

Applied Machine Learning

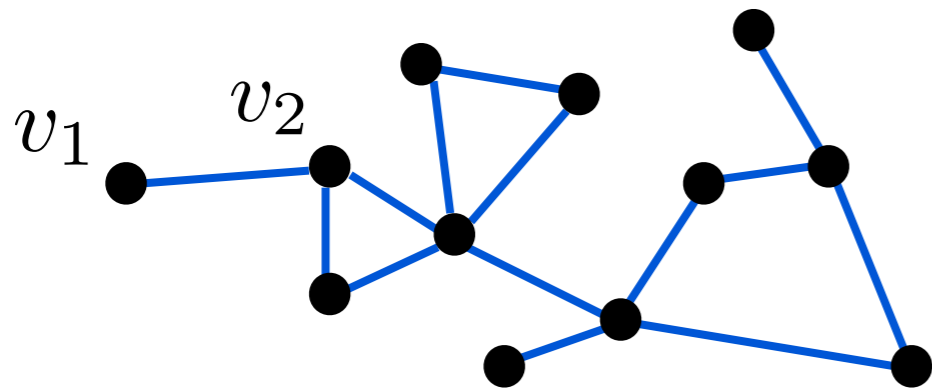
Graph machine learning - Part II

Xiaowen Dong

Department of Engineering Science

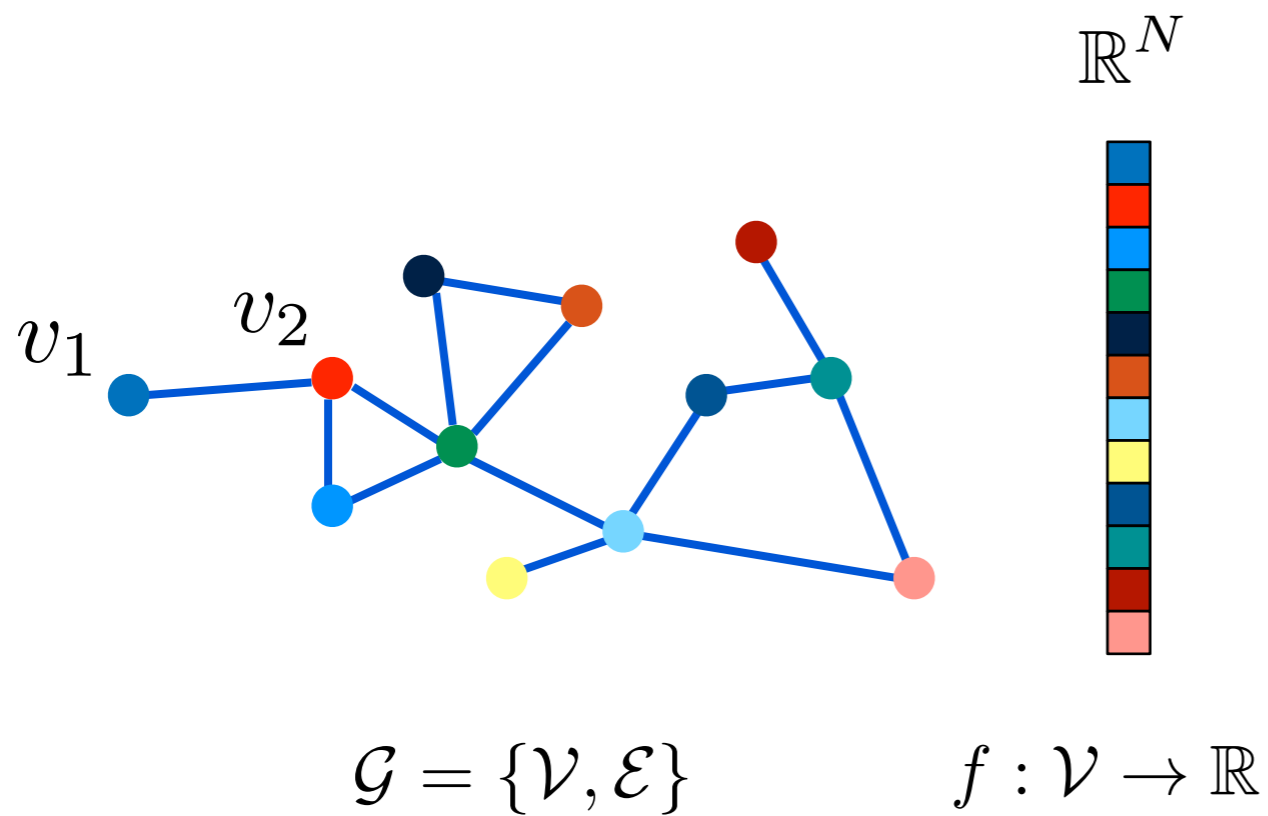


