# Cross-Shape Attention for Part Segmentation of 3D Point Clouds

**Marios Loizou[†1]**  **Siddhant Garg [†2]**  **Dmitry Petrov [†2]**

**Melinos Averkiou[1]**  **Evangelos Kalogerakis[2]**

*[1]University of Cyprus / CYENS CoE*   *[2]University of Massachusetts Amherst*

**† Equal Contribution**

# Goal: learn more coordinated feature representations



**test** shape

# Goal: learn more coordinated feature representations
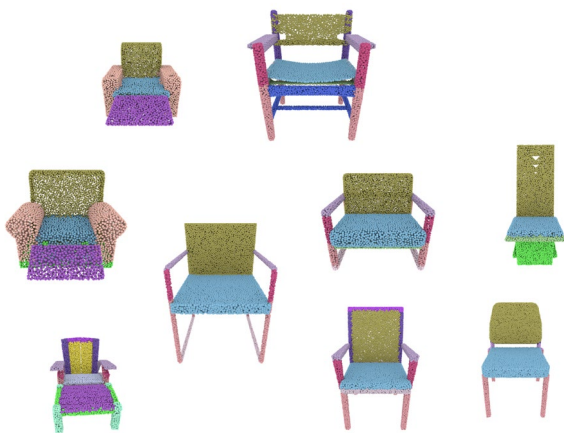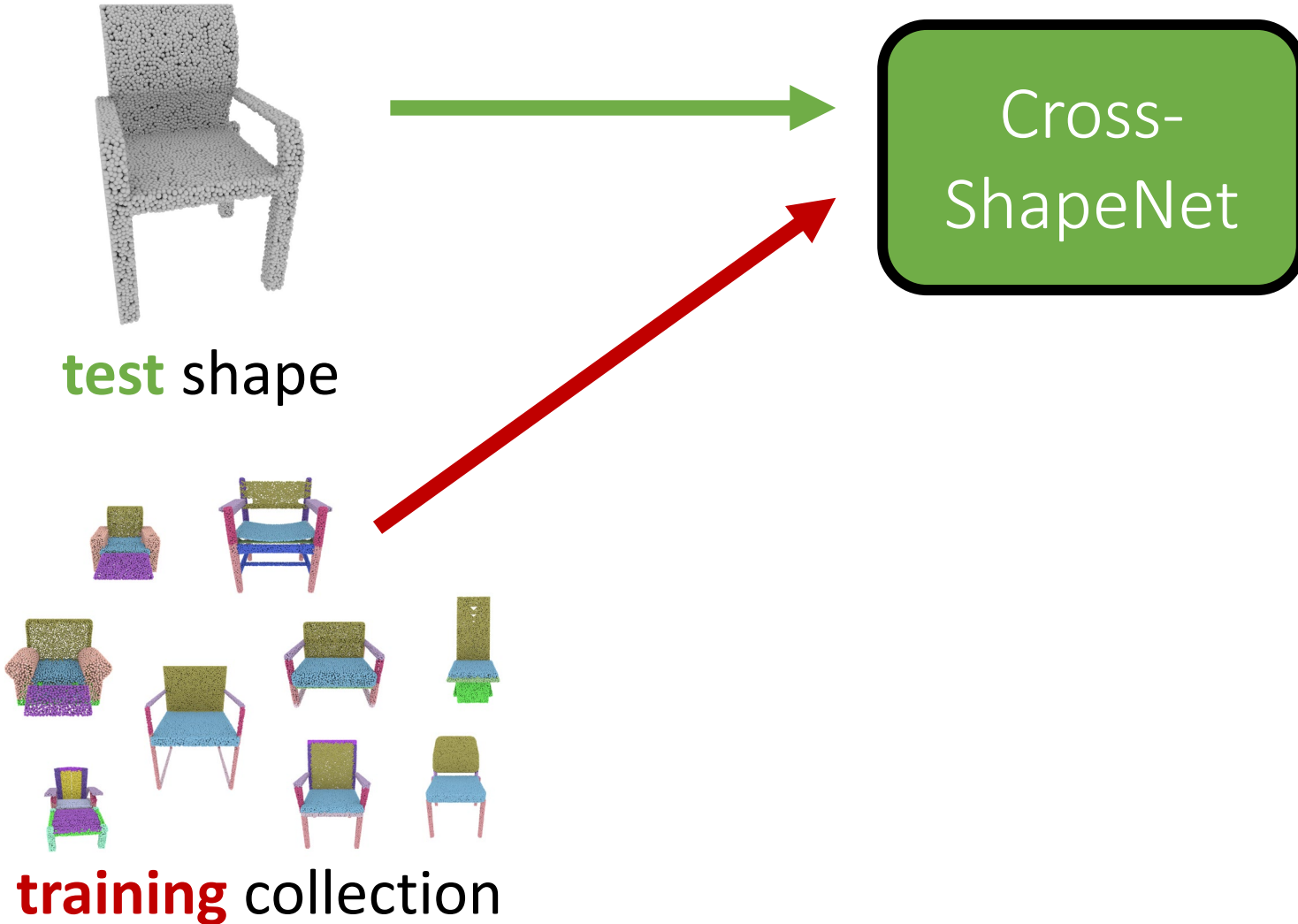


**test** shape

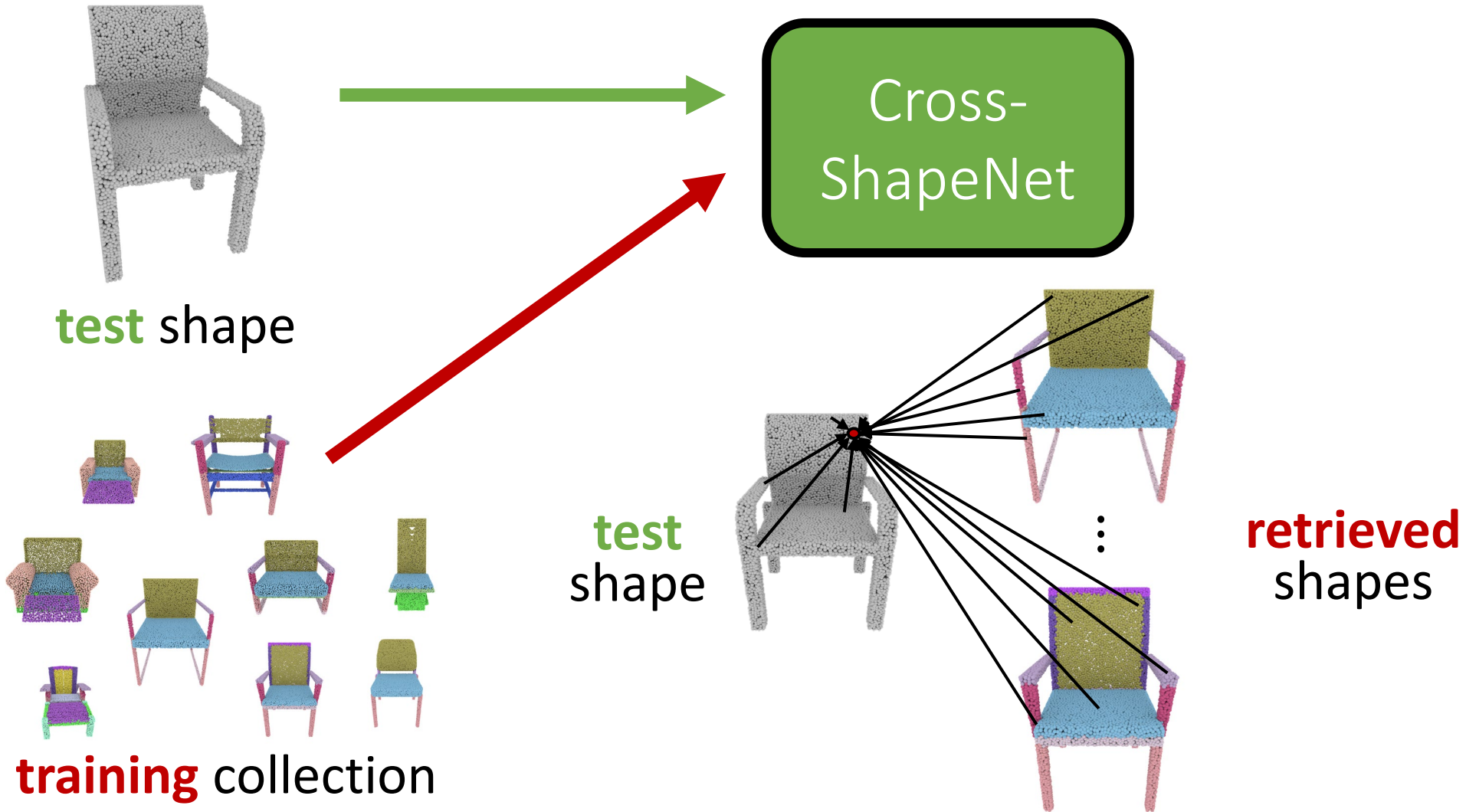**training** collection

# Goal: learn more coordinated feature representations



test shape

training collection

Cross-ShapeNet

**Goal**: learn more coordinated feature representations

Cross-ShapeNet

**test** shape

**training** collection

**test** shape

**retrieved** shapes

**Goal**: learn more coordinated feature representations

**test** shape

Cross-ShapeNet

**test** shape segmentation

- Back
- Armrest top
- Armrest bar
- Seat
- Leg

**training** collection

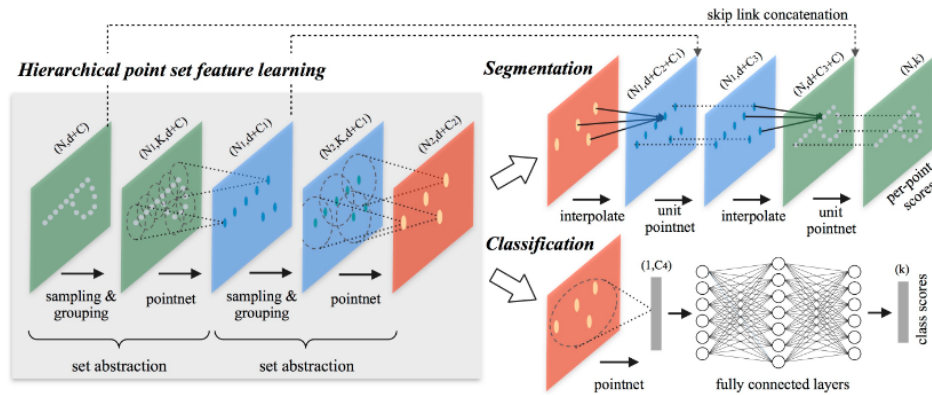**test** shape

**retrieved** shapes

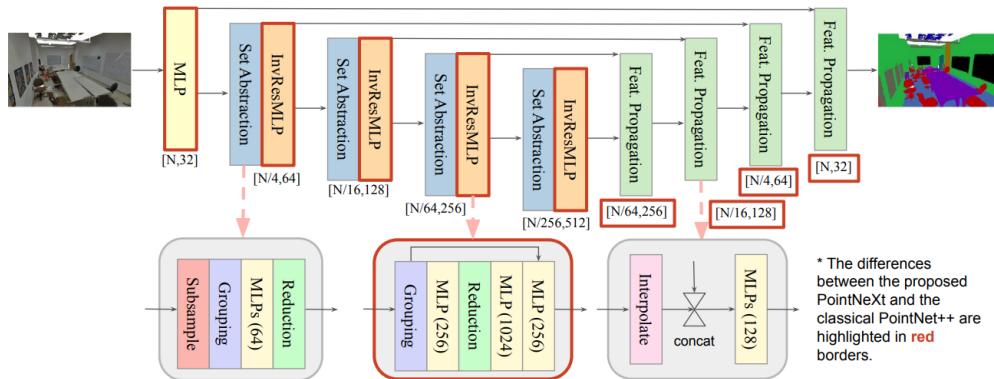# **Prior work**: Point-based networks
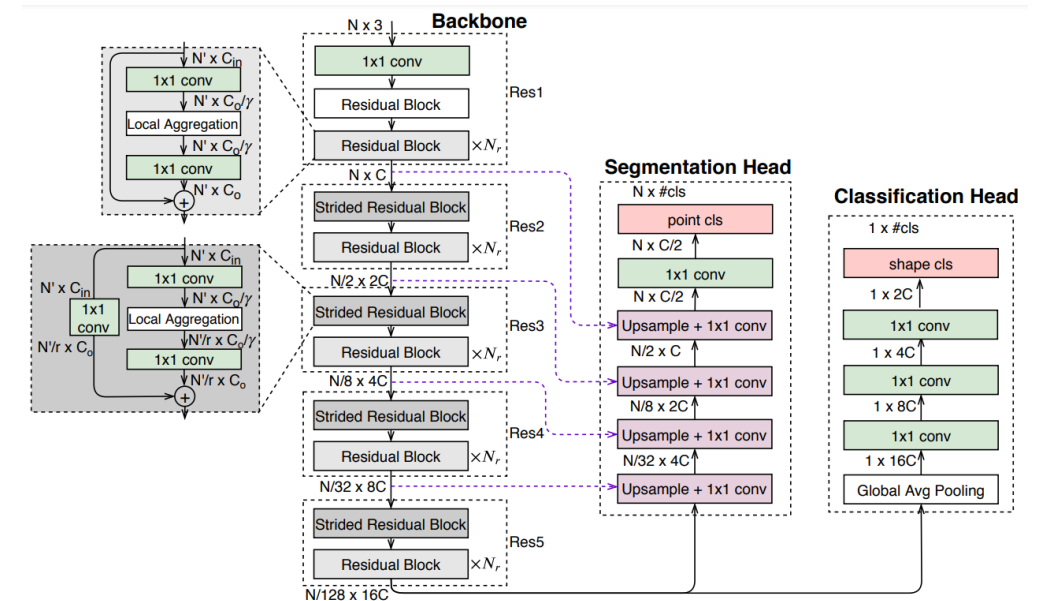


