

# CS6208 : Advanced Topics in Artificial Intelligence

## Graph Machine Learning

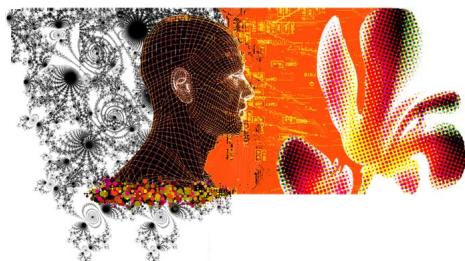
### Lecture 6 : Weisfeiler-Lehman GNNs

Semester 2 2022/23

Xavier Bresson

<https://twitter.com/xbresson>

Department of Computer Science  
National University of Singapore (NUS)



# Outline

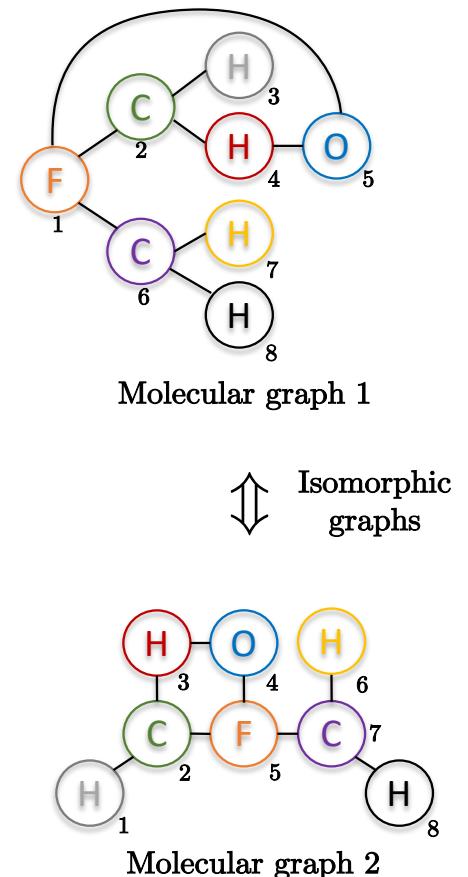
- Weisfeiler-Lehman isomorphism test
- Graph isomorphism networks
- Equivariant graph neural networks
- Expressivity power and generalization
- Equivariant set neural networks
- Conclusion

# Outline

- Weisfeiler-Lehman isomorphism test
- Graph isomorphism networks
- Equivariant graph neural networks
- Expressivity power and generalization
- Equivariant set neural networks
- Conclusion

# Graph isomorphism

- Graph isomorphism : Two graphs are isomorphic if there exists an index permutation between the nodes that preserves node adjacencies.
  - An example of two isomorphic molecular graphs. Node correspondence is given by the node color.
- Determining whether two graphs are isomorphic is NP-intermediate : It is not known if a polynomial time algorithm exists, or the problem is NP-hard.
- Weisfeiler-Lehman (WL) proposed an algorithm<sup>[1]</sup> to test if two graphs are not isomorphic.
  - However, the WL test is not sufficient to guarantee that two graphs are isomorphic (necessary condition only).



[1] Weisfeiler, Lehman, A reduction of a graph to a canonical form and an algebra arising during this reduction, 1968

# Weisfeiler-Lehman test<sup>[1]</sup>

- Design an injective “coloring” function  $f_{WL}$  that takes a pair (node, its neighborhood) as input, and outputs a new node color :

$$f_{WL}(c_i^\ell, \{c_j^\ell\}_{j \in \mathcal{N}_i}) = c_i^{\ell+1}$$

iteration  
Multiset function      Multiset (set of unordered and repetitive elements)

- Properties of  $f_{WL}$  :

- Unique color for the same pair of (node, its neighborhood).

$$f_{WL}\left(\underbrace{\text{(F)}}_{\text{Same color}} \underbrace{\{C, C, O\}}_{=}\right) = \text{blue circle}$$

$$f_{WL}\left(\underbrace{\text{(F)}}_{\text{Same color}} \underbrace{\{O, C, C\}}_{=}\right) = \text{blue circle}$$

- Different colors for distinct pairs of (node, its neighborhood).

$$f_{WL}\left(\underbrace{\text{(F)}}_{\text{Different color}} \underbrace{\{C, O\}}_{\neq}\right) = \text{green circle}$$

$$f_{WL}\left(\underbrace{\text{(F)}}_{\text{Different color}} \underbrace{\{C, C\}}_{\neq}\right) = \text{red circle}$$

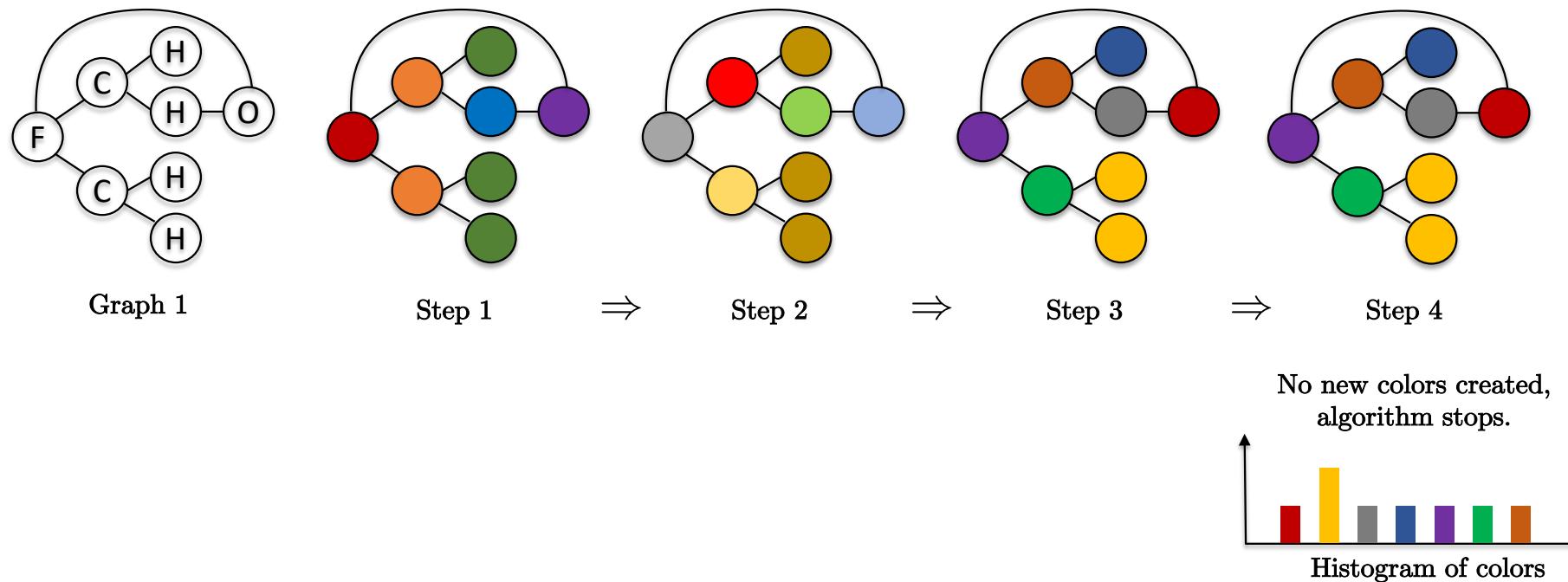
Simply said,  $f_{WL}$  must be injective (mapping different inputs to different outputs).

Injectivity is a discriminative property.

[1] Weisfeiler, Lehman, A reduction of a graph to a canonical form and an algebra arising during this reduction, 1968

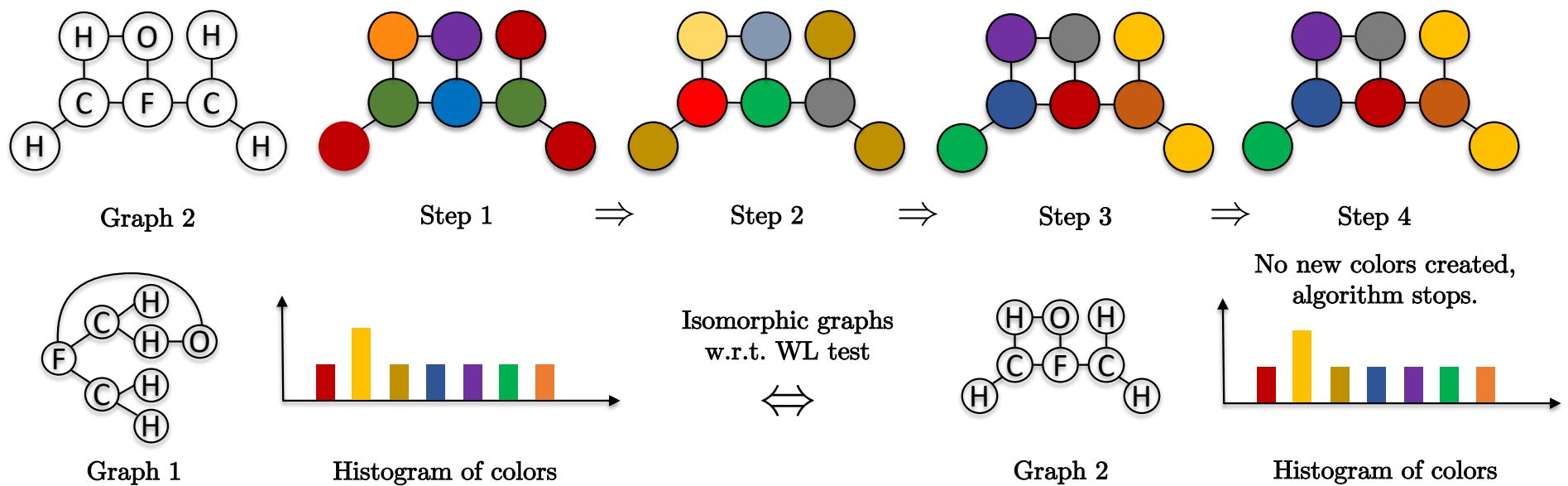
# Weisfeiler-Lehman test

- WL algorithm iteratively applies the coloring function  $f_{WL}$  until no new colors are created.
- This produces a canonical representation of a graph as a histogram of colors.



