

The (expressive) power of graph learning

Floris Geerts (University of Antwerp)

Course

- ❖ Is about recent advances in **graph learning**.
- ❖ With an emphasis on the **expressive power of learning methods**.
 - ❖ **Self-contained** (too some extent).
 - ❖ Mostly **high-level**, but also **low-level**, so basically **all levels**.
 - ❖ **Not all** methods or related works are covered.
 - ❖ Emphasis on **theoretical aspects**.

About the speaker

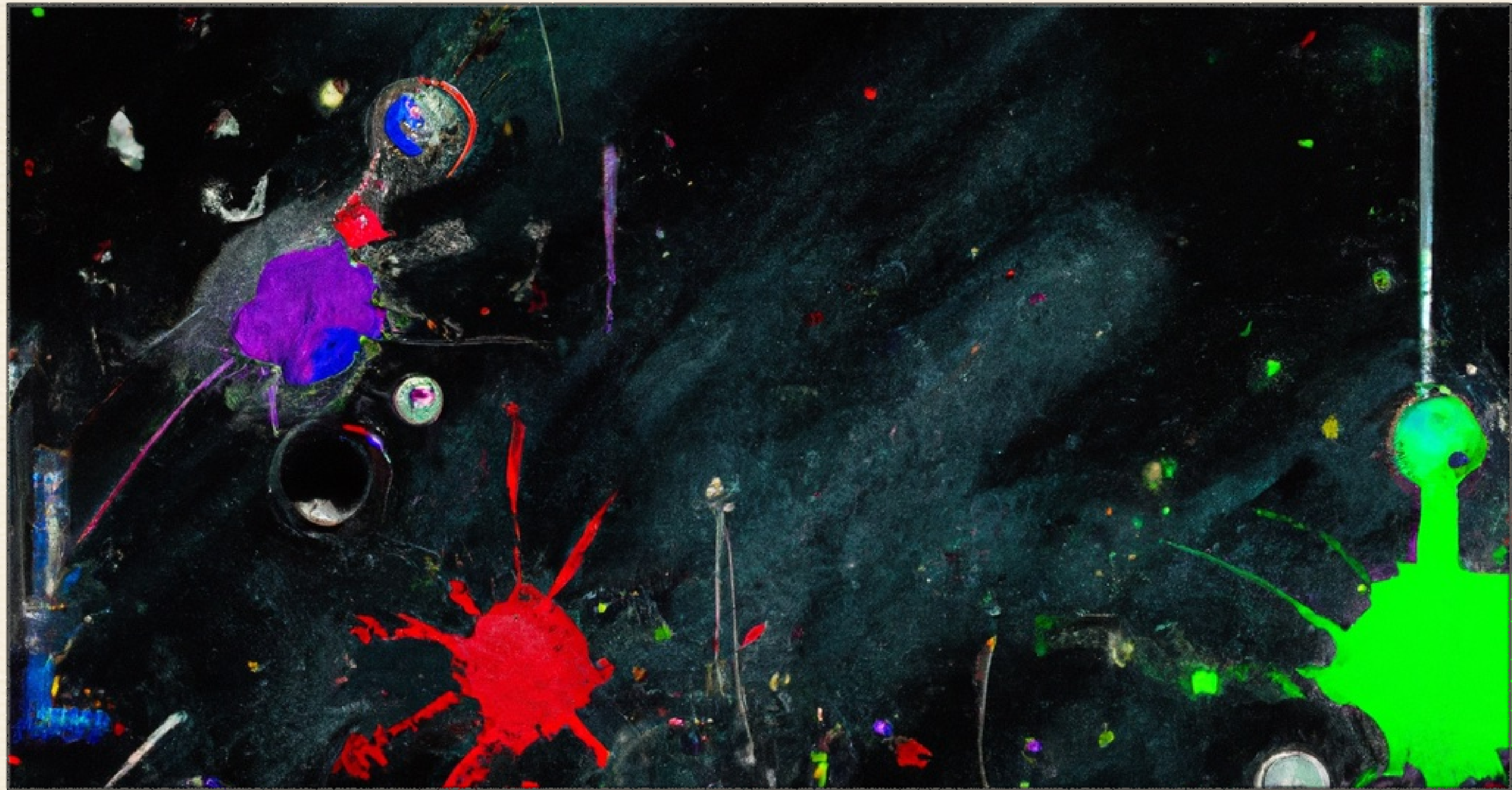
- ❖ Background in mathematics, database theory and expressive power of query languages.
- ❖ Since 2018, expressive power of linear algebra.
- ❖ Natural move to the study of expressive power of graph neural networks.
- ❖ Current focus is on generalisation and relational learning.

Outline

- ❖ Graph learning and expressive power
- ❖ Message Passing Neural Networks
- ❖ Boosting power:
 - ❖ Feature augmentation
 - ❖ Subgraphs
 - ❖ Higher-order message-passing



Ask Questions

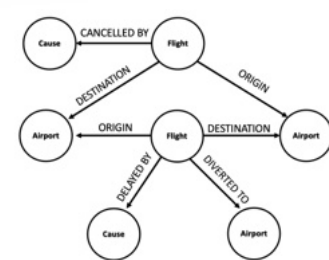


Graph learning

What is it?

Why learning on graphs?

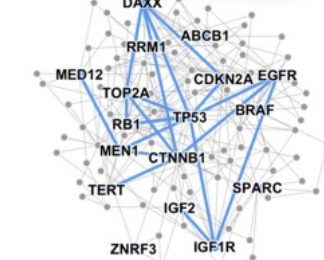
Graphs are everywhere!



Event Graphs



Computer Networks



Disease Pathways

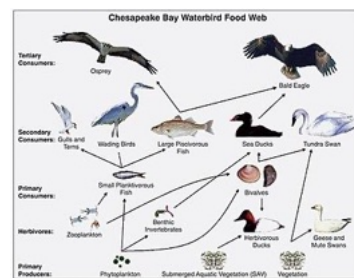


Image credit: Wikipedia

Food Webs



Image credit: Pinterest

Particle Networks

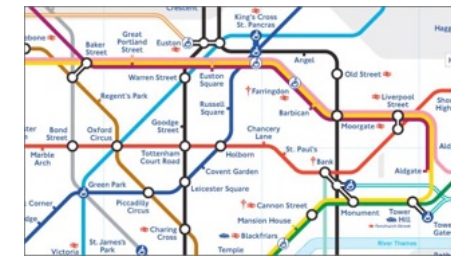


Image credit: visitlondon.com

Underground Networks



Image credit: Medium

Social Networks

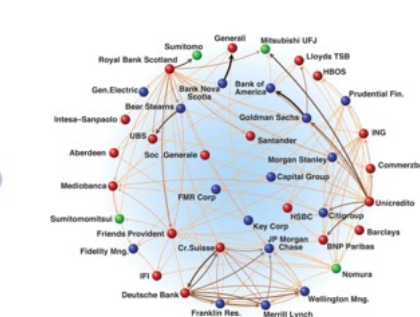


Image credit: Science

Economic Networks

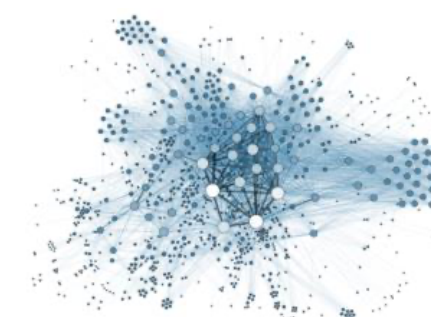


Image credit: Lumen Learning

Communication Networks



Citation Networks



Image credit: Missoula Current News

Internet

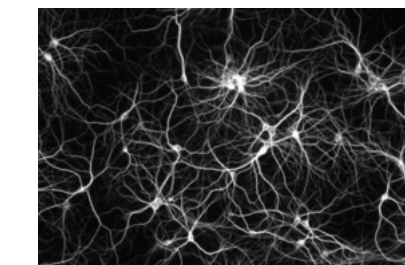


Image credit: The Conversation

Networks of Neurons

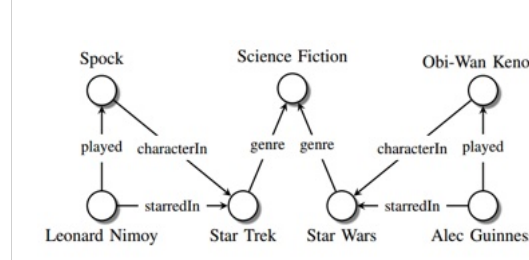


Image credit: Maximilian Nickel et al

Knowledge Graphs

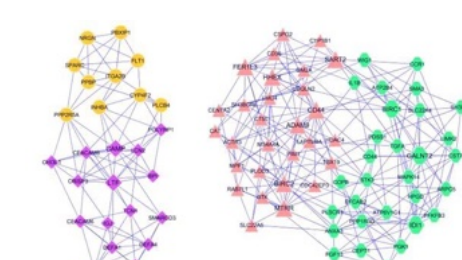


Image credit: ese.wustl.edu

Regulatory Networks

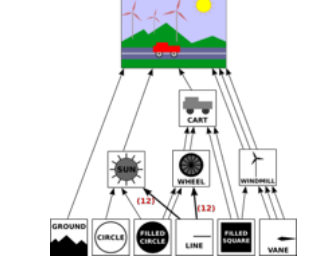


Image credit: math.hws.edu

Scene Graphs

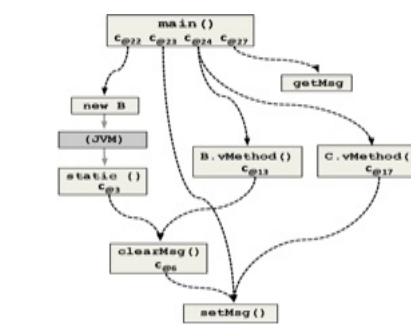


Image credit: ResearchGate

Code Graphs

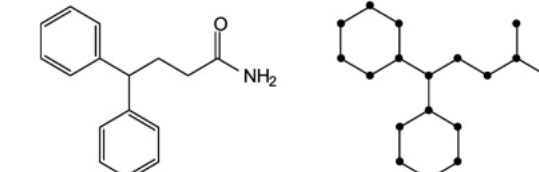


Image credit: MDPI

Molecules

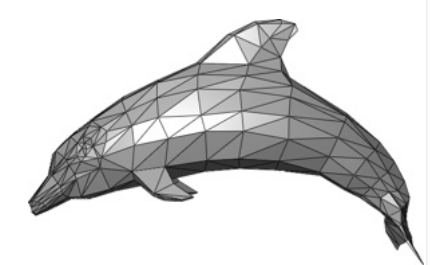


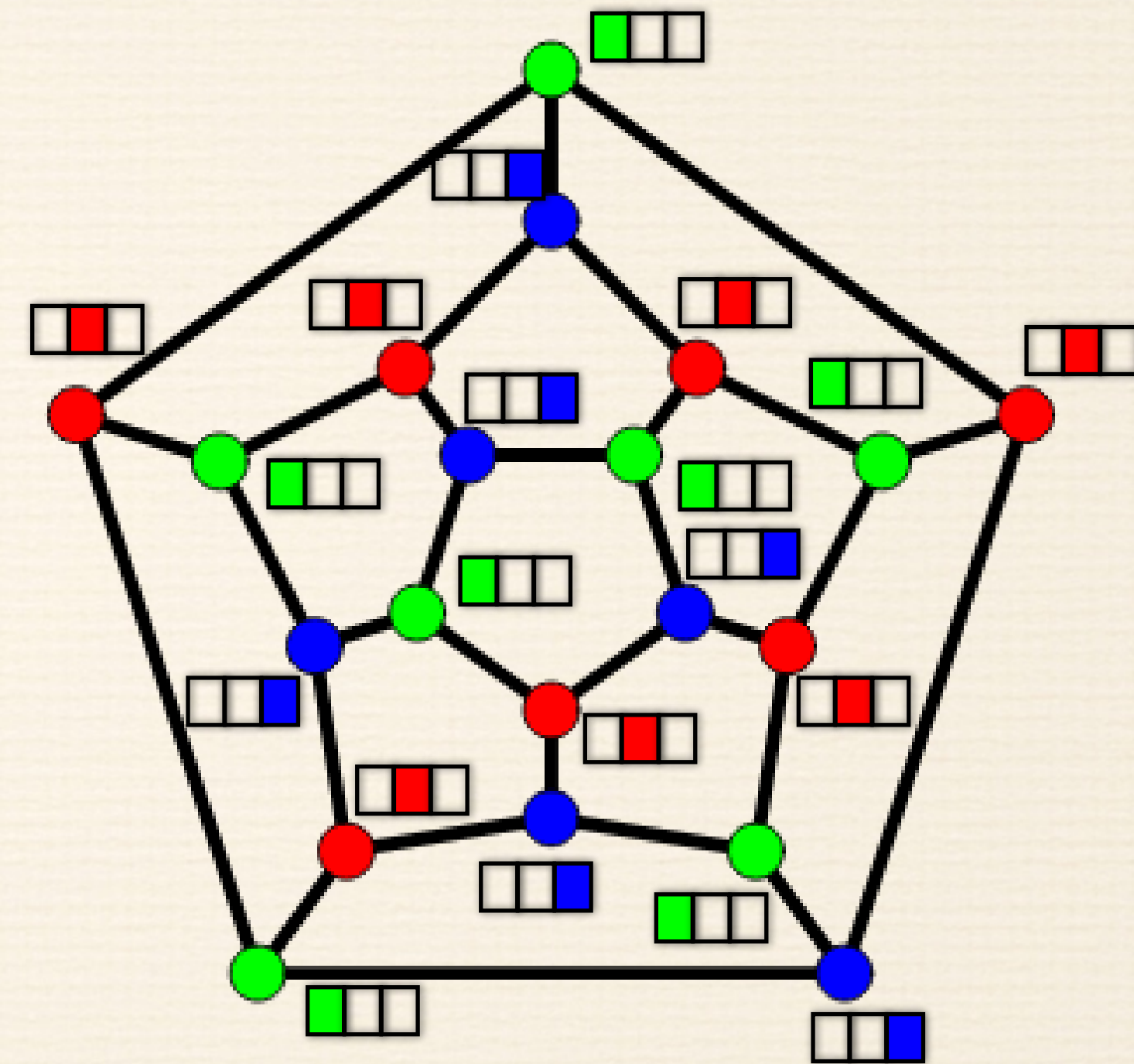
Image credit: Wikipedia

3D Shapes

Graphs: One definition to rule them all

- ❖ Graph $G = (V_G, E_G, L_G)$ with
- ❖ **Vertex** set V_G
- ❖ **Edge** set $E_G \subseteq V_G^2 := V_G \times V_G$
- ❖ **Vertex labels:** $L_G : V_G \rightarrow \Sigma$

Vertex features \mathbb{R}^d



Hot-one encoding

Adjacency matrix representation

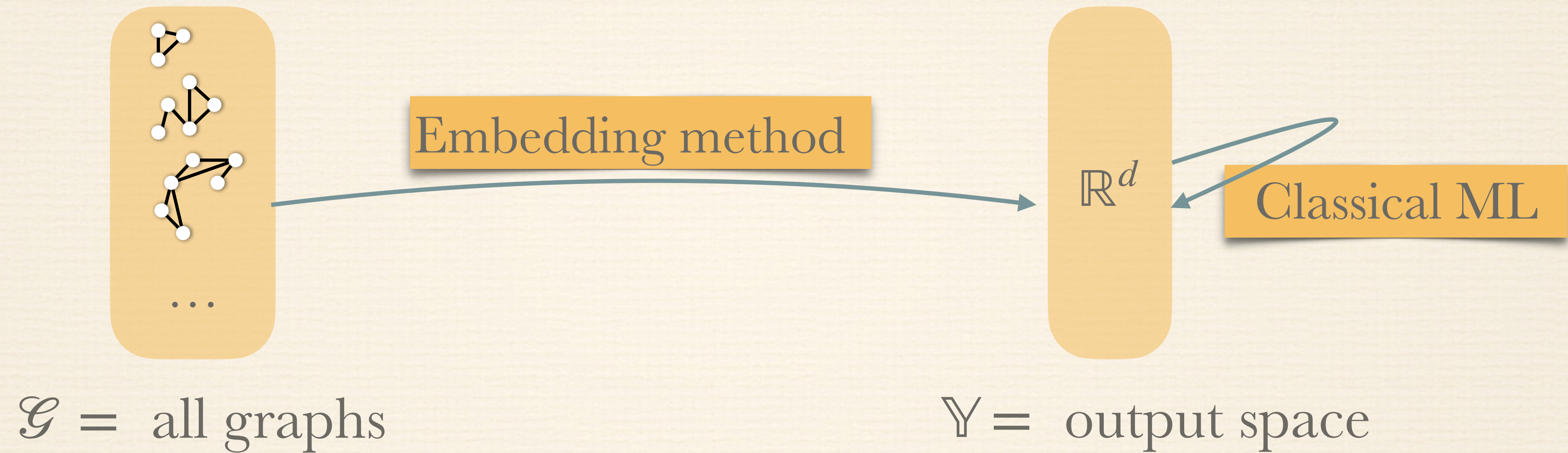
- ❖ Graph $G = (V_G, E_G, L_G)$ can also be represented by **adjacency matrix** A_G and **feature matrix** F_G
- ❖ Let $n = |V_G|$ be the number of vertices. Let $v, w \in [n] := \{1, \dots, n\}$.

adjacency matrix $A_G \in \mathbb{R}^{n \times n} : (v, w) \mapsto \begin{cases} 1 & (v, w) \in E_G \\ 0 & \text{otherwise} \end{cases}$

feature matrix $F_G \in \mathbb{R}^{n \times d} : v \mapsto L_G(v)$

- ❖ To turn graph into matrix, one needs an ordering on the vertices.

Graph learning



Embeddings

\mathcal{G} = all graphs

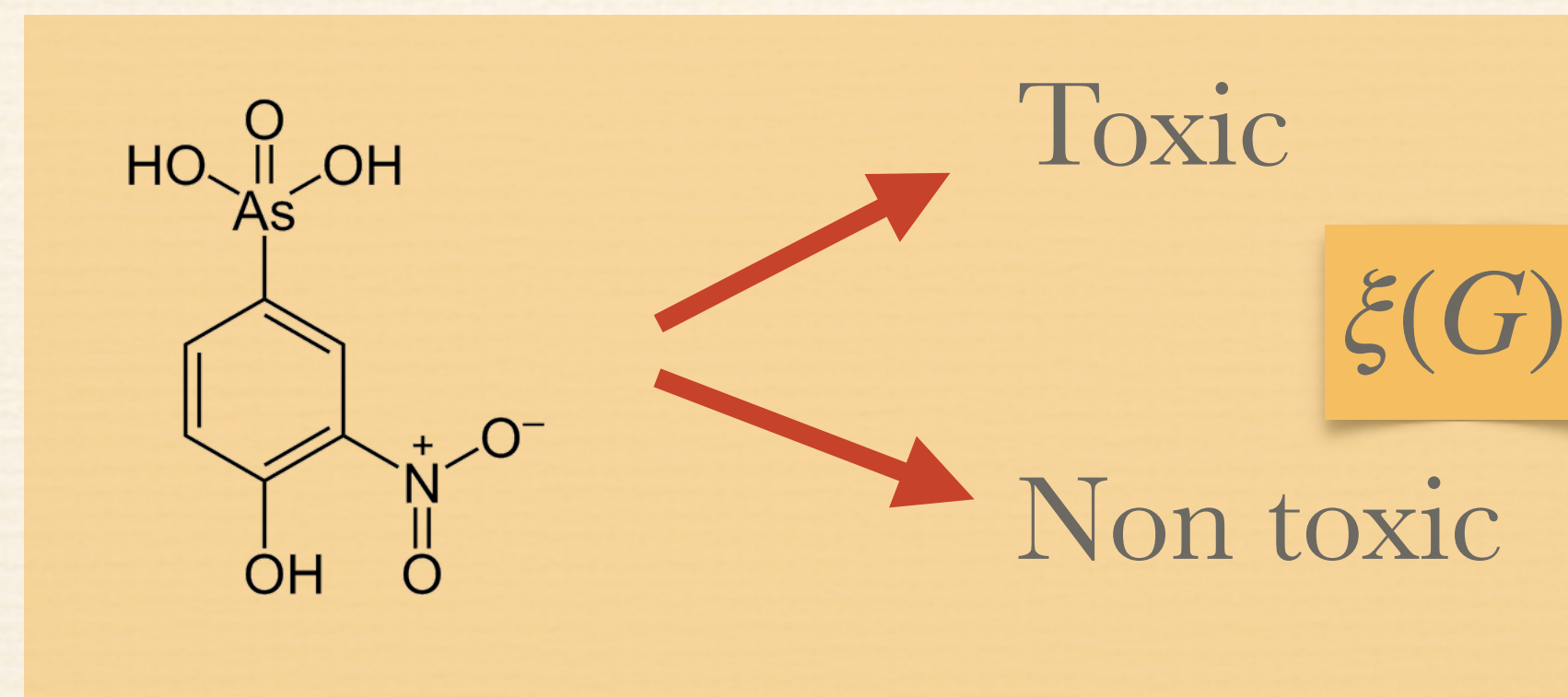
\mathcal{V} = all vertices

\mathbb{Y} = output space

- ❖ Graph embedding: $\xi : \mathcal{G} \rightarrow \mathbb{Y}$
- ❖ Vertex embedding: $\xi : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{Y})$
- ❖ p -Vertex embedding: $\xi : \mathcal{G} \rightarrow (\mathcal{V}^p \rightarrow \mathbb{Y})$

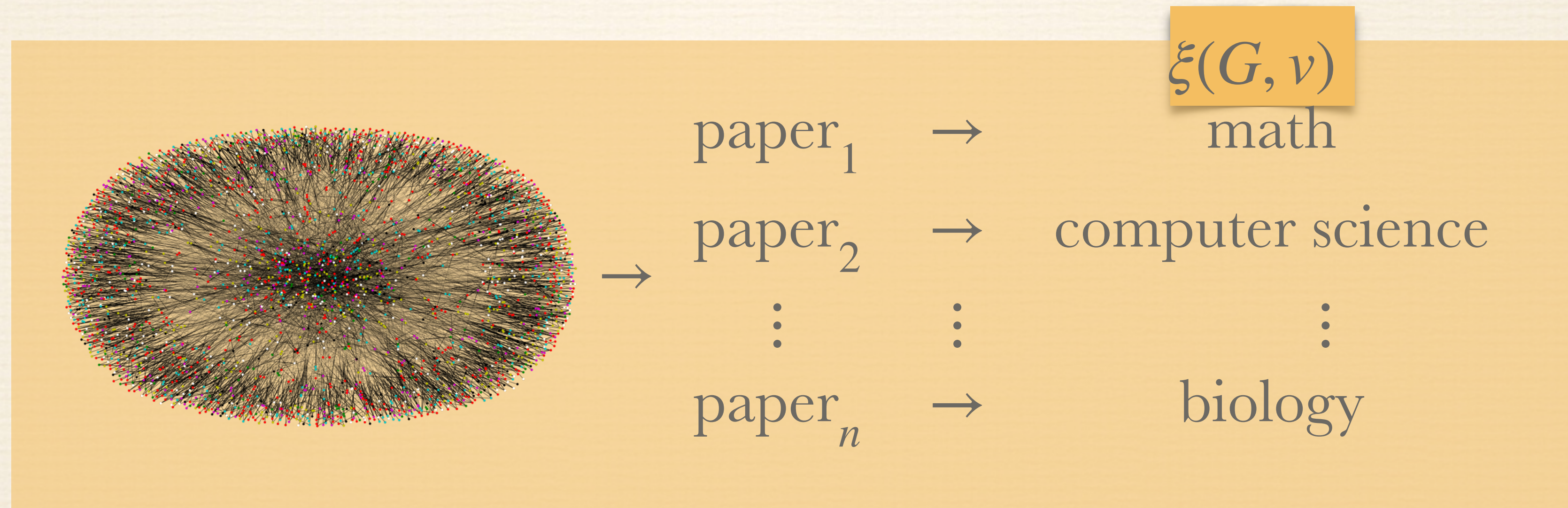
Graph embeddings

- ❖ Graph embedding: $\xi : \mathcal{G} \rightarrow \mathbb{Y}$
- ❖ Graph **classification/regression**



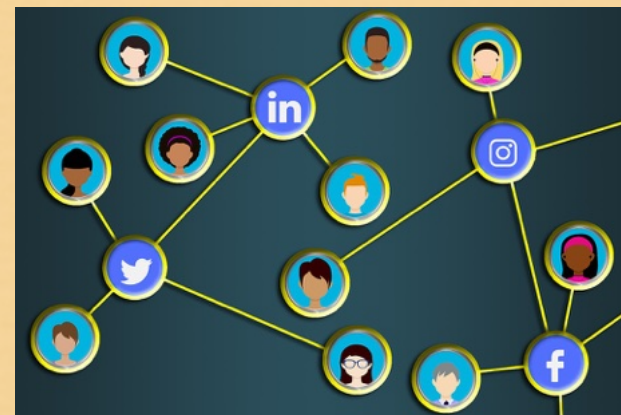
Vertex embeddings

- ❖ Vertex embedding: $\xi : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{Y})$
- ❖ Vertex **classification/regression**. For example, prediction of subject of papers.



p-Vertex embeddings

- ❖ p -Vertex embedding: $\xi : \mathcal{G} \rightarrow (\mathcal{V}^p \rightarrow \mathbb{Y})$
- ❖ For example, 2-vertex embeddings: **link prediction**



→ (Joe, Anna)
(Paolo, Mohammed)

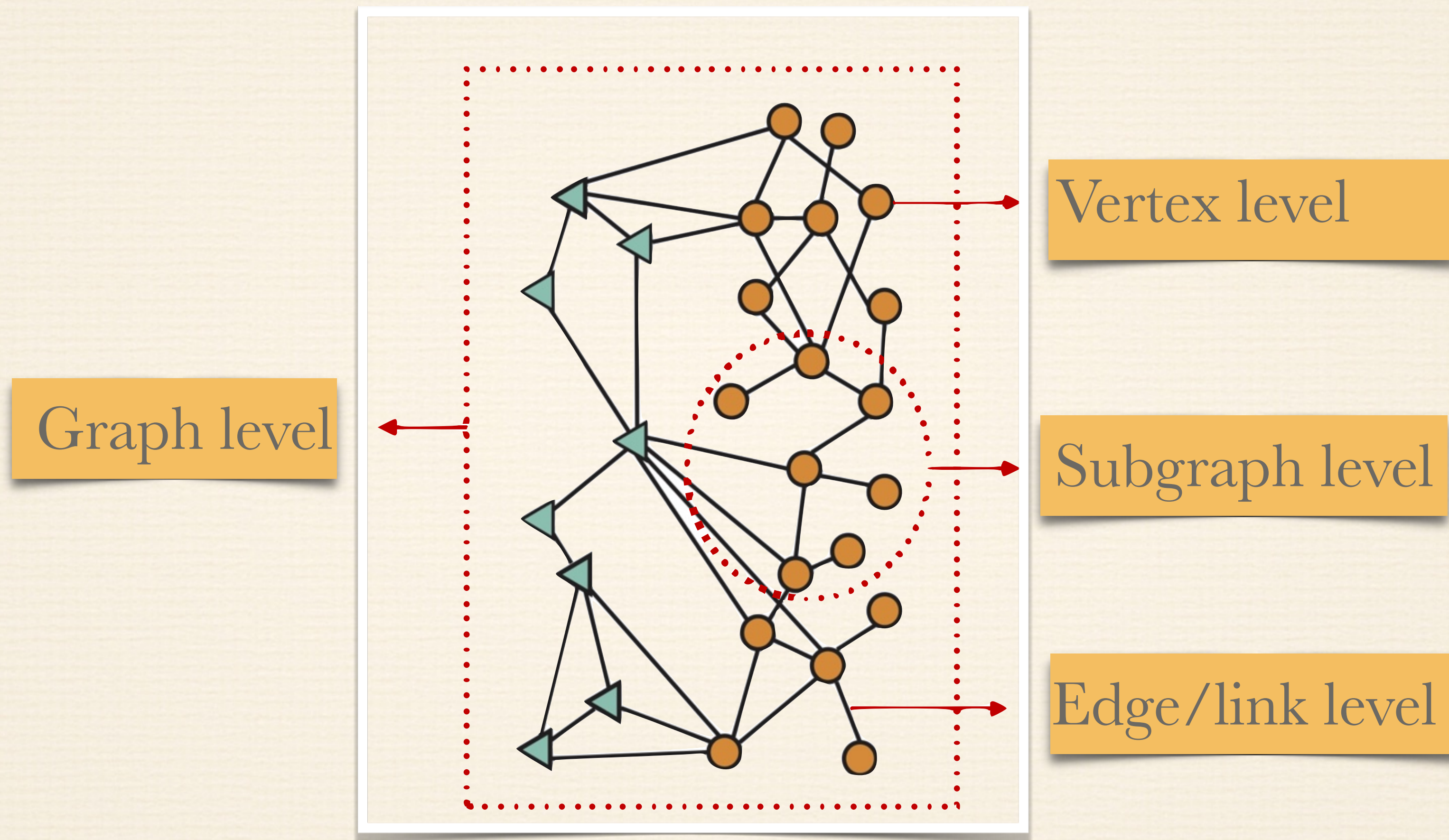
...