PointNet++ [Qi et al. 2017]

# Prior work: Point-based networks



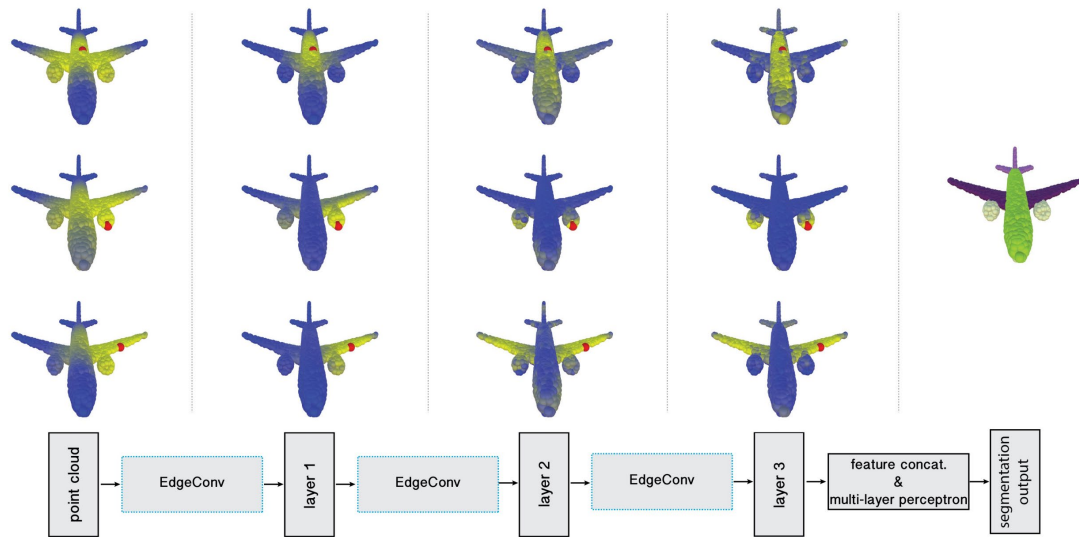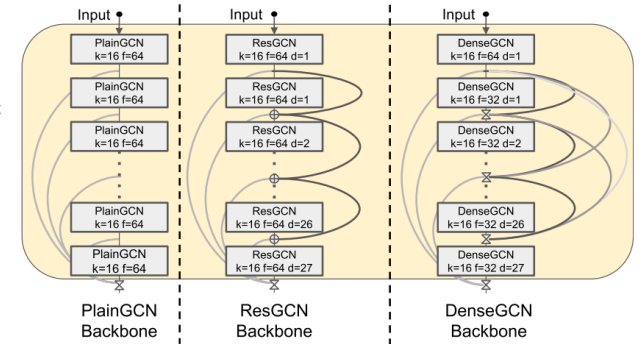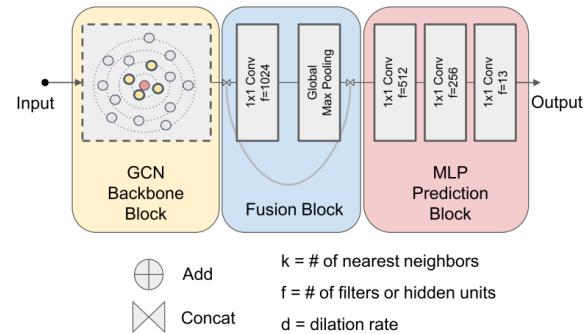PointNet++ [Qi et al. 2017]



PointNeXt [Qian et al. 2022]



CloserLook3D [Liu et al. 2020]
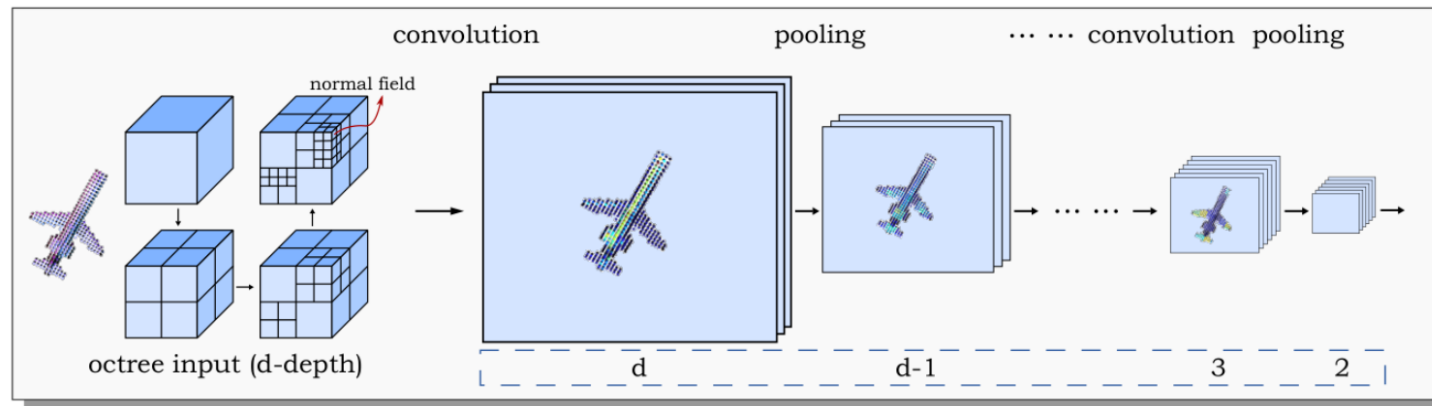
# Prior work: GCNs for non-Eucledian data
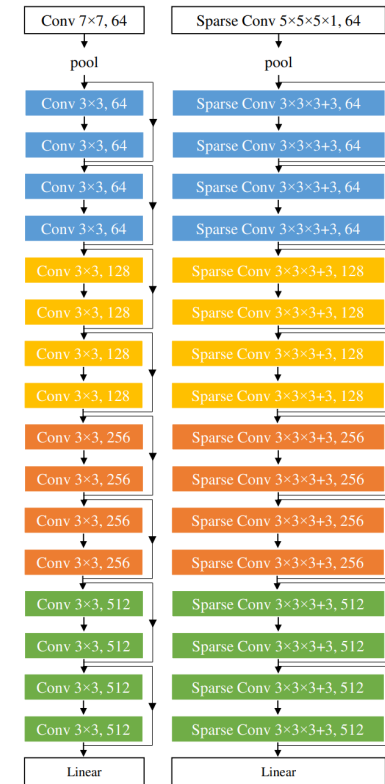


DGCNN [Wang et al. 2019]

DeepGCNs [Li et al. 2023]
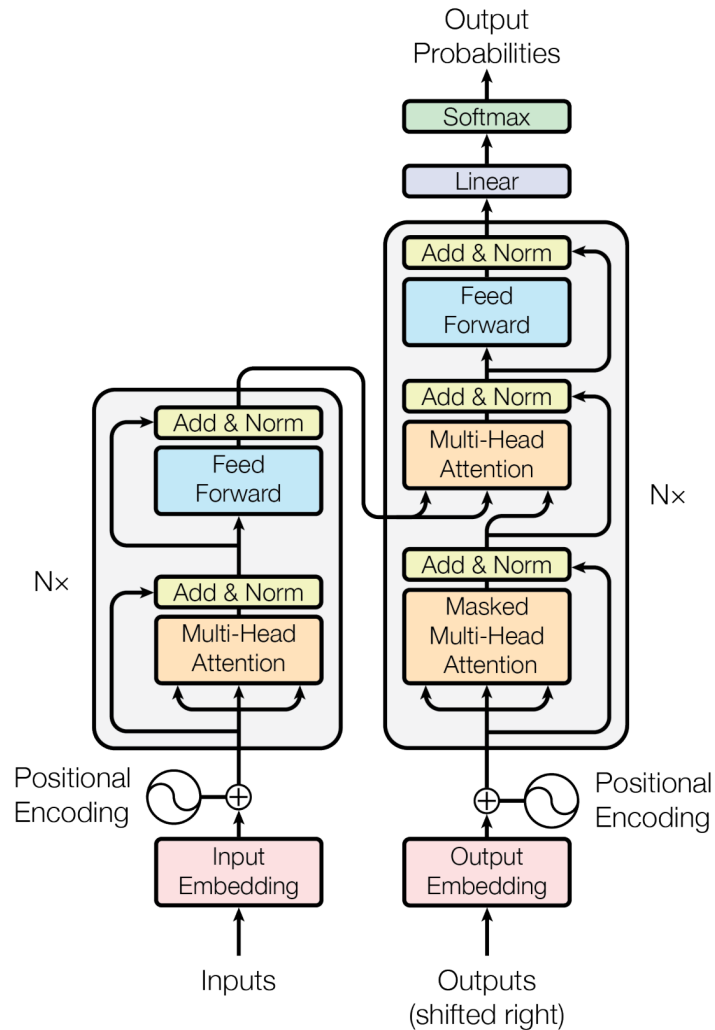
# Prior work: Volumetric networks



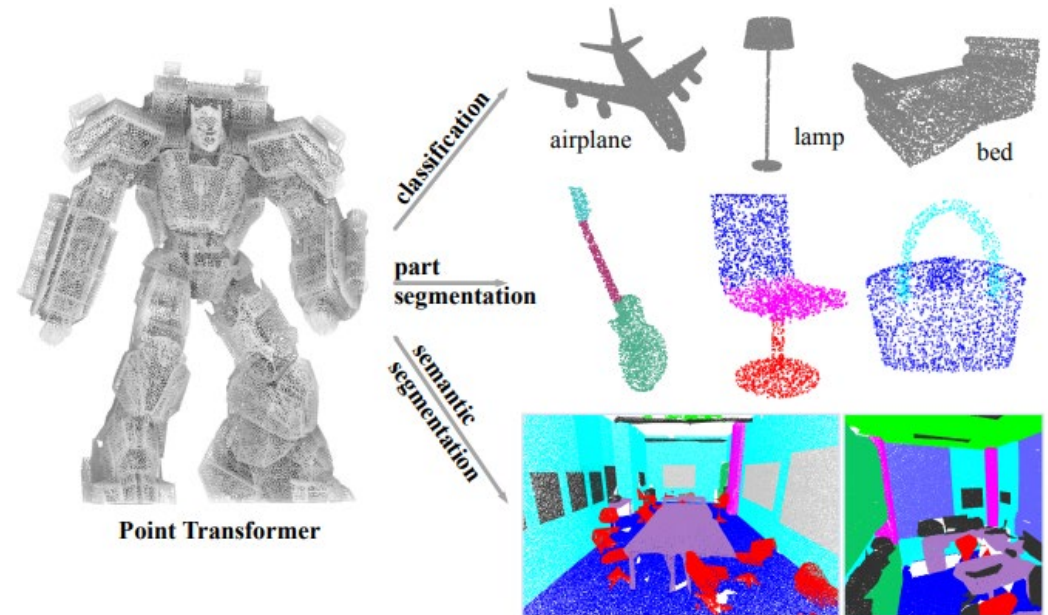O-CNN [Wang et al. 2017]



MinkowskiNet [Choy et al. 2019]

# Prior work: Attention is All You Need



Transformer [Vaswani et al. 2017]



PointTransformer v1/v2 [Zhao et al. 2021, 2022]

