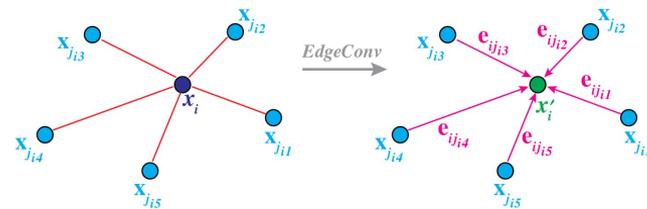
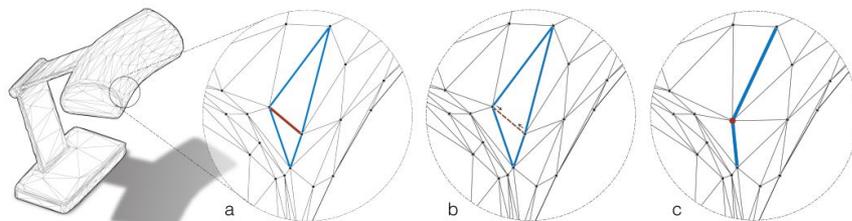
EdgeConv

KPCConv

Point Transformer

Mesh

MeshCNN



Point Cloud-based Architectures

Efficient than voxel based
architectures

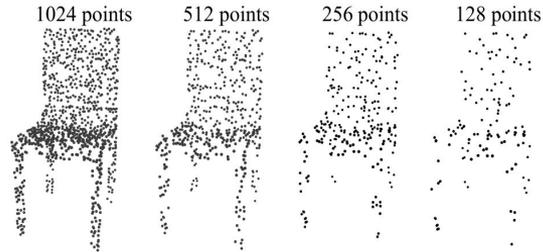
Suitable for point cloud inputs
(LiDAR or RGB-D)

Point Cloud-based Architectures

Efficient than voxel based architectures

Suitable for point cloud inputs (LiDAR or RGB-D)

Point clouds may not be the best way to represent 3D shapes

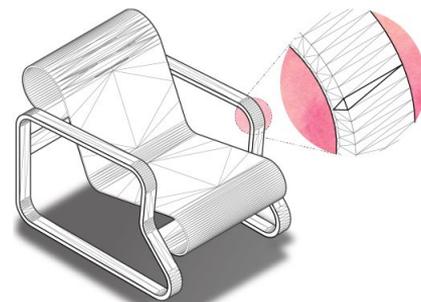
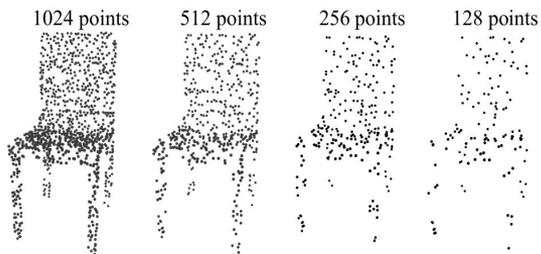


Point Cloud-based Architectures

Efficient than voxel based architectures

Suitable for point cloud inputs (LiDAR or RGB-D)

Point clouds may not be the best way to represent 3D shapes

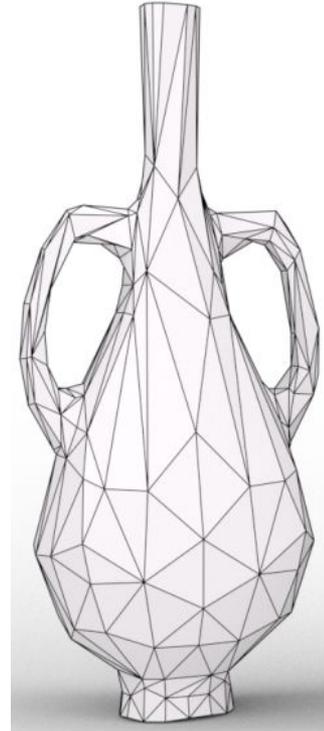


Mesh

Mesh

Mesh Representation

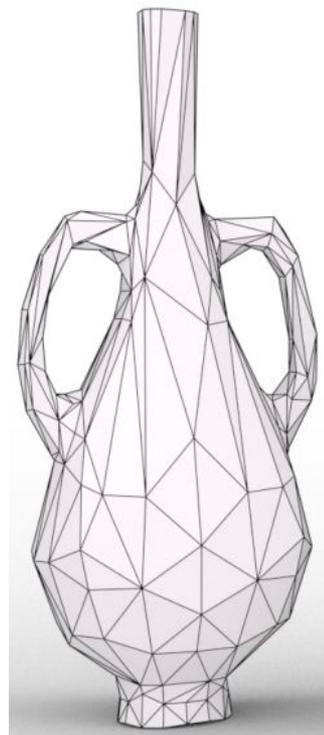
Mesh = Vertices, Faces, Edges



Mesh Representation

Mesh = Vertices, Faces, Edges

3d locations
 $v = (x, y, z)$

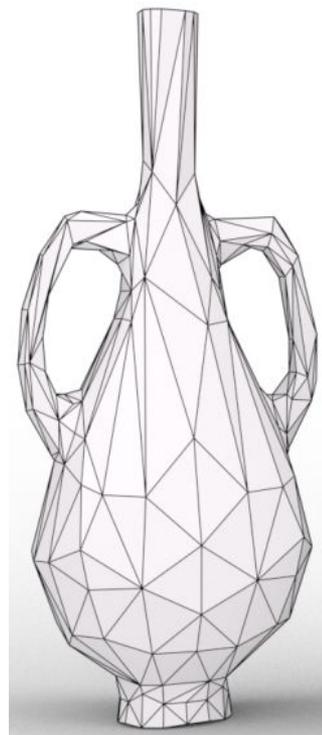


Mesh Representation

Mesh = Vertices, Faces, Edges

3d locations
 $v = (x, y, z)$

Triplet of vertices
 $f = (v_1, v_2, v_3)$



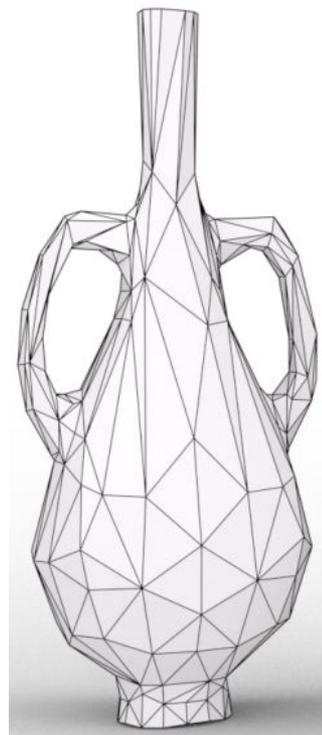
Mesh Representation

Mesh = Vertices, Faces, Edges

3d locations
 $v = (x, y, z)$

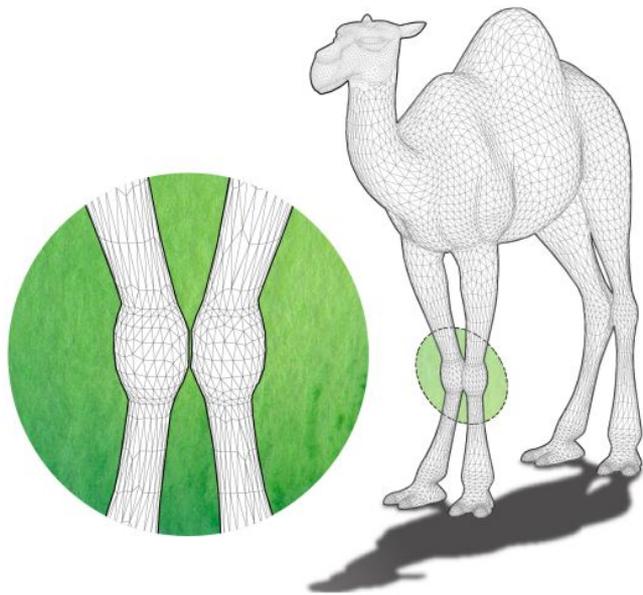
Triplet of vertices
 $f = (v_1, v_2, v_3)$

Pair of vertices
 $e = (v_1, v_2)$

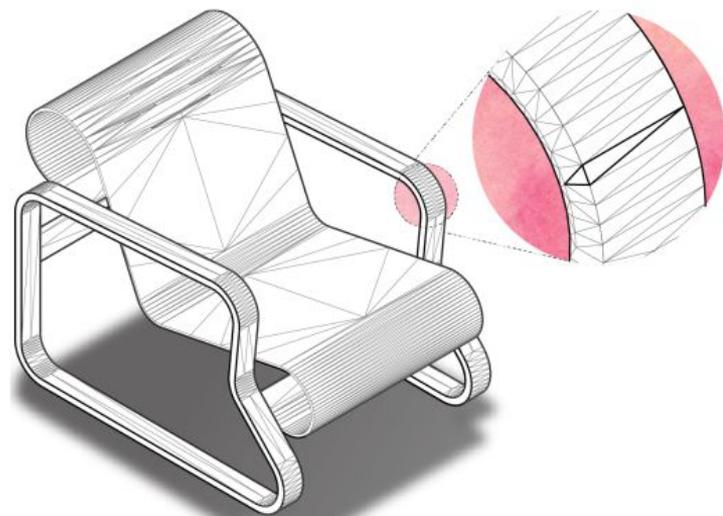


Mesh Connectivity and Local Shape

Conveys distinctness of local shape



Adaptive to non-uniform shape



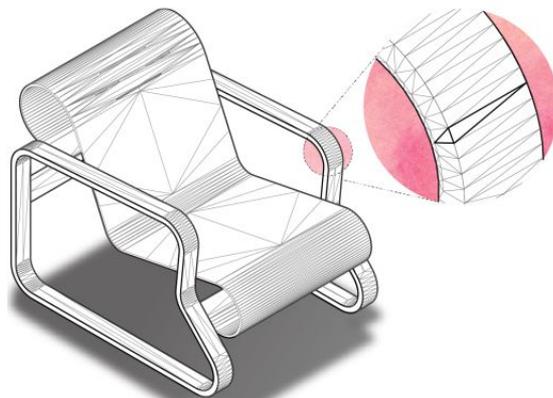
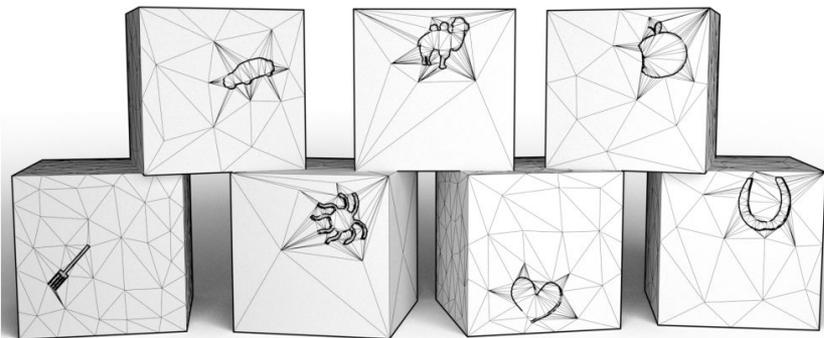
Learning on Meshes