Recap: Graph signal processing

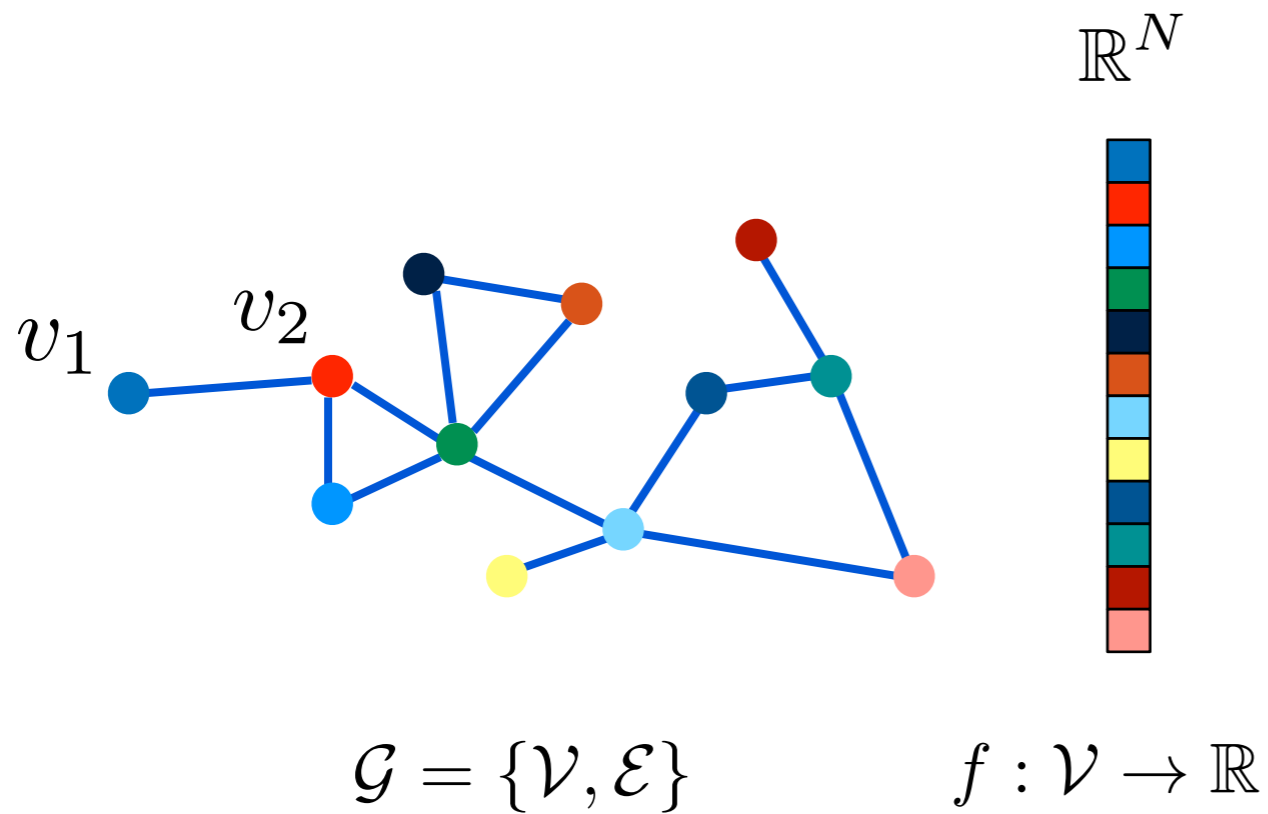


$$\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$$

Recap: Graph signal processing



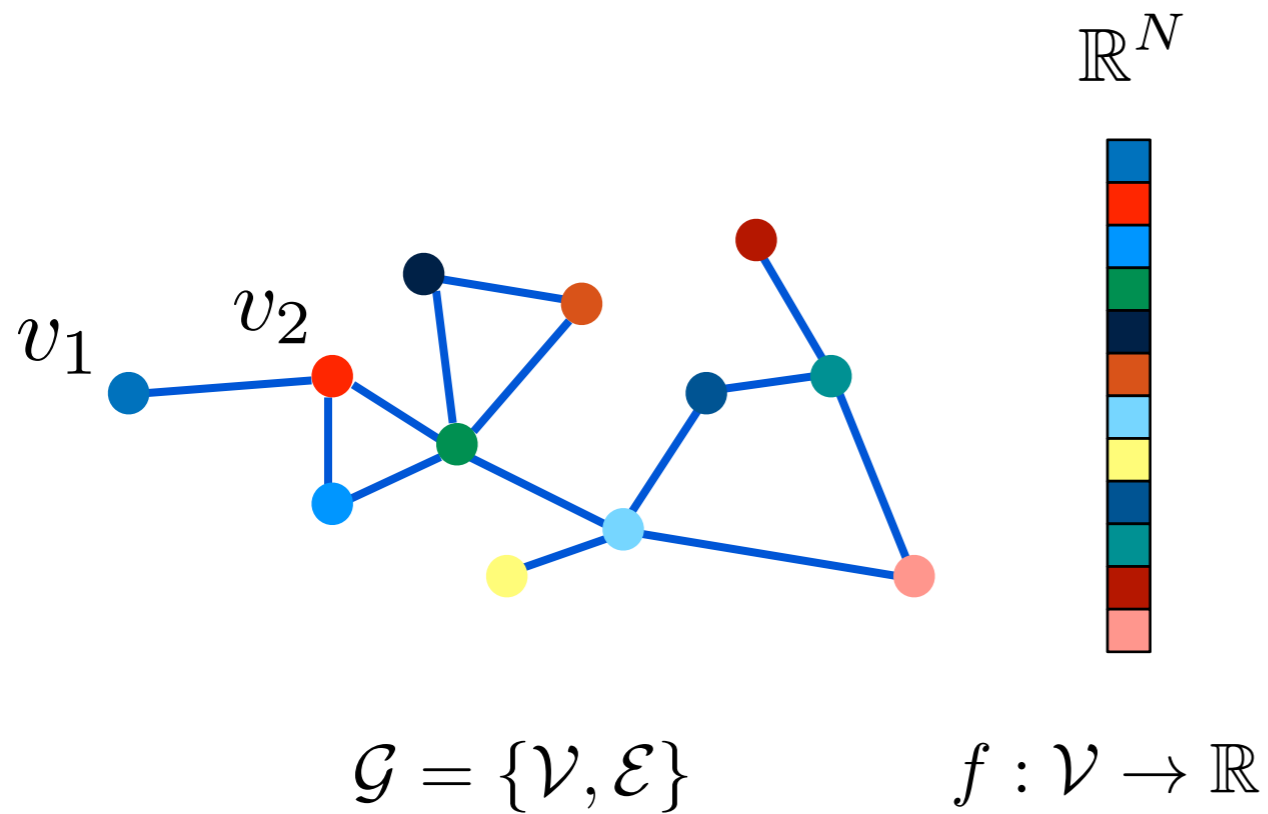