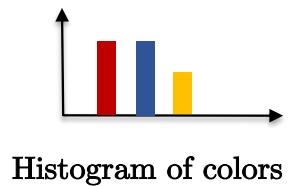
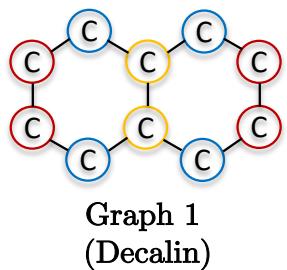
# Weisfeiler-Lehman test

- The WL test is a necessary but not sufficient condition for isomorphism.
  - If the two graphs have different color histograms then the two graphs are guaranteed to be non-isomorphic.
  - If the two graphs produce the same histogram of colors, the graphs are possibly isomorphic (but not guaranteed).

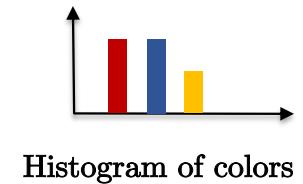
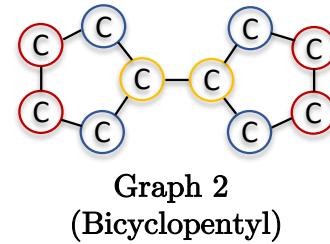


# Weisfeiler-Lehman test

- Example of two (simple) non-isomorphic graphs that fail to be distinguished with the WL test.



Isomorphic graphs  
w.r.t. WL test  
 $\iff$



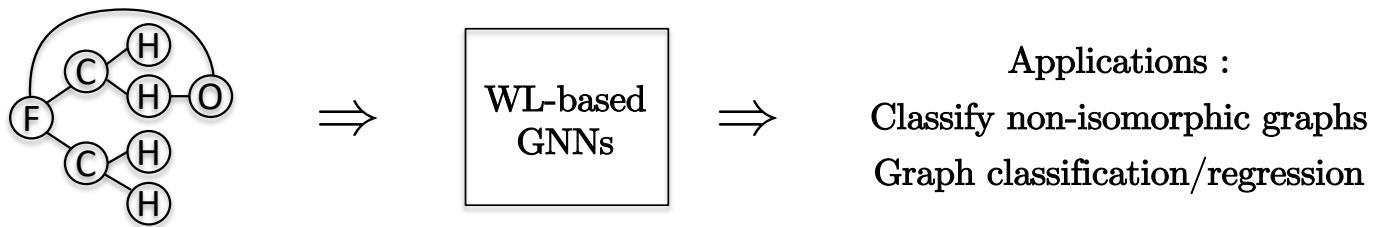
- The original WL test is technically known as the 1-WL test.

# Outline

- Weisfeiler-Lehman isomorphism test
- **Graph isomorphism networks**
- Equivariant graph neural networks
- Expressivity power and generalization
- Equivariant set neural networks
- Conclusion

# WL-based graph neural networks

- Can we turn the WL test into a neural network to classify non-isomorphic graphs ?
- Yes, but is it necessary to have a learning process to implement the WL test ? Yes and no.
  - Yes because designing a constraint-free analytical injective function  $f_{WL}$  is hard (function  $f_{WL}$  as Hash function only provides an approximation).
  - No because neural networks are also limited to provide an approximation of function  $f_{WL}$ . However, their approximation can be better than Hash function because of the feature learning mechanism.
- Observe that distinguishing non-isomorphic graphs is not necessarily useful in practice. What is useful is the “coloring” process that can provide expressive representation of data-structured graphs and can be used for downstream tasks like graph classification or regression.



# Graph isomorphism networks (GINs)<sup>[1]</sup>

- The WL test algorithm iterates the “coloring” formula :

$$f_{WL}(c_i^\ell, \{c_j^\ell\}_{j \in \mathcal{N}_i}) = c_i^{\ell+1}$$

where function  $f_{WL}$  provides :

- A unique color for the same pair  $(c_i, \{c_j\}_{\mathcal{N}_i})$ .
- Different colors for distinct pairs  $(c_i, \{c_j\}_{\mathcal{N}_i})$ .

This implies that  $f_{WL}$  must be injective :

$$c_i \neq c_{i'} \text{ or } \{c_j^\ell\}_{j \in \mathcal{N}_i} \neq \{c_j^\ell\}_{j \in \mathcal{N}_{i'}} \Rightarrow f_{WL}(c_i^\ell, \{c_j^\ell\}_{j \in \mathcal{N}_i}) \neq f_{WL}(c_{i'}^\ell, \{c_j^\ell\}_{j \in \mathcal{N}_{i'}})$$

- Which aggregator functions are injective?

[1] Xu, Hu, Leskovec, Jegelka, How powerful are graph neural networks?, 2019

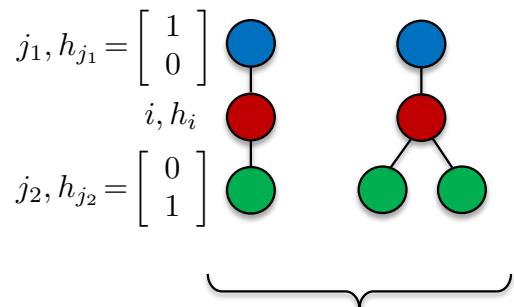
# Graph isomorphism networks

- (Simple) aggregator functions  $f_{WL}$  :

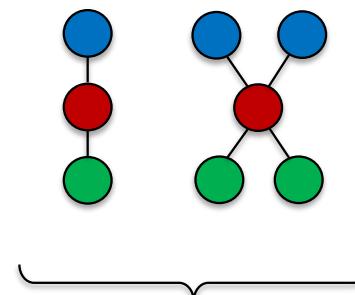
$$h_i = \max_{j \in \mathcal{N}_i} h_j \quad \text{Not injective}$$

$$h_i = \text{mean}_{j \in \mathcal{N}_i} h_j \quad \text{Not injective}$$

$$h_i = \sum_{j \in \mathcal{N}_i} h_j \quad \text{Injective}$$



Max fails to discriminate  
the red node  
Mean, Sum succeed



Max, Mean fail to  
discriminate the red node  
Sum succeeds

Injectivity is a discriminative  
property at the node level  
and by extension at the  
graph level too.

# Graph isomorphism networks

- Can we design a GNN function  $f_{\text{GNN}}$  that is as expressive as the WL test to distinguish non-isomorphic graphs? Yes, this GNN has an aggregation function of the form :

$$f_{\text{GNN}}(h_i^\ell, \{h_j^\ell\}_{j \in \mathcal{N}_i}) = (1 + \varepsilon)g(h_i) + \sum_{j \in \mathcal{N}_i} g(h_j)$$

where we must have

- Function  $g$  injective.
  - Sum aggregator (mean/max are not injective).
  - Constant  $\varepsilon$  must be irrational (a number that cannot be written as a ratio of two integers,  $\pi = 3.1415..$ )
- It is difficult to design an analytical injective function  $g$  (upper bound on the neighborhood size for countable colors, one-hot encoding does not provide meaningful distance between features).
    - An MLP can approximate  $g$  as the universal approximation theorem guarantees the existence of such function (although SGD is not guaranteed to find it):

$$h_i^{\ell+1} = f_{\text{GIN}}(h_i^\ell, \{h_j^\ell\}_{j \in \mathcal{N}_i}) = \text{MLP}^\ell \left( (1 + \varepsilon)h_i^\ell + \sum_{j \in \mathcal{N}_i} h_j^\ell \right)$$

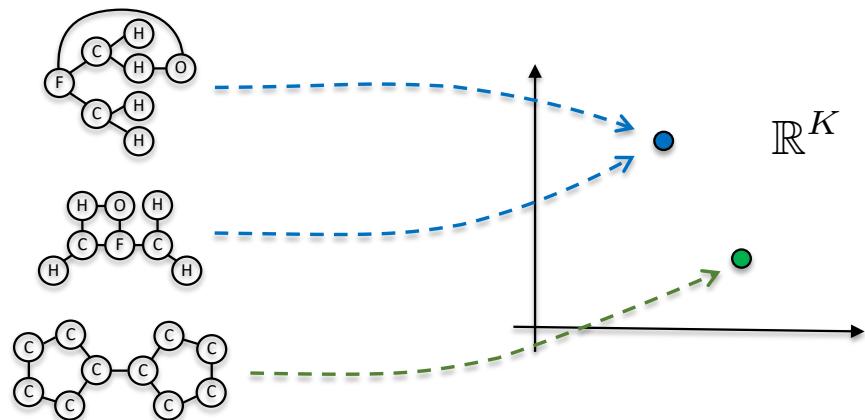
Scalar  $\varepsilon$  is learned (as computer  
cannot represent irrational numbers)

# Graph isomorphism networks

- Graph readout aggregation function must also be injective :

$$h_G = \text{MLP} \left( \sum_{i \in V} h_i^L \right) \in \mathbb{R}^K$$

Sum function



- Summary :
  - Pioneer work<sup>[1]</sup> with<sup>[2]</sup> on the theoretical question of expressivity/representation/discriminative power of GNNs.
  - GINs<sup>[1]</sup> are designed to be as expressive as the original WL test<sup>[3]</sup>, meaning these networks can learn different graph representations in  $\mathbb{R}^K$  for graphs that can be distinguished by the WL test (coloring and injectivity are the key properties).

[1] Xu, Hu, Leskovec, Jegelka, How powerful are graph neural networks?, 2019

[2] Morris, Ritzert, Fey, Hamilton, Lenssen, Rattan, Grohe, Weisfeiler and leman go neural: Higher-order graph networks, 2019

[3] Weisfeiler, Lehman, A reduction of a graph to a canonical form and an algebra arising during this reduction, 1968

# Principal Neighbourhood Aggregation (PNA)<sup>[1]</sup>

- Graph Isomorphism Networks<sup>[2]</sup> introduced GNNs are as expressive as the WL test to distinguish non-isomorphic graphs and have an aggregation function of the form :

$$f_{\text{NN}}(h_i^\ell, \{h_j^\ell\}_{j \in \mathcal{N}_i}) = (1 + \varepsilon)g(h_i^\ell) + \sum_{j \in \mathcal{N}_i} g(h_j^\ell)$$

where function  $g$  is injective and sum aggregator is used (mean/max are not injective).

