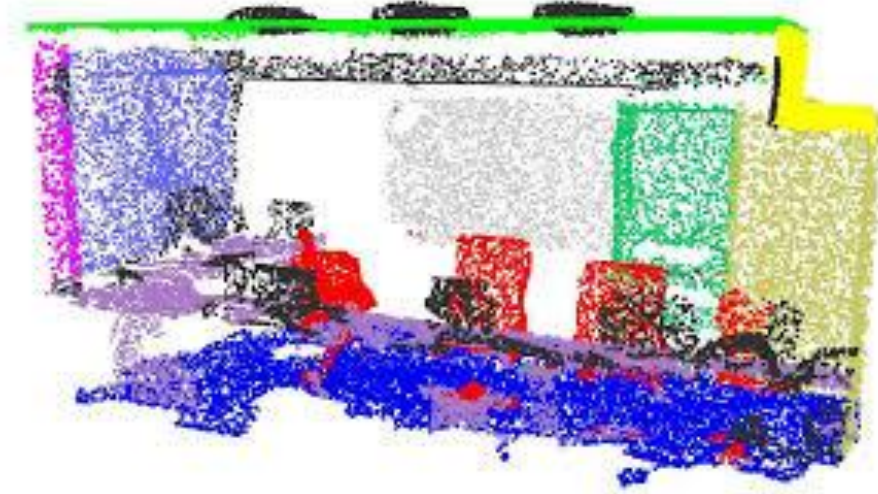


# Set Models



Markus Chardonnet  
*SiDNN* – 29.03.2022

# Outline

- Sets and Multisets
- Motivations for set models
- Setup
  - Permutation Invariance
  - Permutation Equivariance
- Models
  - Equivariant models
  - Invariant models
- Conclusion

# Sets and Multisets

## Definition :

**Set** : Collection of different elements

**Multiset** : Collection of elements, where we allow for multiple instances of the same element

## Example :

$\{a, b, c\}$

$\{a, a, b, a, b, c\}$

# Sets and Multisets

Definition :

**Set** : Collection of different elements

**Multiset** : Collection of elements, where we allow for multiple instances of the same element

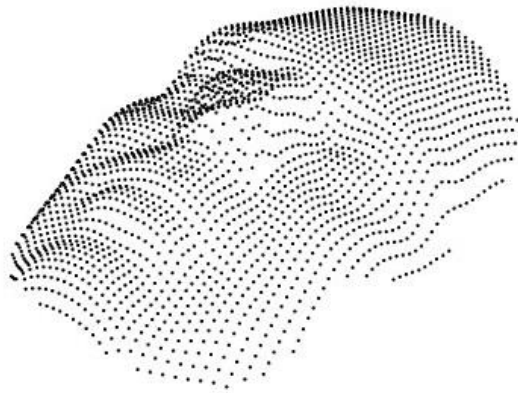
Example :

$\{a, b, c\}$

$\{a, a, b, a, b, c\}$

# Motivation for set models

## Point Clouds

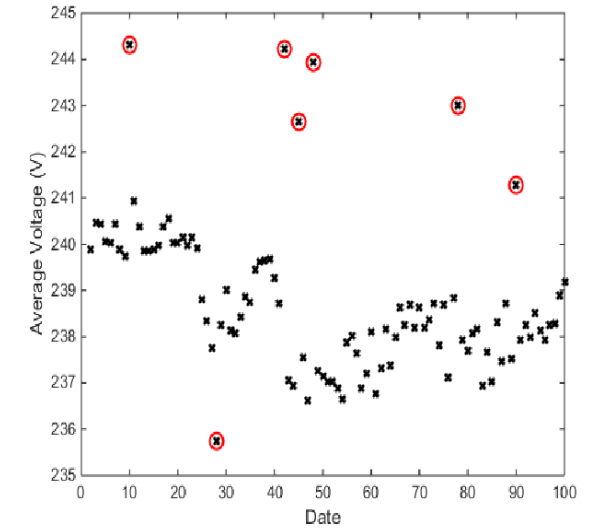


## Set expansion

Please enter a comma separated seed list of terms:

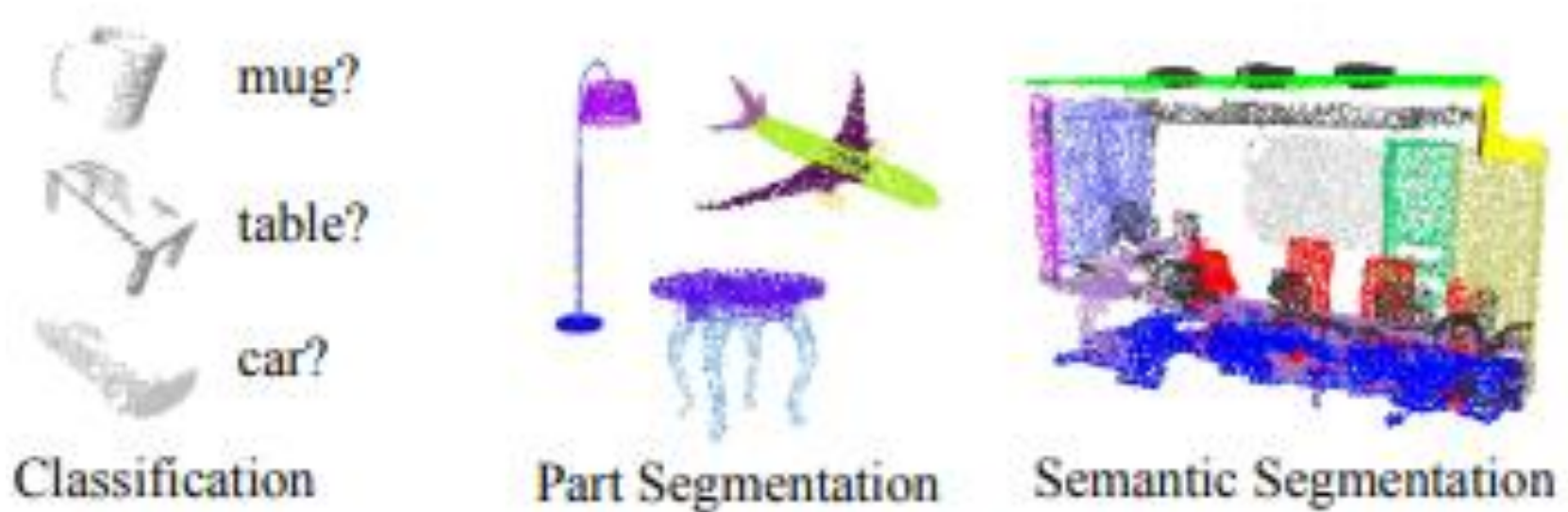
#	Results	Score
0	Germany	0.7273473739624023
1	the USA	0.7269347310066223
2	Lebanon	0.7219724655151367
3	Belgium	0.7119301557540894
4	the United States	0.7025378942489624
5	Saudi Arabia	0.6978805065155029
6	Morocco	0.6954267024993896

## Outlier detection



# Point Clouds

→ Set of point in space, representing 3D shape or object



# Outlier Detection



# What should a set model account for ?

→ The set can have varying number of elements

→ The order of the elements doesn't count



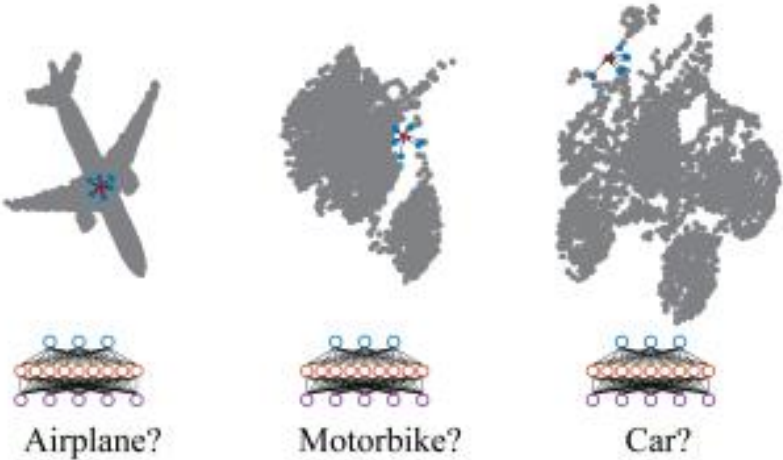
# What should a set model account for ?