# Why use **attention** for 3D representations?

Encode points such that their features capture **relations** wrt the rest of the shape
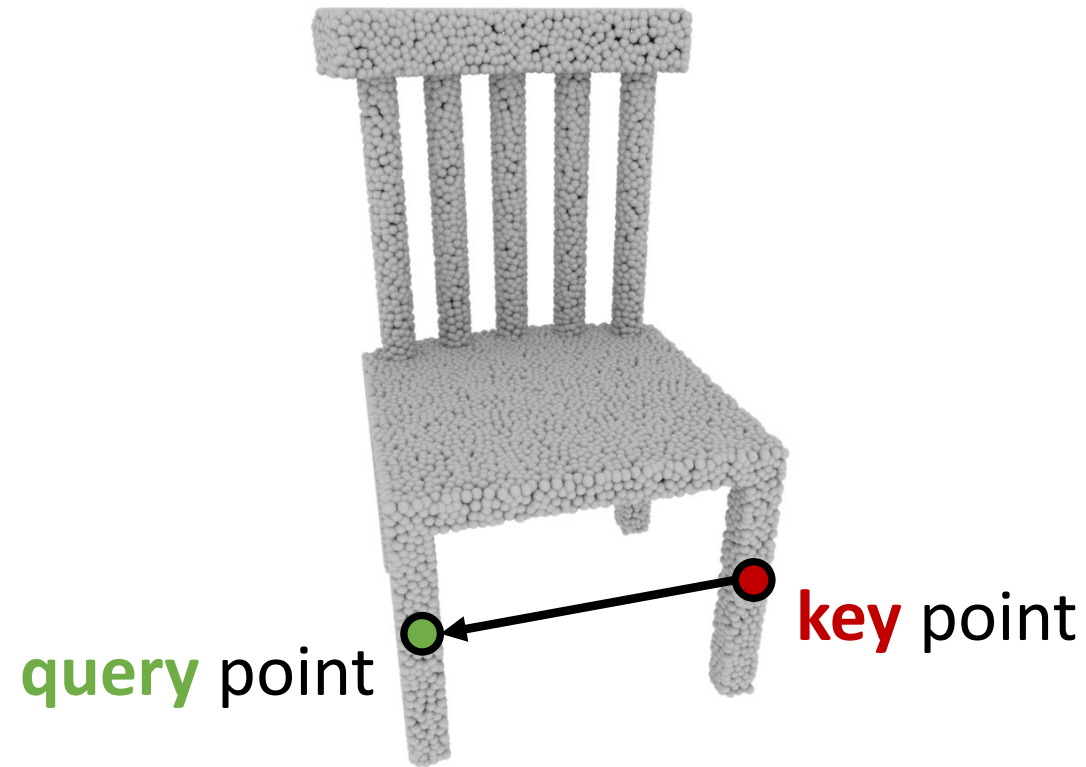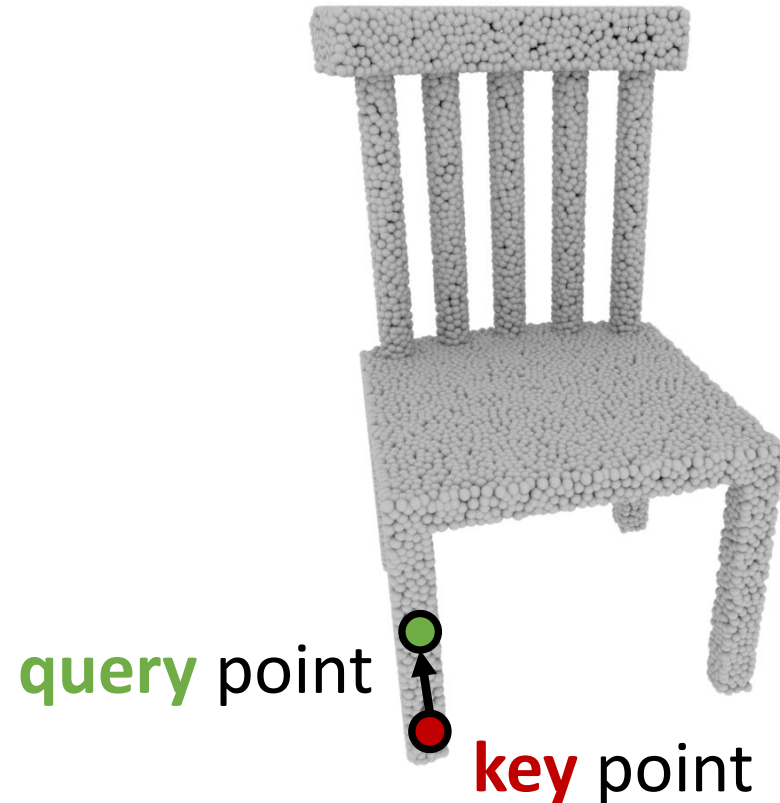


**query** point

# Why use **attention** for 3D representations?

Encode points such that their features capture **relations** wrt the rest of the shape

# Why use **attention** for 3D representations?

Encode points such that their features capture **relations** wrt the rest of the shape



**query** point

**key** point

# Why use **attention** for 3D representations?

Encode points such that their features capture **relations** wrt the rest of the shape
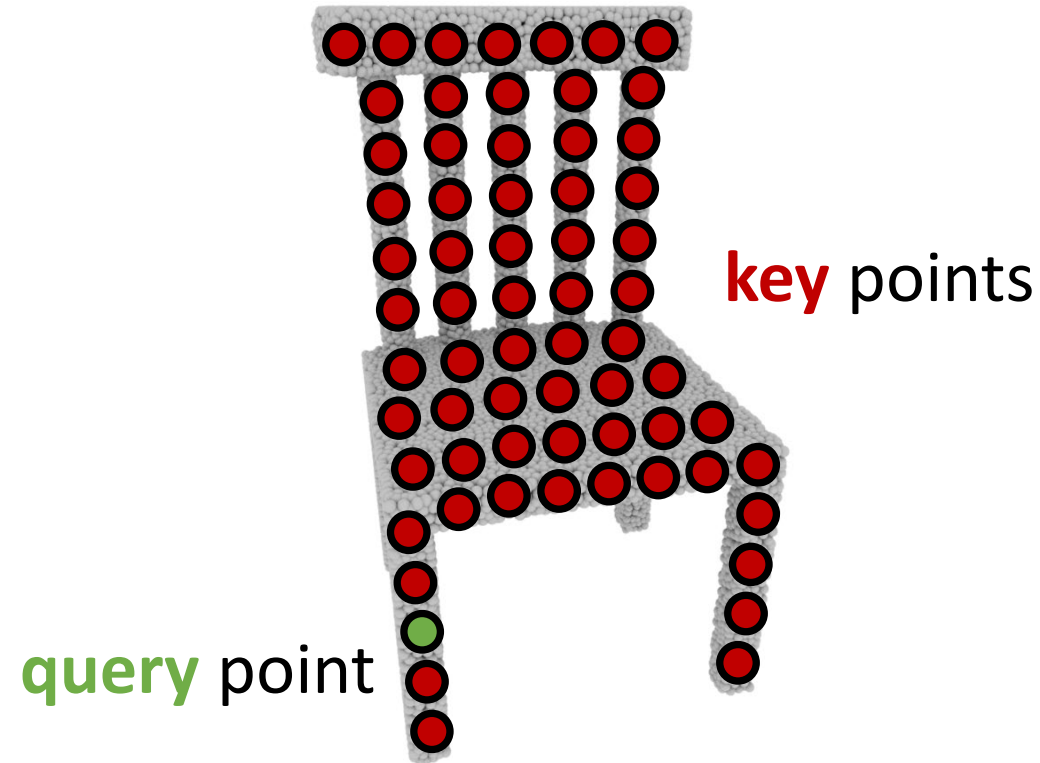


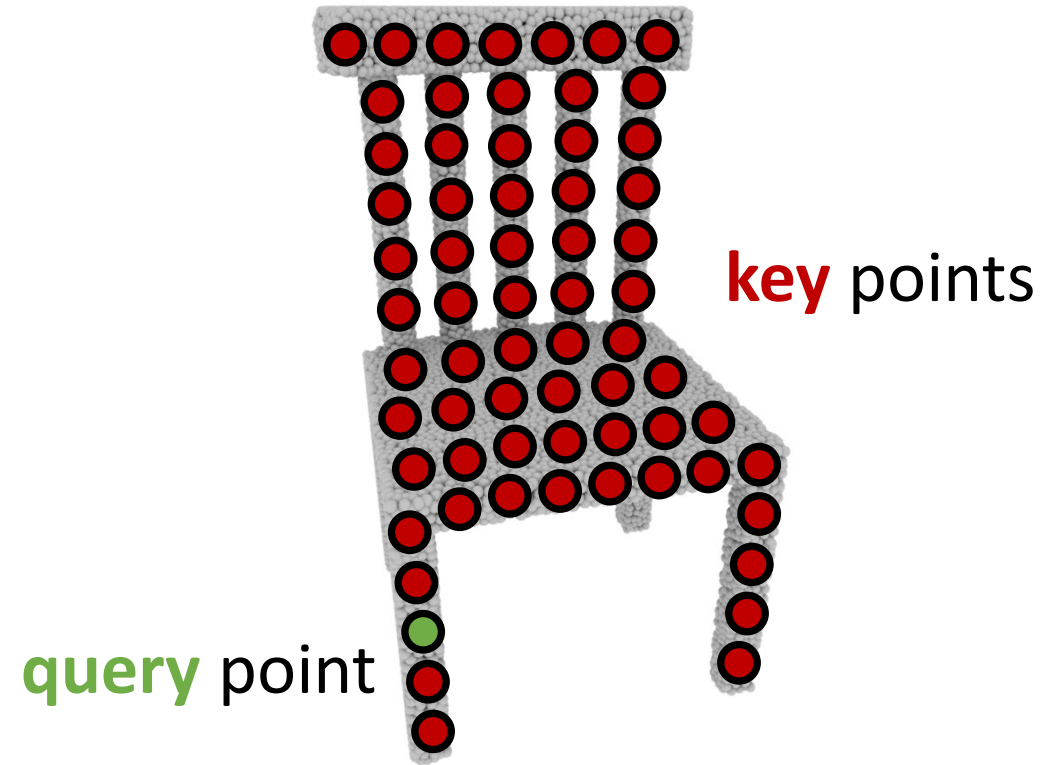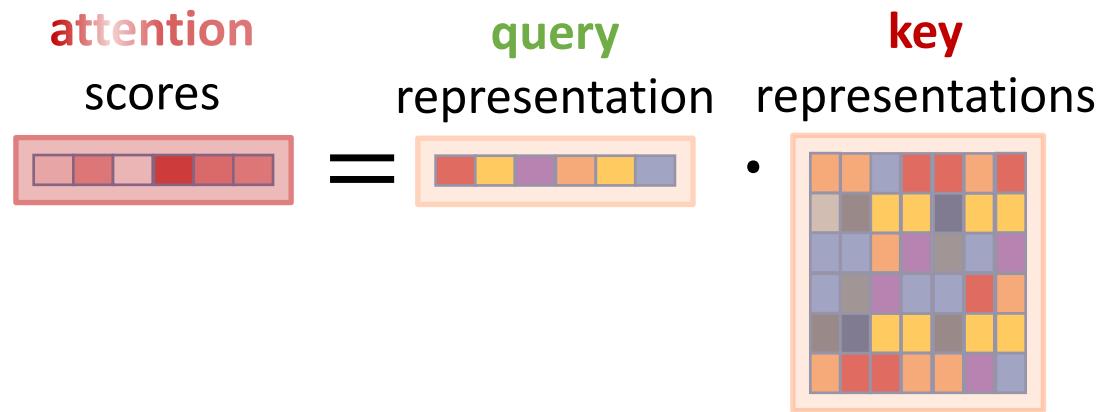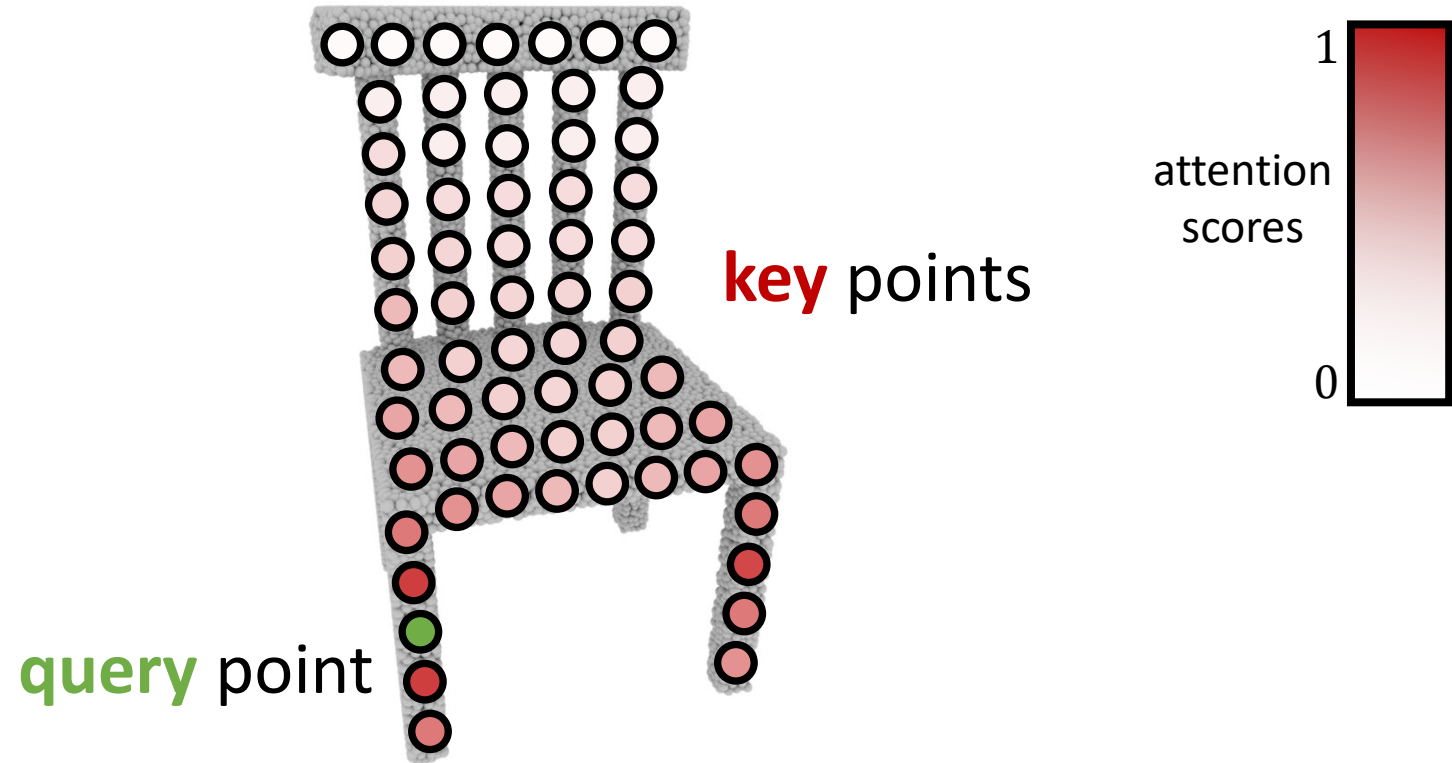**key** points

