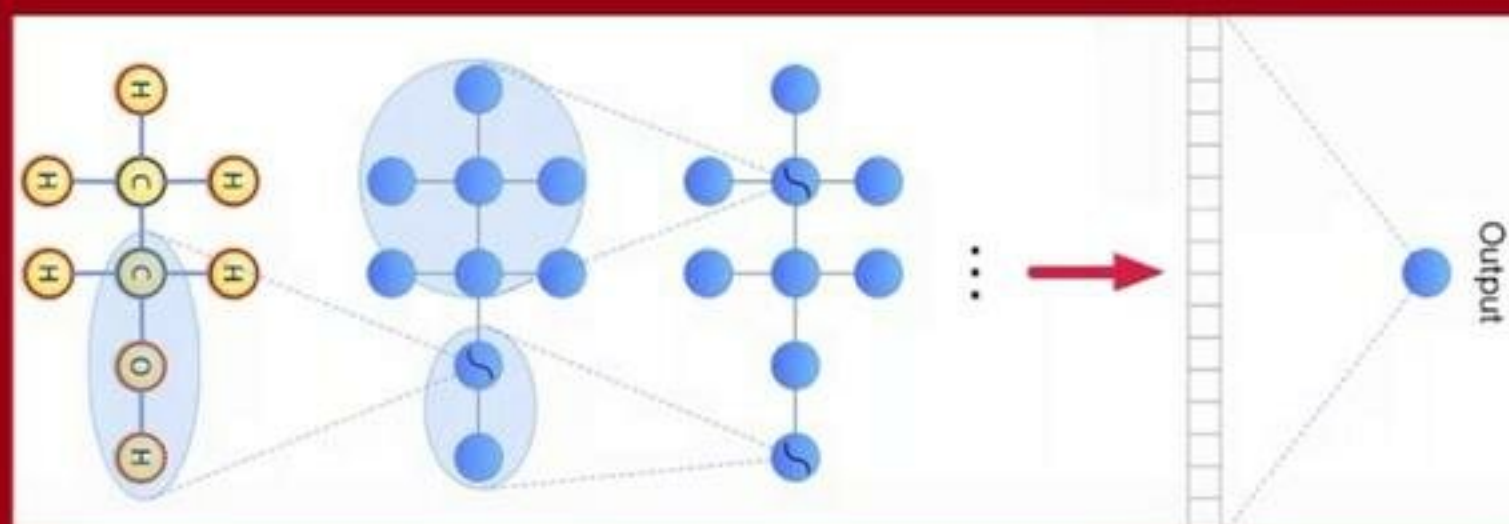


Learning effective representations with Graph Neural Networks



Nicolò Navarin

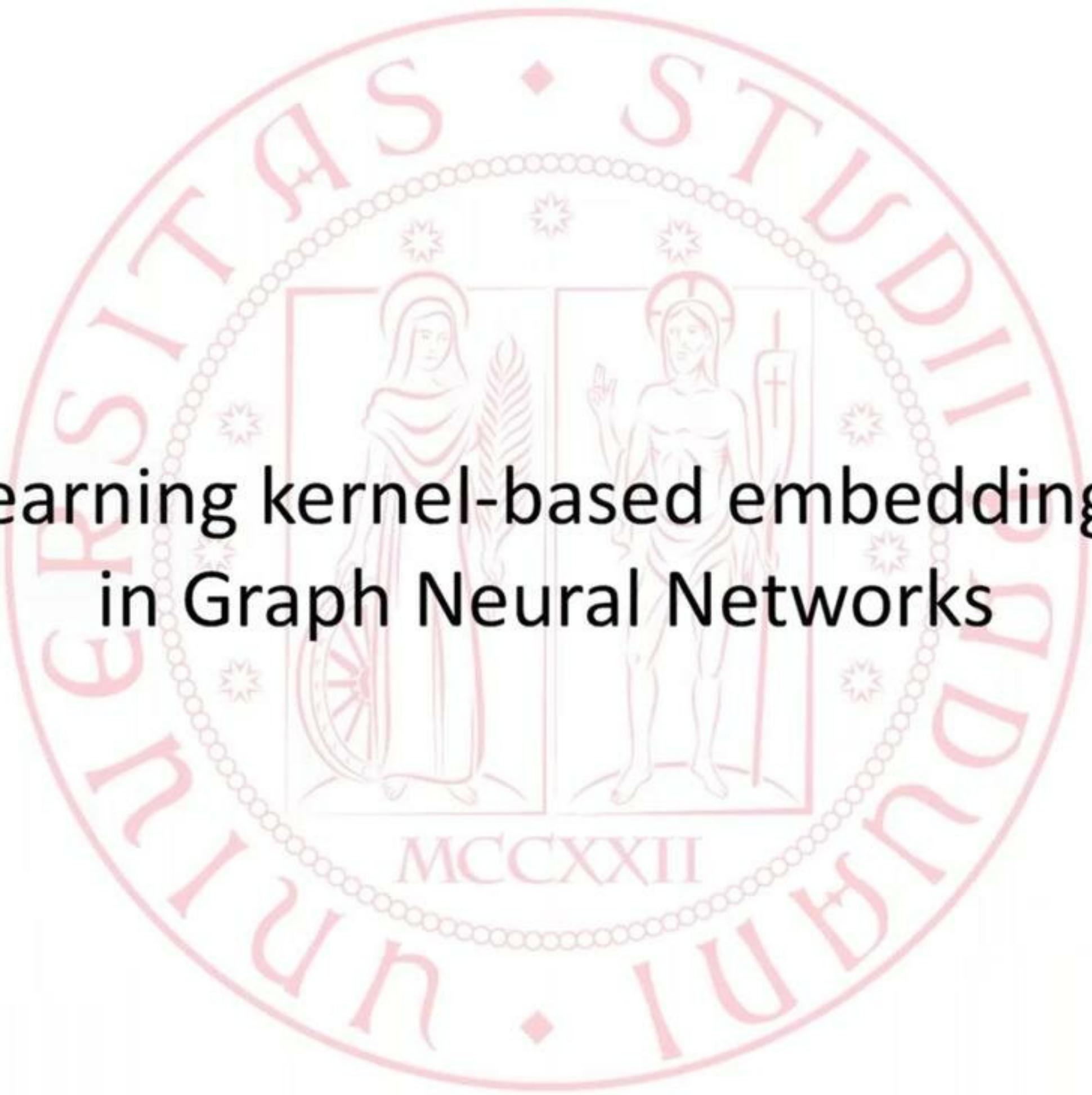
*Assistant Professor
Department of Mathematics
University of Padua*

Talk @ TU Wien
24/11/2020



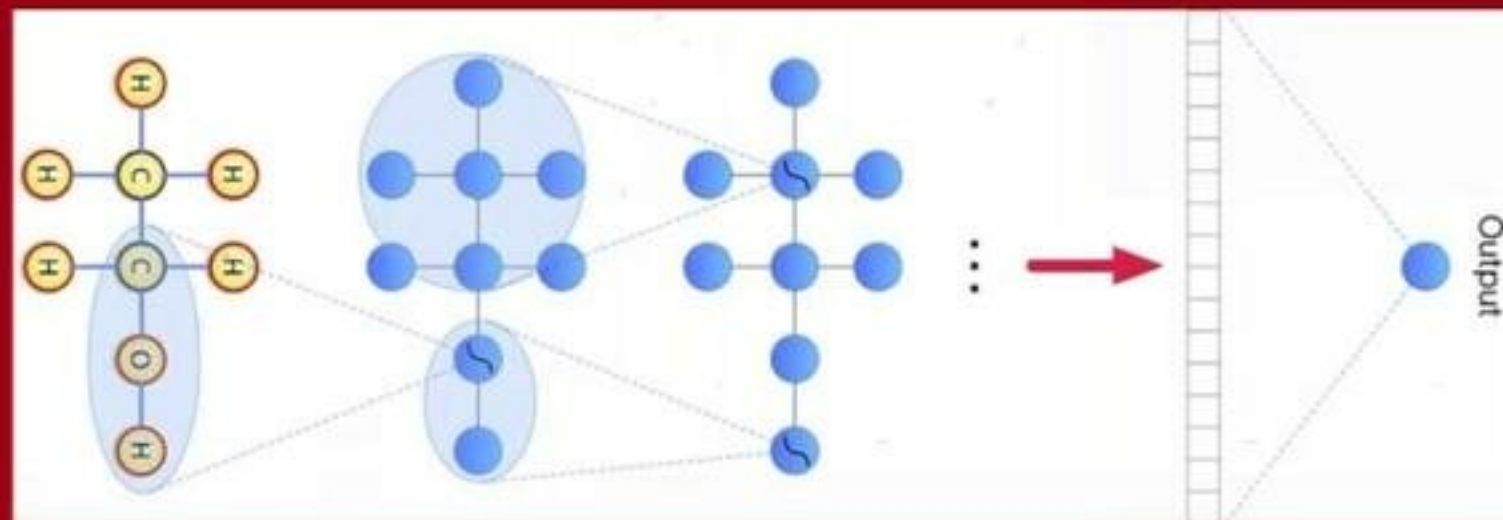
UNIVERSITÀ
DEGLI STUDI
DI PADOVA

DIPARTIMENTO
MATEMATICA

The background of the slide features a large, faint, red circular seal of the University of Vienna. The seal contains a central illustration of two figures, likely saints or scholars, standing under a gothic arch. The text 'UNIVERSITAS STUDII VIENNAE' is inscribed around the perimeter, and 'MCCXXII' is at the bottom.

Learning kernel-based embeddings in Graph Neural Networks

Learning effective representations with Graph Neural Networks



Nicolò Navarin

*Assistant Professor
Department of Mathematics
University of Padua*

Talk @ TU Wien
24/11/2020



UNIVERSITÀ
DEGLI STUDI
DI PADOVA

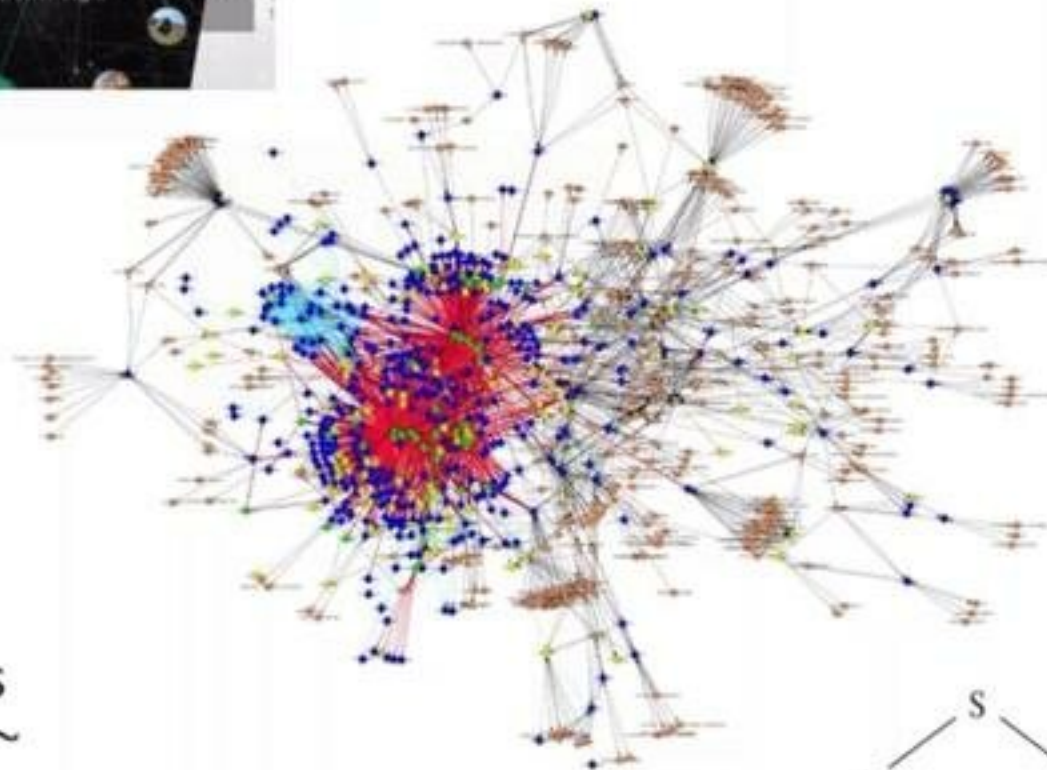
 DIPARTIMENTO
MATEMATICA

Graphs are everywhere..

Knowledge graphs



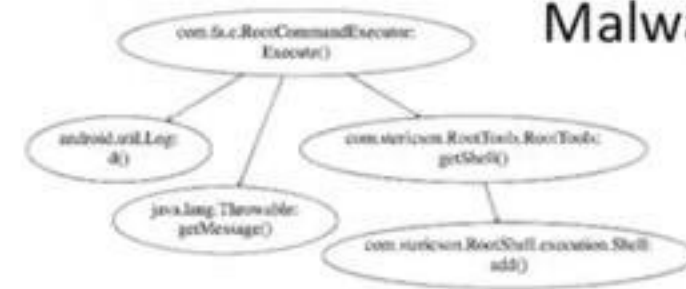
Gene networks



- In several settings it is natural to represent data as graphs.
- We may want to **predict some property:**