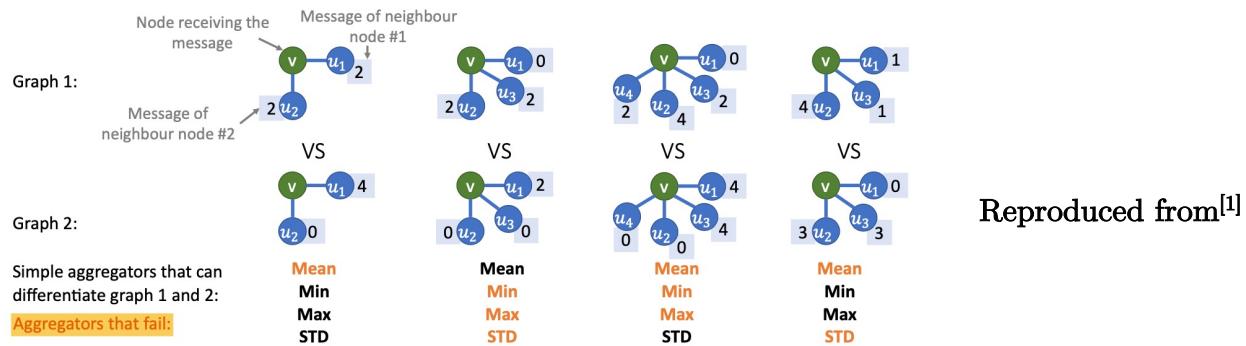
- Principal Neighbourhood Aggregation : Generalization of the injective sum aggregator function to multiple aggregators.
  - Aggregators are mean, max, min and standard deviation.
  - Overall injective function (that guarantees GNN expressivity w.r.t WL test).
  - Increase the expressivity power of standard MP-GNNs while preserving O(N) complexity.

[1] Corso, Cavalleri, Beaini, Lio, and Velickovic, Principal neighbourhood aggregation for graph nets, 2020

[2] Xu, Hu, Leskovec, Jegelka, How powerful are graph neural networks?, 2019

# Principal Neighbourhood Aggregation

- Illustration of aggregators are mean, max, min and standard deviation :



- PNA equation :

$$X_i^{\ell+1} = \text{Concat}\left(X_i^\ell, \tilde{X}_i^\ell, W_1^\ell\right) W_2^\ell \in \mathbb{R}^{d_{\ell+1}}$$

$$\tilde{X}_i^\ell = \bigoplus_{j \in \mathcal{N}_i} (X_j^\ell, X_j^\ell) \in \mathbb{R}^{12 \times d_\ell}$$

$$= \begin{bmatrix} \text{Mean}_{j \in \mathcal{N}_i}(X_j) \\ \text{Std Dev}_{j \in \mathcal{N}_i}(X_j) \\ \max_{j \in \mathcal{N}_i}(X_j) \\ \min_{j \in \mathcal{N}_i}(X_j) \end{bmatrix} \otimes \begin{bmatrix} \text{Id} \\ S(d, 1) \\ S(d, -1) \end{bmatrix}, \quad S(d, \alpha) = \left(\frac{\log(d+1)}{\delta}\right)^\alpha$$

Scalers that amplifies or decreases amplitude of incoming message.

[1] Corso, Cavalleri, Beaini, Lio, and Velickovic, Principal neighbourhood aggregation for graph nets, 2020

# Outline

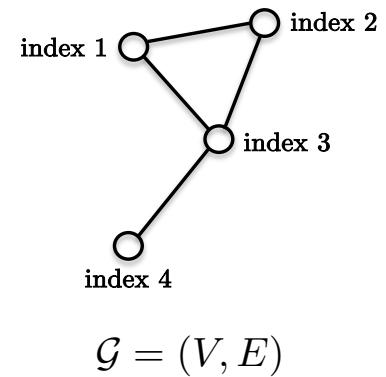
- Weisfeiler-Lehman isomorphism test
- Graph isomorphism networks
- **Equivariant graph neural networks**
- Expressivity power and generalization
- Equivariant set neural networks
- Conclusion

# Graph representations

- Graphs are rank-2 tensors. They can be represented as

- List of edges (short list for sparse graphs) :

$$E = [ \{1,2\}, \{1,3\}, \{2,3\}, \{3,4\} ]$$



- Matrix (dense tensor with n being the #nodes) :

$$A = \begin{matrix} V & 1 & 2 & 3 & 4 \\ 1 & \begin{bmatrix} 0 & 1 & 1 & 0 \end{bmatrix} \\ 2 & \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix} \\ 3 & \begin{bmatrix} 1 & 1 & 0 & 1 \end{bmatrix} \\ 4 & \begin{bmatrix} 0 & 0 & 1 & 0 \end{bmatrix} \end{matrix} \in \mathbb{R}^{n \times n}$$

- Other tensors :

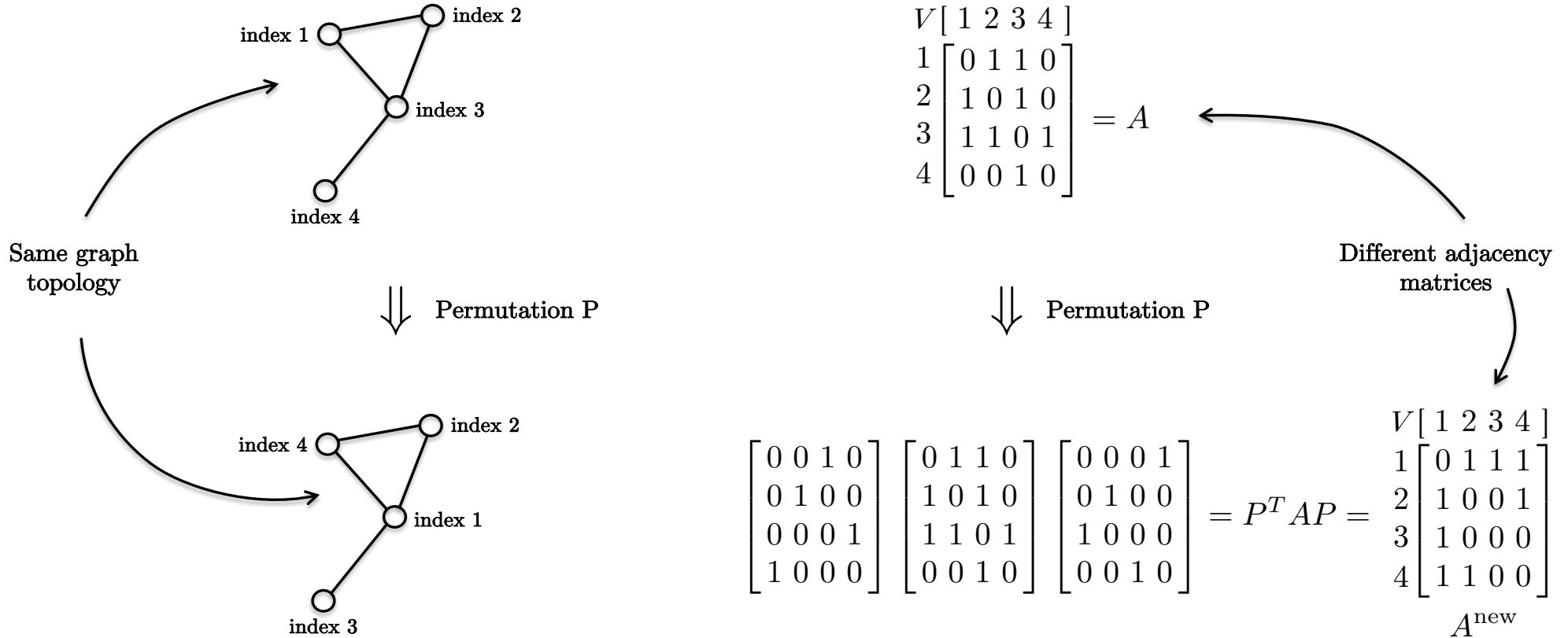
- Vectors (ordered and unordered sets) are rank-1 tensors :  $\begin{bmatrix} \quad \end{bmatrix} \in \mathbb{R}^n$

- General tensors (s.a. 3D-MRIs or hyper-graphs) of rank-m :



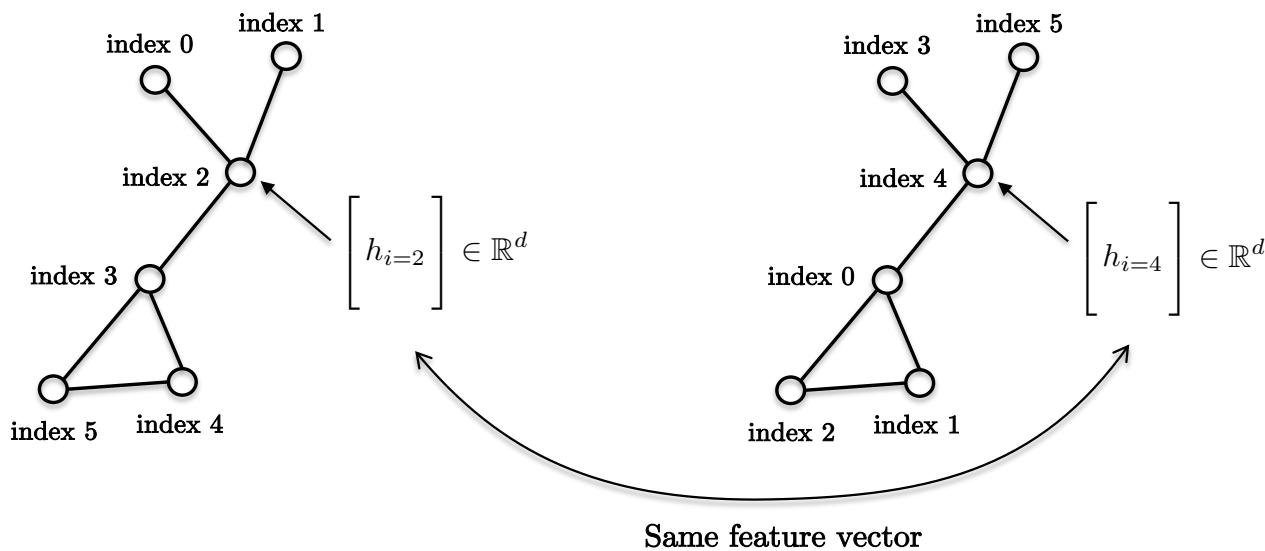
# Graph symmetry

- The most basic graph symmetry is the invariance of the topology w.r.t. index permutation :



# GNNs are permutation invariant

- A built-in property of GNNs is the invariance to index parametrization, i.e. node representation is independent from the choice of node ordering :



- Note that all GNNs s.a. GCNs<sup>[1]</sup>, GAT<sup>[2]</sup> are permutation invariant.

