# AIMS CDT - Signal Processing Michaelmas Term 2025

Xiaowen Dong

Department of Engineering Science

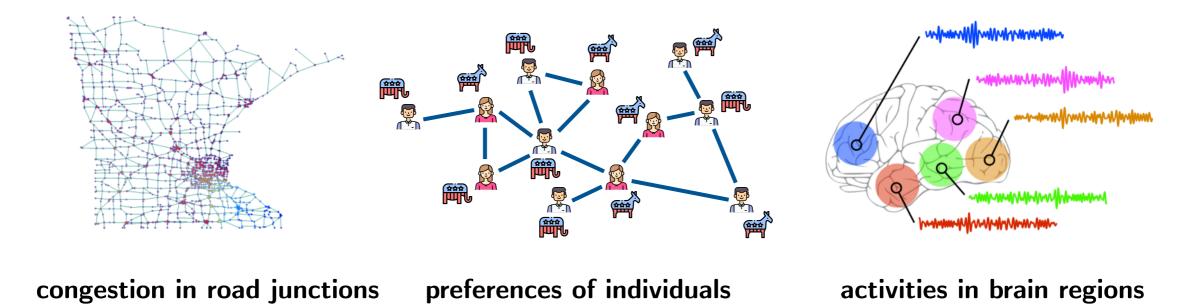


# Deep Learning on Graphs

#### Lecture 3

- Machine learning on graphs: Overview
- Convolutional neural networks on graphs
- Message passing neural networks
- Recent developments & Applications

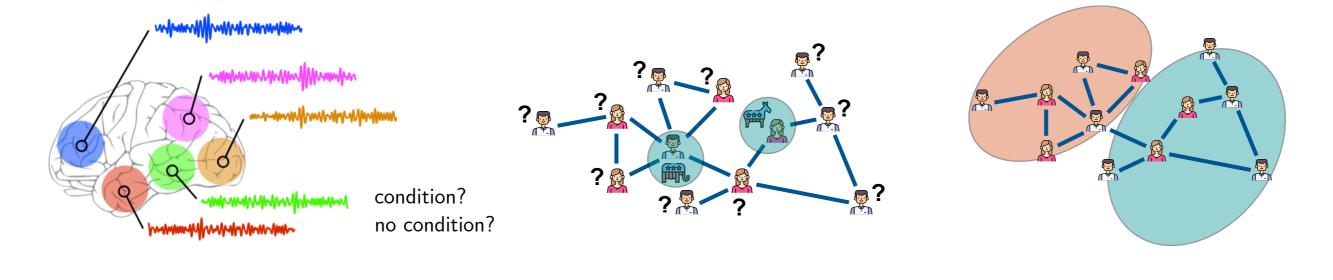
### Graph-structured data are pervasive



preferences of individuals

activities in brain regions

### Learning with graph-structured data

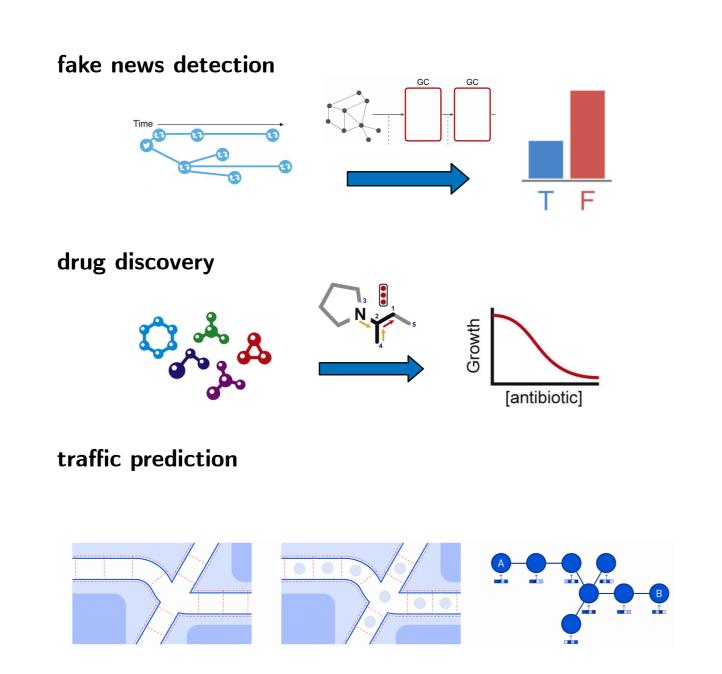


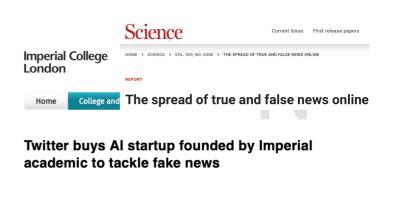
graph-level classification (supervised)

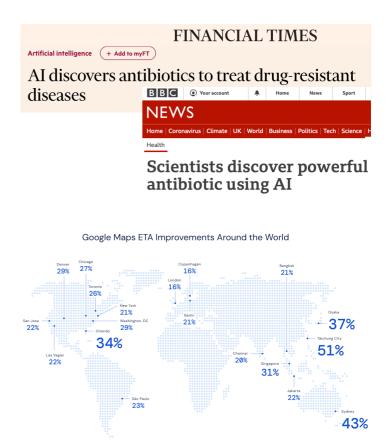
node-level classification (semi-supervised)

graph clustering (unsupervised)

#### Learning with graph-structured data



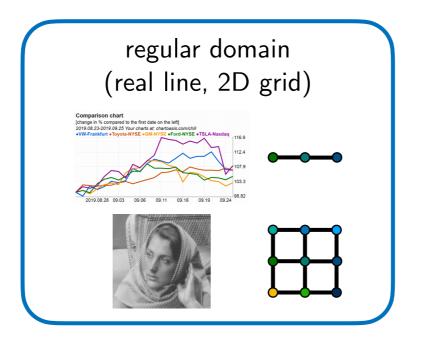




Monti et al., "Fake news detection on social media using geometric deep learning," ICLR Workshop, 2019. Stokes et al., "A deep learning approach to antibiotic discovery," Cell, 2020. Derrow-Pinion et al., "ETA prediction with graph neural networks in Google Maps," CIKM, 2021.

### Classical ML vs Graph ML

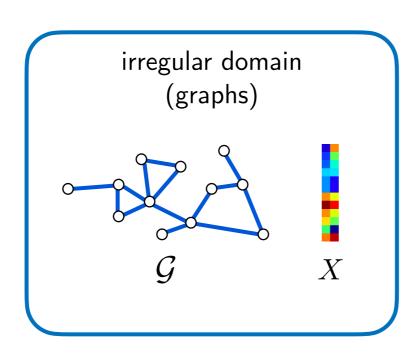
**Classical ML** 



(X) time series forecasting

image classification

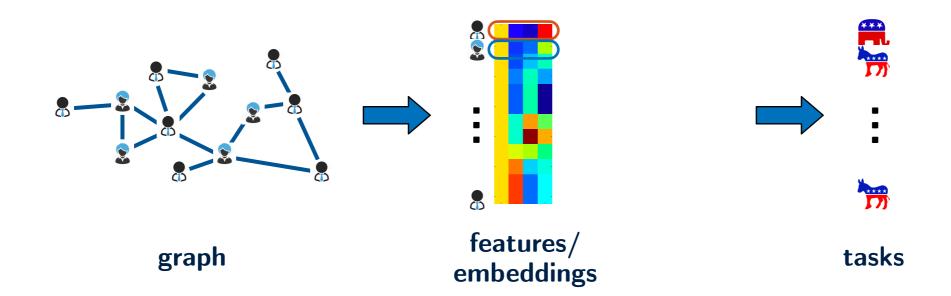
**Graph ML** 



 $f(\mathcal{G},X)$ 

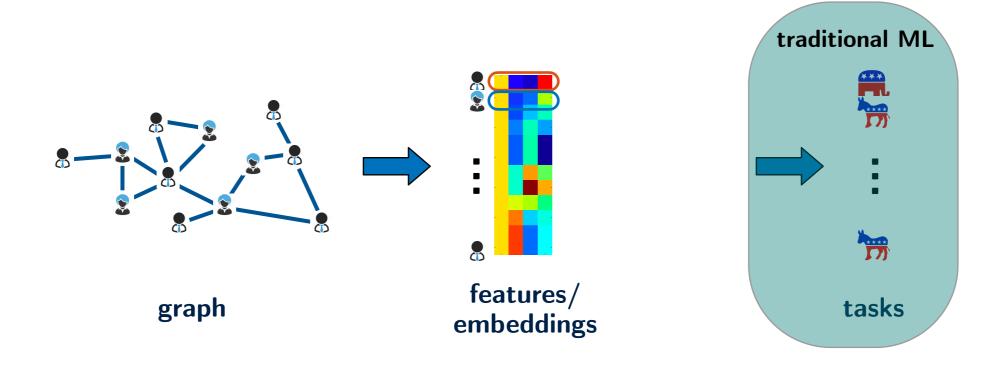
node classification
link prediction
graph classification
graph clustering

Traditional machine learning on graphs



- Limitations
  - hand-crafted features or optimised embeddings, often focused on graph structure

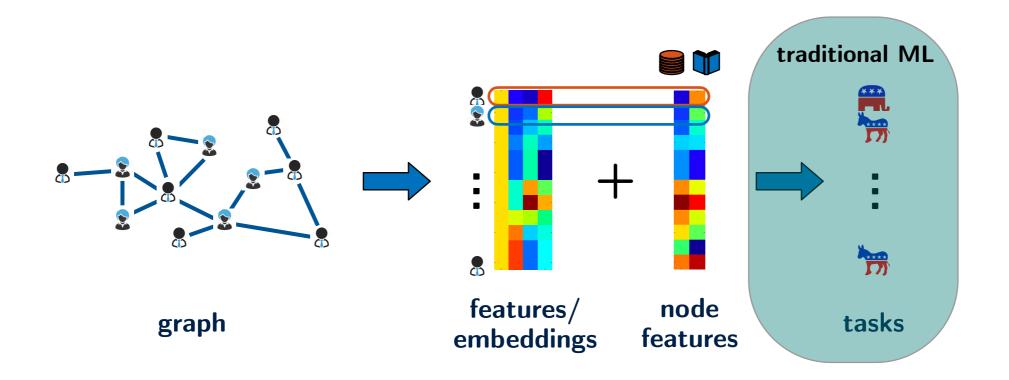
Traditional machine learning on graphs



#### Limitations

- hand-crafted features or optimised embeddings, often focused on graph structure
- respect notion of "closeness" in the graph, but do not adapt to downstream tasks

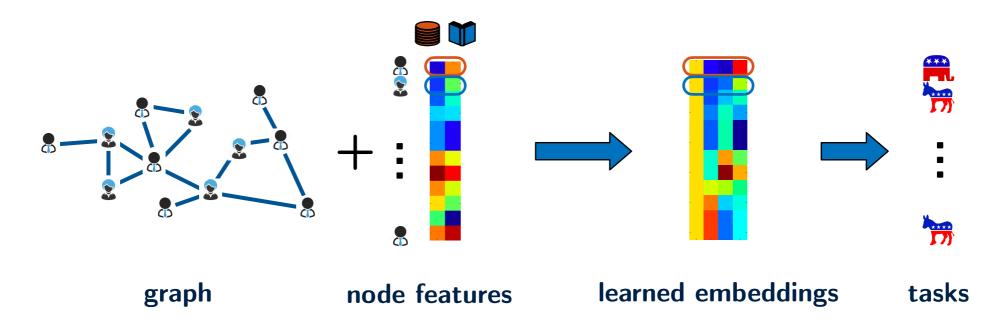
Traditional machine learning on graphs



#### Limitations

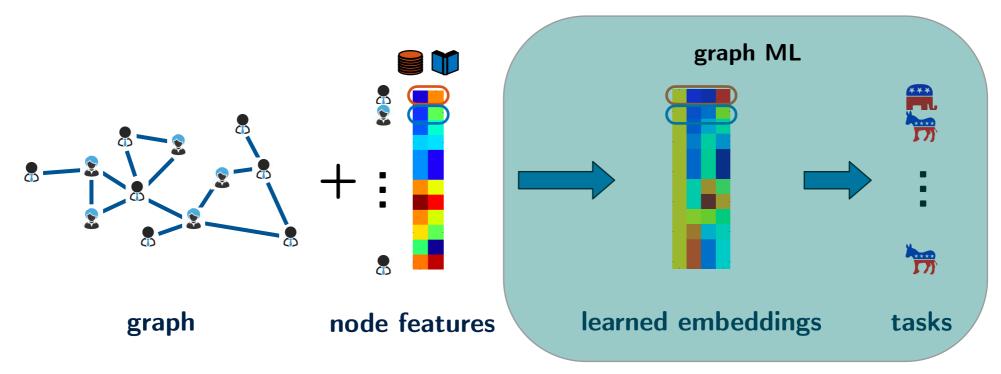
- hand-crafted features or optimised embeddings, often focused on graph structure
- respect notion of "closeness" in the graph, but do not adapt to downstream tasks
- can incorporate additional node features, but in a mechanical way

Graph machine learning



- Advantages
  - naturally combine graph structure and node features in analysis and learning
    - new tools: graph signal processing, graph neural networks

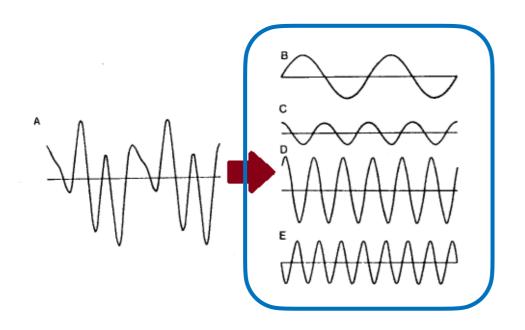
Graph machine learning



#### Advantages

- naturally combine graph structure and node features in analysis and learning
  - new tools: graph signal processing, graph neural networks
- embeddings can adapt to downstream tasks and be trained in end-to-end fashion
- offers more flexibility and enables "deeper" architectures and embeddings

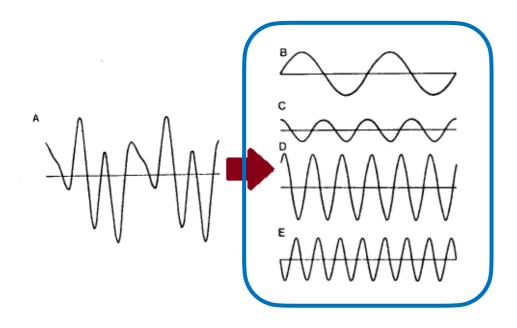
### Graph signal processing

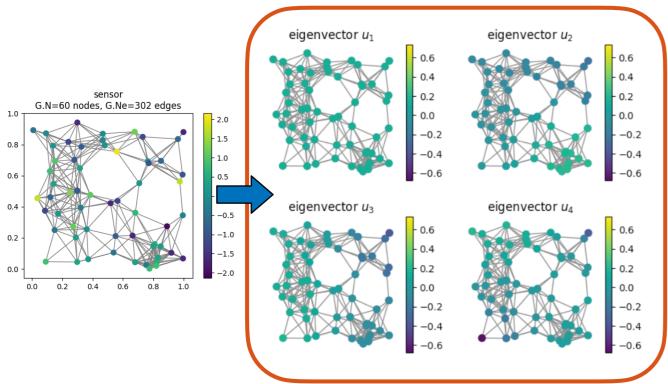


#### classical signal processing

- complex exponentials provide
   "building blocks" of 1D signal
   (different oscillations or frequencies)
- leads to Fourier transform
- enables convolution and filtering

#### Graph signal processing





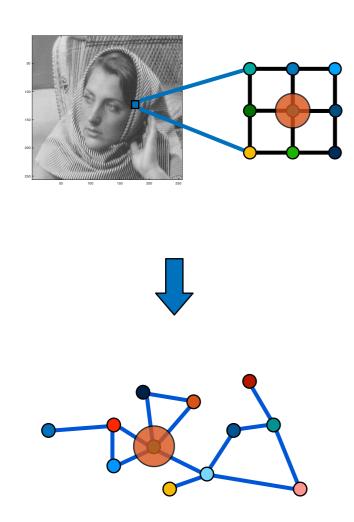
#### classical signal processing

- complex exponentials provide
   "building blocks" of 1D signal
   (different oscillations or frequencies)
- leads to Fourier transform
- enables convolution and filtering