→ The set can have varying number of elements

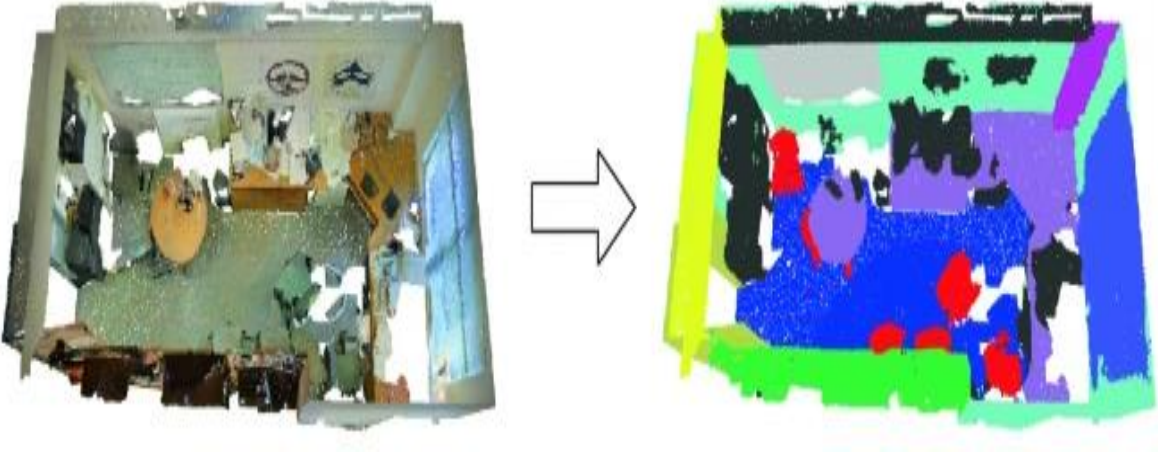
→ The order of the elements doesn't count

# Two types of tasks

Entire set

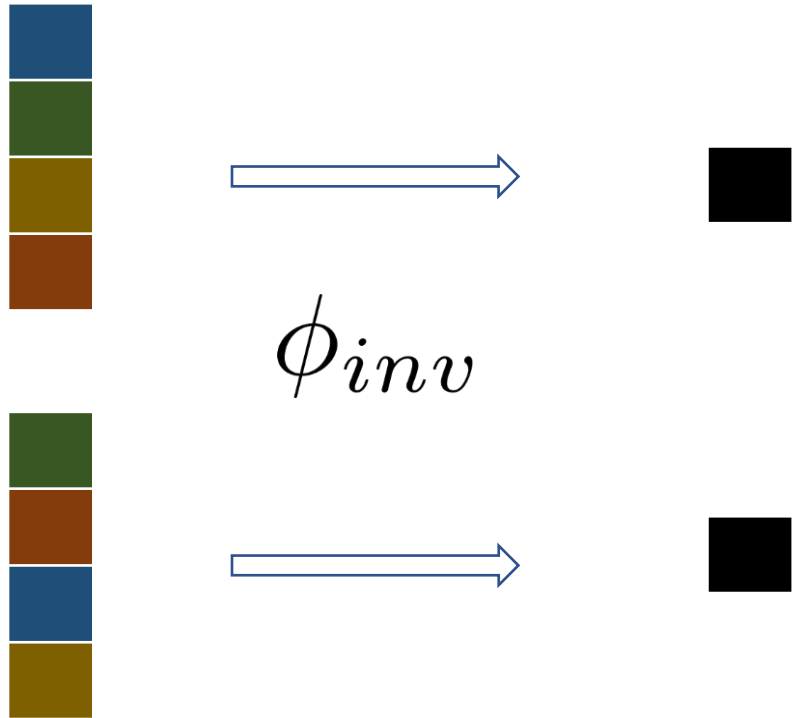


Elementwise

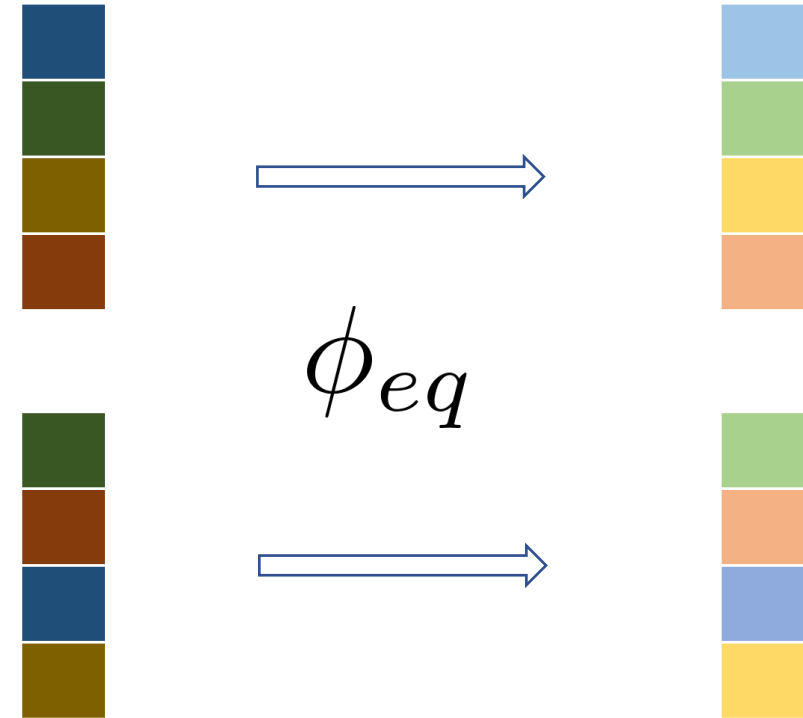


# Two notions arise

## Permutation Invariance



## Permutation Equivariance



# Group Invariance/Equivariance

Let  $\phi : X \rightarrow Y$  be a map and let  $G$  be a group acting on  $X$ . This means that  $G$  is identified by a sub-group of the bijections of  $X$ .

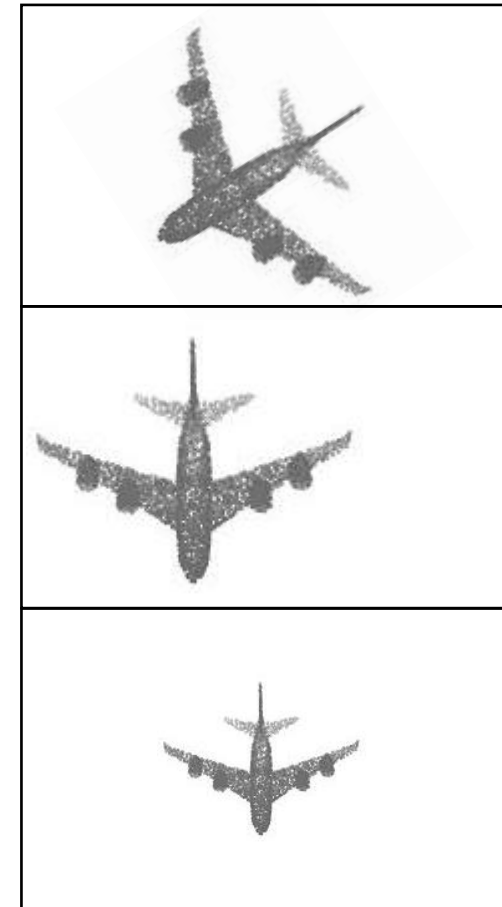
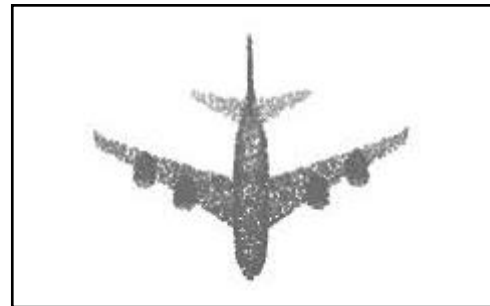
Invariance : we say  $\phi$  is  $G$ -invariant if  $\forall g \in G, x \in X, \phi(g.x) = \phi(x)$ .