[1] Kipf, Welling, Semi-supervised classification with graph convolutional networks, 2017

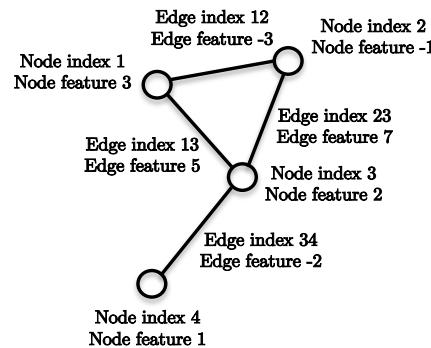
[2] Velickovic, Cucurull, Casanova, Romero, Lio, Bengio, Graph attention networks, 2017

# Equivariant functions

- What are the functions that preserve index permutation ?
- They exist two types of such functions :
  - Equivariant functions  $f$  defined as :

$$\begin{array}{ccc} A \in \mathbb{R}^{n \times n} & \xrightarrow{f} & f(A) \in \mathbb{R}^{n \times n} \\ \downarrow P & & \downarrow P \\ P^T AP \in \mathbb{R}^{n \times n} & \xrightarrow{f} & f(P^T AP) = P^T f(A)P \in \mathbb{R}^{n \times n} \end{array}$$

- Equivariant functions model hidden layers. For example, consider  $A_{ij}$  to be the feature of edge  $ij$  and  $A_{ii}$  to be the feature of node  $i$  :



$$A = \begin{bmatrix} 3 & -3 & 5 & 0 \\ -3 & -1 & 7 & 0 \\ 5 & 7 & 2 & -2 \\ 0 & 0 & -2 & 1 \end{bmatrix} \quad \text{and} \quad f(A) = A11^T$$

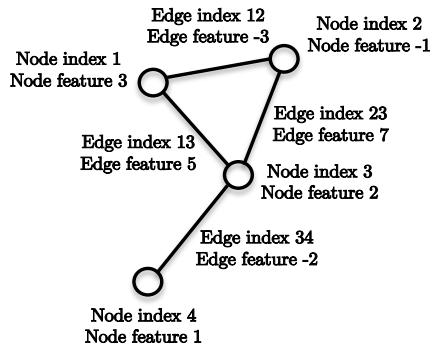
Function  $f$  satisfies :  $f(P^T AP) = P^T f(A)P$

# Invariant functions

- Invariant functions  $g$  defined as :

$$\begin{array}{ccc} A \in \mathbb{R}^{n \times n} & \xrightarrow{g} & g(A) \in \mathbb{R}^K, K \geq 1 \\ \downarrow P & & \downarrow P \\ P^T AP \in \mathbb{R}^{n \times n} & \xrightarrow{g} & g(P^T AP) = g(A) \in \mathbb{R}^K, K \geq 1 \end{array}$$

- Invariant functions are readout layers. For example,



$$A = \begin{bmatrix} 3 & -3 & 5 & 0 \\ -3 & -1 & 7 & 0 \\ 5 & 7 & 2 & -2 \\ 0 & 0 & -2 & 1 \end{bmatrix} \quad \text{and} \quad g(A) = 1^T A 1 = \sum_{ij} A_{ij}$$

Function  $g$  satisfies :  $g(P^T AP) = g(A)$

# Symmetry & equivariance

- The most basic image symmetry is translation.
  - Image content remains unchanged w.r.t. translations.
- What is the function that preserves translations ?
  - Convolution (translation equivariant function)
  - Convolutional neural networks<sup>[1]</sup> are translation equivariant (but not rotation equivariant).
- Generalization<sup>[2,3]</sup> :
  - Group convolution × symmetry (e.g. rotation+translation, mesh convolution)
- Why symmetry/equivariance matter ?
  - Great inductive biases
    - Less parameters to learn.
    - Faster training, better generalization.

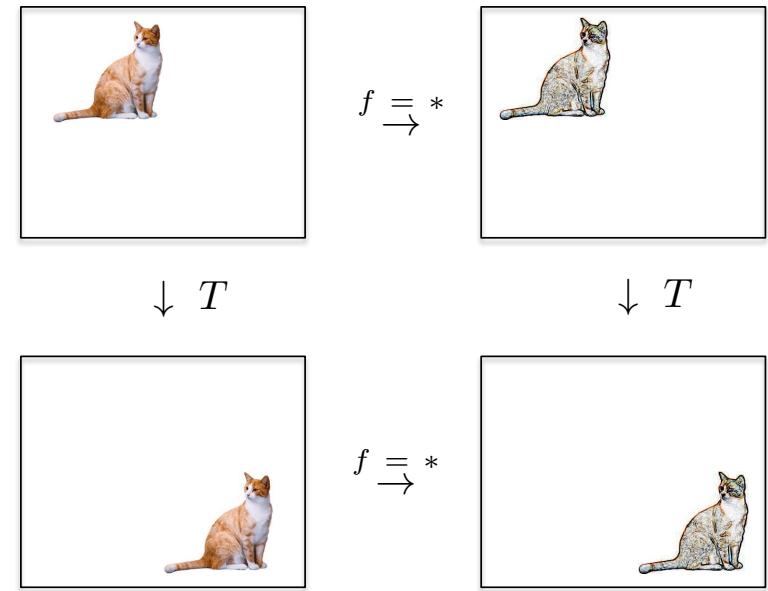


Image symmetry : Translation  
 Equivariant function : Convolution  
 (with a line detector kernel)

$$f(x) = x * k, \quad k = \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

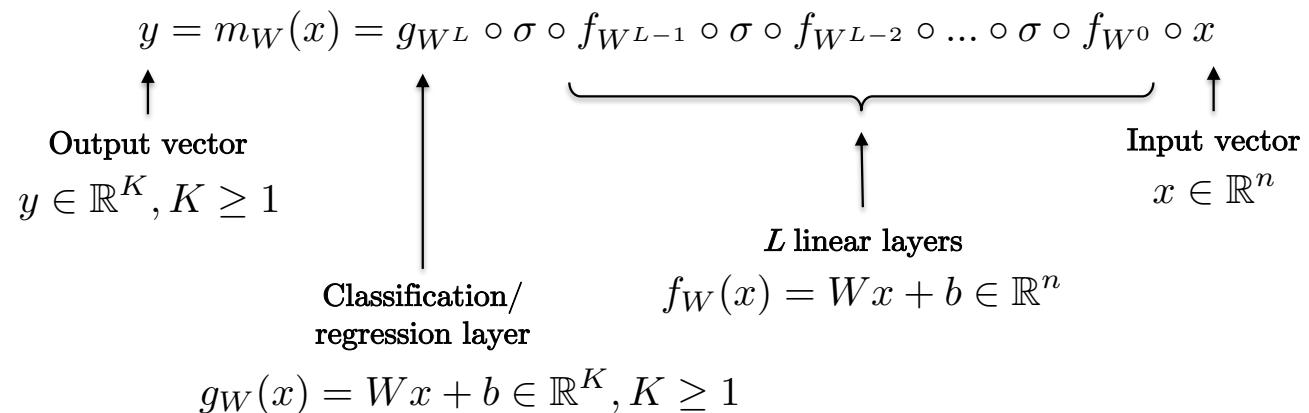
[1] LeCun, Bottou, Bengio, Haffner, Gradient-based learning applied to document recognition, 1998

[2] Cohen, Welling, Group Equivariant Convolutional Networks, 2016

[3] de Haan, Weiler, Cohen, Welling, Gauge Equivariant Mesh CNNs: Anisotropic convolutions on geometric graphs, 2020

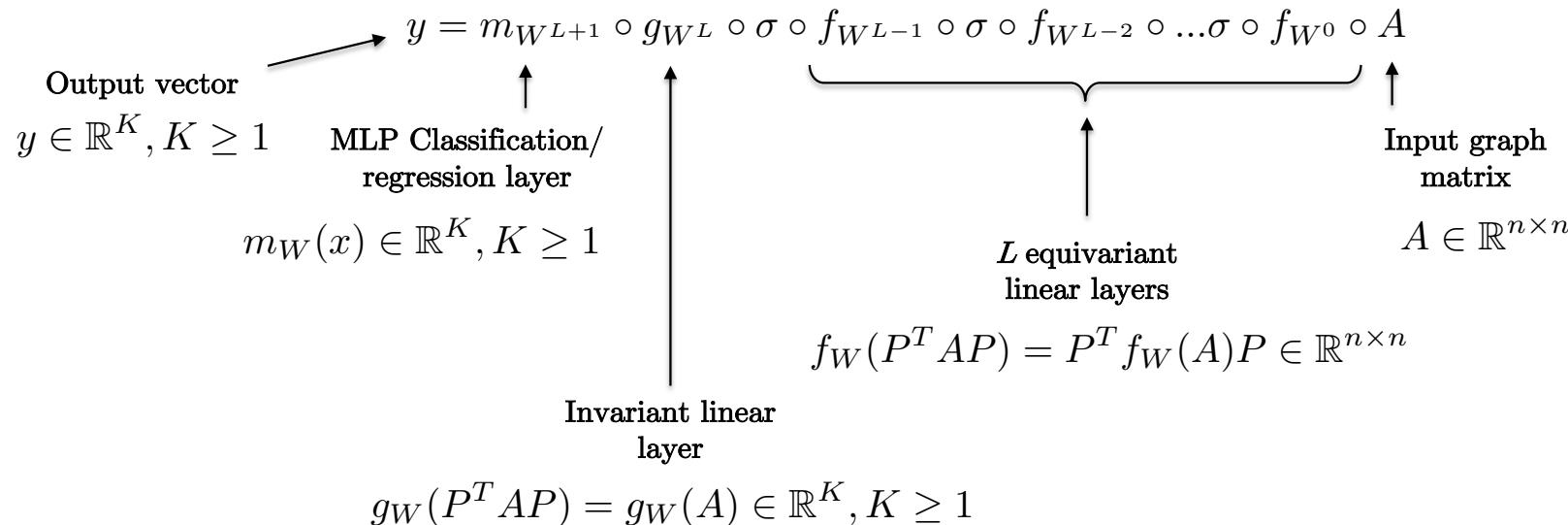
## MLP (for vectors)