Architectures should be able to exploit this property

Learning on Meshes

Architectures should be able to exploit this property

Problem: non-uniformity of the mesh

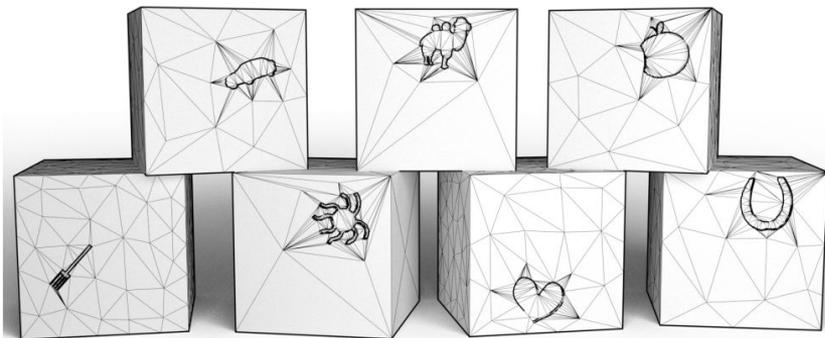


source: Hanocka et al. "MeshCNN: A Network with an Edge" ACM Trans. Graph. 2019

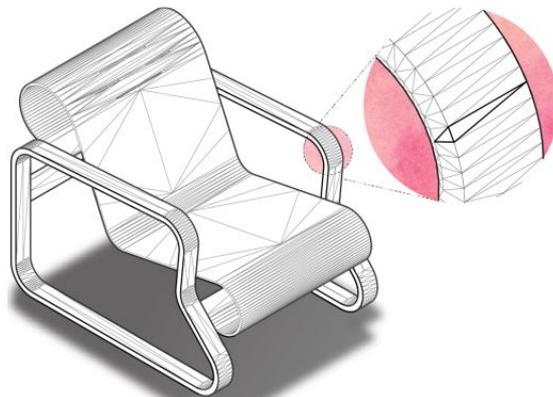
Learning on Meshes

Architectures should be able to exploit this property

Problem: non-uniformity of the mesh



Each vertex has varying number of neighbors

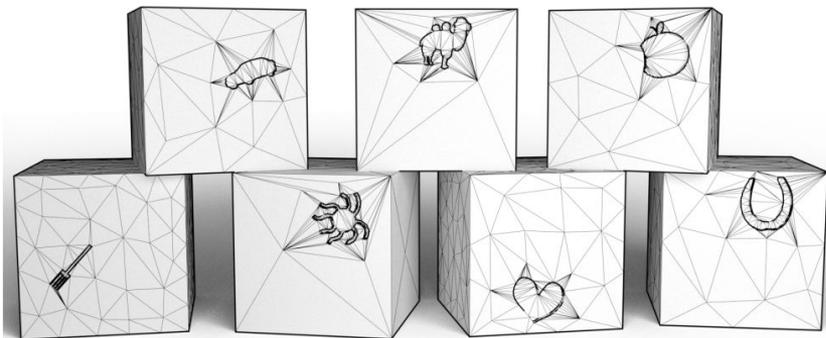


source: Hanocka et al. "MeshCNN: A Network with an Edge" ACM Trans. Graph. 2019

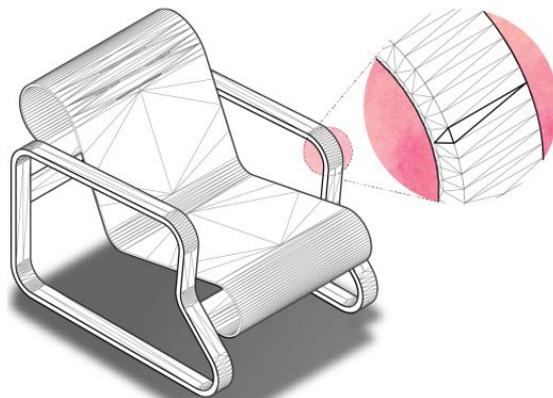
Learning on Meshes

How do we define convolution, pooling, and unpooling on this?

Problem: non-uniformity of the mesh



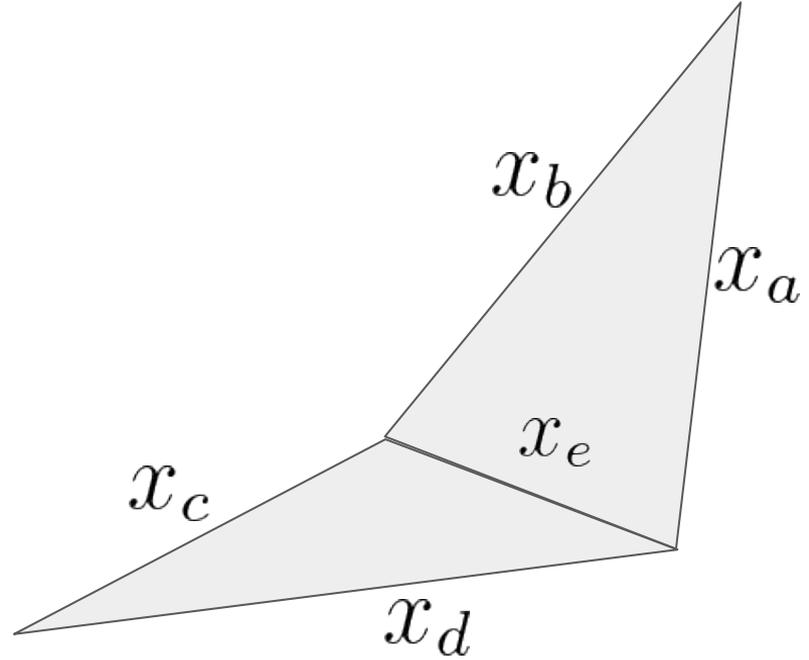
Each vertex has varying number of neighbors



source: Hanocka et al. "MeshCNN: A Network with an Edge" ACM Trans. Graph. 2019

MeshCNN

Every edge has two adjacent faces
and four adjacent edges

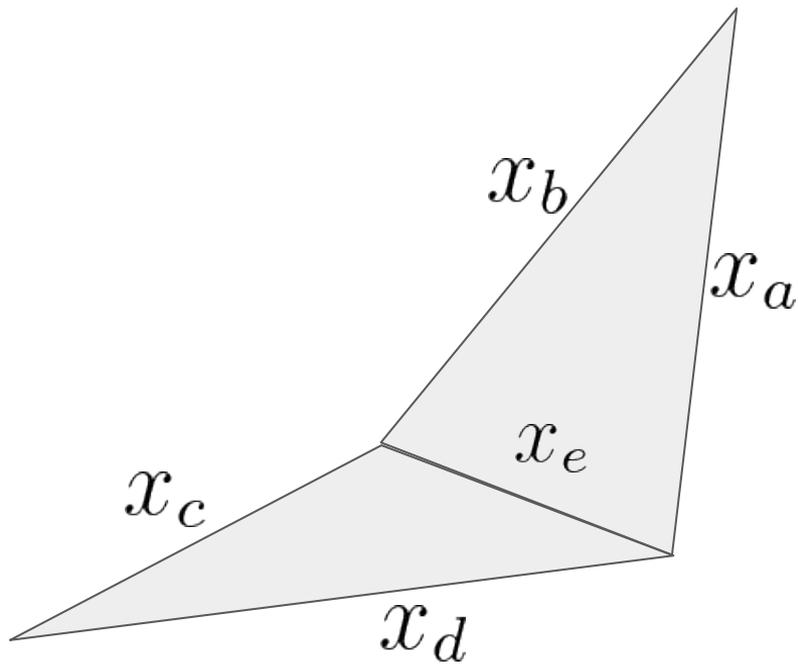


MeshCNN

Every edge has two adjacent faces
and four adjacent edges

Operates over mesh edges

Generates and updates
representation vectors over mesh
edges



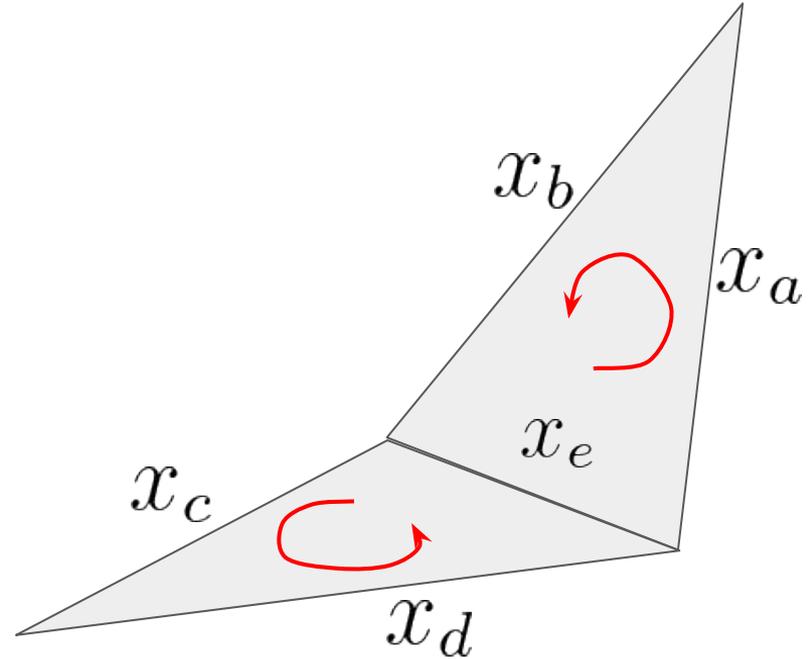
Ordering Neighboring Edges

Always order counter-clockwise

Two possibilities

(x_a, x_b, x_c, x_d)

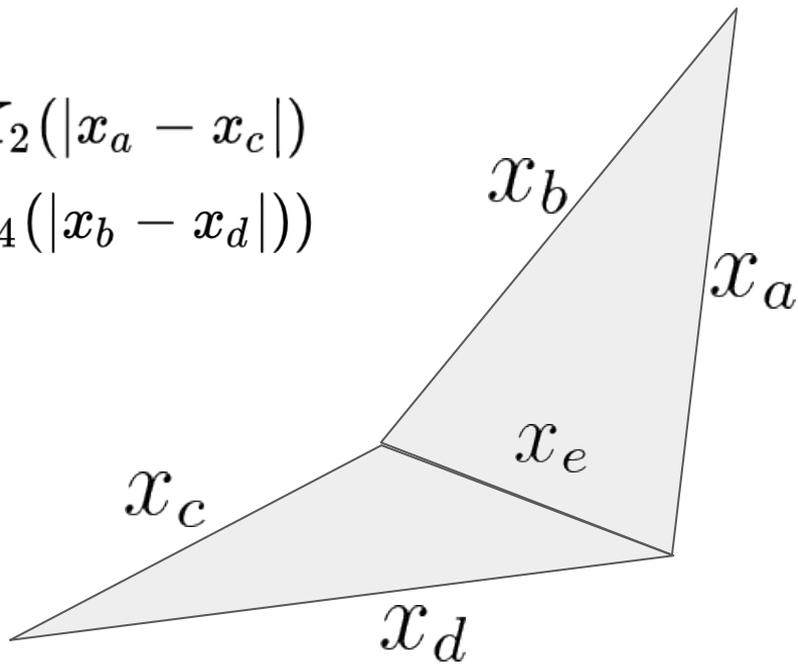
(x_c, x_d, x_a, x_b)



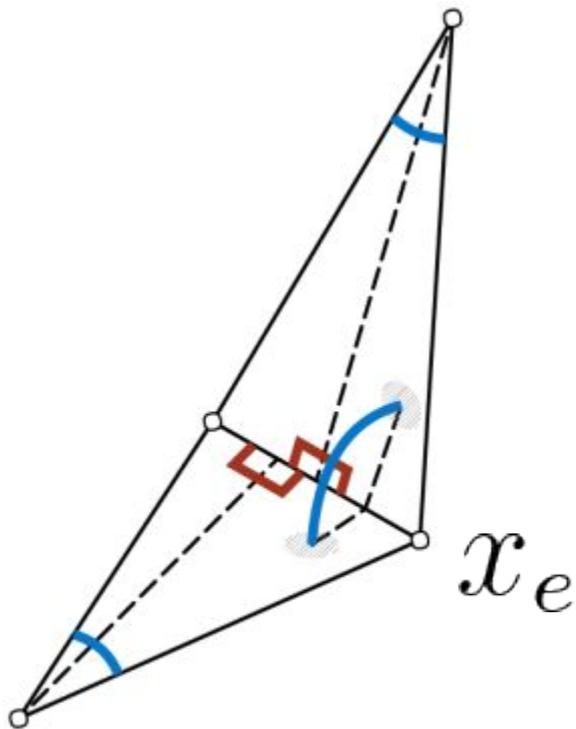
Aggregation should be invariant to these two possibilities

Updating Edge Features

$$\begin{aligned}x'_e = \sigma & (Kx_e + K_1(x_a + x_c) + K_2(|x_a - x_c|) \\ & + K_3(x_b + x_d) + K_4(|x_b - x_d|))\end{aligned}$$



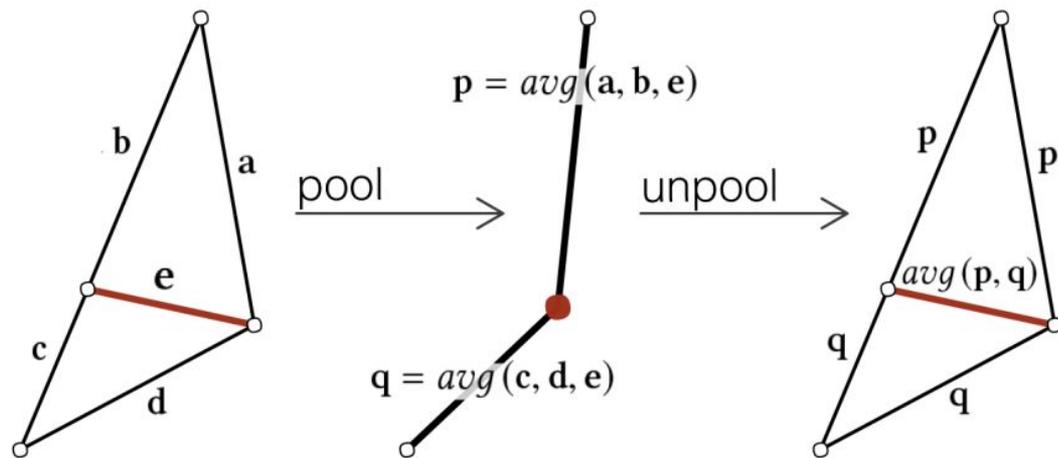
Initial features



Vector of length 5

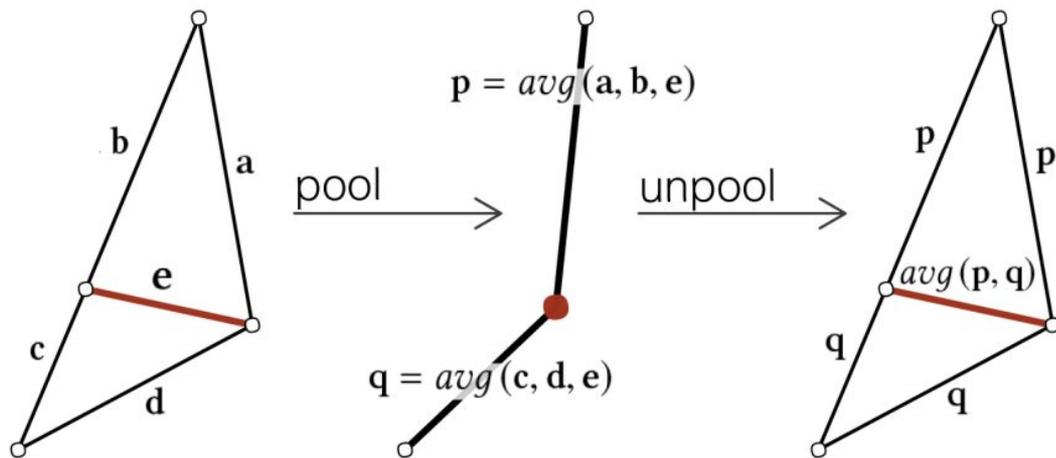
Invariant to translation, rotation,
and uniform scale

Pooling and Unpooling



In the figure a is \mathcal{X}_a ...