**query** point

# Why use **attention** for 3D representations?

Encode points such that their features capture **relations** wrt the rest of the shape

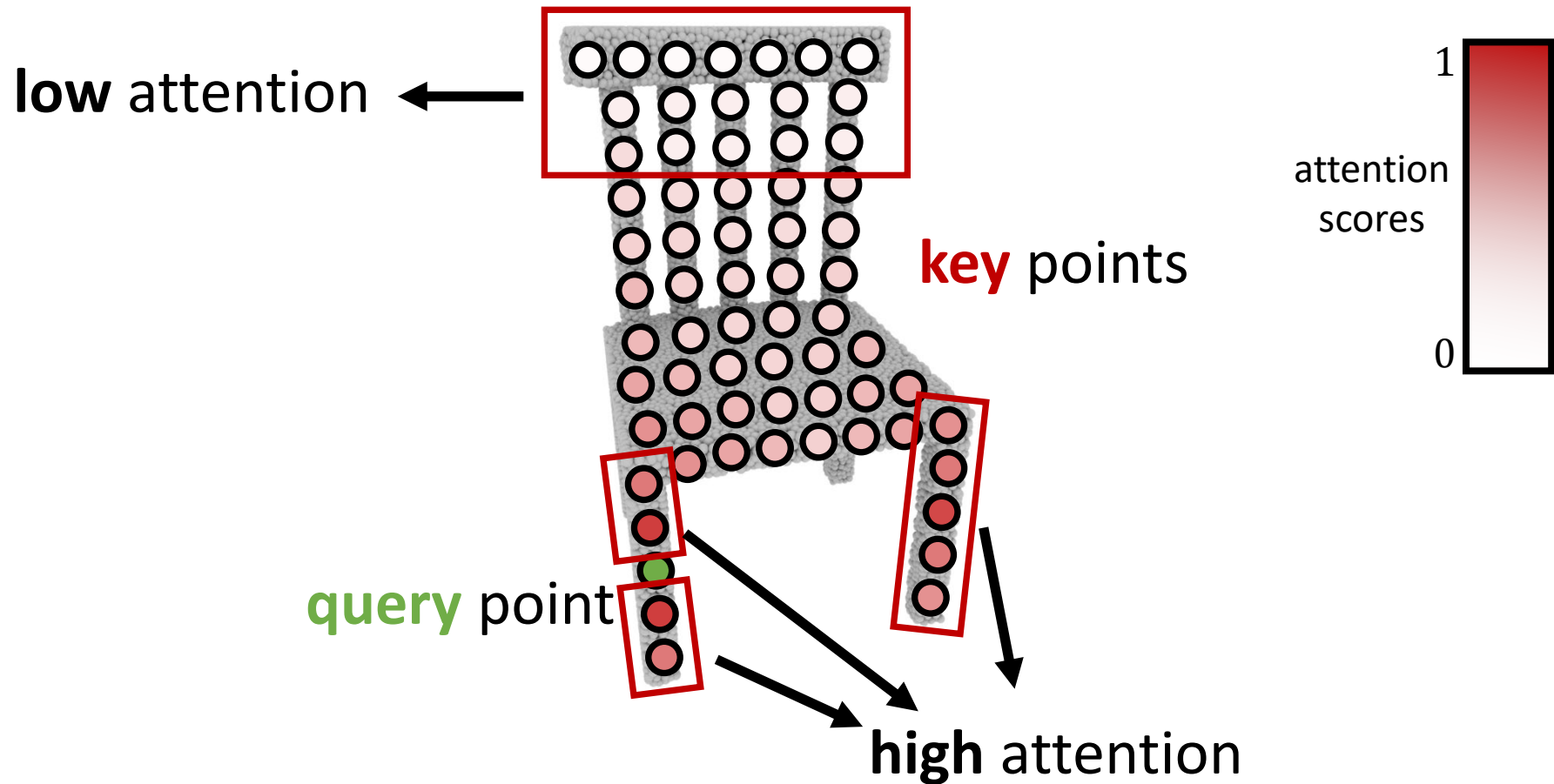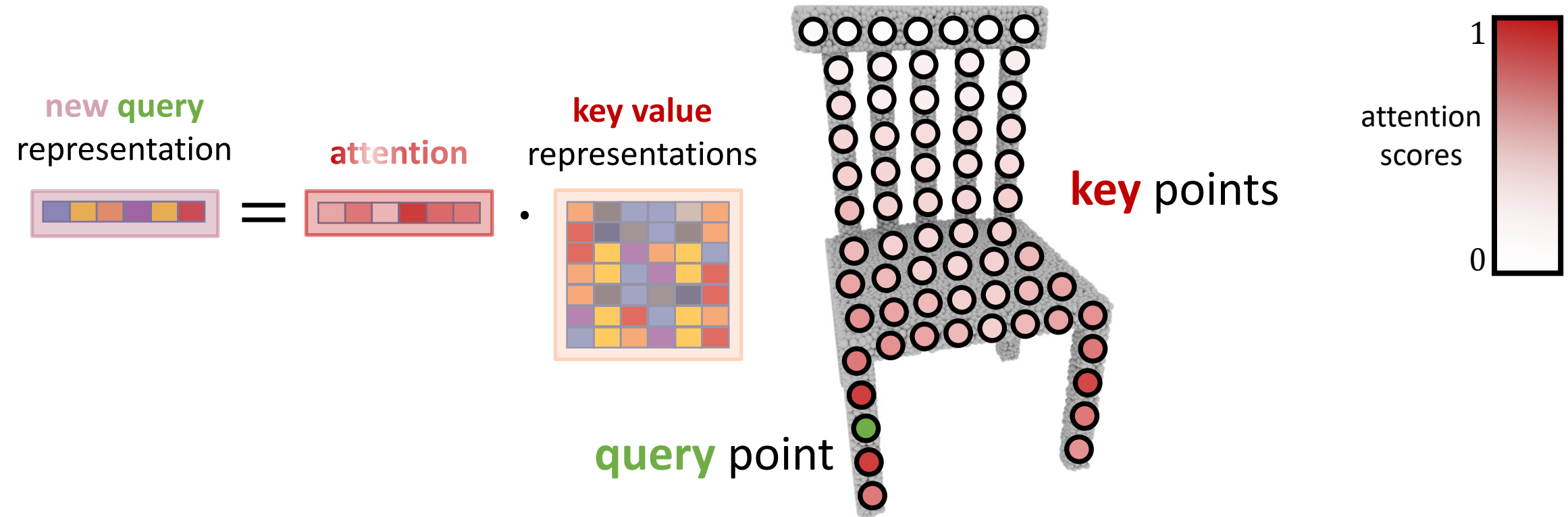# Why use **attention** for 3D representations?
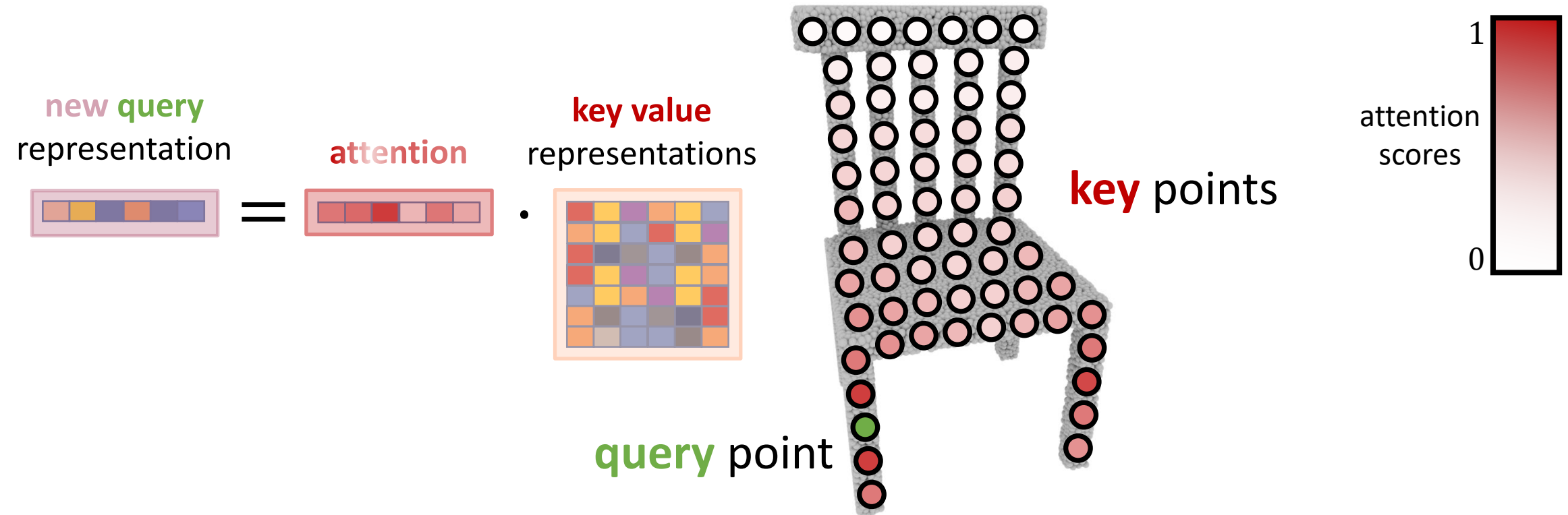
# Why use **attention** for 3D representations?
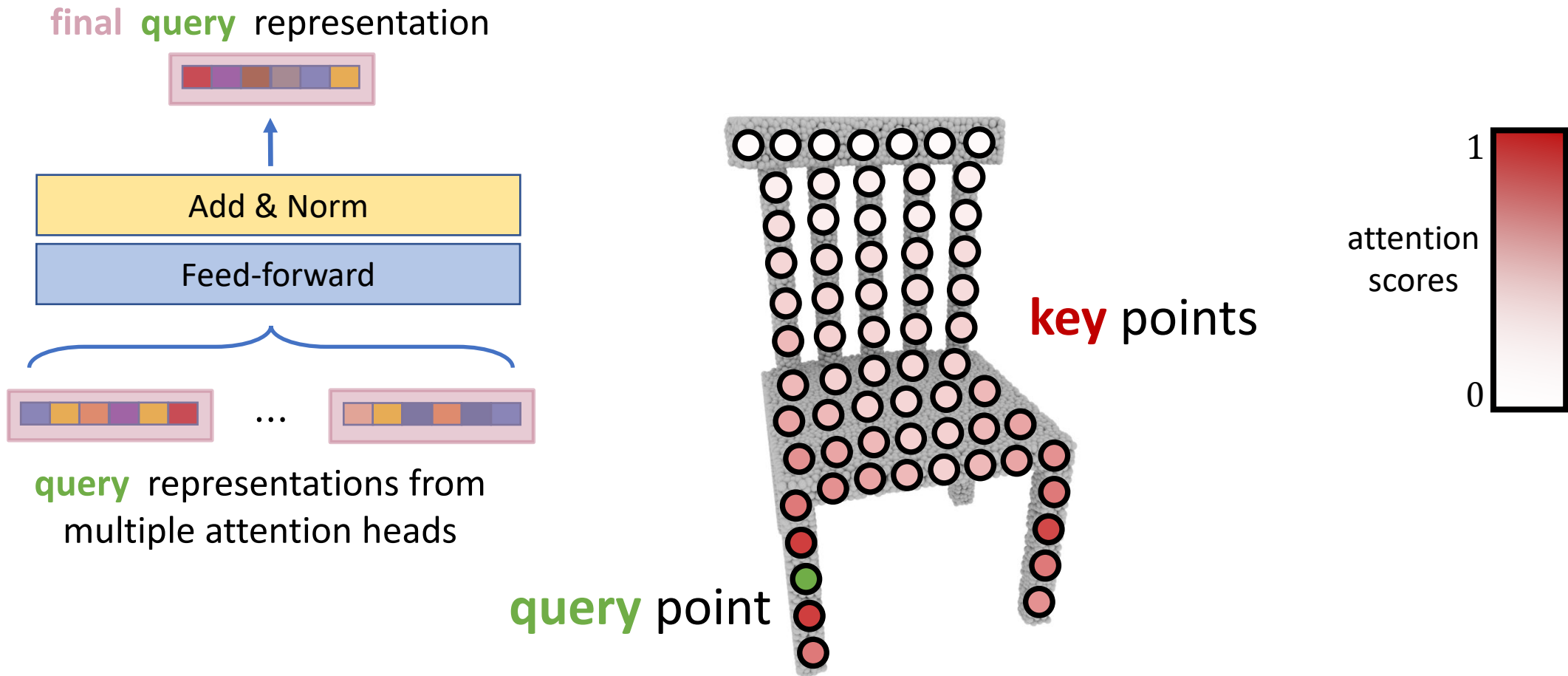