$\xi(G, v, w)$

\mapsto link

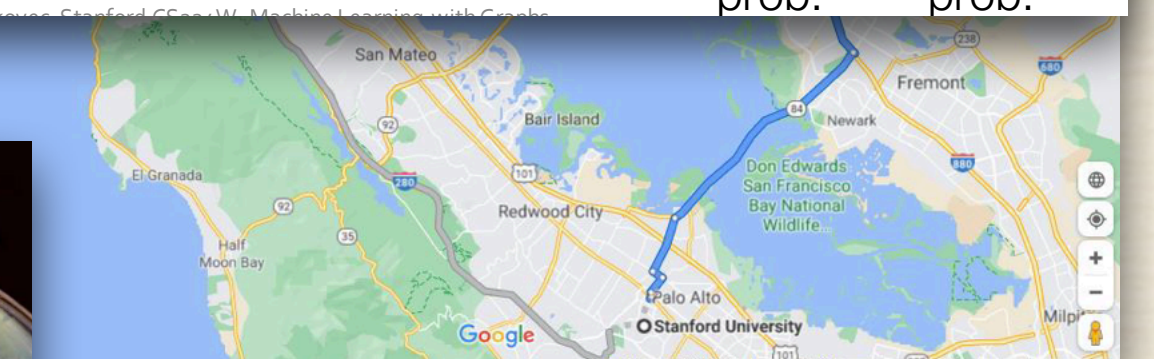
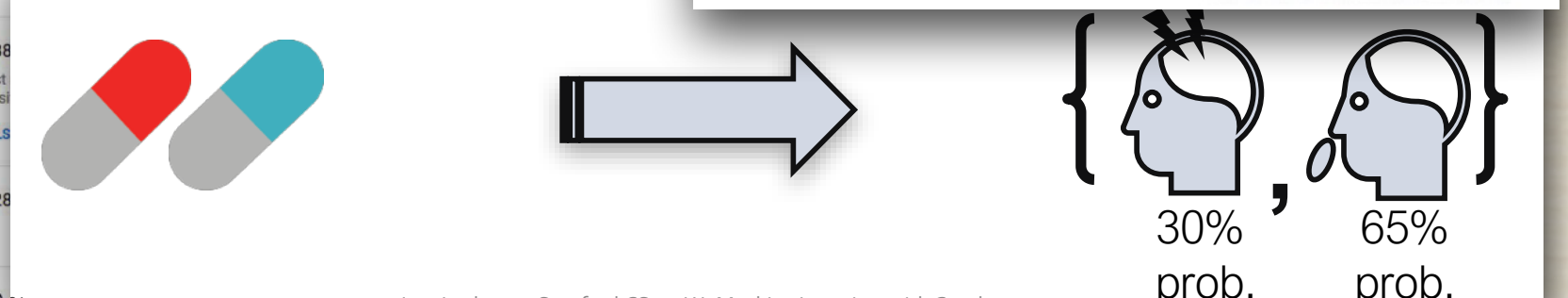
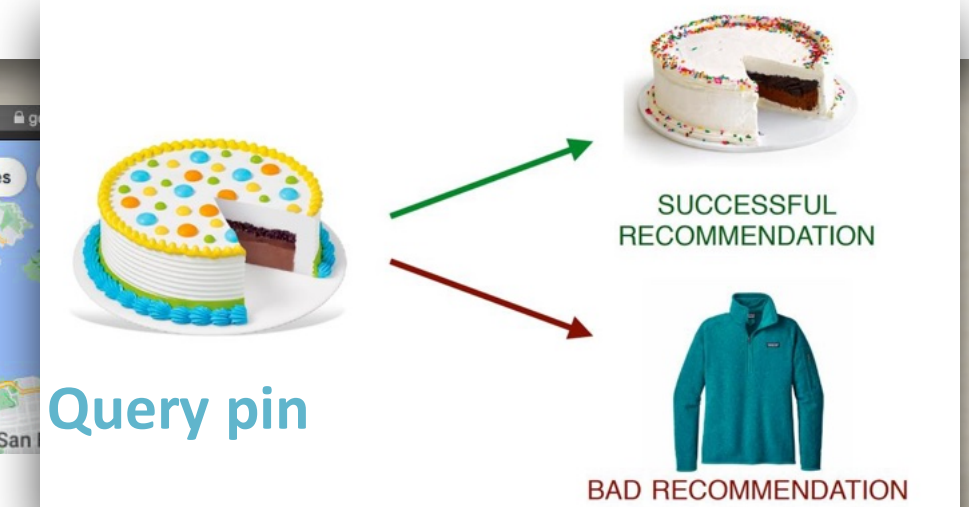
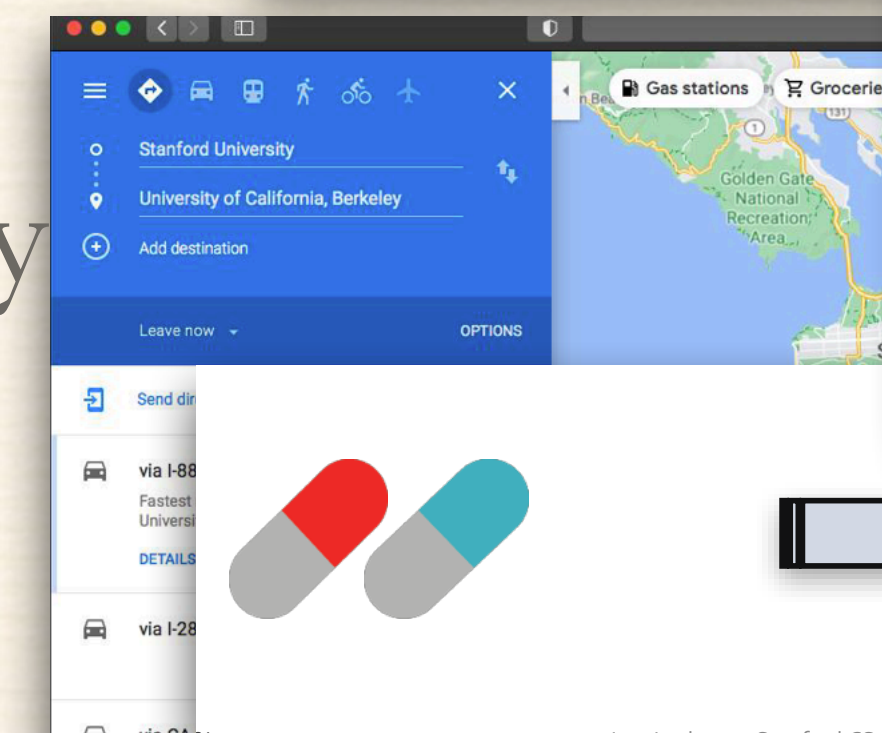
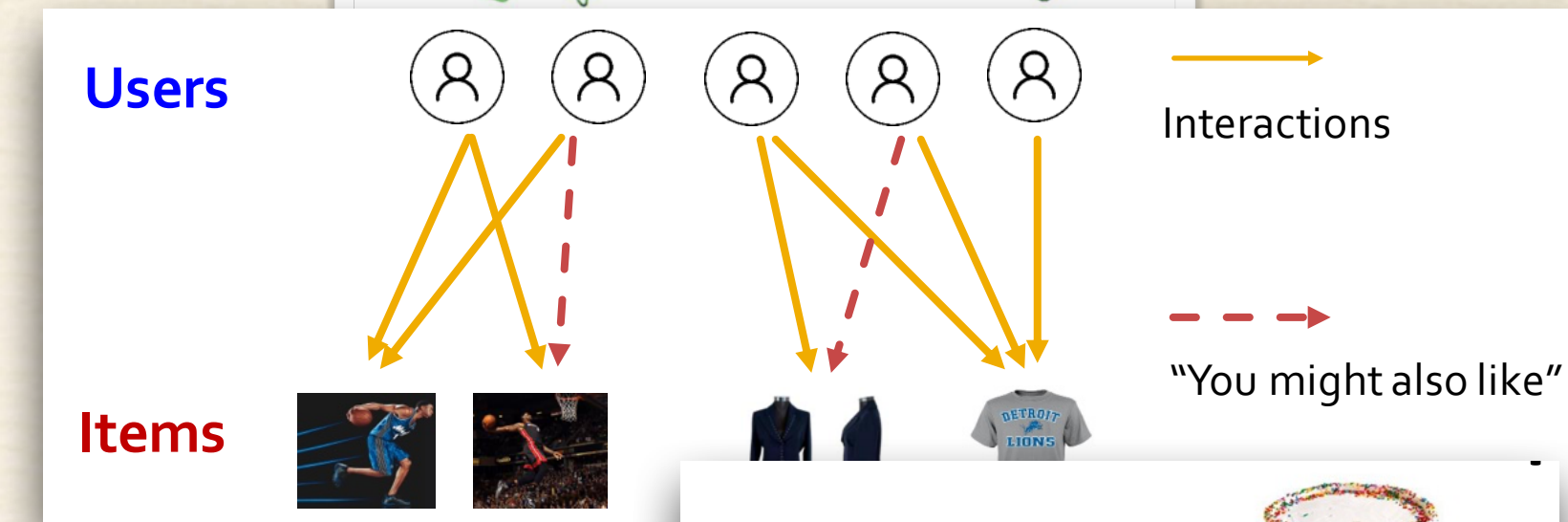
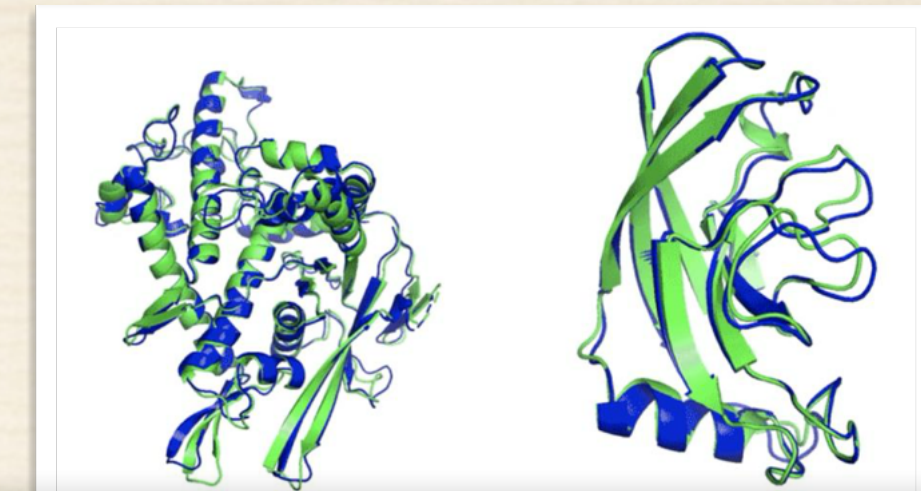
\mapsto no link

Graph learning tasks



Applications

- ❖ **Vertex classification:** categorise online user/items, location amino acids (protein folding, alpha fold)
- ❖ **Link prediction:** graph completion, recommender systems, drug side effect discovery
- ❖ **Graph classification:** molecule property, drug discovery
- ❖ **Subgraph tasks:** traffic prediction



Graph learning

- ❖ We want to learn an unknown embedding $\Xi : \mathcal{G} \rightarrow (\mathcal{V}^p \rightarrow \mathbb{Y})$



What does this mean???

- ❖ The embedding Ξ is partially revealed by means of a **training set**

$$\mathcal{T} := \{(G_1, \mathbf{v}_1, y_1), \dots, (G_\ell, \mathbf{v}_\ell, y_\ell)\} \subseteq \mathcal{G} \times \mathcal{V}^p \times \mathbb{Y}$$

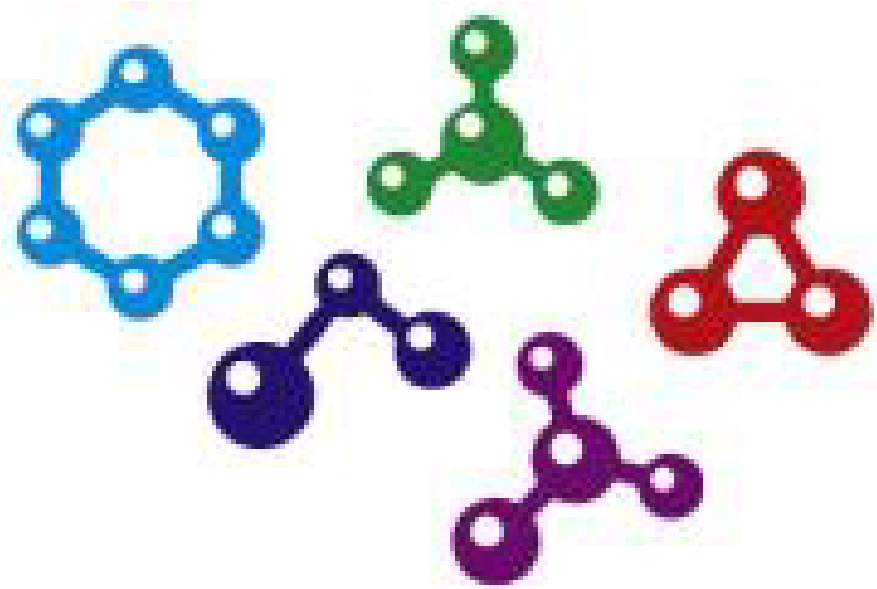


$\Xi(G_1, \mathbf{v}_1)$



$\Xi(G_\ell, \mathbf{v}_\ell)$

Training sets



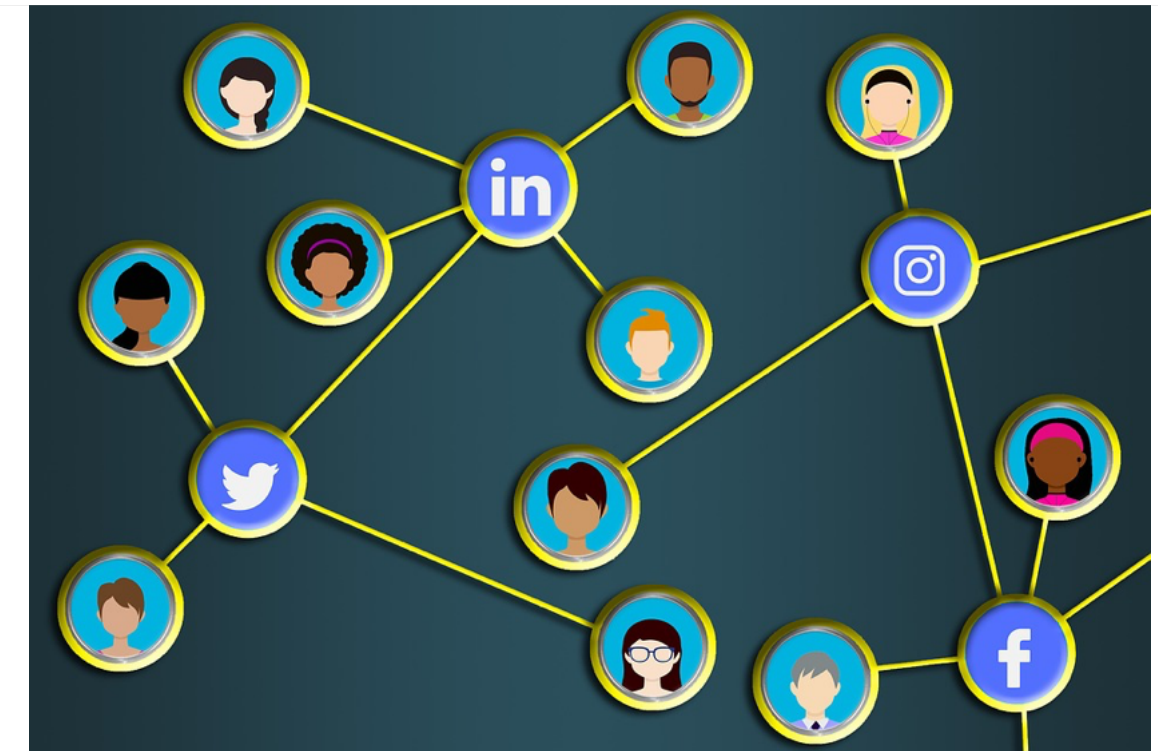
(molecule, yes/no)

Graph classification



(cora, paper, topic)

Vertex classification



(social, p_x, p_y , yes/no)

Link prediction

Graph learning: hypothesis class

- ❖ We want to find the best model consistent with training set \mathcal{T}



What does this mean???

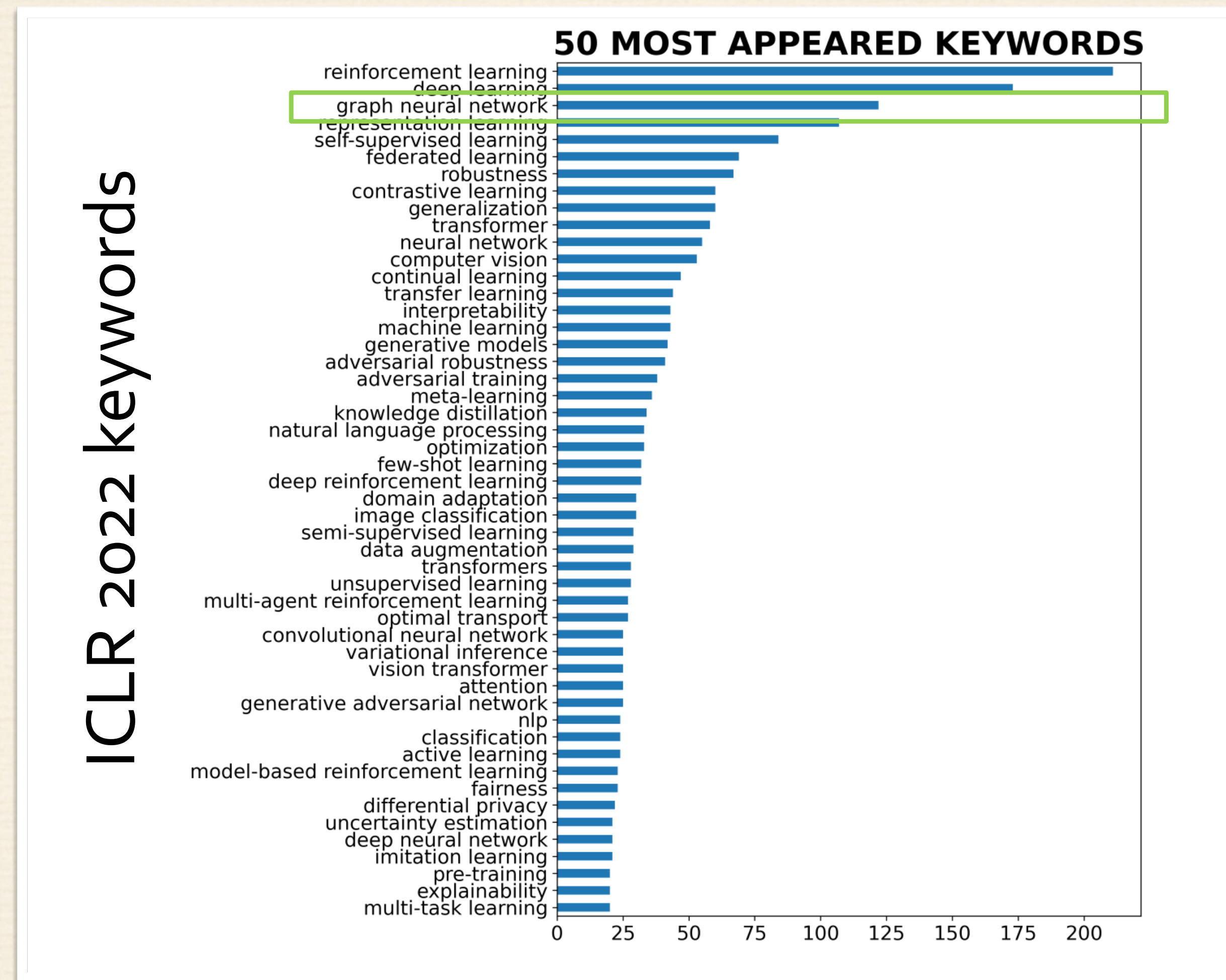
- ❖ Models are selected from an hypothesis class \mathcal{H}
- ❖ In the graph setting \mathcal{H} consists of embeddings

Hypothesis classes

\mathcal{H}

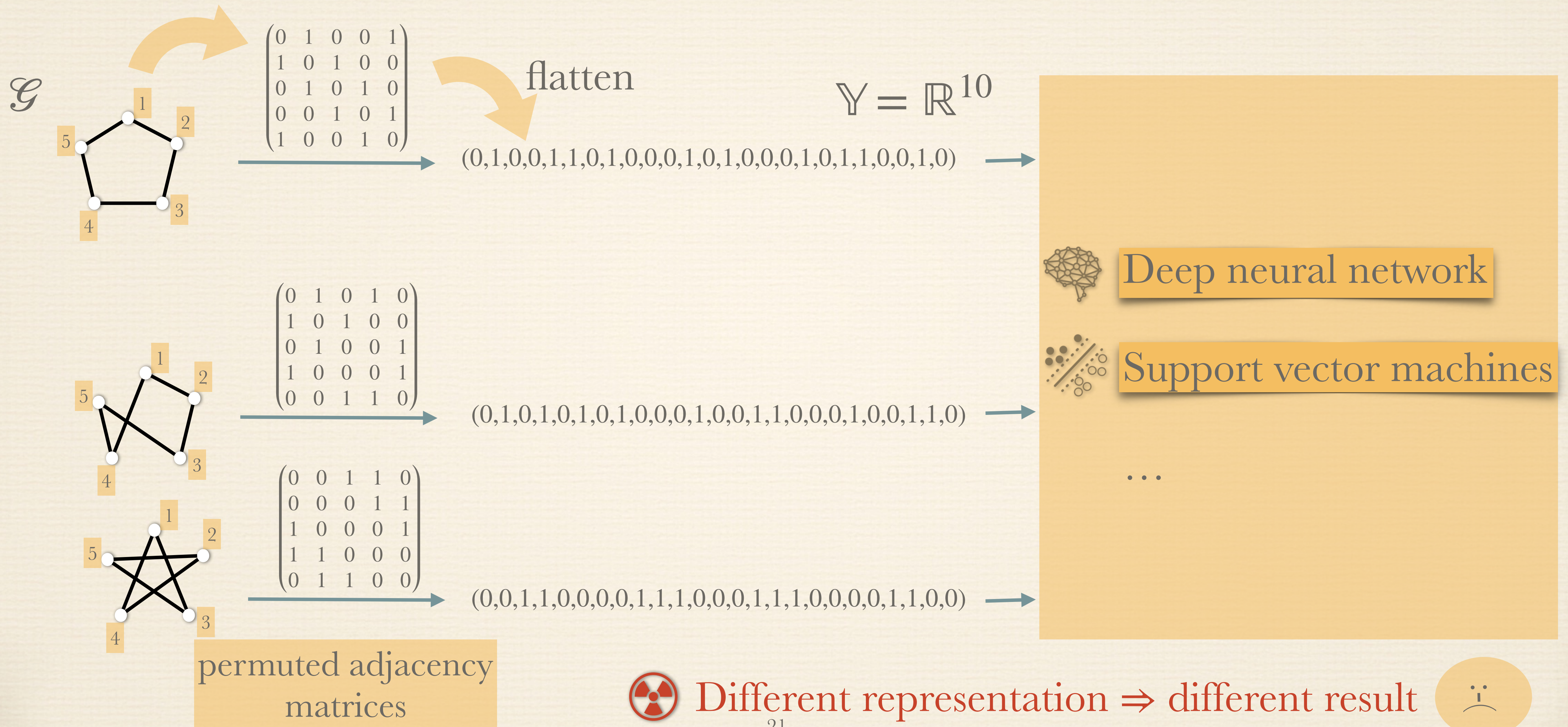
MPNN PPGN
GSN 2-IGN $\times 1000$
Graphormer GATs CayleyNet
CWN $\delta - k - \text{GNNs}$ GIN GCNs
ChebNet Dropout GNN
k-IGNs GraphSage k-SAN k-GNNs
Id-aware GNN

Hypothesis Class Explosion



Count	
Large Language Models	318
Reinforcement Learning	201
Graph Neural Networks	123
Diffusion Models	112
Deep Learning	110
Representation Learning	107

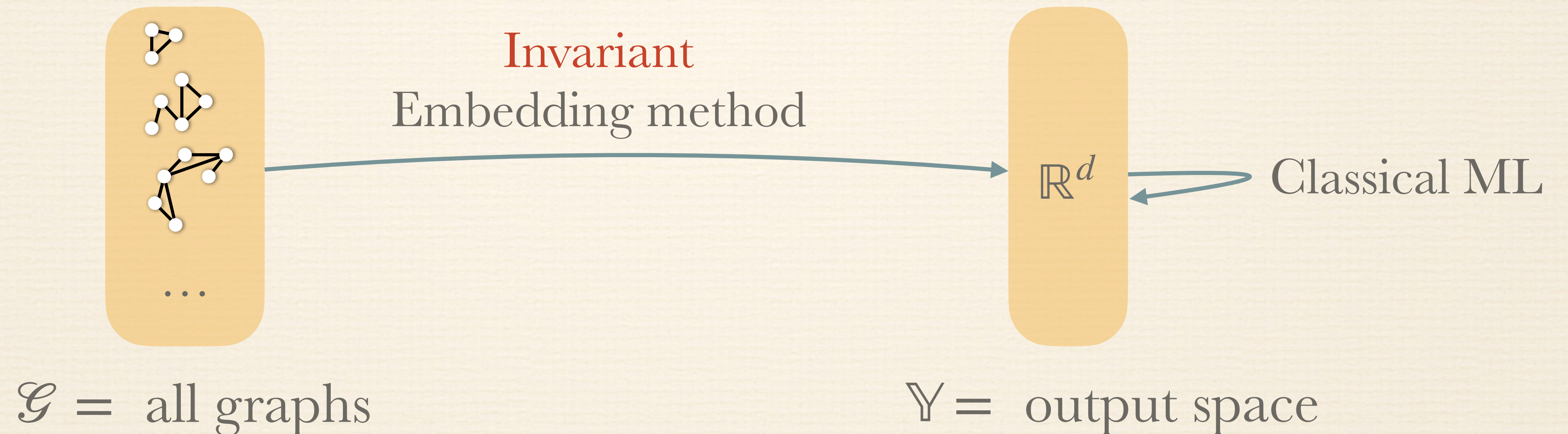
What's new?



Graph learning

Classical embedding methods **depend on representation**

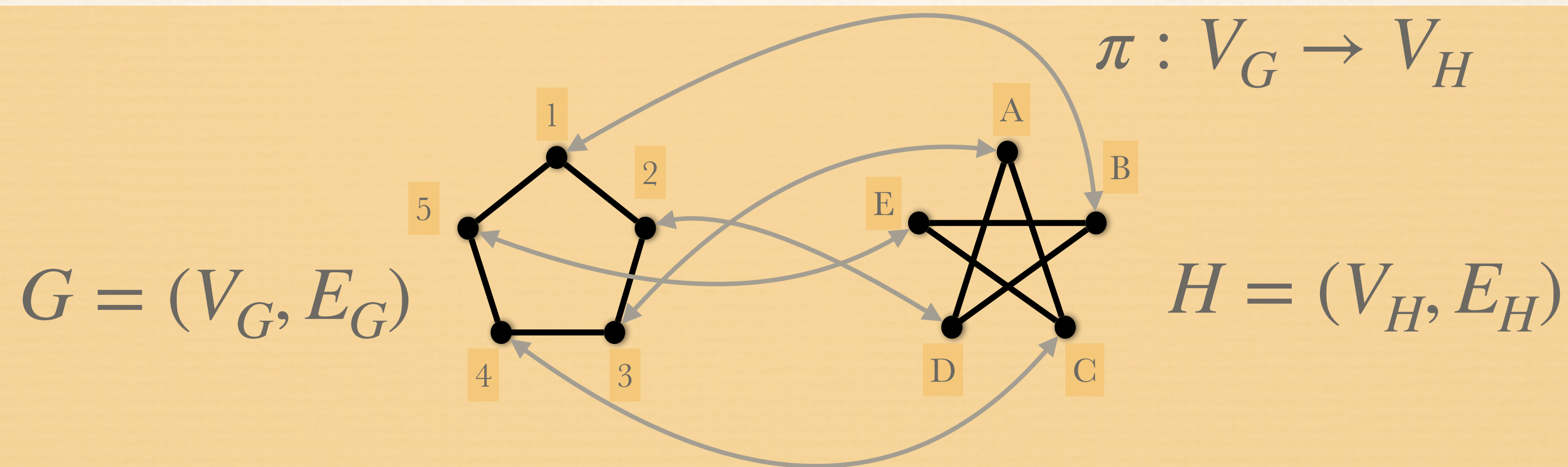
E.g., think of **MultiLayerPerceptron** on **vector representation** of flattened adjacency matrix



A desired property: Invariance

- ❖ Embeddings should be **invariant**, that is, independent of the chosen graph representation.
- ❖ Invariance is defined in terms of **graph isomorphisms**.

$$G \cong H$$

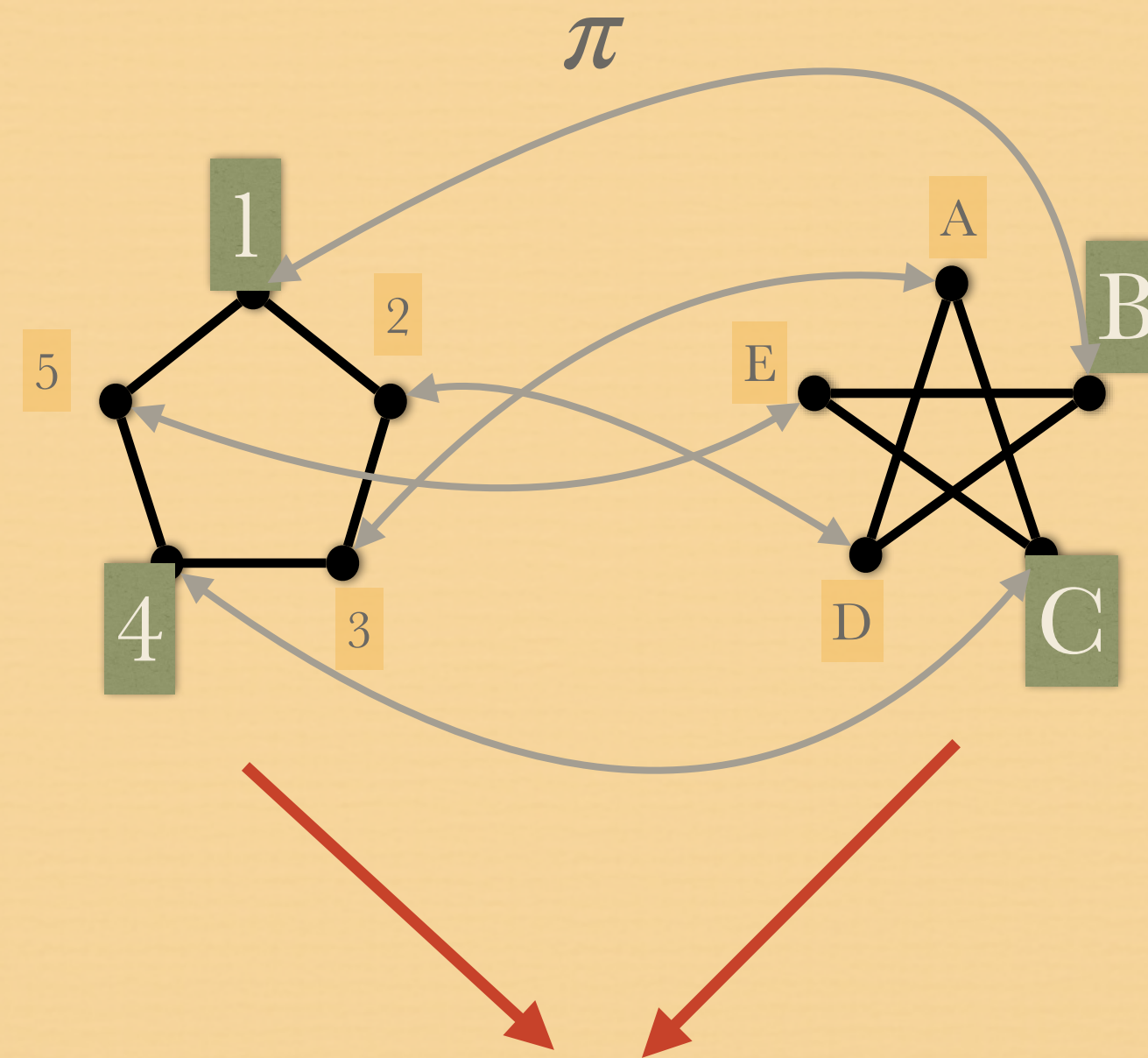


- ❖ The mapping π is a bijective vertex function satisfying $(v, v') \in E_G \iff (\pi(v), \pi(v')) \in E_H$ also $L_G(v) = L_H(\pi(v))$ must hold.