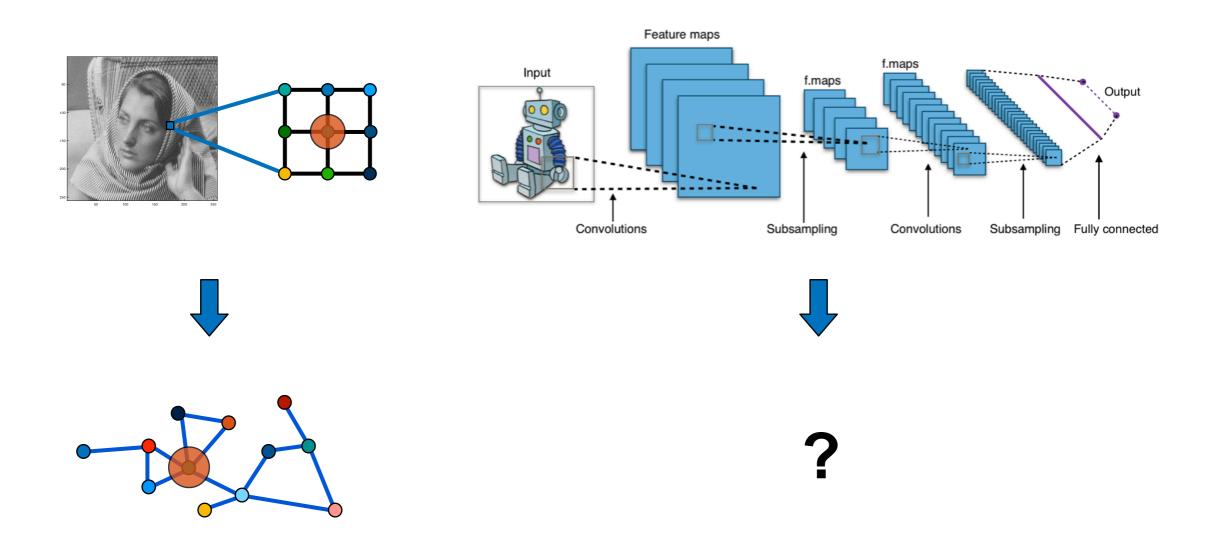
#### graph signal processing

- Laplacian eigenvectors provide
   "building blocks" of graph signal
   (different oscillation or frequencies)
- leads to graph Fourier transform
- enables convolution and filtering on graphs

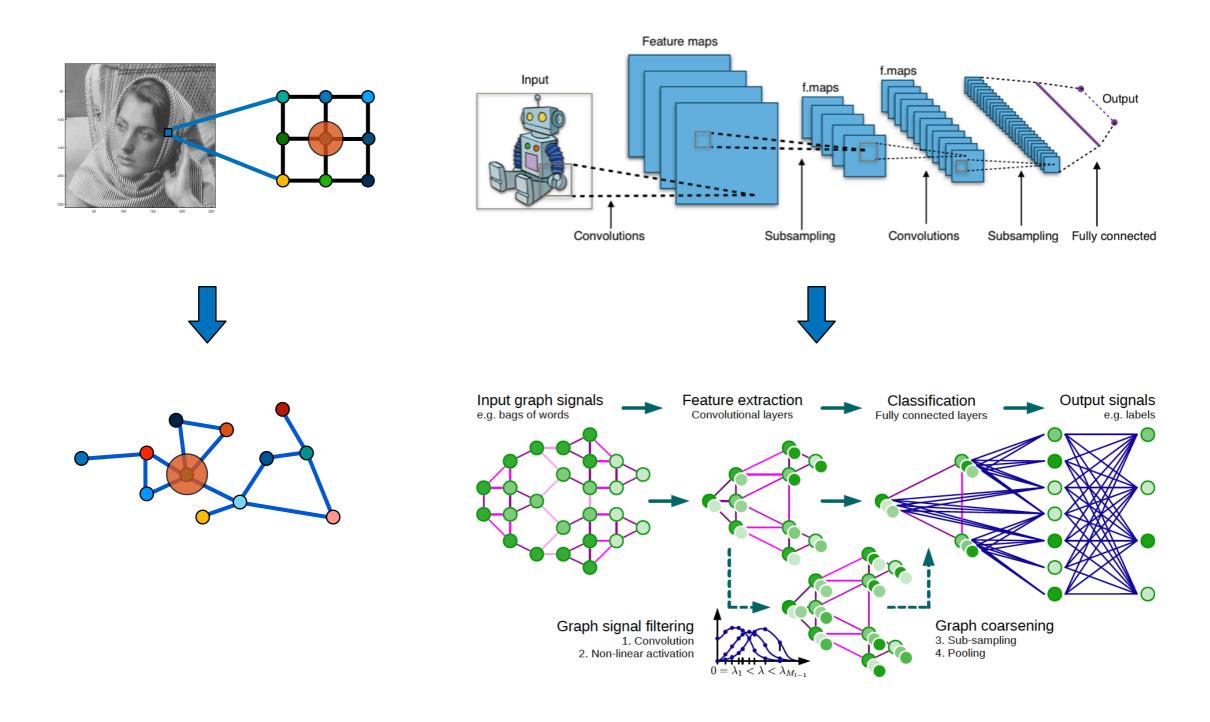
# Convolutional neural networks on graphs



### Convolutional neural networks on graphs



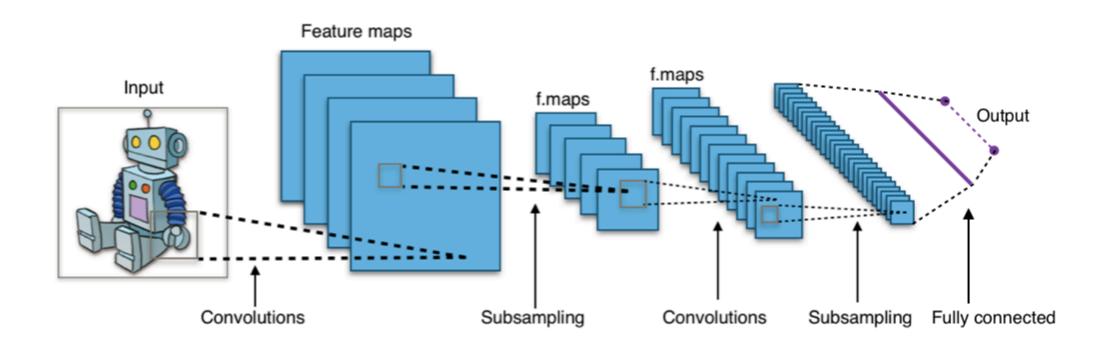
### Convolutional neural networks on graphs



#### Lecture 3

- Machine learning on graphs: Overview
- Convolutional neural networks on graphs
- Message passing neural networks
- Recent developments & Applications

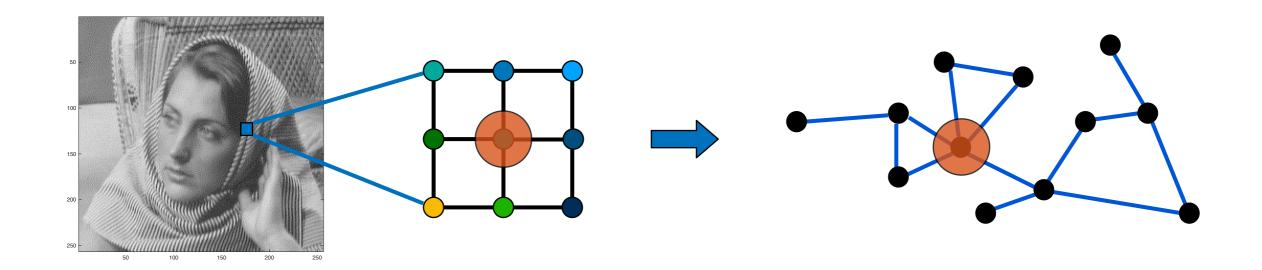
#### CNNs exploit structure within data



#### checklist

- convolution: translation equivariance
- localisation: compact filters (independent of sample dimension)
- multi-scale: compositionality
- **efficiency:**  $\mathcal{O}(N)$  computational complexity

### CNNs on graphs?



#### checklist

- **convolution**: how to define convolution? what about invariance?
- **localisation:** what is the notion of locality?
- multi-scale: how to down-sample on graphs?
- **efficiency:** how to keep the computational complexity low?

#### classical convolution

time domain

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

30	3	$2_2$	1	0
02	$0_2$	$1_0$	3	1
30	1,	2	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

#### classical convolution

time domain

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

30	3	$2_2$	1	0
02	$0_2$	$1_{0}$	3	1
30	1,	2	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

#### classical convolution

time domain

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



frequency domain

$$\widehat{(f * g)}(\omega) = \hat{f}(\omega) \cdot \hat{g}(\omega)$$

30	3	$2_{2}$	1	0
$0_2$	$0_2$	$1_0$	3	1
30	1,	22	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

#### classical convolution

convolution on graphs

time domain

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



frequency domain

$$\widehat{(f * g)}(\omega) = \hat{f}(\omega) \cdot \hat{g}(\omega)$$

graph spectral domain

$$\widehat{(f * g)}(\lambda) = ((\chi^T f) \circ \hat{g})(\lambda)$$

#### classical convolution

time domain

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



frequency domain

$$\widehat{(f * g)}(\omega) = \widehat{f}(\omega) \cdot \widehat{g}(\omega)$$

#### convolution on graphs

spatial (node) domain

$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$



graph spectral domain

$$\widehat{(f * g)}(\lambda) = ((\chi^T f) \circ \hat{g})(\lambda)$$

#### classical convolution

time domain

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



frequency domain

$$\widehat{(f * g)}(\omega) = \widehat{f}(\omega) \cdot \widehat{g}(\omega)$$

#### convolution on graphs

spatial (node) domain

$$f*g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$
 convolution = filtering



graph spectral domain

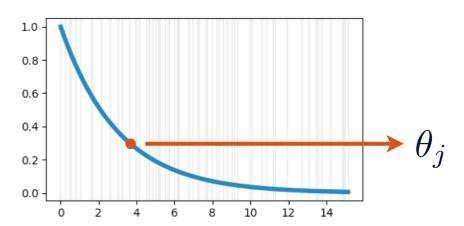
$$\widehat{(f * g)}(\lambda) = ((\chi^T f) \circ \hat{g})(\lambda)$$

$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$



learning a non-parametric filter:

$$\hat{g}_{\theta}(\Lambda) = \operatorname{diag}(\theta), \ \theta \in \mathbb{R}^{N}$$

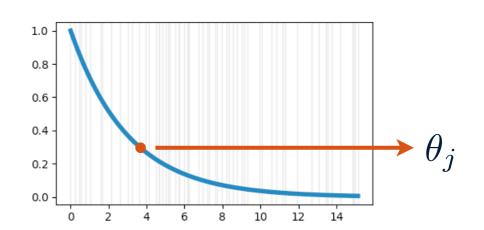


$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$



learning a non-parametric filter:

$$\hat{g}_{\theta}(\Lambda) = \operatorname{diag}(\theta), \ \theta \in \mathbb{R}^{N}$$



- convolution expressed in the graph spectral domain
- no localisation in the spatial (node) domain
- computationally expensive

$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$



parametric filter as polynomial of Laplacian

$$\hat{g}_{\theta}(\lambda) = \sum_{j=0}^{K} \theta_{j} \lambda^{j}, \ \theta \in \mathbb{R}^{K+1} \qquad \qquad \hat{g}_{\theta}(L) = \sum_{j=0}^{K} \theta_{j} L^{j}$$



$$\hat{g}_{\theta}(L) = \sum_{j=0}^{K} \theta_{j} L^{j}$$

$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$



parametric filter as polynomial of Laplacian

$$\hat{g}_{\theta}(\lambda) = \sum_{j=0}^{K} \theta_{j} \lambda^{j}, \ \theta \in \mathbb{R}^{K+1}$$
 
$$\hat{g}_{\theta}(L) = \sum_{j=0}^{K} \theta_{j} L^{j}$$

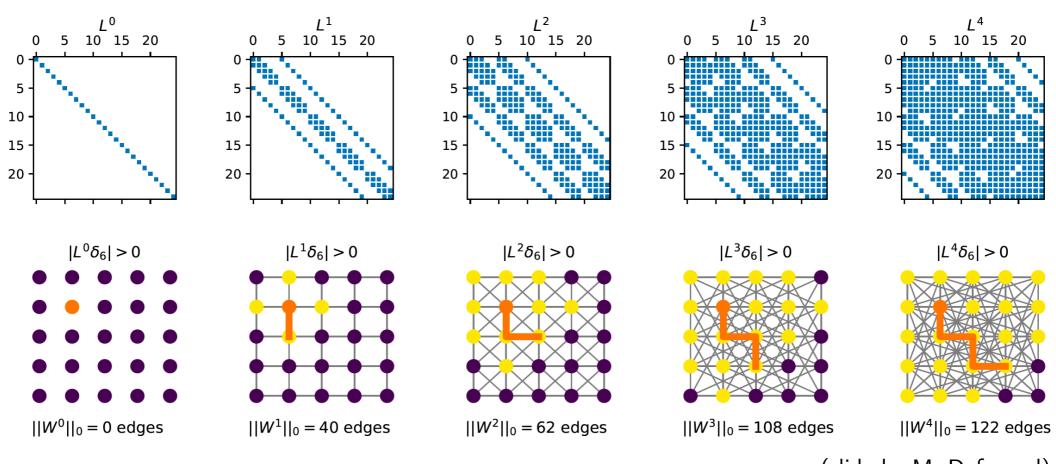


$$\hat{g}_{\theta}(L) = \sum_{j=0}^{K} \theta_{j} L^{j}$$

what do powers of graph Laplacian capture?

### Powers of graph Laplacian

#### $L^k$ defines the k-neighborhood



Localization:  $d_{\mathcal{G}}(v_i, v_j) > K$  implies  $(L^K)_{ij} = 0$ 

(slide by M. Deferrard)

$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$



parametric filter as polynomial of Laplacian

$$\hat{g}_{\theta}(\lambda) = \sum_{j=0}^{K} \theta_{j} \lambda^{j}, \ \theta \in \mathbb{R}^{K+1} \qquad \qquad \hat{g}_{\theta}(L) = \sum_{j=0}^{K} \theta_{j} L^{j}$$



$$\hat{g}_{\theta}(L) = \sum_{j=0}^{K} \theta_{j} L^{j}$$

$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$

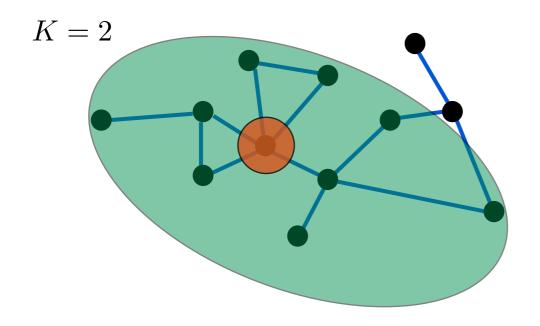


parametric filter as polynomial of Laplacian

$$\hat{g}_{\theta}(\lambda) = \sum_{j=0}^{K} \theta_{j} \lambda^{j}, \ \theta \in \mathbb{R}^{K+1}$$
 
$$\hat{g}_{\theta}(L) = \sum_{j=0}^{K} \theta_{j} L^{j}$$



$$\hat{g}_{\theta}(L) = \sum_{j=0}^{K} \theta_j L^j$$



localisation within K-hop neighbourhood

$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$

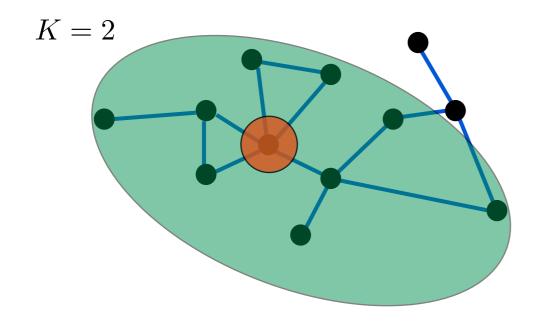


parametric filter as polynomial of Laplacian

$$\hat{g}_{\theta}(\lambda) = \sum_{j=0}^{K} \theta_j \lambda^j, \ \theta \in \mathbb{R}^{K+1} \qquad \qquad \hat{g}_{\theta}(L) = \sum_{j=0}^{K} \theta_j L^j$$



$$\hat{g}_{\theta}(L) = \sum_{j=0}^{K} \theta_j L^j$$



- localisation within K-hop neighbourhood
- Chebyshev approximation enables efficient computation via recursive multiplication with scaled Laplacian

$$\tilde{L} = \frac{2}{\lambda_{N-1}}L - I$$

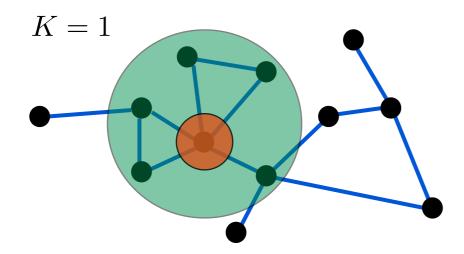
### A simplified parametric filter

$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$



simplified parametric filter

$$\hat{g}_{\theta}(L) = \sum_{j=0}^{K} \theta_j L^j$$



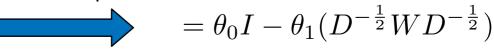
#### normalised Laplacian

$$L_{\text{norm}} = D^{-\frac{1}{2}} L D^{-\frac{1}{2}}$$

$$= D^{-\frac{1}{2}} (D - W) D^{-\frac{1}{2}}$$

$$= I - D^{-\frac{1}{2}} W D^{-\frac{1}{2}} = I - W_{\text{norm}}$$

$$K=1 \label{eq:K}$$
 normalised Laplacian



(localisation within 1-hop neighbourhood)