- **over Graphs or**
- **over Graph Nodes**

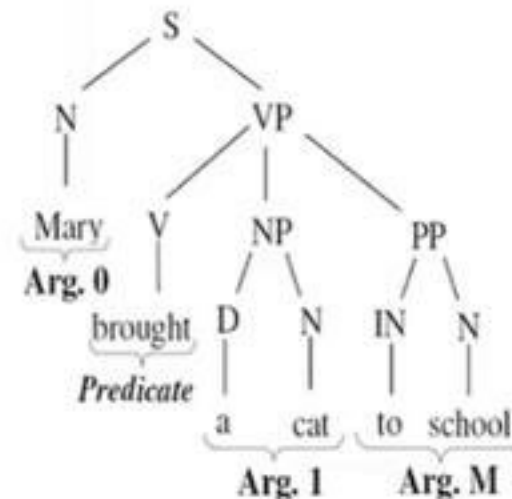
Malware API calls



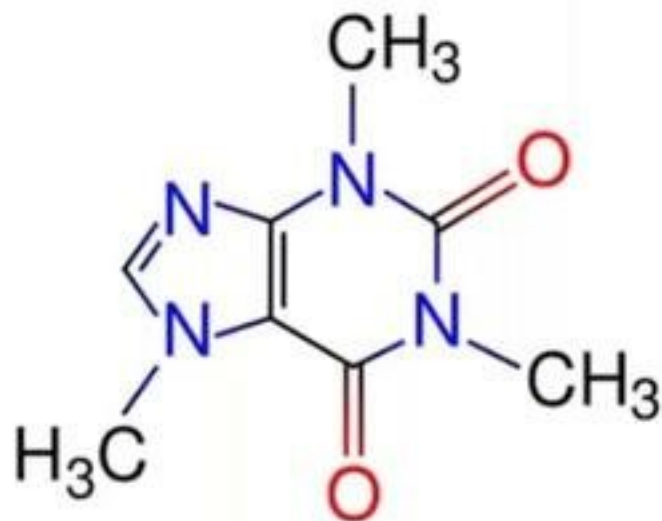
Images



Parse trees

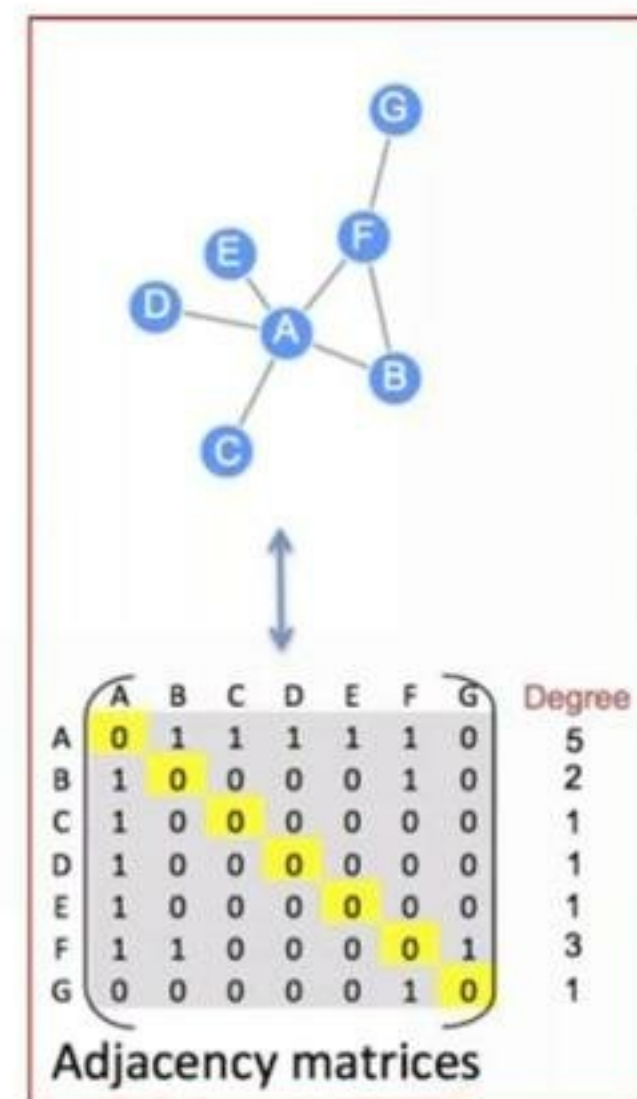
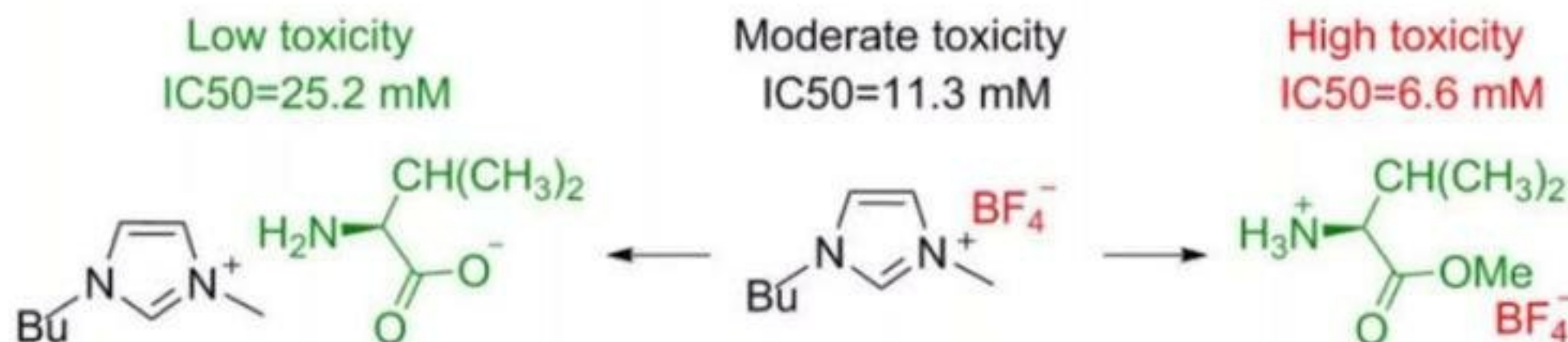


Chemical compounds



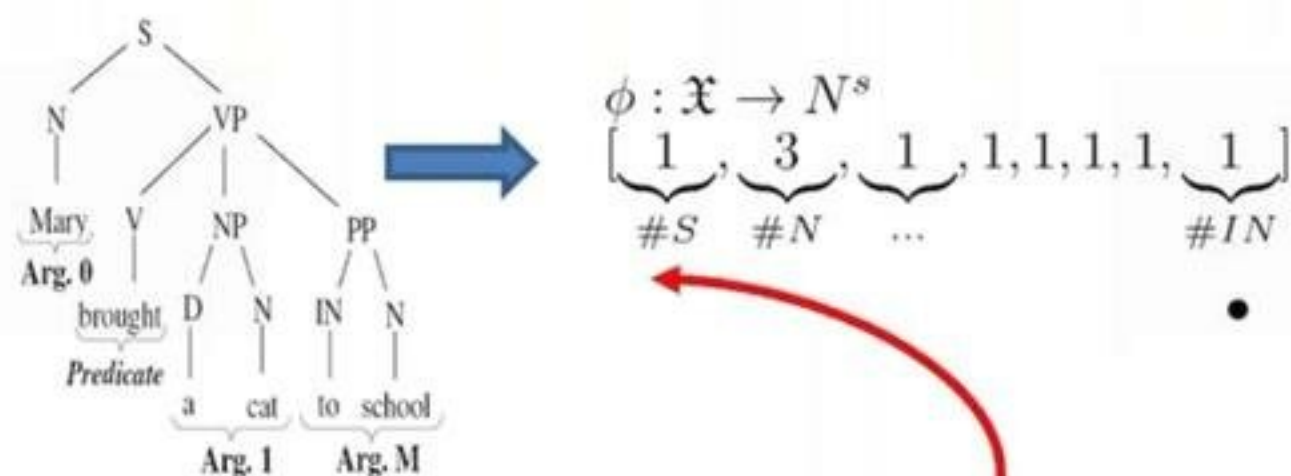
Classification/Regression on Graphs

In this talk, classification/regression on graphs



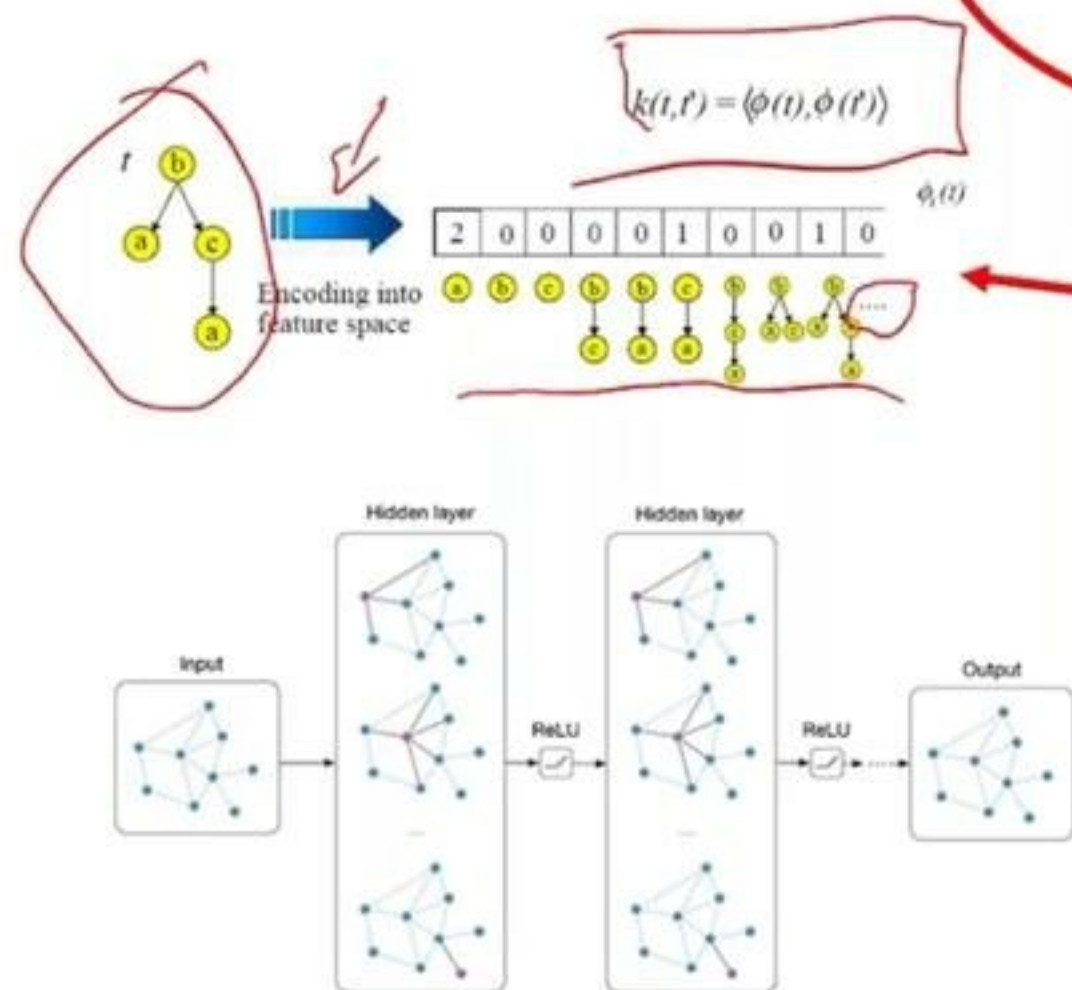
- Dataset composed of N pairs $\{(G_i, y_i), 1 \leq i \leq N\}$
 - Each graph:
 - n_i vertices
 - (possibly) discrete label associated to each node: $l(v)$
 - d vectorial **attributes** associated to each node: $a(v)$ or $X \in \mathbb{R}^{n_i \times d}$
- Given an unseen graph G , the task is to predict the correct target y

Machine learning on Graphs



- It is difficult to cast “mainstream” learning algorithms to work on graphs.

- 3 possible ways:

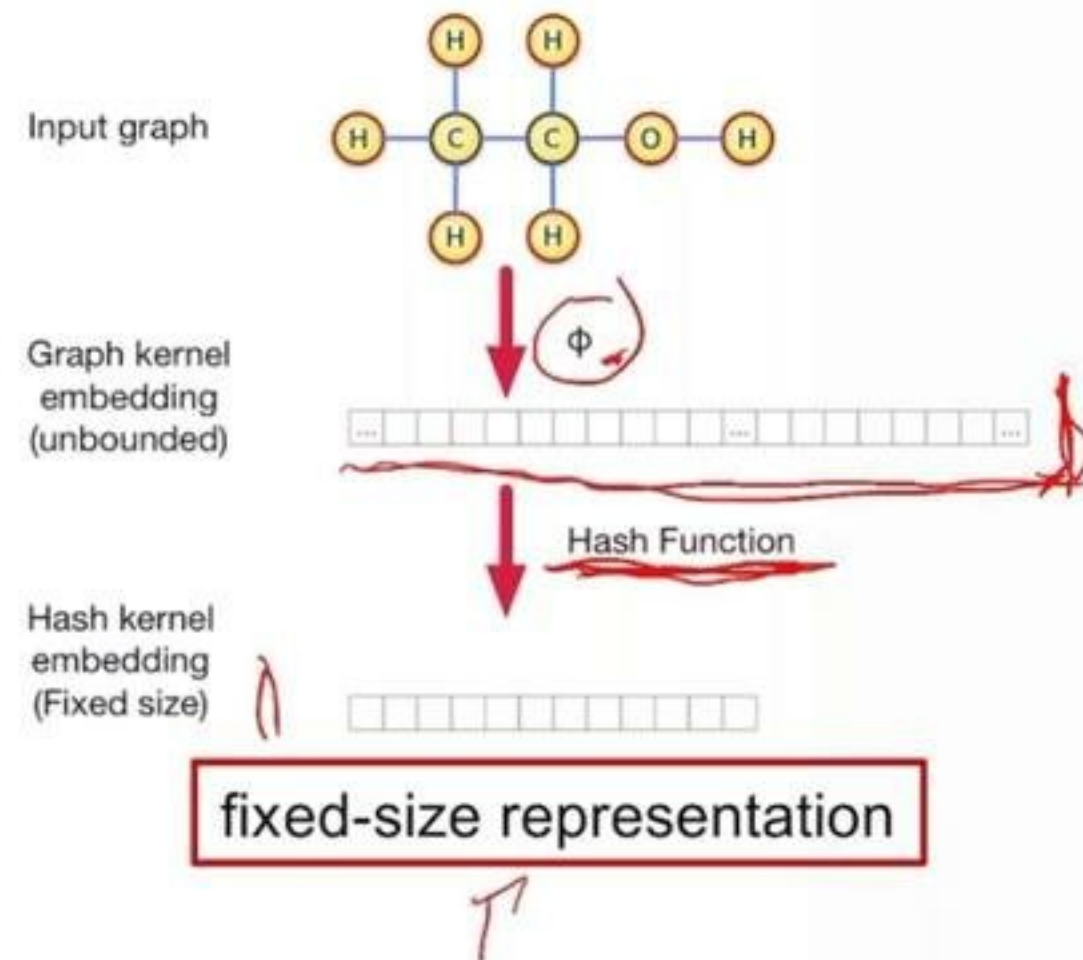
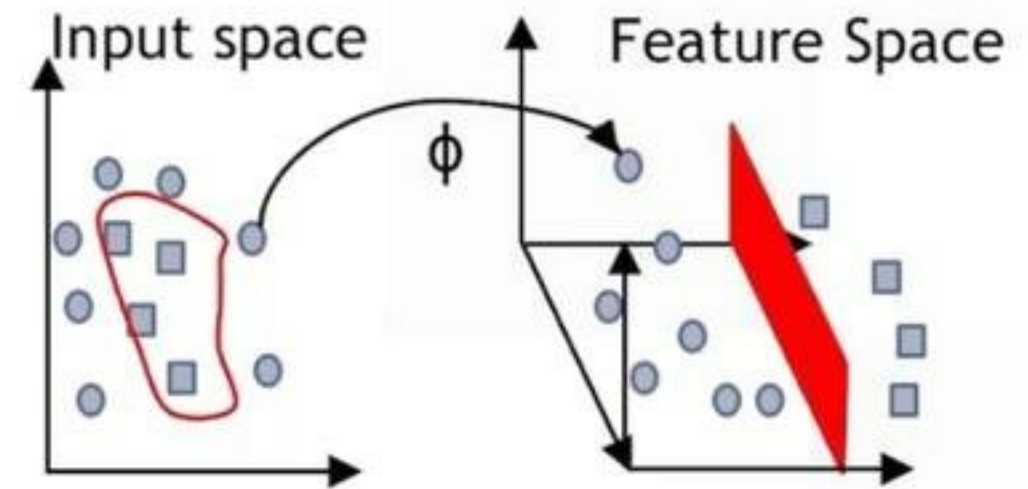