Invariant embeddings

for all π, G and $\mathbf{v} \in V_G^p : \xi(G, \mathbf{v}) = \xi(\pi(G), \pi(\mathbf{v}))$

Isomorphism



$(1,4)$ and (B,C) have same embedding in \mathbb{Y}

We typically assume invariant embedding methods (unless said otherwise)

Graph learning: ERM

Best one! • ξ

\mathcal{H}

- ❖ Given training set \mathcal{T} and hypothesis class \mathcal{H}
- ❖ Empirical risk minimisation:

Find embedding ξ in \mathcal{H} which minimises empirical loss

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \text{loss}(\xi(G_i, \mathbf{v}_i), y_i)$$

Loss function is a mapping from $\mathbb{Y} \times \mathbb{Y} \rightarrow \mathbb{R}$

Loss functions

❖ Choice depends on learning task (regression, classification,...)

❖ L1: $\text{loss}(y_{\text{predicted}}, y_{\text{true}}) := |y_{\text{predicted}} - y_{\text{true}}|$

❖ L2: $\text{loss}(y_{\text{predicted}}, y_{\text{true}}) := (y_{\text{predicted}} - y_{\text{true}})^2$

❖ (Binary) cross entropy:

$$\text{loss}(y_{\text{predicted}}, y_{\text{true}}) := y_{\text{true}} \log(y_{\text{predicted}}) + (1 - y_{\text{true}}) \log(1 - y_{\text{predicted}})$$

Graph learning

- ❖ Graph learning systems solve ERM using **back propagation** and **gradient descent**...

$$\hat{\xi} : \arg \min_{\xi \in \mathcal{H}} \frac{1}{\ell} \sum_{i=1}^{\ell} \text{loss}(\xi(G_i, \mathbf{v}_i), y_i))$$

What one really wants

- ❖ Simply predicting labels for training data is insufficient.
- ❖ We want to predict labels for graphs not in the training data.
- ❖ We will assume the presence of distribution \mathcal{D} over $\mathcal{G} \times \mathbb{Y}$.

We do not know the true distribution \mathcal{D} , as it represents real-world unseen data.

Risk minimisation

- ❖ **Risk minimisation**: Find embedding $\tilde{\xi}$ in \mathcal{H} which minimises expected loss over \mathcal{D} :

$$\tilde{\xi} := \arg \min_{\xi \in \mathcal{H}} \text{Prob}_{(G,y) \sim \mathcal{D}}[\xi(G) \neq y]$$

- ❖ RM focuses on **minimising errors over all the data** according to their distribution.

Generalisation Error

- ❖ We want to find a hypothesis $\xi \in \mathcal{H}$ that does well on the training data (small empirical loss $L_{\mathcal{T}}(\xi)$)
- ❖ But also has small expected loss $L_{\mathcal{D}}(\xi)$
- ❖ We want the **generalisation error** $L_{\mathcal{D}}(\xi) - L_{\mathcal{T}}(\xi)$ to be small.
- ❖ Statistical learning theory **provides bounds** on training data guaranteeing small generalisation error.

Generalisation error

- ❖ We want the generalisation error $L_{\mathcal{D}}(\xi) - L_{\mathcal{T}}(\xi)$ to be small.

Theorem (Vapnik and Chervonenkis 1964)


- ❖ For $\delta > 0$, with probability $1 - \delta$ (in our selection of training data \mathcal{T} of size m) for all $\xi \in \mathcal{H}$:

$$L_{\mathcal{D}}(\xi) - L_{\mathcal{T}}(\xi) \leq \sqrt{\frac{2d \log(\frac{em}{d})}{m}} + \sqrt{\frac{\log(\frac{1}{\delta})}{2m}}$$

- ❖ where d is the **VC dimension** of \mathcal{H}

Graph learning

- ❖ Graph learning systems solve ERM using **back propagation** and **gradient descent**...

$$\hat{\xi} : \arg \min_{\xi \in \mathcal{H}} \frac{1}{\ell} \sum_{i=1}^{\ell} \text{loss}(\xi(G_i, \mathbf{v}_i), y_i)$$


Our focus will be on the **expressive power** of hypothesis classes \mathcal{H}

Expressive power

- ❖ Which embeddings can **be expressed** by embeddings in \mathcal{H} ?
- ❖ Which embeddings can **be approximated** by embeddings in \mathcal{H} ?
- ❖ Which **inputs** can be **separated/distinguished** by embeddings in \mathcal{H} ?
- ❖ What is the relationship between **expressiveness** and **generalisation** of \mathcal{H} ?

Notions of expressivity

❖ Let $\Xi : \mathcal{G} \rightarrow (\mathcal{V}^p \rightarrow \mathbb{Y})$ be a *p*-vertex embedding

\mathcal{H} can *express* Ξ if there exists a $\xi \in \mathcal{H}$ such that for all $G \in \mathcal{C}, \mathbf{v} \in V_G^p :$
$$\xi(G, \mathbf{v}) = \Xi(G, \mathbf{v})$$

❖ Let $\Xi : \mathcal{G} \rightarrow \{0,1\}$ indicator function for *connected* graphs.

❖ Can we find hypothesis in $\xi \in \mathcal{H}$ such that $\xi(G) = \Xi(G)$ for all graphs G

Notions of expressivity

Separation/distinguishing power of \mathcal{H}

$$\rho(\mathcal{H}) := \{(G, \mathbf{v}, H, \mathbf{w}) \mid \forall \xi \in \mathcal{H} : \xi(G, \mathbf{v}) = \xi(H, \mathbf{w})\}$$

- ❖ All pairs of inputs **that cannot be separated** by any embedding in \mathcal{H}
- ❖ Can we find hypothesis in $\xi \in \mathcal{H}$ such that $\xi(G) \neq \xi(H)$ for any connected graph G and disconnected graph H ?

Distinguishing power

- ❖ **Strongest power:** \mathcal{H} powerful enough to detect non-isomorphic graphs
- ❖ **Weakest power:** \mathcal{H} cannot differentiate any two graphs



Distinguishing power

- ❖ Allows for comparing **different classes of** embeddings methods!

$$\rho(\text{methods1}) \subseteq \rho(\text{methods2})$$

Methods1 is more powerful than Methods2
Methods 2 is bounded by Methods 1 in power

$$\rho(\text{methods1}) = \rho(\text{methods2})$$

Both methods are as powerful

- ❖ Allows for comparing embedding methods with **algorithms, logic, ...**

Expressive power in ML community

- ❖ Focus has been on **distinguishing power** of classes \mathcal{H} of embedding methods.
- ❖ Goal is to **characterise $\rho(\mathcal{H})$** in a way to sheds light on what **graph properties** a learning method can detect/use.
- ❖ We see an example shortly for \mathcal{H} = the class of **Message-Passing Neural Networks (MPNNs)**
- ❖ Recent work addresses **uniform expressiveness**.

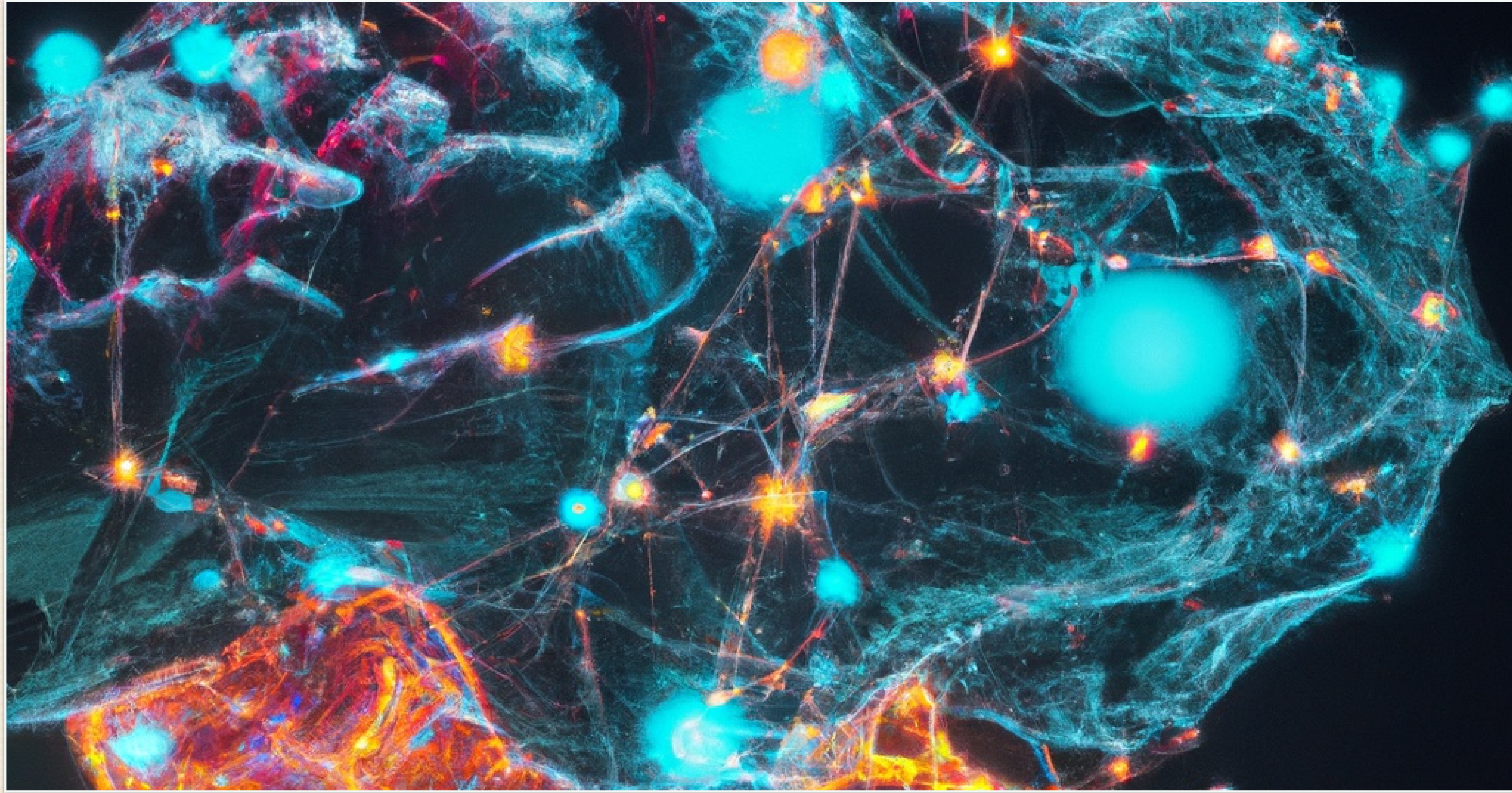
Expressive power in ML community

- ❖ Search for **increase** in expressive power has led to **surge of new methods** of graph learning.
- ❖ Despite theoretical underpinning... still a bit of **alchemy** to find the right method...





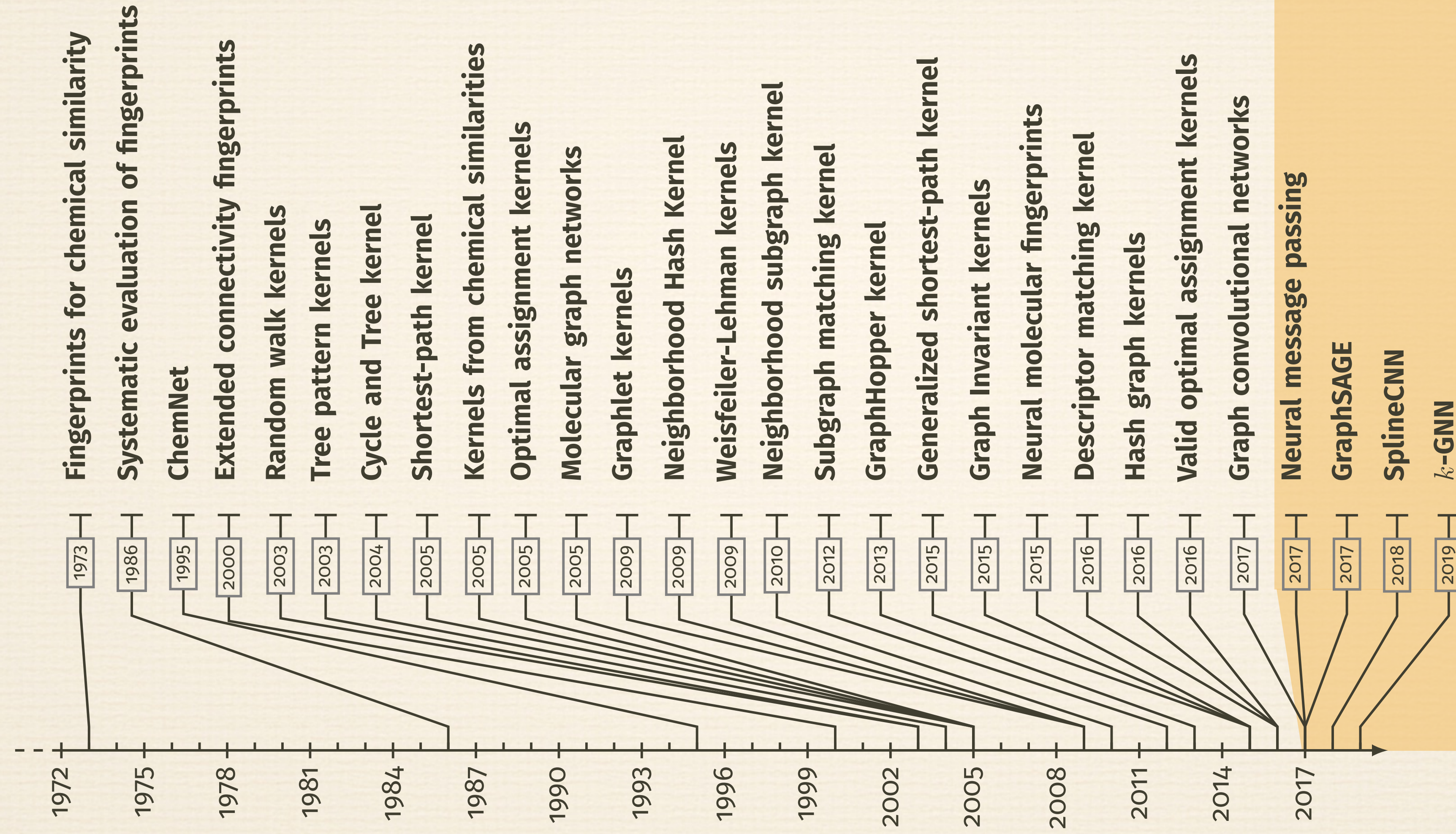
Questions?



Message Passing Neural Networks

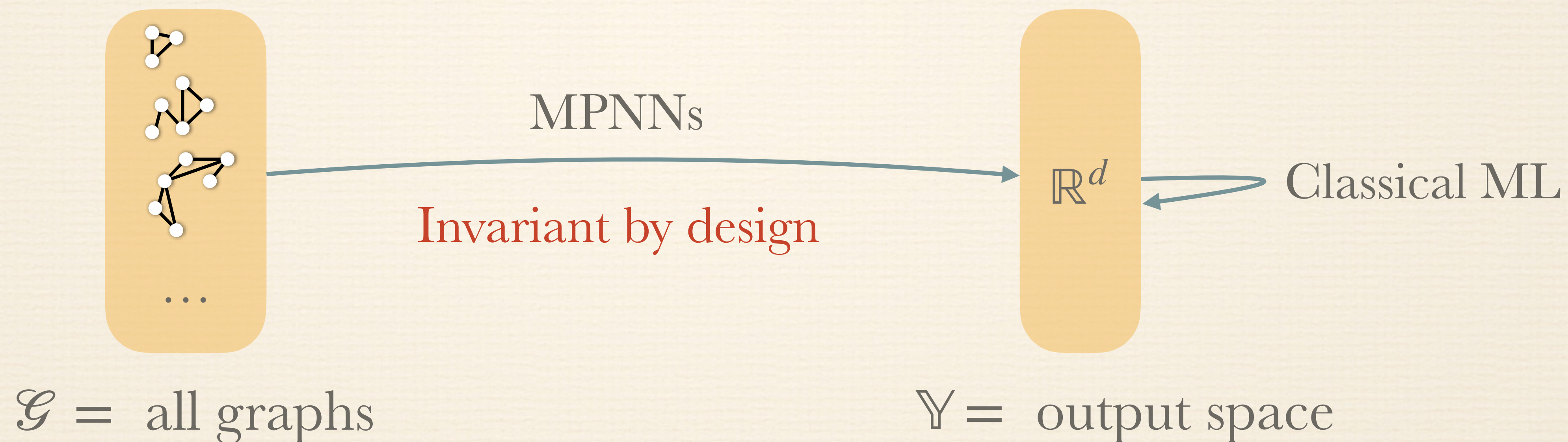
(Still) the most popular type of Graph Learning Architecture

A little graph embedding history

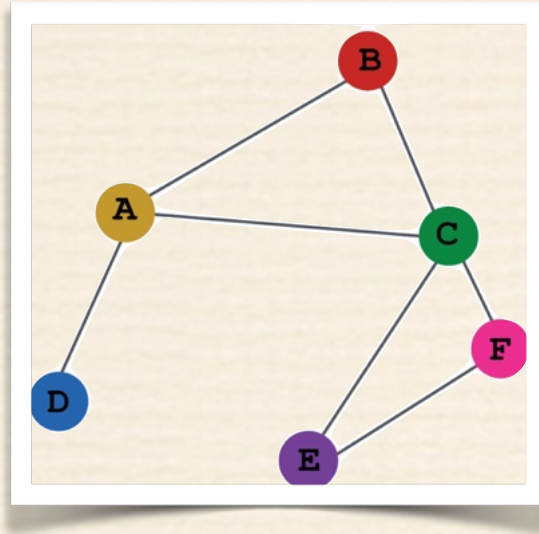


Message passing neural networks

A class of **invariant vertex** and **graph** embedding methods

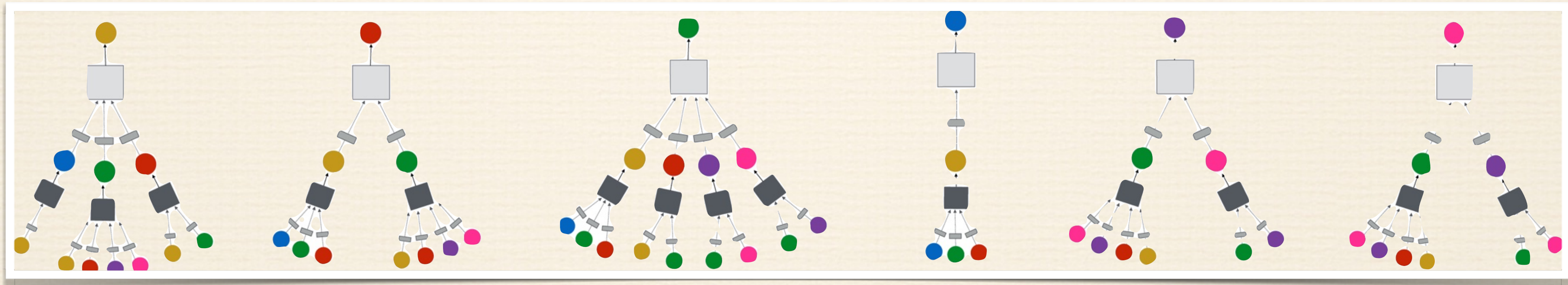


Idea behind MPNNs: Neighbour aggregation



Every vertex defines a computation graph

 Neural networks



MPNNs: Vertex embedding

$$\xi(G, v) := \xi^{(L)} \circ \xi^{(L-1)} \circ \dots \circ \xi^{(0)}(G, v)$$

Message Passing Layers $\xi^{(i)}(G, v) \in \mathbb{R}^d$

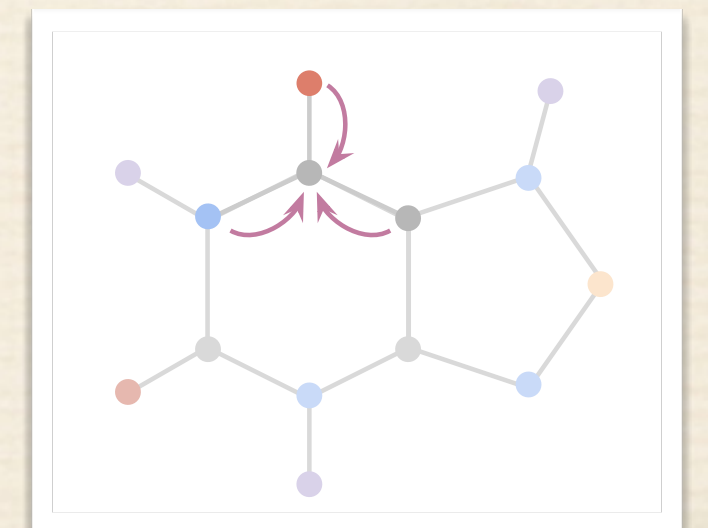
$\xi^{(0)}(G, v) :=$ Hot-one encoding of label of vertex $v \in \mathbb{R}^d$

$$\xi^{(t)}(G, v) := \text{Upd}^{(t)}\left(\xi^{(t-1)}(G, v), \text{Agg}^{(t)}\left(\{\{\xi^{(t-1)}(G, v), \xi^{(t-1)}(G, u) \mid u \in N_G(v)\}\}\right)\right) \in \mathbb{R}^d$$

neighbourhoods

Message Passing between v and its
neighbours $u \in N_G(v)$

Update and aggregate function contain
learnable parameters (NNs)



MPNNs: Graph embedding

$$\rho(G) := \rho \circ \xi^{(L)} \circ \xi^{(L-1)} \circ \dots \circ \xi^{(0)}(G, v)$$

↑
Readout

$$\rho(G) := \text{Readout}\left(\left\{\left\{\xi^{(L)}(G, v) \mid v \in V_G\right\}\right\}\right) \in \mathbb{R}^d$$

↑
Has learnable parameters

↑
Aggregation over all vertices

Typical choices for update, aggregate and readout: **Multilayer Perceptrons**

MPNN example: GNN 101

Neural Network Activation Functions: a small subset!

ReLU $\max(0, x)$	GELU $\frac{x}{2} \left(1 + \tanh \left(\sqrt{\frac{2}{\pi}} (x + ax^3) \right) \right)$	PReLU $\max(0, x)$
ELU $\begin{cases} x & \text{if } x > 0 \\ \alpha(\exp x - 1) & \text{if } x < 0 \end{cases}$	Swish $\frac{x}{1 + \exp -x}$	SELU $\alpha(\max(0, x) + \min(0, \beta(\exp x - 1)))$
SoftPlus $\frac{1}{\beta} \log(1 + \exp(\beta x))$	Mish $x \tanh \left(\frac{1}{2} \log(1 + \exp(\beta x)) \right)$	RReLU $\begin{cases} x & \text{if } x \geq 0 \\ ax & \text{if } x < 0 \text{ with } a \sim \mathcal{U}(l, u) \end{cases}$
HardSwish $\begin{cases} 0 & \text{if } x \leq -3 \\ x & \text{if } x \geq 3 \\ x(x+3)/6 & \text{otherwise} \end{cases}$	Sigmoid $\frac{1}{1 + \exp(-x)}$	SoftSign $\frac{x}{1 + x }$
Tanh $\tanh(x)$	Hard tanh $\begin{cases} a & \text{if } x \geq a \\ b & \text{if } x \leq b \\ x & \text{otherwise} \end{cases}$	Hard Sigmoid $\begin{cases} 0 & \text{if } x \leq -3 \\ 1 & \text{if } x \geq 3 \\ x/6 + 1/2 & \text{otherwise} \end{cases}$
Tanh Shrink $x - \tanh(x)$	Soft Shrink $\begin{cases} x - \lambda & \text{if } x > \lambda \\ x + \lambda & \text{if } x < -\lambda \\ 0 & \text{otherwise} \end{cases}$	Hard Shrink $\begin{cases} x & \text{if } x > \lambda \\ x & \text{if } x < -\lambda \\ 0 & \text{otherwise} \end{cases}$

- ❖ **Non-linear activation** function σ (ReLU, sign, sigmoid, ...)
- ❖ $\mathbf{F}_{v\bullet}^{(t)} \in \mathbb{R}^d$ denotes embedding of vertex v
- ❖ **Weight** matrices $\mathbf{W}_1^{(t)} \in \mathbb{R}^{d \times d}$ and $\mathbf{W}_2^{(t)} \in \mathbb{R}^{d \times d}$ and **bias** vector $\mathbf{b} \in \mathbb{R}^{1 \times d}$

Matrix form

$$\mathbf{F}_{v\bullet}^{(0)} := L_G(v) \quad \leftarrow \text{Embedding vertex labels}$$

$$\mathbf{F}_{v\bullet}^{(t)} := \sigma \left(\mathbf{F}_{v\bullet}^{(t-1)} \mathbf{W}_1^{(t)} + \sum_{u \in N_G(v)} \mathbf{F}_{u\bullet}^{(t-1)} \mathbf{W}_2^{(t)} + \mathbf{b}^{(t)} \right)$$

$$\mathbf{F}^{(t)} := \sigma \left(\mathbf{F}^{(t-1)} \mathbf{W}_1^{(t)} + \mathbf{A} \mathbf{F}^{(t-1)} \mathbf{W}_2^{(t)} + \mathbf{B}^{(t)} \right) \quad \begin{array}{l} \leftarrow \text{Aggregation over} \\ \text{neighbours} \end{array}$$

adjacency matrix

GNN 101: Graph embedding

❖ **Weight** matrix $\mathbf{W} \in \mathbb{R}^{d \times d}$ and **bias vector** $\mathbf{b} \in \mathbb{R}^{1 \times d}$

$$\mathbf{F}^{(t)} := \sigma \left(\sum_{v \in V_G} \mathbf{F}^{(L)} \mathbf{W} + \mathbf{b} \right)$$

Aggregation over all vertices

ERM: Find best parameters $\mathbf{W}_1^{(1)}, \dots, \mathbf{W}_1^{(L)}, \mathbf{W}_2^{(1)}, \dots, \mathbf{W}_2^{(L)}, \mathbf{W}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}, \mathbf{b}$

Two more examples of MPNNs

❖ Graph Isomorphism Networks (**GIN**)

$$\mathbf{F}_{v\bullet}^{(t)} := \text{MLP}^{(t)} \left((1 + \epsilon^{(t)})\mathbf{F}_{v\bullet}^{(t-1)} + \sum_{u \in N_G(v)} \mathbf{F}_{u\bullet}^{(t-1)} \right)$$

❖ Graph Convolution Network (**GCN**)

$$\mathbf{F}_{v\bullet}^{(t)} := \text{MLP}^{(t)} \left(\frac{1}{\sqrt{|N_G(v)| + 1}} \sum_{u \in N_G(v) \cup \{v\}} \frac{1}{\sqrt{|N_G(u)| + 1}} \mathbf{F}_{u\bullet}^{(t-1)} \right)$$

MPNNs: Expressive power

What is $\rho(\text{MPNNs})$?



Recall: All pairs of graphs (G, H) such that all MPNNs return same graph embedding on both graphs.