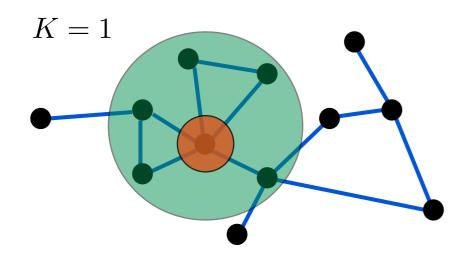
### A simplified parametric filter

$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$



simplified parametric filter

$$\hat{g}_{\theta}(L) = \sum_{j=0}^{K} \theta_j L^j$$



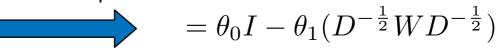
#### normalised Laplacian

$$L_{\text{norm}} = D^{-\frac{1}{2}} L D^{-\frac{1}{2}}$$

$$= D^{-\frac{1}{2}} (D - W) D^{-\frac{1}{2}}$$

$$= I - D^{-\frac{1}{2}} W D^{-\frac{1}{2}} = I - W_{\text{norm}}$$

$$K=1 \\ \label{eq:K}$$
 normalised Laplacian



(localisation within 1-hop neighbourhood)

$$\alpha = \theta_0 = -\theta_1$$

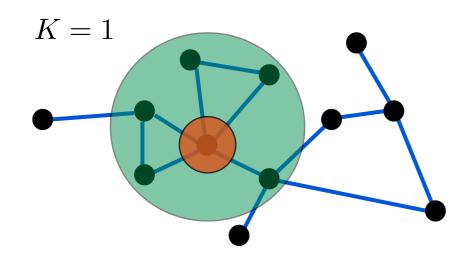
$$= \alpha (I + D^{-\frac{1}{2}} W D^{-\frac{1}{2}})$$

$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$



simplified parametric filter

$$\hat{g}_{\theta}(L) = \sum_{j=0}^{K} \theta_j L^j$$



#### normalised Laplacian

$$L_{\text{norm}} = D^{-\frac{1}{2}} L D^{-\frac{1}{2}}$$

$$= D^{-\frac{1}{2}} (D - W) D^{-\frac{1}{2}}$$

$$= I - D^{-\frac{1}{2}} W D^{-\frac{1}{2}} = I - W_{\text{norm}}$$

$$K=1 \label{eq:K}$$
 normalised Laplacian



$$= \theta_0 I - \theta_1 (D^{-\frac{1}{2}} W D^{-\frac{1}{2}})$$

(localisation within 1-hop neighbourhood)

$$\alpha = \theta_0 = -\theta_1$$



$$= \alpha (I + D^{-\frac{1}{2}} W D^{-\frac{1}{2}})$$

renormalisation



$$\Rightarrow \alpha(\tilde{D}^{-\frac{1}{2}}\tilde{W}\tilde{D}^{-\frac{1}{2}})$$

#### renormalisation

$$\tilde{W} = W + I$$
  $\tilde{D} = D + I$ 

$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$



simplified parametric filter

$$\hat{g}_{\alpha}(L) = \alpha(I + D^{-\frac{1}{2}}WD^{-\frac{1}{2}})$$

$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$

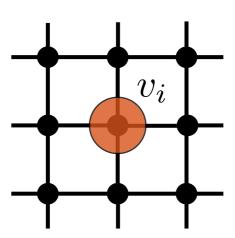


simplified parametric filter

$$\hat{g}_{\alpha}(L) = \alpha(I + D^{-\frac{1}{2}}WD^{-\frac{1}{2}})$$



$$y_i = \alpha f_i + \alpha \frac{1}{\sqrt{d_i}} \sum_{j:(i,j)\in\mathcal{E}} w_{ij} \frac{1}{\sqrt{d_j}} f_j$$



$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$



simplified parametric filter

$$\hat{g}_{\alpha}(L) = \alpha(I + D^{-\frac{1}{2}}WD^{-\frac{1}{2}})$$

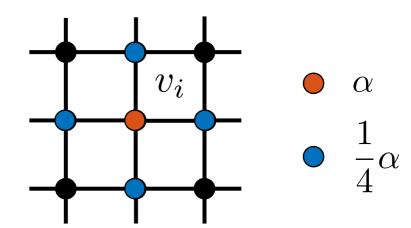


$$y_i = \alpha f_i + \alpha \frac{1}{\sqrt{d_i}} \sum_{j:(i,j)\in\mathcal{E}} w_{ij} \frac{1}{\sqrt{d_j}} f_j$$



unitary edge weights

$$y_i = \alpha f_i + \frac{1}{4} \alpha \sum_{j:(i,j)\in\mathcal{E}} f_j$$



$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$



simplified parametric filter

$$\hat{g}_{\alpha}(L) = \alpha(I + D^{-\frac{1}{2}}WD^{-\frac{1}{2}})$$



$$y_i = \alpha f_i + \alpha \frac{1}{\sqrt{d_i}} \sum_{j:(i,j)\in\mathcal{E}} w_{ij} \frac{1}{\sqrt{d_j}} f_j$$

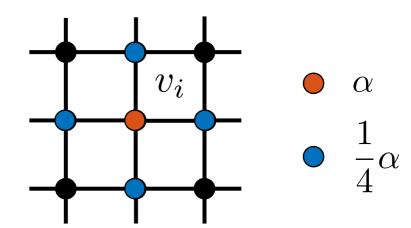


unitary edge weights

$$y_i = \alpha f_i + \frac{1}{4} \alpha \sum_{j:(i,j)\in\mathcal{E}} f_j$$

30	3,	$2_{2}$	1	0
$0_2$	02	$1_0$	3	1
30	1,	22	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0



Convolution is defined via the graph spectral domain..

$$f * g = \chi \hat{g}(\Lambda) \chi^T f = \hat{g}(L) f$$

- ..but can be implemented in the spatial (node) domain
  - simplified filter:  $y=\hat{g}_{\theta}(L)f=lpha( ilde{D}^{-\frac{1}{2}} ilde{W} ilde{D}^{-\frac{1}{2}})f$
  - interpretation: at each layer nodes exchange information in 1-hop neighbourhood
  - more generally: receptive field size determined by degree of polynomial

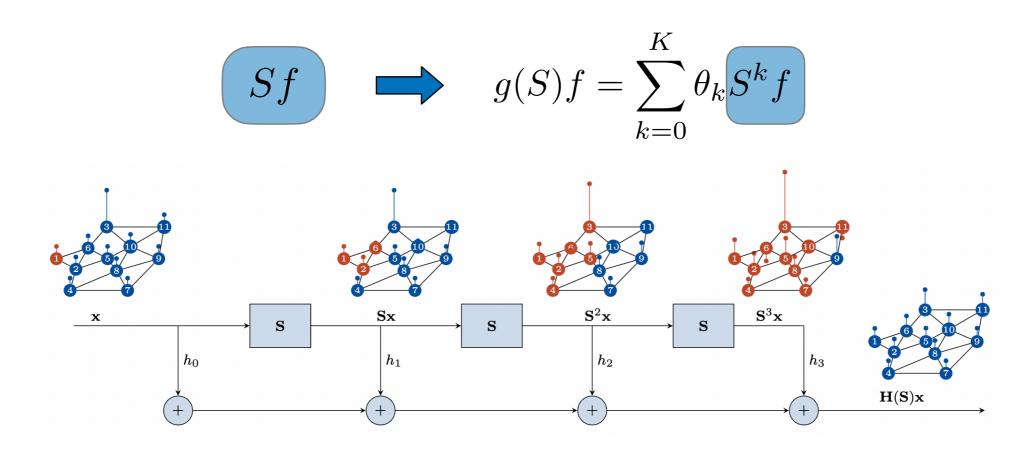
Convolution in classical signal processing relies on the shift operator

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

Convolution in classical signal processing relies on the shift operator

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

Idea: convolution via a graph shift operator (e.g., adjacency matrix)



Convolution can also be interpreted as a weighted summation

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

Convolution can also be interpreted as a weighted summation

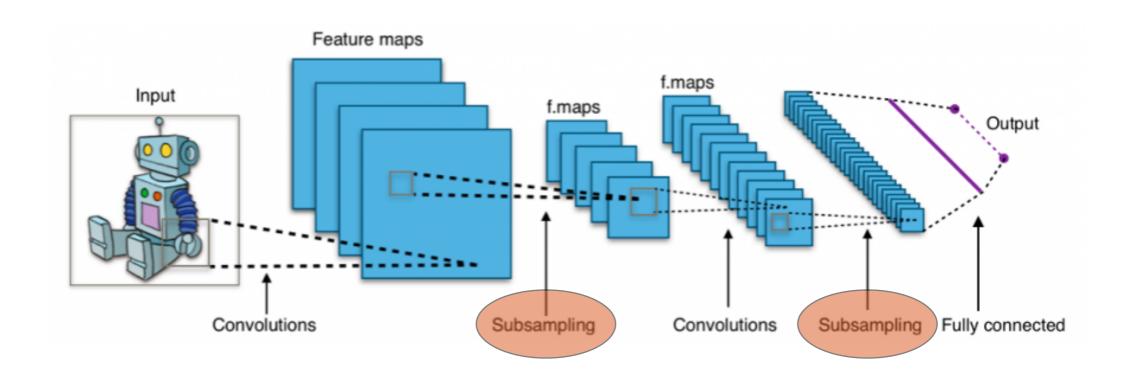
$$(f * g)(t) = \int_{-\infty}^{\infty} f(\tau) g(t - \tau) d\tau$$

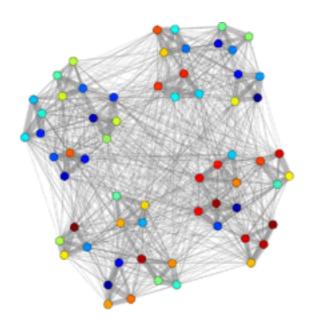
Idea: convolution via a spatial weighted summation in graph domain

$$(f * g)(v) = \sum_{v' \in \mathcal{V}} f(v') g(v, v')$$

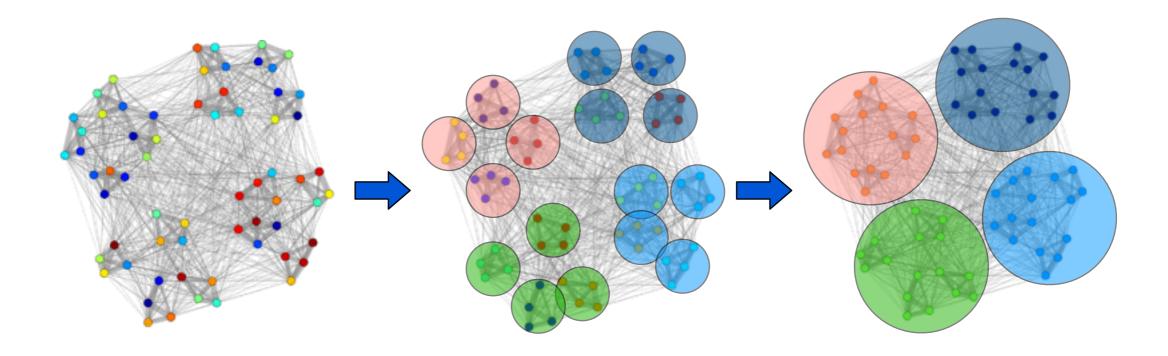


relative importance of each node to v (can be designed to achieve locality)

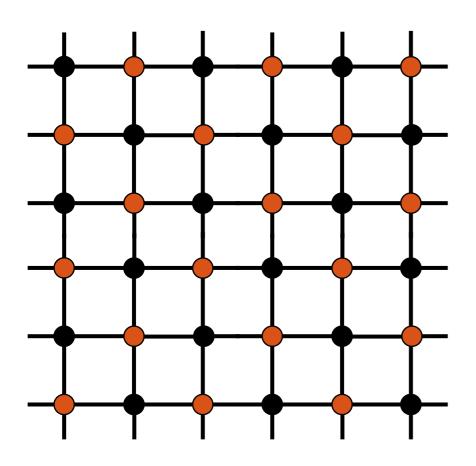




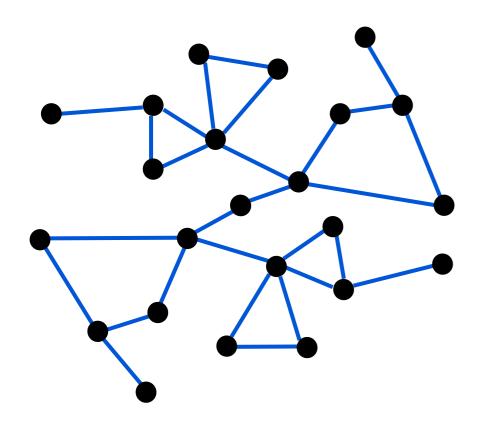
• pooling = downsampling on graphs, but how?



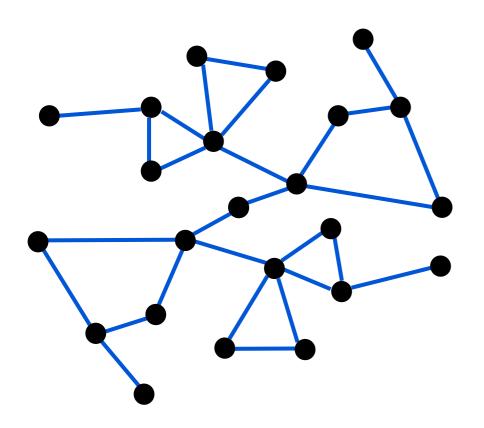
- pooling = downsampling on graphs, but how?
- natural idea: graph coarsening



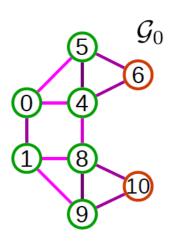
coarsening is straightforward on regular grids



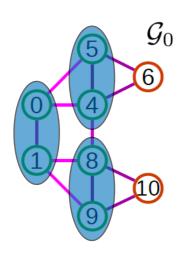
- coarsening is straightforward on regular grids
- not so much on irregular graphs



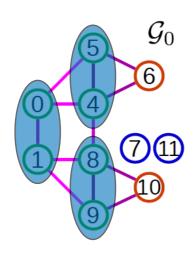
- coarsening is straightforward on regular grids
- not so much on irregular graphs
- can be achieved via node clustering
  - multi-level partitioning
  - roughly fixed downsampling factor (e.g., 2)
  - need for efficiency



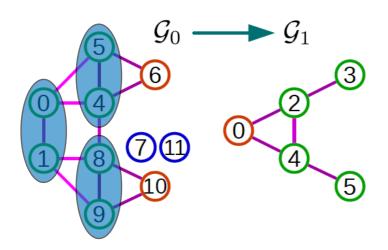
- pooling based on Graclus algorithm
  - local greedy way of merging vertices: maximising  $w_{ij}(1/d_i+1/d_j)$
  - adding artificial vertices to ensure two children for each node



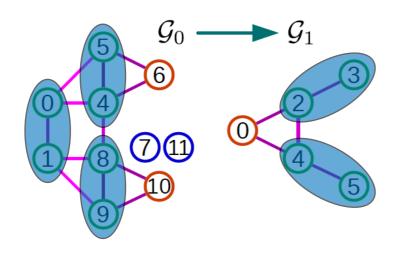
- pooling based on Graclus algorithm
  - local greedy way of merging vertices: maximising  $w_{ij}(1/d_i+1/d_j)$
  - adding artificial vertices to ensure two children for each node



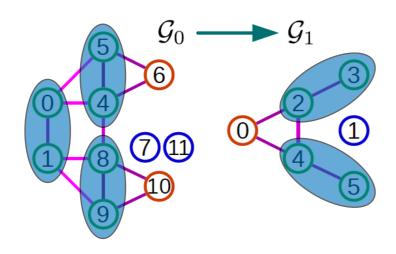
- pooling based on Graclus algorithm
  - local greedy way of merging vertices: maximising  $w_{ij}(1/d_i+1/d_j)$
  - adding artificial vertices to ensure two children for each node



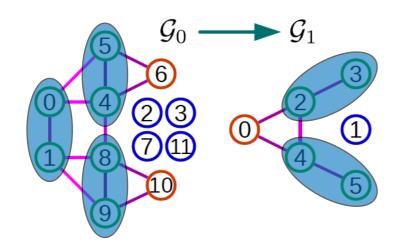
- pooling based on Graclus algorithm
  - local greedy way of merging vertices: maximising  $w_{ij}(1/d_i+1/d_j)$
  - adding artificial vertices to ensure two children for each node



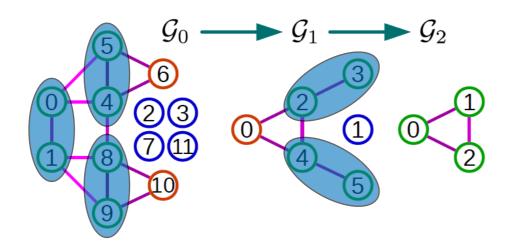
- pooling based on Graclus algorithm
  - local greedy way of merging vertices: maximising  $w_{ij}(1/d_i+1/d_j)$
  - adding artificial vertices to ensure two children for each node



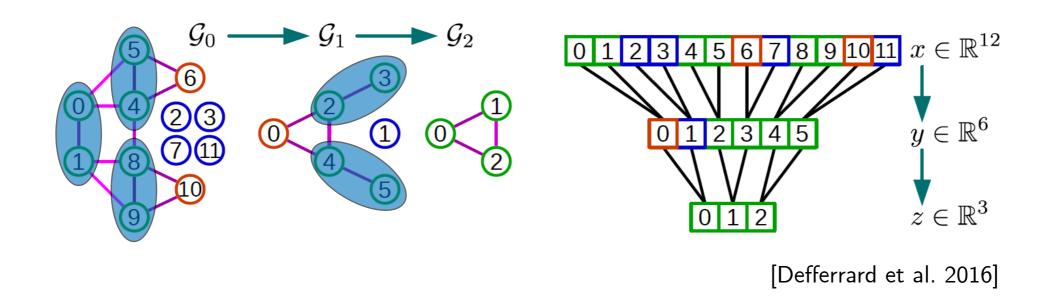
- pooling based on Graclus algorithm
  - local greedy way of merging vertices: maximising  $w_{ij}(1/d_i+1/d_j)$
  - adding artificial vertices to ensure two children for each node