Recap: Graph signal processing



- key concepts
 - smoothness

$$f^T L f = \frac{1}{2} \sum_{i,j=1}^N W_{ij} (f(i) - f(j))^2$$

Recap: Graph signal processing

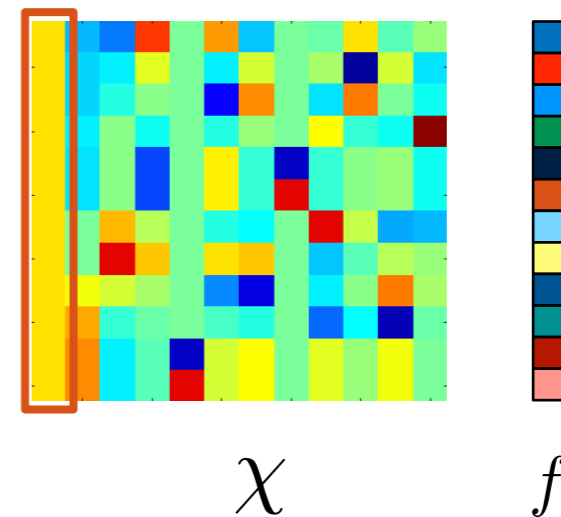


- key concepts

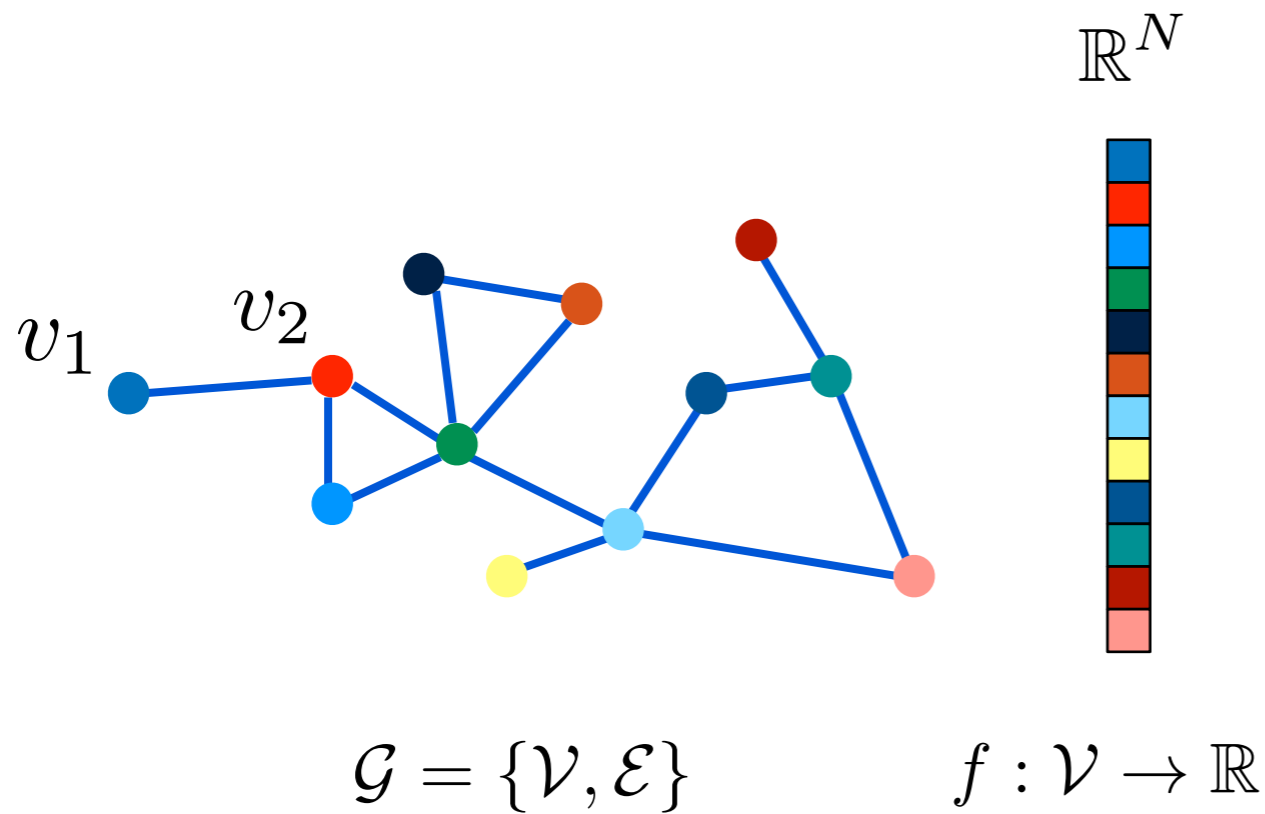
- smoothness

$$f^T L f = \frac{1}{2} \sum_{i,j=1}^N W_{ij} (f(i) - f(j))^2$$

- Fourier-like analysis



Recap: Graph signal processing

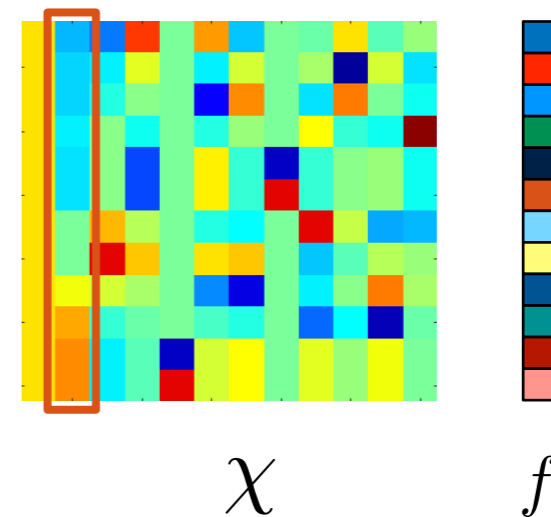