Equivariance : If  $G$  acts also on  $Y$ , we say  $\phi$   $G$ -equivariant if  $\forall g \in G, x \in X, \phi(g.x) = g.\phi(x)$

# Group Invariance/Equivariance

→ there are other relevant groups

Example :



Rotation

Translation

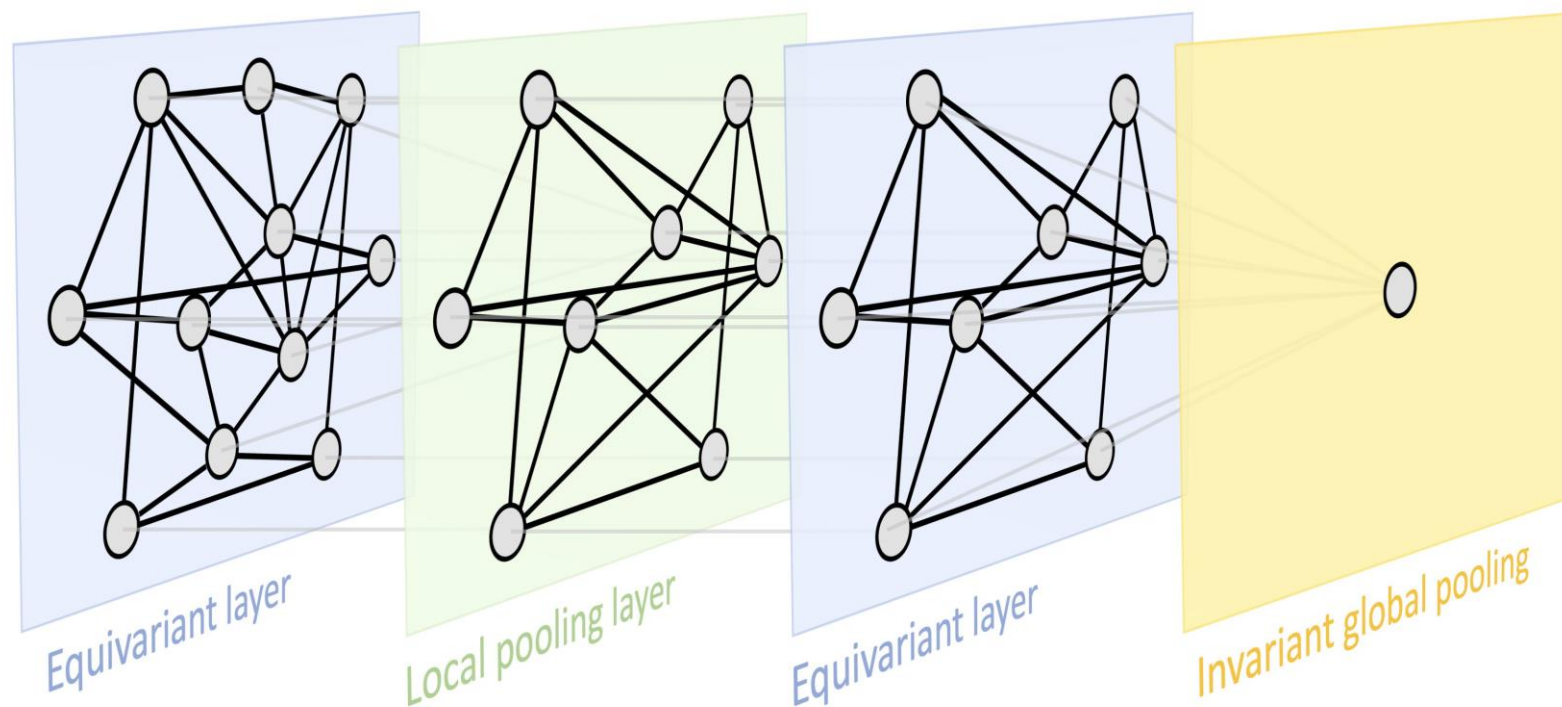
Scale

→ In this example, the group acts on each individual element

→ For permutation inv-/equivariance  
the group acts between elements of the set

# Permutation Invariance / Equivariance

→ two connected notions, often used side by side

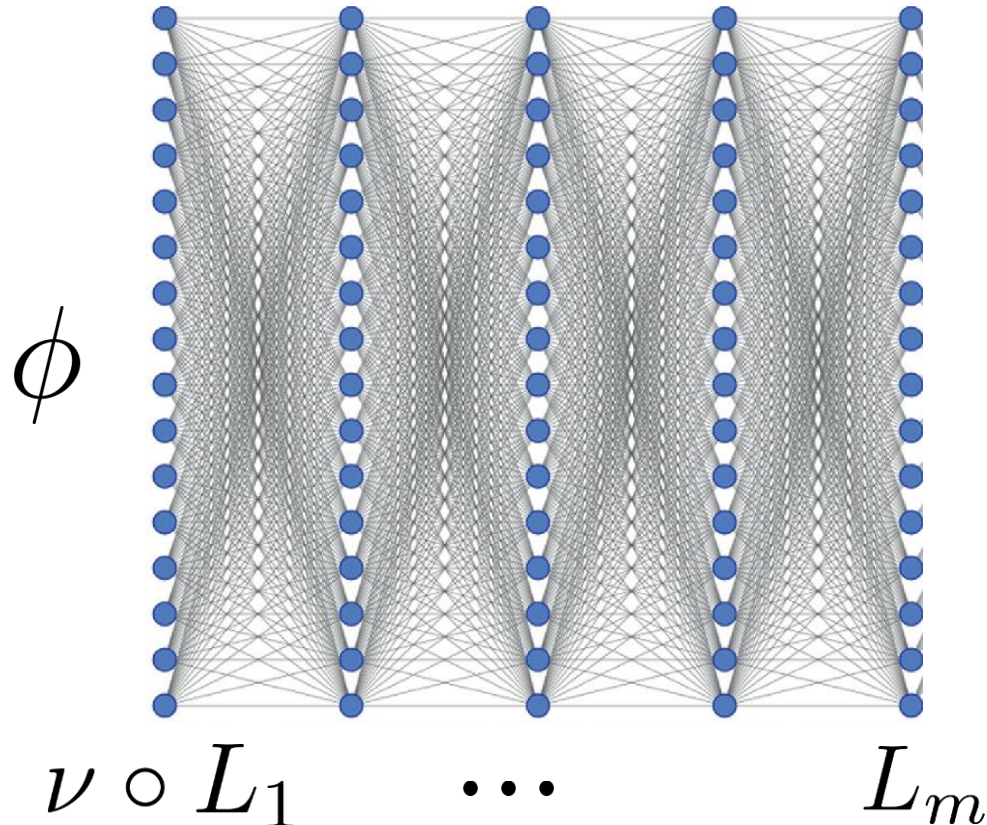


# Framework

We will consider a set of  $n$  elements  $x_1, \dots, x_n$ , each having features in  $\mathbb{R}^k$ , and group them into one matrix  $\mathbb{R}^{n \times k} \ni X = (x_1, \dots, x_n)^T$ .