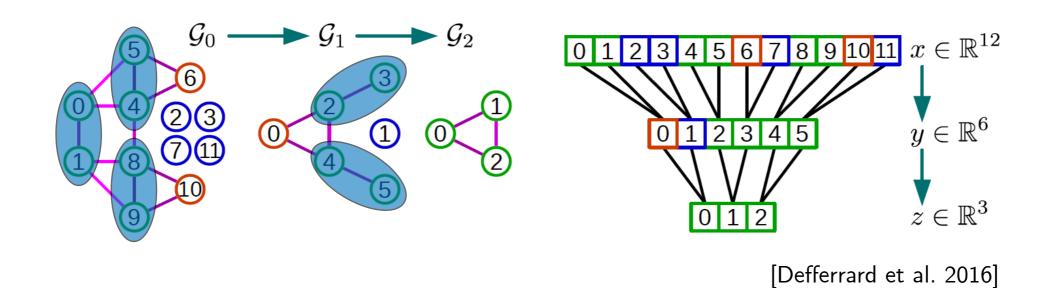
- pooling based on Graclus algorithm
  - local greedy way of merging vertices: maximising  $w_{ij}(1/d_i+1/d_j)$
  - adding artificial vertices to ensure two children for each node



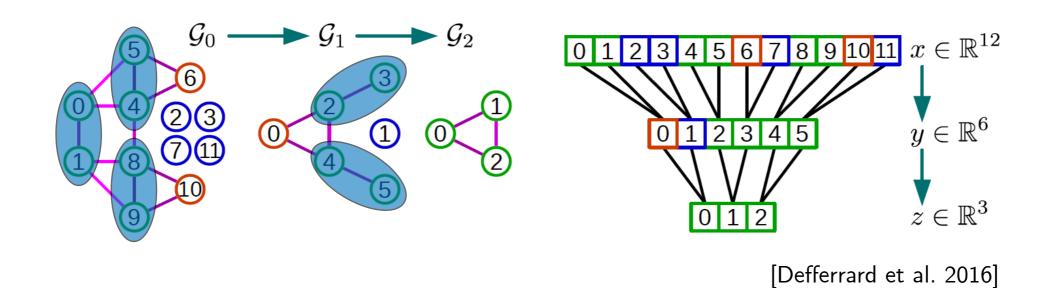
- pooling based on Graclus algorithm
  - local greedy way of merging vertices: maximising  $w_{ij}(1/d_i+1/d_j)$
  - adding artificial vertices to ensure two children for each node



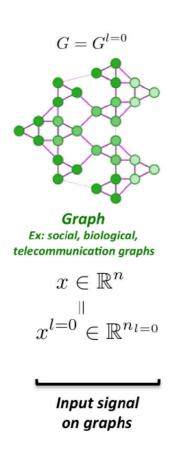
- pooling based on Graclus algorithm
  - local greedy way of merging vertices: maximising  $w_{ij}(1/d_i+1/d_j)$
  - adding artificial vertices to ensure two children for each node

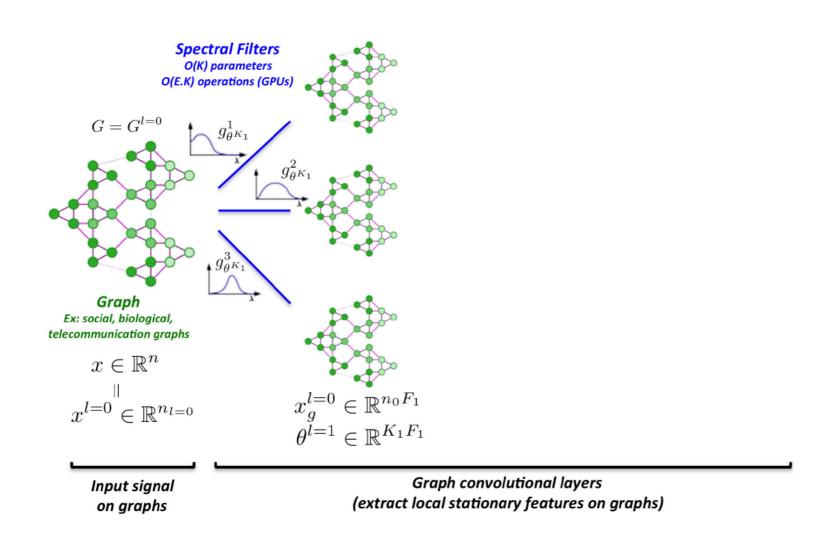


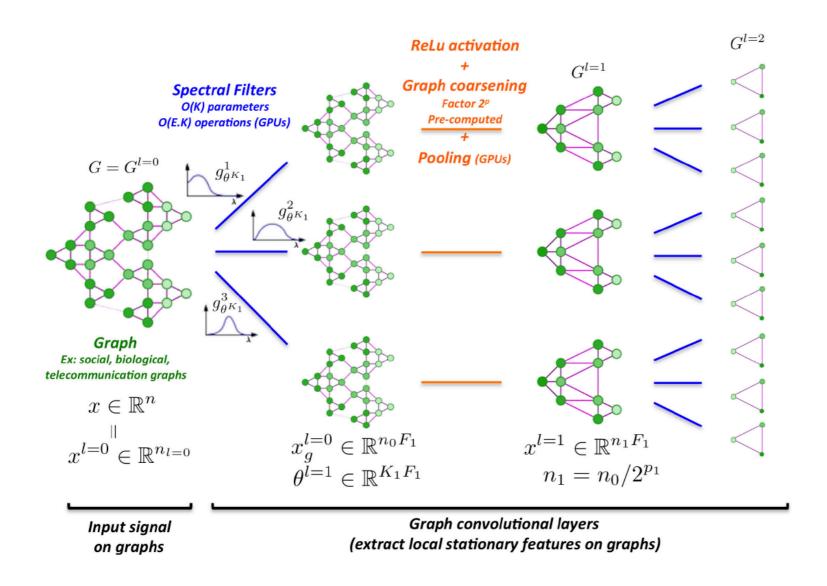
- pooling based on Graclus algorithm
  - local greedy way of merging vertices: maximising  $w_{ij}(1/d_i+1/d_j)$
  - adding artificial vertices to ensure two children for each node
  - 1D grid pooling: [ max(0,1) max(4,5,6) max(8,9,10) ]

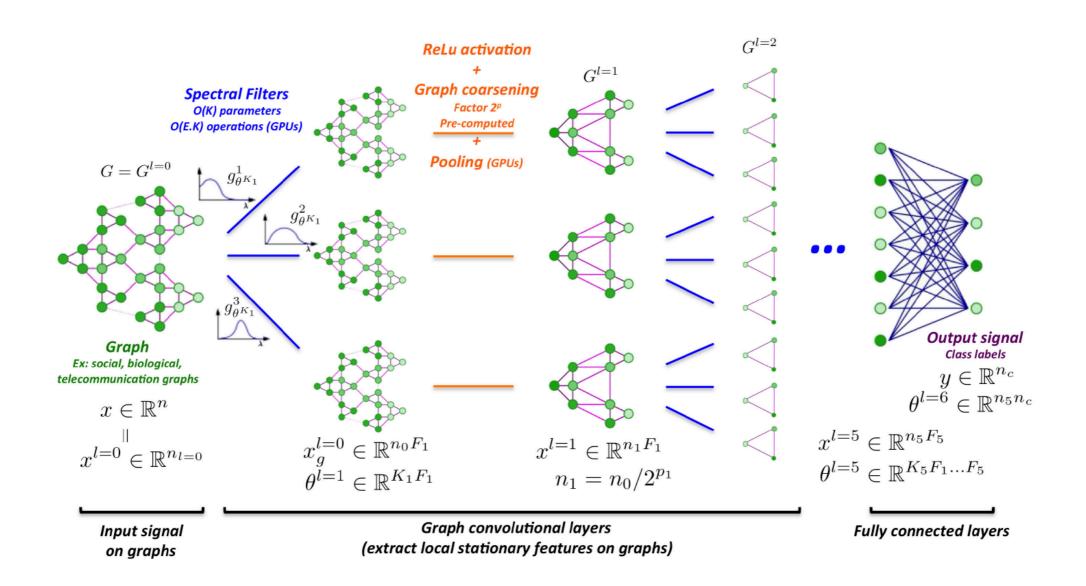


- pooling based on Graclus algorithm
  - local greedy way of merging vertices: maximising  $w_{ij}(1/d_i+1/d_j)$
  - adding artificial vertices to ensure two children for each node
  - 1D grid pooling: [ max(0,1) max(4,5,6) max(8,9,10) ]
  - only based on graph (and no signal) information

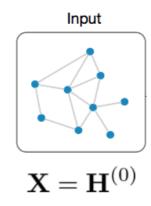




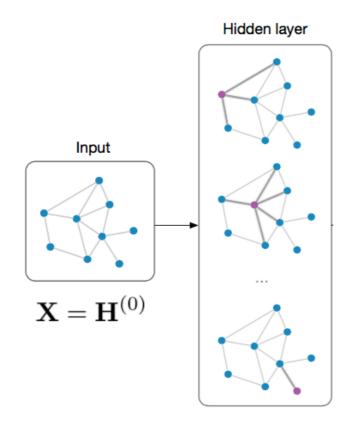




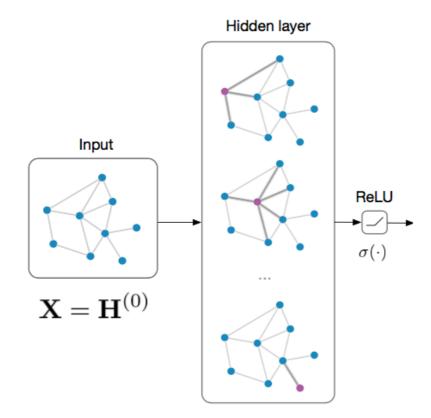
$$\hat{g}_{\theta^{(k+1)}}(L)\Big(\operatorname{ReLU}(\hat{g}_{\theta^{(k)}}(L)f)\Big)$$

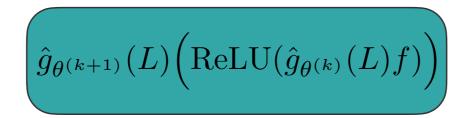


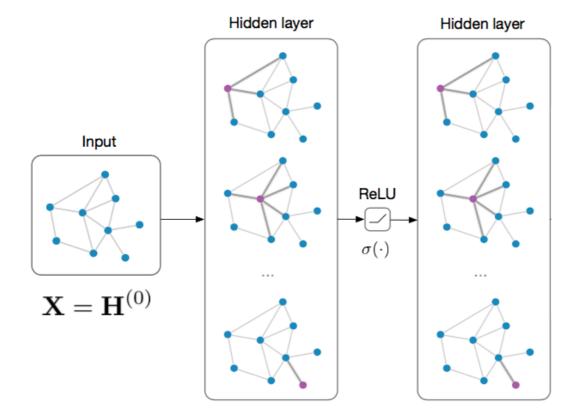
$$\hat{g}_{\theta^{(k+1)}}(L)\Big(\mathrm{ReLU}(\hat{g}_{\theta^{(k)}}(L)f)\Big)$$

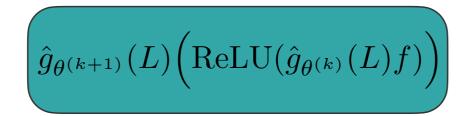


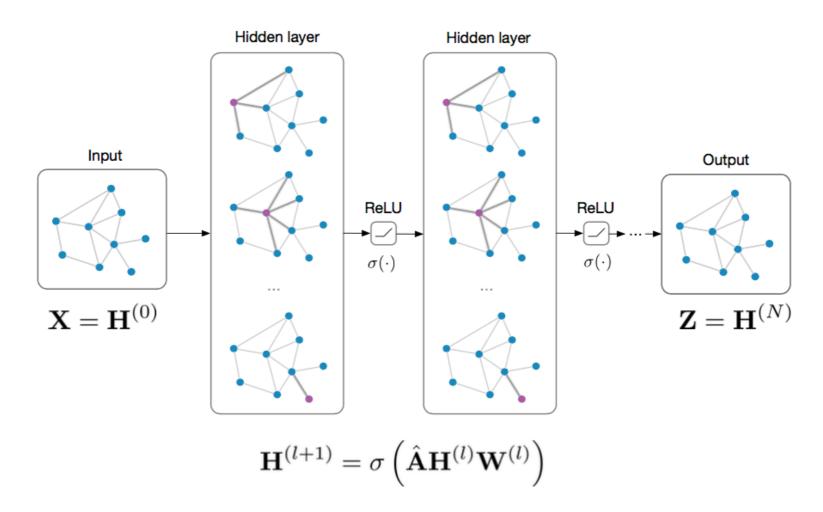








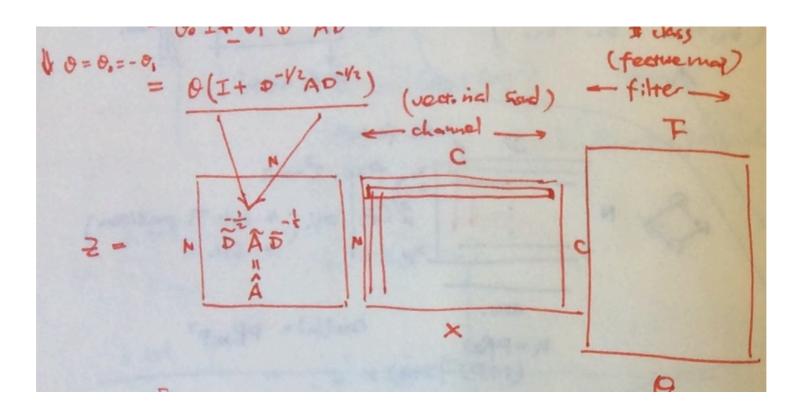




# CNNs on graphs: Node classification

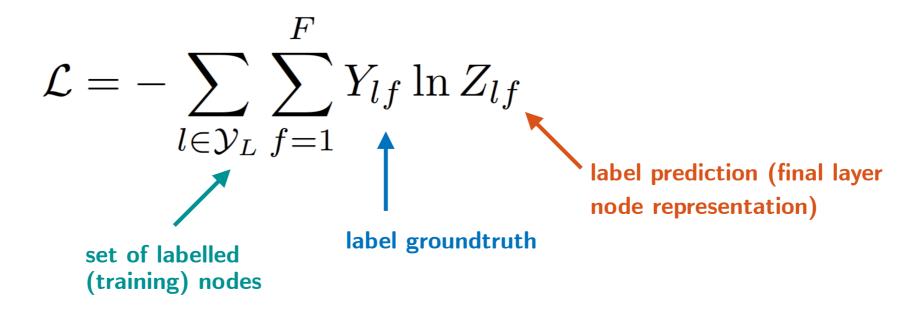
GCN architecture

$$\left(\hat{g}_{\theta^{(k+1)}}(L)\Big(\mathrm{ReLU}(\hat{g}_{\theta^{(k)}}(L)f)\Big)\right)$$



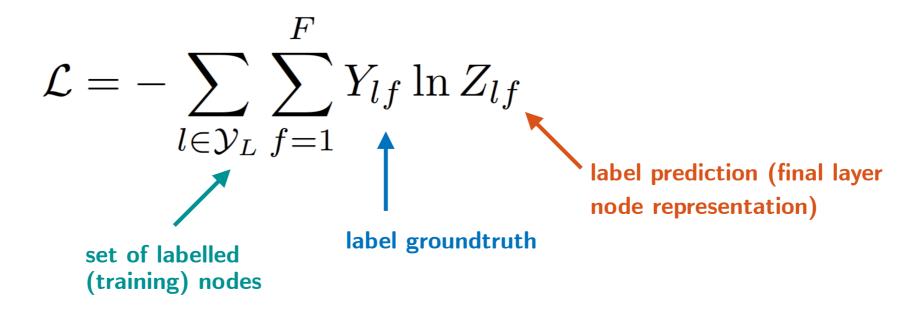
$$\mathbf{H}^{(l+1)} = \sigma \left( \hat{\mathbf{A}} \mathbf{H}^{(l)} \mathbf{W}^{(l)} \right)$$

- Node-level task
  - cross-entropy loss function for (semi-supervised) node classification



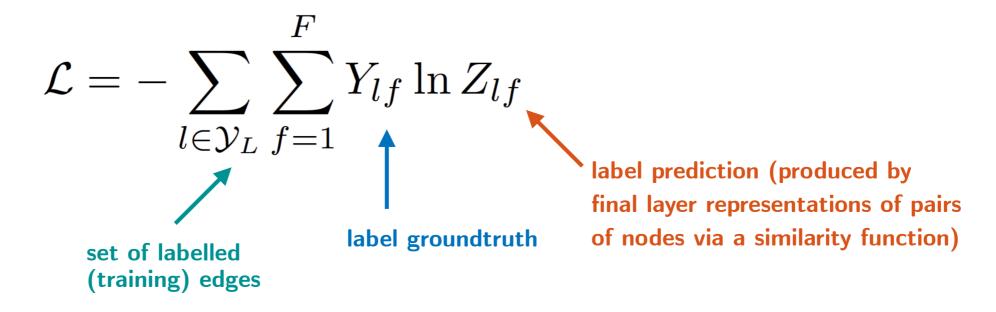
- training by minimising loss function and making predictions on testing nodes

- Node-level task
  - cross-entropy loss function for (semi-supervised) node classification



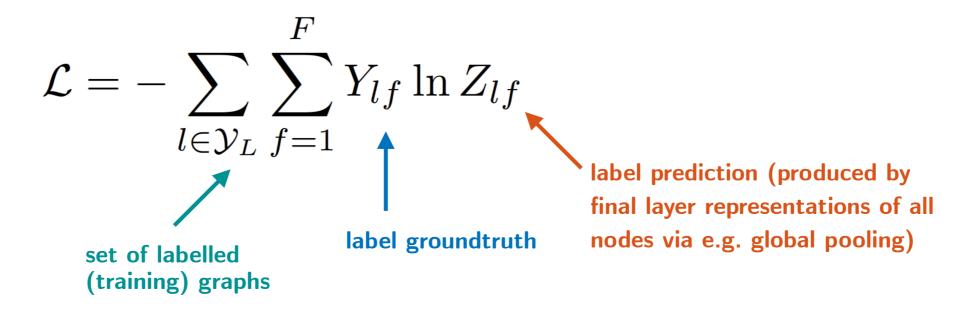
- training by minimising loss function and making predictions on testing nodes
- Factors influencing model behaviour
  - what label distribution favours GCN in this task?
  - what about perturbation of input graph topology?