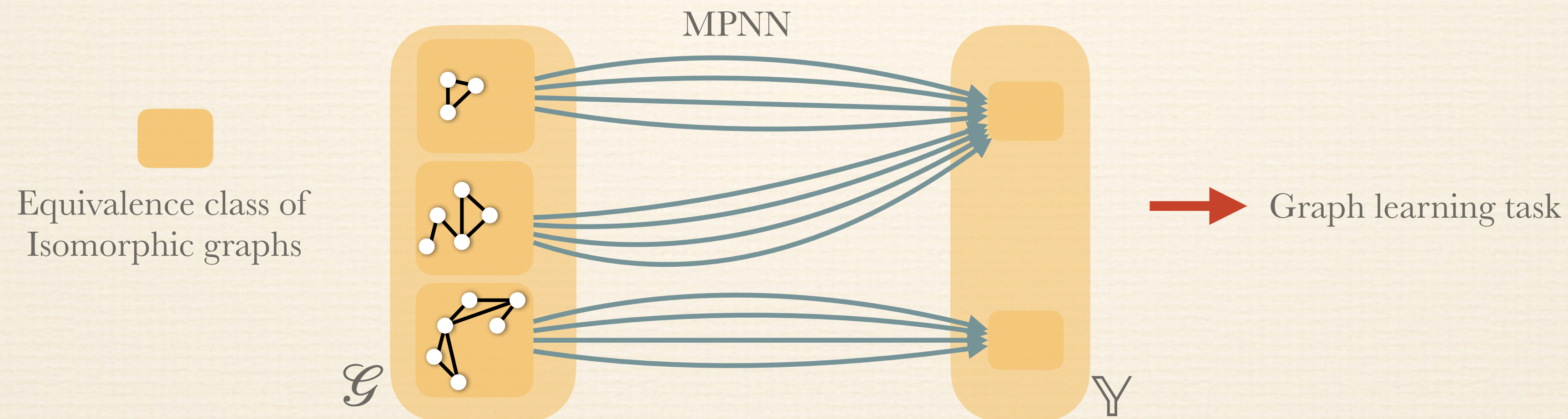


Understanding $\rho(\text{MPNNs})$ translates in understanding power of GNN 101, GCNs, GINs,

A short detour to **graph isomorphism testing**

MPNNs and isomorphic graphs

- ❖ Because of **invariance**: MPNNs embed **isomorphic** graphs in the same way. That is, if $G \cong H \Rightarrow (G, H) \in \rho(\text{MPNN})$
- ❖ Can MPNNs embed **non-isomorphic** graphs differently?



The graph isomorphism problem

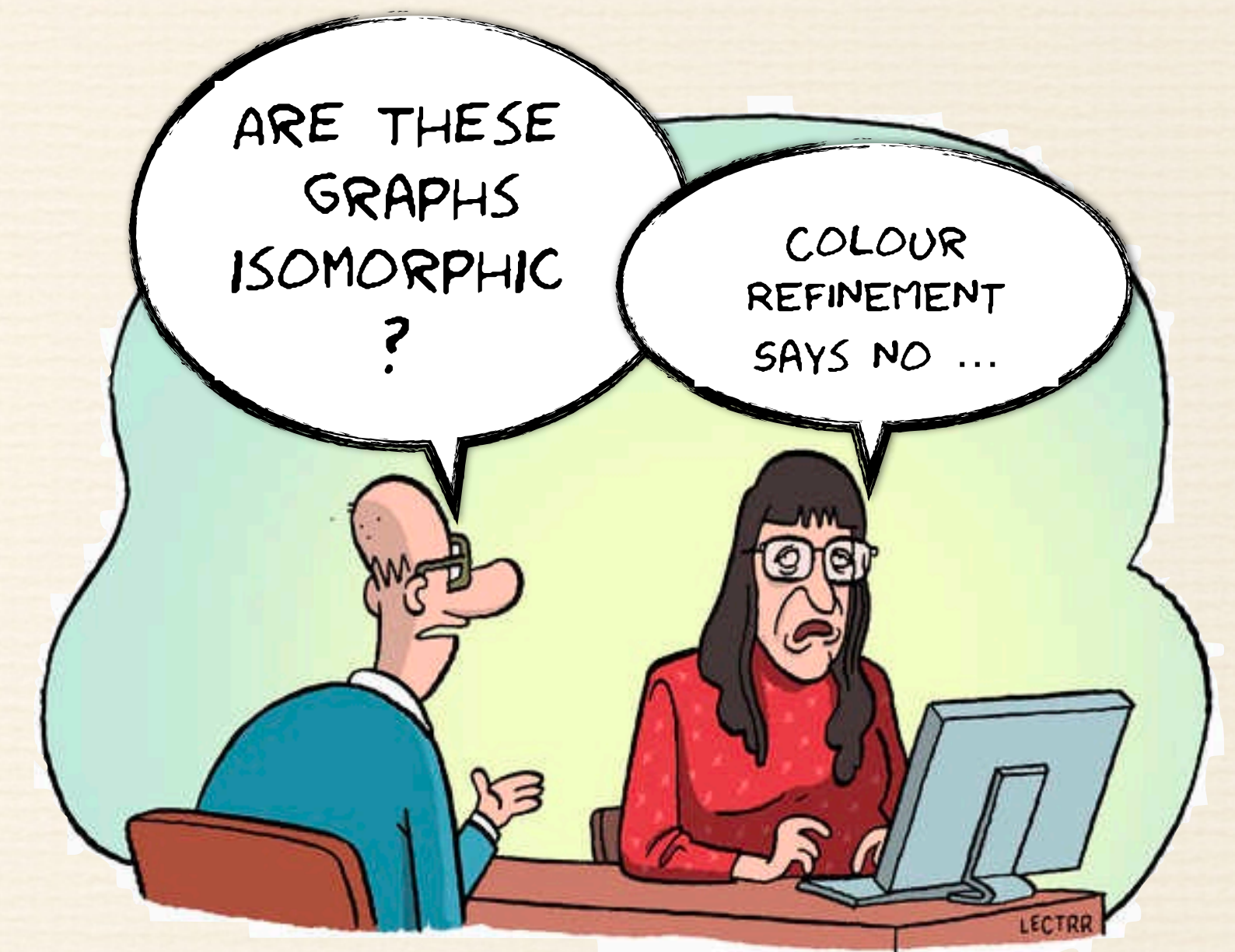
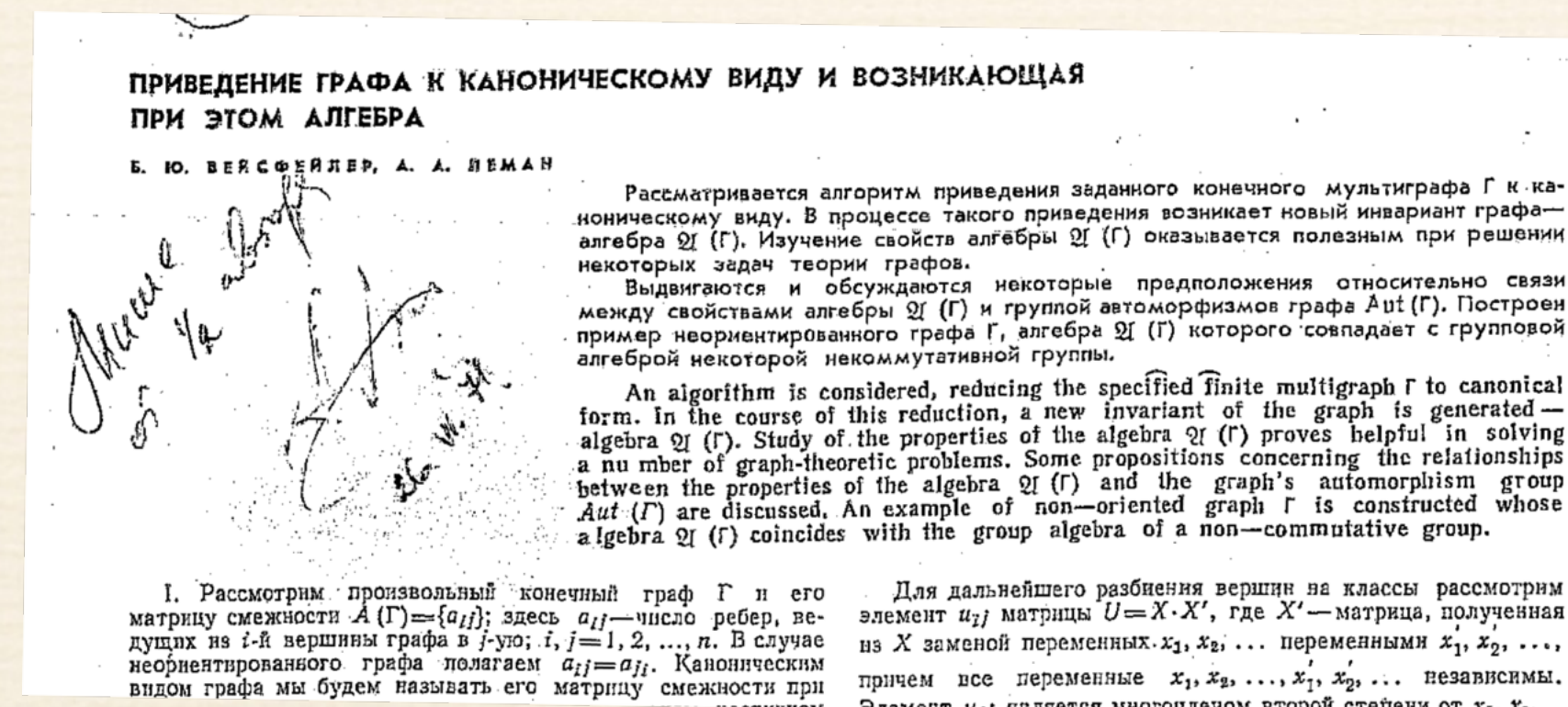
Given two graph $G = (V_G, E_G, L_G)$ and $H = (V_H, E_H, L_H)$: are they isomorphic? Or is $G \cong H$?

- ❖ Does there exist a **graph isomorphism** $\pi : V_G \rightarrow V_H$?
- ❖ **Theory**: computational complexity is open.
- ❖ Quasi-polynomial algorithm $n^{\log(n)^{O(1)}}$ by László Babai (2016).
- ❖ **Practice**: very fast tests.

One-sided test: Colour refinement

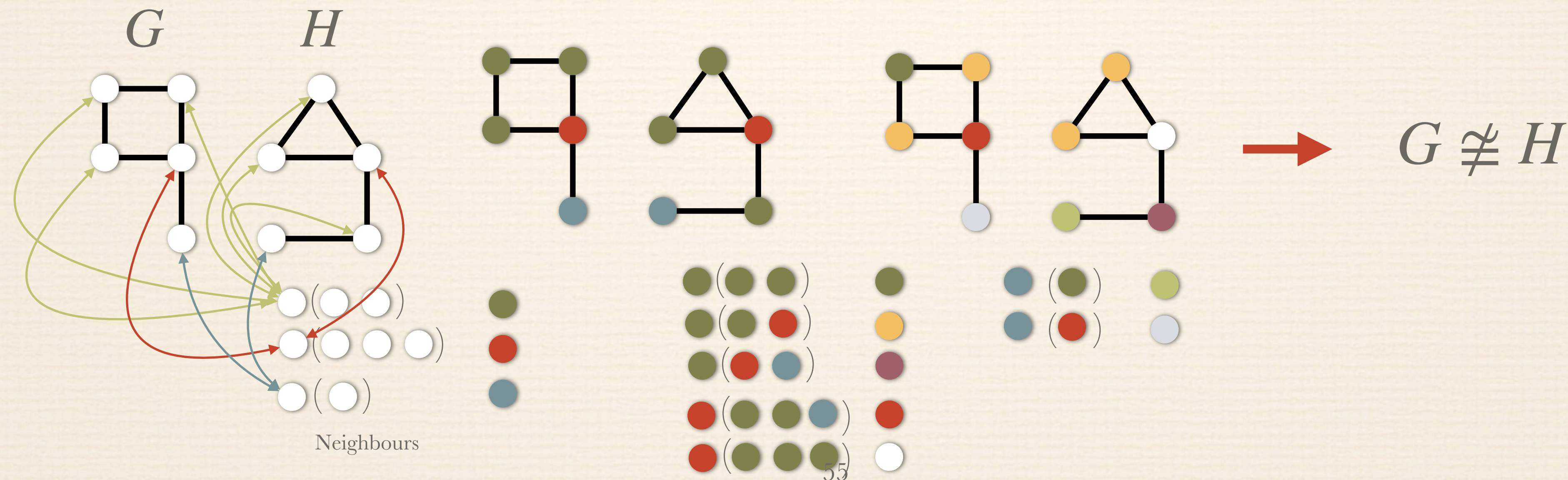
Apply **heuristic** on G and H : If Heuristic say “no” then $G \not\cong H$, otherwise we do not know.

- ❖ Common heuristic is **colour refinement**
- ❖ In a paper by Boris Weisfeiler and Andrei Leman (1968)



Colour refinement

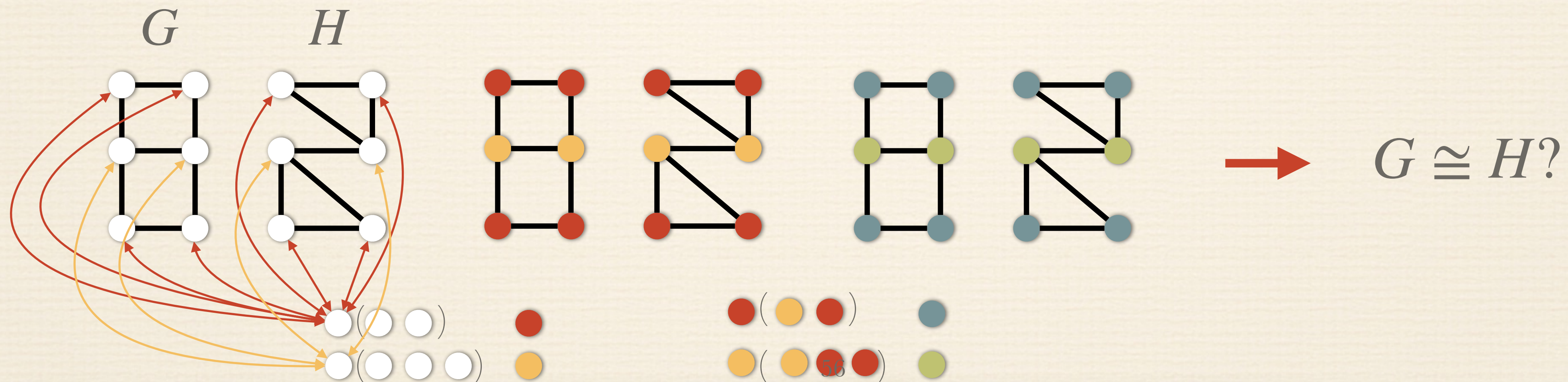
- ❖ Initial: All vertices have their **original colour** (label)
- ❖ Iteration: Separation of identically coloured vertices based on **colour histograms of neighbours**.
- ❖ Two graphs are non-isomorphic if they have **different colour histograms**.



Colour refinement

- ❖ Initial: All vertices have their **original colour** (label)
- ❖ Iteration: Separation of identically coloured vertices based on **colour histograms of neighbours**.
- ❖ Two graphs are non-isomorphic if they have **different colour histograms**.

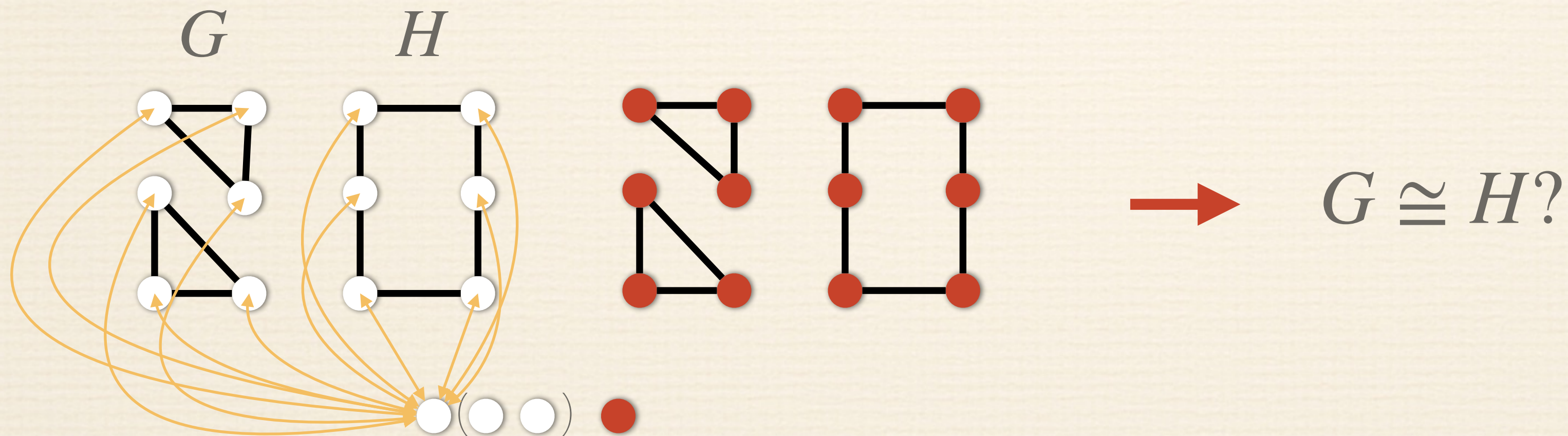
Stops when colour partition does not change (max n iterations)



Colour refinement

- ❖ Initial: All vertices have their **original colour** (label)
- ❖ Iteration: Separation of identically coloured vertices based on **colour histograms of neighbours**.
- ❖ Two graphs are non-isomorphic if they have **different colour histograms**.

Stops when colour partition does not change (max n iterations)



Colour refinement

- ❖ Extensively studied in the theoretical computer science community
- ❖ Many different characterisations of when two graphs have the same colour histograms (equivalent for colour refinement).
- ❖ Successful on random graphs with high probability
- ❖ Weak expressive (distinguishing) power

L. Babai and L. Kucera. *Canonical labelling of graphs in linear average time* (1979)

Cai et al.: *An optimal lower bound on the number of variables for graph identifications*. (1992)

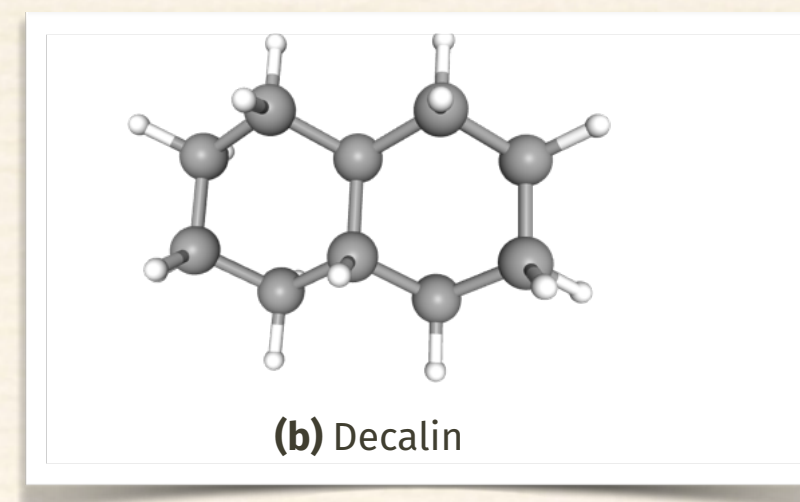
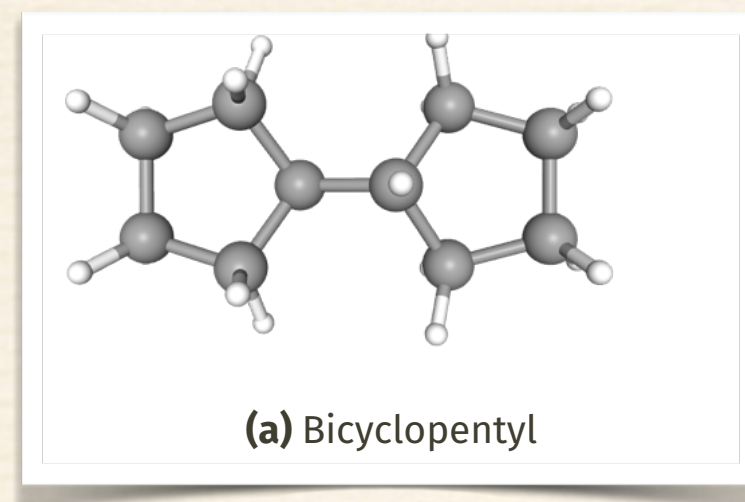
Arvind et al.: *On the power of color refinement* (2015)

M. Grohe: *Descriptive Complexity, Canonisation, and Definable Graph Structure Theory* (2017)

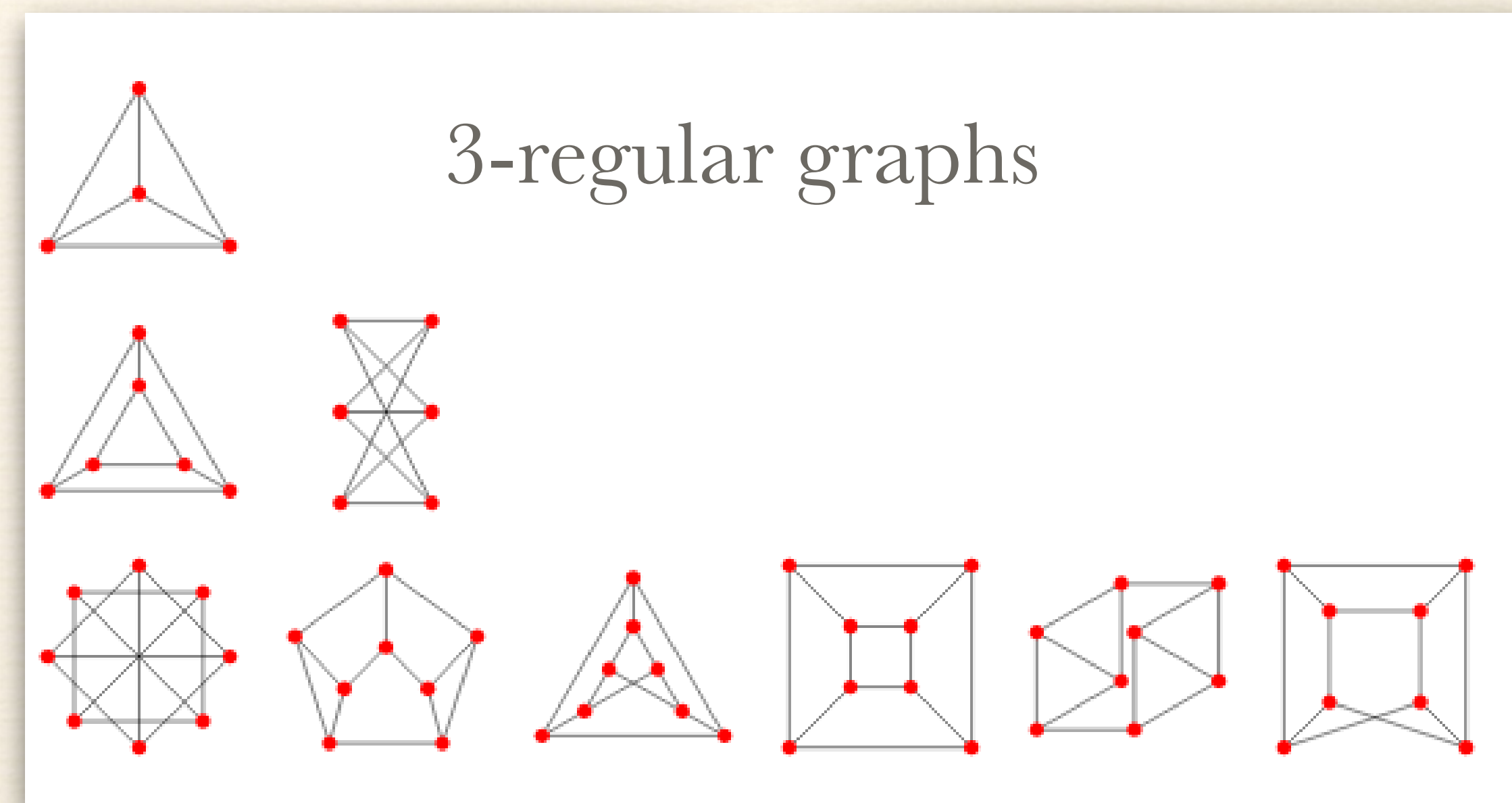
Arvind et al.: *On WL invariance: Subgraph Counts and related properties* (2019)

M. Grohe. *The logic of graph neural networks* (2021)

ρ (colour refinement)



- ❖ Cannot count cycles (triangles)
- ❖ Cannot distinguish d-regular graphs
- ❖ Only tree information



Back to **MPNNs**

MPNNs & Colour refinement

Theorem (Morris et al. 2019, Xu et al. 2019)

If colour refinement cannot tell two graphs apart
then neither can any MPNN!

MPNNs

$\xi^{(0)}(G, v) :=$ Hot-one encoding of label of vertex v

$\xi^{(t)}(G, v) := \text{Upd}^{(t)}\left(\xi^{(t-1)}(G, v), \text{Agg}^{(t)}\left(\{\{\xi^{(t-1)}(G, v), \xi^{(t-1)}(G, u) \mid u \in N_G(v)\}\}\right)\right)$

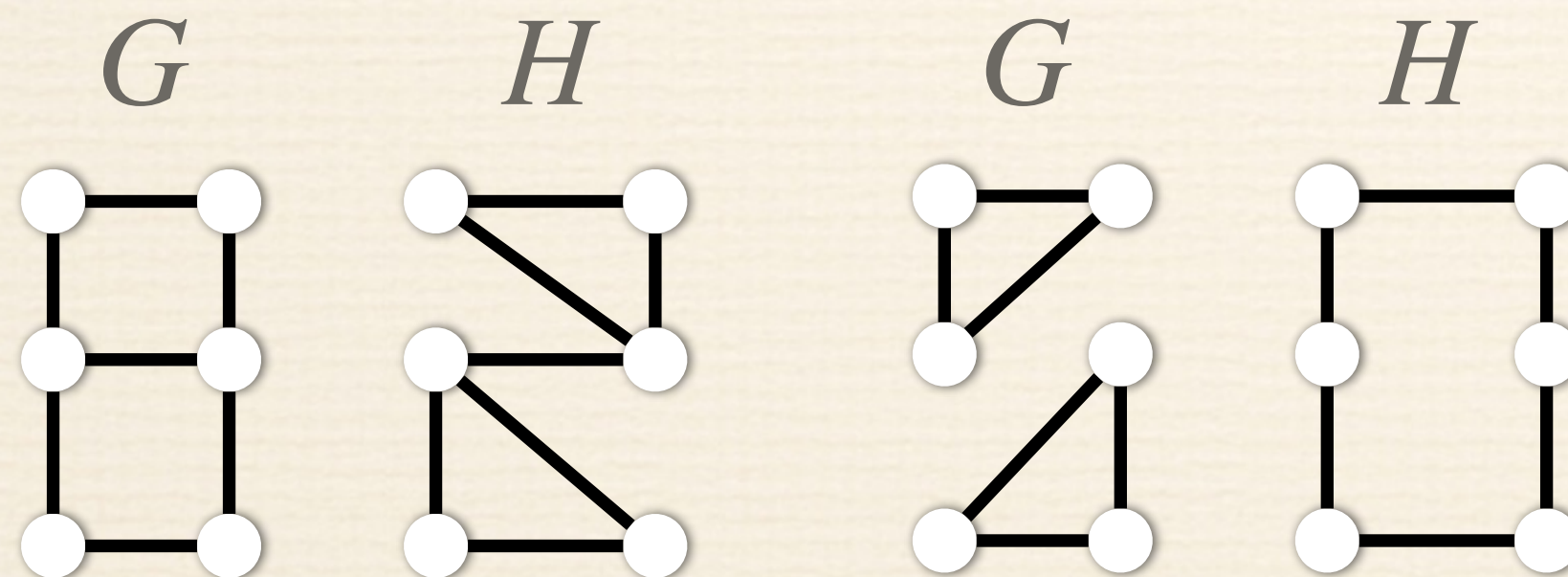
$\rho(G) := \text{Readout}\left(\{\{\xi^{(L)}(G, v) \mid v \in V_G\}\}\right)$

Color refinement

$\text{cr}^{(0)}(G, v) :=$ Initial label of v

$\text{cr}^{(t)}(G, v) := \text{Hash}\left(\text{cr}^{(t-1)}(G, v), \{\{\text{cr}^{(t-1)}(G, u) \mid u \in N_G(v)\}\}\right)$

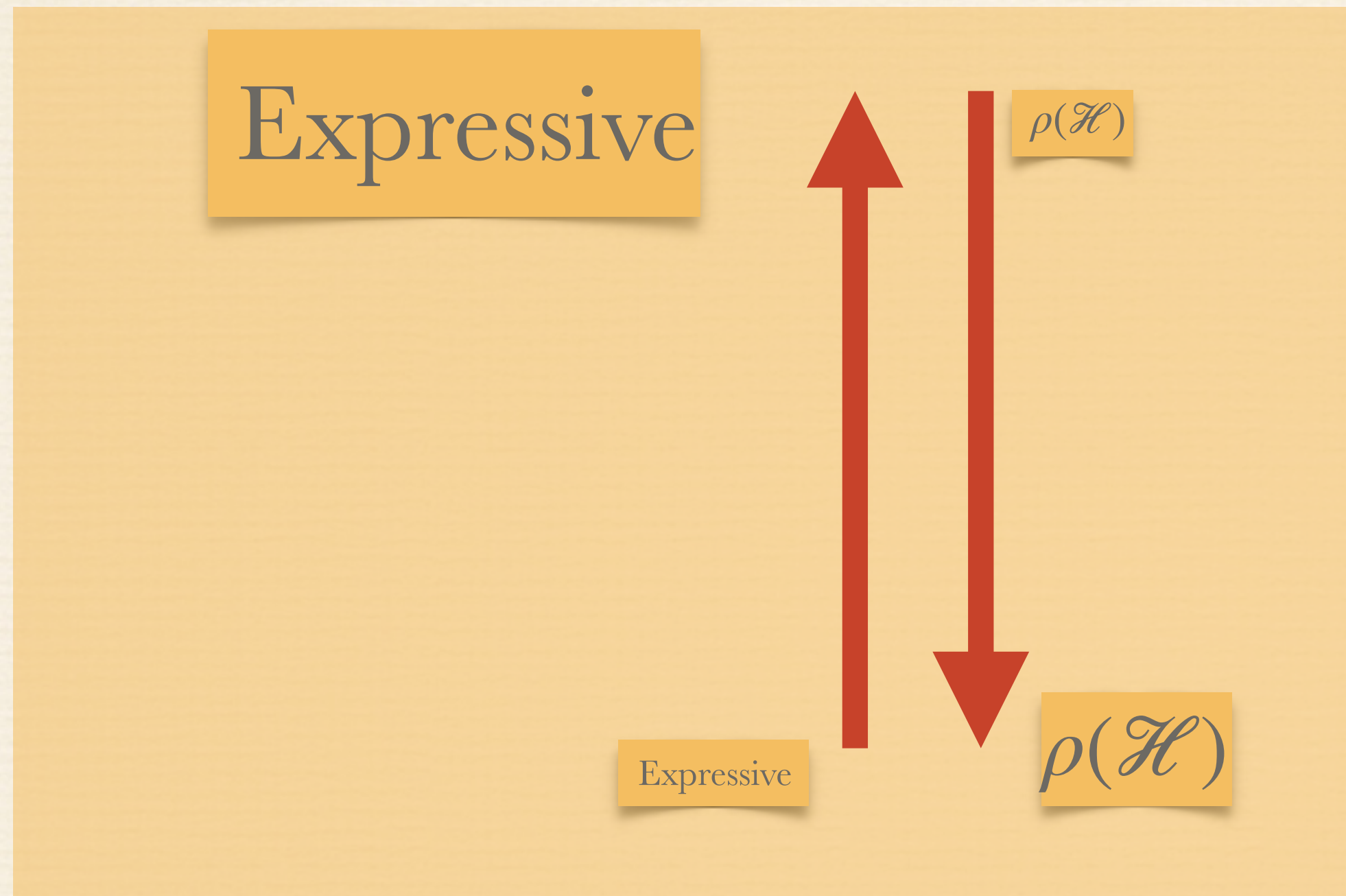
$\rho(G) := \{\{\text{cr}(G, v) \mid v \in V_G\}\}$



→ No MPNN can separate these graphs

MPNNs & Colour refinement

Recall:



We have just shown: $\rho(\text{colour refinement}) \subseteq \rho(\text{MPNNs})$

Expressive power of MPNNs is **upper bounded** by colour refinement

Lower bound?