Graph Kernels

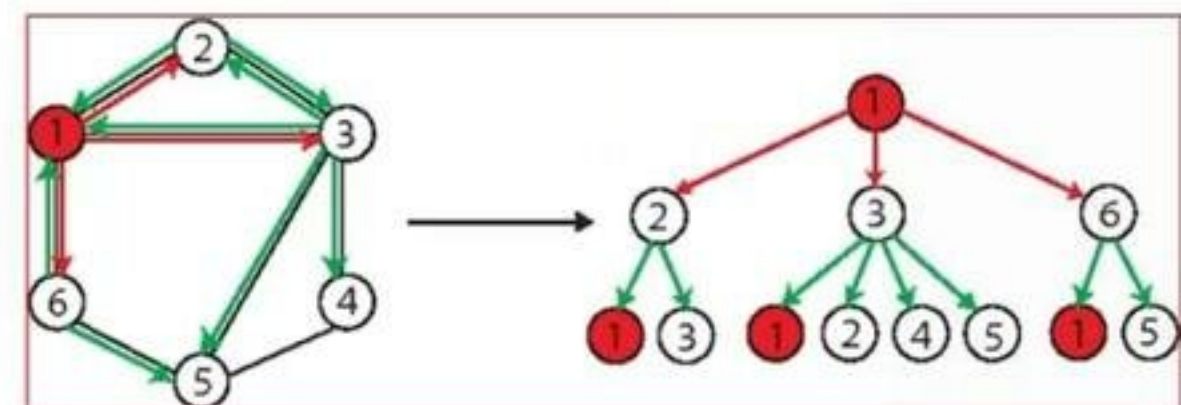
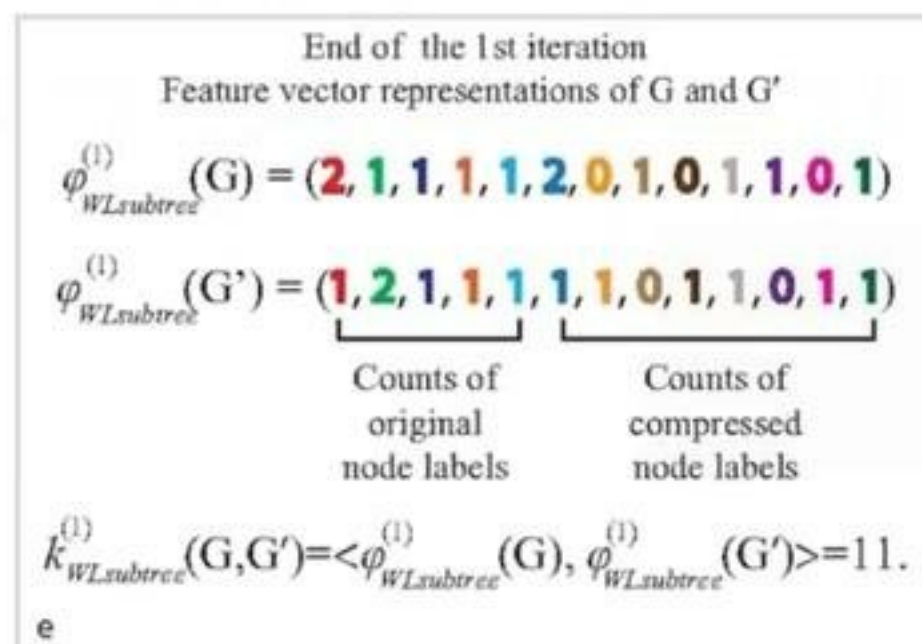
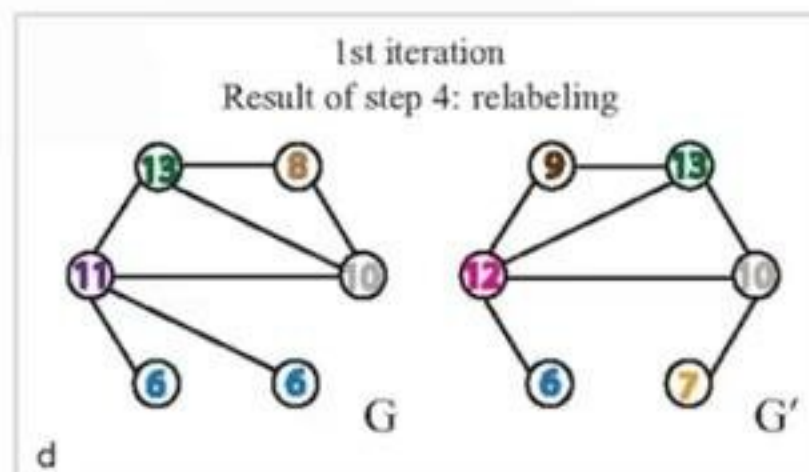
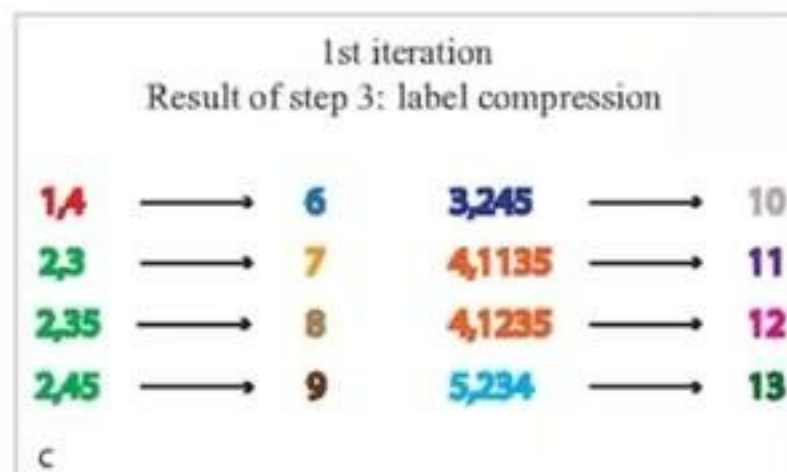
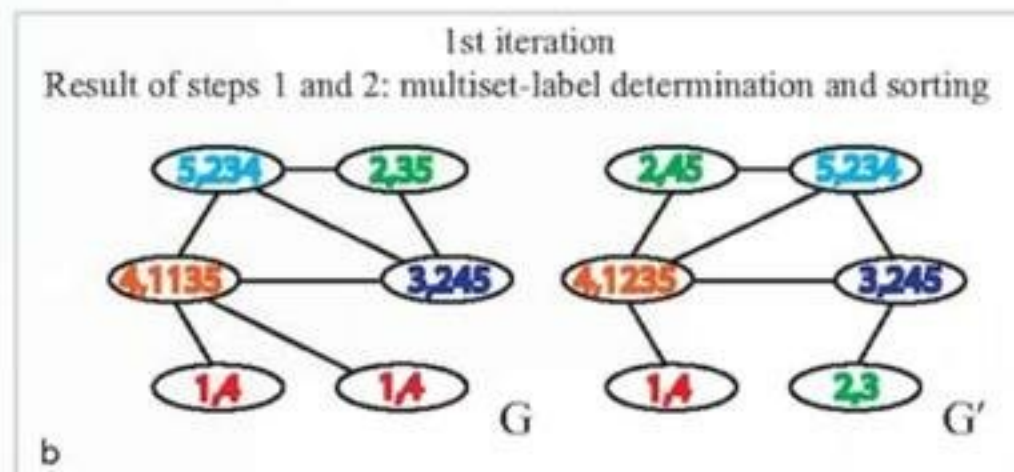
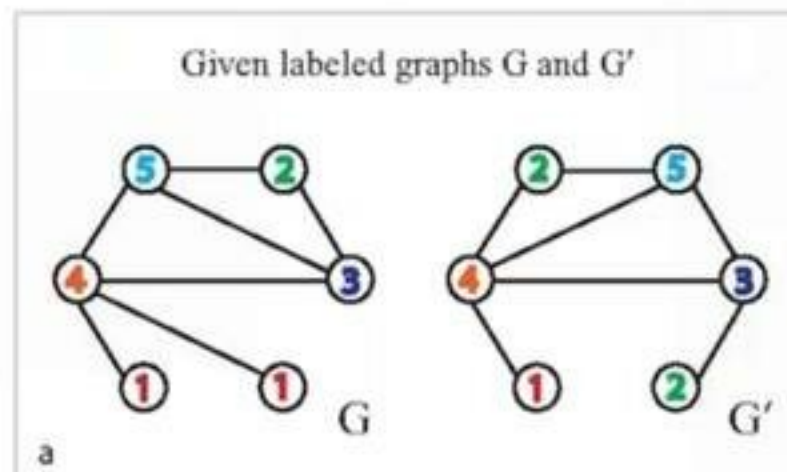


$$k(X, Y) = \phi(X)^T \phi(Y)$$

- Implicitly define a mapping from input space (the space of graphs) to feature space.
- Several proposals in literature:
 - Global measures: Random walks, Shortest paths, Graphlets
 - Local measures:
 - Weisfeiler-Lehman subtree kernel
 - Ordered Decomposition DAGs kernel
 - Explicit feature space, hashing to obtain a fixed-size representation

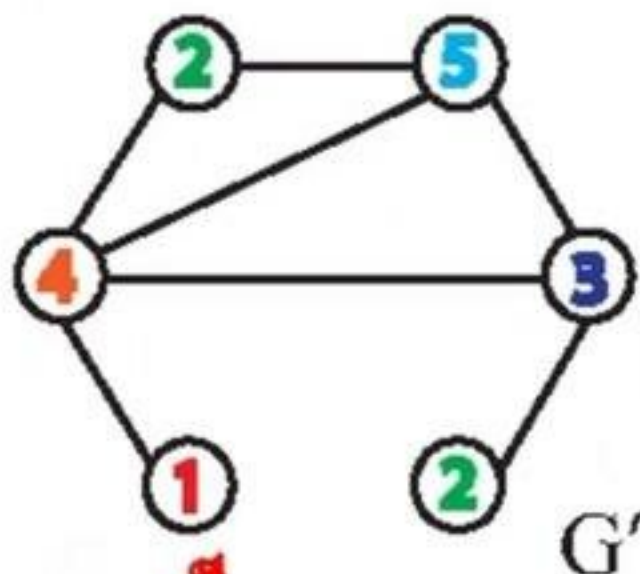
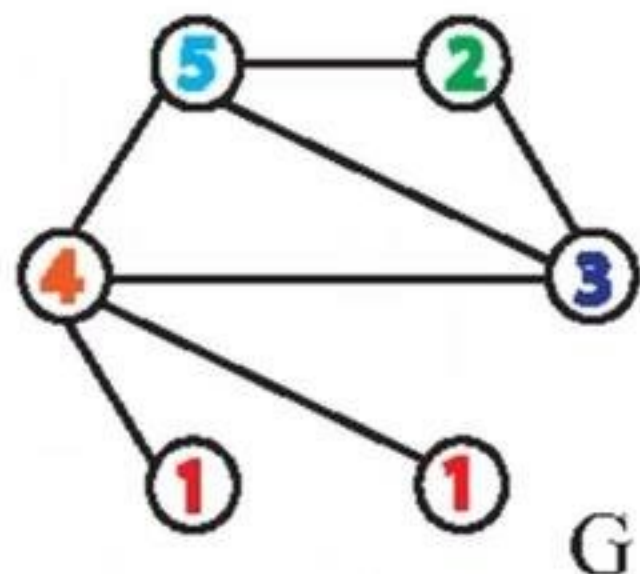


Weisfeiler-Lehman Subtree kernel

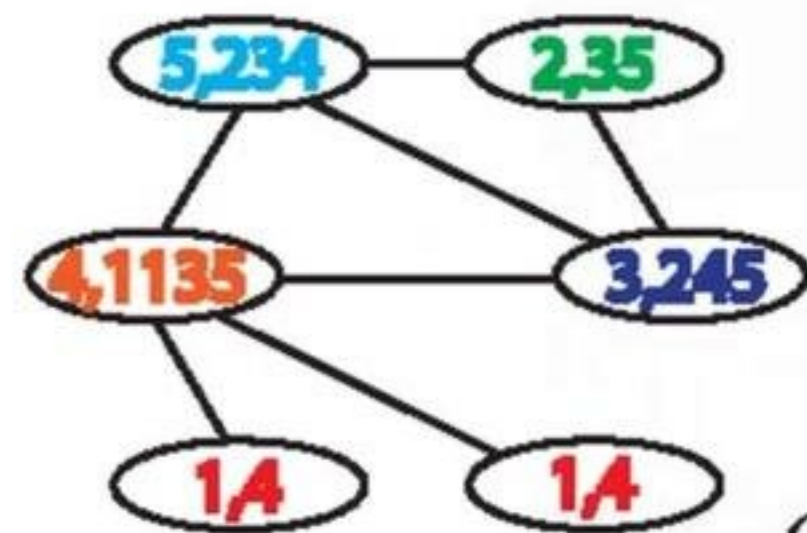


Weisfeiler-Lehman Subtree

Given labeled graphs G and G'



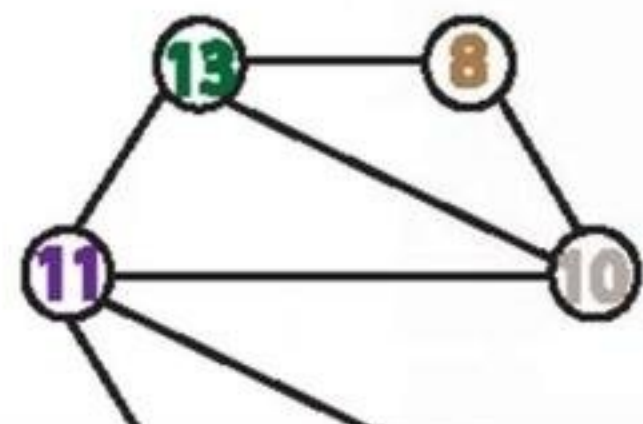
1st iteration
Result of steps 1 and 2: multiset



1st iteration
Result of step 3: label compression

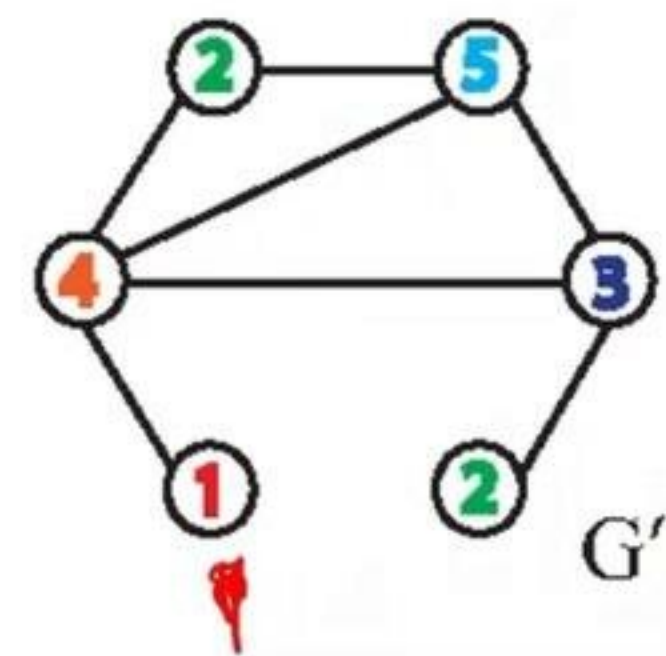
1,4	→	6	3,245	→	10
2,3	→	7	4,1135	→	11
2,35	→	8	4,1235	→	12

1st iteration
Result of step

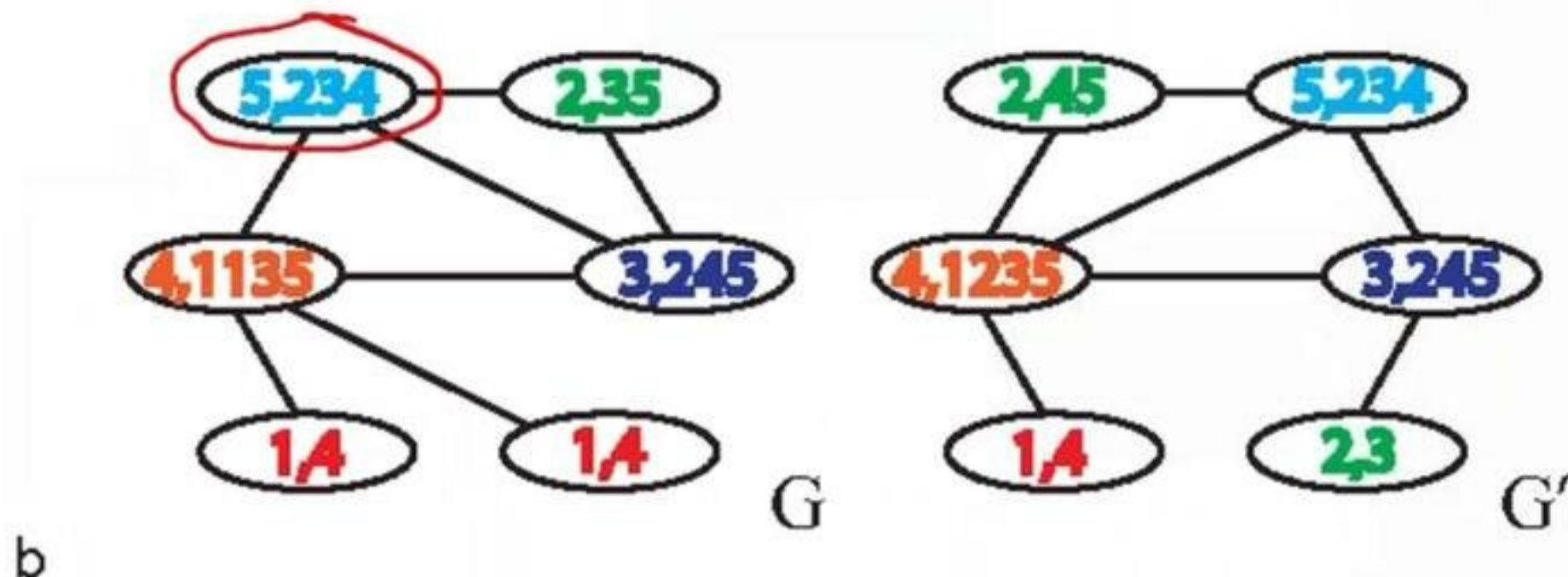


er-Lehman Subtree kernel

Graphs G and G'