- Edge-level task
  - cross-entropy loss function for link prediction



- training by minimising loss function and making predictions on testing edges

- Graph-level task
  - cross-entropy loss function for graph classification

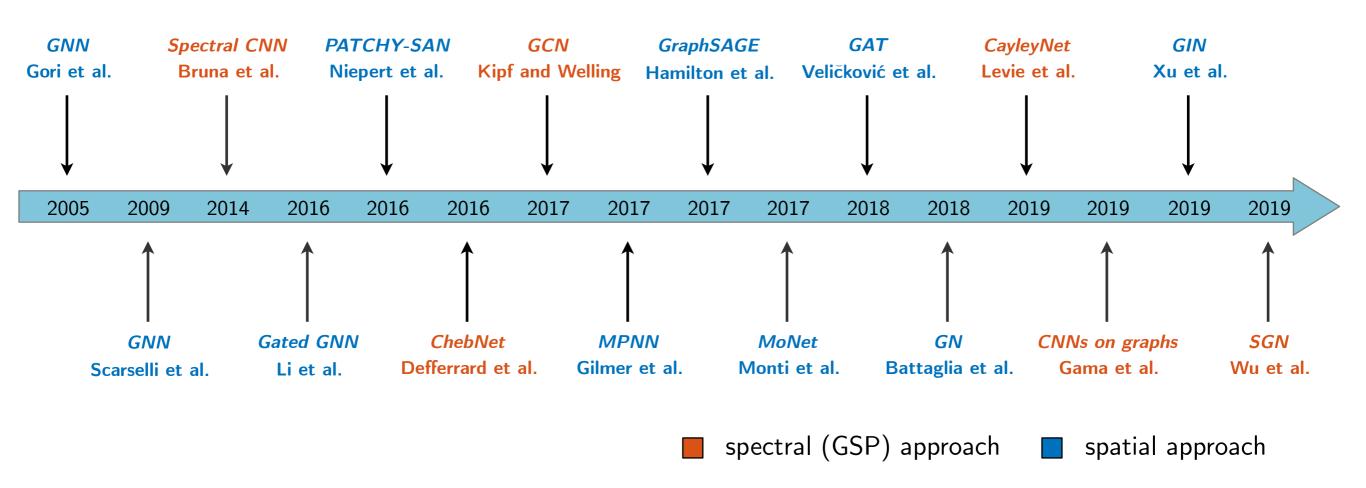


- training by minimising loss function and making predictions on testing graphs

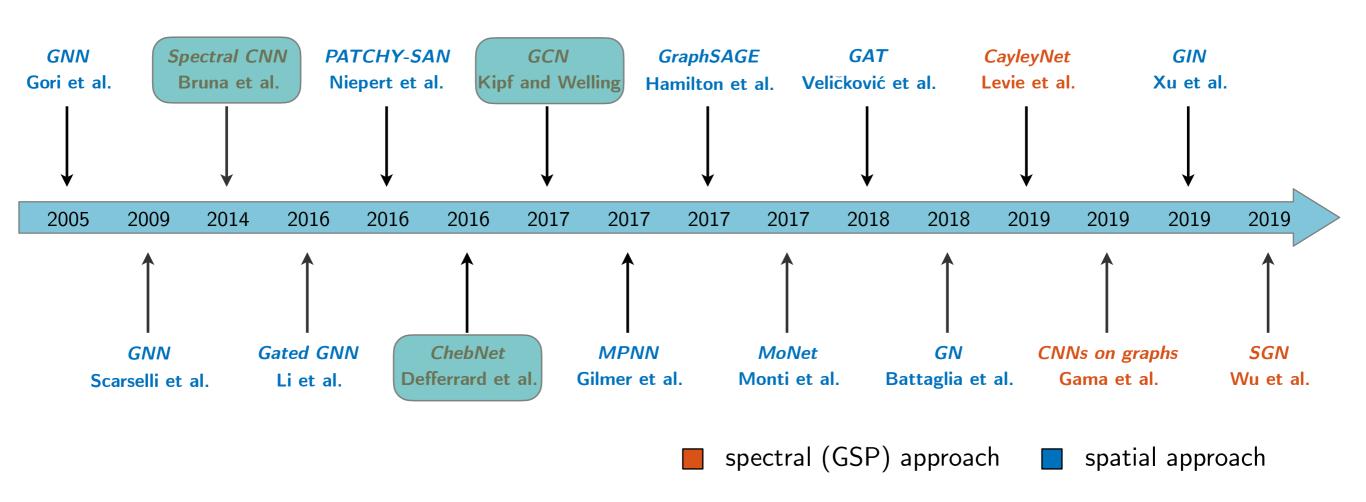
#### Lecture 3

- Machine learning on graphs: Overview
- Convolutional neural networks on graphs
- Message passing neural networks
- Recent developments & Applications

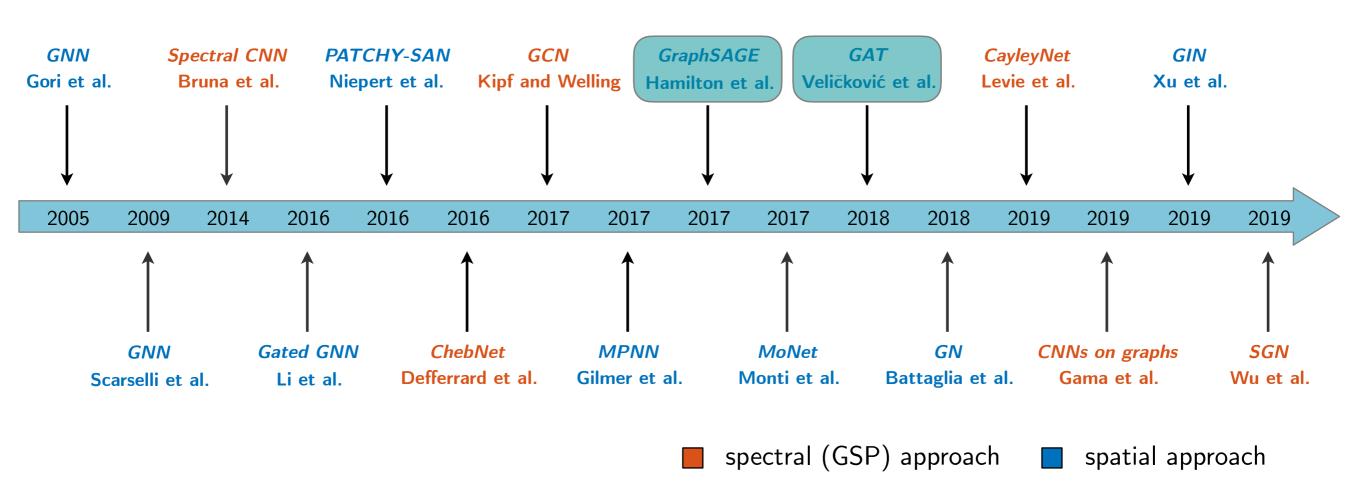
### GNNs - A historical timeline



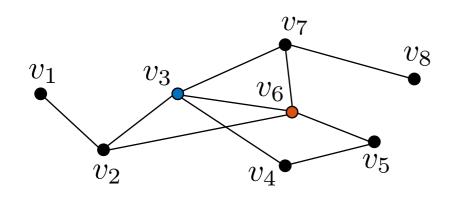
### GNNs - A historical timeline

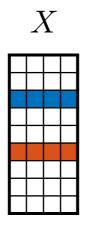


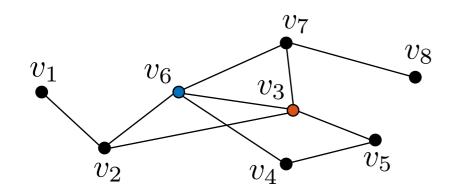
### GNNs - A historical timeline

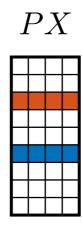


Permutation invariance: function invariant w.r.t. permutation (node re-ordering)

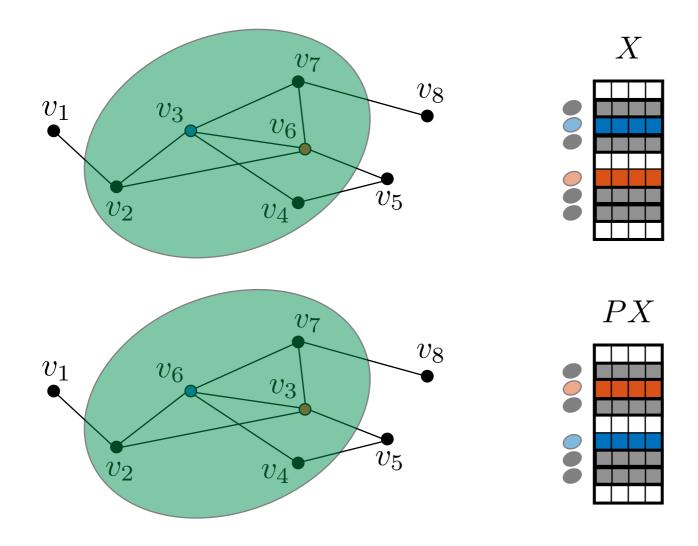




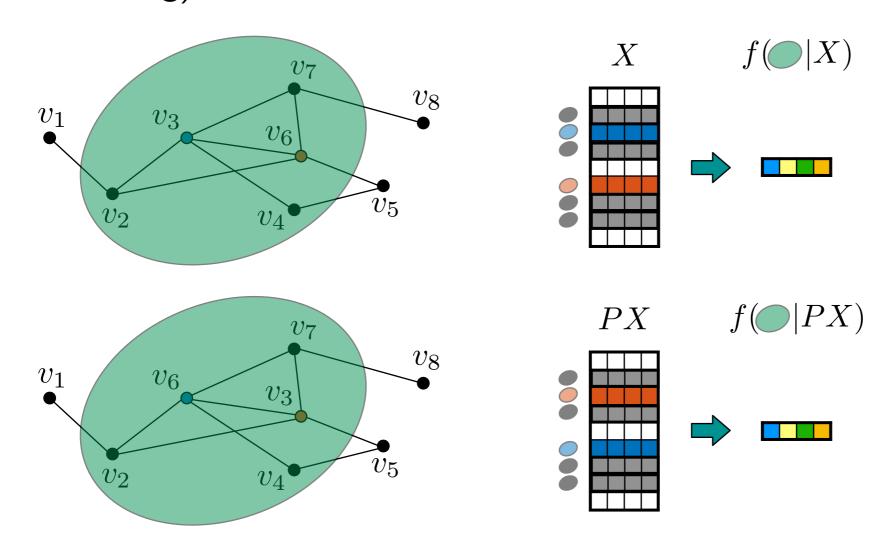




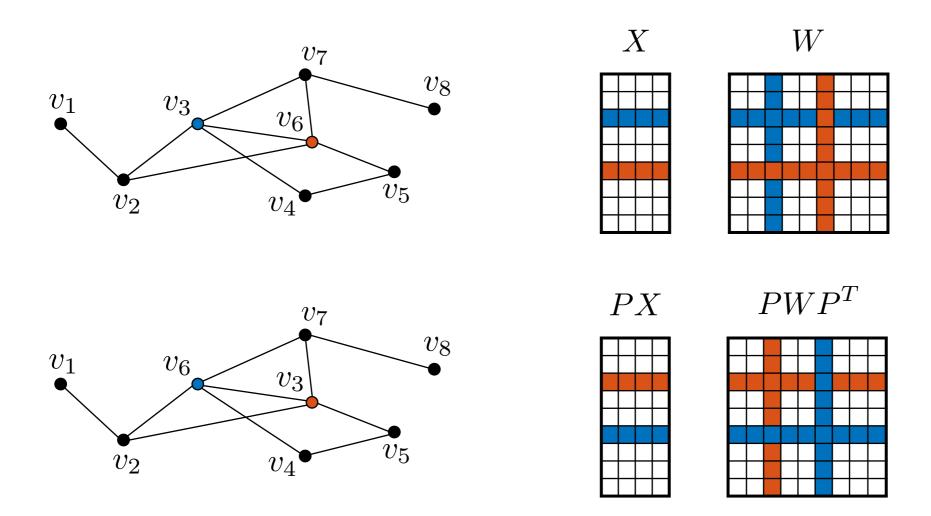
Permutation invariance: function invariant w.r.t. permutation (node re-ordering)



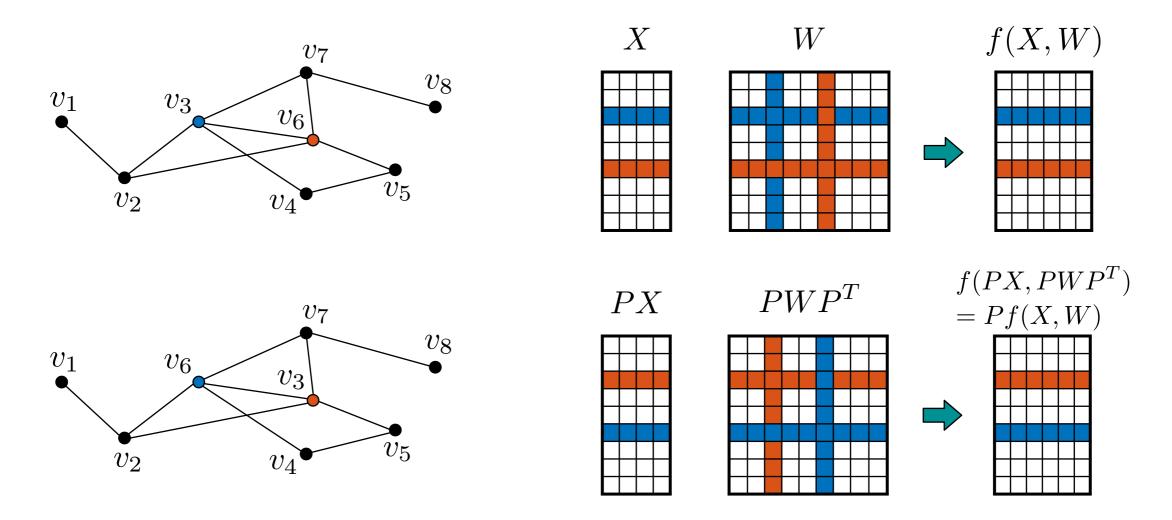
Permutation invariance: function invariant w.r.t. permutation (node re-ordering)



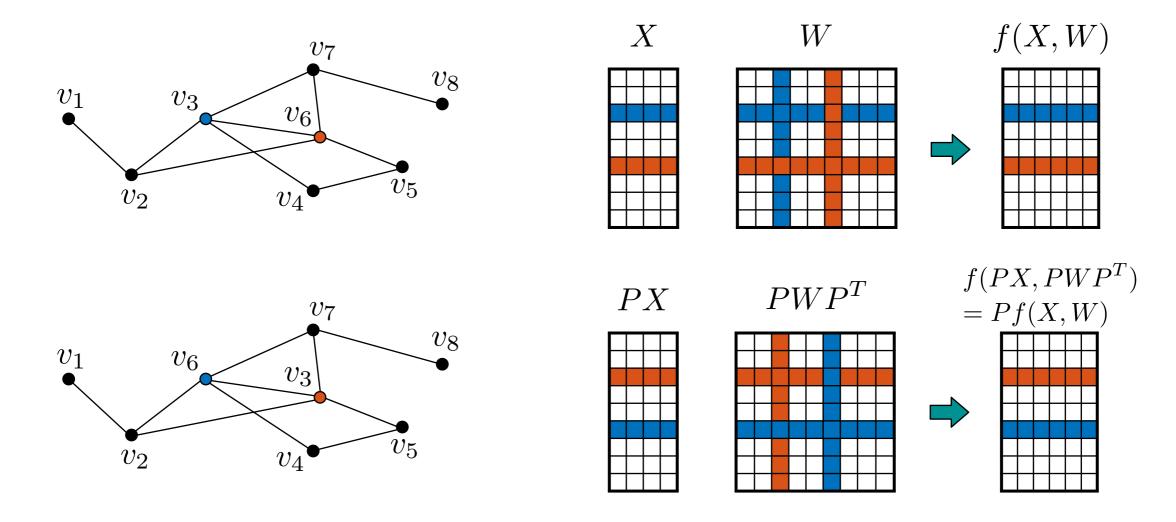
• Permutation equivariance: function equivariant w.r.t. permutation (permutation of input  $\Rightarrow$  same permutation of output)