low attention

**key** points

attention scores

**query** point

**high** attention

# Why use **attention** for 3D representations?

# Why use **attention** for 3D representations?

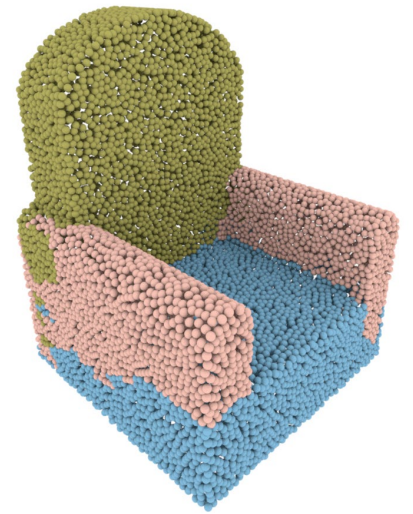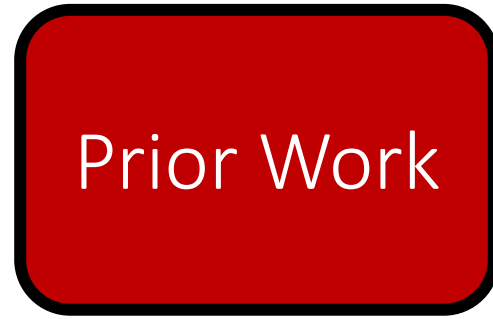# Why use **attention** for 3D representations?

# Motivation: Long-range interactions **across** shapes



test shape

**Prior Work**

output part labels

**No interactions across shapes**

# Motivation: Long-range interactions across shapes



test shape

Cross-ShapeNet

output part labels

shapes from input collection
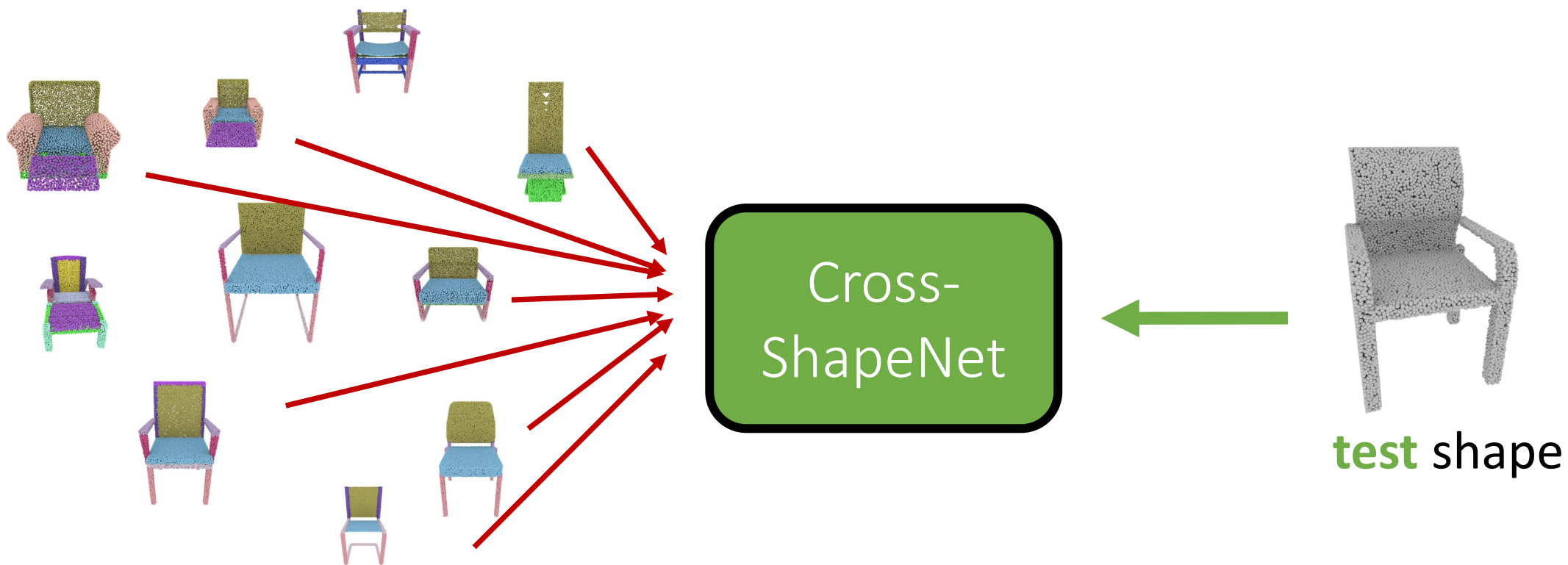
# Key challenge: Retrieve compatible shapes



**Shape Collection**

**test** shape

# Key challenge: Retrieve compatible shapes



**Shape Collection**
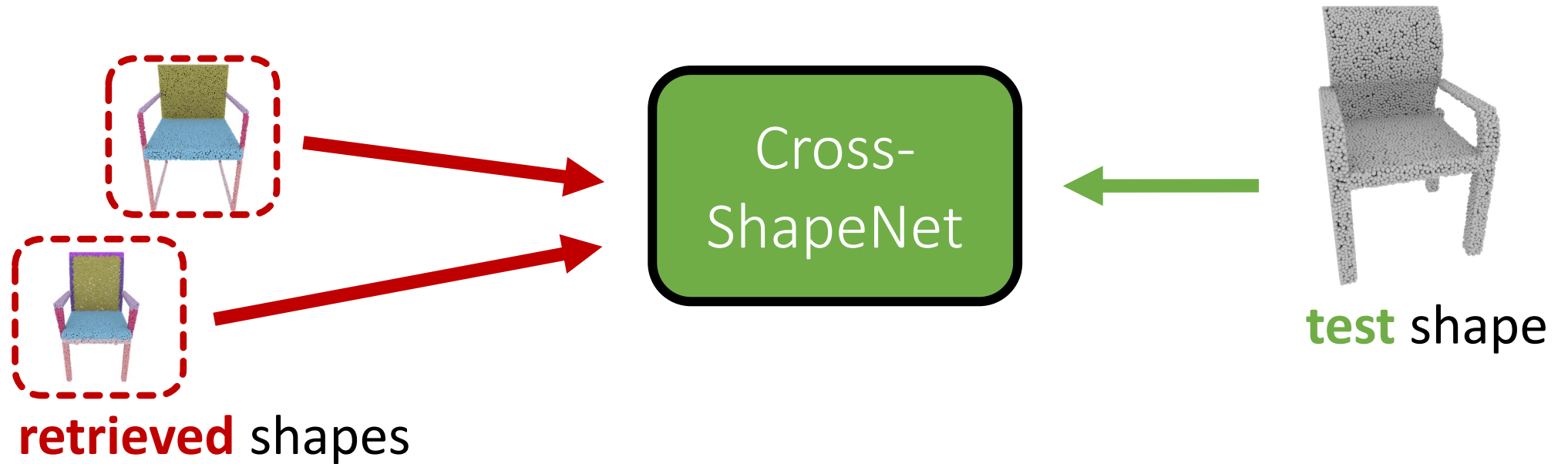
Cross-ShapeNet

**test** shape

# Key challenge: Retrieve compatible shapes



**Shape Collection**

Cross-ShapeNet

**test** shape

# Key challenge: Combine multiple shapes



**retrieved** shapes

Cross-ShapeNet

**test** shape

# Key challenge: Combine multiple shapes



**retrieved** shapes

**test** shape

**Cross-shape attention**

# Key challenge: Combine multiple shapes

# Pipeline



**Shape Collection**

# Pipeline



**Shape Collection**

# Pipeline



**Shape Collection**

# Pipeline



**Shape Collection**

**Cross-shape attention**

# Pipeline



**Shape Collection**

**Cross-shape attention**

# Pipeline



**Shape Collection**

**Cross-shape attention**

# Pipeline



**Shape Collection**

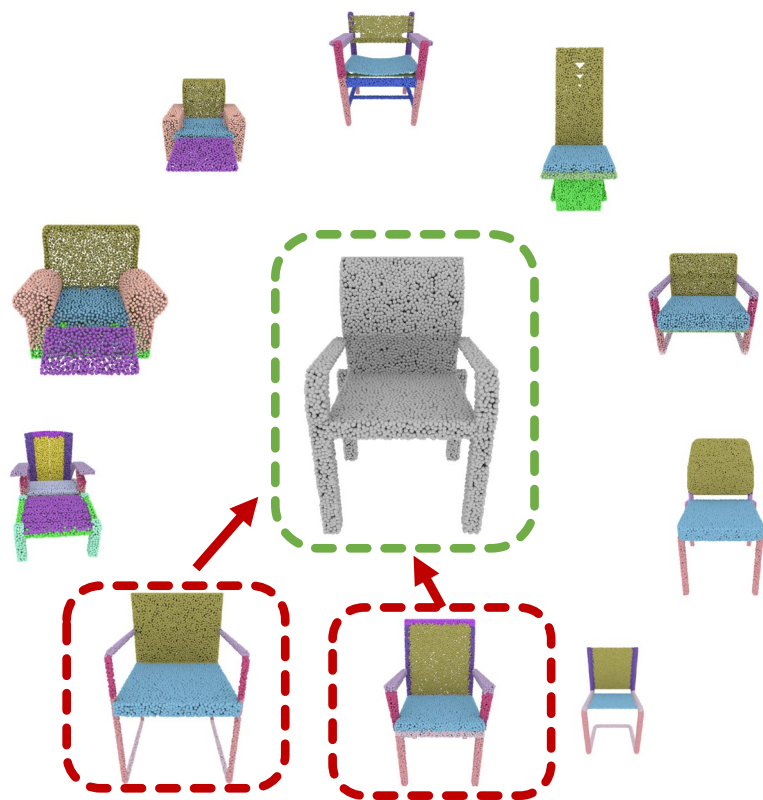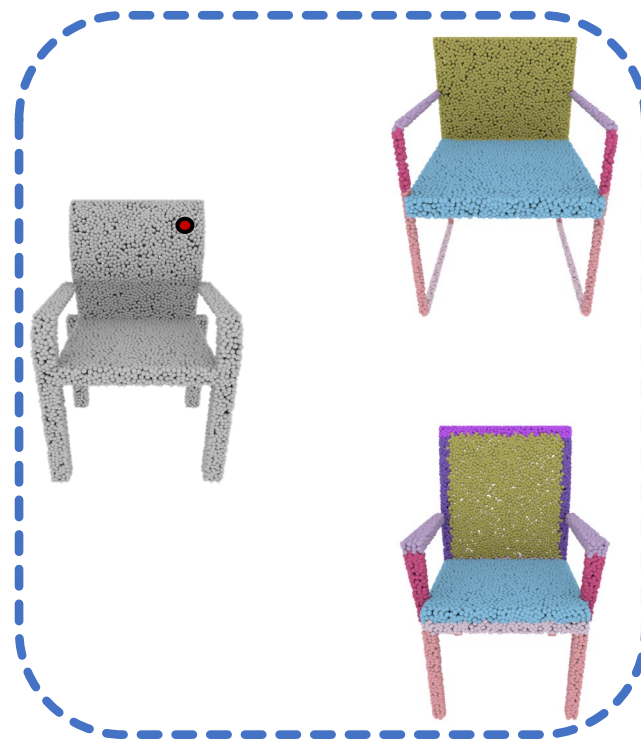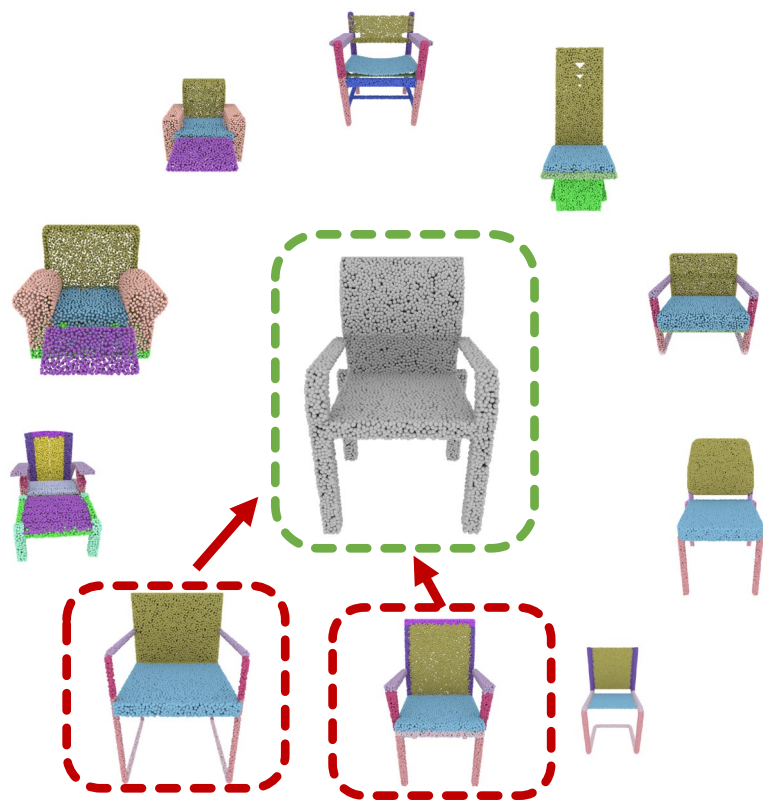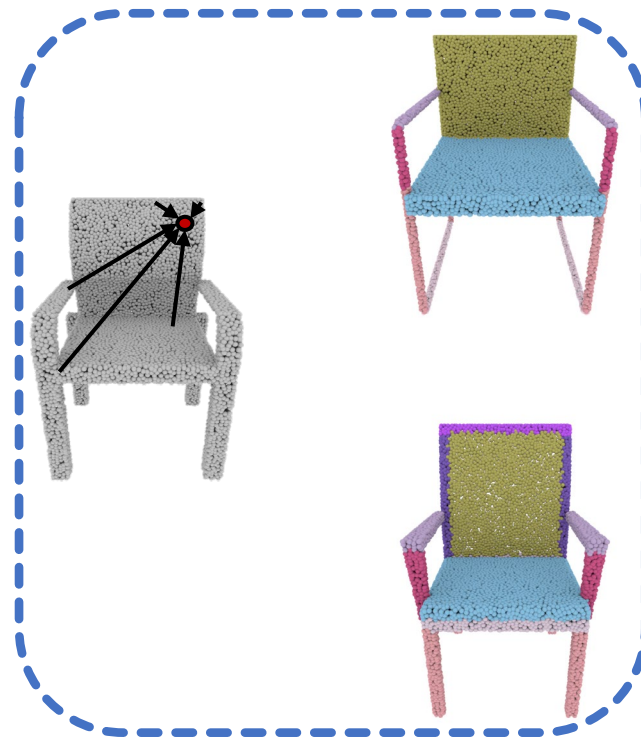**Cross-shape attention**