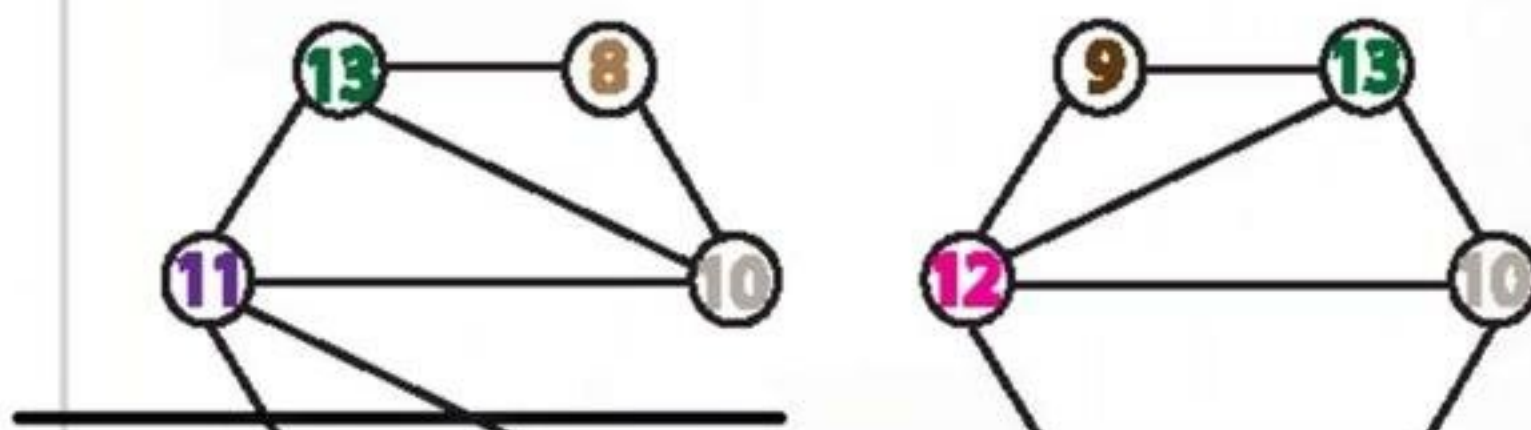
1st iteration
Result of steps 1 and 2: multiset-label determination and sorting



Label compression



1st iteration
Result of step 4: relabeling



Brief History of GNNs



- The idea of Neural Networks for structured data dates back to '97 [Sperduti & Starita, 1997]
- In the '00s, two proposals:
 - Graph Neural Network Model [Scarselli, Gori et al., 2009]
 - Recurrent model, contraction mapping
 - Neural Networks for Graphs [Micheli, 2009] ~~scribbled~~
 - Convolutional model, layer-wise training
- Recently, many works proposing slight modifications, e.g. :
 - [Tarlow et al., 2016] Extends [Scarselli, Gori et al., 2009] ~~scribbled~~
 - no contraction mapping, GRUs
 - [Kipf & Welling, 2017] proposes an approach similar to [Micheli, 2009] for node classification
 - end-to-end

Sperduti & Starita (1997). *Supervised neural networks for the classification of structures*. IEEE TNNs.

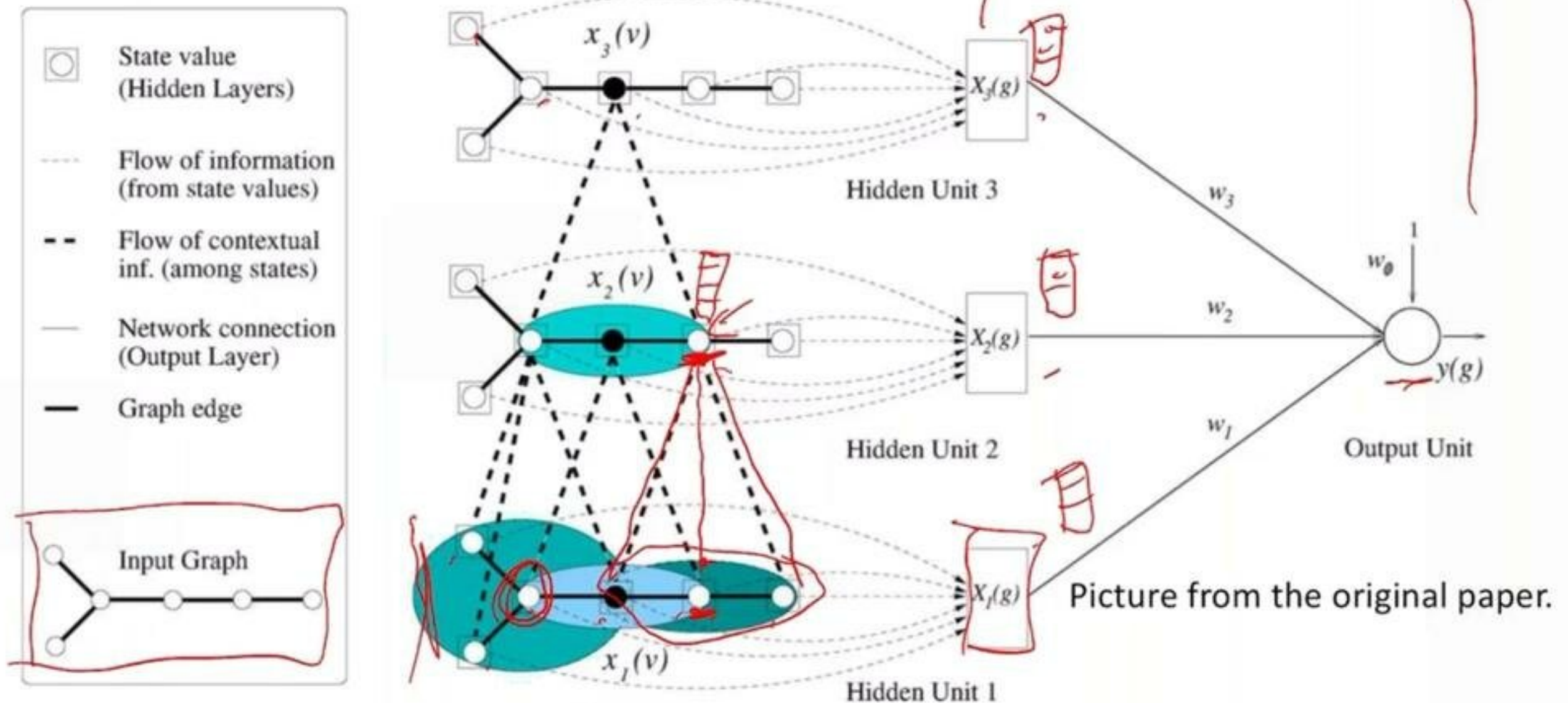
Scarselli, Gori et al. (2009). *The Graph Neural Network Model*. IEEE TNNs.

Micheli (2009). *Neural network for graphs: A contextual constructive approach*. IEEE TNNs.

Kipf & Welling(2017). *Semi-Supervised Classification with Graph Convolutional Networks*. In ICLR.

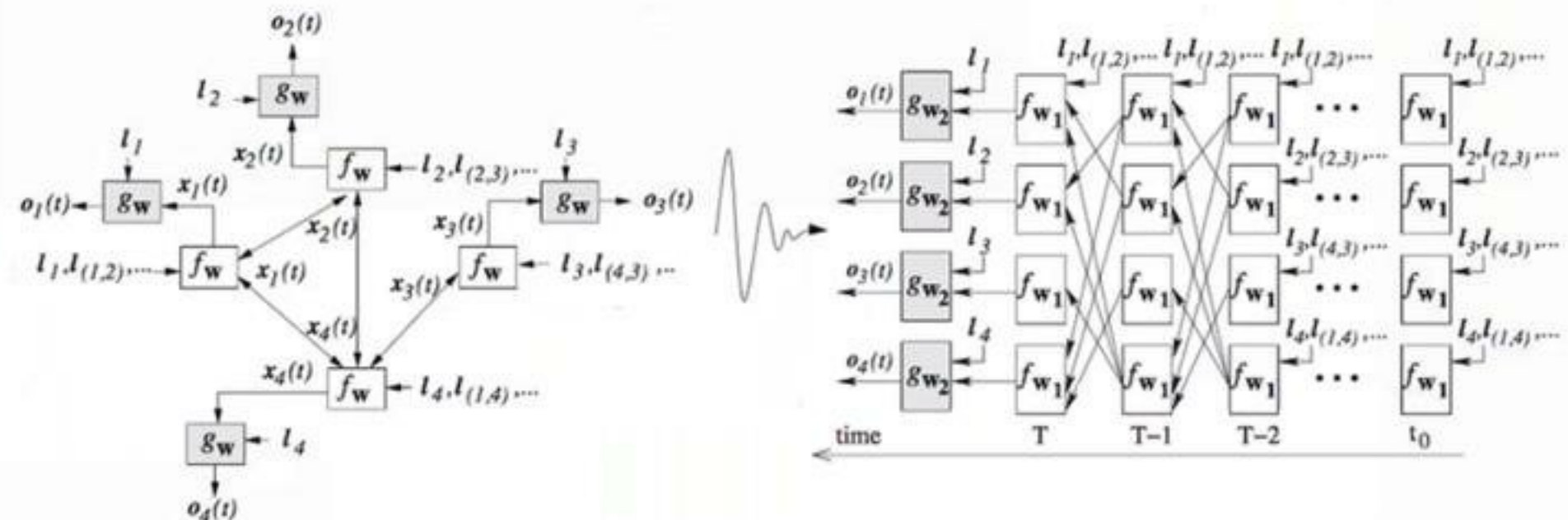
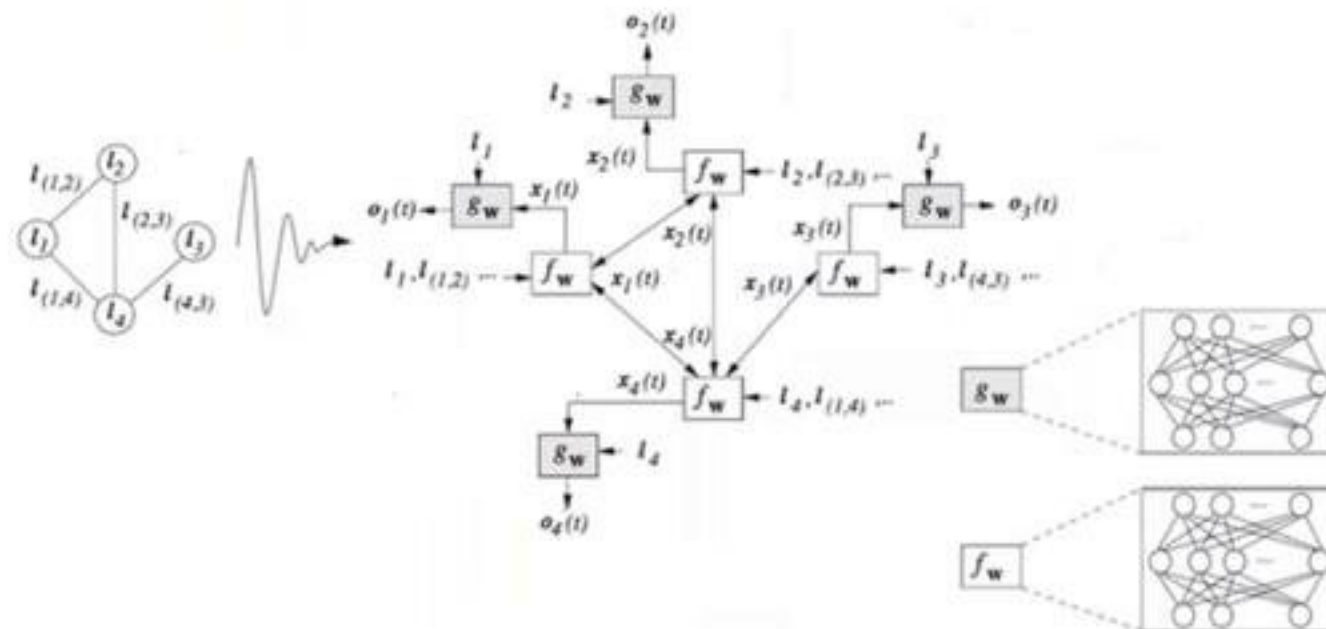
Tarlow et al. (2016). *Gated Graph Sequence Neural Networks*. In ICLR.

NN4G by Micheli



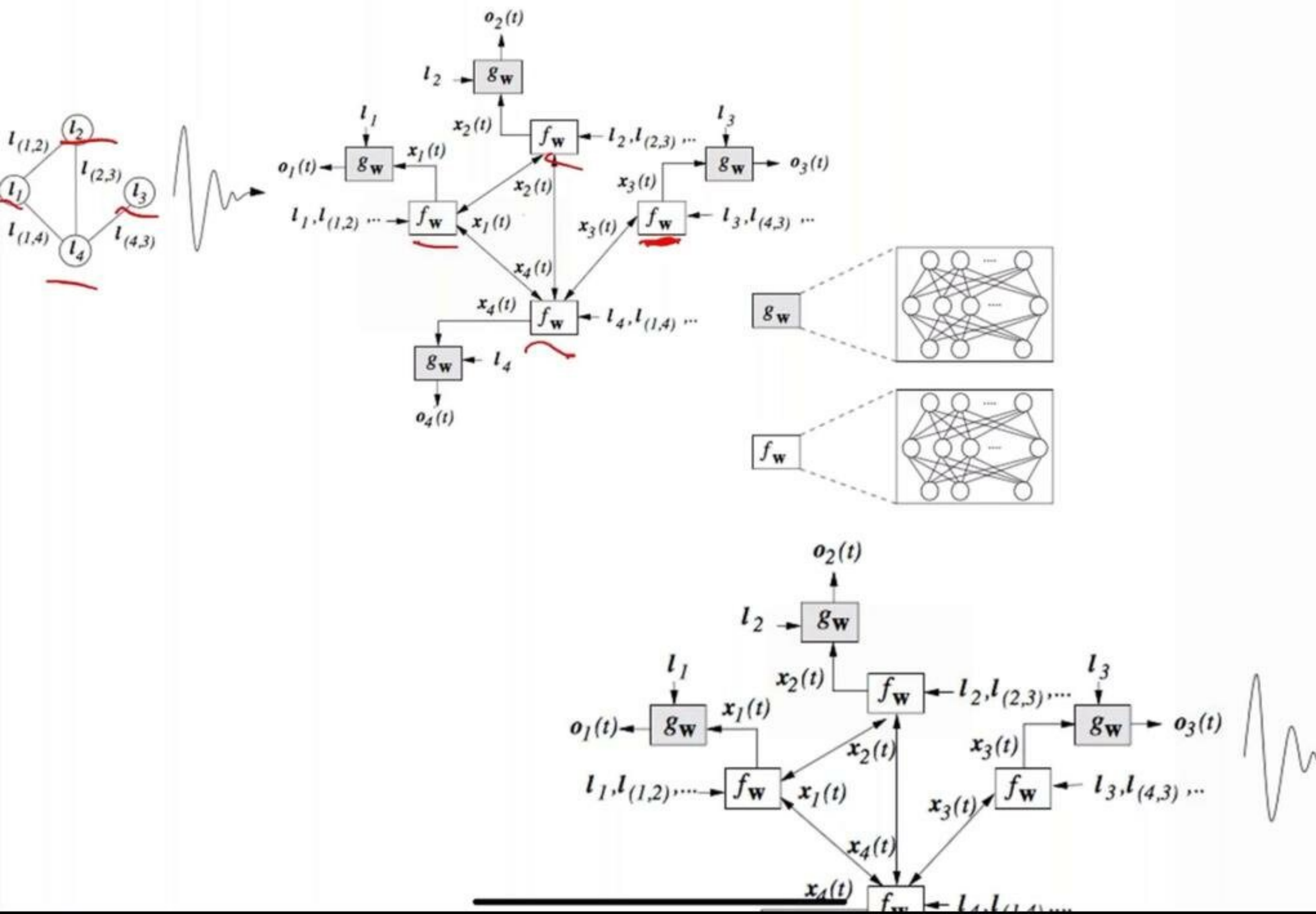
- Each convolution takes as input the representation of all previous layers
- Trained layer-wise (cascade correlation)
- Readout: a representation per-graph per-layer is computed using the sum

GNN by Scarselli, Gori et al.