# Equivariance

How do we model this with neural networks ?



$$\phi : \mathbb{R}^{n \times d_{in}} \rightarrow \mathbb{R}^{n \times d_{out}}$$

$$X \rightarrow L_m \circ \nu \circ \dots \circ \nu \circ L_1(X)$$

$$L_i : \mathbb{R}^{n \times d_i} \rightarrow \mathbb{R}^{n \times d_{i+1}}$$

→ the activation  $\nu$  is equivariant  
(elementwise operations)

→ need to pay attention to  $L_i$



# Equivariance

One Natural way :

→ apply an elementwise function :

$$x_1, \dots, x_n \longrightarrow (\pi(x_1), \dots, \pi(x_n))$$

# Equivariance

One Natural way :

→ apply an elementwise function :

$$x_1, \dots, x_n \longrightarrow (\pi(x_1), \dots, \pi(x_n))$$

Leads to :

$$L_i(X) = \lambda_i IX + \mathbf{1}c_i$$

$I \in \mathbb{R}^{n \times n}$  is the identity

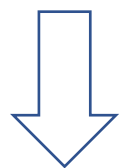
$\mathbf{1} = [1, \dots, 1]^T \in \mathbb{R}^n$

$\lambda_i \in \mathbb{R}^{d_i \times d_{i+1}}$ ,  $c_i \in \mathbb{R}^{d_{i+1}}$

→ Drawback : no interaction between elements of the set

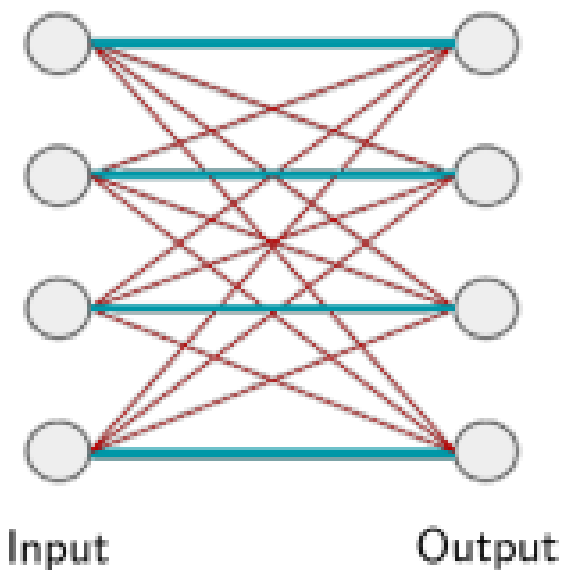
# Equivariance

$$L_i(X) = \lambda_i IX + \mathbf{1}c_i$$



$$L'_i(X) = \lambda_i IX + \gamma_i(\mathbf{1}^T \mathbf{1}) + \mathbf{1}c_i$$

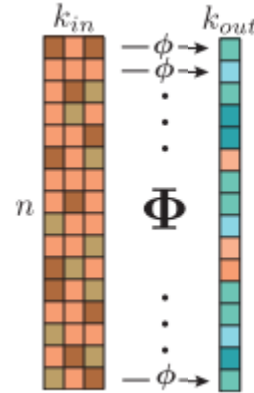
→ no interactions



# Equivariance

Deepsets :

$$X \rightarrow L'_m \circ \nu \circ \dots \circ \nu \circ L'_1(X)$$

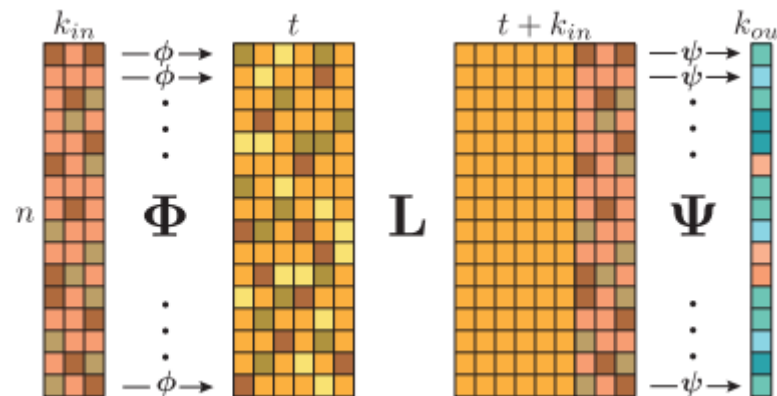


$$L_i(X) = \lambda_i IX + \mathbf{1}c_i$$

$$L'_i(X) = \lambda_i IX + \gamma_i(\mathbf{1}^T \mathbf{1}) + \mathbf{1}c_i$$

PointNetST :

$$X \rightarrow L_m \circ \nu \circ \dots \circ \nu \circ L'_i \circ \nu \circ \dots \circ \nu \circ L_1(X)$$



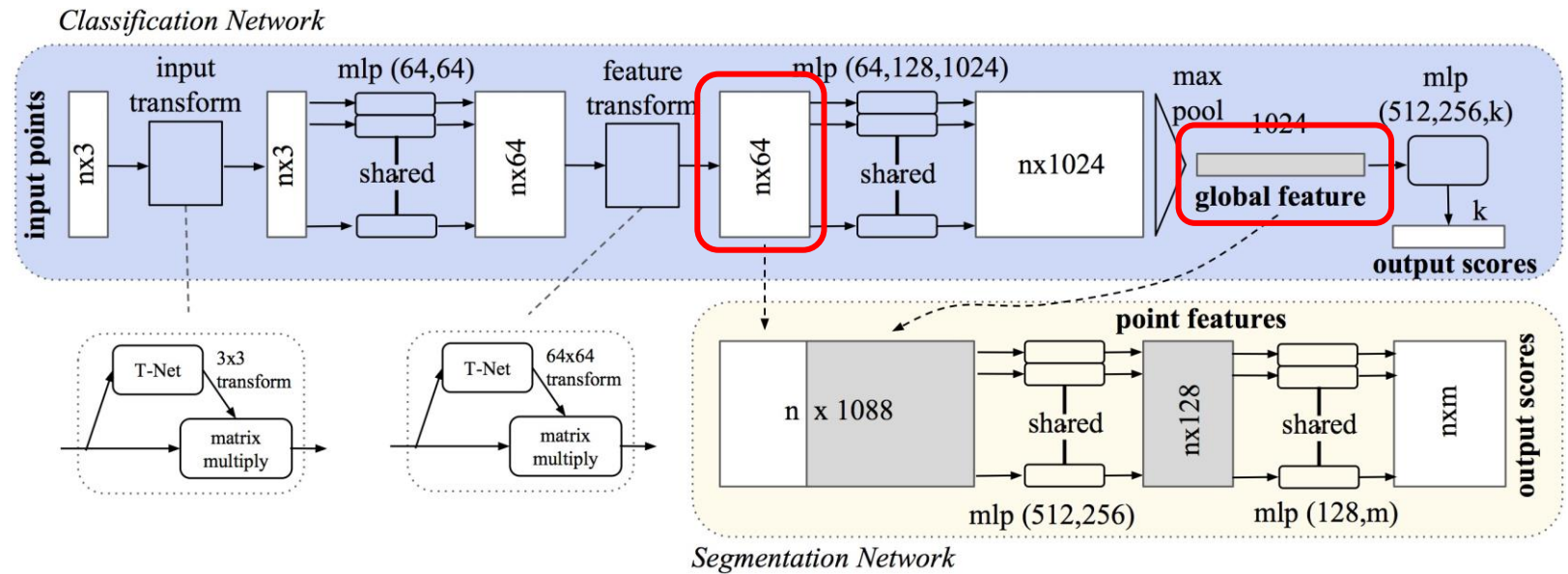
**→ both are universal approximator**

# Equivariance

For each element, concatenate :

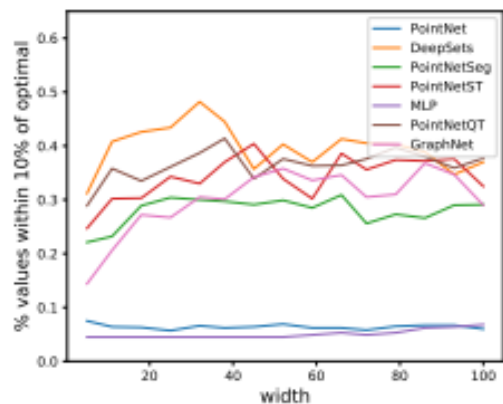
- features obtained by applying elementwise function
- the output of an invariant function

PointNetSeg :