# Pipeline



**Shape Collection**

**Cross-shape attention**

Cross-ShapeNet

# Pipeline



**Shape Collection**

**Cross-shape attention**

Cross-ShapeNet

# Cross-Shape Attention

**query** shape $\mathcal{S}_m = \{\boldsymbol{p}_i\}_{i=1}^M$



**key** shape $\mathcal{S}_n = \{\boldsymbol{p}_j\}_{j=1}^N$

# Cross-Shape Attention

**query** shape $\mathcal{S}_m = \{\boldsymbol{p}_i\}_{i=1}^{M}$

$\boldsymbol{X}_m \in R^{M \times D}$



**Backbone**

Backbone point representations

**Backbone**

**key** shape $\mathcal{S}_n = \{\boldsymbol{p}_j\}_{j=1}^{N}$

$\boldsymbol{X}_n \in R^{N \times D}$

# Cross-Shape Attention

$$\boldsymbol{X}_m \in R^{M \times D}$$



Backbone point
representations



$$\boldsymbol{X}_n \in R^{N \times D}$$

Transformer [Vaswani et al. 2017]

# Cross-Shape Attention

$$\boldsymbol{X}_m \in R^{M \times D} \qquad \boldsymbol{W}_Q \in R^{D \times D}$$

$$\boldsymbol{Q}_m = \boldsymbol{W}_Q \cdot \boldsymbol{X}_m$$



Backbone point representations

Query Transformation

Intermediate representations

Key Transformation

$$\boldsymbol{X}_n \in R^{N \times D} \qquad \boldsymbol{W}_K \in R^{D \times D} \qquad \boldsymbol{K}_n = \boldsymbol{W}_K \cdot \boldsymbol{X}_n$$

Transformer [Vaswani et al. 2017]

# Cross-Shape Attention

$$\boldsymbol{W}_Q \in R^{D \times D}$$

$$\boldsymbol{X}_m \in R^{M \times D}$$

$$\boldsymbol{Q}_m \in R^{M \times D}$$

Query Transformation

Backbone point
representations

Intermediate
representations

Key Transformation

$$\boldsymbol{X}_n \in R^{N \times D}$$

$$\boldsymbol{W}_K \in R^{D \times D}$$

$$\boldsymbol{K}_n \in R^{N \times D}$$

Transformer [Vaswani et al. 2017]

# Cross-Shape Attention



$\boldsymbol{X}_m \in R^{M \times D}$

$\boldsymbol{W}_Q \in R^{D \times D}$

$\boldsymbol{Q}_m \in R^{M \times D}$

Query Transformation

Backbone point representations

Intermediate representations

$\boldsymbol{W}_V \in R^{D \times D}$

Key-value representations

Key Transformation

$\boldsymbol{X}_n \in R^{N \times D}$

$\boldsymbol{W}_K \in R^{D \times D}$

$\boldsymbol{K}_n \in R^{N \times D}$

Value Transformation

$\boldsymbol{V}_n = \boldsymbol{W}_V \cdot \boldsymbol{X}_n$

Transformer [Vaswani et al. 2017]

# Cross-Shape Attention



$\boldsymbol{X}_m \in R^{M \times D}$    $\boldsymbol{W}_Q \in R^{D \times D}$    $\boldsymbol{Q}_m \in R^{M \times D}$

Query Transformation

Backbone point representations

Intermediate representations

Key-value representations

$\boldsymbol{X}_n \in R^{N \times D}$    $\boldsymbol{W}_K \in R^{D \times D}$    $\boldsymbol{K}_n \in R^{N \times D}$

Key Transformation

$\boldsymbol{W}_V \in R^{D \times D}$

Value Transformation

$\boldsymbol{V}_n \in R^{N \times D}$

Transformer [Vaswani et al. 2017]

# Cross-Shape Attention

Query
representations

$Q_m$

Key
representations

$\mathbf{K}_n^T$

Transformer [Vaswani et al. 2017]

# Cross-Shape Attention

Query
representations

Key
representations

$$\frac{\boldsymbol{Q}_m \quad \mathbf{K}_n^T}{\sqrt{D}}$$

$\boldsymbol{Q}_m \in R^{M \times D}$

$\boldsymbol{K}_n \in R^{N \times D}$

Transformer [Vaswani et al. 2017]

# Cross-Shape Attention

Query
representations

Key
representations

$$\boldsymbol{Q}_m \qquad \mathbf{K}_n^T$$



$$\boldsymbol{Q}_m \in R^{M \times D}$$
$$\boldsymbol{K}_n \in R^{N \times D}$$

$$\frac{\boxed{\boldsymbol{Q}_m} \quad \boxed{\mathbf{K}_n^T}}{\sqrt{D}}$$

$$Var\left(\frac{\boldsymbol{Q}_{i,:}\boldsymbol{K}_{:,j}^T}{\sqrt{D}}\right) = 1,$$
$$\forall i = 1, \cdots, M$$
$$\forall j = 1, \cdots, N$$

# Cross-Shape Attention

$$softmax \left( \frac{Q_m \quad K_n^T}{\sqrt{D}} \right) = A_{m,n} \in R^{M \times N}$$

Query representations $Q_m$

Key representations $K_n^T$

Attention matrix

Transformer [Vaswani et al. 2017]

# Cross-Shape Attention

$$\boldsymbol{A}_{m,n} \cdot \boldsymbol{V}_n = \boldsymbol{X'}_m^{(CSA)} \in R^{M \times D}$$

**Cross-shape** attention matrix

**Key shape** value representations

**Cross-shape** attention representations

# Cross-Shape Attention

$$\boldsymbol{A}_{m,n} \cdot \boldsymbol{V}_n = \boldsymbol{X'}_m^{(CSA)} \in R^{M \times D}$$

**Cross-shape** attention matrix

**Key shape** value representations

**Cross-shape** attention representations

**query** shape

**key** shape

# Cross-Shape Attention for multiple shapes



**key** shape

**query** shape

# Cross-Shape Attention for multiple shapes



**query** shape

**key** shapes $\mathcal{C}(m)$

# Cross-Shape Attention for multiple shapes



**query** shape

**key** shapes $\mathcal{C}(m)$

# Cross-Shape Attention for multiple shapes



**query** shape

**key** shapes $\mathcal{C}(m)$

- $\mathcal{C}(m)$: set of compatible key shapes
- $c(m, n)$: compatibility function between query shape $S_m$ and key shape $S_n$

**Cross-shape attention output**

$$X'_m = \sum_{n \in \{\mathcal{C}(m), m\}} c(m, n) A_{m,n} V_n$$

# Compatibility function

$$\boldsymbol{X'}_m^{(SSA)} \in R^{M \times D}$$



$$\boldsymbol{X'}_n^{(SSA)} \in R^{N \times D}$$

# Compatibility function

$$\boldsymbol{X'}_m^{(SSA)} \in R^{M \times D}$$



$$\underset{i}{\text{avg}} \, \boldsymbol{X'}_{m,i}^{(SSA)}$$

$$\boldsymbol{y}_m^{(SSA)} \in R^D$$





$$\underset{i}{\text{avg}} \, \boldsymbol{X'}_{n,i}^{(SSA)}$$

$$\boldsymbol{y}_n^{(SSA)} \in R^D$$



$$\boldsymbol{X'}_n^{(SSA)} \in R^{N \times D}$$

# Compatibility function

# Compatibility function



$\boldsymbol{u}_m \in R^D$

$\widehat{\boldsymbol{u}}_m = \boldsymbol{u}_m / ||\boldsymbol{u}_m||$

**Cosine similarity**

$$s(m,n) = \widehat{\boldsymbol{u}}_m \cdot \widehat{\boldsymbol{u}}_n$$

$\boldsymbol{u}_n \in R^D$

$\widehat{\boldsymbol{u}}_n = \boldsymbol{u}_n / ||\boldsymbol{u}_n||$

# Compatibility function

# Compatibility function

Compatibility function

query shape

$\hat{\boldsymbol{u}}_m$

$\hat{\boldsymbol{u}}_{n_1}$

$\hat{\boldsymbol{u}}_{n_k}$

$s(m, n_1)$

$s(m, n_k)$

$\boldsymbol{s(m, m\,)}$

key shapes $\mathcal{C}(m)$

# Compatibility function



query shape

$\widehat{\boldsymbol{u}}_m$

$\widehat{\boldsymbol{u}}_{n_1}$

$\widehat{\boldsymbol{u}}_{n_k}$

$c(m, n_1)$

$c(m, n_k)$

$c(\boldsymbol{m}, \boldsymbol{m})$

key shapes $\mathcal{C}(m)$

compatibility
$$c(m, n) = \frac{e^{s(m,n)}}{\sum_{n \in \{\mathcal{C}(m), m\}} e^{s(m,n)}}$$