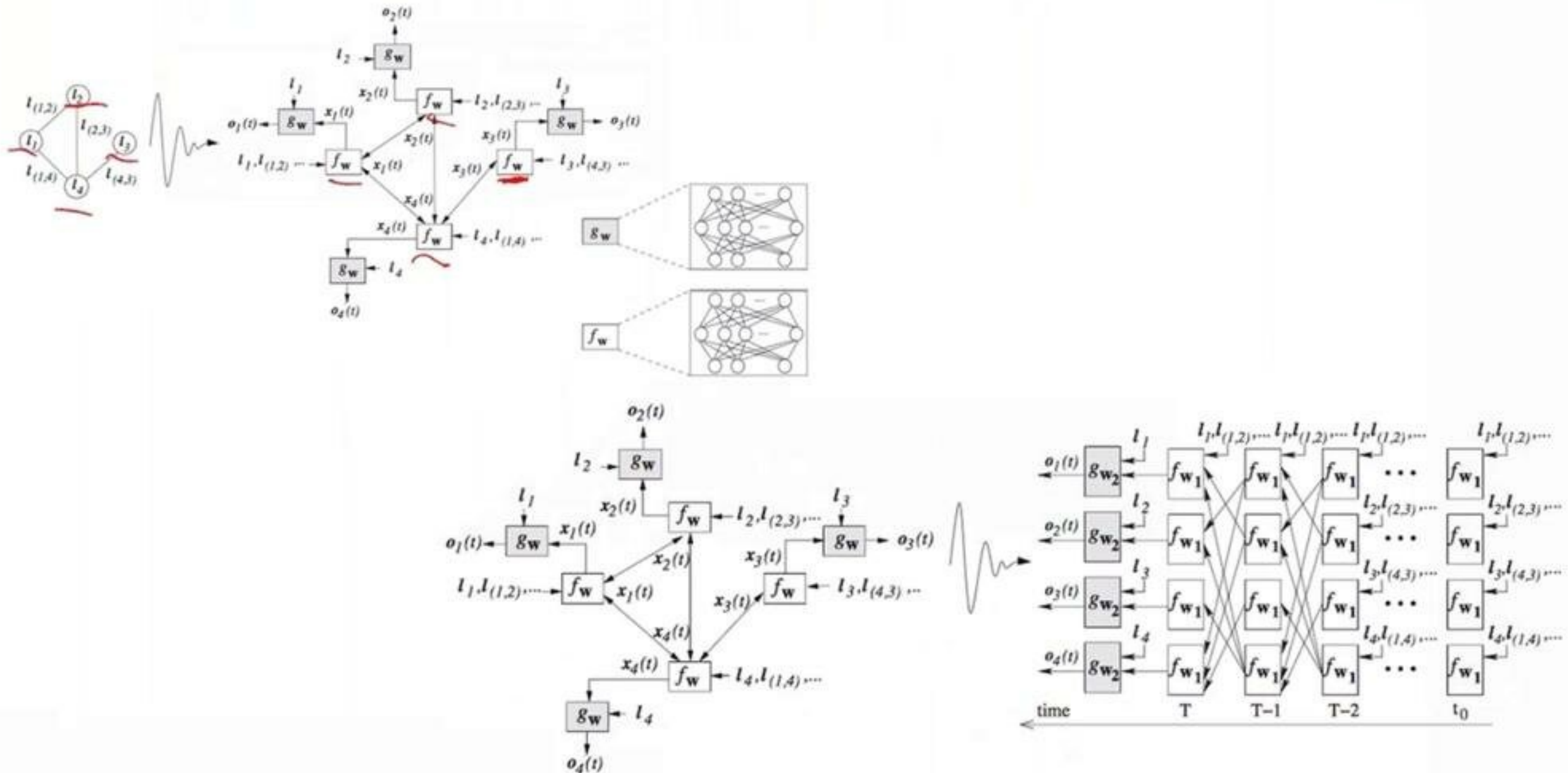


- Recurrent neurons, but similar in spirit to NN4G.
- The unrolled network is close to NN4G, but:
 - tied parameters
 - defined as a contraction mapping

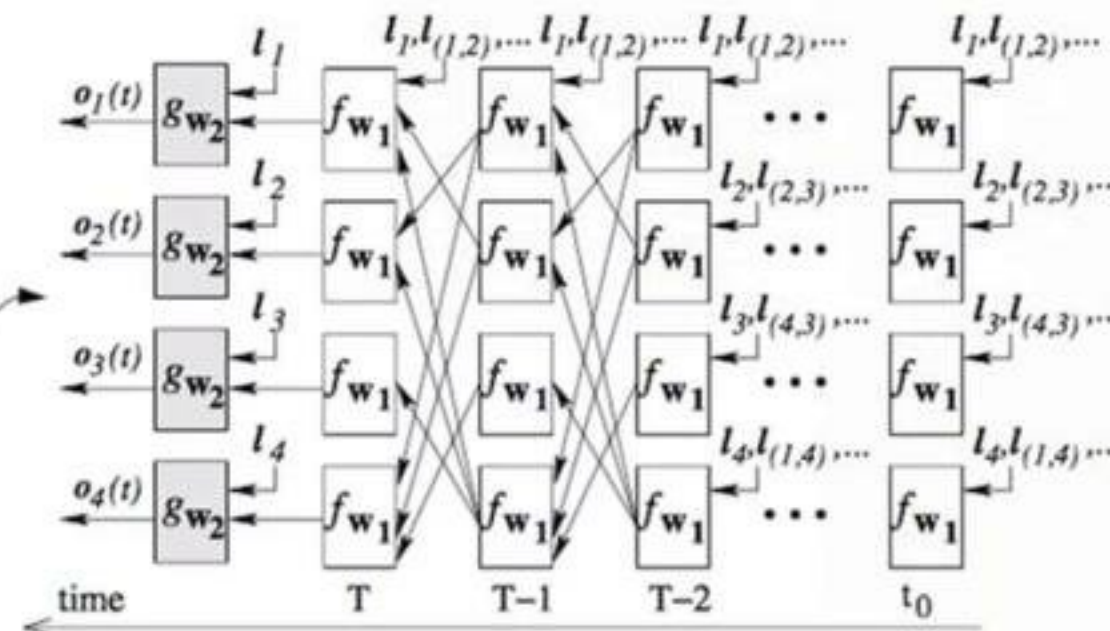
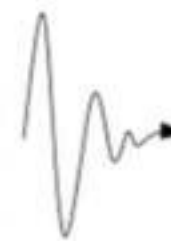
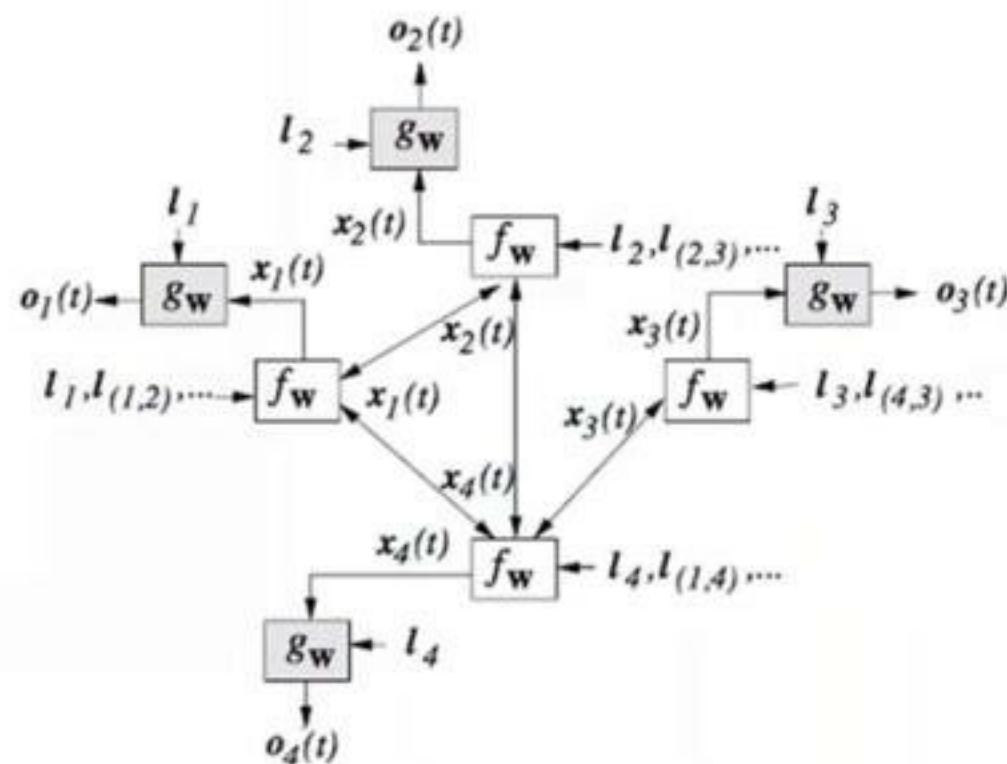
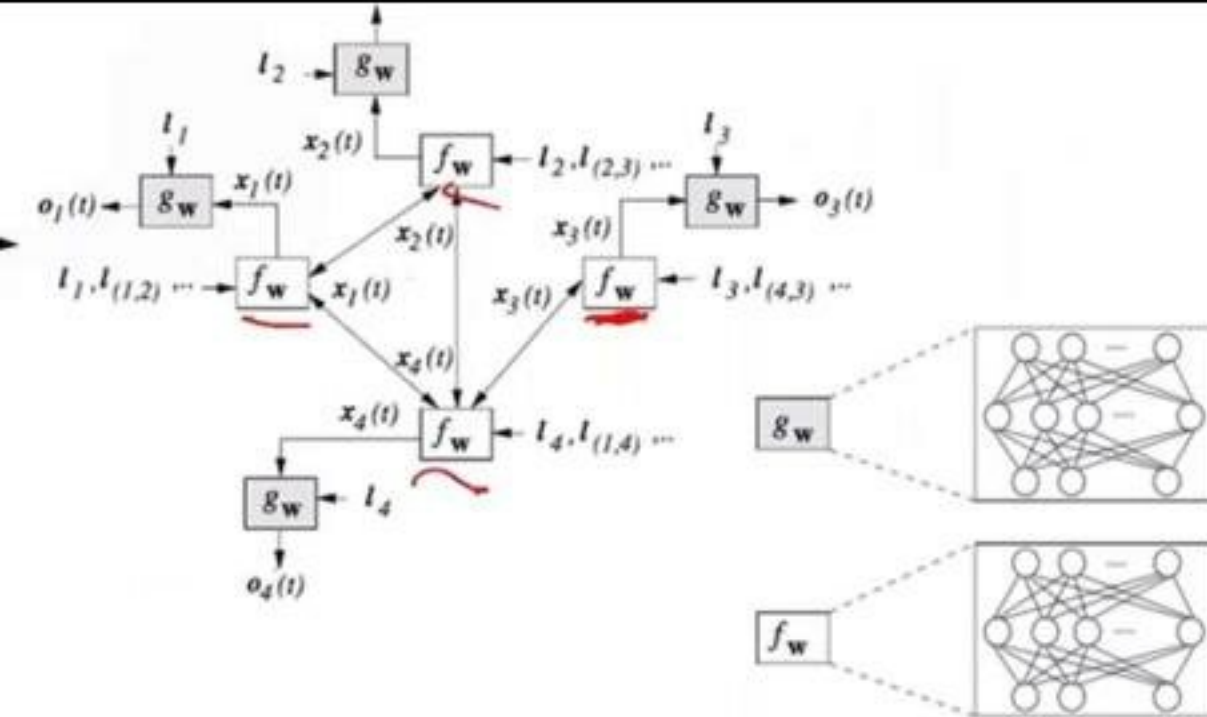
UNN by Scarsen, Gorn et al.



GNN by Scarselli, Gori et al.