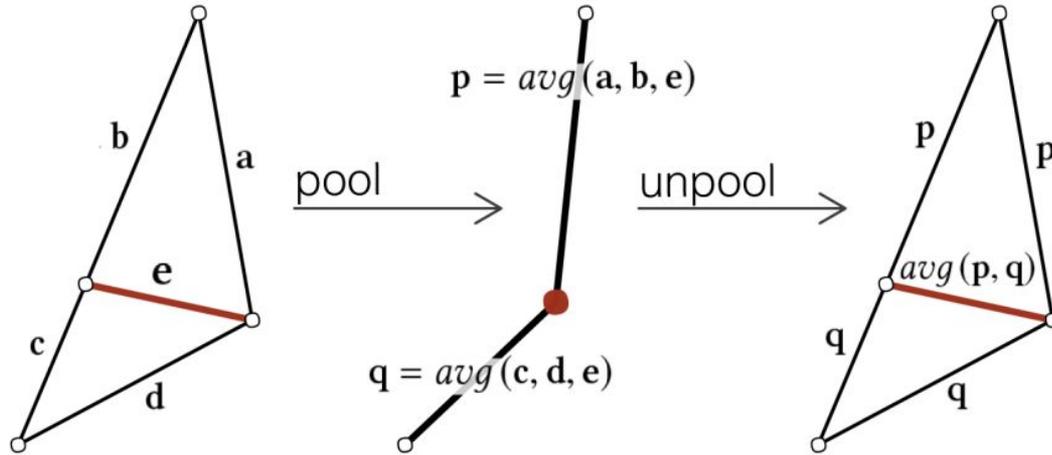
Pooling and Unpooling



edges with N smallest feature vector are collapsed at each layer

In the figure a is x_a ...

Pooling and Unpooling

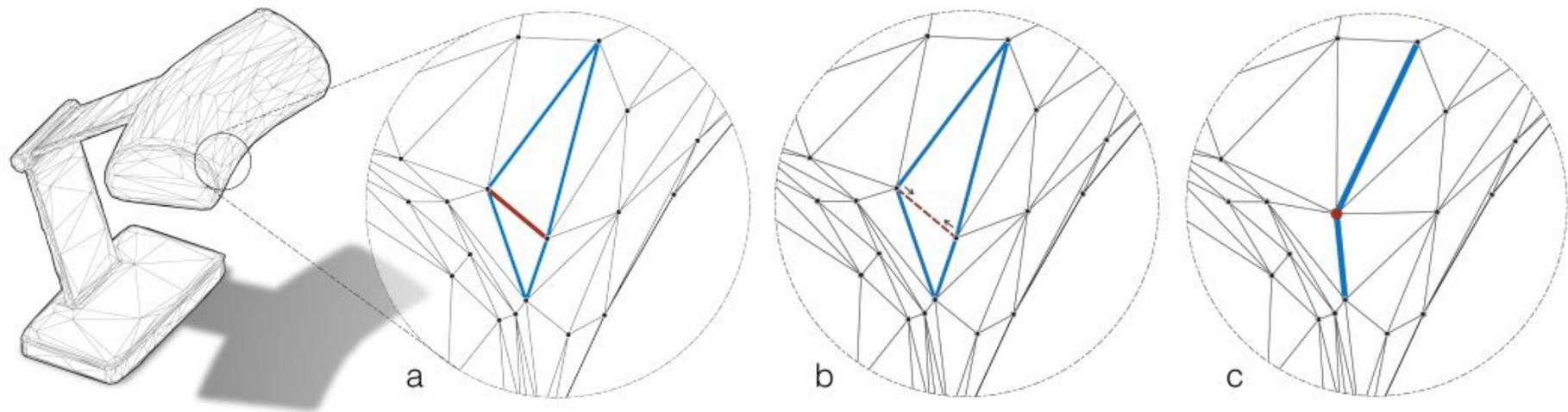


edges with N **smallest feature vector** are collapsed at each layer

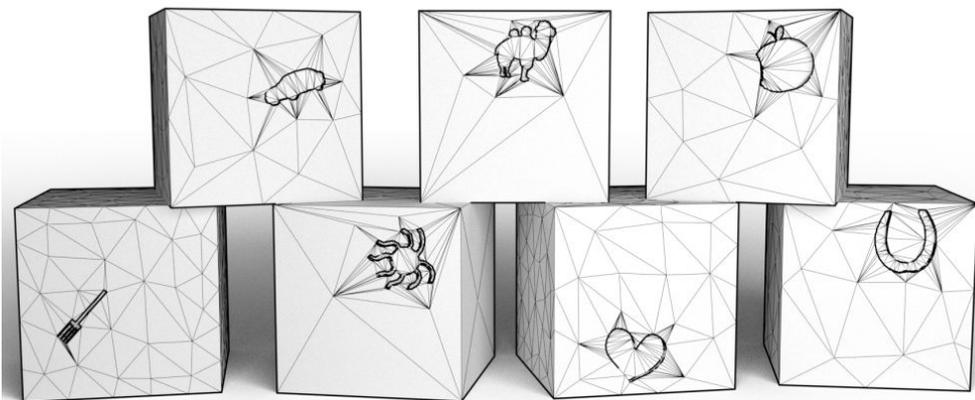
In the figure a is x_a ...

in L2 norm $\|e\|_2$

Pooling and Unpooling



MeshCNN: Interesting Results

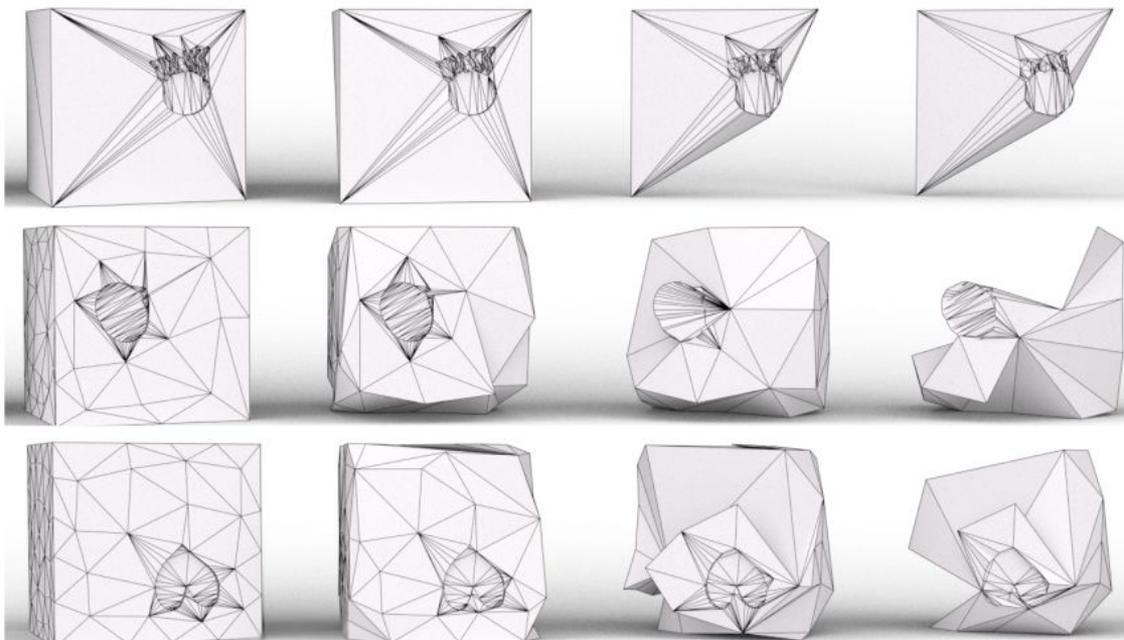


Classifying fine engraved cubes

Cube Engraving Classification

| method | input res | test acc |
|------------|-----------|---------------|
| MeshCNN | 750 | 92.16% |
| PointNet++ | 4096 | 64.26% |

MeshCNN: Interesting Results

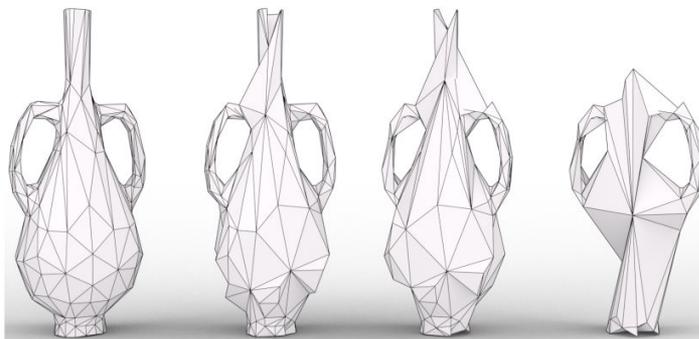


—————→
depth

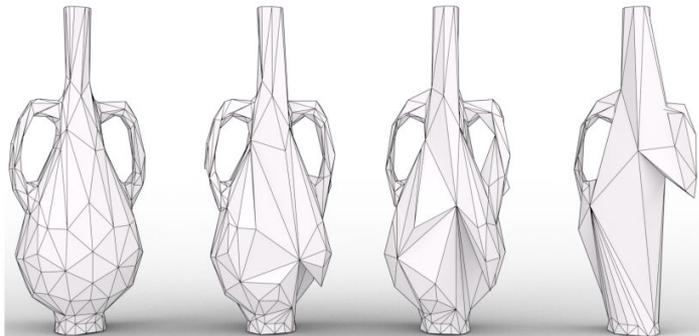
preserves important
edges required for
the task

source: Hanocka et al. "MeshCNN: A Network with an Edge" ACM Trans. Graph. 2019

MeshCNN: Interesting Results



Task 1: Vaze has a handle?



Task 2: Vaze has a neck?

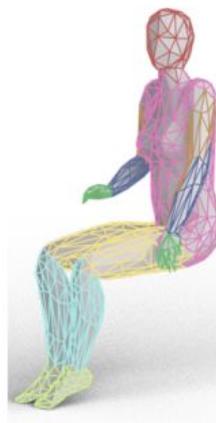
—————→
depth

source: Hanocka et al. "MeshCNN: A Network with an Edge" ACM Trans. Graph. 2019

MeshCNN: Human Shape Segmentation

| Human Body Segmentation | | |
|-------------------------|------------|---------------|
| Method | # Features | Accuracy |
| MeshCNN | 5 | 92.30% |
| SNGC | 3 | 91.02% |
| Toric Cover | 26 | 88.00% |
| PointNet++ | 3 | 90.77% |
| DynGraphCNN | 3 | 89.72% |
| GCNN | 64 | 86.40% |
| MDGCNN | 64 | 89.47% |

} [2018]



Mesh based Architectures

More structure. Opportunity for the architecture to be more expressive.

Computationally expensive than Point Cloud based architectures.