- key concepts

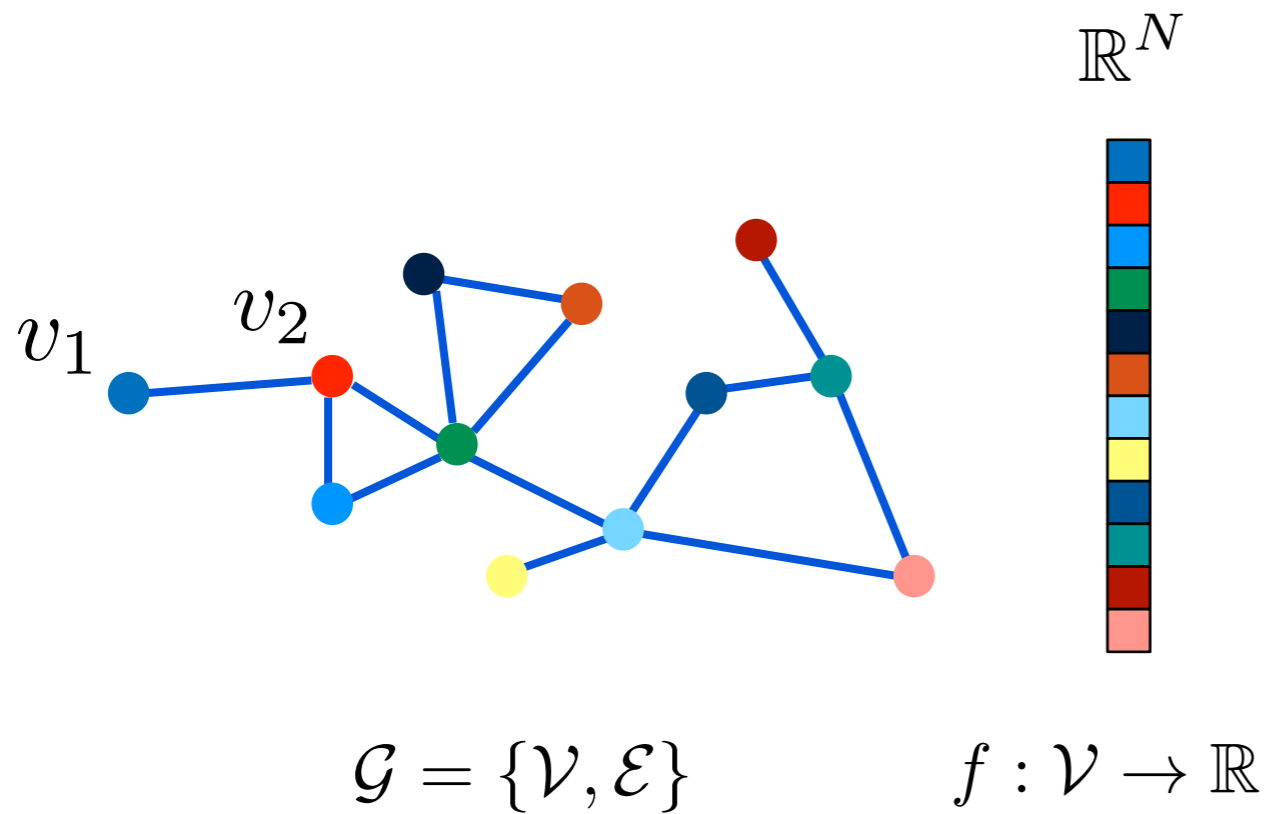
- smoothness

$$f^T L f = \frac{1}{2} \sum_{i,j=1}^N W_{ij} (f(i) - f(j))^2$$

- Fourier-like analysis



Recap: Graph signal processing