• Permutation equivariance: function equivariant w.r.t. permutation (permutation of input  $\Rightarrow$  same permutation of output)



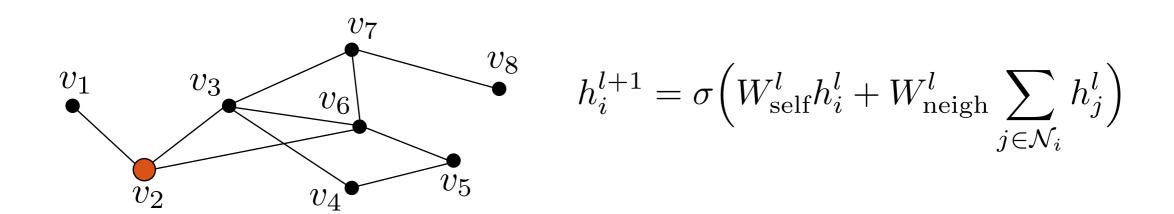
• Permutation equivariance: function equivariant w.r.t. permutation (permutation of input  $\Rightarrow$  same permutation of output)



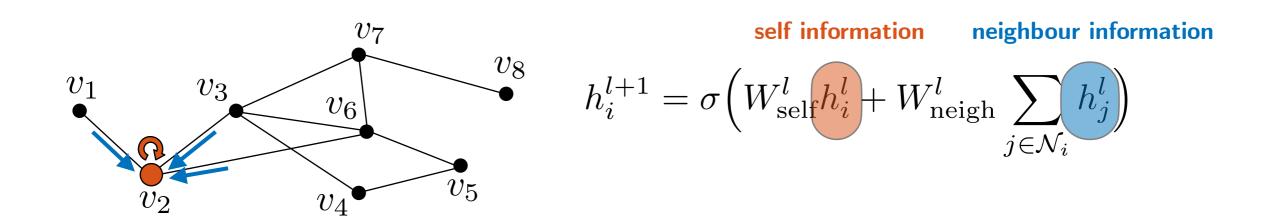
• Spectral GNNs:  $\hat{g}_{\theta}(L)$  is permutation equivariant because it acts on local neighbourhood and its behaviour is permutation invariant

Recall graph convolution can also be defined in spatial (node) domain

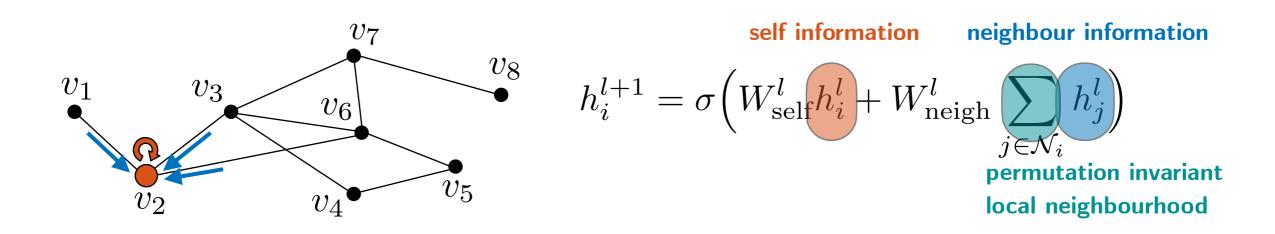
Recall graph convolution can also be defined in spatial (node) domain



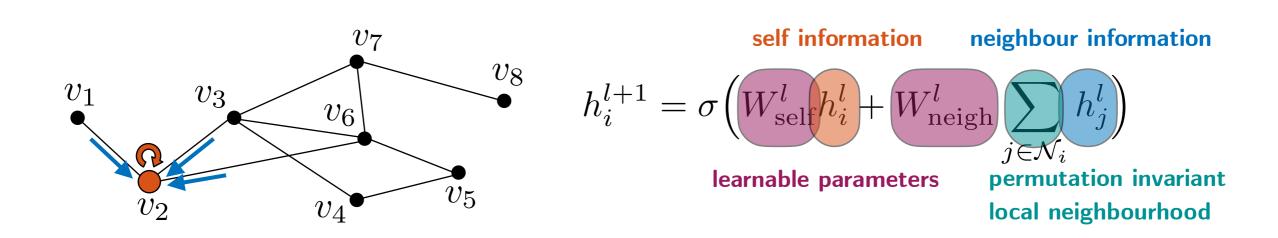
Recall graph convolution can also be defined in spatial (node) domain

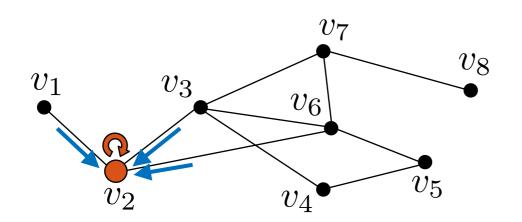


Recall graph convolution can also be defined in spatial (node) domain

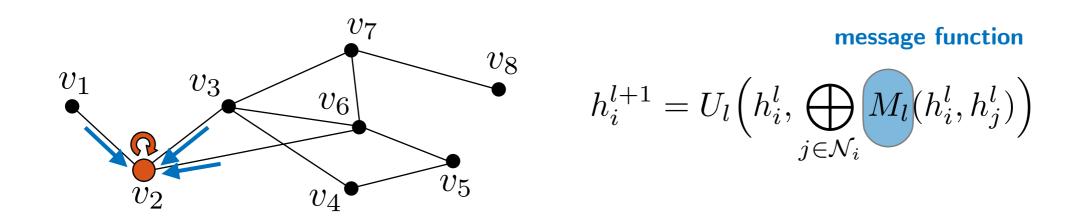


Recall graph convolution can also be defined in spatial (node) domain

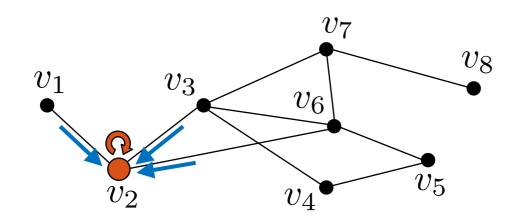




$$h_i^{l+1} = U_l \left( h_i^l, \bigoplus_{j \in \mathcal{N}_i} M_l(h_i^l, h_j^l) \right)$$



- nodes exchange **messages** with local neighbours



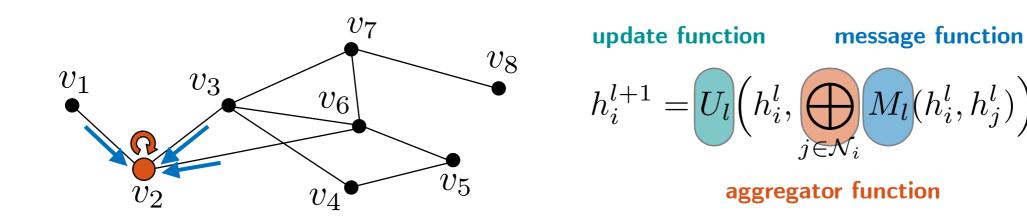
update function

message function

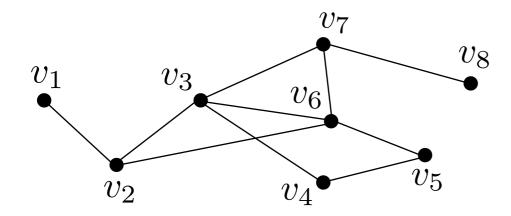
$$h_i^{l+1} = \underbrace{U_l} \Big( h_i^l, \underbrace{M_l} (h_i^l, h_j^l) \Big)$$

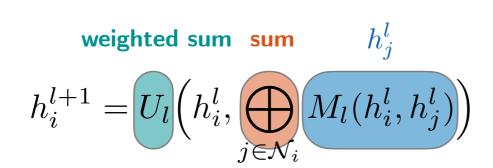
aggregator function

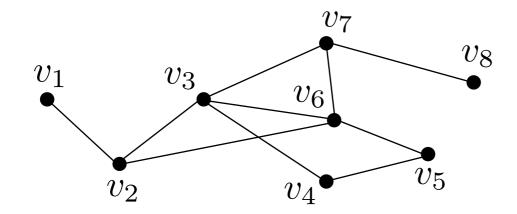
- nodes exchange **messages** with local neighbours
- each node updates its representation by aggregating messages from neighbours

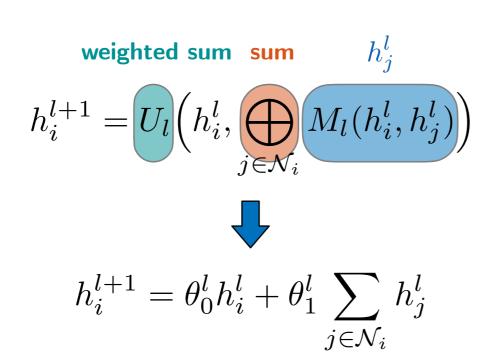


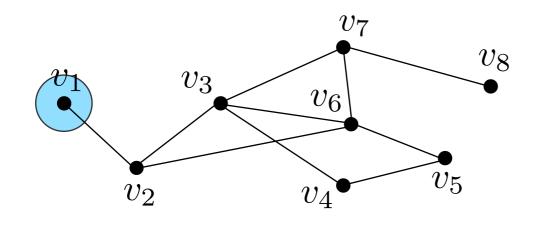
- nodes exchange messages with local neighbours
- each node updates its representation by aggregating messages from neighbours
- functions are differentiable and parameters are learned by minimising loss of downstream task
- key difference between architectures: how nodes aggregate information from neighbours and across layers







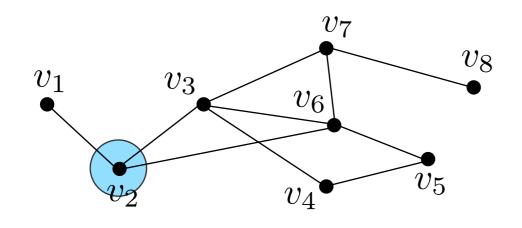




#### one message passing layer

$$h_1^{l+1} = \theta_0^l h_1^l + \theta_1^l h_2^l$$

weighted sum sum 
$$h_j^l$$
 
$$h_i^{l+1} = U_l \left( h_i^l, \bigoplus_{j \in \mathcal{N}_i} M_l(h_i^l, h_j^l) \right)$$
 
$$\downarrow h_i^{l+1} = \theta_0^l h_i^l + \theta_1^l \sum_{i \in \mathcal{N}_i} h_j^l$$

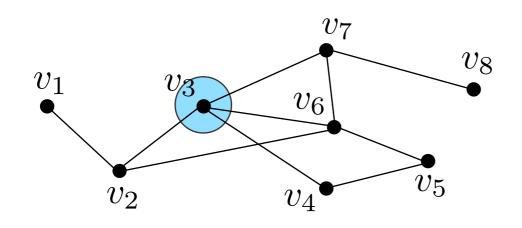


one message passing layer

$$h_1^{l+1} = \theta_0^l h_1^l + \theta_1^l h_2^l$$

$$h_2^{l+1} = \theta_0^l h_2^l + \theta_1^l (h_1^l + h_3^l + h_6^l)$$

weighted sum sum 
$$h_j^l$$
 
$$h_i^{l+1} = U_l \left( h_i^l, \bigvee_{j \in \mathcal{N}_i} M_l (h_i^l, h_j^l) \right)$$
 
$$h_i^{l+1} = \theta_0^l h_i^l + \theta_1^l \sum_{i \in \mathcal{N}_i} h_j^l$$



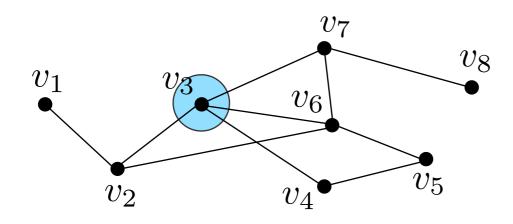
#### one message passing layer

$$h_1^{l+1} = \theta_0^l h_1^l + \theta_1^l h_2^l$$

$$h_2^{l+1} = \theta_0^l h_2^l + \theta_1^l (h_1^l + h_3^l + h_6^l)$$

$$h_3^{l+1} = \theta_0^l h_3^l + \theta_1^l (h_2^l + h_4^l + h_6^l + h_7^l)$$

$$\begin{aligned} \text{weighted sum} & \quad \text{sum} \quad \quad h_j^l \\ h_i^{l+1} = & U_l \Big( h_i^l, \bigoplus_{j \in \mathcal{N}_i} M_l (h_i^l, h_j^l) \Big) \\ & \qquad \qquad \qquad \\ h_i^{l+1} = & \quad \theta_0^l h_i^l + \theta_1^l \sum_{i \in \mathcal{N}_i} h_j^l \end{aligned}$$



#### weighted sum sum

$$h_i^{l+1} = U_l \left( h_i^l, \underbrace{M_l(h_i^l, h_j^l)} \right)$$



$$h_i^{l+1} = \theta_0^l h_i^l + \theta_1^l \sum_{j \in \mathcal{N}_i} h_j^l$$

#### one message passing layer

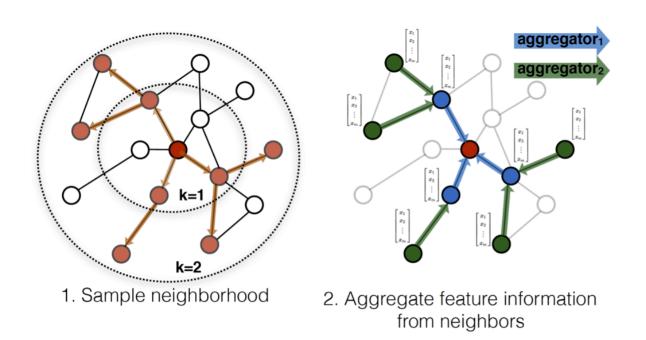
$$\begin{split} h_1^{l+1} &= \theta_0^l h_1^l + \theta_1^l h_2^l \\ h_2^{l+1} &= \theta_0^l h_2^l + \theta_1^l (h_1^l + h_3^l + h_6^l) \\ \hline h_3^{l+1} &= \theta_0^l h_3^l + \theta_1^l (h_2^l + h_4^l + h_6^l + h_7^l) \\ \hline h_4^{l+1} &= \theta_0^l h_3^l + \theta_1^l (h_2^l + h_4^l + h_6^l + h_7^l) \\ \hline h_4^{l+1} &= \theta_0^l h_4^l + \theta_1^l (h_3^l + h_5^l) \\ \hline h_5^{l+1} &= \theta_0^l h_5^l + \theta_1^l (h_4^l + h_6^l) \\ \hline h_6^{l+1} &= \theta_0^l h_6^l + \theta_1^l (h_2^l + h_3^l + h_5^l + h_7^l) \\ \hline h_7^{l+1} &= \theta_0^l h_7^l + \theta_1^l (h_3^l + h_6^l + h_8^l) \\ \hline h_8^{l+1} &= \theta_0^l h_8^l + \theta_1^l h_7^l \end{split}$$

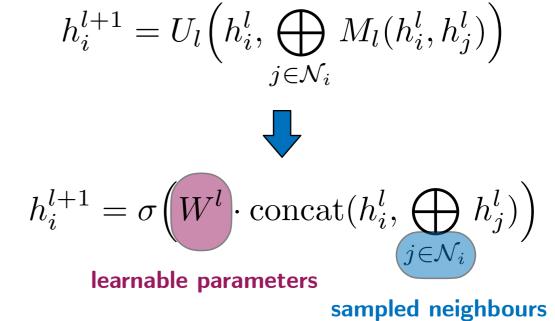
#### final output after L layers

$$Z = \{h_1^L, h_2^L, h_3^L, h_4^L, h_5^L, h_6^L, h_7^L, h_8^L\}^T$$

# GraphSAGE

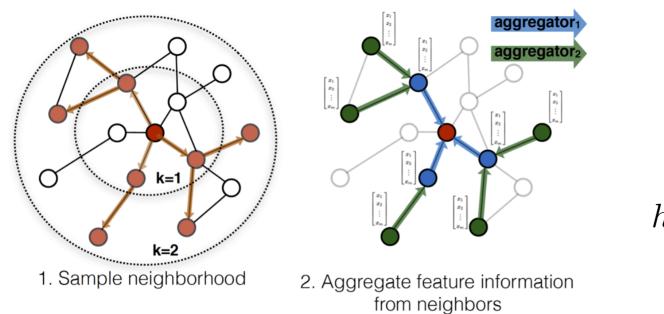
 Idea: sample neighbours from each hop of neighbourhood (i.e., building a "computational" graph) for improved scalability and robustness