- Pooling, unpooling, manifoldness



Mesh based Architectures

Different meshes can represent the same thing

Data Augmentation



Choice of Representation

Expressive Power

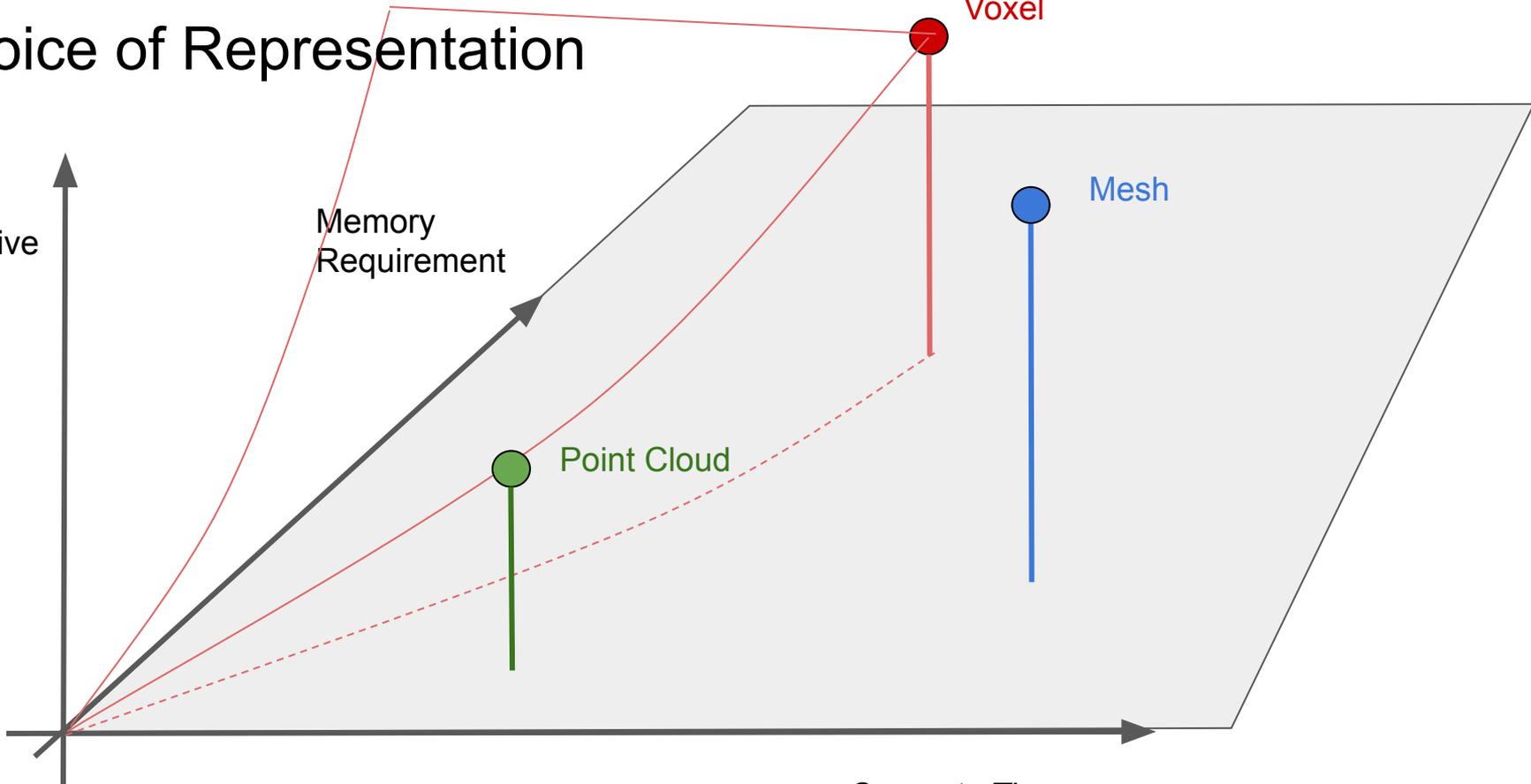
Memory Requirement

Point Cloud

Voxel

Mesh

Compute Time



Choice of Representation

Expressive Power

Memory Requirement

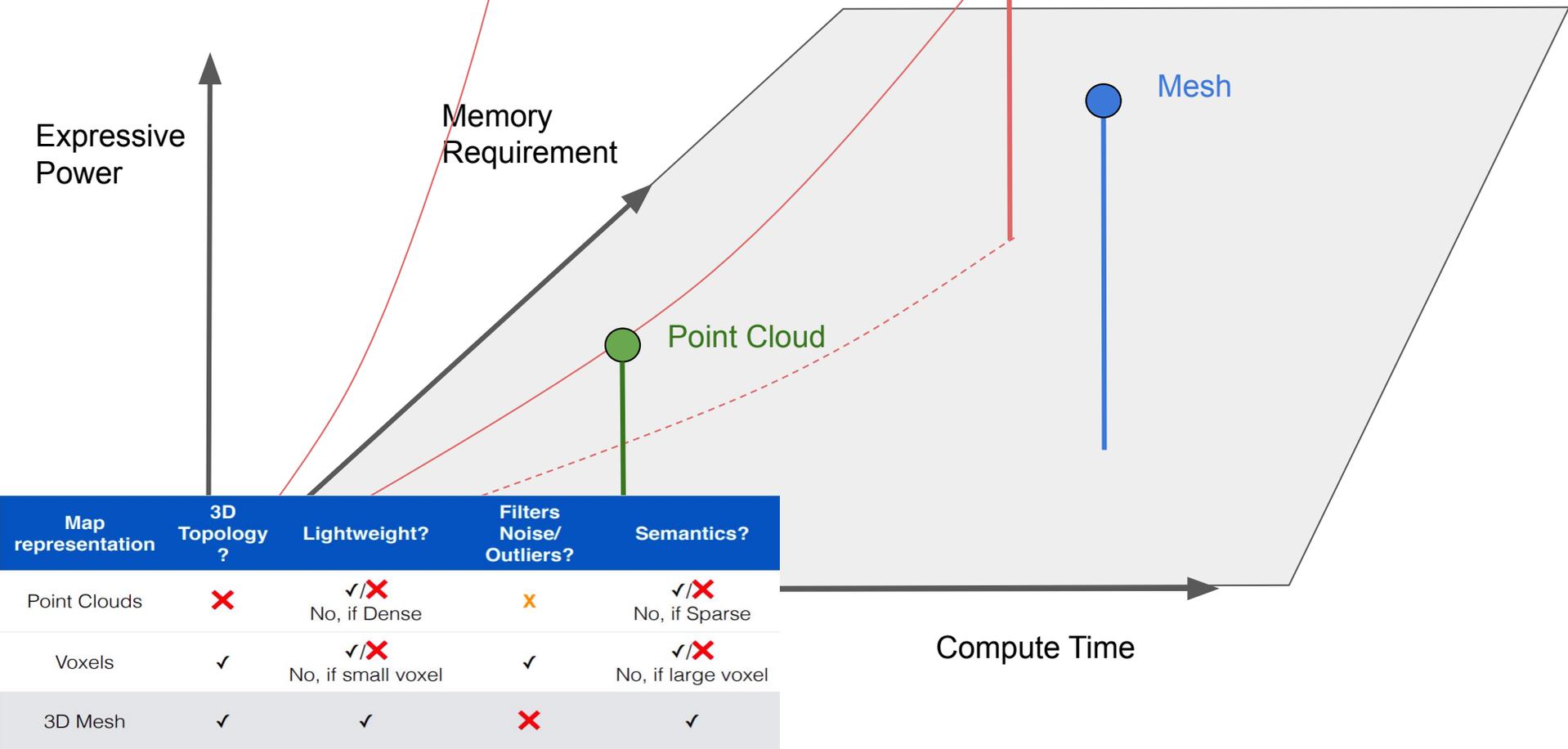
Point Cloud

Voxel

Mesh

| Map representation | 3D Topology ? | Lightweight? | Filters Noise/Outliers? | Semantics? |
|--------------------|---------------|---------------------------|-------------------------|---------------------------|
| Point Clouds | ✗ | ✓/✗ No, if Dense | ✗ | ✓/✗ No, if Sparse |
| Voxels | ✓ | ✓/✗ No, if small voxel | ✓ | ✓/✗ No, if large voxel |
| 3D Mesh | ✓ | ✓ | ✗ | ✓ |

Compute Time



Choice of Representation

Expressive Power

Memory Requirement

Point Cloud

Voxel

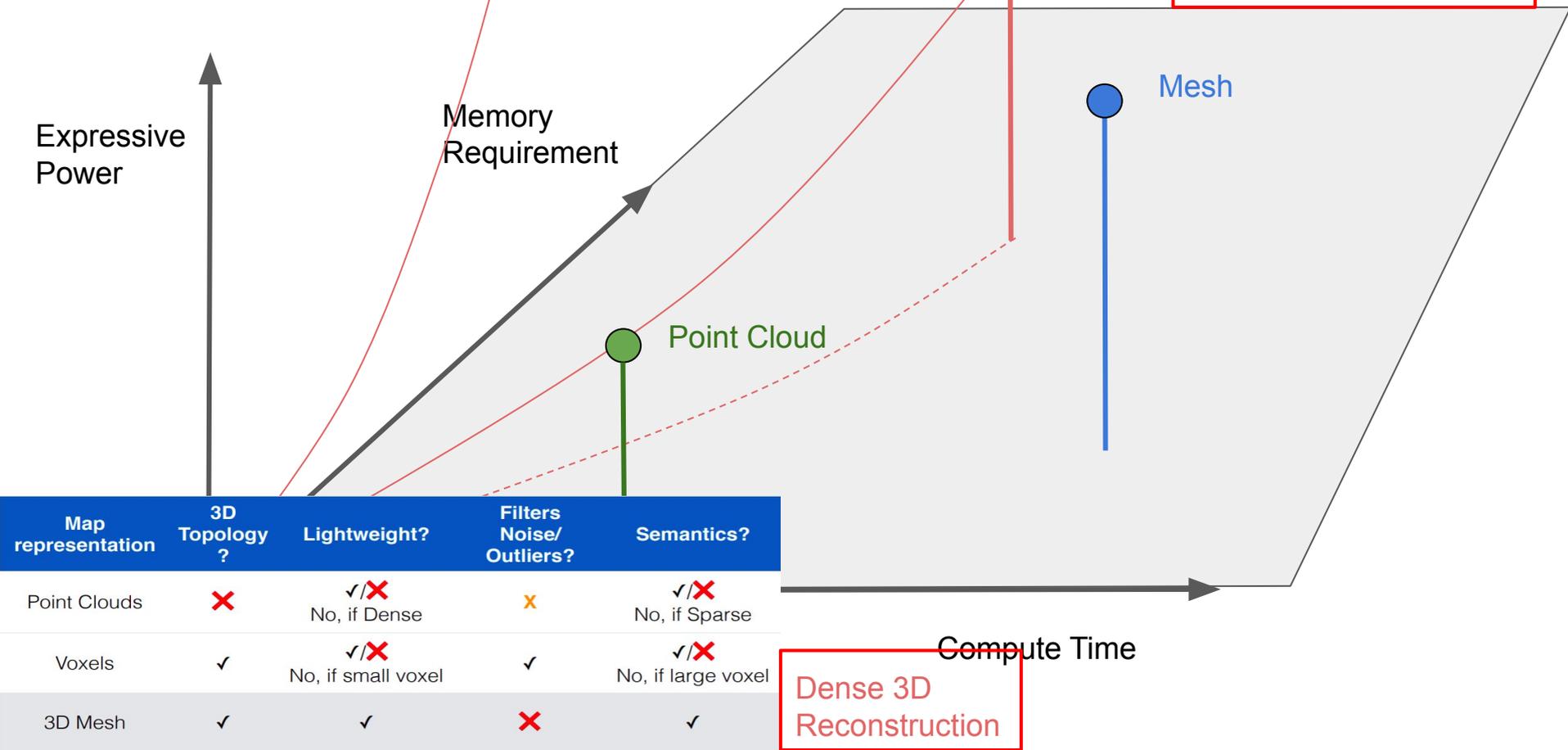
Mesh

Learning Architectures

Compute Time

Dense 3D Reconstruction

| Map representation | 3D Topology ? | Lightweight? | Filters Noise/Outliers? | Semantics? |
|--------------------|---------------|---------------------------|-------------------------|---------------------------|
| Point Clouds | ✗ | ✓/✗ No, if Dense | ✗ | ✓/✗ No, if Sparse |
| Voxels | ✓ | ✓/✗ No, if small voxel | ✓ | ✓/✗ No, if large voxel |
| 3D Mesh | ✓ | ✓ | ✗ | ✓ |



Second Part

2

Geometric Deep Learning

- Unifying view of developing architectures on all data
- Symmetry
- Equivariance, Invariance, Convolutions
- Unified Blueprint

Provide a unifying framework for developing deep learning architectures

References:

Bronstein et al. "Geometric Deep Learning: Grids, Groups, Graphs, Geodesics, and Gauges" 2021.
Bronstein et al. "Geometric Deep Learning" Lectures for AMMI, 2021.