- ❖ We have seen that MPNNs cannot separate more graphs than colour refinement.
- ❖ Can colour refinement separate **more graphs than MPNNs?** No!

Theorem (Morris et al. 2019)

There exists a GNN 101 which can embed G and H differently when colour refinement assigns them different colours

- ❖ The class of MPNNs is as powerful (or weak) as colour refinement

What else can we say?

$$\rho(\text{colour refinement}) = \rho(\text{MPNNs})$$

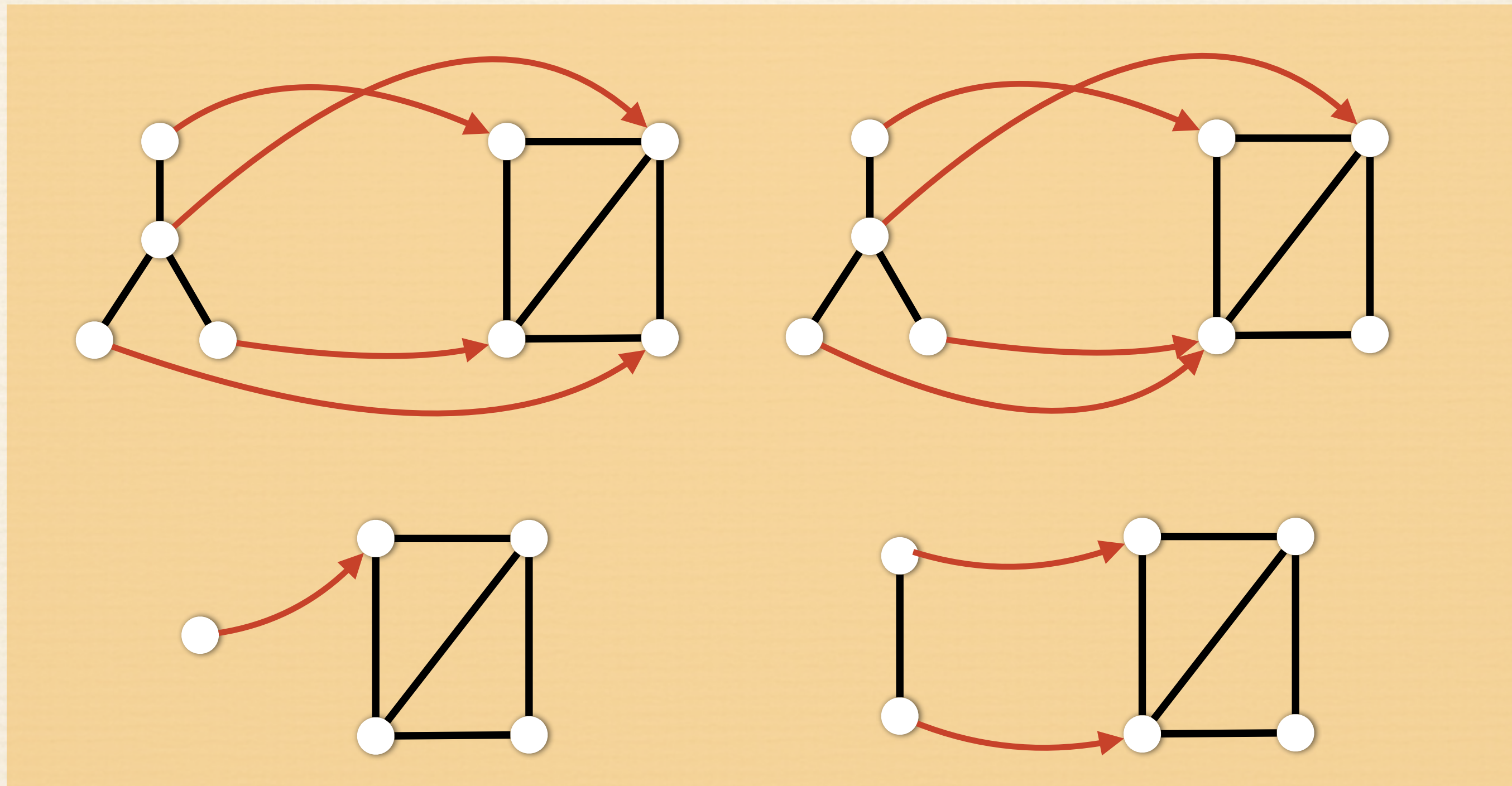


Other - more insightful - characterisations?

A detour to homomorphism counts

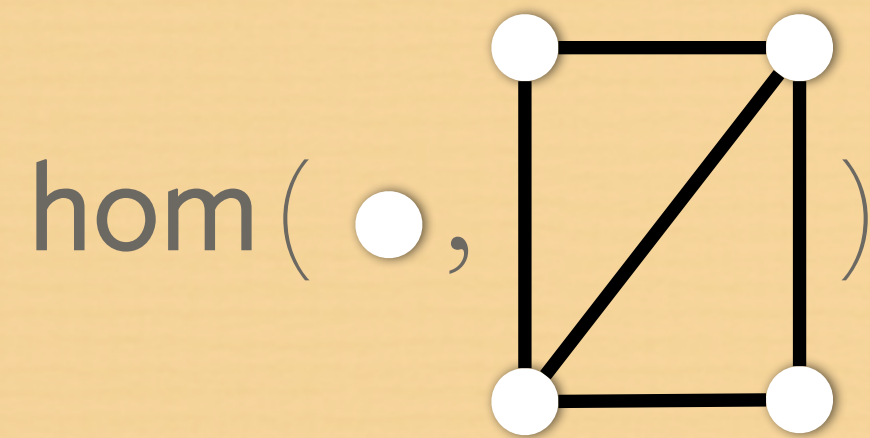
Homomorphisms

- ❖ Let $P = (V_P, E_P, L_P)$ and $G = (V_G, E_G, L_G)$ be graphs.
- ❖ A function $h : V_P \rightarrow V_G$ is a **homomorphism** if it is **edge preserving** $(v, w) \in E_P \Rightarrow (h(v), h(w)) \in E_G$ and **label preserving**.

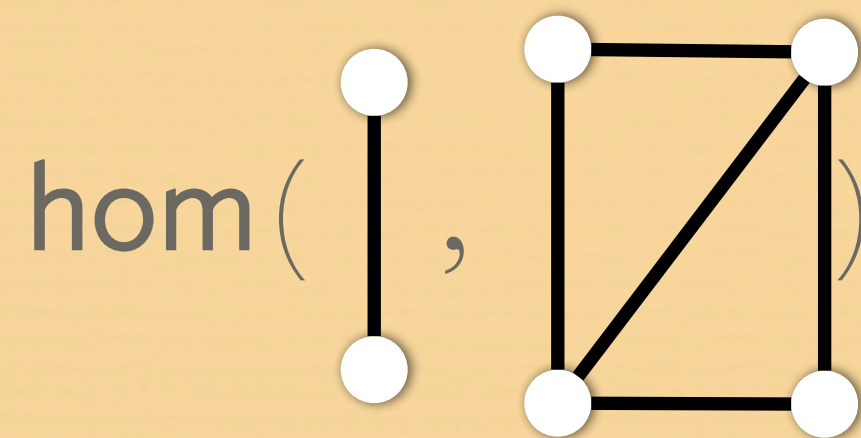


Homomorphism counts

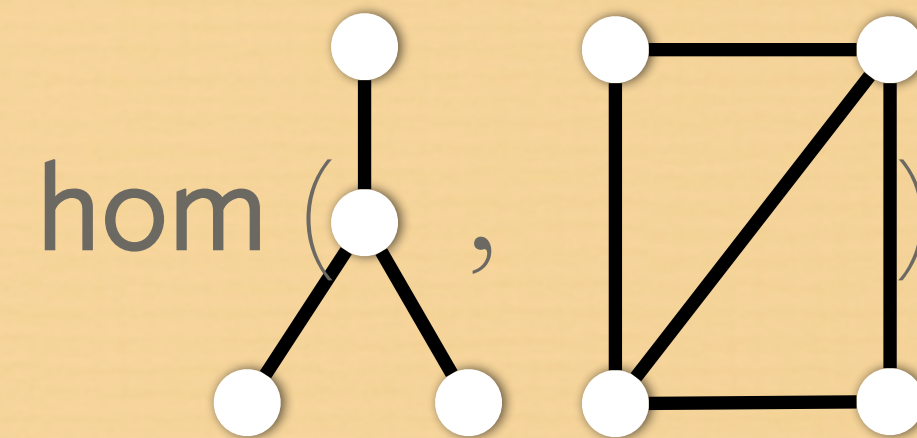
- ❖ Define $\text{HOM}(P, G) := \{ \text{all homomorphisms from } P \text{ to } G \}$
- ❖ Define $\text{hom}(P, G) := |\text{HOM}(P, G)|$.



#vertices = 4



2#edges=10



$$70 = 2 \cdot 2^3 + 2 \cdot 3^3$$

Homomorphisms

- ❖ Weaker notion than subgraph isomorphism (see later).
- ❖ Underlies semantics of many graph query languages.
- ❖ Algebra of homomorphism counts: A rich and active area of research.

MPNNs and hom counts

Theorem (Dell et al. 2019, Dvorák 2010)

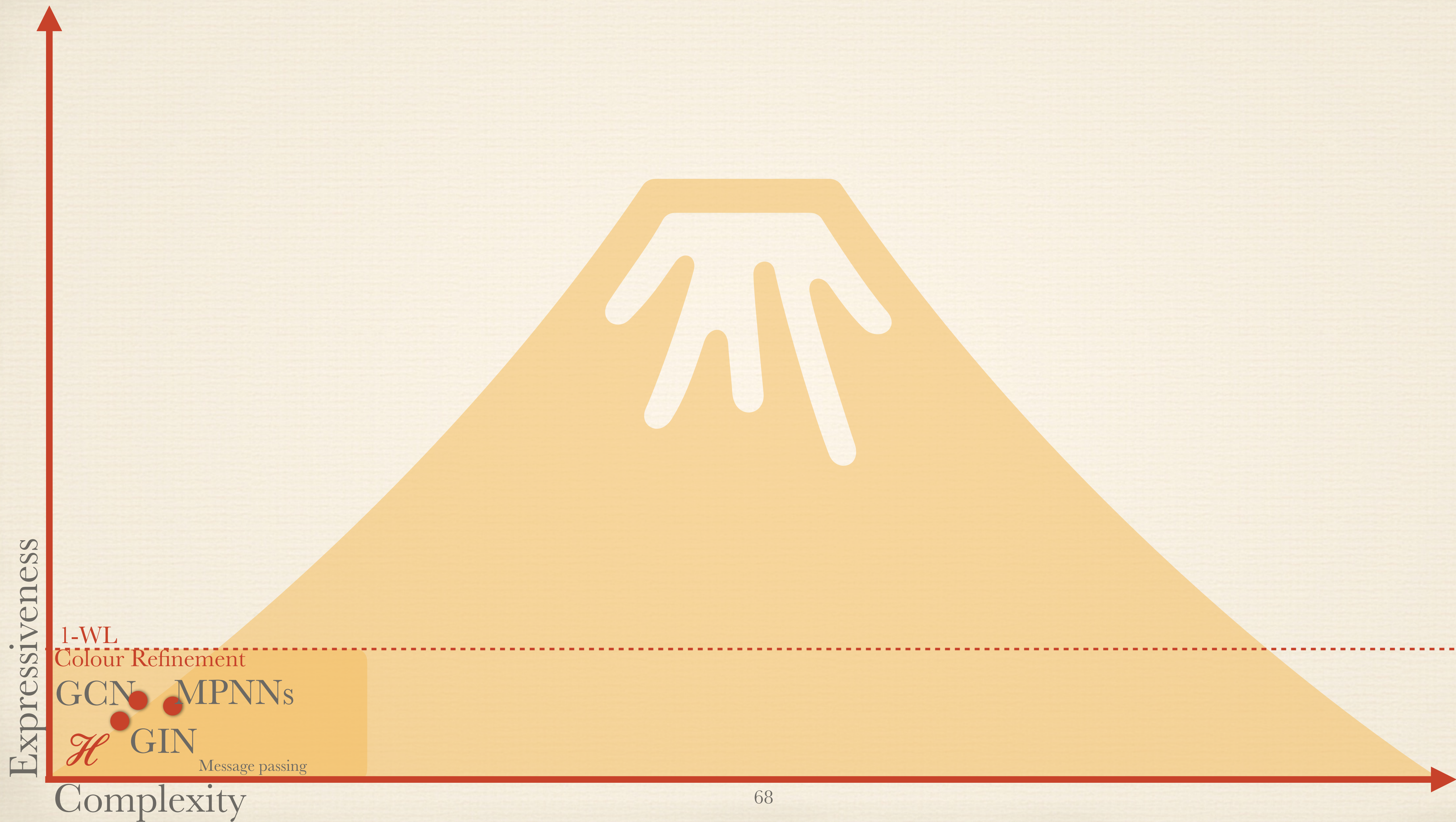
$\text{hom}(T, G) = \text{hom}(T, H)$ for all *trees* T
if and only if
colour refinement cannot distinguish G from H .

Corollary

$\text{hom}(T, G) = \text{hom}(T, H)$ for all *trees* T if and only if no
MPNN can distinguish G from H .

Follows from $\rho(\text{cr}) = \rho(\text{MPNN})$

❖ MPNNs can only detect **tree information** from a graph!



Beyond distinguishing power?

- ❖ Approximation properties (universality)
- ❖ Logical expressiveness
- ❖ Generalization

Approximation properties

- ❖ Equip set of graphs \mathcal{G} with a **topology** and assume that \mathcal{H} consists of **continuous graph embeddings** from \mathcal{G} to \mathbb{R} .
- ❖ Let $\mathcal{C} \subseteq \mathcal{G}$ be a **compact set** of graphs.

Stone-Weierstrass

Theorem (Azizian & Lelarge 2021, G. and Reutter 2022)

If \mathcal{H} is **closed under linear combinations and products**, then \mathcal{H} can approximate any continuous function $\Xi : \mathcal{C} \rightarrow \mathbb{R}$ satisfying

$$\rho(\mathcal{H}) \subseteq \rho(\{\Xi\}).$$

- ❖ Can be generalised to general embeddings with output space \mathbb{R}^d

MPNNs: Approximation

Theorem (Azizian & Lelarge 2021, G. and Reutter 2022)

On compact set of graphs, MPNNs can approximate any continuous graph embedding $\Xi : \mathcal{G} \rightarrow \mathbb{R}$ satisfying $\rho(\text{colour refinement}) \subseteq \rho(\{\Xi\})$

- ❖ We know $\rho(\text{MPNNs}) = \rho(\text{colour refinement})$
- ❖ Update functions can be used to approximate product and take linear combinations of MPNNs
- ❖ Intricate relation between distinguishing power and approximation properties

Universality and graph isomorphism

Theorem (Chen et al. (2019))

In order for a class of methods to be able to approximate **any (invariant) continuous functions**, the class of methods should be able to distinguish any two **non-isomorphic** graphs.

Proof

Minimal size $\rho(\mathcal{H}) \subseteq \rho(\{\Xi\})$



$$(G, H) \in \rho(\mathcal{H}) \Leftrightarrow G \cong H$$

Logical expressiveness

- ❖ Finite variable logics.
- ❖ Extension with Presburger quantifiers.

Colour refinement (again)

I mentioned that ρ (colour refinement) has **many characterisations**.

Of interest is also a **logical one**, in particular **First-order logic with 2 variables and counting quantifiers (C_2)**.

$$\varphi(x) = \exists^{\leq 5} y \left(E(x, y) \wedge \exists^{\geq 2} x \left(E(y, x) \wedge L_a(x) \right) \right)$$

↑
binary edge predicate

↑
unary label predicate

Given graph G , vertex $v \in V_G$ **satisfies** φ : It has at most 5 neighbours
 $(G, v) \models \varphi$ each with at least to neighbours labeled “a”

Colour refinement and C_2

Theorem (Cai et al. 1992)

Two graphs have the same colour histogram after t iterations of colour refinement *if and only if* they satisfy the same C_2 sentences of quantifier depth t

$$\rho(\text{colour refinement}) = \rho(\text{MPNNs}) = \rho(C_2)$$

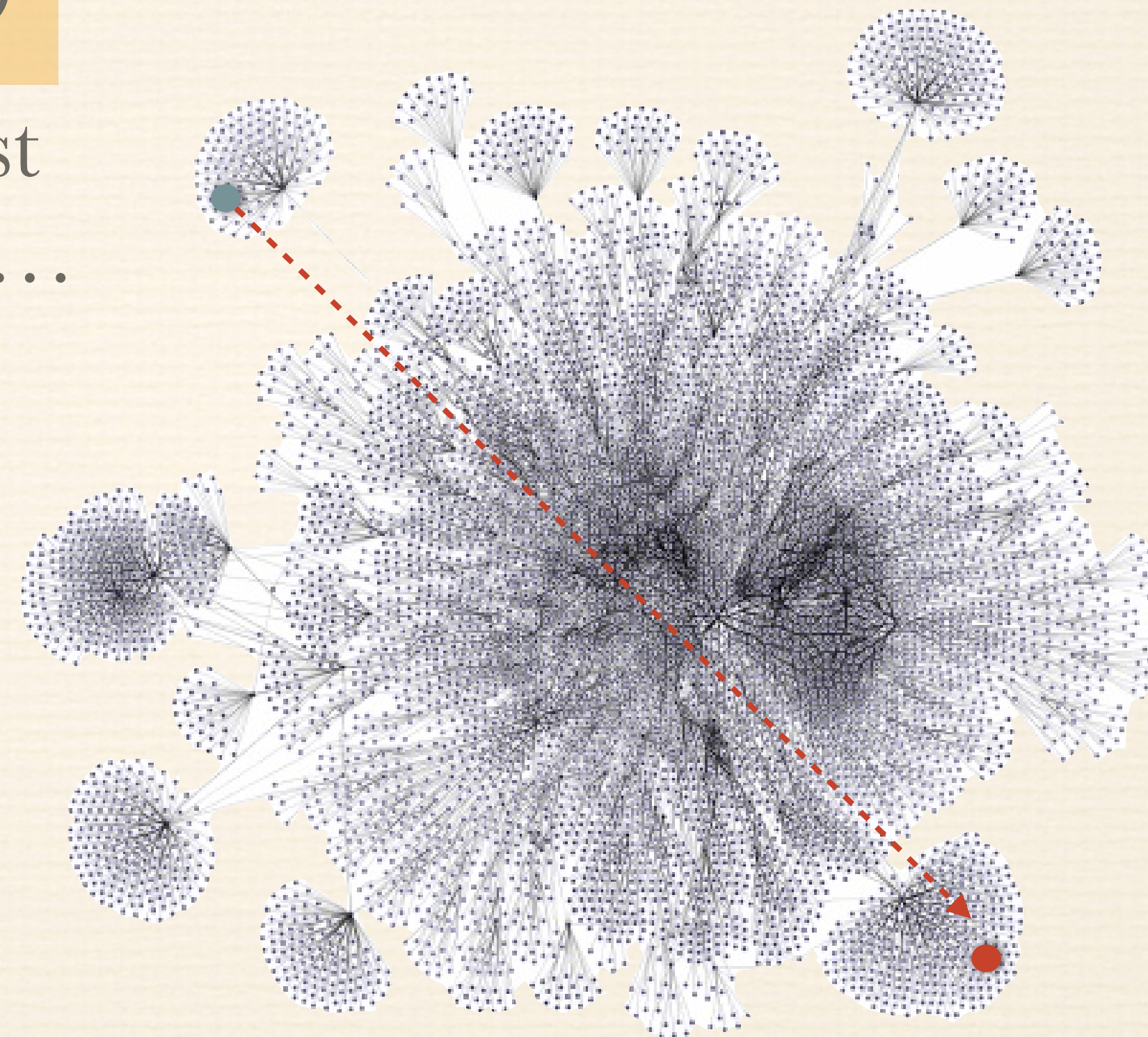
What about vertices?

Which unary C_2 formulas can MPNNs express?

\mathcal{H} can \mathcal{C} -express Ξ if there exists a $\xi \in \mathcal{H}$ such that for all $G \in \mathcal{C}, \mathbf{v} \in V_G^p : \xi(G, \mathbf{v}) = \Xi(G, \mathbf{v})$

❖ Not all: $\varphi(x) := L_b(x) \wedge \exists y L_r(y)$

I am blue and there exist
a red vertex somewhere...



Cannot be reached by message passing!

Which unary C_2 formulas can MPNNs express?

- ❖ Not all: $\varphi(x) := L_b(x) \wedge \exists y L_r(y)$
- ❖ **Graded modal logic**: syntactical fragment of C_2 in which quantifiers are of the form $\exists^{\geq N} (E(x, y) \wedge \varphi'(y))$

Theorem (Barceló et al. 2020)

Let $\varphi(x)$ be a unary FO formula. Then, $\varphi(x)$ is equivalent to a graded modal logic formula *if and only if* $\varphi(x)$ is expressible by the class of MPNNs.

$$\exists \xi \in \text{MPNNs} : \forall G \in \mathcal{G}, \forall v \in V_G : (G, v) \models \varphi \Leftrightarrow \xi(G, v) = 1$$

Role of activation functions

Theorem

Let $\varphi(x)$ be a unary FO formula. Then, $\varphi(x)$ is equivalent to a graded modal logic formula *if and only if* $\varphi(x)$ is expressible by the class of MPNNs.

❖ Proof relies on sign, ReLU, trReLU activation function.

Theorem (Sammy Khalife 2023)

There is a $\varphi(x)$ in graded modal logic that is not expressible by MPNNs using *polynomial activation functions*.

MPNN+: Extended MPNNs

- ❖ Can we extend MPNNs such that **all C_2 formulas** (including $\varphi(x) := L_b(x) \wedge \exists y L_r(y)$) can be expressed?

$$\xi^{(t)}(G, v) := \text{Upd}^{(t)}\left(\xi^{(t-1)}(G, v), \text{Agg}^{(t)}\left(\{\{\xi^{(t-1)}(G, v), \xi^{(t-1)}(G, u) \mid u \in N_G(v)\}\}\right)\right)$$

Add global aggregation in every layer



$$\xi^{(t)}(G, v) := \text{Upd}^{(t)}\left(\xi^{(t-1)}(G, v), \text{Agg}^{(t)}\left(\{\{\xi^{(t-1)}(G, v), \xi^{(t-1)}(G, u) \mid u \in N_G(v)\}\}, \text{Global}^{(t)}\left(\{\{\xi^{(t-1)}(G, u) \mid u \in V_G\}\}\right)\right)\right)$$

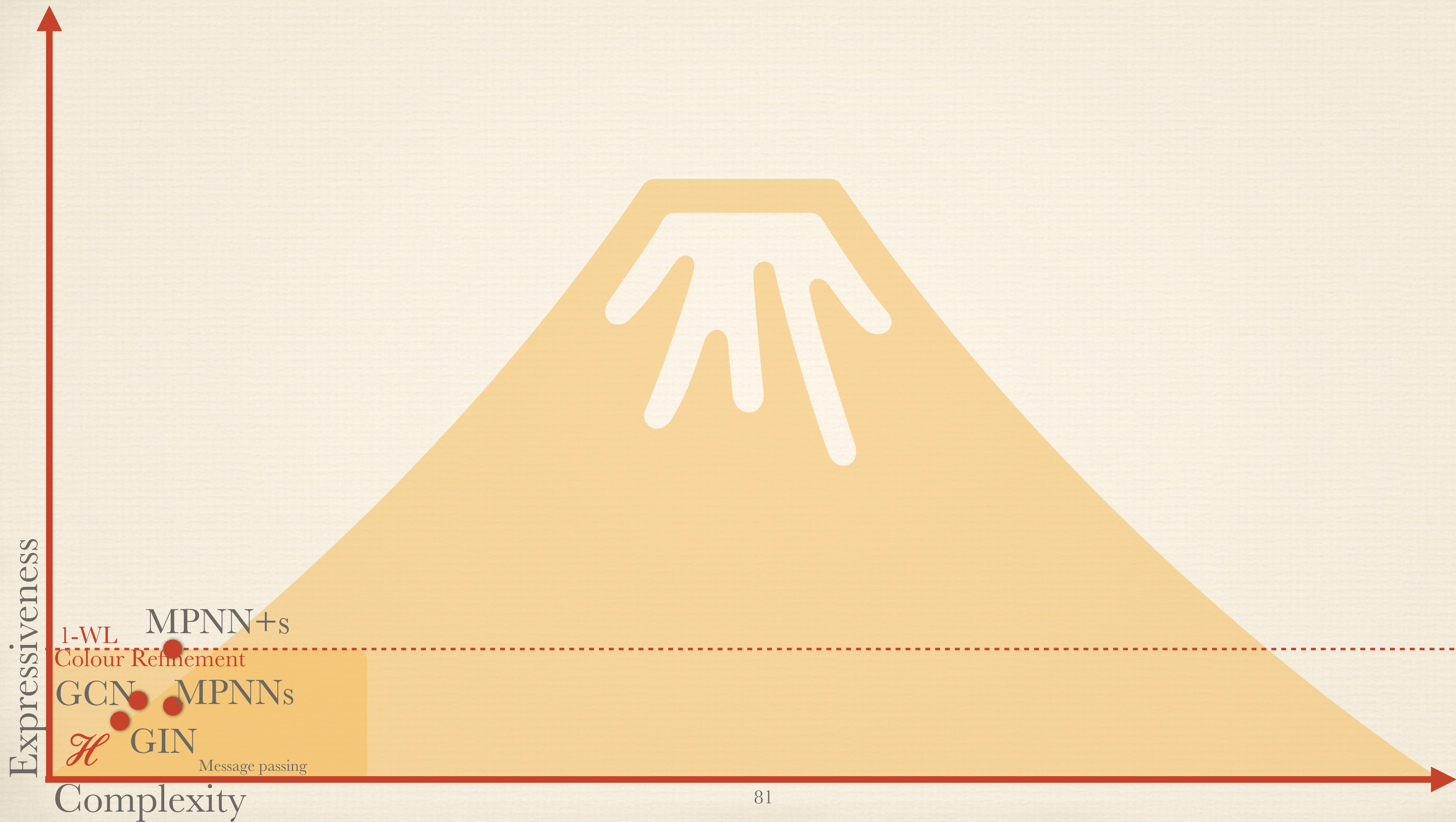
MPNN_s+

Theorem (Barceló et al. 2020)

Every unary C_2 formula $\varphi(x)$ is expressible by the class of MPNN_s+

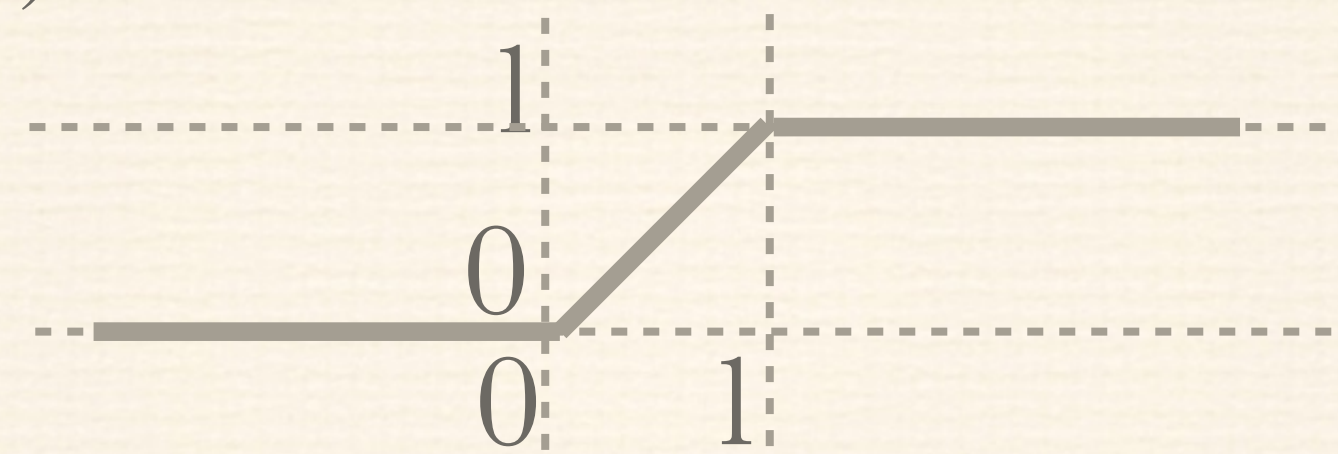
- ❖ The corresponding colour refinement version is known as the **one-dimensional Weisfeiler-Leman algorithm** or **1-WL** on vertices.

$$\rho(1\text{-WL}) = \rho(\text{MPNN}_{\text{s}}^+)$$



Wait a moment! MPNNs go easily beyond FO!

$$\text{trRelU}\left(\sum_{u \in N_G(v)} P_r(u) - \sum_{u \in N_G(v)} P_b(u)\right) = \begin{cases} 1 & v \text{ had more red than blue neighbors} \\ 0 & \text{otherwise} \end{cases}$$



This property is known **not to be expressible** as an FO formula $\varphi(x)$

Proof

By means of *locality* of FO

or

By playing so-called *Erhenfeucht-Fraisse* game.