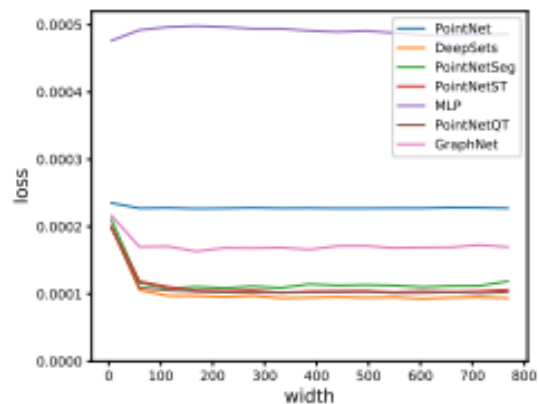


# Comparison

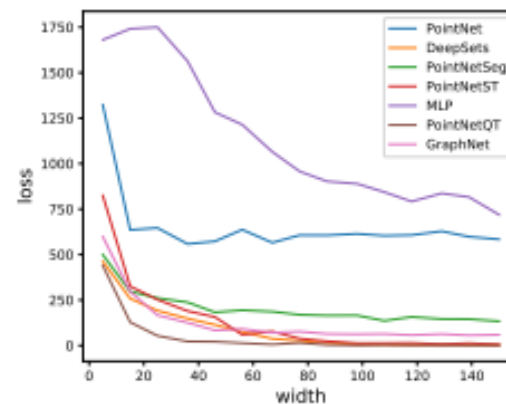
Knapsack test



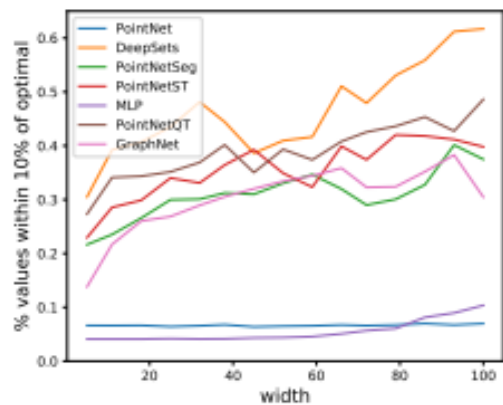
Fiedler test



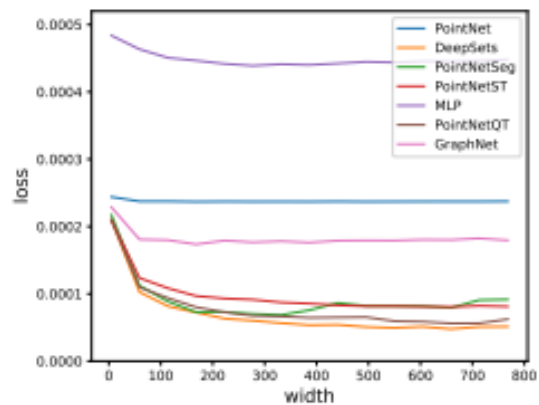
$\sum_{x \in \mathcal{X}} (x - \frac{1}{2})^2$  test



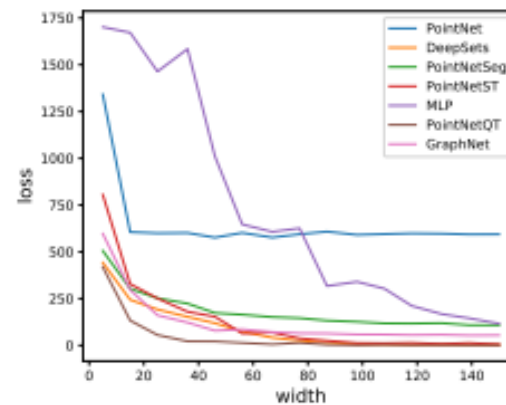
Knapsack train



Fiedler train



$\sum_{x \in \mathcal{X}} (x - \frac{1}{2})^2$  train

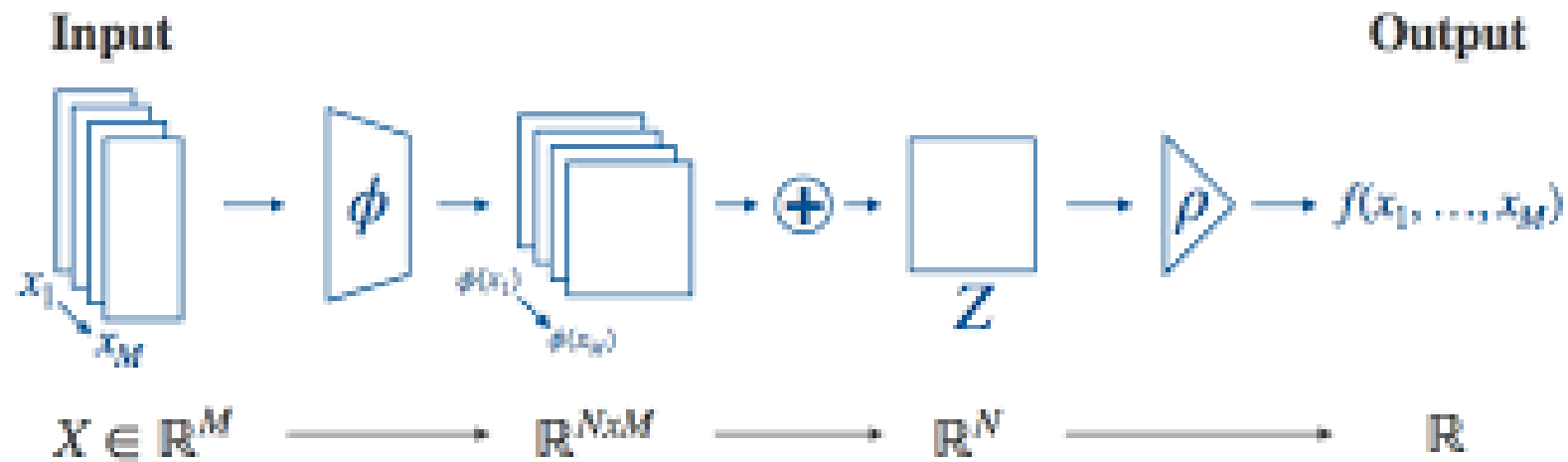


# Invariance

Permutation Invariance :  $\phi$  is permutation invariant if for all permutation  $\pi$ , we have :

$$\phi(x_{\pi(1)}, \dots, x_{\pi(n)}) = \phi(x_1, \dots, x_n)$$

# Deep Sets



→ **Universal Approximator of invariant functions**



# Deep Sets

$$f(X) = \rho \left( \sum_{i=1}^n \phi(x_i) \right)$$

- Encoder
- Aggregation / Pooling
- Decoder

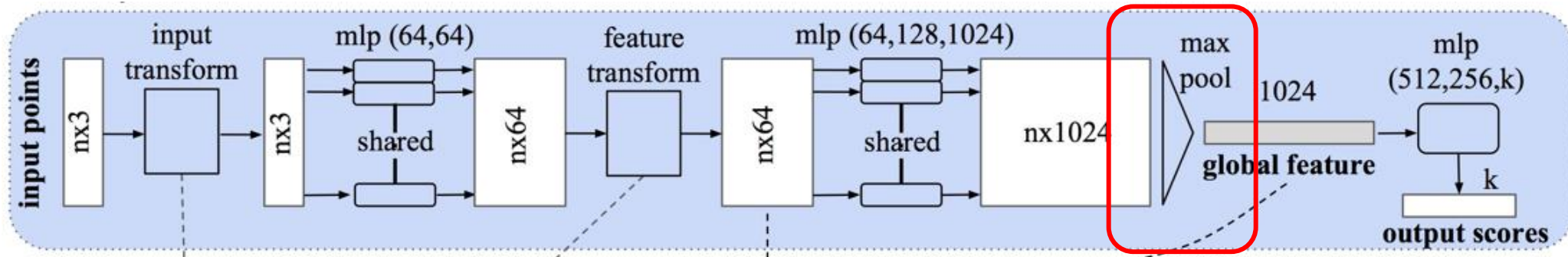
# Deep Sets

Other choices of aggregation functions :

- Mean

- Max (ex : PointNet ) :

$$f(X) = \rho \left( \max_{i=1..n} \phi(x_i) \right)$$



→ advantage of better generalization on varying set sizes

# Deep Sets

Some information on interactions is discarded during the encoding

$$f(X) = \rho \left( \sum_{i=1}^n \phi(x_i) \right)$$

→ Use previously discussed equivariant layers

→ Higher order Janossy pooling

# Janossy pooling

**Idea : consider each possible permutation on the elements and pass them into the same function, then average**

$$f(X) = \rho \left( \frac{1}{S_n} \sum_{\pi \in S_n} \phi(\pi(X)) \right)$$

Encoder

Aggregation  
/ Pooling

Decoder

# Janossy pooling

Obvious drawback : computational complexity (number of permutations)