- What is the simplest neural network that can be defined for (ordered) vectors ?
  - An MLP : Composition of { linear functions + non-linear activations/ReLU } + final linear layer for classification/regression :



# MLP for graphs<sup>[1]</sup>

- What is the simplest neural network that can be defined for graph adjacency matrices  $A \in \mathbb{R}^{n \times n}$  ?
  - Equivariant GNNs (E-GNNs)<sup>[1]</sup> : Composition of { equivariant linear functions + non-linear activations/ReLU } + invariant linear function + final linear function for classification/regression :



[1] Maron, Ben-Hamu, Shamir, Lipman, Invariant and equivariant graph networks, 2019

# Linear equivariant/invariant spaces

- Given a rank-2 tensor  $A \in R^{n \times n}$ , can we fully characterize all equivariant/invariant linear functions  $f : R^{n \times n} \rightarrow R^{n \times n}$  /  $g : R^{n \times n} \rightarrow R^K$ ?
  - In other words, what is the biggest collection of equivariant/invariant linear layers ?
  - This is to achieve maximal expressivity of linear functions invariant by permutation.
- Theorem<sup>[1]</sup> : The space of all equivariant linear functions is fully characterizable and of dimension  $\text{Bell}(4)=15$ , and the space of all invariant linear functions is of dimension  $\text{Bell}(2)=2$ .
- Note that the number of functions is independent of the number of nodes  $n$ .
  - These functions can be used with datasets of graphs with variable sizes.
- General case<sup>[1]</sup> : Given an order- $m$  tensor  $A \in R^{n^m}$ , the space of all equivariant linear functions  $f : R^{n^m} \rightarrow R^{n^{m'}}$  is of dimension  $\text{Bell}(m+m')$  and the space of all invariant linear functions  $g : R^{n^m} \rightarrow R^K$  is of dimension  $K \cdot \text{Bell}(m)$ .

[1] Maron, Ben-Hamu, Shamir, Lipman, Invariant and equivariant graph networks, 2019

# Linear equivariant functions

- The 15 equivariant linear functions  $f_k : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$  for a given tensor  $A \in \mathbb{R}^{n \times n}$  are defined as :

- The identity and transpose operations:  $f_1(A) = A, f_2(A) = A^T$
- The diag operation:  $f_3(A) = \text{diag}(\text{diag}(A))$
- Sum of rows replicated on rows/ columns/ diagonal:  

$$f_4(A) = A\mathbf{1}\mathbf{1}^T, f_5(A) = \mathbf{1}(A\mathbf{1})^T, f_6(A) = \text{diag}(A\mathbf{1}), \mathbf{1} \in \mathbb{R}^n$$

$$A = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

- Sum of columns replicated on rows/ columns/ diagonal:  

$$f_7(A) = A^T\mathbf{1}\mathbf{1}^T, f_8(A) = \mathbf{1}(A^T\mathbf{1})^T, f_9(A) = \text{diag}(A^T\mathbf{1})$$

$$f_1(A) = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix} \quad f_2(A) = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

- Sum of all elements replicated on all matrix/ diagonal:  

$$f_{10}(A) = (\mathbf{1}^T A \mathbf{1}) \cdot \mathbf{1}\mathbf{1}^T, f_{11}(A) = (\mathbf{1}^T A \mathbf{1}) \cdot \text{diag}(\mathbf{1})$$

$$f_3(A) = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad f_4(A) = \begin{bmatrix} 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 3 & 3 & 3 & 3 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

- Sum of diagonal elements replicated on all matrix/diagonal:  

$$f_{12}(A) = (\mathbf{1}^T \text{diag}(A)) \cdot \mathbf{1}\mathbf{1}^T, f_{13}(A) = (\mathbf{1}^T \text{diag}(A)) \cdot \text{diag}(\mathbf{1})$$

$$f_5(A) = \begin{bmatrix} 2 & 2 & 3 & 1 \\ 2 & 2 & 3 & 1 \\ 2 & 2 & 3 & 1 \\ 2 & 2 & 3 & 1 \end{bmatrix} \quad f_6(A) = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

- Replicate diagonal elements on rows/columns:  

$$f_{14}(A) = \text{diag}(A)\mathbf{1}^T, f_{15}(A) = \mathbf{1}\text{diag}(A)^T$$

$$f_7(A) = \begin{bmatrix} 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 3 & 3 & 3 & 3 \\ 1 & 1 & 1 & 1 \end{bmatrix} \quad f_8(A) = \begin{bmatrix} 2 & 2 & 3 & 1 \\ 2 & 2 & 3 & 1 \\ 2 & 2 & 3 & 1 \\ 2 & 2 & 3 & 1 \end{bmatrix}$$

etc

# Equivariant layer parametrization

- From 1-dim input to 1-dim output :

$$A \in \mathbb{R}^{n \times n} \rightarrow f_W(A) = \sum_{k=1}^{15+2} W_k f_k(A) \in \mathbb{R}^{n \times n}$$

15 equivariant functions      2 bias functions  
 ↓                          ↓  
 15+2                      17 (learnable) parameters  
 ↓                          ↓

- From d-dim input to 1-dim output :

$$A \in \mathbb{R}^{n \times n \times d} \rightarrow f_W(A) = \sum_{k=1}^{15+2} \sum_{q=1}^d W_{k,q} f_k(A_{\cdot,\cdot,q}) \in \mathbb{R}^{n \times n}$$

17 x d (learnable) parameters      Sliced tensor  $A$  along the third dimension.  
 ↓                          ↓  
 17 x d                      17 (learnable) parameters

- From d-dim input to d'-dim output :

$$A \in \mathbb{R}^{n \times n \times d} \rightarrow (f_W(A))_{\cdot,\cdot,q'} = \sum_{k=1}^{15+2} \sum_{q=1}^d W_{k,q,q'} f_k(A_{\cdot,\cdot,q}) \in \mathbb{R}^{n \times n \times d'}$$

17 x d x d' (learnable) parameters  
 ↓                          ↓

- Memory/speed complexity is  $O(n^2)$ .

# Equivariant GNNs<sup>[1]</sup>

- Start with :

$$A \in \mathbb{R}^{n \times n \times (1+d+d_e)}$$

↓  
 Adjacency matrix      ↓  
 d-dim node features

↑  
 d<sub>e</sub>-dim edge features

- Forward pass of E-GNNs :

$$y = \underbrace{m_{W^{L+1}}}_{\text{MLP layer}} \circ \underbrace{g_{W^L}}_{\text{Invariant layer}} \circ \sigma \circ \underbrace{f_{W^{L-1}} \circ \sigma \circ f_{W^{L-2}} \circ \dots \sigma \circ \circ f_{W^0}}_{L \text{ Equivariant layers}} \circ A$$

with  $f_{W^\ell} : \mathbb{R}^{n \times n \times d^\ell} \rightarrow \mathbb{R}^{n \times n \times d^{\ell+1}}$   
 $g_{W^L} : \mathbb{R}^{n \times n \times d^L} \rightarrow \mathbb{R}^K$

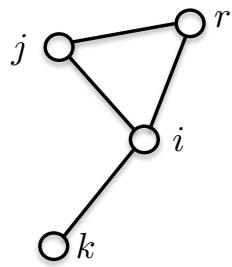
- #Parameters/layer is :  $17 \times d^\ell \times d^{\ell+1}$
- Memory/speed complexity is  $O(n^2)$ .
- E-GNNs are the most expressive linear networks for graphs w.r.t. the original 1-WL test.
  - However, the 1-WL test fails to distinguish simple non-isomorphic graphs.

[1] Maron, Ben-Hamu, Shamir, Lipman, Invariant and equivariant graph networks, 2019

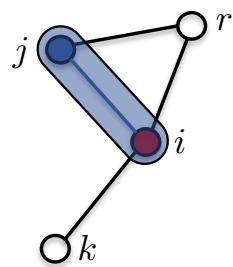
[2] Gilmer, Schoenholz, Riley, Vinyals, Dahl, Neural message passing for quantum chemistry, 2017

## Boosting the expressivity power

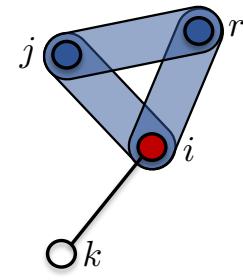
- Can we design more expressive GNNs, i.e. GNNs that are better than the original WL test ?
  - k-WL tests : Original test uses 2-tuple of nodes to produce colors. To produce more colors, and improve expressivity power of WL tests, higher-order interactions between nodes with k-tuple of nodes with  $k \geq 3$  can be used.



$$\mathcal{G} = (V, E)$$



$$\text{Edges} = \{i, j\}, \{i, r\}, \{i, k\}$$



$$\begin{aligned} \text{Hyper-edges} &= \{i, j, r\}, \\ &\{i, k, r\}, \{i, j, k\} \end{aligned}$$

2-tuple of nodes can distinguish non-isomorphic graphs with the 1-WL/2-WL tests.

3-tuple of nodes can distinguish non-isomorphic graphs with the 3-WL test.

# Higher expressivity power

- What is the expressivity power of equivariant GNNs in terms of WL tests ?
  - Let us define a k-order Equivariant GNNs :

$$y = m_{WL+1} \circ g_{WL} \circ \sigma \circ f_{WL-1} \circ \sigma \circ f_{WL-2} \circ \dots \sigma \circ f_{W^0} \circ A$$

$$\text{with } k = \max_{\ell \in [0, L-1]} k_\ell$$

$$\text{and } f_{W^\ell} : \mathbb{R}^{n^{k_\ell} \times d_\ell} \rightarrow \mathbb{R}^{n^{k_\ell+1} \times d_{\ell+1}}$$

- Theorem<sup>[1]</sup> : There exists a k-order E-GNN that can distinguish non-isomorphic graphs with the k-WL test.
  - Existence result (although SGD is not guaranteed to find such network).
  - It is not a practical result as k-order E-GNNs require  $O(n^k)$  memory/speed complexities.
    - Note that 1-WL/2-WL GNNs (s.a. GINs) have the same discriminative power<sup>[2]</sup>. It implies that we must have at least  $k=3$  to be more powerful than GINs. But then, it requires  $O(n^3)$ .