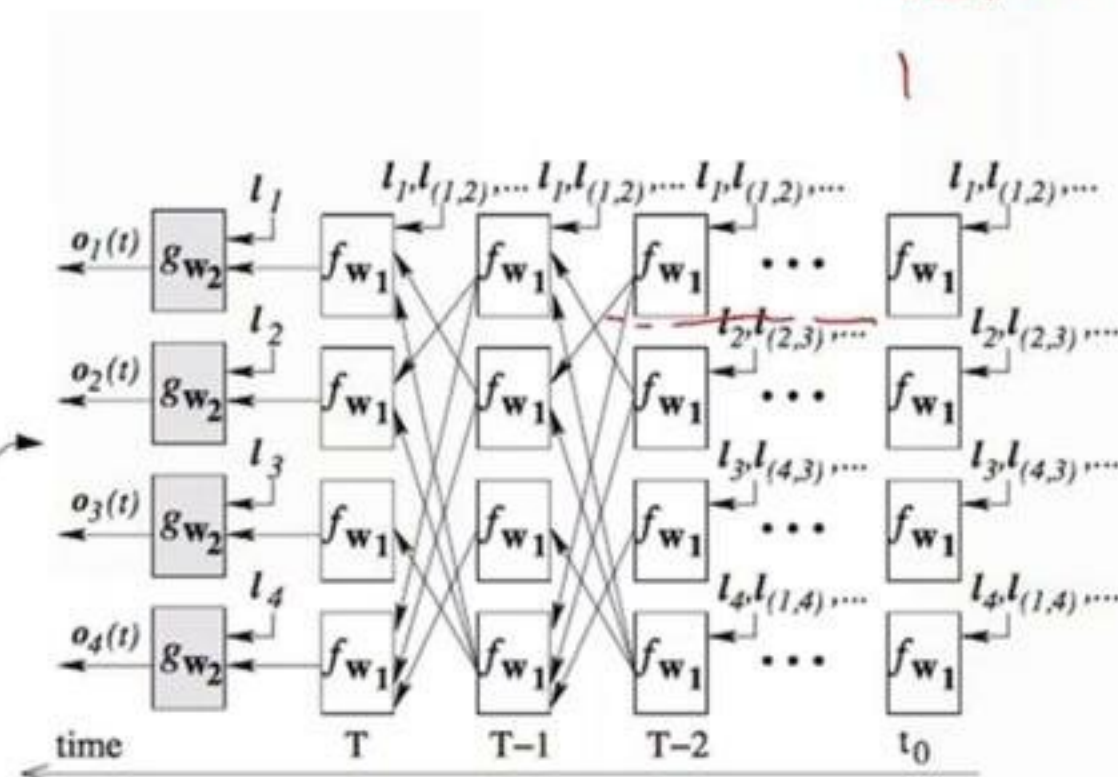
- Recurrent neurons, but similar in spirit to NN4G.
- The unrolled network is close to NN4G, but:
 - tied parameters
 - defined as a contraction mapping



urrent neurons, but similar in spirit to NN4G.

unrolled network is close to NN4G, but:

- tied parameters
- defined as a contraction mapping



INPUT

Convolution operator (on images)

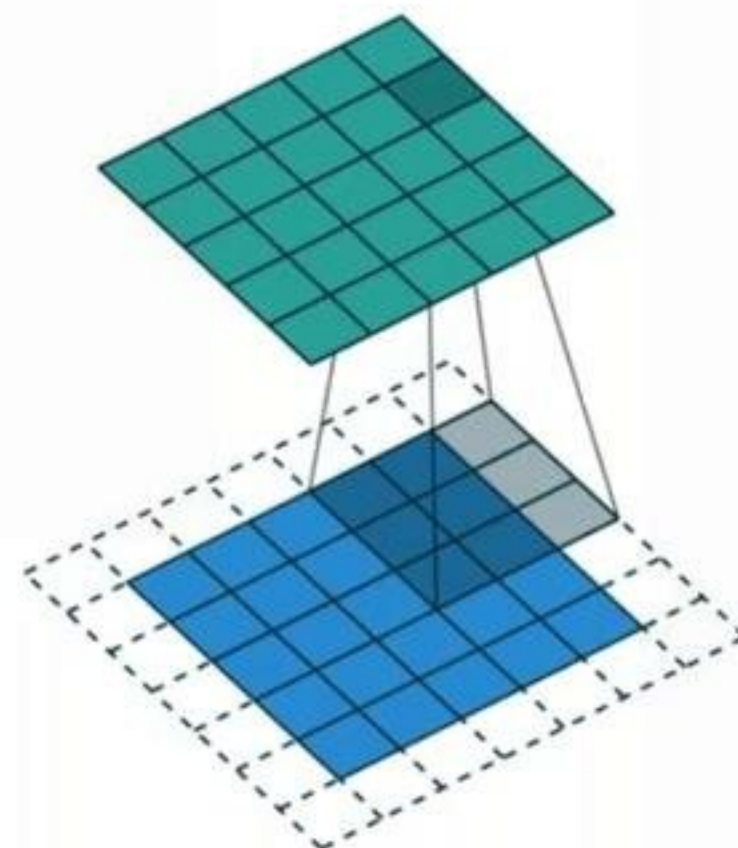
General definition for filter f and signal x :

$$(f * x)(t) = \int_{-\infty}^{\infty} f(\tau)x(t - \tau) d\tau$$

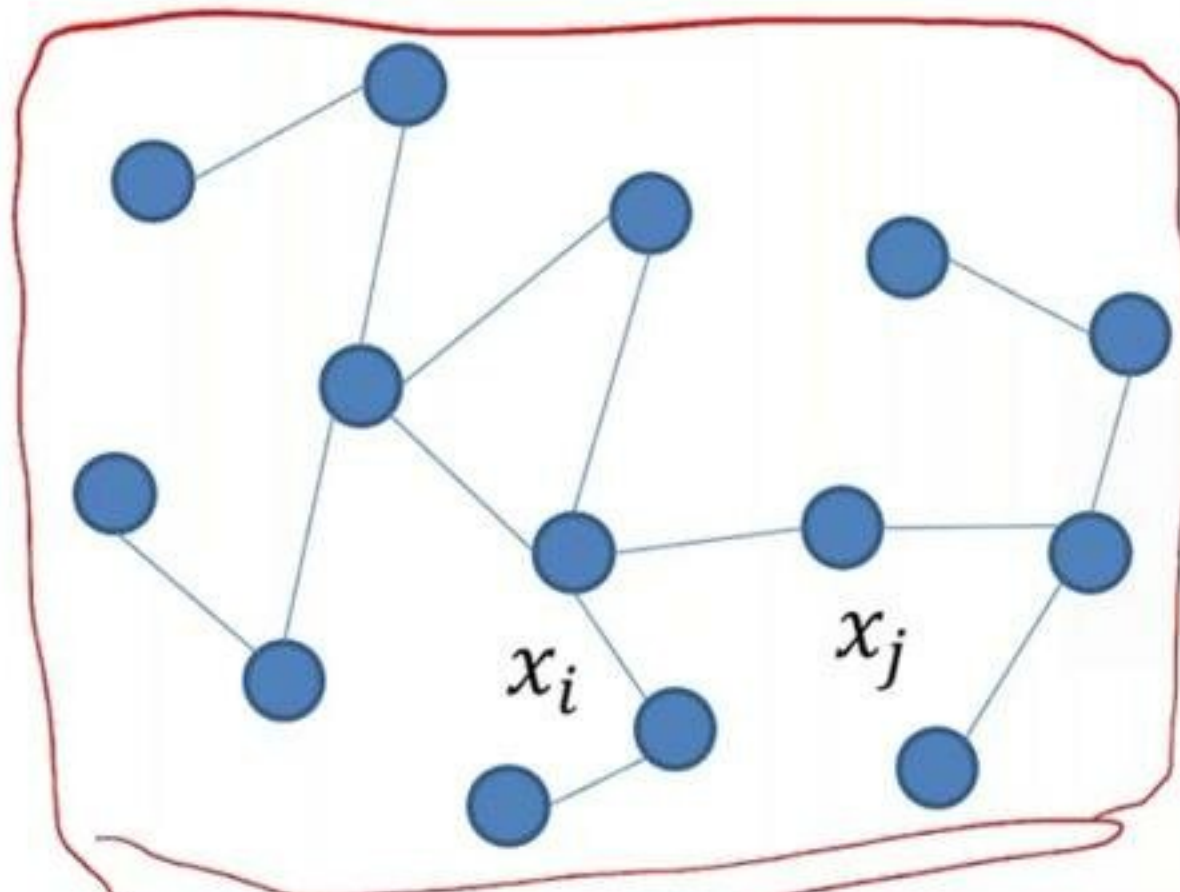
In images, it corresponds to

$$(f * x)(i, j) = \sum_{s=-a}^a \sum_{t=-b}^b f(s, t)x(i - s, j - t)$$

Where f is a $2a \times 2b$ filter and x an image



How to define convolution on graphs?



Consider a simple setting:

- single undirected graph
- $x: V \rightarrow \mathbb{R}$: a signal on the nodes of a graph
 - Represented as vector $\underline{x} \in \mathbb{R}^n$

The convolution operator is difficult to define in the vertex domain

Convolution Theorem:

- Convolution in one domain (time, space)
- corresponds to pointwise multiplication in frequency domain

$$\widehat{f * g} = \widehat{f} \odot \widehat{g}$$

\widehat{f} : Fourier transform of f

- \odot Hadamard (element-wise) product

How to define convolution on graphs?



Main steps:

- Graph Fourier Transform
 - Fourier Basis are eigenvectors of normalized Graph Laplacian

$$L = U \Lambda U^T$$

- We can then define the graph convolution in the frequency domain

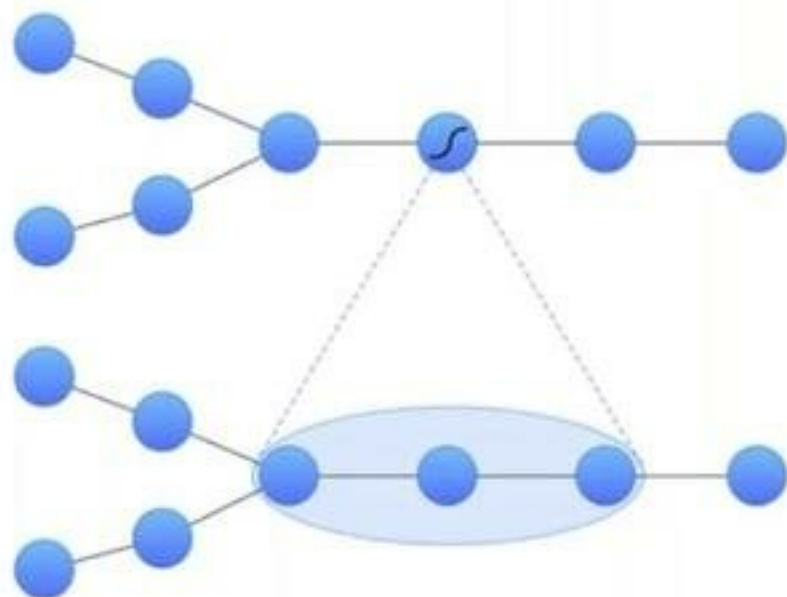
$$f *_G x = U \hat{F} U^T x$$

where $\hat{F} = \text{diag}(\hat{f})$

- For some choice of filters, e.g. polynomials of the spectral matrix
 - The convolution can be computed **in the node space** directly

$$\hat{F}_\Theta = \sum_{k=0}^K \theta_k \Lambda^k \rightarrow f *_G x = \sum_{k=0}^K \theta_k L^k x$$

Summary of Graph Convolution



- **1-localized** GCN maps multisets of representations (node and neighbours at the previous layer) to a new one:

$$H_v^{l+1} = f(\{H_v^l, H_u^l, \forall u \in ne(v)\})$$

- where $H^0 = X$

- f : linear mapping & non-linear activation function, e.g.

$$H^{l+1} = \sigma(\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2} H^l \Theta^l)$$

The convolution operator can be generalized to be more expressive than 1-WL [6]

If f is expressive enough (and with an injective readout), a multilayer 1-localized GCN is as expressive as the 1-dim WL isomorphism test

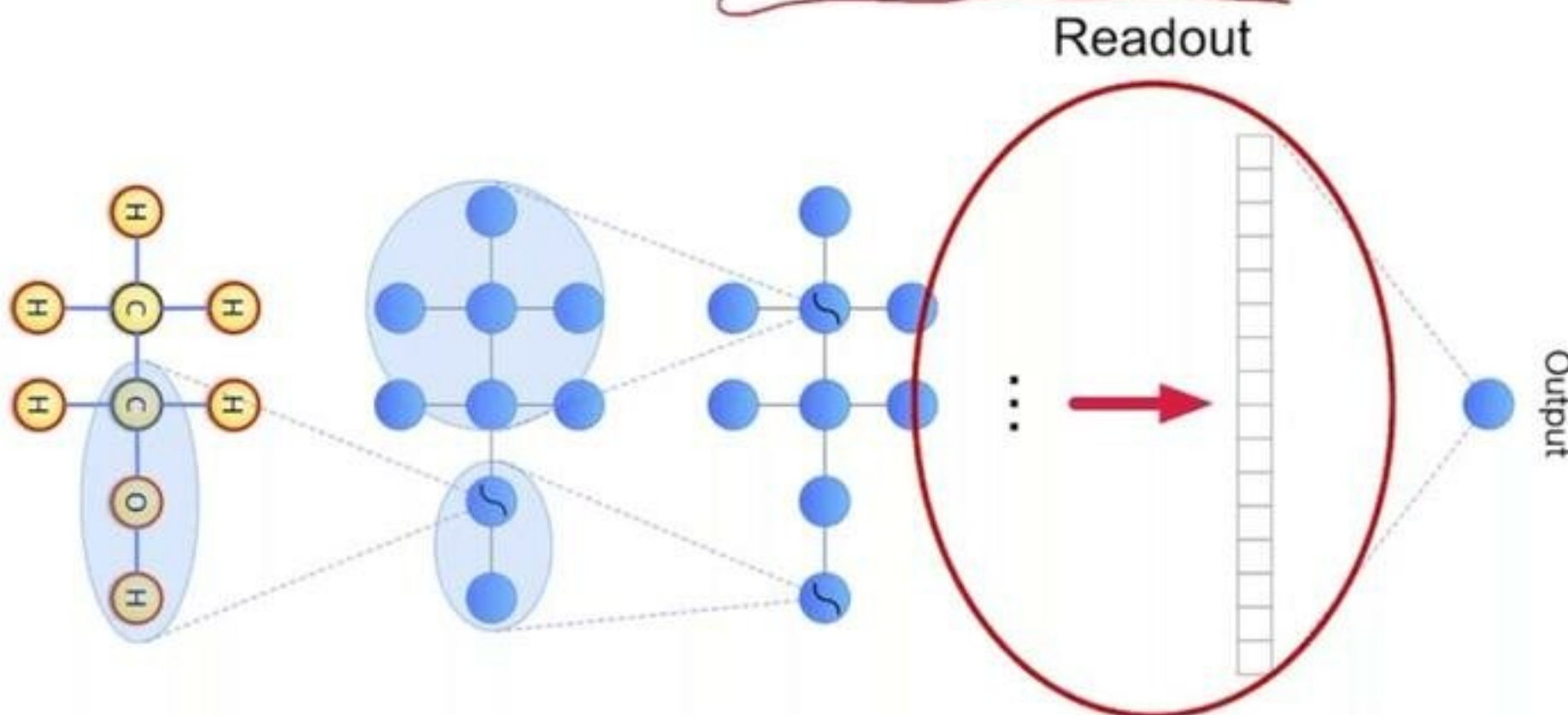
Readout Layer (back to graph classification)



With GC we have a representation for each graph node.

How can we map node representations to a graph-level representation?

- **Readout** function:
 - Maps a (multi) **set** of node representations to a **graph-level representation**
 - Differentiable
- Naïve solutions:
 - sum (or average) of node representations
 - More complex alternatives: Universal readout (DeepSets) [7]



[7] Navarin, N., Tran, D. Van, & Sperduti, A. (2019). Universal Readout for Graph Convolutional Neural Networks. *International Joint Conference on Neural Networks*. Budapest, Hungary.