# Cross-Shape Attention for multiple shapes



**key** shapes

**Cross-shape attention**

# Cross-Shape Attention for multiple shapes

# Retrieve compatible shapes



**Shape Collection**

Cross-ShapeNet

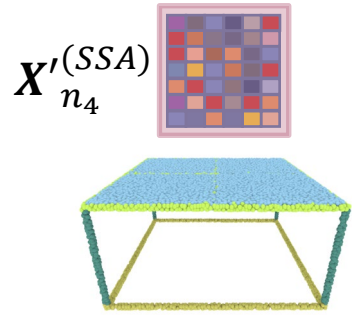**query** shape
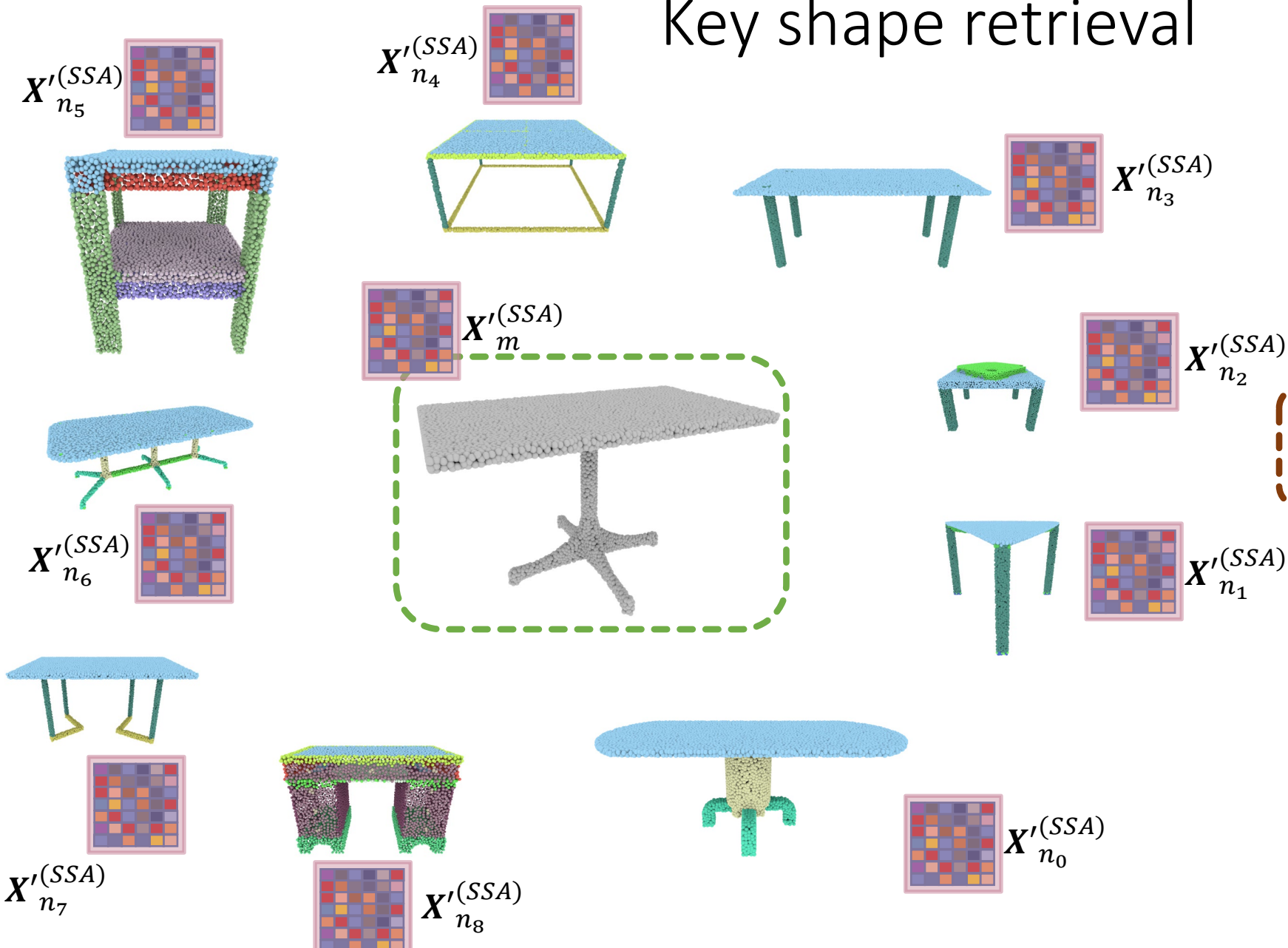
# Retrieve compatible shapes



**Shape Collection**

Cross-ShapeNet

**query** shape

Key shape retrieval

Key shape retrieval

$\boldsymbol{X'}_{n_5}^{(SSA)}$

$\boldsymbol{X'}_{n_4}^{(SSA)}$
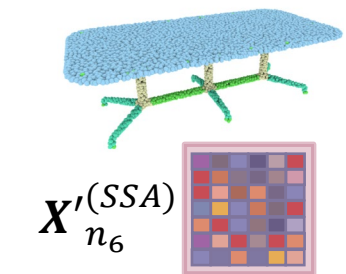
$\boldsymbol{X'}_{n_3}^{(SSA)}$

$\boldsymbol{X'}_{n_2}^{(SSA)}$

$\boldsymbol{X'}_{n_6}^{(SSA)}$

$\boldsymbol{X'}_{n_1}^{(SSA)}$

$\boldsymbol{X'}_{n_7}^{(SSA)}$

$\boldsymbol{X'}_{n_8}^{(SSA)}$

$\boldsymbol{X'}_{n_0}^{(SSA)}$

Key shape retrieval

$X'^{(SSA)}_{n_5}$

$X'^{(SSA)}_{n_4}$

$X'^{(SSA)}_{n_3}$

$X'^{(SSA)}_{m}$

$X'^{(SSA)}_{n_2}$

$X'^{(SSA)}_{n_6}$

$X'^{(SSA)}_{n_1}$

$X'^{(SSA)}_{n_7}$

$X'^{(SSA)}_{n_8}$

$X'^{(SSA)}_{n_0}$

Key shape retrieval

$X'^{(SSA)}_{n_5}$

$X'^{(SSA)}_{n_4}$

$X'^{(SSA)}_{n_3}$

$X'^{(SSA)}_{m}$

$X'^{(SSA)}_{n_2}$

$X'^{(SSA)}_{n_6}$

$X'^{(SSA)}_{n_1}$

$X'^{(SSA)}_{n_7}$

$X'^{(SSA)}_{n_8}$

$X'^{(SSA)}_{n_0}$

**Cosine similarity**

$$S_{m,n_k} = X'^{(SSA)}_m \cdot \left( X'^{(SSA)}_{n_k} \right)^{\top}$$

# Key shape retrieval



$$r_i(m, n_k) = \max_j S_{m,n_k}[i, j]$$

# Key shape retrieval



$X'^{(SSA)}_{n_5}$

$X'^{(SSA)}_{n_4}$

$X'^{(SSA)}_{n_3}$

$X'^{(SSA)}_{m}$

$X'^{(SSA)}_{n_2}$

$X'^{(SSA)}_{n_6}$

$X'^{(SSA)}_{n_1}$

$X'^{(SSA)}_{n_7}$

$X'^{(SSA)}_{n_8}$

$X'^{(SSA)}_{n_0}$

**retrieval measure**

$$r(m, n_k) = \operatorname*{avg}_{i} r_i(m, n_k)$$

# Key shape retrieval

$X'^{(SSA)}_{n_5}$

$X'^{(SSA)}_{n_4}$

$X'^{(SSA)}_{n_3}$

$X'^{(SSA)}_m$

$X'^{(SSA)}_{n_2}$

$X'^{(SSA)}_{n_6}$

$X'^{(SSA)}_{n_1}$

$X'^{(SSA)}_{n_7}$

$X'^{(SSA)}_{n_8}$

$X'^{(SSA)}_{n_0}$

**retrieval measure**

$$r(m, n_k) = \underset{i}{\mathrm{avg}}\, r_i(m, n_k)$$

# Key shape retrieval

$X'^{(SSA)}_{n_5}$

$X'^{(SSA)}_{n_4}$

$X'^{(SSA)}_{n_3}$

$X'^{(SSA)}_{m}$

$X'^{(SSA)}_{n_2}$

**retrieval measure**

$$r(m, n_k) = \underset{i}{\mathrm{avg}}\, r_i(m, n_k)$$

$X'^{(SSA)}_{n_6}$

$X'^{(SSA)}_{n_1}$

$X'^{(SSA)}_{n_7}$

$X'^{(SSA)}_{n_8}$

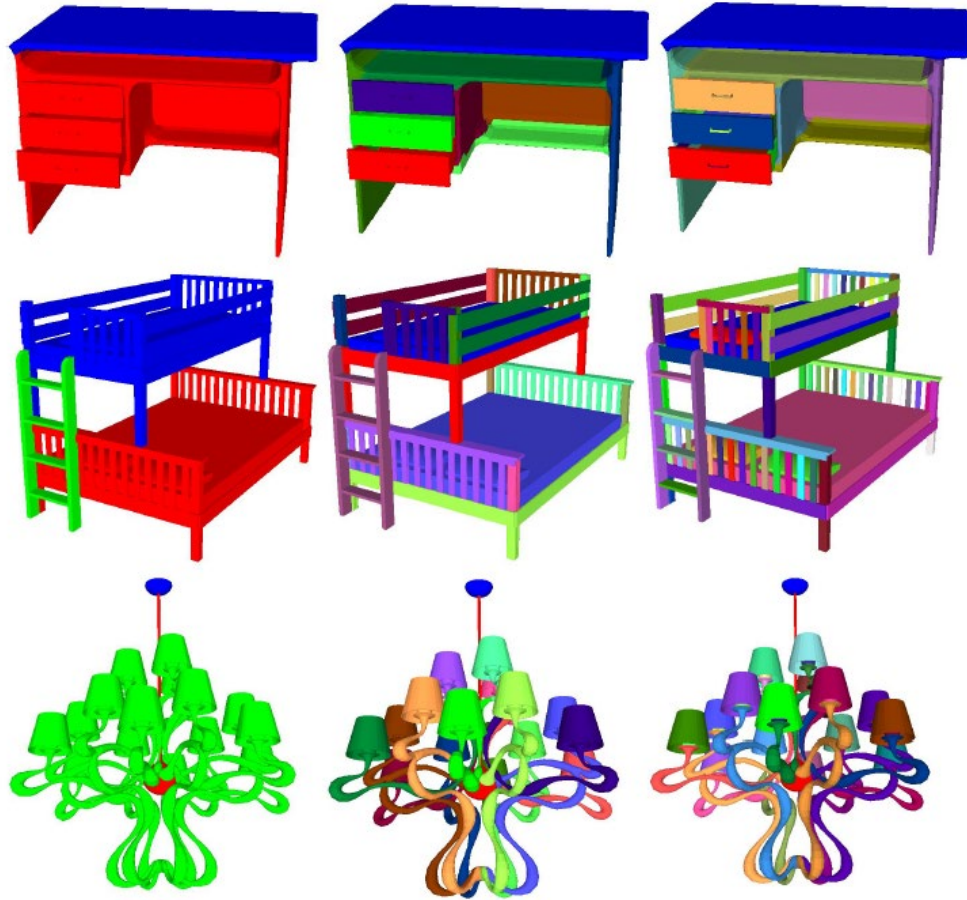$X'^{(SSA)}_{n_0}$

# Key shape retrieval: Examples
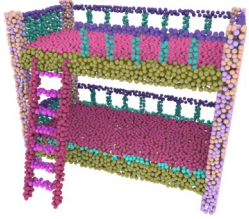
**query** shapes
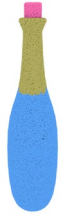


**key** shapes

# PartNet dataset
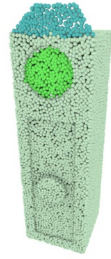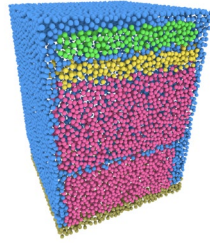

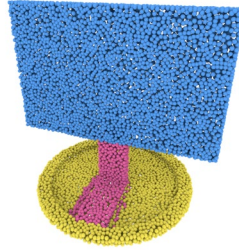
Coarse → Fine-grained

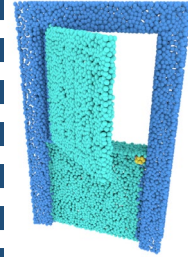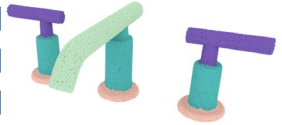[Mo et al. 2019]
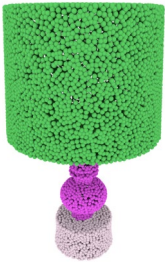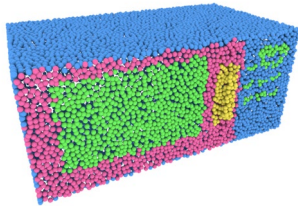
# PartNet dataset



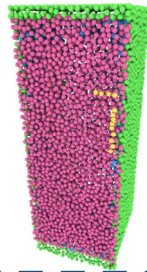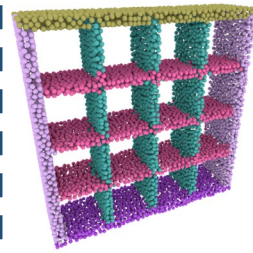Bed  Bottle  Chair  Clock  Dishwasher  Display  Door  Earphone  Faucet

Knife  Lamp  Microwave  Refrigerator  Storage Furn.  Table  Trashcan  Vase

[Mo et al. 2019]

# Examples of shape collections



Table category

Chair category

**5,707** training shapes

**4,489** training shapes

# Training details: Loss

$$L_{CE} = -\sum_{\boldsymbol{p}_i \in \mathcal{S}_k} \widehat{\boldsymbol{q}}_i \log \boldsymbol{q}_i$$

$\mathcal{S}_k$: shape $k = \{\boldsymbol{p}_i\}_{i=1}^{P_k}$
$\widehat{\boldsymbol{q}}_i$: ground-truth one-hot label vector for point $\boldsymbol{p}_i$
$\boldsymbol{q}_i$: predicted label probabilities for point $\boldsymbol{p}_i$

training
data

# Training details: Backbones



MID-FC [Wang et al. 2021]

MinkowskiNet [Choy et al. 2019]

# Training details: Backbones



HRNet [Wang et al. 2021]

**Shape Collection**

**Collection graph**

**Collection graph**

**Collection graph**

# Inference: Collection graph

**Collection graph**

**test** shape

# Inference: Collection graph

**Collection graph**

**test** shape

# Results

| Method | Part IoU |
|--------|----------|

# Results: MinkowskiNet variants

| Method | Part IoU |
|--------|----------|
| MinkHRNet | 48.0 |

# Results: MinkowskiNet variants

| Method | Part IoU |
|--------|----------|
| MinkHRNet | 48.0 |
| MinkHRNetCSN-SSA | 48.7 |

+0.7%

# Results: MinkowskiNet variants

| Method | Part IoU |
|---|---|
| MinkHRNet | 48.0 |
| MinkHRNetCSN-SSA | 48.7 |
| MinkHRNetCSN-K1 | **49.9** |
| MinkHRNetCSN-K2 | 49.7 |