[1] Maron, Ben-Hamu, Shamir, Lipman, Invariant and equivariant graph networks, 2019

[2] Cai, Furer, Immerman, An optimal lower bound on the number of variables for graph identification, 1992

## 3-WL GNNs<sup>[1]</sup>

- How to design GNNs that are provable 3-WL but do not require  $O(n^3)$  memory/speed complexities ?
- 3-WL GNNs<sup>[1]</sup> : To achieve higher interactions between nodes, it was proposed to multiply matrices (that generates non-local interactions).
  - Theorem : There exists a 3-WL GNN as expressive as the 3-WL test.
  - Memory is quadratic  $O(n^2)$  but matrix-matrix multiplication implies  $O(n^3)$  speed.
    - It improves memory complexity compared to 3-order E-GNNs<sup>[2]</sup>.
    - Matrix-matrix multiplication densifies sparse matrix.

[1] Maron, Ben-Hamu, Serviansky, Lipman, Provably powerful graph networks, 2019

[2] Maron, Ben-Hamu, Shamir, Lipman, Invariant and equivariant graph networks, 2019

## 3-WL GNNs

- 3-WL GNN update layer :

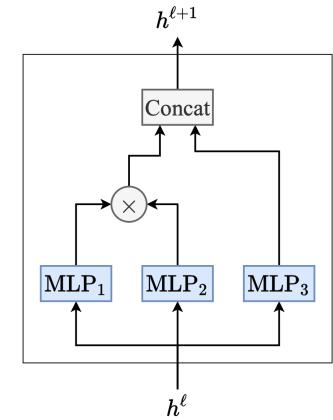
$$h^{\ell+1} = \text{Concat}(m_{W_1^\ell}(h^\ell) \cdot m_{W_2^\ell}(h^\ell), m_{W_3^\ell}(h^\ell)),$$

where  $h^\ell, h^{\ell+1} \in \mathbb{R}^{n \times n \times d}, W_a, W_b \in \mathbb{R}^{d \times d}$ ,

where  $m_W$  are 2-layer MLPs applied along the feature dimension. The matrix-matrix multiplication is carried out along the first and second dimensions as :

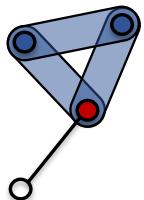
$$(M_{W_1}(h) \cdot M_{W_2}(h))_{i,j,k} = \sum_{p=1}^n (M_{W_1}(h))_{i,p,k} \cdot (M_{W_2}(h))_{p,j,k}$$

with complexity  $O(n^3)$ .



## Sparse k-WL-GNNs<sup>[1]</sup>

- Original k-WL GNNs are provable as expressive as the k-WL test, but they are not practical because of  $O(n^k)$  memory/speed requirements.
  - Dense rank-k tensors represent all possible k-tuple of nodes (existing or not).
  - These GNNs do not leverage graph sparsity, which is known to be a good inductive bias for generalization. Also, dense tensors increase over-fitting.
- Sparse k-WL-GNNs<sup>[1]</sup> solely consider the k-tuple of nodes present in the given graphs.
  - Theorem<sup>[1]</sup> : Sparse k-WL-GNNs are strictly more powerful than dense k-WL-GNNs.
  - Theorem<sup>[1,2]</sup> : Sparse k-WL-GNNs exhibit a smaller generalization error than dense k-WL-GNNs.
  - In practice, the memory and speed complexities are improved, and benefit from the graph sparsity. However, the numerical time still remains large as all k-tuple of nodes present in the graph must be considered at each layer.

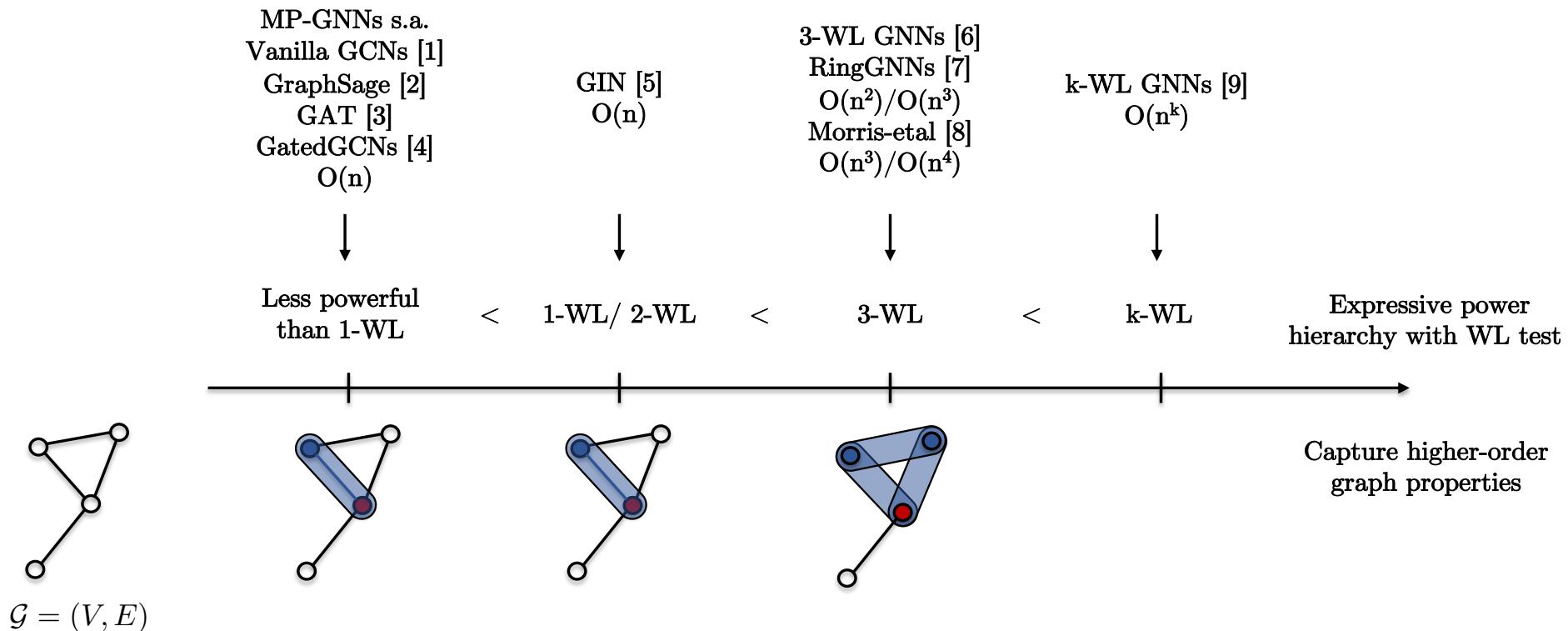


[1] Morris, Rattan, Mutzel, Weisfeiler and Leman go sparse: Towards scalable higher-order graph embeddings, 2020  
[2] Garg, Jegelka, Jaakkola, Generalization and representational limits of graph neural networks, 2020

# Outline

- Weisfeiler-Lehman isomorphism test
- Graph isomorphism networks
- Equivariant graph neural networks
- **Expressivity power and generalization**
- Equivariant set neural networks
- Conclusion

# Expressivity power w.r.t WL test



- [1] Kipf, Welling, Semi-supervised classification with graph convolutional networks, 2017
- [2] Hamilton, Ying, Leskovec, Inductive representation learning on large graphs, 2017
- [3] Velickovic, Cucurull, Casanova, Romero, Lio, Bengio, Graph attention networks, 2017
- [4] Bresson, Laurent, Residual gated graph convnets, 2017
- [5] Xu, Hu, Leskovec, Jegelka, How powerful are graph neural networks?, 2019
- [6] Maron, Ben-Hamu, Serviansky, Lipman, Provably powerful graph networks, 2019
- [7] Chen, Villar, Chen, Bruna, On the equivalence between graph isomorphism testing and function approximation with gnn, 2019
- [8] Morris, Ritzert, Fey, Hamilton, Lenssen, Rattan, Grohe, Weisfeiler and leman go neural: Higher-order graph neural networks, 2019
- [9] Maron, Ben-Hamu, Shamir, Lipman, Invariant and equivariant graph networks, 2019

# Universal approximator

- What functions can be approximated by MLP ?
  - Universal approximation theorem<sup>[1,2]</sup> : Any continuous function can be arbitrarily approximated with a MLP, with the necessary condition that the number of hidden neurons goes to infinity (or very large).
- What functions can be approximated by Equivariant GNNs ?
  - Theorem<sup>[3]</sup> : Any continuous function invariant by permutation can be arbitrarily approximated by E-GNNs, with the necessary condition the network has a tensor of order  $m=\text{poly}(n)=n(n-1)/2$ , where  $n$  is the number of nodes.
  - High-order tensors are not practical as memory/speed complexities grow as  $O(n^m)$ .
- Result in [3] is for the worst case and was improved in [4] to be independent of  $n$  (but no upper bound on the order is known).