Domains and Architectures

| | | | | | |
|-------------------|---------|--------------|-------|---------|-------------------------------|
| Image | 2D Grid | CNN | Voxel | 3D Grid | VoxNet, OctNet |
| Point Cloud | Sets | PointNet ... | Time | 1D Grid | RNN, LSTM |
| Mesh, Manifold | Graph | MeshCNN | Graph | Graph | EdgeConv, MeshCNN, GNN ... |

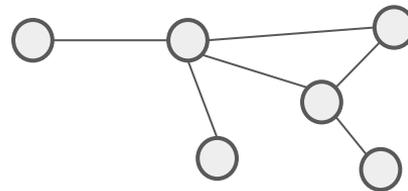
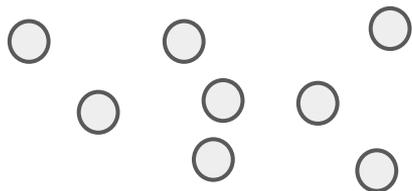
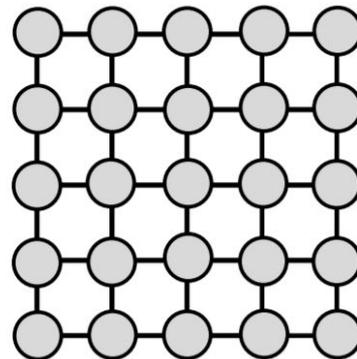
Domains and Architectures

| | | | | | |
|-------------------|---------|--------------|-------|---------|-------------------------------|
| Image | 2D Grid | CNN | Voxel | 3D Grid | VoxNet, OctNet |
| Point Cloud | Sets | PointNet ... | Time | 1D Grid | RNN, LSTM |
| Mesh, Manifold | Graph | MeshCNN | Graph | Graph | EdgeConv, MeshCNN, GNN ... |

Need an abstraction

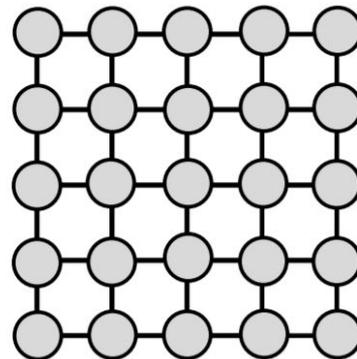
Domain

Ω

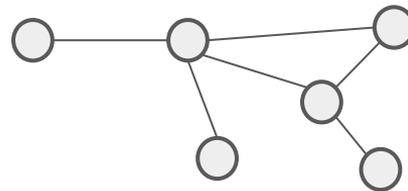


Signals on the Domain

Ω



$$\mathcal{X}(\Omega) = \{x : \Omega \rightarrow \mathbb{R}\}$$



Functions

Ω

$$\mathcal{F}_C = \{f : \mathcal{X}(\Omega) \rightarrow \mathbb{R}\}$$

$$\mathcal{F}_S = \{f : \mathcal{X}(\Omega) \rightarrow \mathcal{X}(\Omega)\}$$

Functions

Ω

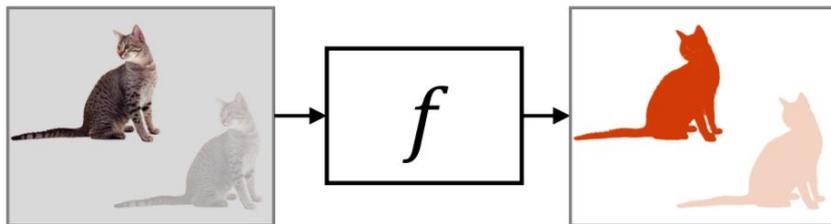
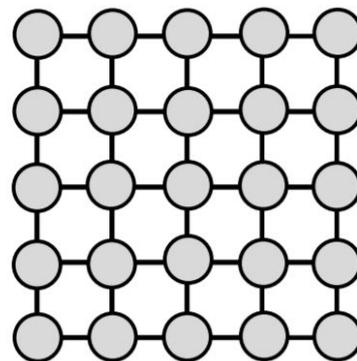
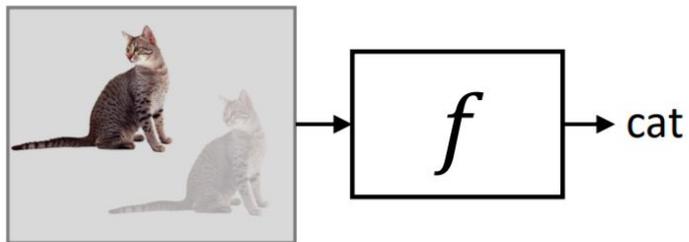


$$\mathcal{F}_C = \{f : \mathcal{X}(\Omega) \rightarrow \mathbb{R}\} \quad \text{Classification}$$

$$\mathcal{F}_S = \{f : \mathcal{X}(\Omega) \rightarrow \mathcal{X}(\Omega)\} \quad \text{Segmentation}$$

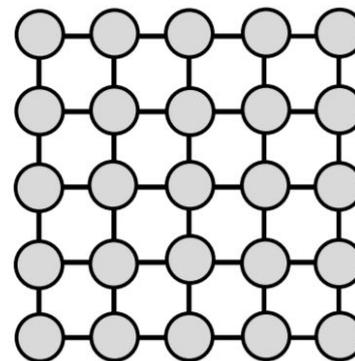
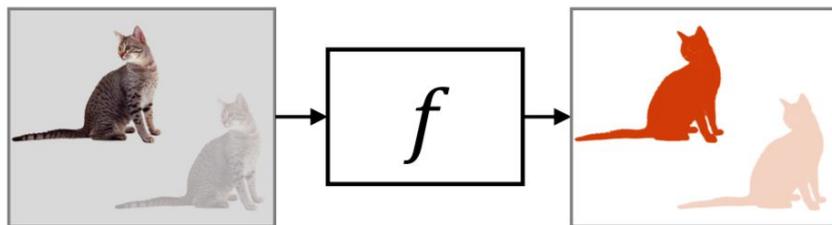
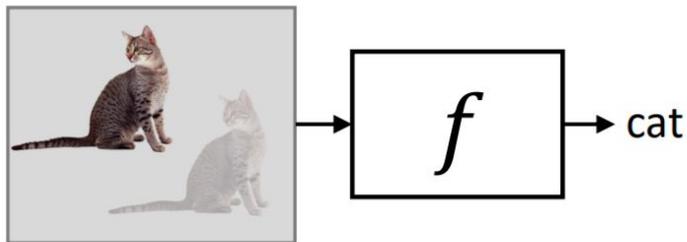
Symmetries

Ω



Symmetries

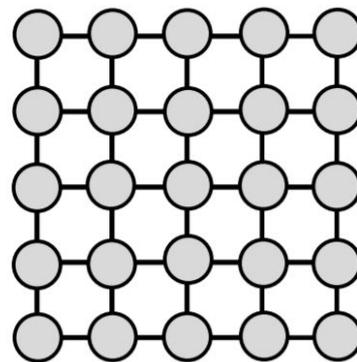
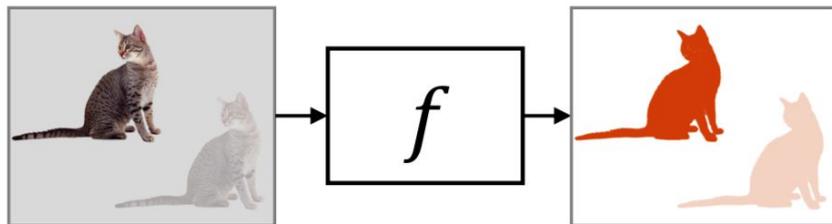
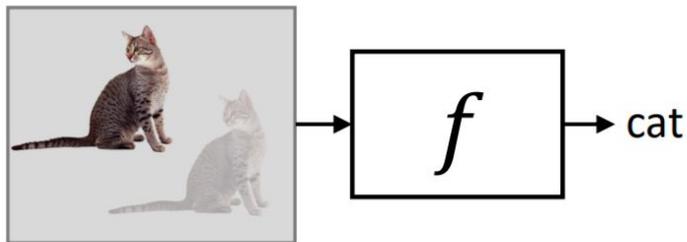
Ω



$S_{a,b}$ = shift or translation operator

Symmetries

Ω



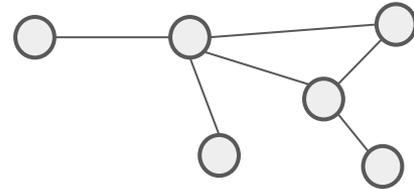
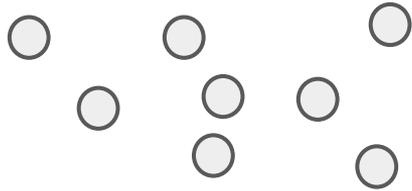
$$f(x) = f(S_{a,b} \cdot x)$$

$$f(x) = S_{a,b} \cdot f(x)$$

Symmetries

Ω

P = permutation of the domain



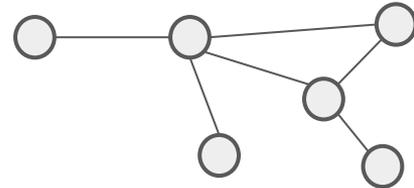
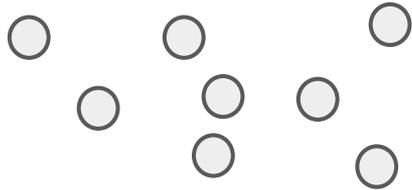
Symmetries

Ω

P = permutation of the domain

$$f(x) = f(P \cdot x)$$

$$f(x) = P \cdot f(x)$$



Space of Symmetries S

What structure does the set of all symmetries possess?

Space of Symmetries \mathcal{S}

What structure does the set of all symmetries possess?

1. Identity:

$$f(x) = f(i \cdot x) \quad i \cdot x = x$$

2. Inverse:

If $f(x) = f(g \cdot x)$ and $g \cdot x = y$ then there should exist $g^{-1} \cdot y = x$

3. Composition

If $f(x) = f(g \cdot x)$ and $f(y) = f(h \cdot y)$ then $f(x) = f((g \cdot h) \cdot x)$

Space of Symmetries \mathbb{S}

What structure does the set of all symmetries possess?

1. Identity: $i \in \mathbb{S}$

$$f(x) = f(i \cdot x) \quad i \cdot x = x$$

2. Inverse: $g \in \mathbb{S} \Rightarrow g^{-1} \in \mathbb{S}$

If $f(x) = f(g \cdot x)$ and $g \cdot x = y$ then there should exist $g^{-1} \cdot y = x$

3. Composition

If $f(x) = f(g \cdot x)$ and $f(y) = f(h \cdot y)$ then $f(x) = f((g \cdot h) \cdot x)$

$$g, h \in \mathbb{S} \Rightarrow g \cdot h \in \mathbb{S}$$

Space of Symmetries \mathbb{S}

What is this structure?

What structure does the set of all symmetries possess?

1. Identity: $i \in \mathbb{S}$

$$f(x) = f(i \cdot x) \quad i \cdot x = x$$

2. Inverse: $g \in \mathbb{S} \Rightarrow g^{-1} \in \mathbb{S}$

If $f(x) = f(g \cdot x)$ and $g \cdot x = y$ then there should exist $g^{-1} \cdot y = x$

3. Composition

If $f(x) = f(g \cdot x)$ and $f(y) = f(h \cdot y)$ then $f(x) = f((g \cdot h) \cdot x)$

$$g, h \in \mathbb{S} \Rightarrow g \cdot h \in \mathbb{S}$$

Space of Symmetries \mathcal{S}

What structure does the set of all symmetries over a domain possess?

Space of Symmetries = Group

$$\mathcal{S} = G$$

Space of Symmetries ~~S~~ G
Notation

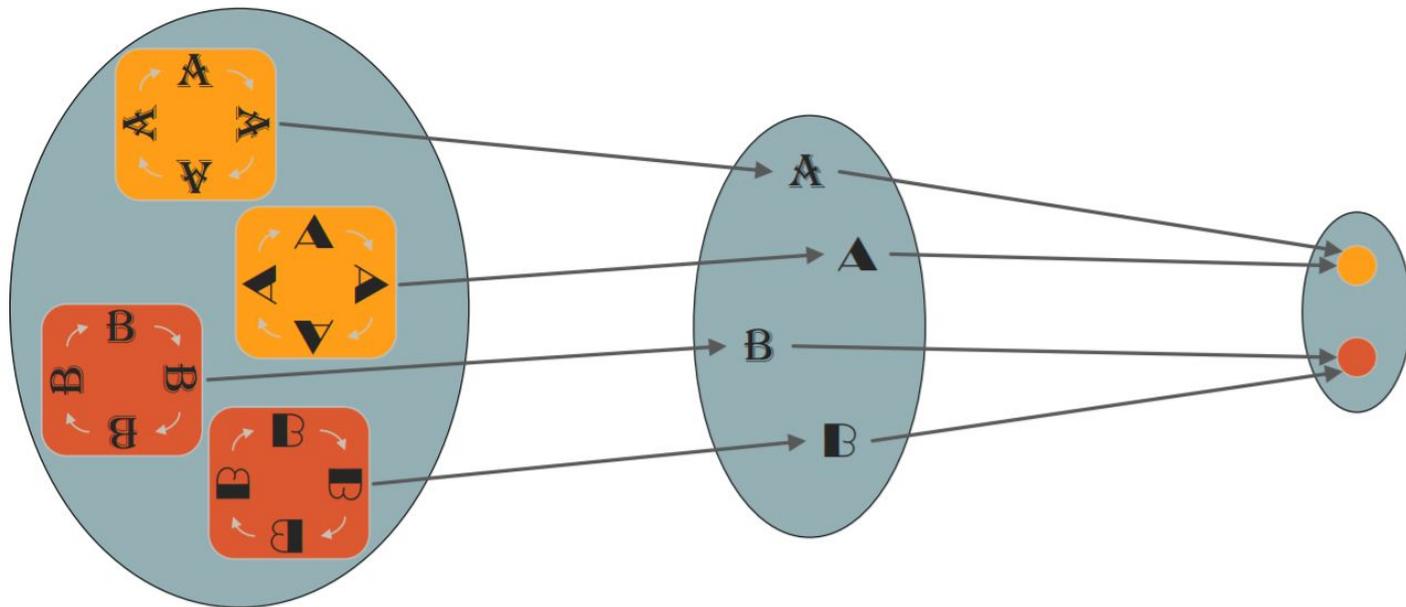
What structure does the set of all symmetries possess?