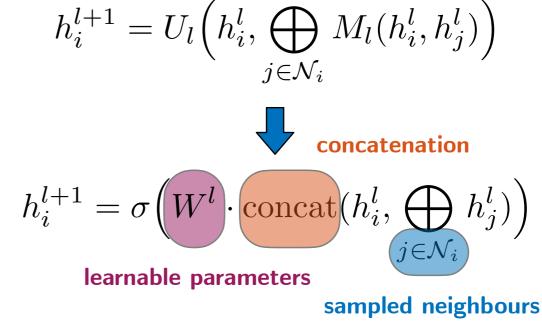




# GraphSAGE

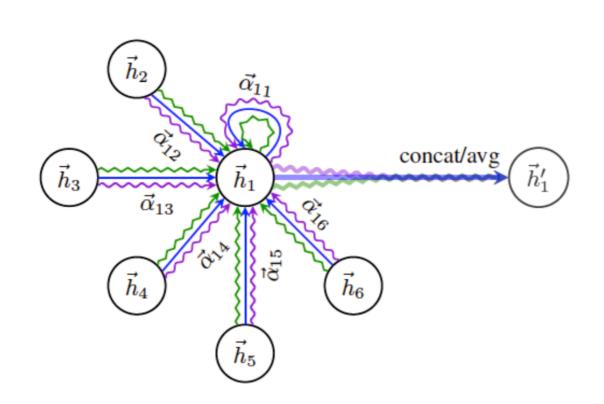
- Idea: sample neighbours from each hop of neighbourhood (i.e., building a "computational" graph) for improved scalability and robustness
- Concatenation of self and neighbour embeddings acts as implicit "skip connections" to prevent loss of self information





# Graph attention network (GAT)

• Idea: learn relative importance of neighbours in aggregation

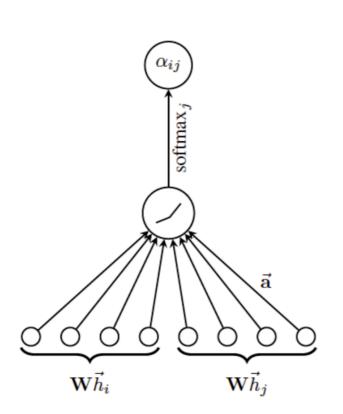


$$h_i^{l+1} = U_l\Big(h_i^l, \bigoplus_{j \in \mathcal{N}_i} M_l(h_i^l, h_j^l)\Big)$$
 relative importance 
$$h_i^{l+1} = \sigma\Big(\alpha_{il}W^lh_i^l + \sum_{i \in \mathcal{N}_i} \alpha_{ij}W^lh_j^l\Big)$$

learnable parameters

# Graph attention network (GAT)

- Idea: learn relative importance of neighbours in aggregation
- An attention function weighs importance of neighbours and computes attention scores



$$h_i^{l+1} = U_l\Big(h_i^l, \bigoplus_{j \in \mathcal{N}_i} M_l(h_i^l, h_j^l)\Big)$$
 relative importance 
$$h_i^{l+1} = \sigma\Big(\alpha_{ij}W^lh_i^l + \sum_{j \in \mathcal{N}_i} \alpha_{ij}W^lh_j^l\Big)$$
 learnable parameters

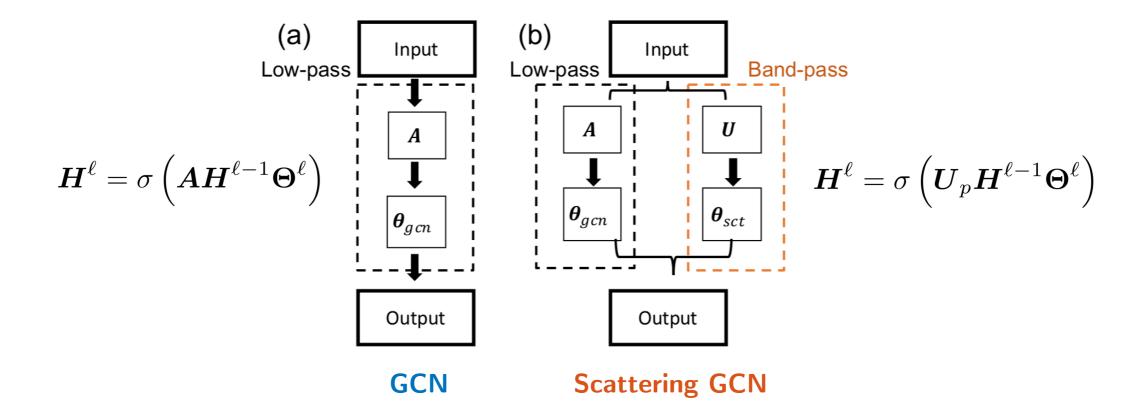
attention function

#### Lecture 3

- Machine learning on graphs: Overview
- Convolutional neural networks on graphs
- Message passing neural networks
- Recent developments & Applications

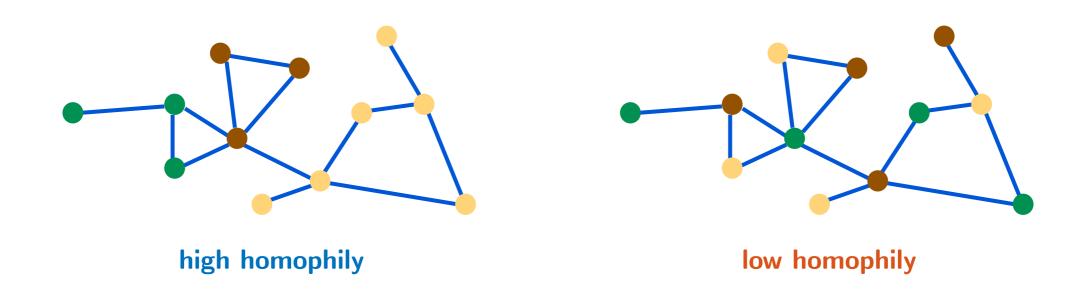
# Over-smoothing

- GCN implements low-pass filtering (which can lead to "over-smoothing")
- Idea: combine low-pass and band-/high-pass filtering



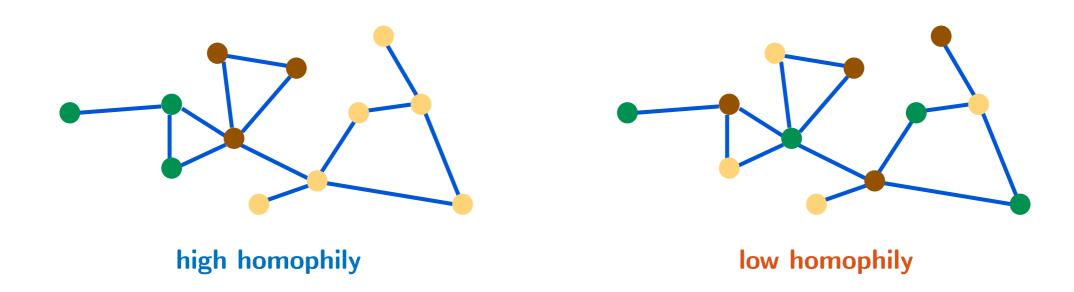
# Label homophily

Heterophily of node labels poses a challenge



# Label homophily

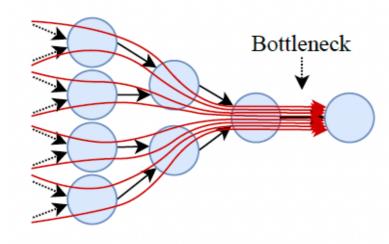
Heterophily of node labels poses a challenge



- Idea:
  - ego- and neighbour-embedding separation:  $\mathbf{r}_v^{(k)} = \texttt{COMBINE}\left(\mathbf{r}_v^{(k-1)}, \texttt{AGGR}(\{\mathbf{r}_u^{(k-1)}: u \in \bar{N}(v)\})\right)$
  - $\quad \text{higher-order neighbourhoods:} \quad \mathbf{r}_v^{(k)} = \text{combine}\left(\mathbf{r}_v^{(k-1)}, \text{ aggr}(\{\mathbf{r}_u^{(k-1)}: u \in N_1(v)\}), \text{ aggr}(\{\mathbf{r}_u^{(k-1)}: u \in N_2(v)\}), \ldots\right)$
  - combination of intermediate representations:  $\mathbf{r}_v^{(\text{final})} = \text{combine}\left(\mathbf{r}_v^{(1)}, \mathbf{r}_v^{(2)}, \dots, \mathbf{r}_v^{(K)}\right)$

# Over-squashing

• Input graph may not be ideal for message passing (e.g., "over-squashing")

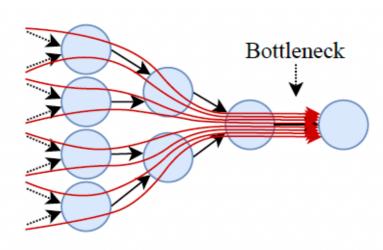


(b) The bottleneck of graph neural networks

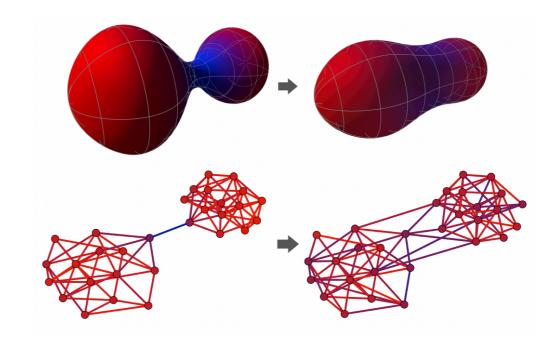
over-squashing caused by **bottlenecks** 

# Over-squashing

- Input graph may not be ideal for message passing (e.g., "over-squashing")
- Idea: "rewiring" as pre-processing step to mitigate over-squashing



(b) The bottleneck of graph neural networks

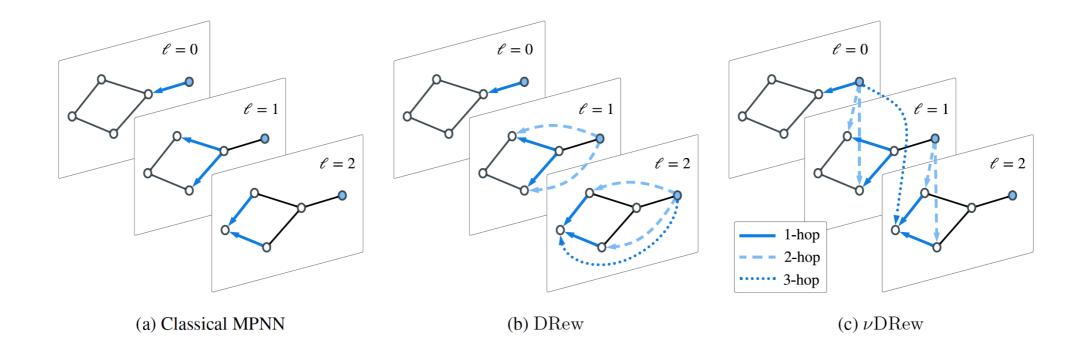


over-squashing caused by **bottlenecks** 

bottlenecks are linked to negatively curved edges

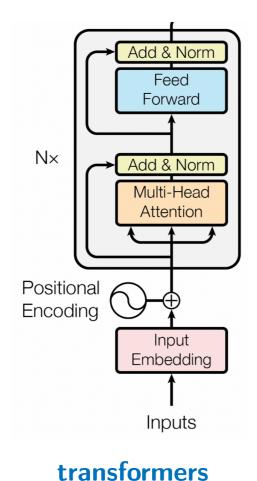
## Adaptive message passing

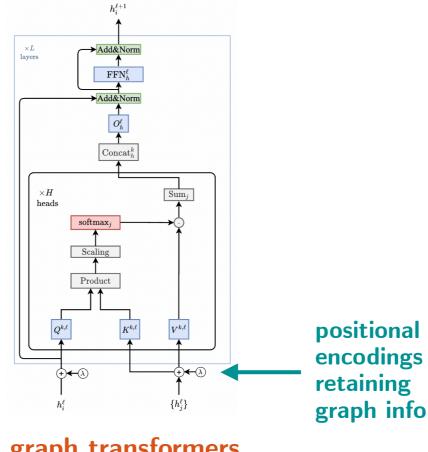
• Idea: modifying the "computational" graph for adaptive message passing



# Graph transformers

Idea: build transformers on graphs by generalising GAT to global attention





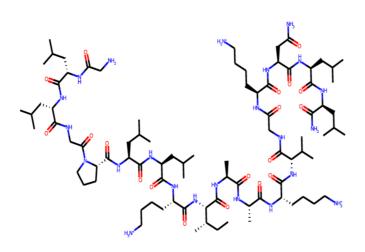
graph transformers

### Long-range interactions

Which tasks require long-range interactions and how to measure them?

### Long-range interactions

Which tasks require long-range interactions and how to measure them?



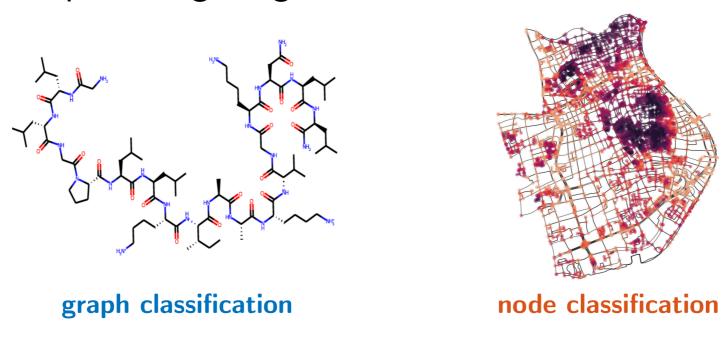
graph classification

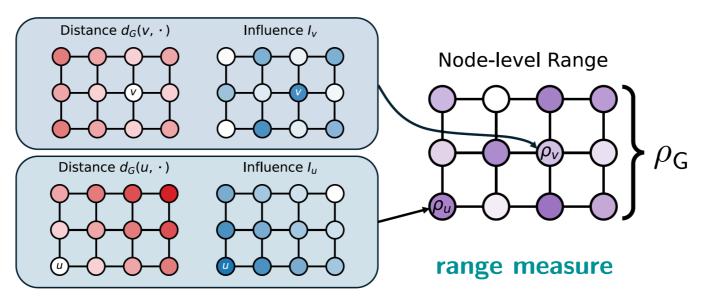


node classification

### Long-range interactions

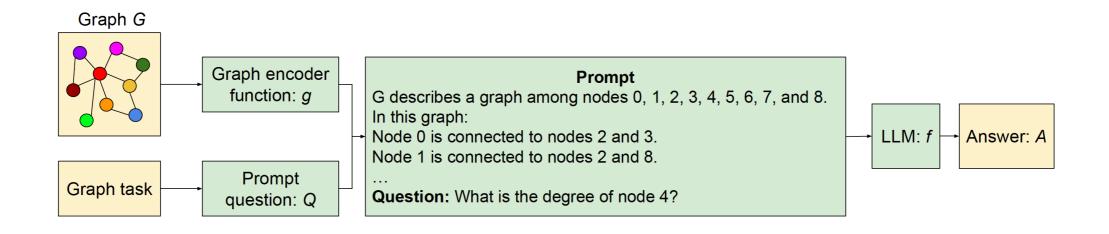
Which tasks require long-range interactions and how to measure them?





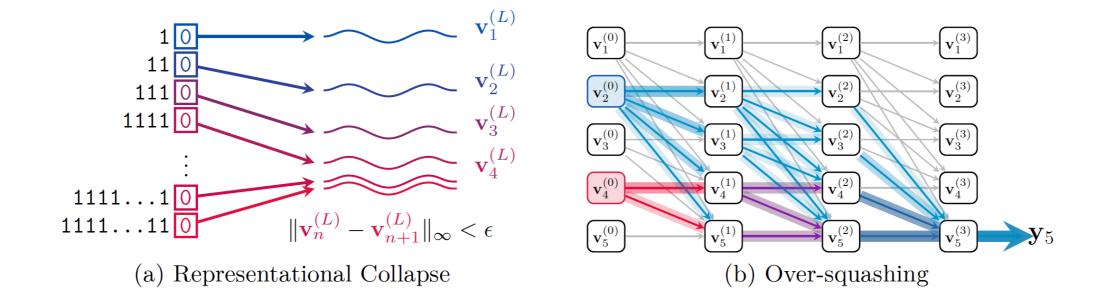
# LLMs for graph learning

• Can LLMs understand graph-structured data?