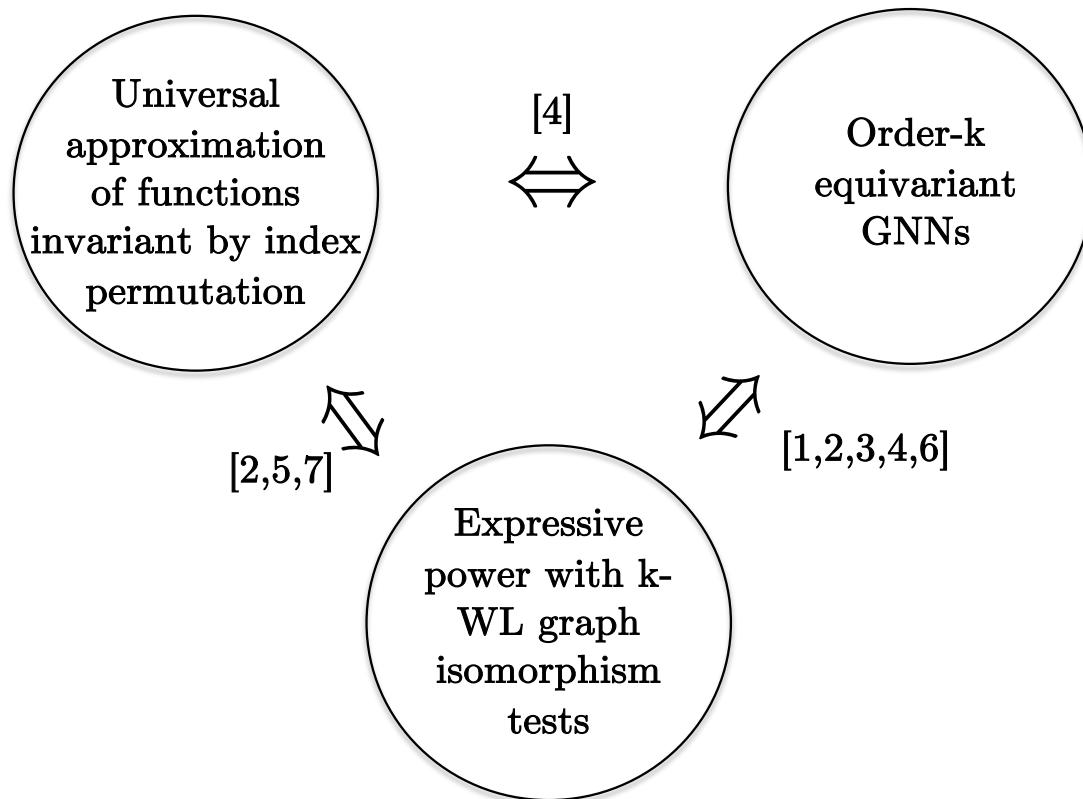
[1] Cybenko, Approximation by superpositions of a sigmoidal function, 1989

[2] Hornik, Approximation capabilities of multilayer feedforward networks, 1991

[3] Maron Fetaya, Segol, Lipman, On the universality of invariant networks, 2019

[4] Keriven, Peyré, Universal invariant and equivariant graph neural networks, 2019

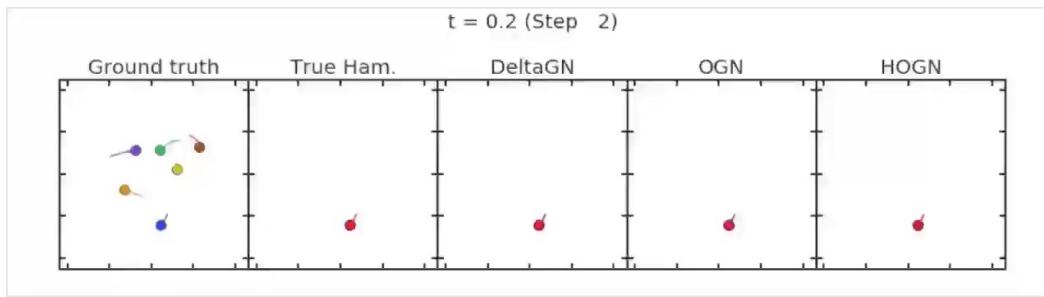
# Relationships



- [1] Xu, Hu, Leskovec, Jegelka, How powerful are graph neural networks?, 2019
- [2] Maron Fetaya, Segol, Lipman, On the universality of invariant networks, 2019
- [3] Keriven, Peyré, Universal invariant and equivariant graph neural networks, 2019
- [4] Maron, Ben-Hamu, Serviansky, Lipman, Provably powerful graph networks, 2019
- [5] Chen, Villar, Chen, Bruna, On the equivalence between graph isomorphism testing and function approximation with gnn, 2019
- [6] Morris, Ritzert, Fey, Hamilton, Lenssen, Rattan, Grohe, Weisfeiler and leman go neural: Higher-order graph neural networks, 2019
- [7] Loukas, What graph neural networks cannot learn: depth vs width, 2019

# Universal approximator and generalization

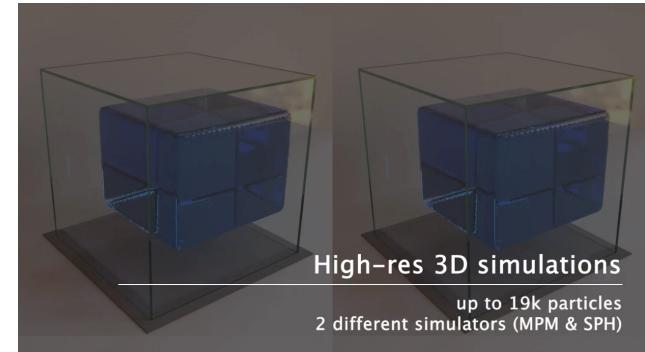
- Universal approximator for functions defined on graphs and expressiveness/representation power to distinguish non-isomorphic graphs are necessary conditions to design good GNNs.
- However, these properties do not guarantee the best generalization performance. It is important to make the distinction between representation power and generalization!
  - MLP are universal approximators for general functions but they do not generalize well.
  - Equivariant GNNs are MLP for permutation-invariant functions but they do not generalize as well as MP-GNNs s.a. GAT/GatedGCNs<sup>[1]</sup>.
- MP-GNNs which benefit from tailored inductive biases can lead to strong generalization performance.
  - For example, combining MP-GNNs with physical constraints like Hamiltonian mechanics<sup>[2,3]</sup>.



[1] Dwivedi, Joshi, Laurent, Bengio, Bresson, Benchmarking graph neural networks, 2020

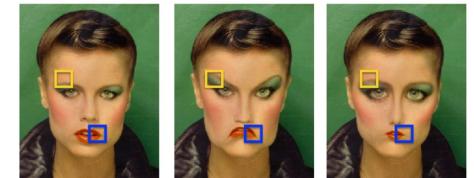
[2] Sanchez-Gonzalez, Bapst, Cranmer, Battaglia, Hamiltonian Graph Networks with ODE Integrators, 2019

[3] Sanchez-Gonzalez, Godwin, Pfaff, Ying, Leskovec, Battaglia, Learning to simulate complex physics with graph networks, 2020

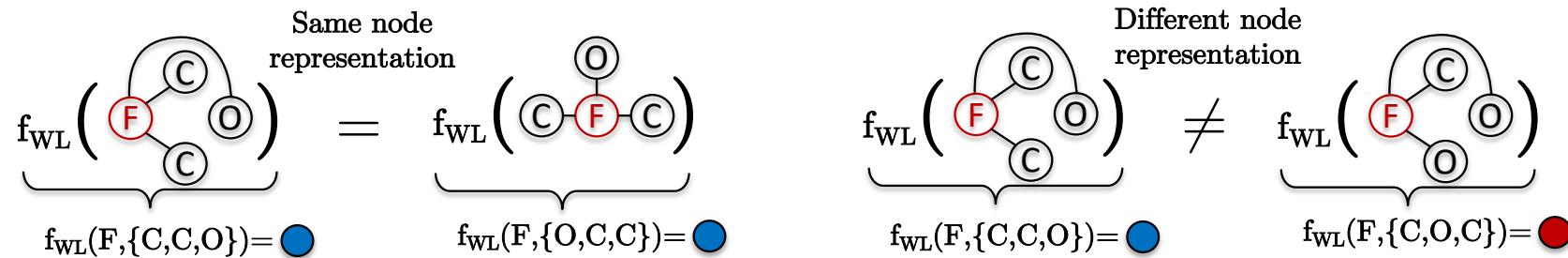


# Question on isomorphism expressivity

- CNNs are robust to local deformation invariance (diffeomorphism on grids). This property is critical for generalization.
- Are WL-GNNs robust to local graph deformations?
  - The WL test focused on finding distinct node representation for different neighborhoods.
  - Two nodes must have two different representations if their neighborhoods are distinct just by a single neighbor.
  - But a representation robust to perturbations (i.e. unchanged or very close) is desirable for generalization.



Reproduced from Bruna's NYU lectures

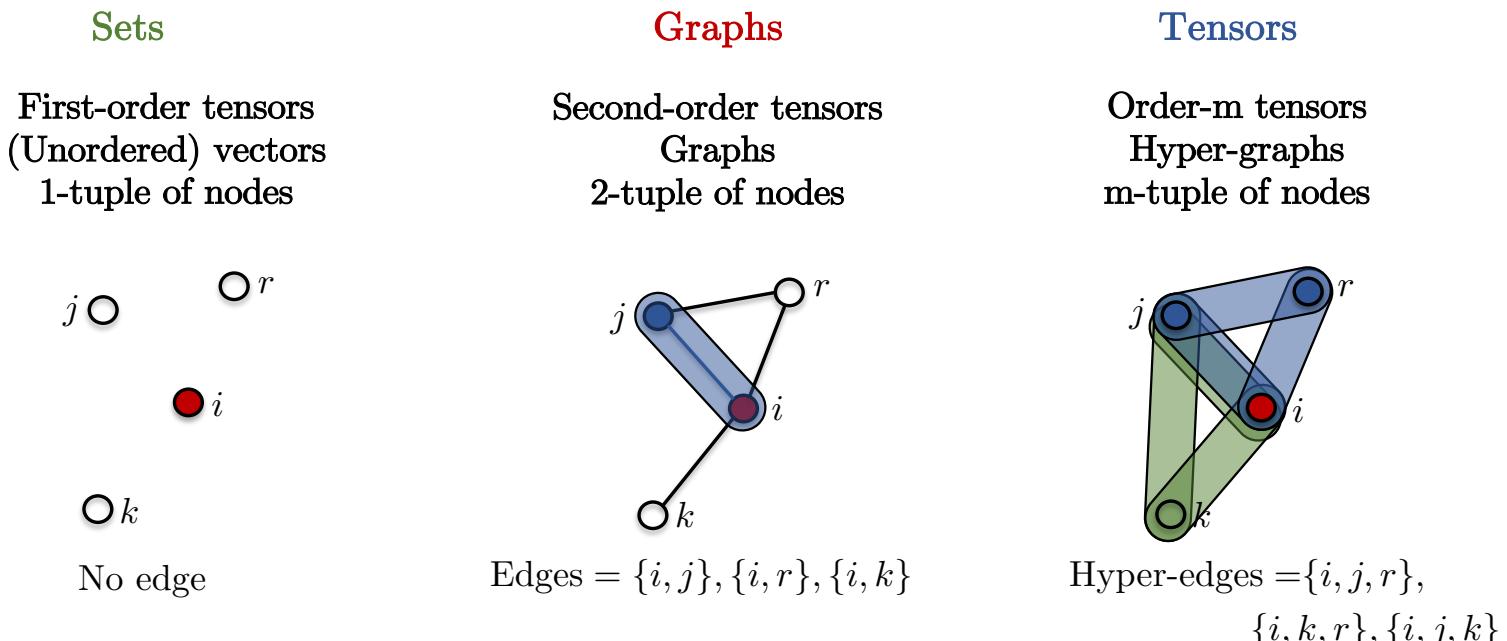


# Outline

- Weisfeiler-Lehman isomorphism test
- Graph isomorphism networks
- Equivariant graph neural networks
- Expressivity power and generalization
- **Equivariant set neural networks**
- Conclusion