Space of Symmetries = Group

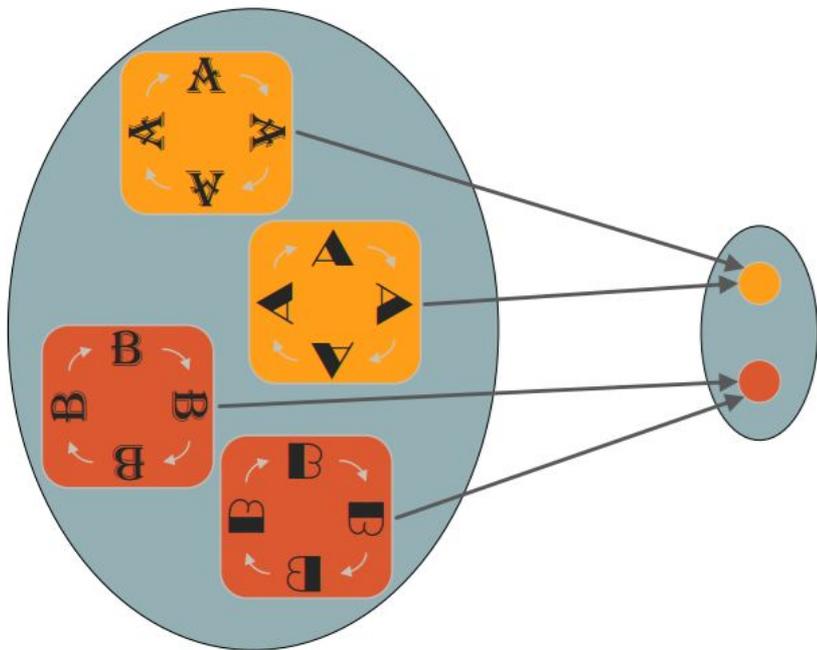
$$\del{S} = G$$

Notation

Equivalence Relation on $\mathcal{X}(\Omega)$



Equivalence Relation on $\mathcal{X}(\Omega)$



G-Equivalence

$$x \sim_G y \Leftrightarrow \exists g \in G : gx = y$$

Satisfies the axioms of an equivalence relation:

1. Reflexivity: $x \sim_G x$

- (Because G contains the identity)

2. Transitivity:

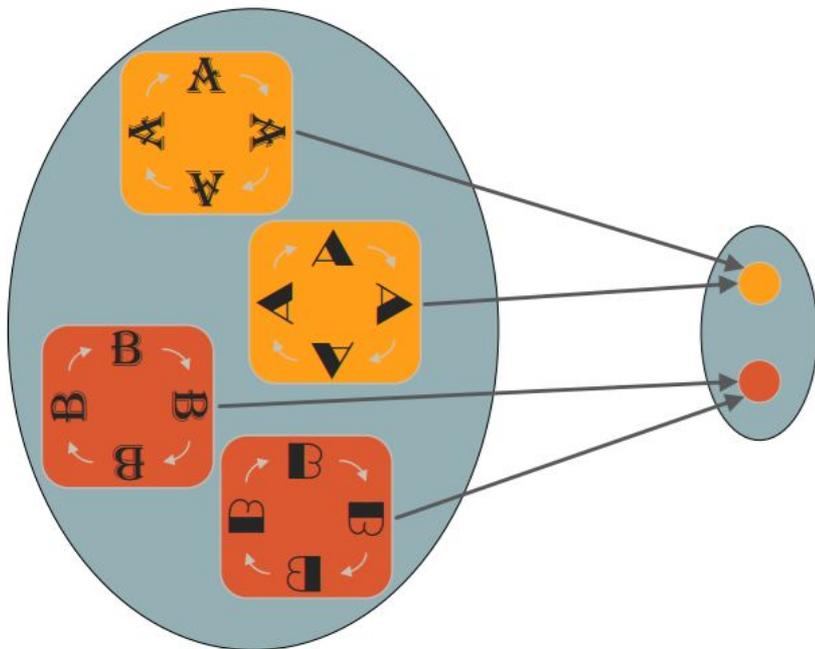
$$x \sim_G y \wedge y \sim_G z \Leftrightarrow x \sim_G z$$

- (Because G is closed under composition)

3. Symmetry: $x \sim_G y \Leftrightarrow y \sim_G x$

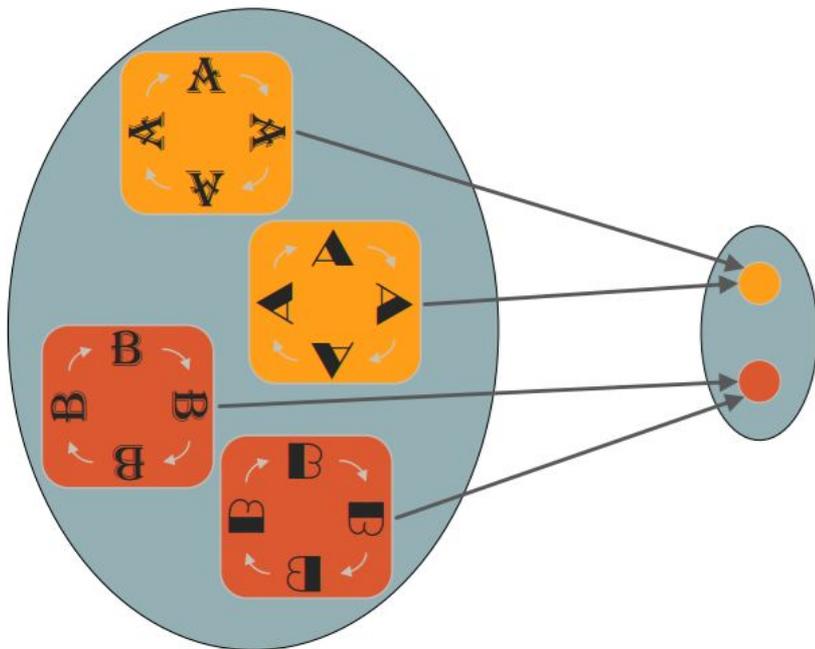
- (Because G is closed under inverses)

Equivalence Relation on $\mathcal{X}(\Omega)$



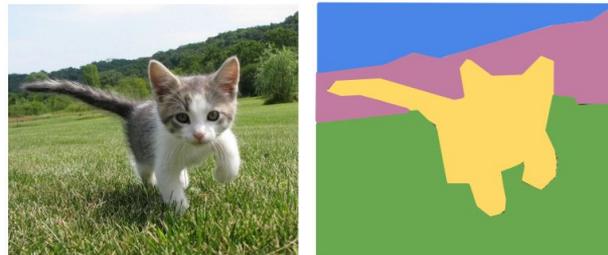
If we knew all the symmetries in the input signal space, we wouldn't need any data

Equivalence Relation on $\mathcal{X}(\Omega)$



*More symmetries we exploit,
the less data and parameters
we will need*

Adding Structure using Symmetry G



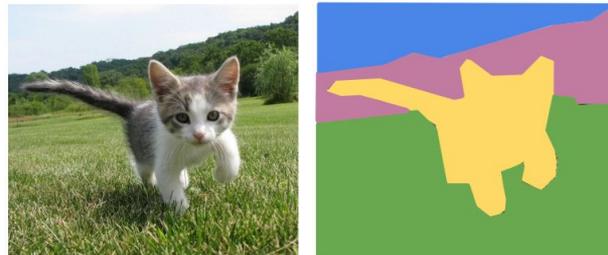
Classification

$$\mathcal{F}_C = \{f : \mathcal{X}(\Omega) \rightarrow \mathbb{R}\}$$

Segmentation

$$\mathcal{F}_S = \{f : \mathcal{X}(\Omega) \rightarrow \mathcal{X}(\Omega)\}$$

Adding Structure using Symmetry G



Classification

$$\mathcal{F}_C = \{f : \mathcal{X}(\Omega) \rightarrow \mathbb{R}\}$$

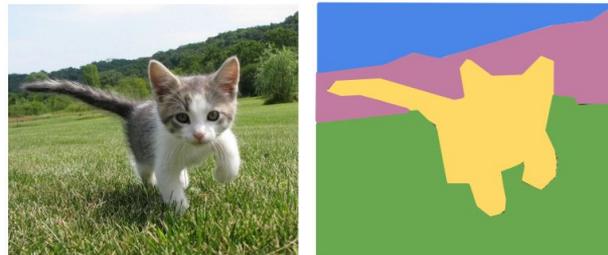
+ invariance to group action

Segmentation

$$\mathcal{F}_S = \{f : \mathcal{X}(\Omega) \rightarrow \mathcal{X}(\Omega)\}$$

+ equivariance to group action

Adding Structure using Symmetry G



Classification

$$\mathcal{F}_C = \{f : \mathcal{X}(\Omega) \rightarrow \mathbb{R} \mid f(g \cdot x) = f(x) \forall g \in G\}$$

Segmentation

$$\mathcal{F}_S = \{f : \mathcal{X}(\Omega) \rightarrow \mathcal{X}(\Omega) \mid f(g \cdot x) = g \cdot f(x) \forall g \in G\}$$

A Simple Example

Permutation invariant single layer perceptron

Translation Equivariance and Convolution

$$\Omega = [d] \times [d]$$

$$\mathcal{X}(\Omega) = \mathbb{R}^{d \times d}$$

Translation Equivariance and Convolution

$$\Omega = [d] \times [d]$$

$$G = \{S_{k,l} \mid S_{k,l} = \text{shift by } (k, l)\}$$

$$\mathcal{X}(\Omega) = \mathbb{R}^{d \times d}$$

Translation Equivariance and Convolution

$$\Omega = [d] \times [d]$$

$$\mathcal{X}(\Omega) = \mathbb{R}^{d \times d}$$

$$G = \{S_{k,l} \mid S_{k,l} = \text{shift by } (k, l)\}$$

$$(S_{k,l} \cdot x)(i, j) = x(i \oplus k, j \oplus l)$$

Translation Equivariance and Convolution

$$\Omega = [d] \times [d]$$

$$G = \{S_{k,l} \mid S_{k,l} = \text{shift by } (k, l)\}$$

$$\mathcal{X}(\Omega) = \mathbb{R}^{d \times d}$$

$$(S_{k,l} \cdot x)(i, j) = x(i \oplus k, j \oplus l)$$

$L : \mathcal{X}(\Omega) \rightarrow \mathcal{X}(\Omega)$ is a linear map

Translation Equivariance and Convolution

$$\Omega = [d] \times [d]$$

$$G = \{S_{k,l} \mid S_{k,l} = \text{shift by } (k, l)\}$$

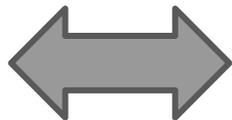
$$\mathcal{X}(\Omega) = \mathbb{R}^{d \times d}$$

$$(S_{k,l} \cdot x)(i, j) = x(i \oplus k, j \oplus l)$$

$L : \mathcal{X}(\Omega) \rightarrow \mathcal{X}(\Omega)$ is a linear map

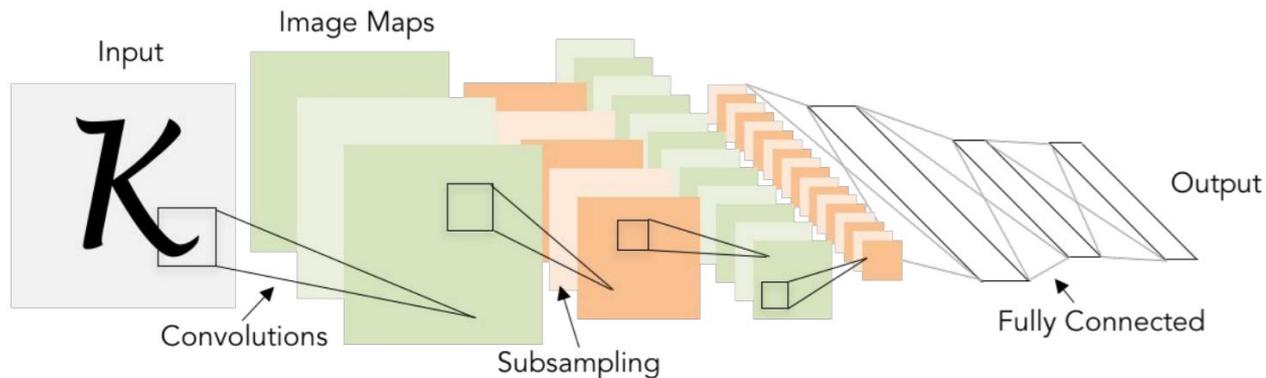
Theorem

$f(x) = \sigma(L(x))$ is
 G -equivariant



$L(x)$ is a convolution

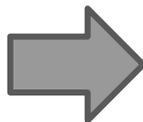
What does this mean?