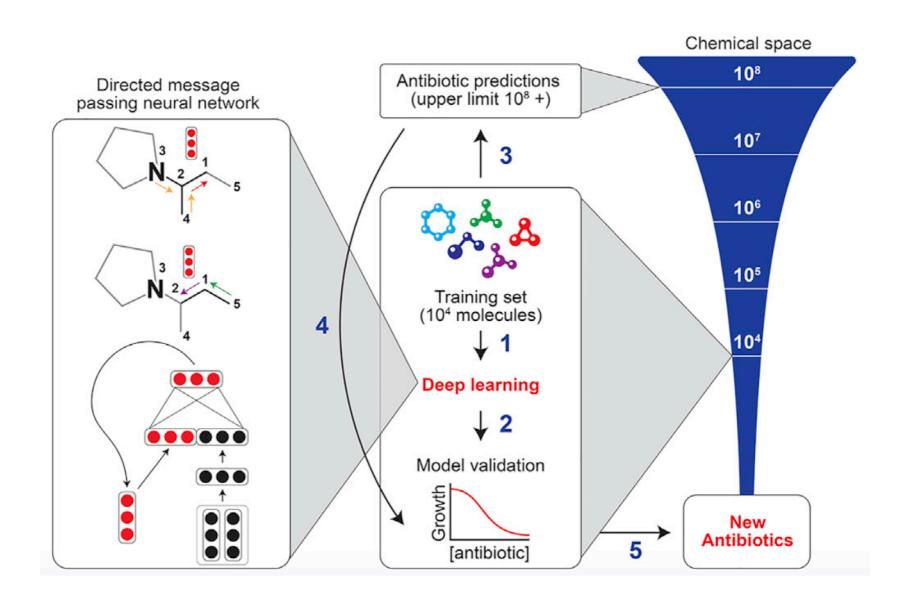


# Graph learning for LLMs

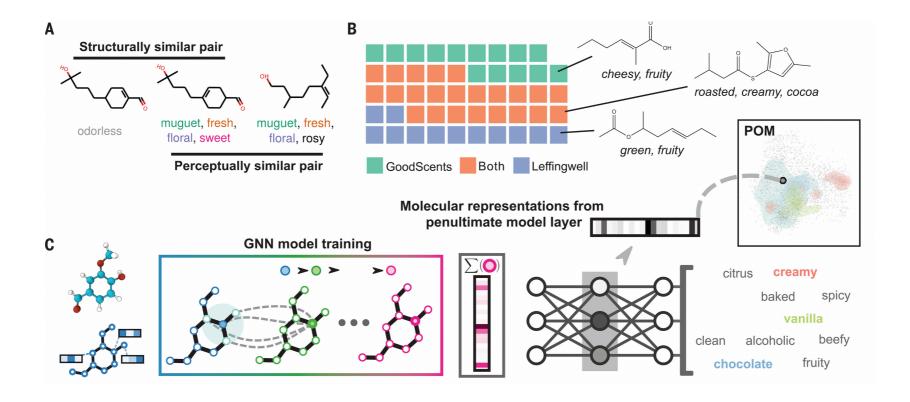
Can insights from graph learning help understand LLMs?



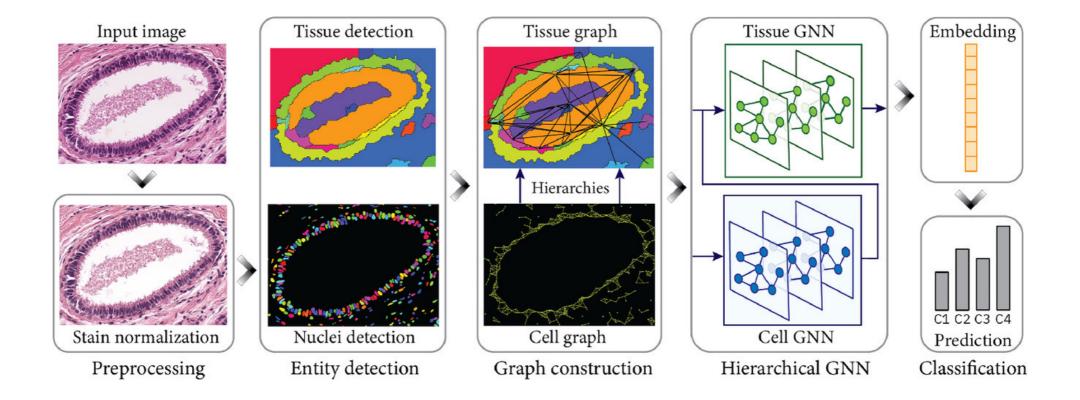
### Application I: Drug discovery



# Application II: Odour perception



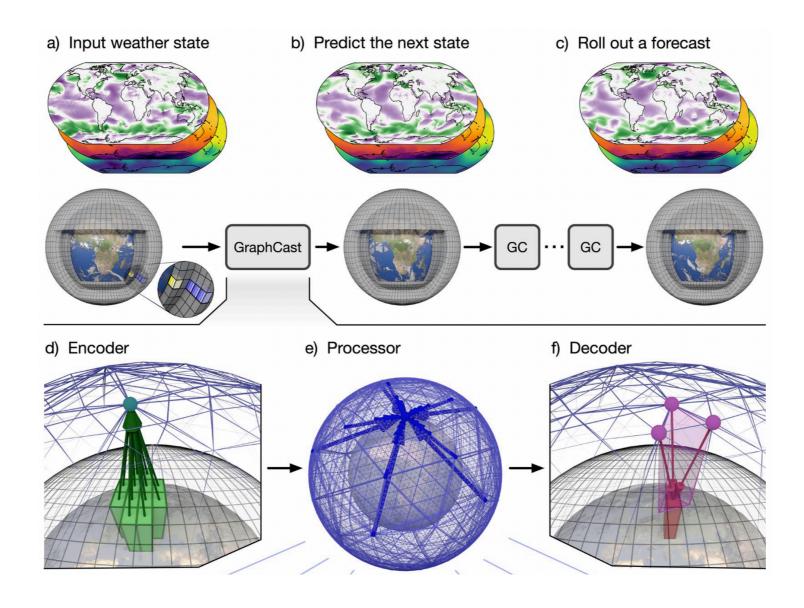
# Application III: Medical imaging



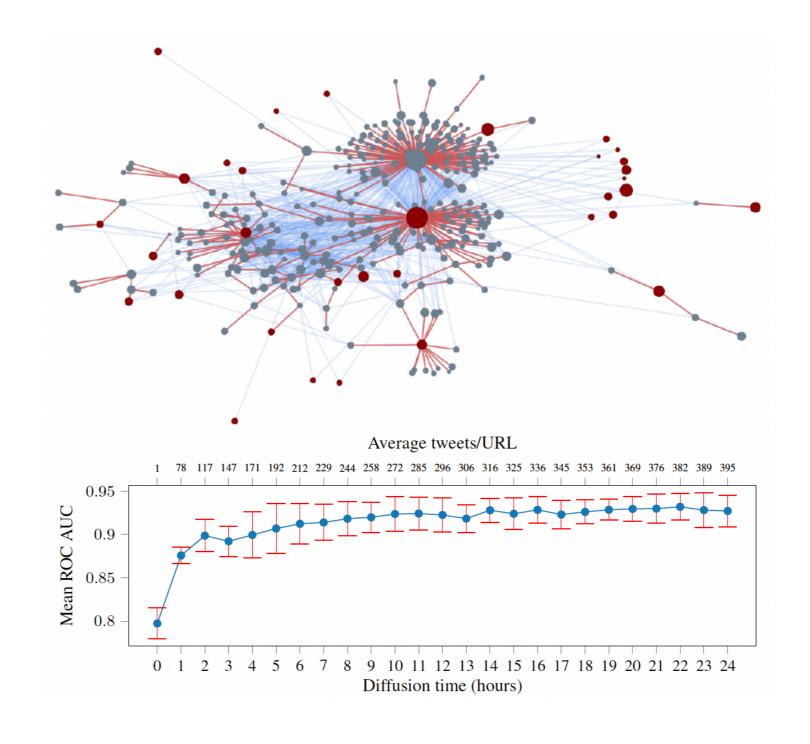
# Application IV: Traffic prediction

# Application IV: Traffic prediction

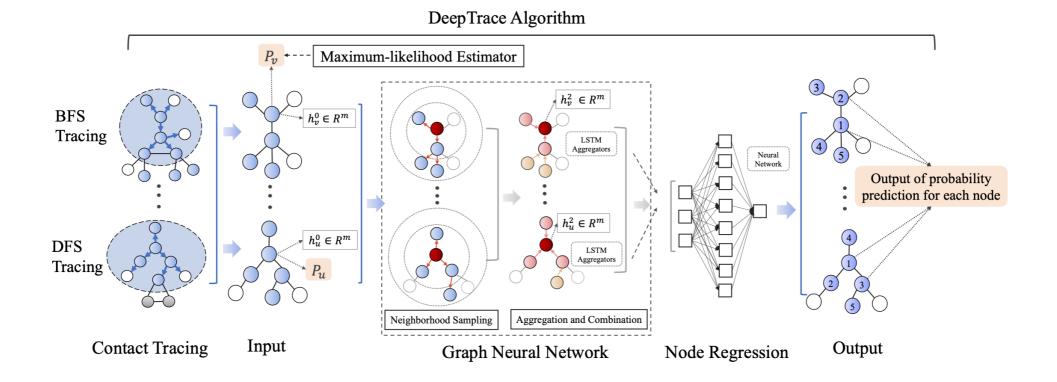
## Application V: Weather forecasting



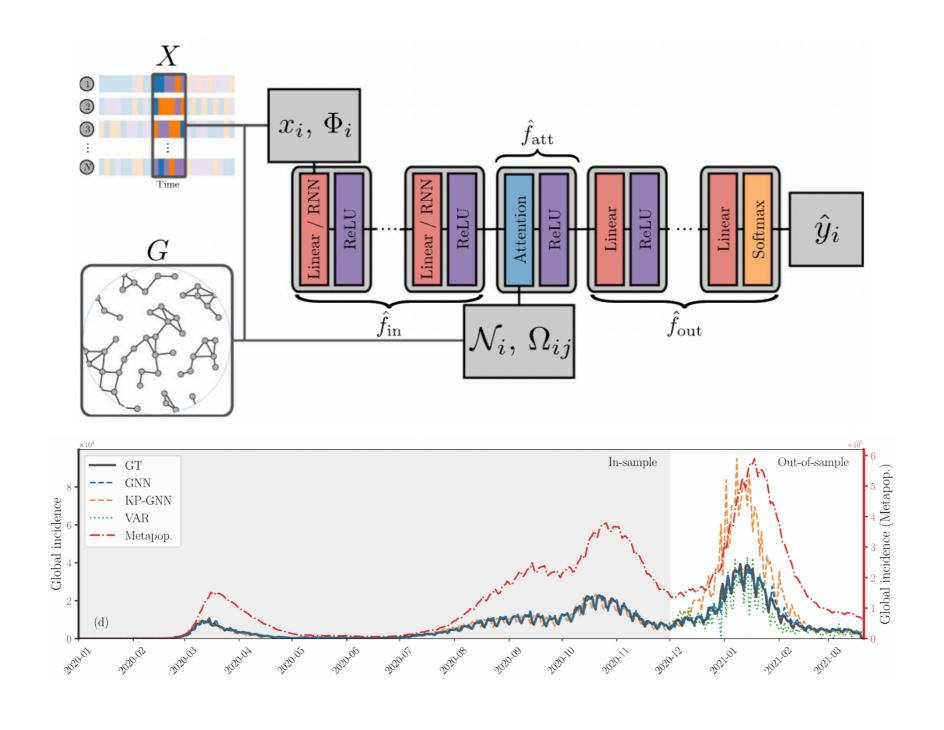
## Application VI: Fake news detection



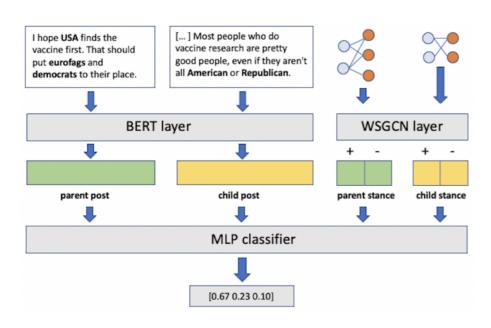
# Application VII: Contact tracing



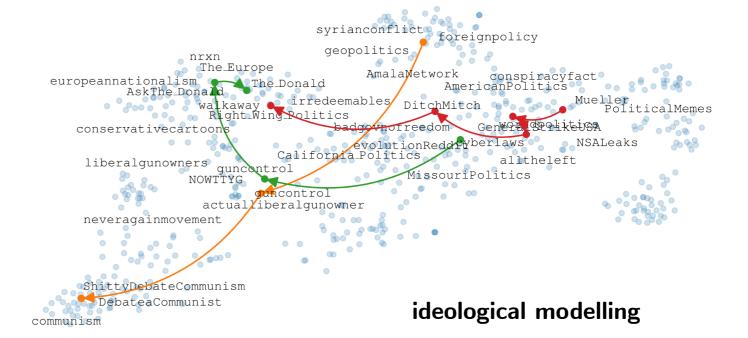
# Application VIII: Contagion dynamics



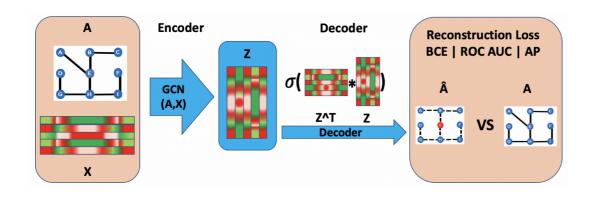
# Application IX: Language modelling

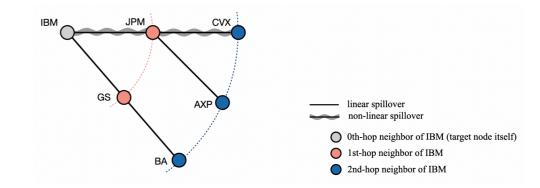


disagreement prediction



## Application X: Stock market analysis





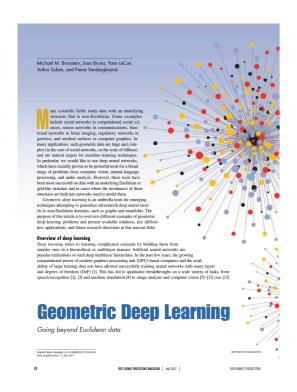
market instability

volatility forecasting

# Deep learning on graphs - Summary

- Fast-growing field that extends data analysis to non-Euclidean domain
- Highly interdisciplinary: machine learning, signal processing, harmonic analysis, network science, differential geometry, applied statistics
- Promising directions
  - beyond convolutional models or MPNNs
  - expressive power of graph ML models
  - robustness & generalisation & scalability
  - interpretability & causal inference
  - construction & refinement of initial graphs
  - optimisation & implementation issues
  - foundation models for graph-structured data

#### References



#### Representation Learning on Graphs: Methods and Applications

Department of Computer Science Stanford University Stanford, CA, 94305

Machine learning on graphs is an important and subsigations task with applications ranging from drug design to friendship recommendation is nocial networks. The primary challenge in this domain is finding a way to represent, or encode, graph structure so that it can be easily explored by machine learning models. Traditionally, machine learning approaches relied on user-defined heuristics to extract features modelly. Traditionally, for a graph (e.g., deeper statistics or beart functions). However, recent years have seen a surge in approaches that automatically learn to encode graph structure into two discontinuous desired interest in the second graph structure into two-dimensional medicalings, using exchanges based on deep learning and nonlinear dimensionally reduction. Here we provide a conceptual review of key advancements in this area of propresentation learning on graphs, including matrix factorization-based embods, random-wank based algorithms, and graph comodational networks. We review methods to embed individual nodes as well as approaches to embed enter (subsyraphs. In doing so, we develop a unified framework to describe these recent approaches, and we highlight a number of important applications and directions for future work.

Graphs are a ubiquitous data structure, employed extensively within computer science and related fields. Social networks, molecular graph structures, biological protein-protein networks, recommender systems—all of these domains and many more can be readily modeled as graphs, which capture interactions (i.e., edges) between individual units (i.e., nodes). As a consequence of their ubiquity, graphs are the backbone of countiess systems, allowing relational knowledge about interacting entiries to be efficiently stored and accessed [2].

However, graphs are not only useful as structured knowledge repositaries: they also play a key role in summing the control of the protein of the structured of the statement of the structured of the statement of the

Copyright 2017 IEEE Personal use of this material is permitted. However, permission to advertising or promotional purposes or for creating new collective works for resale or redistribution copyrighted component of this work in other works must be obtained from the IEEE.

Balletin of the IEEE Computer Society Technical Committee on Data Engineering

#### A Comprehensive Survey on Graph Neural Networks

Zonghan Wu<sup>©</sup>, Shirui Pan<sup>©</sup>, Member, IEEE, Fengwen Chen, Guodong Long<sup>©</sup>

#### Machine Learning on Graphs: A Model and Comprehensive Taxonomy

Sami Abu-El-Haija USC Information Sciences Institute Marina Del Rey, CA, 90292, USA

Bryan Perozzi Google Research New York, NY, 10011, USA

Christopher Ré Stanford University Stanford, CA, 94305, USA

Abstract

There has been a surge of recent interest in graph representation learning (GRL), GRL methods have generally fallen into three main categories, based on the availability of labeled data. The first, network emboding, focuses on learning unsupervised representations of relational structure. The second, graph regularized neural networks, leverages graphs to aggreat normal networks on surface regularization objective for semi-supervised learning. The third, graph neural networks, sims to learn differentiable function over thereot topological structures of the properties of the formal properties of the properties of the properties of the formal properties of the pr

**Keywords:** Network Embedding, Graph Neural Networks, Geometric Deep Learning, Manifold Learning, Relational Learning

License: CC-BY 4.0, see https://creativecommons.org/licenses/by/4.0/. Attribution requirements are provided at http://jmir.org/papers/v23/20-852.html.