## Results: MinkowskiNet variants

| Method | Part IoU |
|---|---|
| MinkHRNet | 48.0 |
| MinkHRNetCSN-SSA | 48.7 |
| **MinkHRNetCSN-K1** | **49.9** |
| MinkHRNetCSN-K2 | 49.7 |

+1.2%

# Results: MinkowskiNet variants

**Ground truth**



**MinkHRNet**

- ■ Pillow
- ■ Mattress
- ■ Stretcher
- ■ Leg
- ■ Horizontal bar
- ■ Vertical bar
- ■ Bed post
- ■ Ladder vertical bar
- ■ Rung

# Results: MinkowskiNet variants

Ground truth

MinkHRNet

**MinkHRNetCSN-SSA**

Pillow
Mattress
Stretcher
Leg
Horizontal bar
Vertical bar
Bed post
Ladder vertical bar
Rung

# Results: MinkowskiNet variants

**Ground truth**

**MinkHRNetCSN-SSA**

**MinkHRNetCSN-K1**

Pillow
Mattress
Stretcher
Leg
Horizontal bar
Vertical bar
Bed post
Ladder vertical bar
Rung

# Results: MID-FC variants

| Method | Part IoU |
|--------|----------|
| MID-FC | 60.8 |

# Results: MID-FC variants

| Method | Part IoU |
|:---:|:---:|
| MID-FC | 60.8 |
| MID-FC-CSN-SSA | 61.8 |

+1.0%

# Results: MID-FC variants

| Method | Part IoU |
|---|---|
| MID-FC | 60.8 |
| MID-FC-CSN-SSA | 61.8 |
| MID-FC-CSN-K1 | 61.9 |
| MID-FC-CSN-K2 | 61.9 |
| MID-FC-CSN-K3 | 62.0 |
| MID-FC-CSN-K4 | **62.1** |
| MID-FC-CSN-K5 | 62.0 |

+0.3%

# Ground truth

# Results: MID-FC variants

## MID-FC

Bar
Leg
Board
Shelf

Ground truth

Results: MID-FC variants

MID-FC

**MID-FC-CSN-SSA**

Bar
Leg
Board
Shelf

# Ground truth

# Results: MID-FC variants

# MID-FC-CSN-SSA

# **MID-FC-CSN-K4**

Bar
Leg
Board
Shelf

# Results: Comparison with other methods

| Method | Part IoU |
|---|---|
| ResGCN-28 (Li et al. 2023) | 45.1 |
| CloserLook3D (Liu et al. 2020) | 53.8 |
| MinkResUNet (Choy et al. 2019) | 46.8 |
| MinkHRNetCSN-K1 (ours) | 49.9 |
| MID-FC (Wang et al. 2021) | 60.8 |
| MID-FC-CSN-K4 (ours) | **62.1** |

# Results: Comparison with other methods

| Method | Part IoU |
|---|---|
| ResGCN-28 (Li et al. 2023) | 45.1 |
| CloserLook3D (Liu et al. 2020) | 53.8 |
| MinkResUNet (Choy et al. 2019) | 46.8 |
| MinkHRNetCSN-K1 (ours) | 49.9 |
| MID-FC (Wang et al. 2021) | 60.8 |
| MID-FC-CSN-K4 (ours) | **62.1** |

+3.1%

# Results: Comparison with other methods

| Method | Part IoU |
|---|---|
| ResGCN-28 (Li et al. 2023) | 45.1 |
| CloserLook3D (Liu et al. 2020) | 53.8 |
| MinkResUNet (Choy et al. 2019) | 46.8 |
| MinkHRNetCSN-K1 (ours) | 49.9 |
| MID-FC (Wang et al. 2021) | 60.8 |
| MID-FC-CSN-K4 (ours) | **62.1** |

+1.3%

# Results: Comparison with other methods

| Method | Part IoU |
|---|---|
| ResGCN-28 (Li et al. 2023) | 45.1 |
| CloserLook3D (Liu et al. 2020) | 53.8 |
| MinkResUNet (Choy et al. 2019) | 46.8 |
| MinkHRNetCSN-K1 (ours) | 49.9 |
| MID-FC (Wang et al. 2021) | 60.8 |
| MID-FC-CSN-K4 (ours) | **62.1** |

+1.3%

**SOTA performance on the PartNet dataset**

# Summary



Cross-shape convolution

Cross-ShapeNet

Shape Collection

- Enable long range point feature interactions **across shapes**

# Summary



Shape Collection

Cross-shape convolution

Cross-ShapeNet

- Enable long range point feature interactions **across shapes**
- Introduce a **novel cross-shape attention** mechanism

# Summary



Cross-shape convolution

Cross-ShapeNet

Shape Collection

- Enable long range point feature interactions **across shapes**

- Introduce a **novel cross-shape attention** mechanism

- Retrieve **compatible shapes** for cross-shape attention

# Summary



Shape Collection

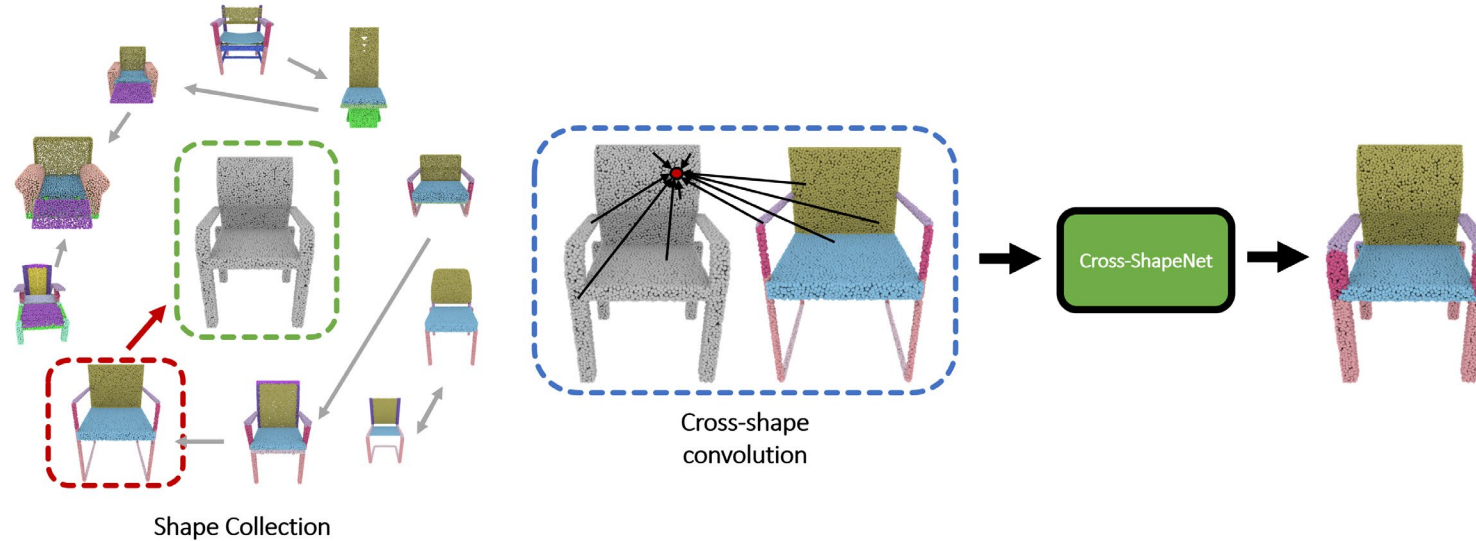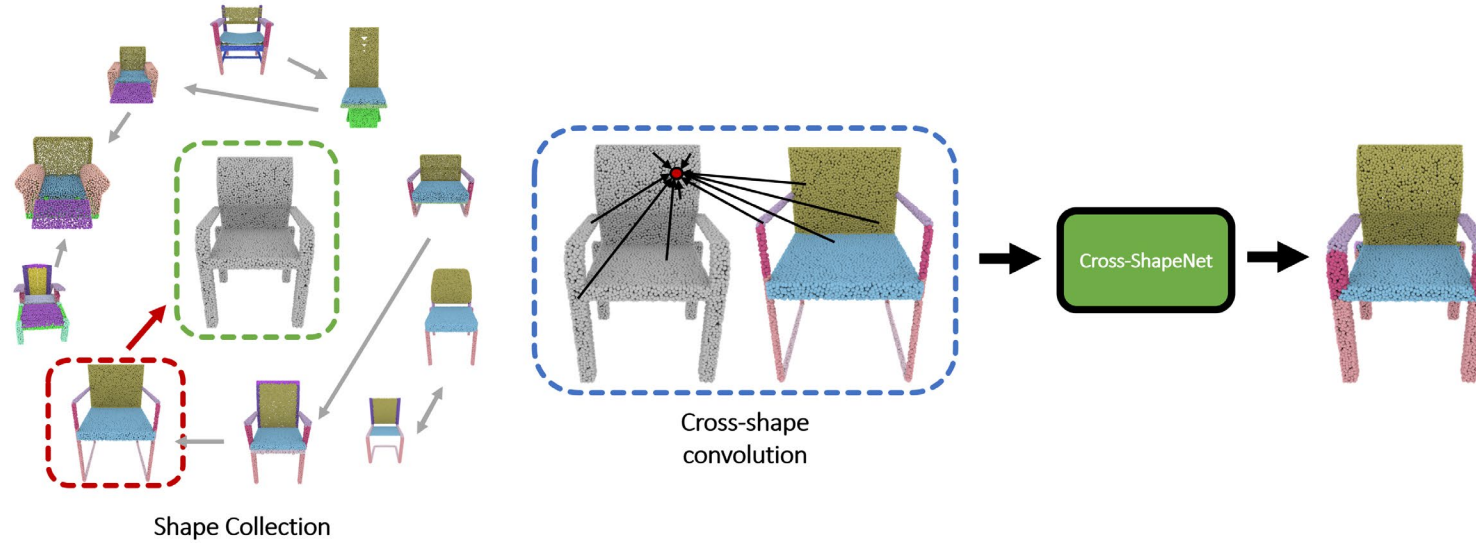Cross-shape convolution

Cross-ShapeNet

- Enable long range point feature interactions **across shapes**
- Introduce a **novel cross-shape attention** mechanism
- Retrieve **compatible shapes** for cross-shape attention
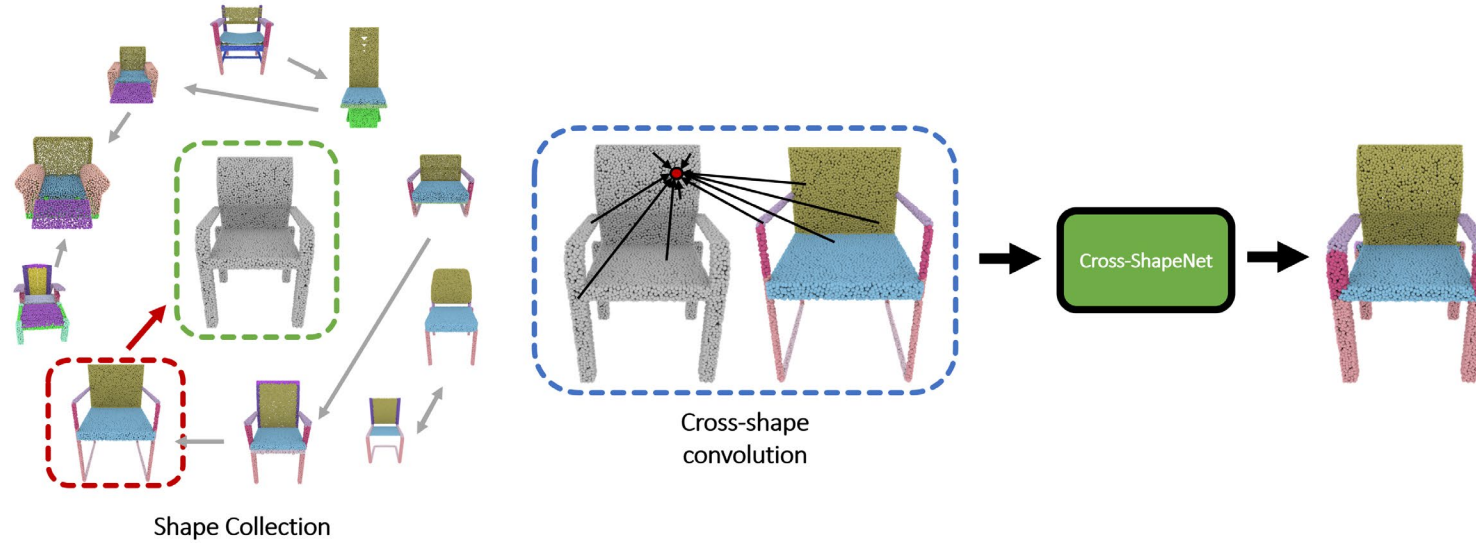- **SOTA performance** on PartNet

# Summary



Shape Collection

Cross-shape convolution

Cross-ShapeNet

Limitations:

- **Increased computational cost** due to shape retrieval

# Summary



Cross-shape convolution

Cross-ShapeNet

Shape Collection

Limitations:

- **Increased computational cost** due to shape retrieval
- Currently no support for **multi-object scenes**

# Thank you!

## Acknowledgements:

**Our project web page:**
https://marios2019.github.io/CSN/