Solution? Add more complex quantifiers!

$$\varphi(x) = \exists^{\geq 2} y (E(x, y) \wedge L_r(y)) \longrightarrow \varphi(x) = (\#_y[E(x, y) \wedge L_r(y)] \geq 2)$$

true when x has more than two red neighbours

$$\varphi(x) = (\#_y[E(x, y) \wedge L_r(y)] - \#_y[E(x, y) \wedge L_b(y)] \geq 0)$$

true when x has more red than blue neighbours

$$\varphi(x) = \left(\sum_{i=1}^k a_i \#_y[E(x, y) \wedge \psi_i(y)] \leq \delta \right)$$

true when the neighbours of x satisfy the linear inequality

Presburger quantifier

Local and Global

$$\varphi(x) = (\#_y[E(x, y) \wedge L_r(y)] \geq 2)$$

true when x has more than two
red neighbours

$$\varphi(x) = (\#_y[L_r(y)] - \#_y[E(x, y) \wedge L_r(y)] = 0)$$

true when x has all red nodes
in graph as neighbours



global counting



local counting

- ❖ Extending Two-variable FO with Presburger quantifiers:
- ❖ New logics: **MP** (global) and **L-MP** (local)

Uniform characterisation of sum-GNNs

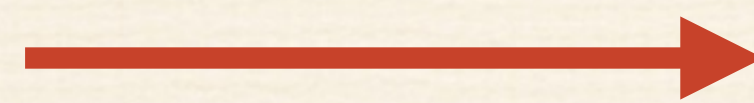
Theorem (Benedikt et al. 2024)

- ❖ The language L-MP is equivalent to sum-GNNs using **eventually constant activation functions**
- ❖ Allowing global aggregation in sum-GNNs, equivalence to MP.

a vertex has more red than blue
neighbours



L-MP expressible



$$\text{trReLU}\left(\sum_{u \in N_G(v)} P_r(u) - \sum_{u \in N_G(v)} P_b(u)\right)$$

sum-GNN expressible

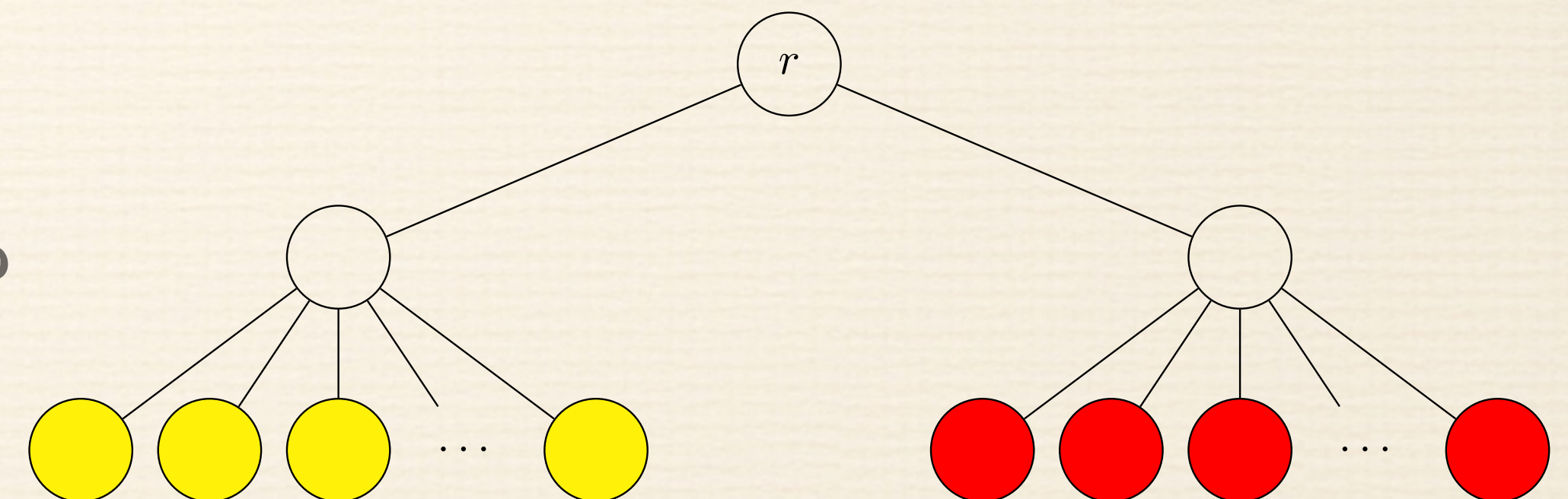
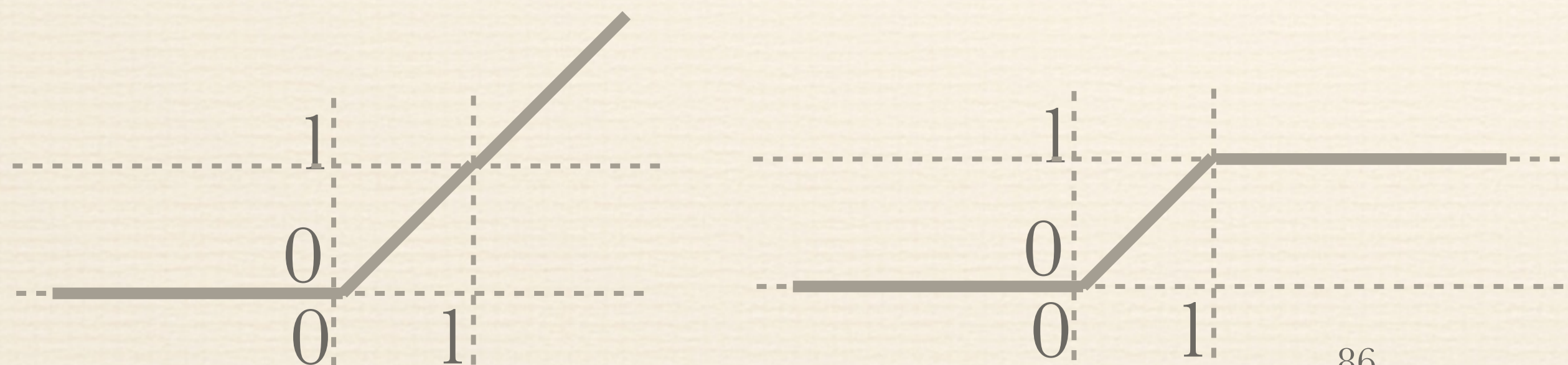
$$\varphi(x) = (\#_y[E(x, y) \wedge L_r(y)] - \#_y[E(x, y) \wedge L_b(y)] \geq 0)$$

Uniform characterisation of sum-GNNs

Theorem (Benedikt et al. 2024)

- ❖ The language L-MP is equivalent to sum-GNNs using **eventually constant activation functions**
- ❖ Allowing global aggregation in sum-GNNs, equivalence to MP.

- ❖ What about **ReLU**?? This activation is not eventually constant. Theorem fails.
- ❖ What about **max**, **avg** as aggregation functions?



Generalisation and expressivity

- ❖ We want the generalisation error $L_{\mathcal{D}}(\xi) - L_{\mathcal{T}}(\xi)$ to be small.

Theorem (Vapnik and Chervonenkis 1964)

- ❖ For $\delta > 0$, with probability $1 - \delta$ (in our selection of training data \mathcal{T} of size m) for all $\xi \in \mathcal{H}$:

$$L_{\mathcal{D}}(\xi) - L_{\mathcal{T}}(\xi) \leq \sqrt{\frac{2d \log(\frac{em}{d})}{m}} + \sqrt{\frac{\log(\frac{1}{\delta})}{2m}}$$

- ❖ where d is the **VC dimension** of \mathcal{H}

VC dimension

❖ A set of graphs G_1, \dots, G_d is **shattered** by an embedding class \mathcal{H}

if, for **any** labeling $y_1, \dots, y_d \in \{0,1\}^d$

we can find an embedding $\xi \in \mathcal{H}$ (which may depend on the labeling)

such that $\xi(G_1) = y_1, \dots, \xi(G_d) = y_d$

❖ **VC dimension** = maximal number of graphs that can be shattered.

VC dimension

❖ **VC dimension** = maximal number of graphs that can be shattered.

Let us assume we consider graphs up to size n : \mathcal{G}_n

Theorem

The VC dimension of MPNNs on \mathcal{G}_n is **bounded** by the number of graphs in \mathcal{G}_n that can be distinguished by MPNNs.

We can also show matching lower bound.

Corollary

The VC dimension of MPNNs on all graphs is unbounded.

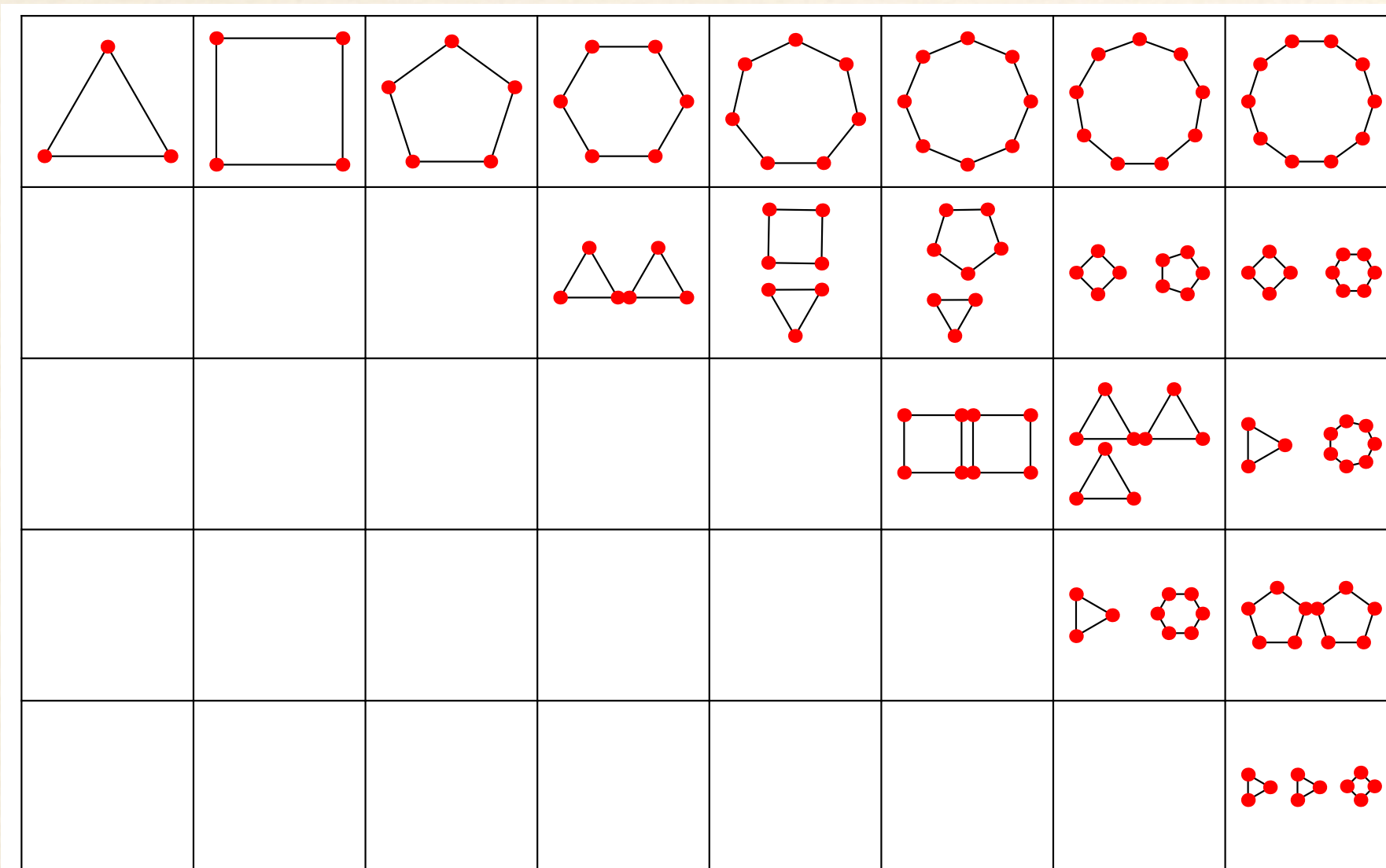
Colour complexity

If color refinement does not need many colours for a graph: **low colour complexity**

We can get smaller bounds on number of distinguished graphs.



VC dimension is small  Less training data needed



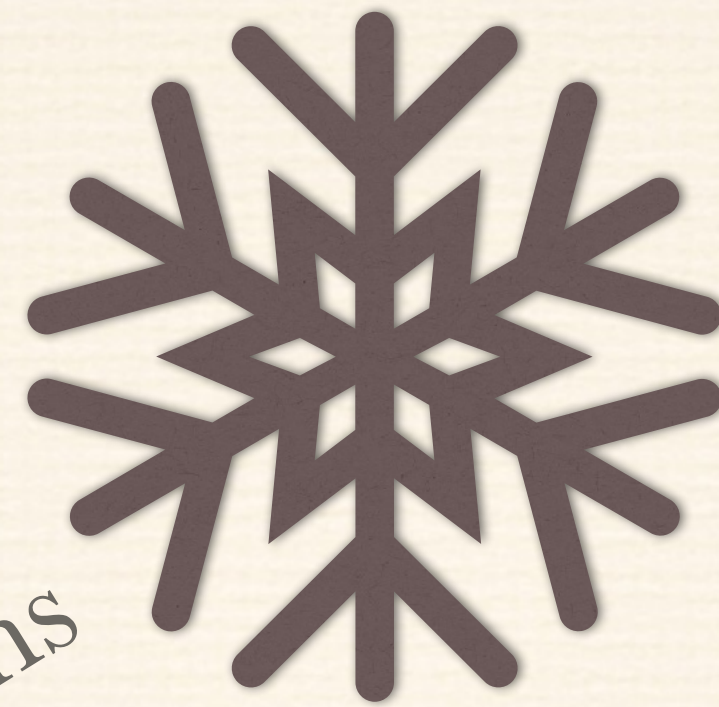
on 2-regular graphs only one colour is needed.

MPNNs only distinguish based on size

Generalisation and expressivity

Only tip of iceberg

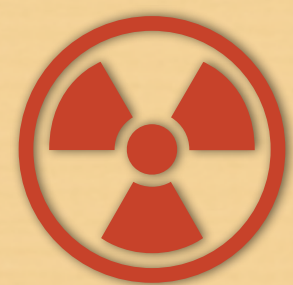
Continuity and covering numbers



Distance measures

graphons

graph neural tangent kernels



Understanding precise impact of expressiveness on generalization, not well understood yet

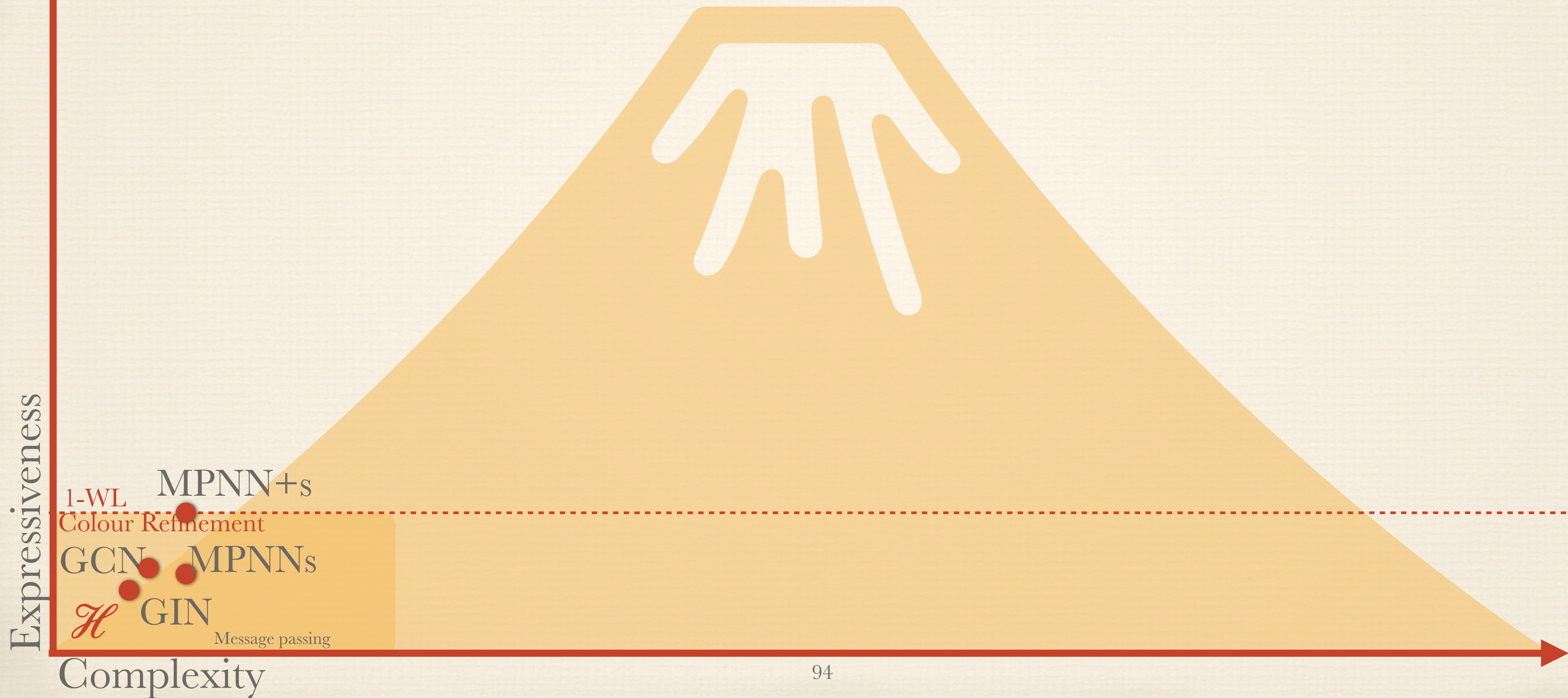
Questions?



More powerful methods

Boosting expressive power

More expressive MPNNs?



How to beyond MPNNs?

Theoretical research guides architecture design!

- ❖ Feature augmentation
- ❖ subgraph GNNs.
- ❖ Higher-order MPNNs

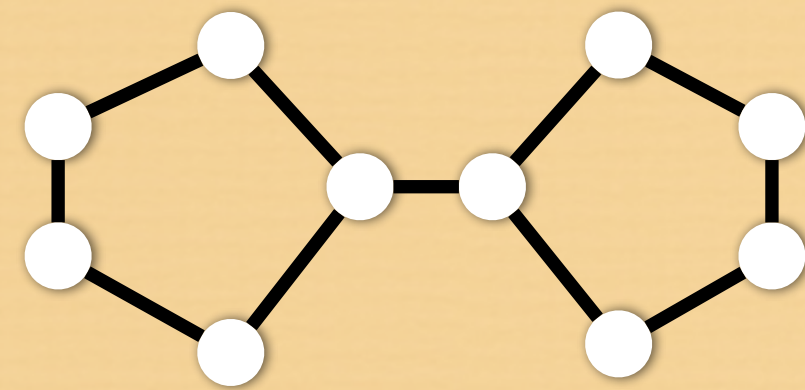


Feature Augmentation

Boost the expressive power by adding information

Feature engineering

- ❖ Deep learning and MPNNs have replaced “old school” **feature engineering approach**.



- ❖ Number of edges
- ❖ Number of cycles of length 5
- ❖ Centrality measures



\mathbb{R}^d



SVM

- ❖ MPNNs were supposed to learn such **features automatically** ...

Idea #1: Adding expressive features

Theorem

Recall:

$\text{hom}(T, G) = \text{hom}(T, H)$ for all *trees* T if and only if
no MPNN can distinguish G from H .

- ❖ What if we add **subgraph information** before doing message-passing?



More than trees

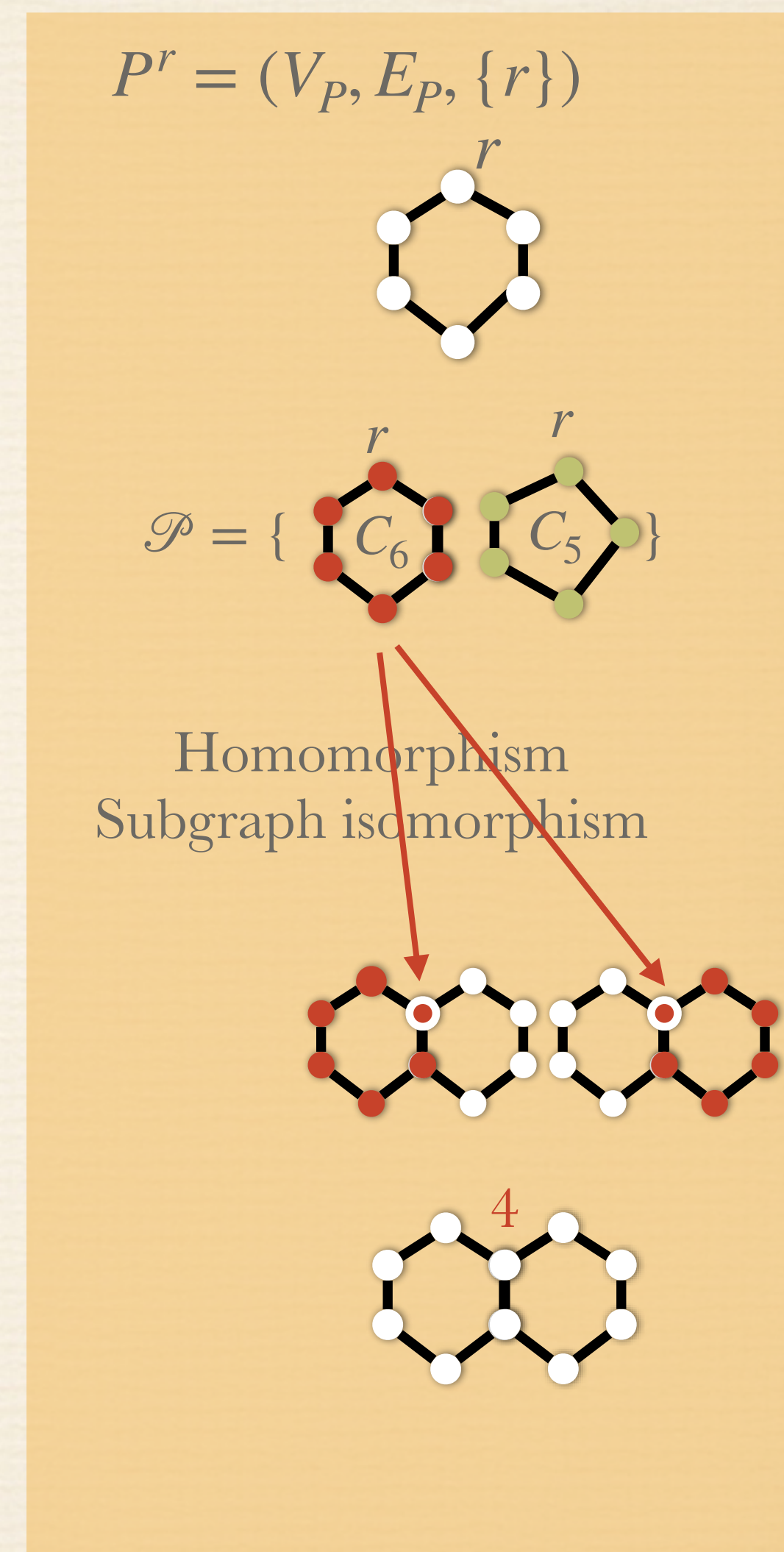
Structural encodings

1. Choose collection of rooted graph patterns/motifs

$$\mathcal{P} := \{P_1^r, \dots, P_\ell^r\}$$

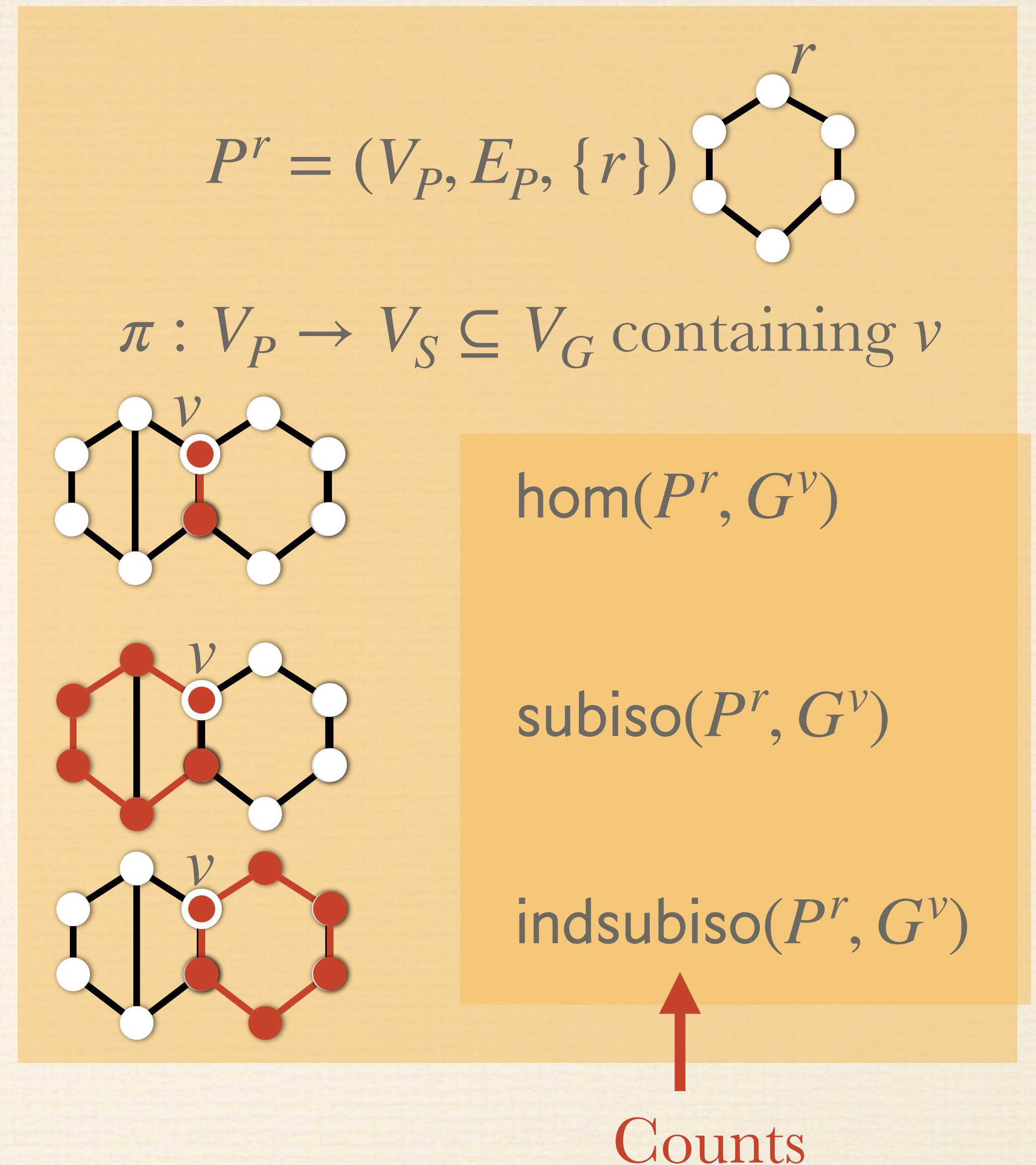
2. Choose how to match subgraphs in \mathcal{P} with data graph G

3. Add count of matches to vertices as extended features.



Matches

- ❖ **Homomorphism**: edge preserving
- ❖ **Subgraph isomorphism**: bijection , edge preserving
- ❖ **Induced subgraph isomorphism**:
bijection, edge preserving (both ways)



\mathcal{P} -MPNNs

- ❖ Add structural encoding as **vertex features** and run MPNN

$$\mathcal{P} := \{P_1^r, \dots, P_\ell^r\}$$

\mathcal{P} -MPNNs

$$\xi^{(0)}(G, v) := \text{Hot-one encoding of label of vertex } v + \text{hom}(P_1^r, G^v), \dots, \text{hom}(P_\ell^r, G^v)$$

$$\xi^{(t)}(G, v) := \text{Upd}^{(t)}\left(\xi^{(t-1)}(G, v), \text{Agg}^{(t)}\left(\{\{\xi^{(t-1)}(G, v), \xi^{(t)}(G, u), \text{hom}(P_1^r, G^u), \dots, \text{hom}(P_\ell^r, G^u) \mid u \in N_G(v)\}\}\right)\right)$$

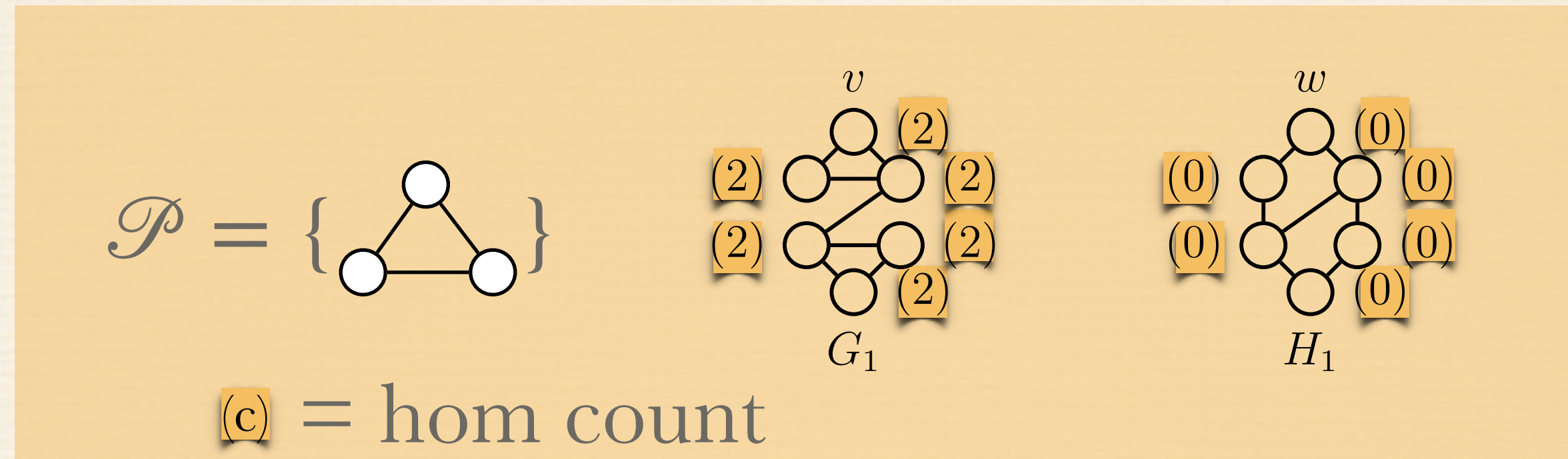
$$\rho(G) := \text{Readout}\left(\{\{\xi^{(L)}(G, v) \mid v \in V_G\}\}\right)$$

hom counts of patterns



- ❖ Did we **increase expressive power**?