Solutions :

- **Sorting** the elements
- **Sampling** among the permutations
- **Restricting** permutations to k-tuples

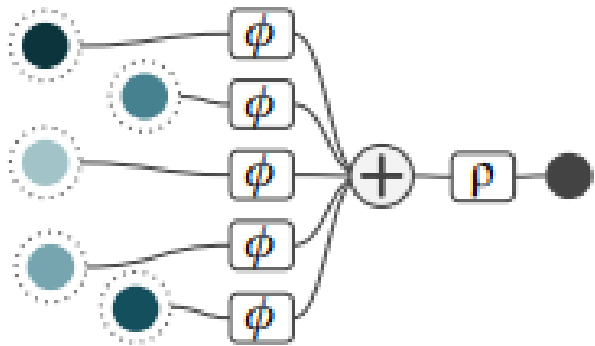
# Janossy pooling with k-tuples

$$f(X) = \rho \left( \frac{(n-k)!}{n!} \sum_{X_k} \phi(X_k) \right)$$

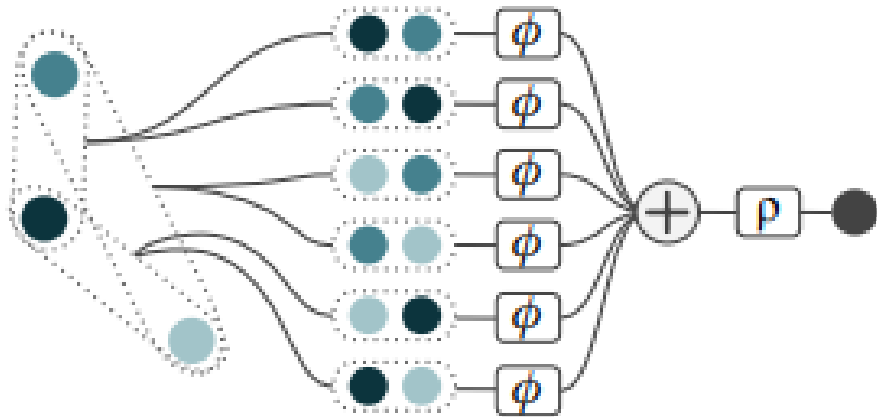
Sum over  $X_k$  formed by any k-tuple of elements in any order.

- Reduces complexity to  $\mathcal{O}(n^k)$
- Tuples incorporate interactions between elements
- Deep Set :  $k = 1$

# Janossy pooling

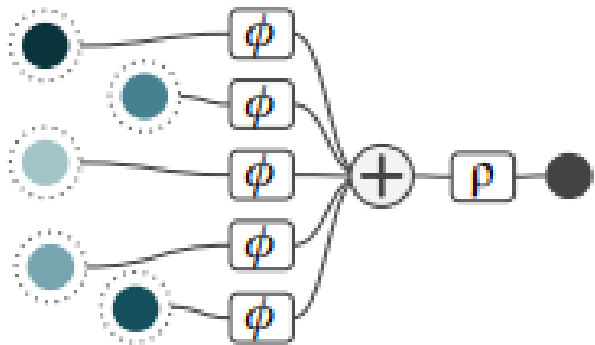


(a) Janossy pooling with  $k = 1$  (*Deep Sets*)

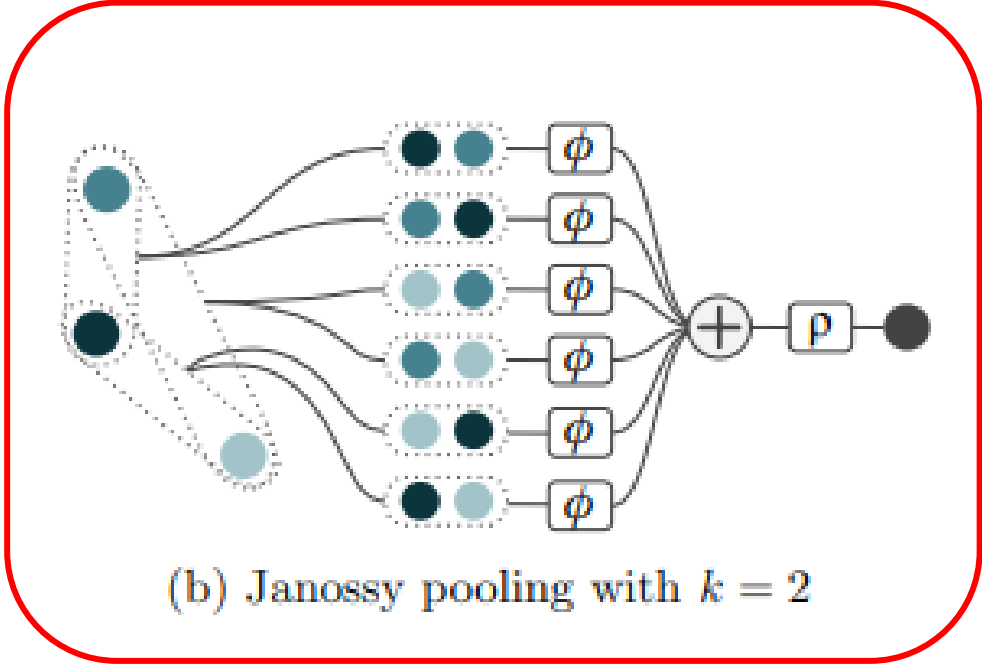


(b) Janossy pooling with  $k = 2$

# Janossy pooling



(a) Janossy pooling with  $k = 1$  (*Deep Sets*)