- In many cases, **graph kernels perform better** (or comparably) to GCNs
- In particular in the domain of molecules

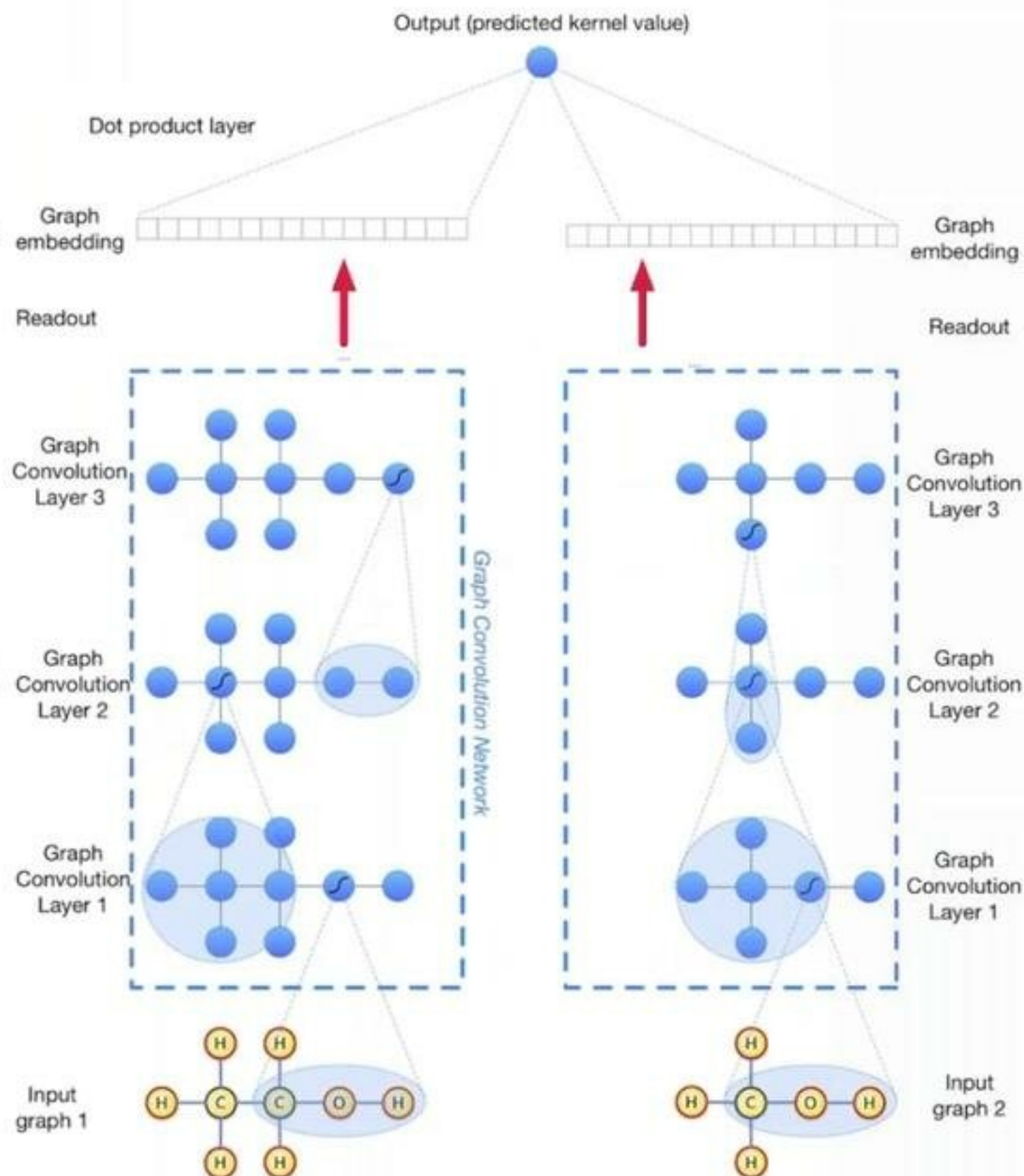
Our idea is to integrate the knowledge carried by the kernel in the representation learned by the GCN

Method/Dataset		MUTAG	PTC	NCI1	PROTEINS	D&D
Graph kernels	RW	79.17±2.07	55.91±0.32	>3 days	59.57±0.09	>3 days
	PK	76.00±2.69	59.50±2.44	82.54±0.47	73.68±0.68	78.25±0.51
	WL	84.11±1.91	57.97±2.49	84.46±0.45	74.68±0.49	78.34±0.62
Graph Neural Networks	PSCN	-	-	76.34±1.68	75.00±2.51	76.27±2.64
	CapsGCN	86.67±6.88	-	78.35±1.55	76.28±3.63	75.38±4.17
	DGCNN	82.48±1.49	57.14±2.19	72.97±0.92	73.96±0.41	78.09±0.72

Pre-training GNNs with kernels

First attempt: **pre-training**

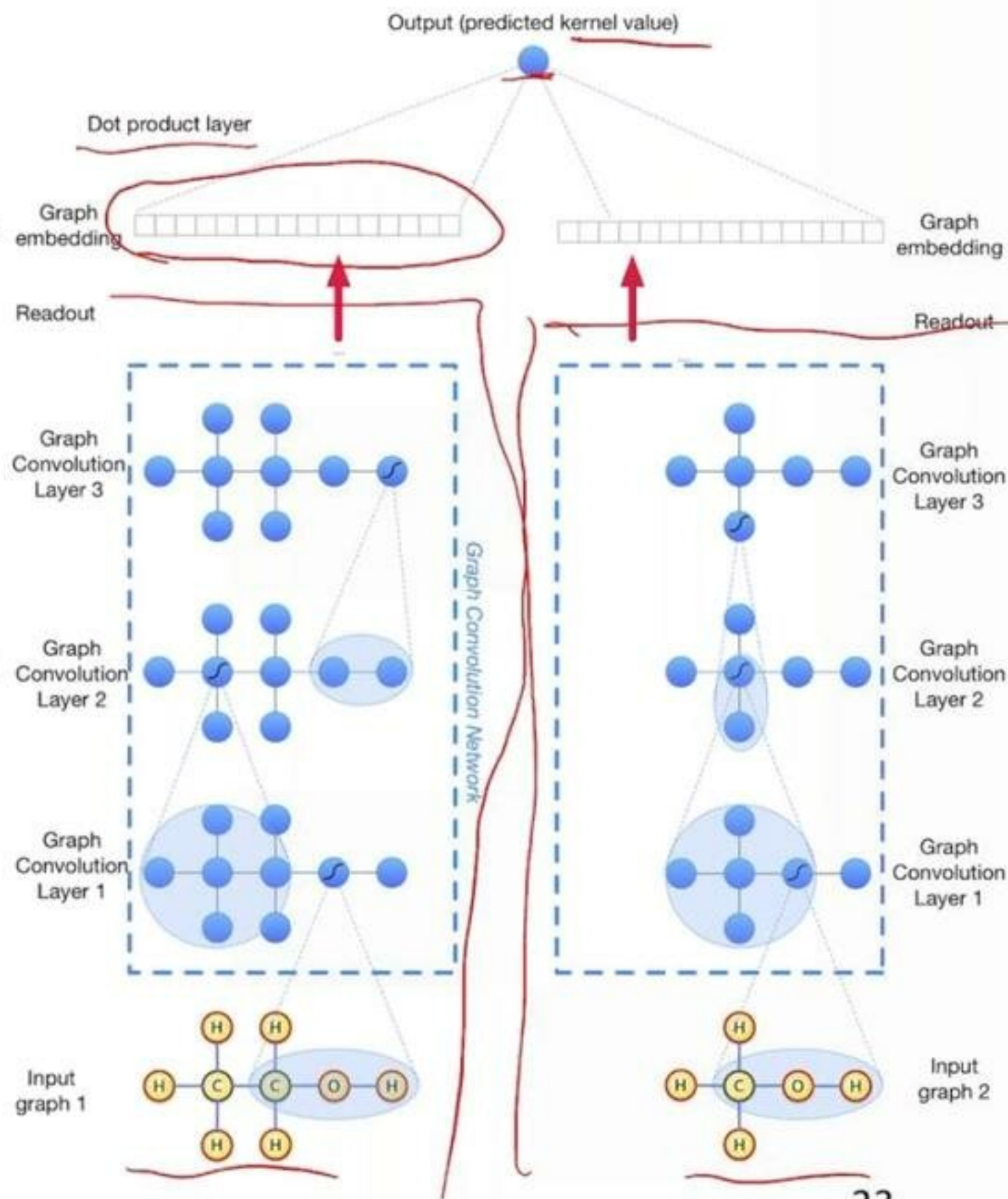
- Kernels are defined over **pairs** of graphs, so we define a siamese GCN architecture
- The target is the kernel value
- After pre-training, we can train a single network using supervised labels as usual



Pre-training GNNs with kernels

First attempt: **pre-training**

- Kernels are defined over **pairs** of graphs, so we define a siamese GCN architecture
- The target is the kernel value
- After pre-training, we can train a single network using supervised labels as usual



Pre-training GNNs with kernels



Pros:

- Promising results
- Unsupervised
- Pre training acts as a bias toward good representations
- bound on (training) error depending on reconstruction loss

Cons:

- Computational time: quadratic in the number of examples

Method / Dataset	MUTAG	PTC	NCI1
GK	81.39±1.74	55.65±0.46	62.49±0.27
RW	79.17±2.07	55.91±0.32	>3 days
PK	76.00±2.69	59.50±2.44	82.54±0.47
WL	84.11±1.91	57.97±2.49	84.46±0.45
DGCNN	85.83±1.66	58.59±2.47	74.44±0.47
Pre-trained DGCNN	88.10±1.05	61.03±2.86	77.13±0.45

Multi-task training approach

Second approach: **Multi-task training**

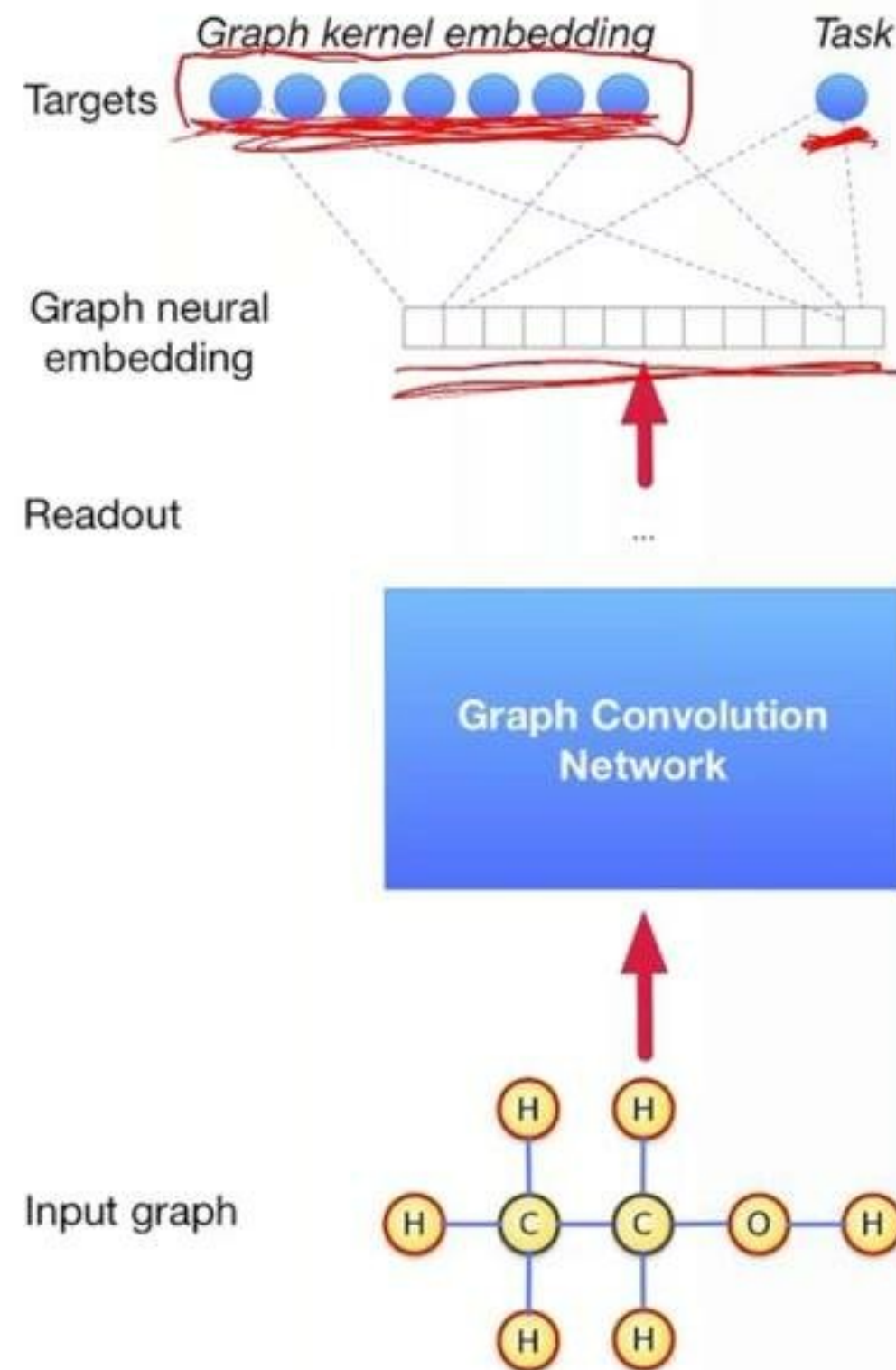
- Single network, two outputs
- **Hashing** to map the kernel embedding to a fixed, small vector
- We try to reconstruct the kernel embedding as a secondary output

Pros:

- No added complexity (pre-processing of kernel computation)
- Any graph embedding can be used (i.e. domain knowledge)

Cons:

- We don't have control in where the network will store information from the embedding



Third approach: **Layer-wise Multi-task training**

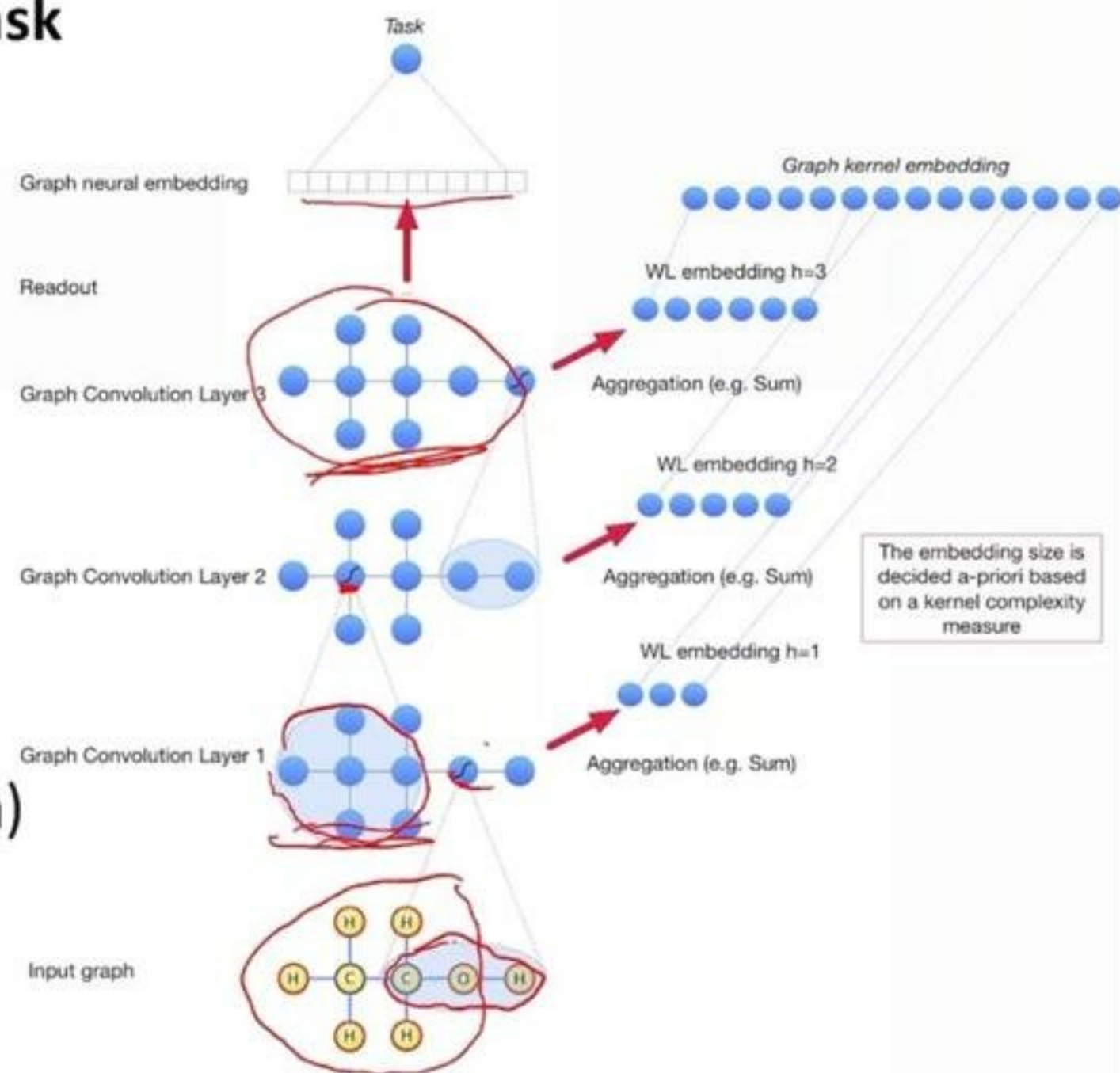
- Single network, multiple outputs
- **Split features** in the embedding according to their **complexity**
- Deeper layers reconstruct more complex features

Pros:

- No added complexity (pre-processing of kernel computation)

Cons (future works):

- Even lower-grained supervision possible?