\mathcal{P} -MPNNs



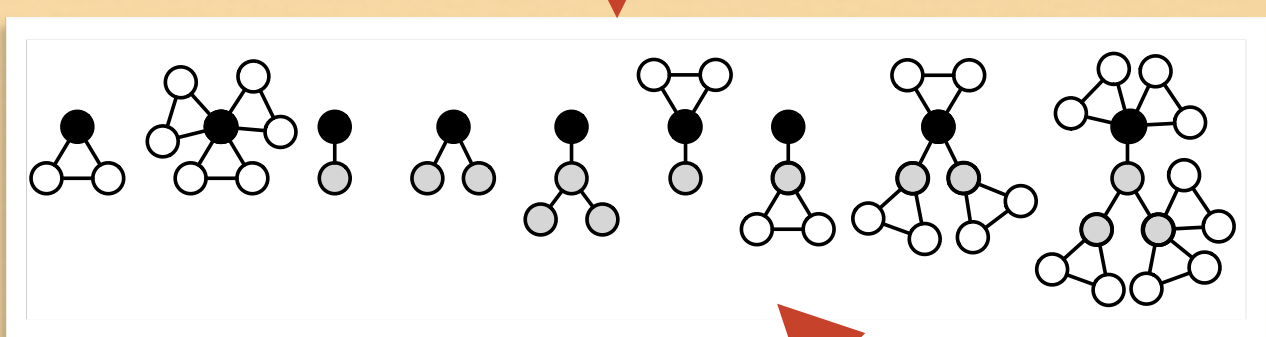
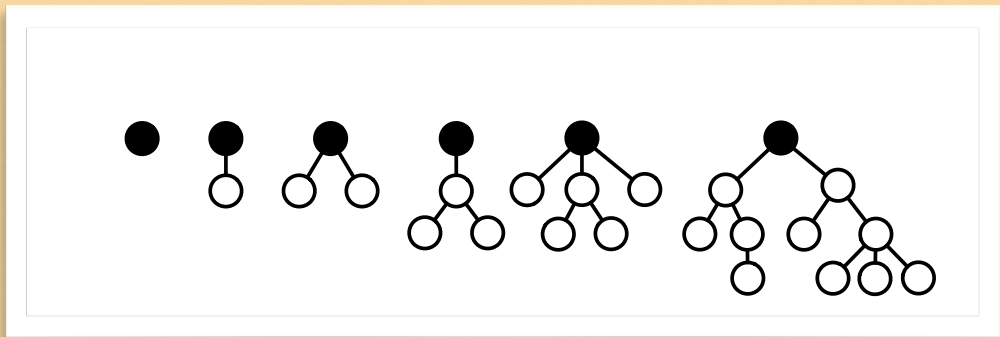
- ❖ We have seen that these graphs equivalent for colour refinement but clearly not for \triangle -MPNNs.
- ❖ So, **increase** in power!
- ❖ What is their **precise expressive power**?

\mathcal{P} -MPNNs: Expressive power

Theorem

$\text{hom}(T, G) = \text{hom}(T, H)$ for all \mathcal{P} -pattern trees T if and only if no P-MPNN can distinguish G from H .

$$\mathcal{P} = \left\{ \begin{array}{c} \text{triangle} \end{array} \right\}$$

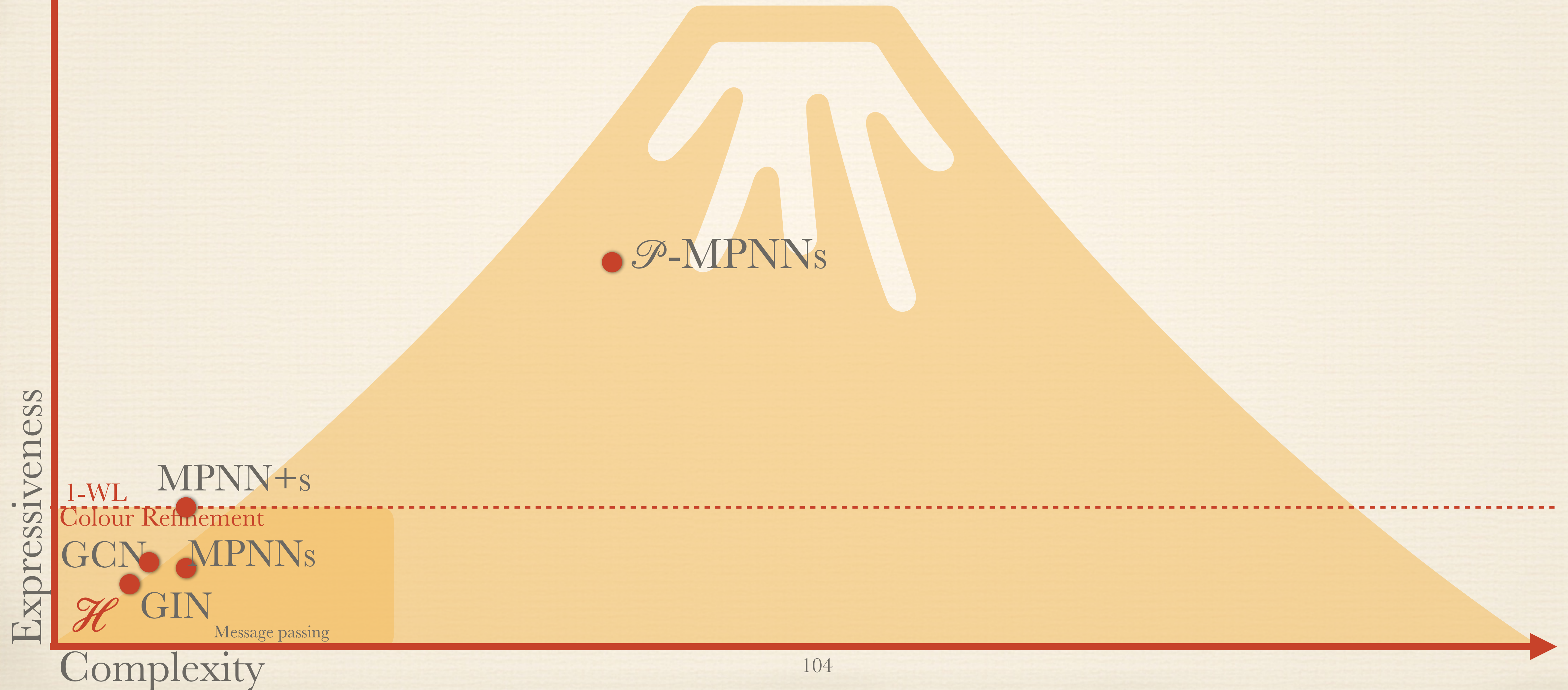


Zinc dataset

SET (\mathcal{F})	MAE
NONE	0.47 ± 0.02
$\{C_3\}$	0.45 ± 0.01
$\{C_4\}$	0.34 ± 0.02
$\{C_6\}$	0.31 ± 0.01
$\{C_5, C_6\}$	0.28 ± 0.01
$\{C_3, \dots, C_6\}$	0.23 ± 0.01
$\{C_3, \dots, C_{10}\}$	0.22 ± 0.01

Take tree: add in each tree vertex
copies of rooted patterns

The larger and complex $\mathcal{P} \Rightarrow$ more complexity counting
 \Rightarrow more expressive power



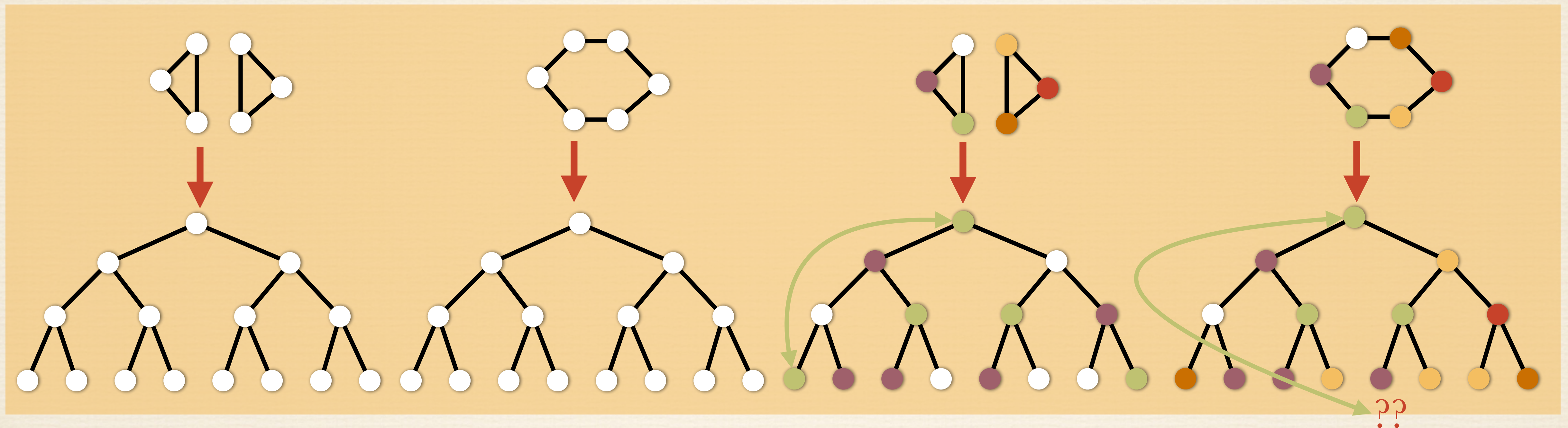
Idea #2: (Random) Vertex identifiers

- ❖ Message-Passing is only based on **vertex features** and **adjacency** information.
- ❖ **Two different** vertices with the **same vertex features** will be treated the **same** (if they have the same colour in colour refinement).

What if we add **vertex identifiers**?

Vertex identifiers

Self identification: useful for cycle detection



In terms of **colour refinement**: every vertex has a **unique colour**

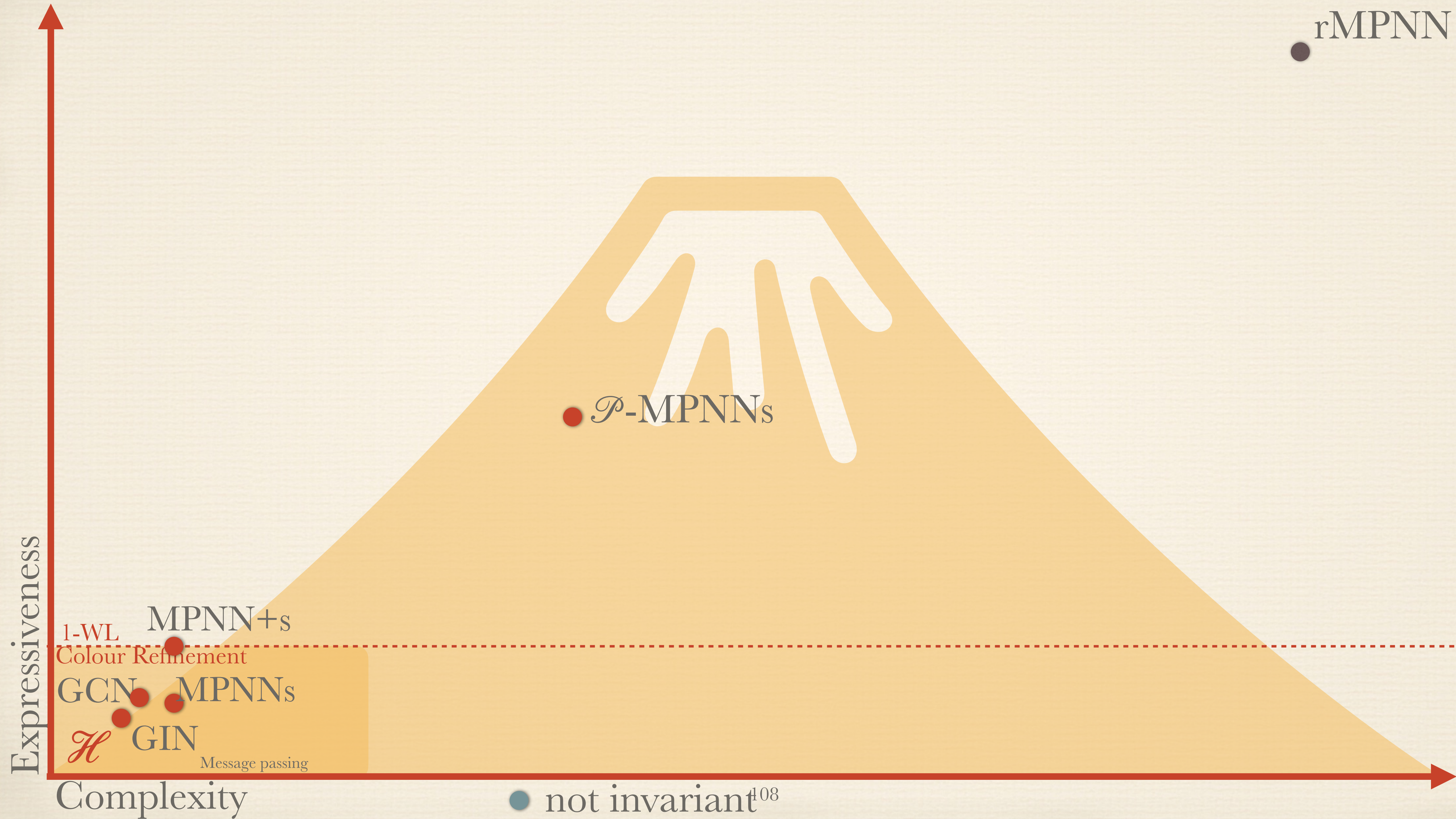
rMPNNs

- ❖ How to choose identifiers? Common choice is **at random!**
- ❖ With high probability random features are vertex identifiers

Theorem

rMPNNs approximate any invariant graph/vertex embedding with high probability

- ❖ Invariance of computed embedding only **in expectation!**

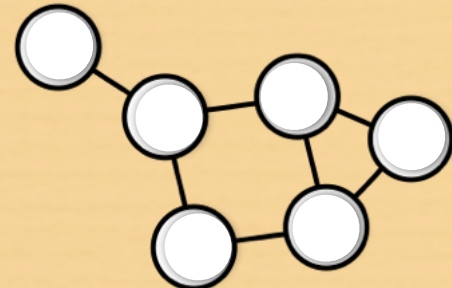


Idea #3: Use global information

- ❖ Extract **global graph information** and use it as positional encodings of vertices
- ❖ **Spectral** information
- ❖ **Shortest paths** (distance information)
- ❖ **Biconnectivity** (connectivity information)

Spectral graph theory

- ❖ **Eigenvalues/vector**: $\mathbf{M} \cdot \mathbf{v} = \lambda \mathbf{v}$
- ❖ For adjacency matrices: Eigenvalues and eigenvectors of **Laplacian**
 $\mathbf{L}_G = \mathbf{D}_G - \mathbf{A}_G$

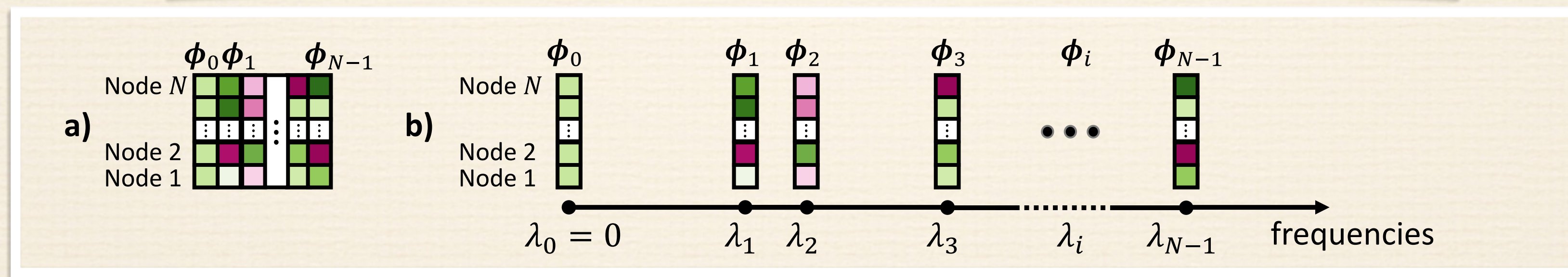
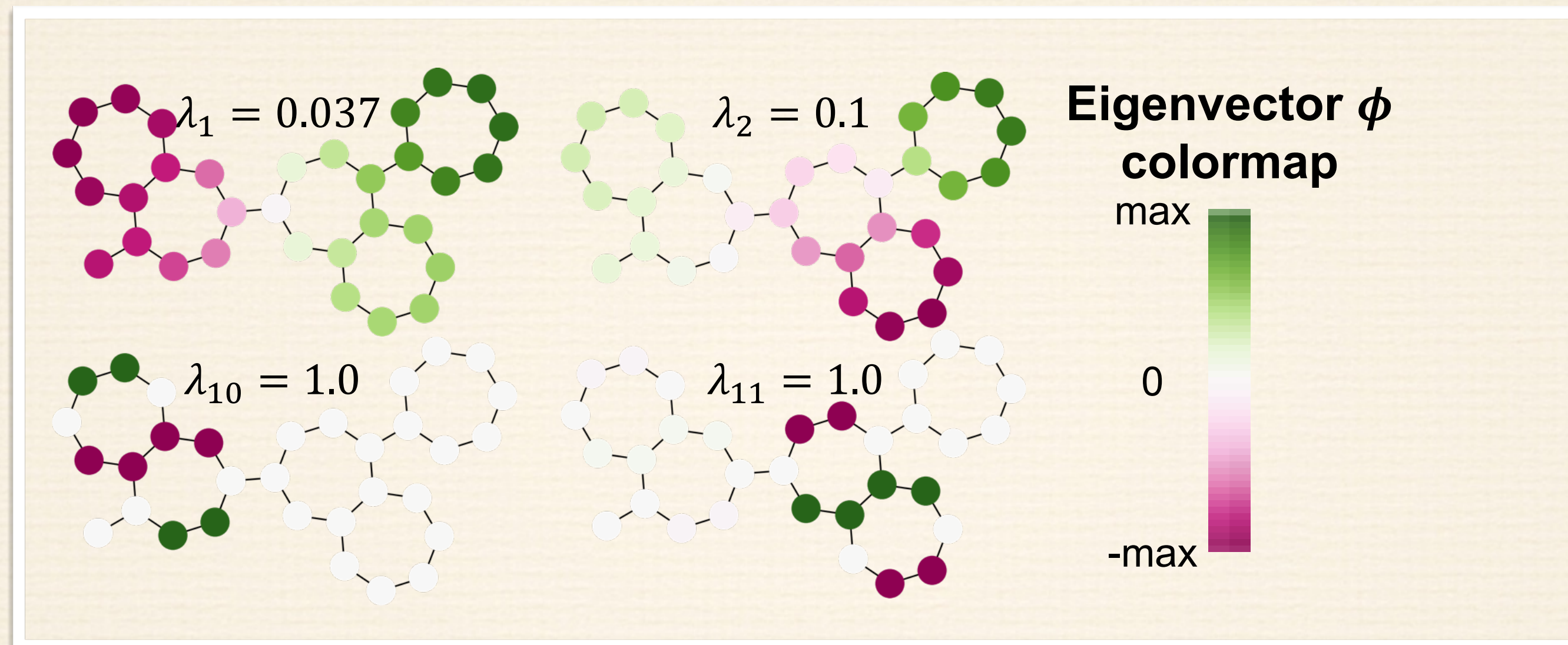


$$\begin{array}{c}
 \mathbf{L}_G \\
 \left(\begin{array}{cccccc}
 2 & -1 & 0 & 0 & -1 & 0 \\
 -1 & 3 & -1 & 0 & -1 & 0 \\
 0 & -1 & 2 & -1 & 0 & 0 \\
 0 & 0 & -1 & 3 & -1 & -1 \\
 -1 & -1 & 0 & -1 & 3 & 0 \\
 0 & 0 & 0 & -1 & 0 & 1
 \end{array} \right)
 \end{array}
 =
 \begin{array}{c}
 \mathbf{D}_G \\
 \left(\begin{array}{cccccc}
 2 & 0 & 0 & 0 & 0 & 0 \\
 0 & 3 & 0 & 0 & 0 & 0 \\
 0 & 0 & 2 & 0 & 0 & 0 \\
 0 & 0 & 0 & 3 & 0 & 0 \\
 0 & 0 & 0 & 0 & 3 & 0 \\
 0 & 0 & 0 & 0 & 0 & 1
 \end{array} \right)
 \end{array}
 -
 \begin{array}{c}
 \mathbf{A}_G \\
 \left(\begin{array}{cccccc}
 0 & 1 & 0 & 0 & 1 & 0 \\
 1 & 0 & 1 & 0 & 1 & 0 \\
 0 & 1 & 0 & 1 & 0 & 0 \\
 0 & 0 & 1 & 0 & 1 & 1 \\
 1 & 1 & 0 & 1 & 0 & 0 \\
 0 & 0 & 0 & 1 & 0 & 0
 \end{array} \right)
 \end{array}$$

- ❖ Laplacian eigenvalues and vectors contain **connectivity information**
 - ❖ multiplicity 1st eigenvalue \sim connected components.

Spectral MPNNs

Add eigenvectors as vertex features



Spectral invariant

$$\mathbf{A} = \sum_{\lambda} \lambda \mathbf{P}_{\lambda} \quad \mathbf{P}_{\lambda} = \begin{pmatrix} p_{11}^{\lambda} & p_{12}^{\lambda} & \dots & p_{1n}^{\lambda} \\ \vdots & \vdots & \ddots & \vdots \\ p_{n1}^{\lambda} & p_{n2}^{\lambda} & \dots & p_{nn}^{\lambda} \end{pmatrix}$$

Multiset

Spectral invariant

$$v \mapsto \text{specinv}(v) := (\lambda, p_{vv}^{\lambda}, \{\{p_{vu}^{\lambda} \mid u \in V_G\}\})_{\lambda \in \Lambda}$$

Graph properties

Number of length 3, 4, or 5 cycles, whether a graph is connected and the number of length k closed walks from any vertex to itself



Beyond 1-WL/Colour Refinement

SpecMPNN

Spectral invariant

$$v \mapsto \text{specinv}(v) := (\lambda, p_{vv}^\lambda, \{\{p_{vu}^\lambda \mid u \in V_G\}\}_{\lambda \in \Lambda})$$

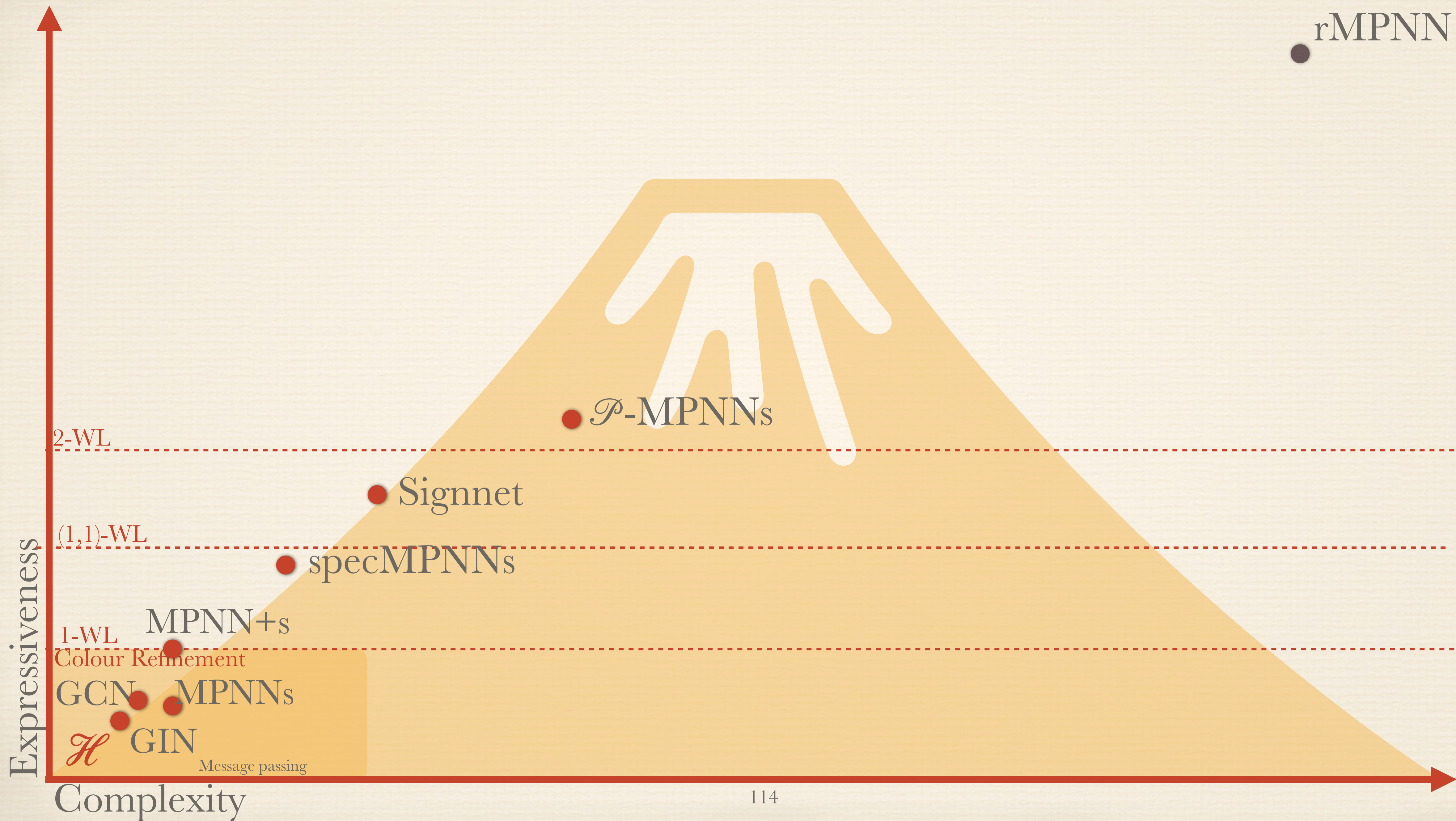
Variation used in Signet and BasisNet  2-WL bound