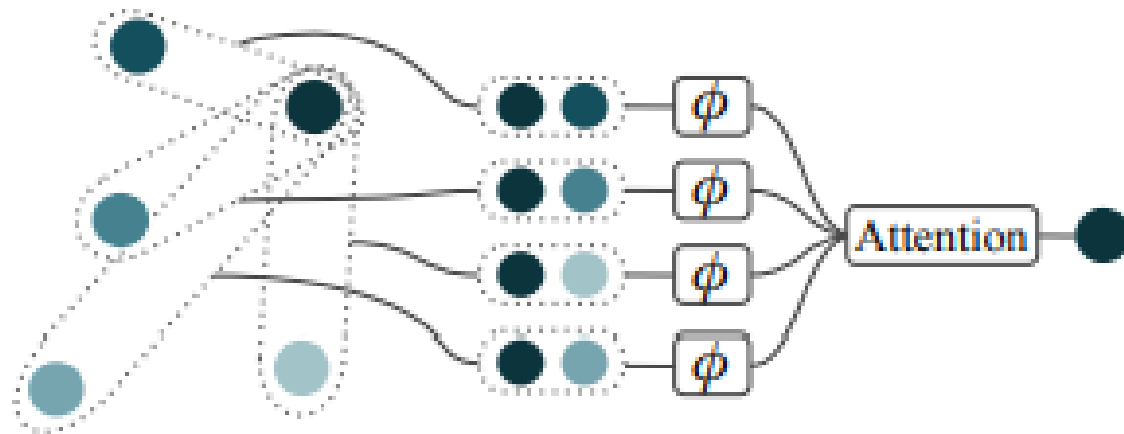


(b) Janossy pooling with  $k = 2$

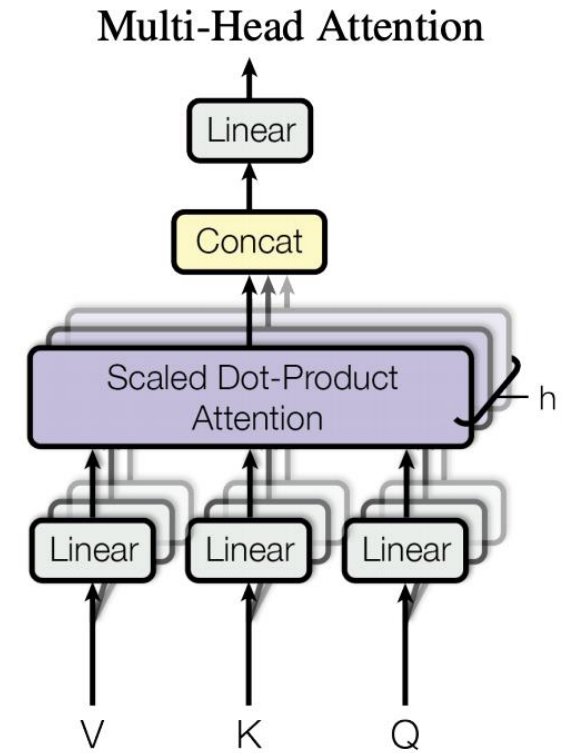


# Set transformer

Deal with interaction ?  
→ Using self-attention



(c) Self-attention

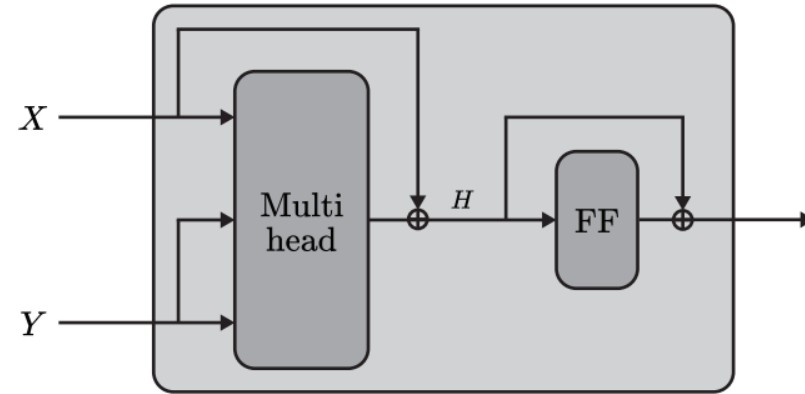


$$\text{MultiHead}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = [\text{head}_1; \dots; \text{head}_h] \mathbf{W}^O$$

where  $\text{head}_i = \text{Attention}(\mathbf{Q}\mathbf{W}_i^Q, \mathbf{K}\mathbf{W}_i^K, \mathbf{V}\mathbf{W}_i^V)$

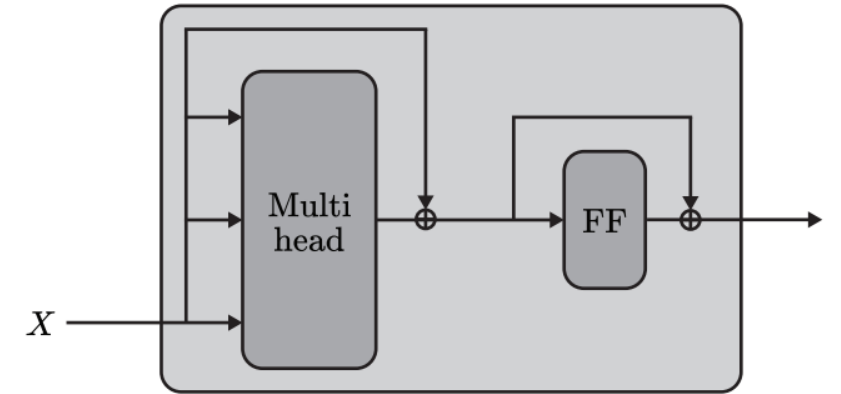
# Set transformer

Deal with interaction ?  
→ Using self-attention



(b) MAB

Multi Attention Block



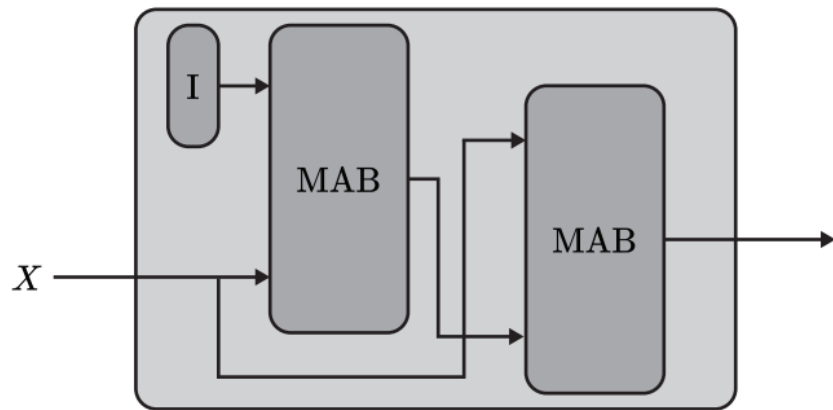
(c) SAB

Set Attention Block

Complexity is  $\mathcal{O}(n^2)$

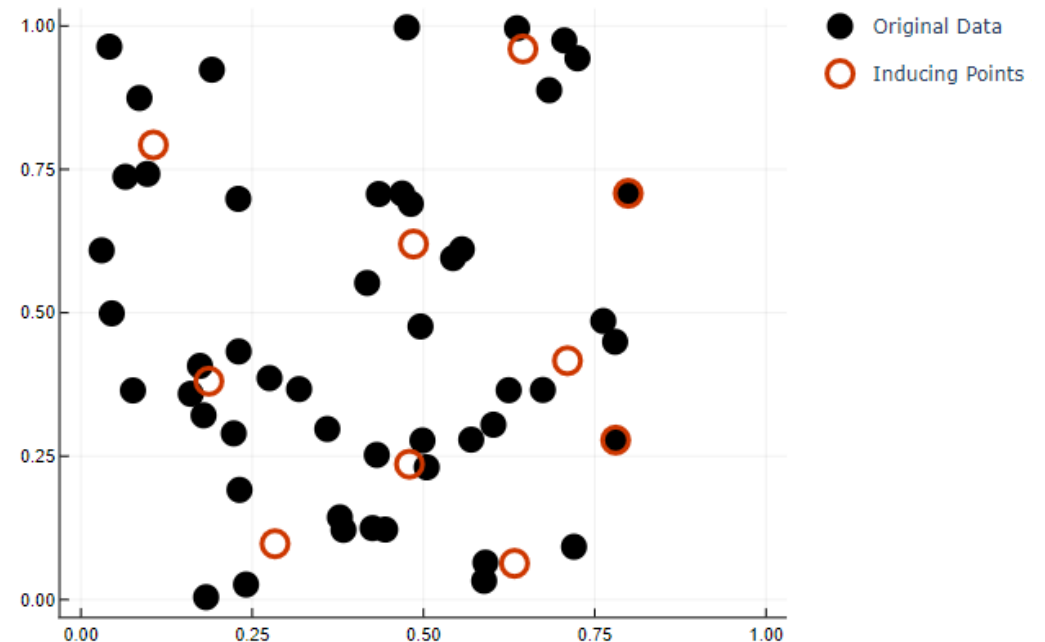
# Set transformer

→ Further reduce the complexity using inducing points



(d) ISAB

Complexity is now :  $\mathcal{O}(nm)$



# Set transformer

**Improve the aggregation function ?**

**→ differentiate influence of each instance**

# Set transformer

Idea :

- **Aggregate the encodings of each element, using attention to weight them according to their respective influence**
- **Output a set of k elements (usually k=1)**

Pooling Multihead Attention :

$$\text{PMA}_k = \text{MAB}(S, \text{rFF}(Z))$$

encoder output  $Z \in \mathbb{R}^{n \times d}$

$k$  seed vectors  $S \in \mathbb{R}^{k \times d}$

# Recap of invariant models

Deep Sets :

$$\boxed{\text{FF}} \circ \boxed{\text{Pooling}} \circ \boxed{\text{rFF}}(X)$$

Set Transformer :

$$\boxed{\text{FF}} \circ \boxed{\text{SAB} \circ \text{PMA}_1} \circ \boxed{\text{ISAB}_m \circ \dots \text{ISAB}_m}(X)$$

- Encoder
- Aggregation / Pooling
- Decoder

→ It is possible to mix those methods

# Comparison

Table 4. Test accuracy for the point cloud classification task using 100, 1000, 5000 points.