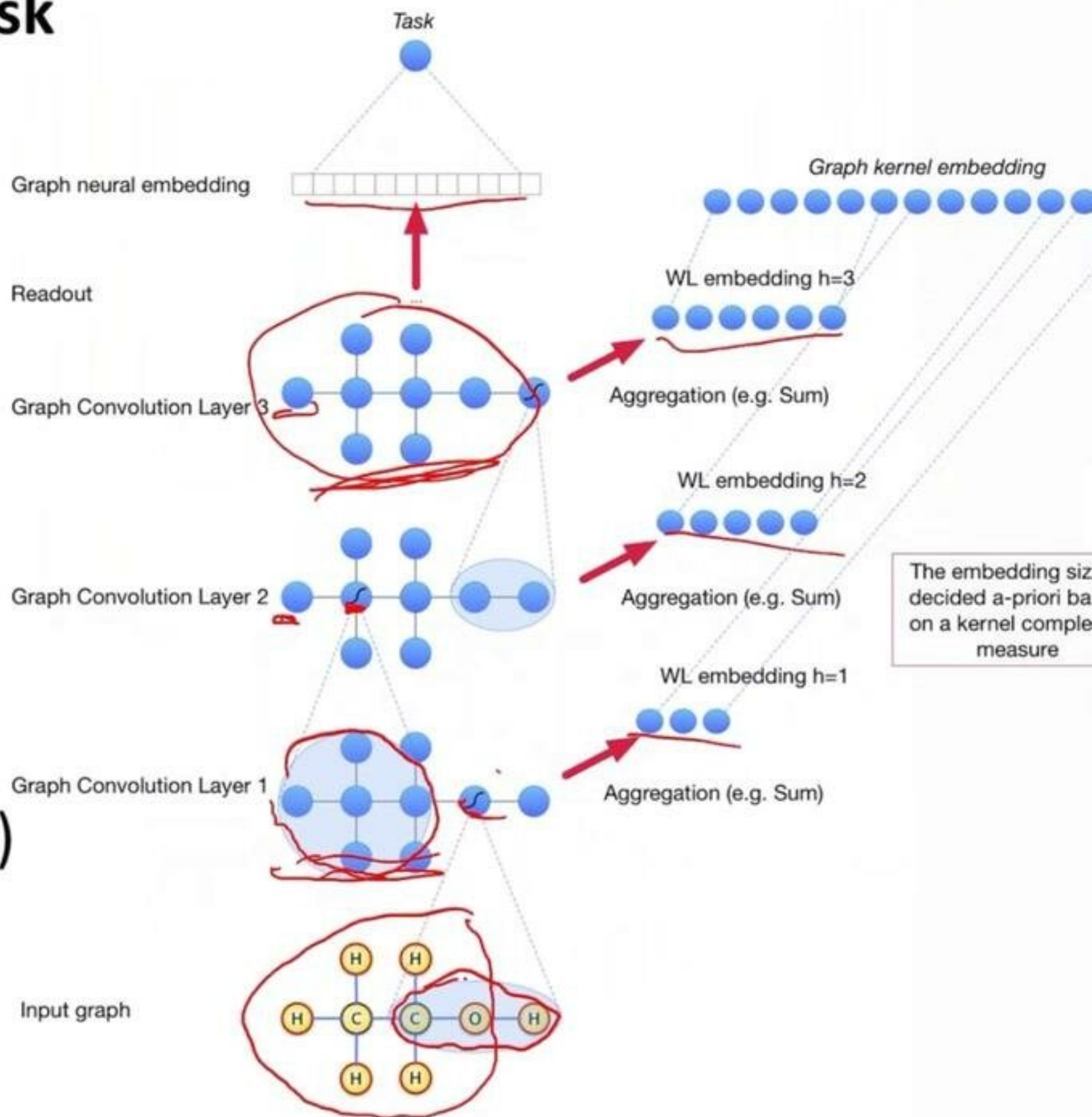
Layer-wise Multi-task

multiple outputs
in the embedding
their **complexity**
reconstruct more
res

complexity (pre-
kernel computation)

ks):

ained supervision



Results for Kernel-based embeddings



- All kernel-based training approaches improve over standard training
- In general, $MT < PT < LMT$
- The techniques are applicable to all GNNs
 - The improvement may vary depending on the architecture

	Method/Dataset	MUTAG	PTC	NCI1	PROTEINS	D&D
	WL (h=3)	76.79*±3.17	57.48*±1.36	82.13±2.17	69.63*±1.22	73.64*±2.56
	DGCNN	82.48*±1.49	57.14*±2.19	72.97*±0.87	73.96*±0.41	78.09*±0.72
	LMT-FGCNN	86.81 ±1.75	59.04±0.94	82.20 ±0.54	76.03 ±0.68	80.14 ±0.76
Ablation study	FGCNN	84.49±1.90	58.82±1.80	81.50±0.39	74.57*±0.80	77.47*±0.86
	LMT-DGCNN	85.00±1.15	59.39 ±0.51	77.02*±0.48	74.61*±0.89	78.11*±0.61
	MT-DGCNN	83.68±1.29	58.39±1.11	76.55*±0.40	74.42*±0.36	78.17*±0.57
	MT-FGCNN	85.81±1.62	59.23±2.35	81.86±0.41	75.18±0.66	79.90±0.39
	PT-DGCNN	85.38±1.47	58.48±1.92	75.20*±0.87	75.19±0.42	78.38*±0.55

Results for Kernel-based embeddings



- All kernel-based training approaches improve over standard training
- In general, $MT < PT < LMT$
- The techniques are applicable to all GNNs
 - The improvement may vary depending on the architecture

	Method/Dataset	MUTAG	PTC	NCI1	PROTEINS	D&D
gASU	WL (h=3)	76.79*±3.17	57.48*±1.36	82.13±2.17	69.63*±1.22	73.64*±2.56
	DGCNN	82.48*±1.49	57.14*±2.19	72.97*±0.87	73.96*±0.41	78.09*±0.72
	LMT-FGCNN	86.81 ±1.75	59.04±0.94	82.20±0.54	76.03±0.68	80.14 ±0.76
	FGCNN	84.49±1.90	58.82±1.80	81.50±0.39	74.57*±0.80	77.47*±0.86
	LMT-DGCNN	85.00±1.15	59.39 ±0.51	77.02*±0.48	74.61*±0.89	78.11*±0.61
	MT-DGCNN	83.68±1.29	58.39±1.11	76.55*±0.40	74.42*±0.36	78.17*±0.57
Ablation study	MT-FGCNN	85.81±1.62	59.23±2.35	81.86±0.41	75.18±0.66	79.90±0.39
	PT-DGCNN	85.38±1.47	58.48±1.92	75.20*±0.87	75.19±0.42	78.38*±0.55

Results for Kernel-based embeddings



- All kernel-based training approaches improve over standard training
- In general, $MT < PT < LMT$
- The techniques are applicable to all GNNs
 - The improvement may vary depending on the architecture

	Method/Dataset	MUTAG	PTC	NCI1	PROTEINS	D&D
gAsu	WL (h=3)	76.79*±3.17	57.48*±1.36	82.13±2.17	69.63*±1.22	73.64*±2.56
	DGCNN	82.48*±1.49	57.14*±2.19	72.97*±0.87	73.96*±0.41	78.09*±0.72
	<u>LMT-FGCNN</u>	86.81 ±1.75	59.04±0.94	82.20±0.54	76.03±0.68	80.14 ±0.76
Ablation study	FGCNN	84.49±1.90	58.82±1.80	81.50±0.39	74.57*±0.80	77.47*±0.86
	LMT-DGCNN	85.00±1.15	59.39 ±0.51	77.02*±0.48	74.61*±0.89	78.11*±0.61
	MT-DGCNN	83.68±1.29	58.39±1.11	76.55*±0.40	74.42*±0.36	78.17*±0.57
	MT-FGCNN	85.81±1.62	59.23±2.35	81.86±0.41	75.18±0.66	79.90±0.39
	PT-DGCNN	85.38±1.47	58.48±1.92	75.20*±0.87	75.19±0.42	78.38*±0.55

Results for Kernel-based embeddings



- All kernel-based training approaches improve over standard training
- In general, $MT < PT < LMT$
- The techniques are applicable to all GNNs
 - The improvement may vary depending on the architecture

	Method/Dataset	MUTAG	PTC	NCI1	PROTEINS	D&D
Ablation study	WL (h=3)	76.79*±3.17	57.48*±1.36	82.13±2.17	69.63*±1.22	73.64*±2.56
	DGCNN	82.48*±1.49	57.14*±2.19	72.97*±0.87	73.96*±0.41	78.09*±0.72
	<u>LMT-FGCNN</u>	86.81 ±1.75	59.04±0.94	82.20 ±0.54	76.03 ±0.68	80.14 ±0.76
	<u>FGCNN</u>	84.49±1.90	58.82±1.80	81.50±0.39	74.57*±0.80	77.47*±0.86
	LMT-DGCNN	85.00±1.15	59.39 ±0.51	77.02*±0.48	74.61*±0.89	78.11*±0.61
	<u>MT-DGCNN</u>	83.68±1.29	58.39±1.11	76.55*±0.40	74.42*±0.36	78.17*±0.57
	MT-FGCNN	85.81±1.62	59.23±2.35	81.86±0.41	75.18±0.66	79.90±0.39
	<u>PT-DGCNN</u>	85.38±1.47	58.48±1.92	75.20*±0.87	75.19±0.42	78.38*±0.55

gASU

- We can add more (unlabelled) data in the training
- Semi-supervised learning
- The more data we add, the higher the performance

Dataset/ Method	DGCNN	FGCNN	LMT-FGCNN		
			+0	+1	+2
NCI1B	72.92 ± 0.56	79.27 ± 0.70	81.01 ± 0.56	81.19 ± 0.46	82.07 ± 0.21
NCI33B	75.00 ± 0.42	81.75 ± 0.67	81.81 ± 0.20	82.60 ± 0.39	82.69 ± 0.56
NCI41B	70.94 ± 0.53	78.30 ± 0.67	79.02 ± 0.17	79.10 ± 0.40	79.54 ± 0.15

- We can add more (unlabelled) data in the training
- Semi-supervised learning
- The more data we add, the higher the performance

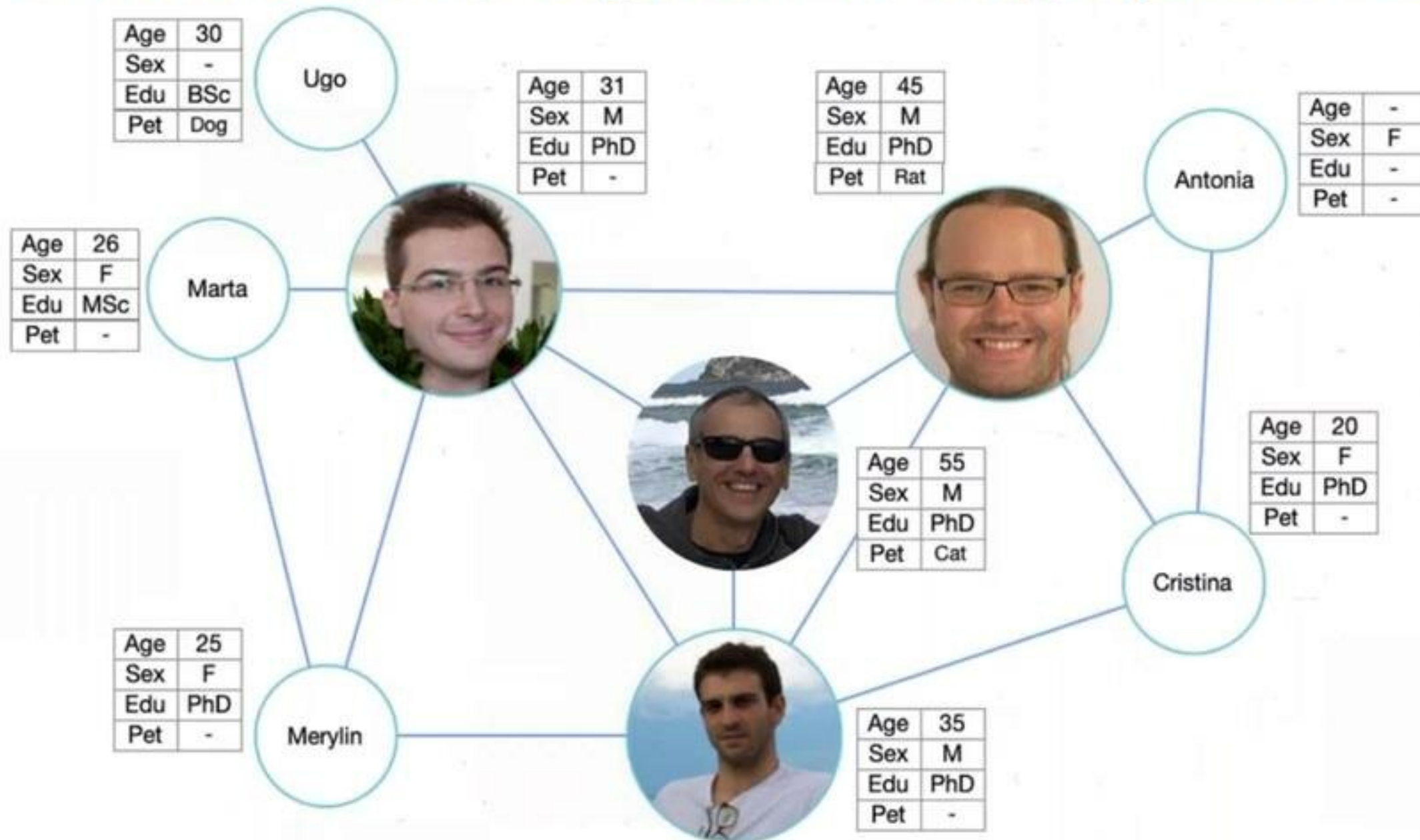
Dataset/ Method	DGCNN	FGCNN	LMT-FGCNN		
			+0	+1	+2
NCI1B	72.92 ± 0.56	79.27 ± 0.70	81.01 ± 0.56	81.19 ± 0.46	82.07 ± 0.21
NCI33B	75.00 ± 0.42	81.75 ± 0.67	81.81 ± 0.20	82.60 ± 0.39	82.69 ± 0.56
NCI41B	70.94 ± 0.53	78.30 ± 0.67	79.02 ± 0.17	79.10 ± 0.40	79.54 ± 0.15

- In this talk, I presented different ways to merge:
 - Kernel methods for graphs
 - Neural Networks for Graphs
- Future work:
 - Other types of (pre) training
 - Recurrent GNNs
 - Work on relational graphs and graphs modelling more complex tasks

The background of the slide features a large, faint, red circular seal of the University of Twente. The seal contains a central illustration of two figures, likely saints or religious figures, standing under a gothic arch. The text 'UNIVERSITAS STUDII TWENTENSIS' is written around the perimeter of the seal, and 'MCCXXII' is at the bottom. The title 'Linear Graph Convolutional Networks' is centered over the seal.

Linear Graph Convolutional Networks

Classification/Regression on graph nodes



- Dataset: a single huge graph:
 - n vertices
 - d **attributes** associated to nodes: $X \in \mathbb{R}^n \times d$
 - target** associated to a small subset of nodes
- Given an unseen node, the task is to predict the correct target y