# Beyond graphs

- Tensors invariant by permutation :



- Equivariant GNNs<sup>[1]</sup> can be run on any tensor of order- $m$ , with memory/speed complexities  $O(n^{\max(m,m')})$ .

$$f_W : \mathbb{R}^{n^m} \rightarrow \mathbb{R}^{n^{m'}}$$

$$g_W : \mathbb{R}^{n^m} \rightarrow \mathbb{R}^K, K \geq 1$$

[1] Maron, Ben-Hamu, Shamir, Lipman, Invariant and equivariant graph networks, 2019

## Equivariant/invariant linear layer for sets

- For rank-1 tensors/vectors, the equivariant linear transformation  $f_W : \mathbb{R}^n$  to  $\mathbb{R}^n$  is fully characterized by  $\text{bell}(2)=2$  equivariant functions, which are  $I$  and  $11^T - I$ . The parametrized equivariant linear transformation is<sup>[1,2,3]</sup> :

$$f_W(x) = [w_1 I + w_2(11^T - I)]x \in \mathbb{R}^n$$

$$f_W(x)_i = w_1 x_i + w_2 \sum_{j \neq i} x_j$$

$$f_W \left( P \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) = f_W \left( \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) = f_W \left( \begin{bmatrix} x_3 \\ x_1 \\ x_2 \end{bmatrix} \right) = \begin{bmatrix} w_1 x_3 + w_2(x_1 + x_2) \\ w_1 x_1 + w_2(x_2 + x_3) \\ w_1 x_2 + w_2(x_1 + x_3) \end{bmatrix} = P f_W \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right)$$

- The invariant linear transformation  $g_W : \mathbb{R}^n$  to  $\mathbb{R}^K$  is fully characterized by  $K \cdot \text{bell}(1)=K$  equivariant function, which is  $1_{K \times n}$ . The parametrized invariant linear transformation is

$$g_W(x) = w_1 1_{K \times n} x \in \mathbb{R}^K, K \geq 1$$

- Linear equivariant functions for sets have limited performance (only 2 parameters to learn per layer), and thus cannot be performant for NLP tasks (sentences are rank-1 tensors).

[1] Zaheer, Kottur, Ravanbakhsh, Poczos, Salakhutdinov, Smola, Deep sets, 2017

[2] Qi, Su, Mo, Guibas, Pointnet: Deep learning on point sets for 3d classification and segmentation, 2017

[3] Maron, Ben-Hamu, Shamir, Lipman, Invariant and equivariant graph networks, 2019

# Equivariant/invariant linear layer for sets

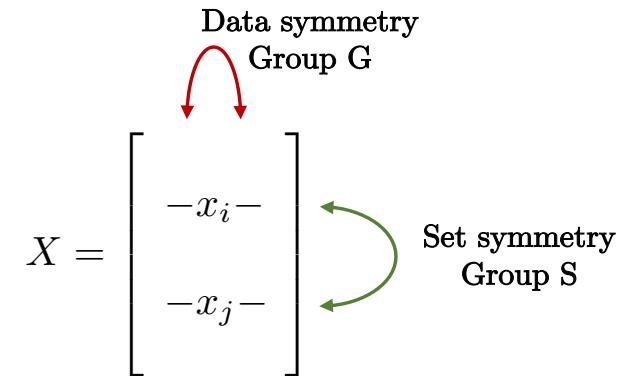
- How to improve the performance of linear equivariant functions for sets ?
  - Consider non-linear equivariant functions.
    - Transformers<sup>[1]</sup> use (non-linear) attention functions to produce node representations that adapt to the context.
    - Breakthrough in NLP
  - Consider additional symmetry/invariance<sup>[2,3]</sup>.

$$h_i^{\ell+1} = f_W(A) = \sum_{j \in V} \frac{\exp(W_1 h_j^\ell)}{\sum_{k \in V} \exp(W_2 h_k^\ell)} W_3 h_j^\ell$$

[1] Vaswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser, Polosukhin, Attention is all you need, 2017  
[2] Zaheer, Kottur, Ravanbakhsh, Poczos, Salakhutdinov, Smola, Deep sets, 2017  
[3] Maron, Litany, Chechik, Fetaya, On Learning Sets of Symmetric Elements, 2020

# On Learning Sets of Symmetric Elements<sup>[1]</sup>

- Neural networks for sets of symmetric data.
  - Set symmetry : order invariance
  - Data symmetry : translation invariance for example
- A set of data can be represented as a matrix  $X \in \mathbb{R}^{n \times d}$
- DeepSets<sup>[2]</sup>/PointNets<sup>[3]</sup> : Neural network for arbitrary sets of data points (order invariance only).
  - These set networks are universal for permutation invariant functions. DeepSets is also universal for equivariant functions<sup>[4]</sup>.
- ConvNets<sup>[5]</sup> : Neural network for symmetric data (translation invariance only)



$$f_W^S(X)_i = f_{W_1}^L(x_i) + f_{W_2}^L(\sum_{j \neq i} x_j)$$

$$f_{W_1}^L(z) = w_1 z, \quad f_{W_2}^L(z) = w_2 z$$

$$f_W^G(x_i) = W * x_i$$

[1] Maron, Litany, Chechik, Fetaya, On Learning Sets of Symmetric Elements, 2020

[2] Zaheer, Kottur, Ravanbakhsh, Poczos, Salakhutdinov, Smola, Deep sets, 2017

[3] Qi, Su, Mo, Guibas, Pointnet: Deep learning on point sets for 3d classification and segmentation, 2017

[4] Segol, Lipman, On universal equivariant set networks, 2019

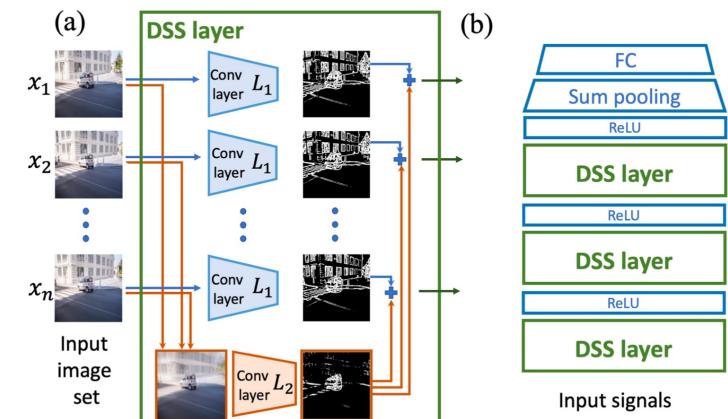
[5] Maron, Litany, Chechik, Fetaya, On Learning Sets of Symmetric Elements, 2020

# On Learning Sets of Symmetric Elements

- Deep Sets for Symmetric elements (DSS)<sup>[1]</sup> : Combining symmetries from DeepSets and ConvNets into a single framework.
- Any linear equivariant layer for the symmetry group  $S \times G$  is of the form :
- Expressivity power : If  $G$ -equivariant networks are universal, then  $S \times G$ -equivariant networks are also universal, that is it can represent any function with the  $S \times G$  symmetry.
  - Note that DSS expressivity is lower than DeepSets (symmetry group is larger) but DSS are as expressive as they can be for this symmetry group.
  - Experimental generalization performance is better with this symmetry group.
- Applications : set of 1D-signals, 3D cloud points and set of images.

$$f_W^{S \times G}(X)_i = f_{W_1}^G(x_i) + f_{W_2}^G(\sum_{j \neq i} x_j),$$

$$f_W^G(x) = W * x$$



# Outline

- Weisfeiler-Lehman isomorphism test
- Graph isomorphism networks
- Equivariant graph neural networks
- Expressivity power and generalization
- Equivariant set neural networks
- Conclusion

# Conclusion

- WL test is used as expressivity measure of representation power of GNNs, but it is not the only expressivity measure. For example, we can define expressivity as the capacity of identifying/counting the number of elementary sub-structures like cycles and cliques<sup>[1]</sup>.
- WL-GNNs are computationally more expensive<sup>[2]</sup> than standard MP-GNNs, i.e.  $O(n^k)$  vs  $O(E)$ .
- Expressivity does not necessarily imply generalization<sup>[3]</sup>, like MLP.
- Equivariance/invariance w.r.t. group of transformations can be powerful to generalize if it captures the right data properties, otherwise it will (strongly) overfit.

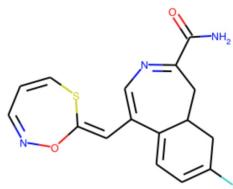


Fig 1. Graph with cycles (rings)

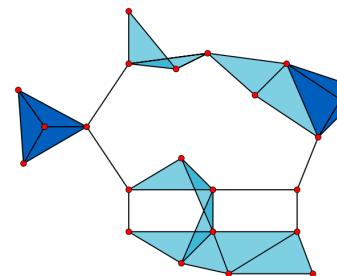


Fig 2. Graph with  
19 × 3-vertex cliques (light and dark blue triangles)  
and 2 × 4-vertex cliques (dark blue areas)

[1] Chen, Chen, Villar, Bruna, Can graph neural networks count substructures? 2020

[2] Maron, Ben-Hamu, Serviansky, Lipman, Provably powerful graph networks, 2019

[3] Dwivedi, Joshi, Laurent, Bengio, Bresson, Benchmarking graph neural networks, 2020



Questions?