	Architecture	100 pts	1000 pts	5000 pts
Deep Sets	rFF + Pooling (Zaheer et al., 2017)	-	$0.83 \pm 0.01$	-
	rFFp-max + Pooling (Zaheer et al., 2017)	$0.82 \pm 0.02$	$0.87 \pm 0.01$	<b><math>0.90 \pm 0.003</math></b>
Set Transformer	rFF + Pooling	$0.7951 \pm 0.0166$	$0.8551 \pm 0.0142$	$0.8933 \pm 0.0156$
	rFF + PMA (ours)	$0.8076 \pm 0.0160$	$0.8534 \pm 0.0152$	$0.8628 \pm 0.0136$
	ISAB (16) + Pooling (ours)	$0.8273 \pm 0.0159$	<b><math>0.8915 \pm 0.0144</math></b>	<b><math>0.9040 \pm 0.0173</math></b>
	ISAB (16) + PMA (ours)	<b><math>0.8454 \pm 0.0144</math></b>	$0.8662 \pm 0.0149$	$0.8779 \pm 0.0122$

**rFF** : Row-wise Feed Forward layer  
**rFFp-max** : rFF with permutation equivariant variants  
**ISAB** : Inducing Set Attention Block  
**Pooling** : Usual pooling without attention  
**PMA** : Pooling by Multihead Attention

# Comparison

Table 5. Meta set anomaly results. Each architecture is evaluated using average of test AUROC and test AUPR.

	Architecture	Test AUROC	Test AUPR
Deep Sets	Random guess	0.5	0.125
	rFF + Pooling	$0.5643 \pm 0.0139$	$0.4126 \pm 0.0108$
	rFFp-mean + Pooling	$0.5687 \pm 0.0061$	$0.4125 \pm 0.0127$
	rFFp-max + Pooling	$0.5717 \pm 0.0117$	$0.4135 \pm 0.0162$
	rFF + Dotprod	$0.5671 \pm 0.0139$	$0.4155 \pm 0.0115$
Set Transformer	SAB + Pooling (ours)	$0.5757 \pm 0.0143$	$0.4189 \pm 0.0167$
	rFF + PMA (ours)	$0.5756 \pm 0.0130$	$0.4227 \pm 0.0127$
	SAB + PMA (ours)	<b><math>0.5941 \pm 0.0170</math></b>	<b><math>0.4386 \pm 0.0089</math></b>



# References

- [1] Zaheer M., Kottur S., Ravanbakhsh S., Póczos B., Salakhutdinov R. R., and Smola A. J. **Deep sets**. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [2] Charles R. Q., Su H., Kaichun M., and Guibas L. J. **PointNet: Deep learning on point sets for 3D classification and segmentation**. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [3] Segol N., Lipman Y. **On Universal Equivariant Set Networks**, 2020.
- [4] Wagstaff E., Fuchs F. B., Engelcke M., Osborne M. A., Posner I. **Universal Approximation of Functions on Sets**, 2021.
- [5] Lee J., Lee Y., Kim J., Kosiorek A. R., Choi S., Teh Y. W. **Set Transformer: A Framework for Attention-based Permutation-Invariant Neural Networks**, In *International Conference on Machine Learning (ICML)*, 2019.

# Conclusion

Questions ?

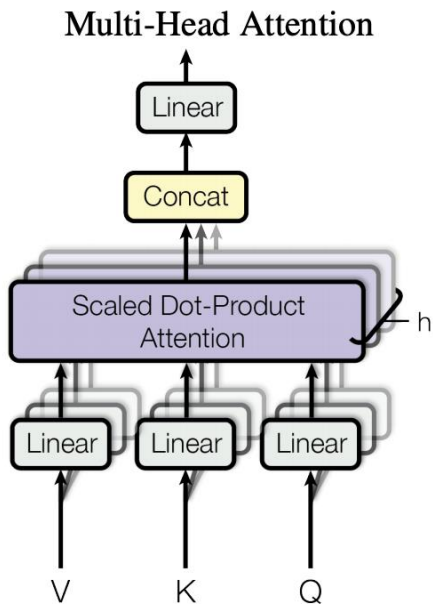
# Set transformer

## Multihead-Attention :

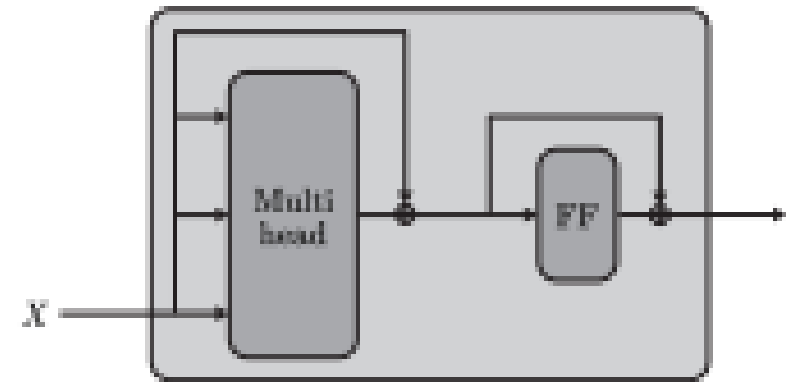
$$\text{Multihead}(Q, K, V; \lambda, \omega) = \text{concat}(O_1, \dots, O_h)W^O$$

$$O_j = \text{Att}(QW_j^Q, KW_j^K, VW_j^V; w_j)$$

$$\text{Att}(Q, K, V; \omega) = \omega(QK^T)V$$



## Attentions blocks :



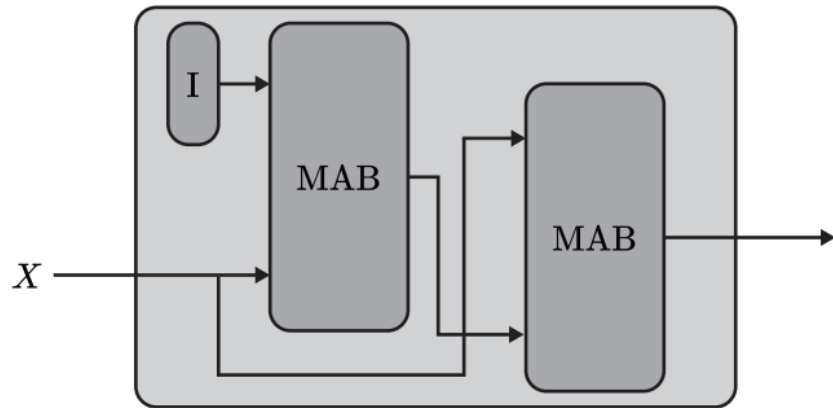
(c) SAB

$$\text{SAB}(X) = \text{MAB}(X, X) \in \mathbb{R}^{n \times d}$$

$$\text{MAB}(X, Y) := \text{LayerNorm}(H + r\text{FF}(H))$$

$$\text{where } H = \text{LayerNorm}(X + \text{Multihead}(X, Y, Y; \omega))$$

# Set transformer



(d) ISAB

$$\text{MAB}(X, Y) := \text{LayerNorm}(H + \text{rFF}(H))$$
  
where  $H = \text{LayerNorm}(X + \text{Multihead}(X, Y, Y; \omega))$

$$\text{ISAB}_m(X) := \text{MAB}(X, H) \in \mathbb{R}^{n \times d}$$
  
where  $H = \text{MAB}(I, X) \in \mathbb{R}^{m \times d}$

# Set transformer

Improve the aggregation function ?

→ differentiate influence of each instance

## Pooling Multihead Attention :

Idea :

- Aggregate the encodings of each element, using attention to weight them according to their respective influence
- Output a set of  $k$  elements

$\text{MAB}(X, Y) := \text{LayerNorm}(H + \text{rFF}(H))$   
where  $H = \text{LayerNorm}(X + \text{Multihead}(X, Y, Y; \omega))$

$\text{PMA}_k = \text{MAB}(S, \text{rFF}(Z))$   
encoder output  $Z \in \mathbb{R}^{n \times d}$   
 $k$  seed vectors  $S \in \mathbb{R}^{k \times d}$