Theorem

$$f^{\text{inv}} : \mathcal{X} \rightarrow \mathbb{R}$$

$$f_1^{\text{eq}}, f_2^{\text{eq}}, \dots, f_L^{\text{eq}} : \mathcal{X} \rightarrow \mathcal{X}$$



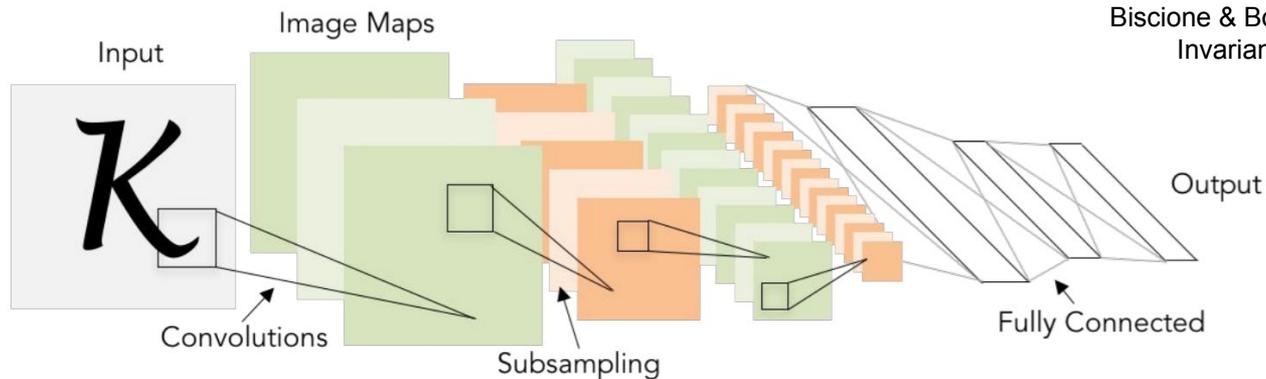
$$f^{\text{inv}} \cdot f_L^{\text{eq}} \cdot f_{L-1}^{\text{eq}} \cdots f_1^{\text{eq}} : \mathcal{X} \rightarrow \mathbb{R}$$

is G-invariant

What does this mean?

CNNs are not strictly translation invariant!

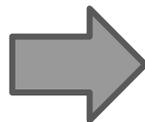
Biscione & Bowers "Learning Translation Invariance in CNNs" NeurIPS 2020



Theorem

$$f^{\text{inv}} : \mathcal{X} \rightarrow \mathbb{R}$$

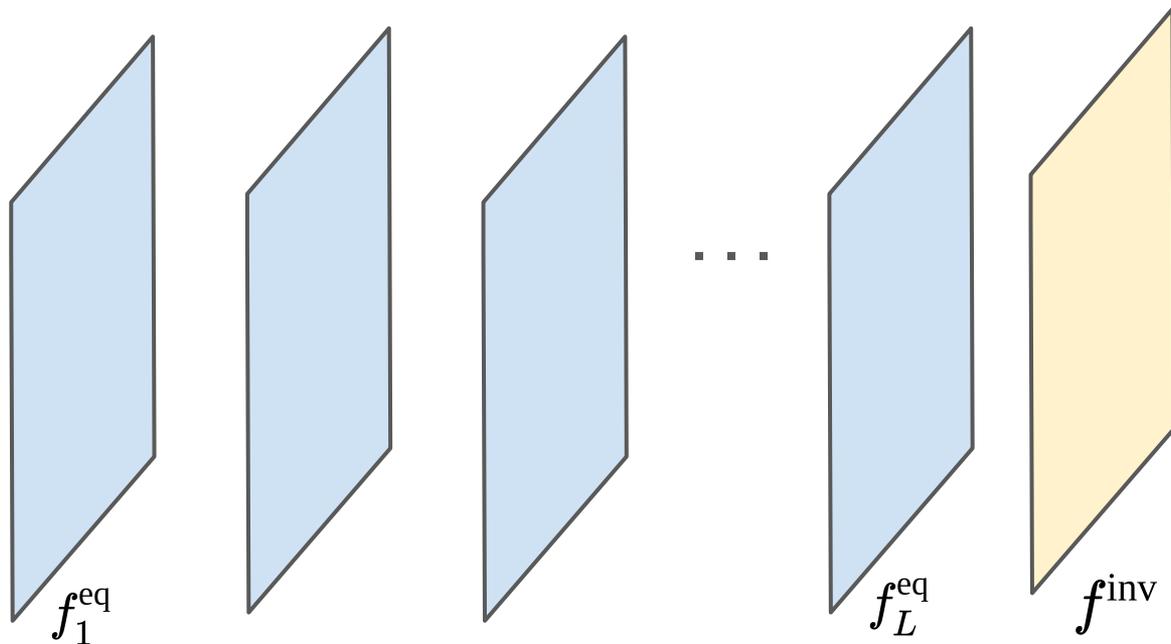
$$f_1^{\text{eq}}, f_2^{\text{eq}}, \dots, f_L^{\text{eq}} : \mathcal{X} \rightarrow \mathcal{X}$$



$$f^{\text{inv}} \cdot f_L^{\text{eq}} \cdot f_{L-1}^{\text{eq}} \cdots f_1^{\text{eq}} : \mathcal{X} \rightarrow \mathbb{R}$$