Can be using combination with any MPNN

Theorem (Seppelt and Rattan (2023))

specMPNN bounded in power by (1,1)-WL and strictly lower than 2-WL

We discuss these WL's later



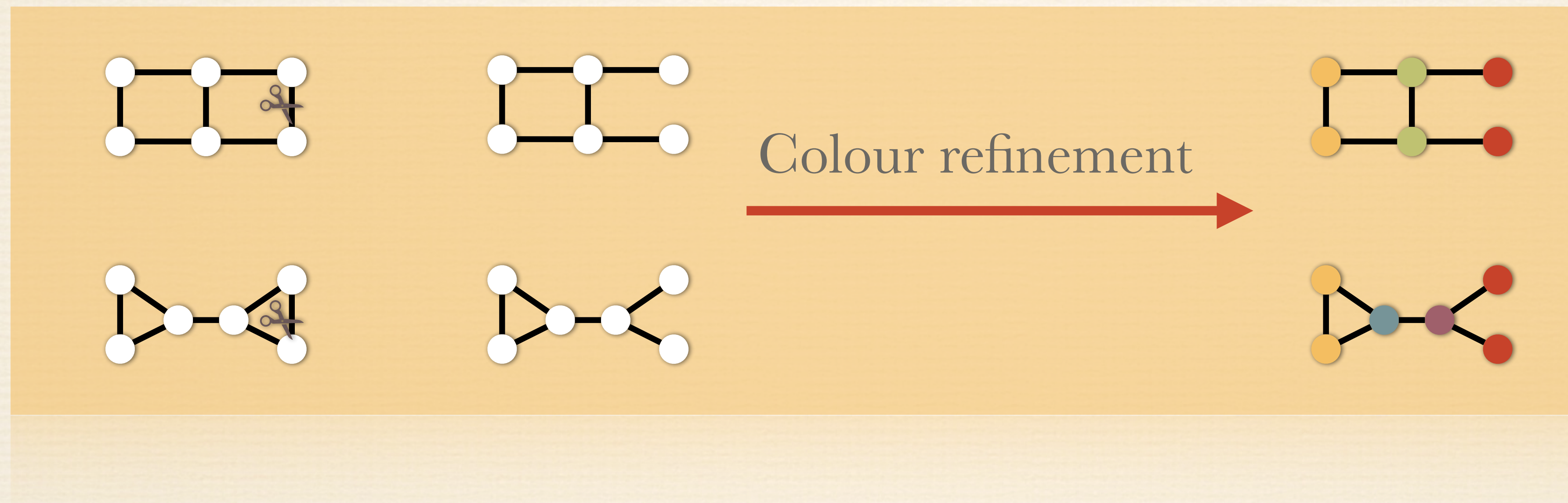


Subgraph GNNs

Turning one graph into many

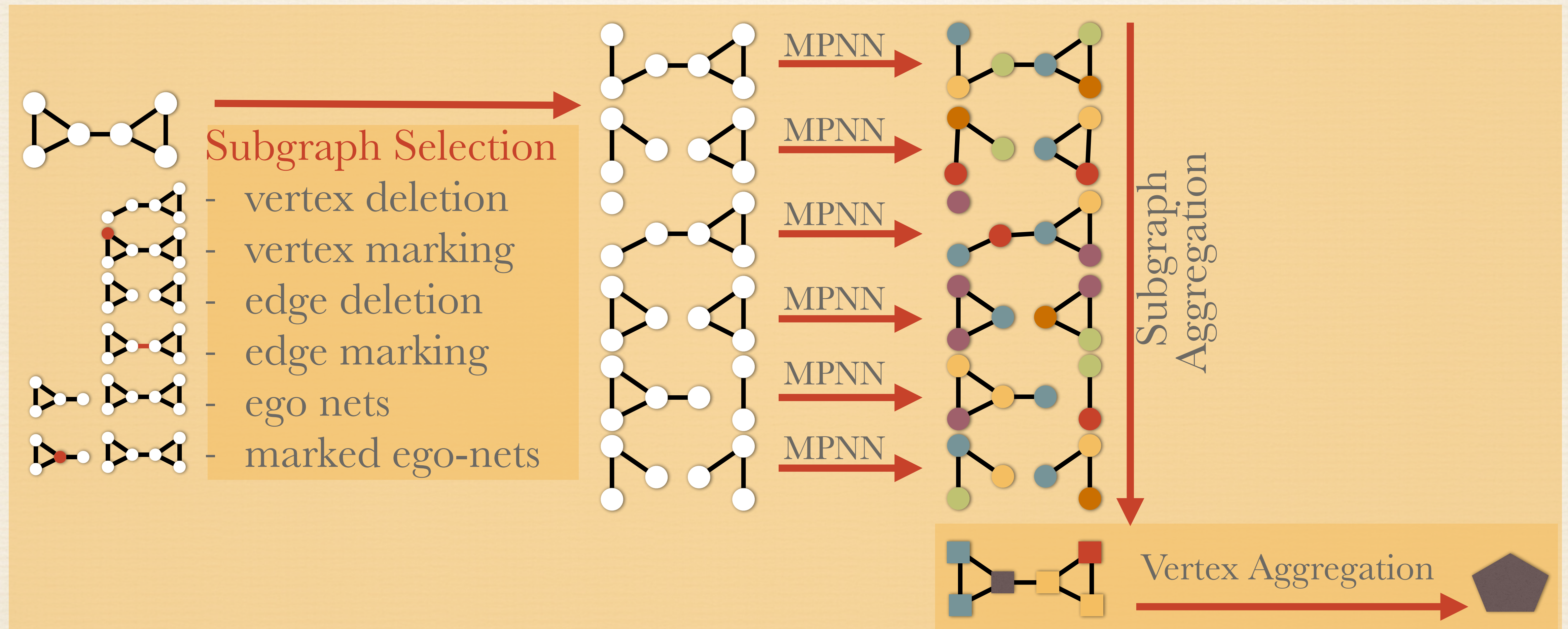
General idea

- ❖ Colour refinement equivalent graphs may contain **colour refinement inequivalent subgraphs**.

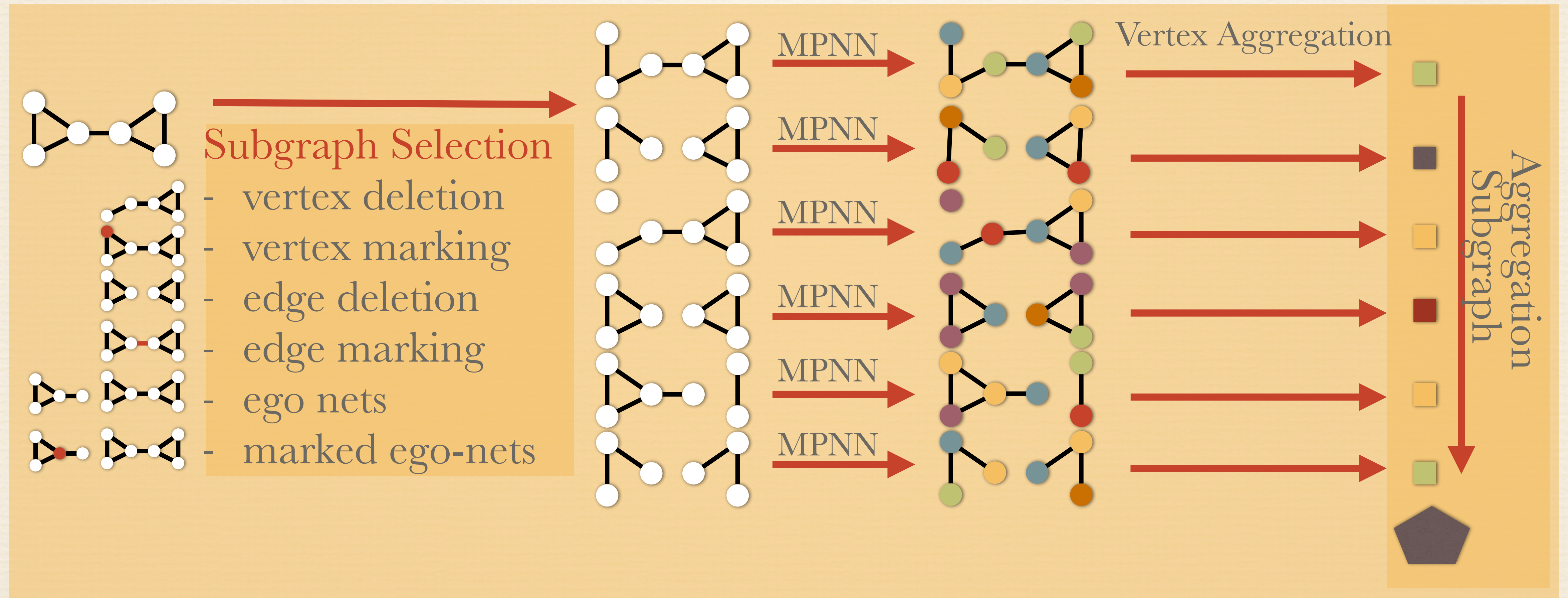


- ❖ View graphs as **a collection of subgraphs** then run MPNN

Subgraph \mapsto Vertex Aggregation



Vertex \rightarrow Subgraph Aggregation



The subgraph GNN “wave”

Vertex \mapsto Subgraph Aggregation

k-OSAN^T

Reconstruction GNN

DS-GNN

DSS-GNN

Subgraph \mapsto Vertex Aggregation

GNN-AK

k-OSAN

NGNN

ID-GNN

DropoutGNN

All provably **more expressive** than MPNNs*

Bevilacqua et al.: *Equivariant subgraph aggregation network* (2022)

Cotta et al.: *Reconstruction for powerful graph representations* (2021)

Bevilacqua et al.: *Understanding and extending subgraph GNNs by rethinking their symmetries* (2022)

Huang et al.: *Boosting the cycle counting power of graph neural networks with I2-GNNs* (2022)

Papp et al.: *DropGNN: Random dropouts increase the expressiveness of graph neural networks.* (2021)

Qian et al.: *Ordered subgraph aggregation networks.* (2022)

You et al.: *Identity-aware graph neural networks.* (2021)

Zhang and P. Li. *Nested graph neural networks* (2021)

Zhao et al.: *From stars to subgraphs: Uplifting any GNN with local structure awareness* (2022)

Selection policies

DS-GNN

- vertex deletion
- edge deletion
- ego nets
- marked ego-nets

ID-GNNs

- marked ego-nets

GNNs-AK

- ego-nets

k-OSAN

- size k subgraph marking

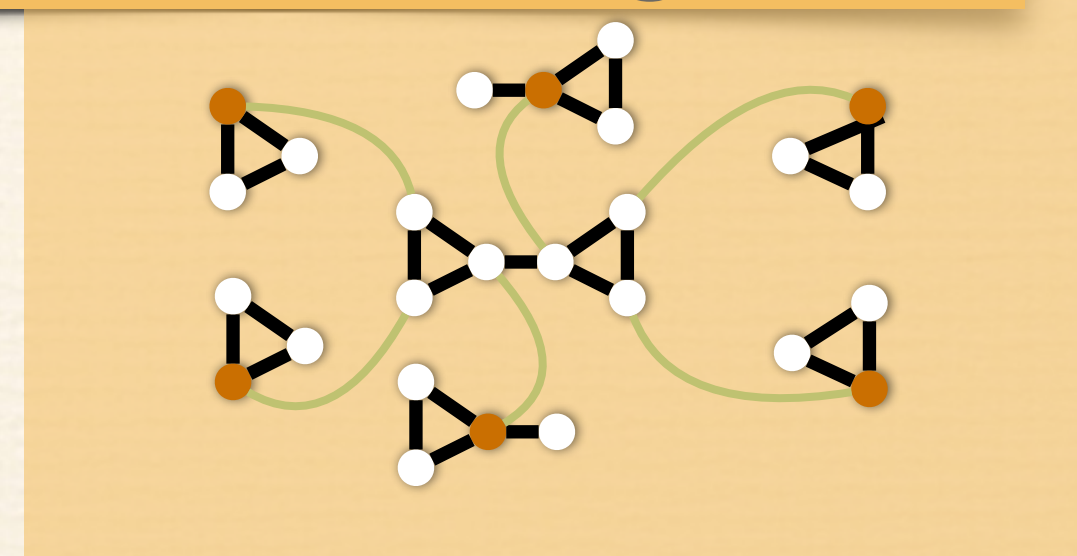
Rec-GNN

- k-vertex deletion

NGNN

- ego-nets

Popular/effective: ego-nets



k-OSAN

Theorem (Qian et al. 2022)

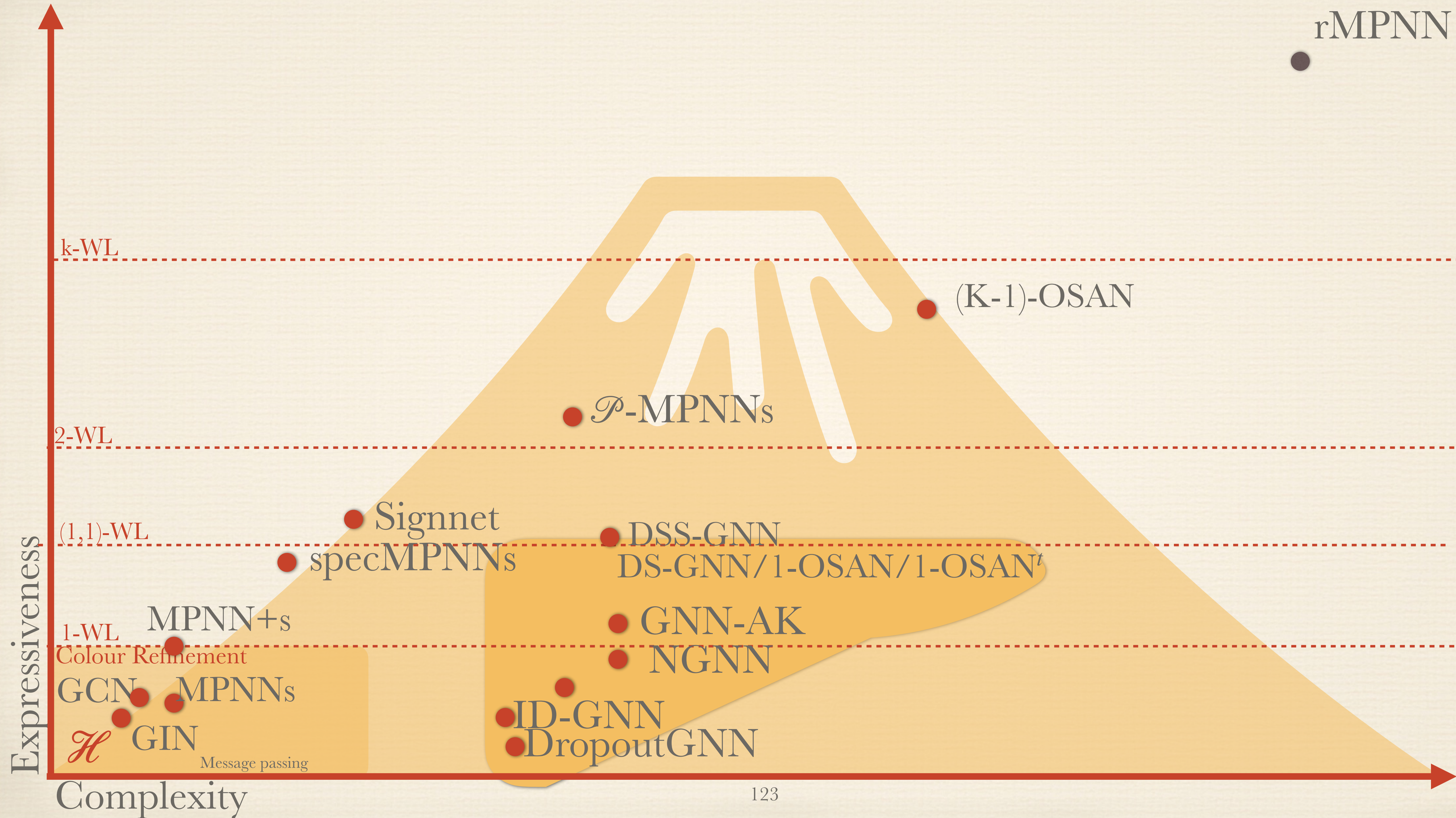
- k-OSANs and k-OSANs^t encompass *almost all* subgraph methods with selection policy involving k vertices.
- *Strictly bounded* in expressive power by $(k+1)$ -WL
- *Incomparable* to k-WL.

k=2

- ❖ if 2-WL cannot distinguish graphs, then **neither can 1-OSANs**
- ❖ 2-WL can distinguish **more graphs** than 1-OSANs
- ❖ There exists graphs than can be distinguished by 1-OSANs but not by MPNNs, and vice versa, there exists graphs that can be distinguished by MPNNs but not by 1-OSANs

Subgraph GNNs

- ❖ Can always ensure to be strictly **more expressive than MPNNs** by including **original graph in batch**.
- ❖ **Tractability** only when **easy subgraph policies** are used, i.e., leading to a small number (linear) of subgraphs.
- ❖ Seems a **good balance** between **complexity** and **expressiveness**





K-dimensional Weisfeiler-Leman

Boosting expressive power by higher-order message-passing

More powerful heuristic

Apply **heuristic** on G and H : If Heuristic say “no” then $G \not\cong H$, otherwise we do not know.

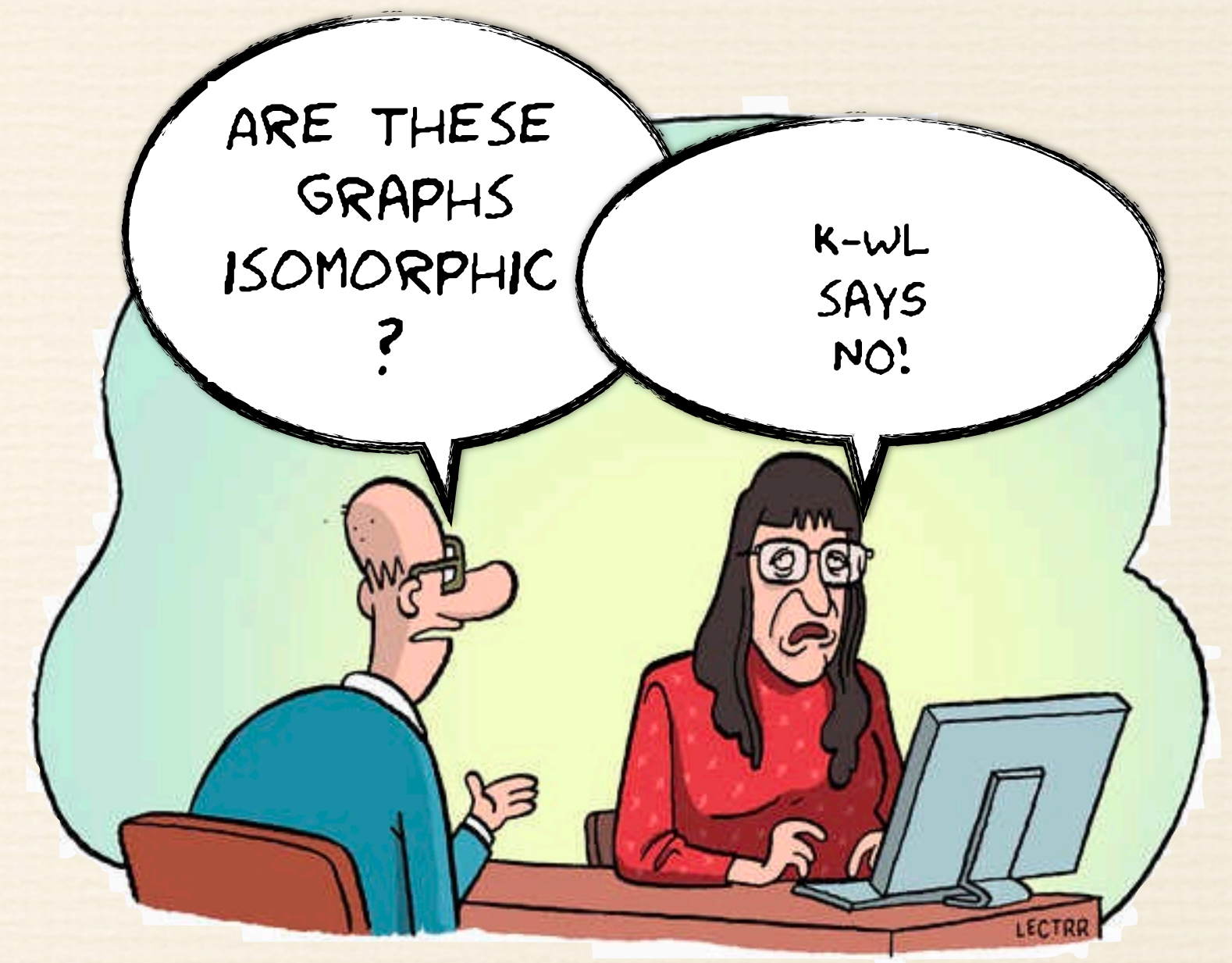
$G \cong H?$

Colour refinement \rightarrow No $\rightarrow G \not\cong H$

\downarrow
1-WL \rightarrow No $\rightarrow G \not\cong H$

\downarrow
2-WL \rightarrow No $\rightarrow G \not\cong H$

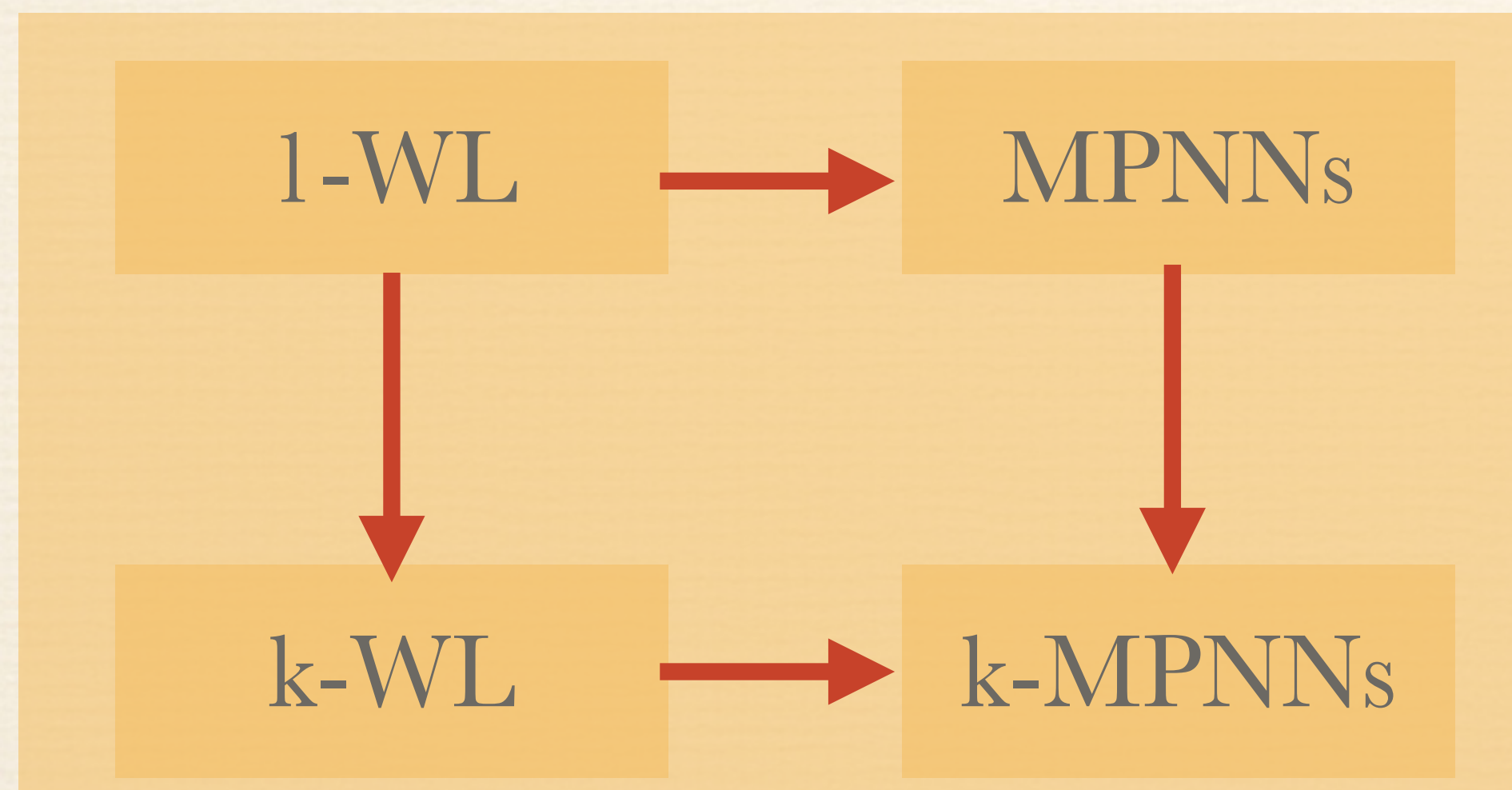
\downarrow
 \dots



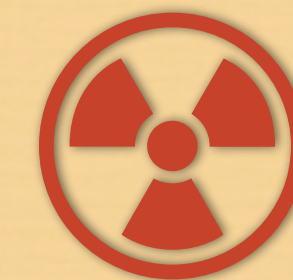
Idea: higher-order GNNs

Theorem (Dell et al. 2018, ...)

$\text{hom}(T, G) = \text{hom}(T, H)$ for all graphs T of tree width k
if and only if
 k -WL cannot tell apart G from H



k-MPNNs will detect more graph information than MPNNs



k-Folklore GNNs (k-FGNNs)

$$\xi^{(t)}(G, v_1, \dots, v_k) := \text{MLP}_1^{(t)}\left(\sum_{u \in V_G} \prod_{i=1}^k \text{MLP}_2^{(t)}(\xi^{(t-1)}(G, v_1, \dots, v_{i-1}, u, v_{i+1}, \dots, v_k))\right)$$

k -vertex embedding

Global aggregation

Uses multiplication

Expressive power?

Theorem (Maron et al. 2019), Azizian and Lelarge 2021)

$$\rho(k\text{-FGNN}) = \rho(k\text{-WL})$$

k-GNNs

A simpler architecture:

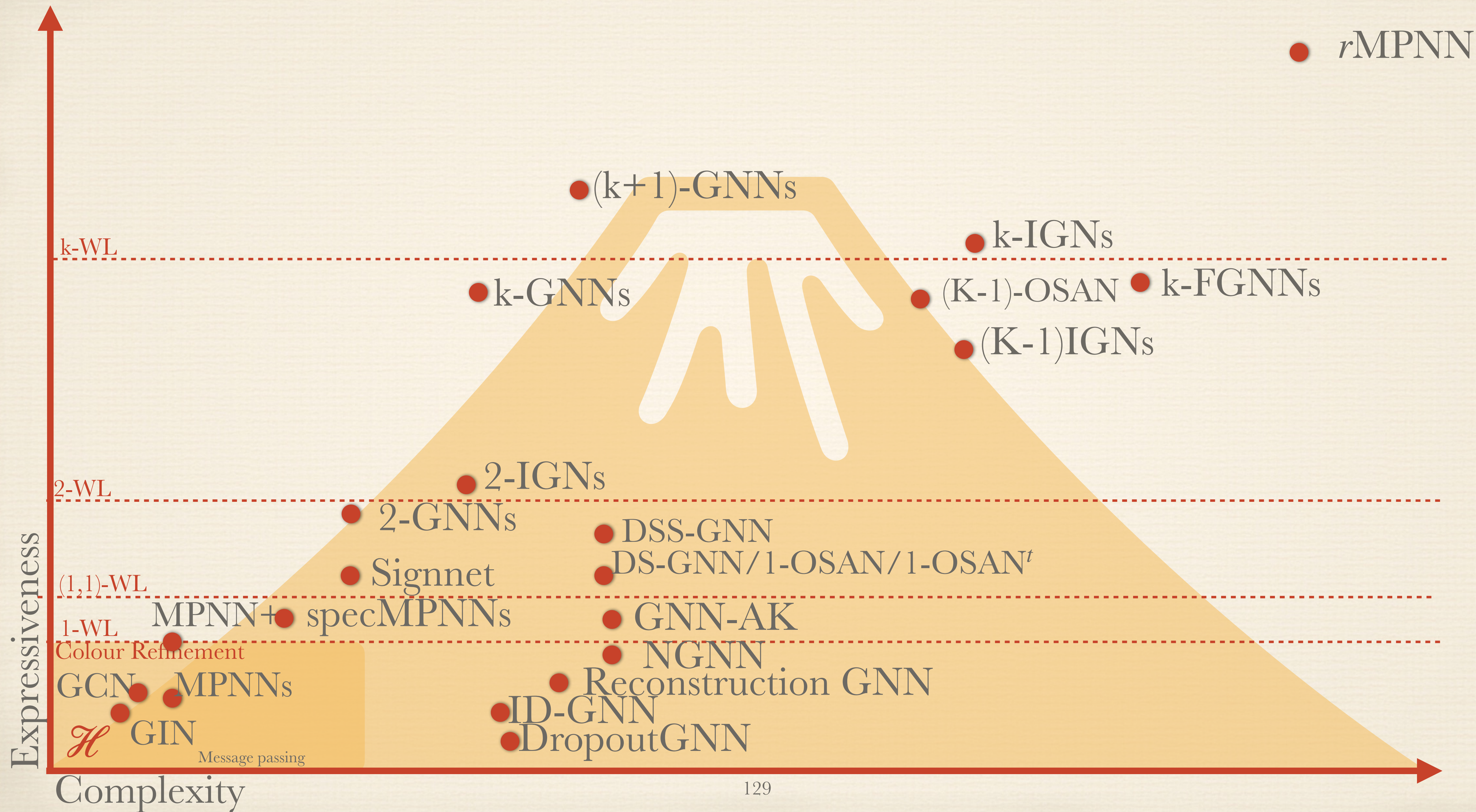
$$\xi^{(t)}(G, v_1, \dots, v_k) := \sigma \left(\xi^{(t-1)}(G, v_1, \dots, v_k) \mathbf{W}_1^{(t)} + \left(\sum_{i=1}^k \sum_{u \in V_G} \xi^{(t)}(G, v_1, \dots, v_{i-1}, u, v_{i+1}, \dots, v_k) \right) \mathbf{W}_2^{(t)} \right)$$

Global aggregation

Expressive power?

Theorem (Morris et al. 2019)

$$\rho(k\text{-GNN}) = \rho(k\text{-WL})$$



Questions?



Conclusions

And look ahead

What to use?

Subgraph

- ❖ Small graphs
- ❖ Good compromise in general

Feature Augmentation

- ❖ Large training datasets
- ❖ Invariance not important
- ❖ Preprocessing ok

Higher-order

- ❖ Graphs are small
- ❖ Efficiency not essential
- ❖ Expressivity guarantee needed

Road ahead

Expressiveness

- ❖ A lot of recent progress
- ❖ WL hierarchy needs better reconciliation with practice
- ❖ Hom count characterisations
- ❖ Relational

Connection with Learning??

- ❖ Optimisation and training unexplored

- ❖ Generalisation properties
- ❖ Sample efficiency?

Conclusion

- ❖ Study of expressive is beautiful area of research for ML researchers
- ❖ Combines theory and practice in an elegant way
- ❖ Many unresolved questions ...