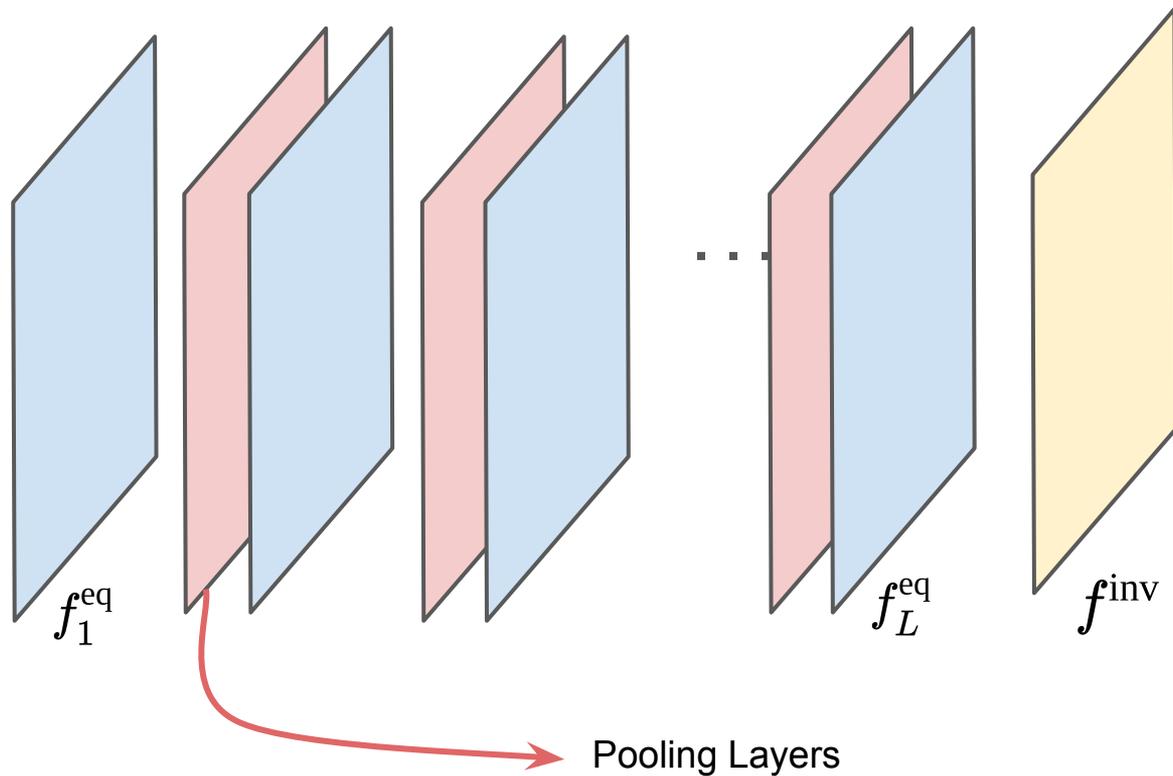
is G-invariant

Geometric Deep Learning Blueprint



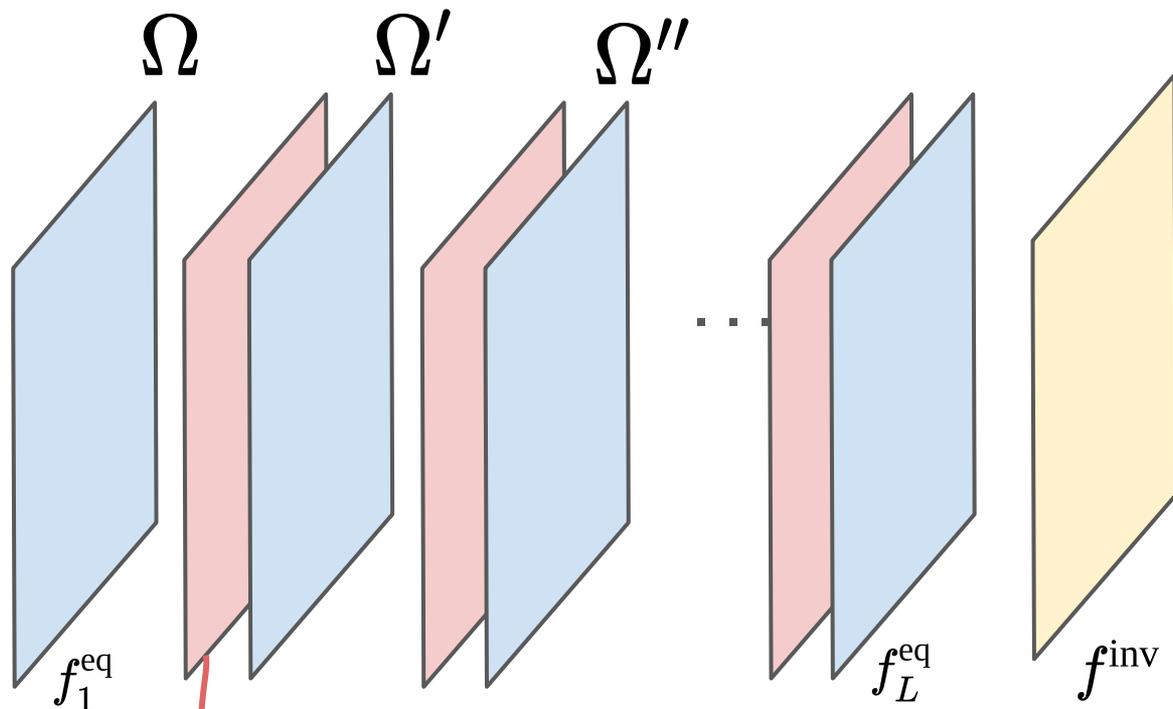
Classification

Geometric Deep Learning Blueprint



Classification

Geometric Deep Learning Blueprint

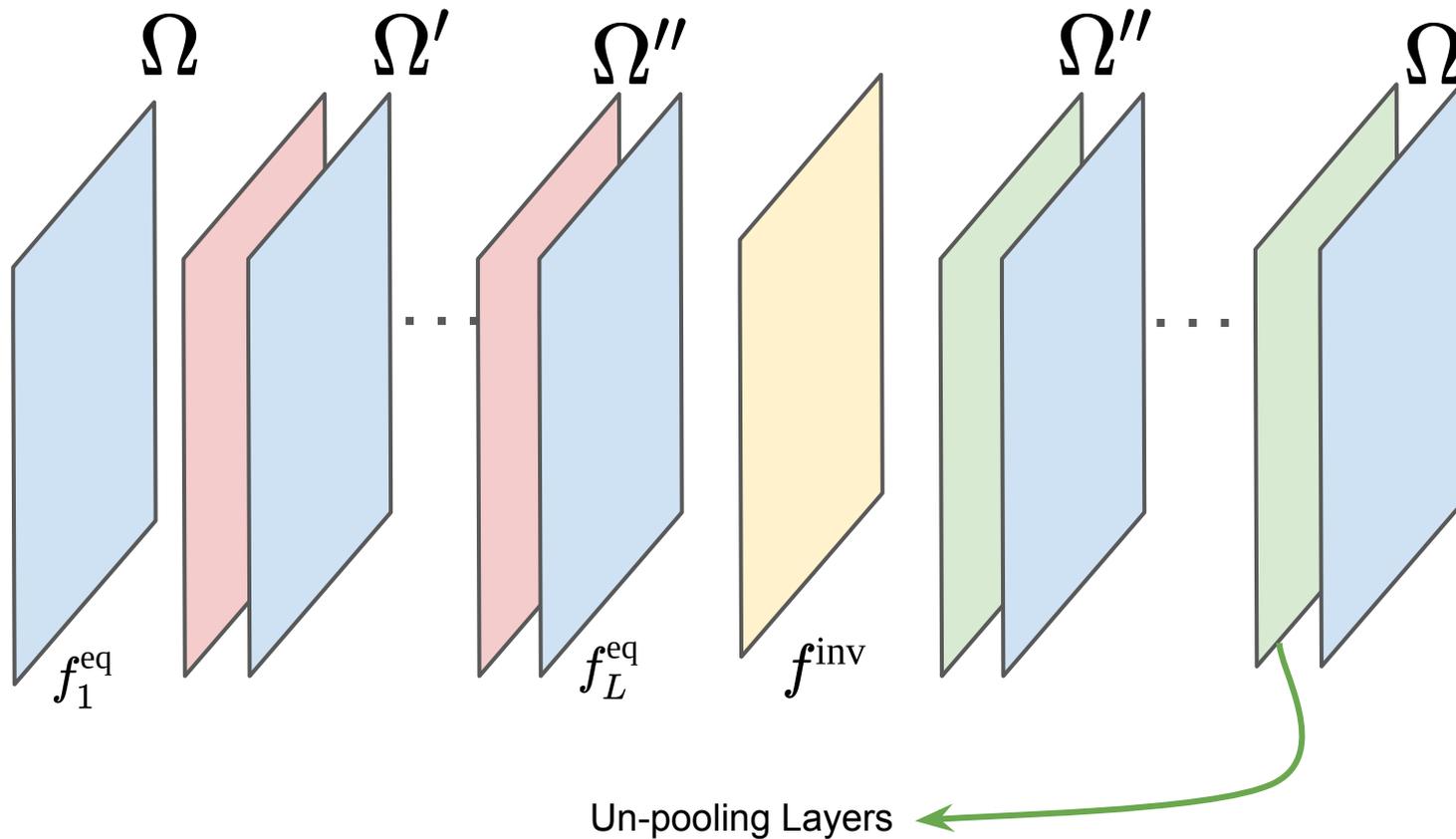


Classification

 Pooling Layers

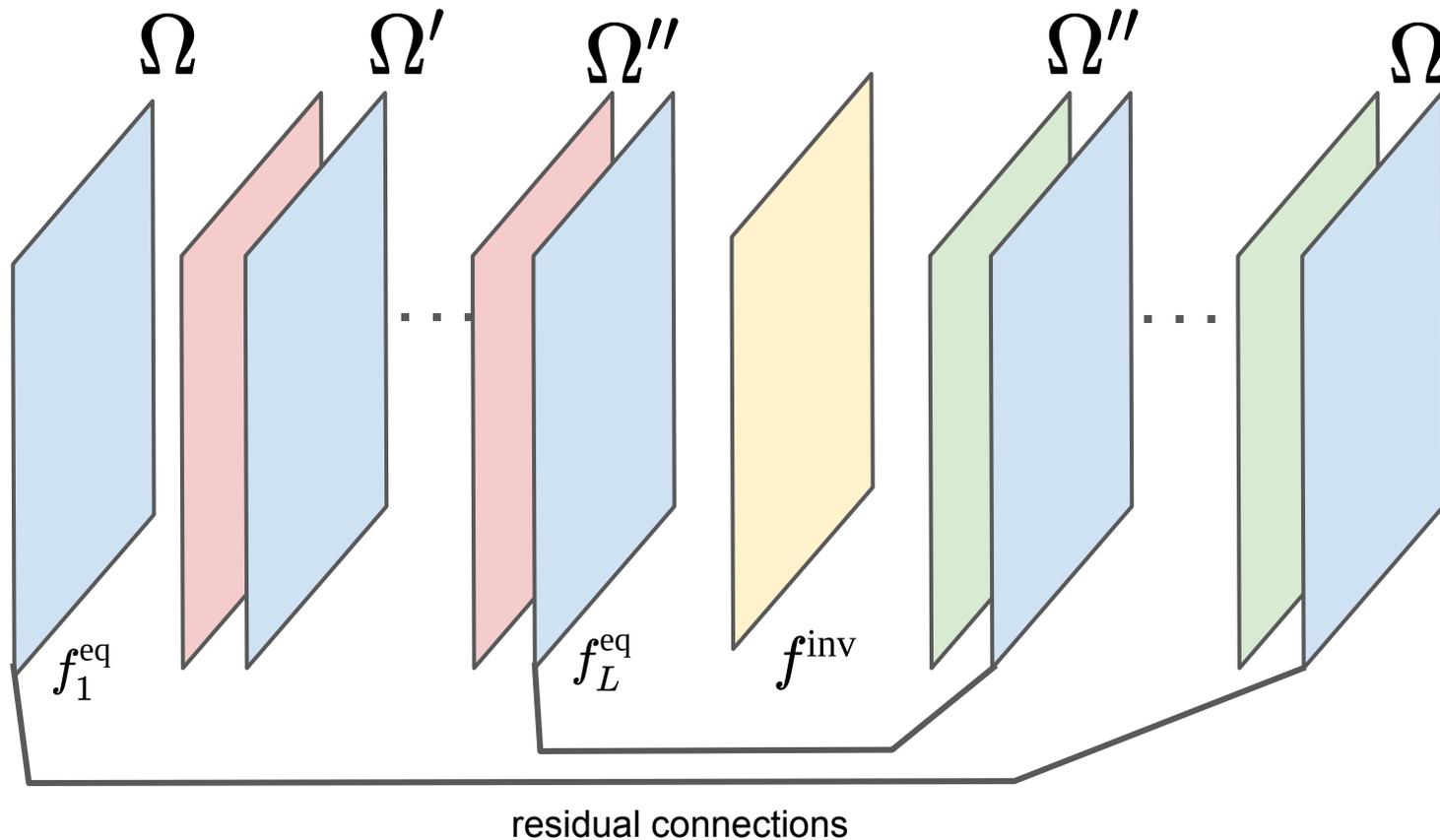
Geometric Deep Learning Blueprint

Segmentation



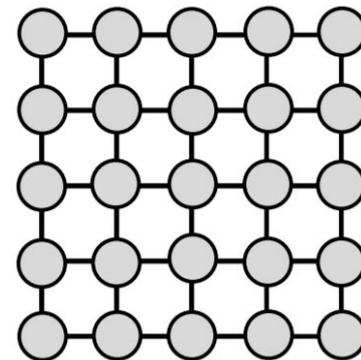
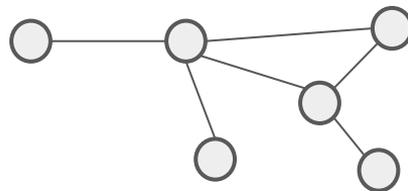
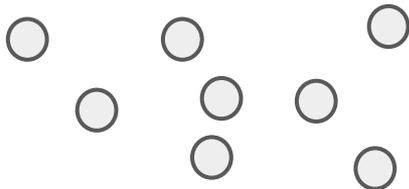
Geometric Deep Learning Blueprint

Segmentation



Using the Blueprint

Suffices to find invariant and equivariant functions on different domains



Sets

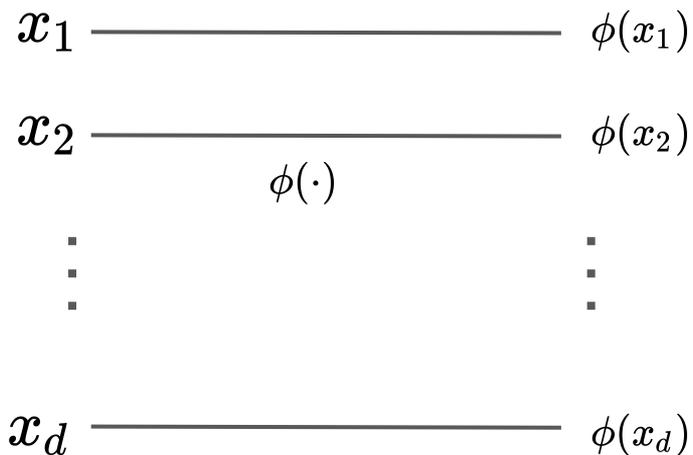
Equivariance over Sets

$$\Omega = [d] \quad \mathcal{X}(\Omega) = \mathbb{R}^d \quad G = \{P \mid P = d \times d \text{ permutation matrix}\}$$

$$f(P \cdot x) = P \cdot f(x)$$

Equivariance over Sets

$$\Omega = [d] \quad \mathcal{X}(\Omega) = \mathbb{R}^d \quad G = \{P \mid P = d \times d \text{ permutation matrix}\}$$



Equivariance over Sets

$$\Omega = [d]$$

$$\mathcal{X}(\Omega) = \mathbb{R}^d$$

$$G = \{P \mid P = d \times d \text{ permutation matrix}\}$$

$$x_1 \text{ ————— } \phi(x_1)$$

$$x_2 \text{ ————— } \phi(x_2)$$

⋮

⋮

$$x_d \text{ ————— } \phi(x_d)$$

$\phi(\cdot)$

*Can be a permutation
invariant function of all
the inputs*



Recall: Point Transformer

Basic version

$$x'_j = \sum_{i \in N(x_j)} \rho(\phi(x_j)^T \psi(x_i)) \cdot \alpha(x_i)$$

Incorporating point feature + location; and using vector for attention

$$x'_j = \sum_{i \in N(x_j)} \rho[\beta(\phi(x_j), \psi(x_i)) + \delta(p_j - p_i)] \odot \alpha(x_i)$$

function other than dot product

position of points

Invariance over Sets

$$\Omega = [d] \quad \mathcal{X}(\Omega) = \mathbb{R}^d \quad G = \{P \mid P = d \times d \text{ permutation matrix}\}$$

$$f(P \cdot x) = f(x)$$

Invariance over Sets

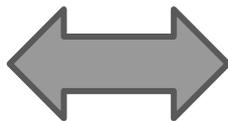
$$\Omega = [d] \quad \mathcal{X}(\Omega) = \mathbb{X}^d \quad G = \{P \mid P = d \times d \text{ permutation matrix}\}$$

 *countable set*

Theorem:

$$f : \mathcal{X}(\Omega) \rightarrow \mathbb{R}$$

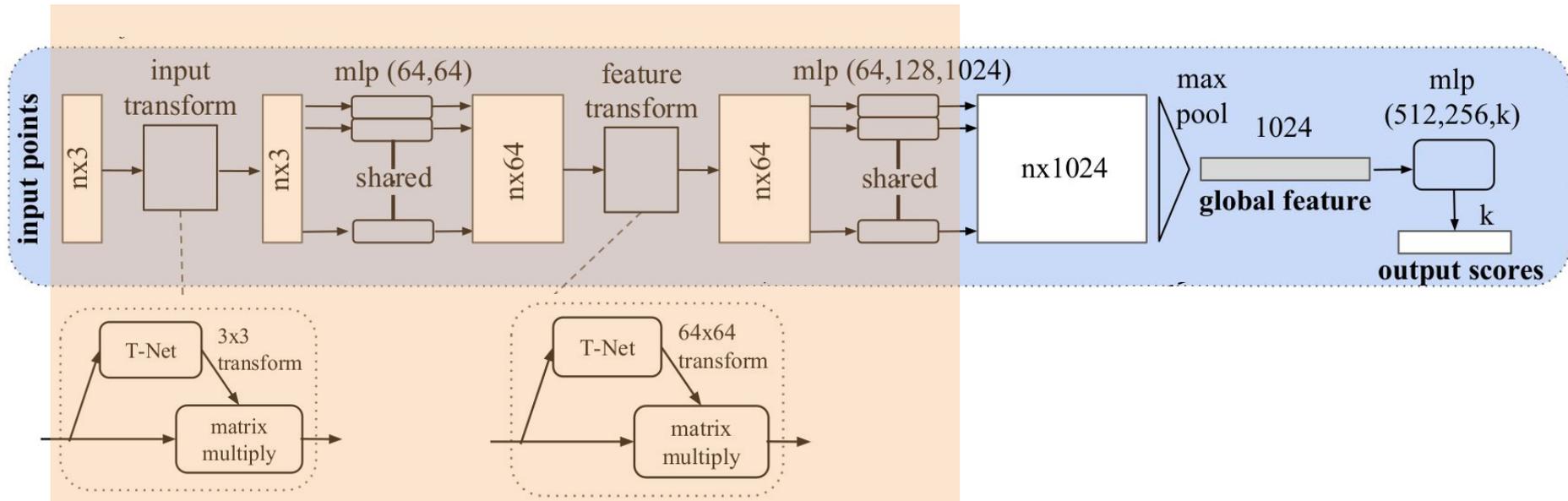
is G -invariant



$$f(\mathbf{x}) = \psi \left(\sum_{i=1}^d \phi(x_i) \right)$$

$\exists \phi, \psi$

Recall: PointNet Architecture



$$f(\{x_1, x_2, \dots, x_n\}) = \max\{h(x_1), h(x_2), \dots, h(x_n)\}$$

Simple Example

Permutation equivariant single layer perceptron

Question

Why not use the blueprint with permutation invariant and equivariant single layer perceptron?

Backup

Error Decomposition