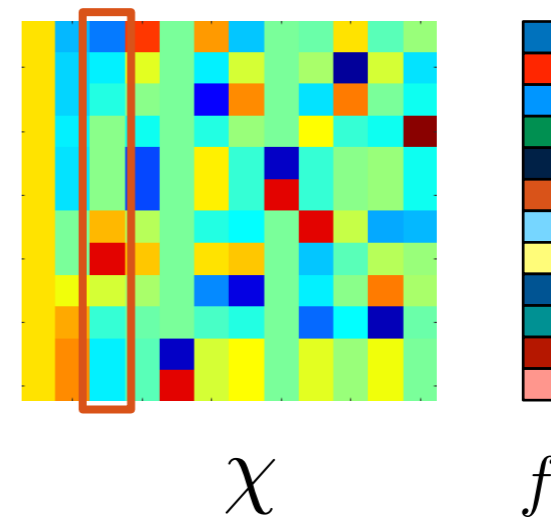


- key concepts

- smoothness

$$f^T L f = \frac{1}{2} \sum_{i,j=1}^N W_{ij} (f(i) - f(j))^2$$

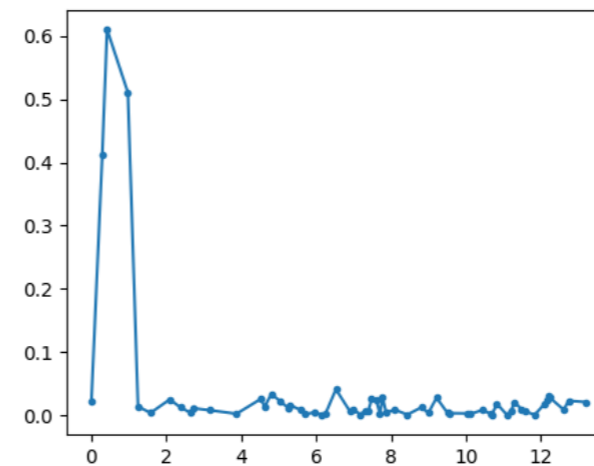
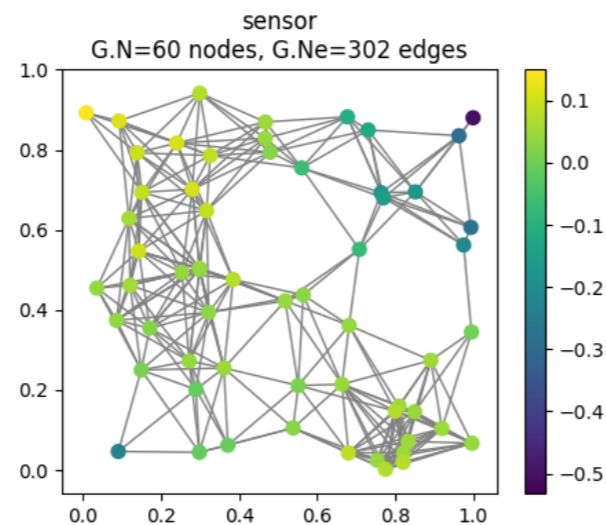
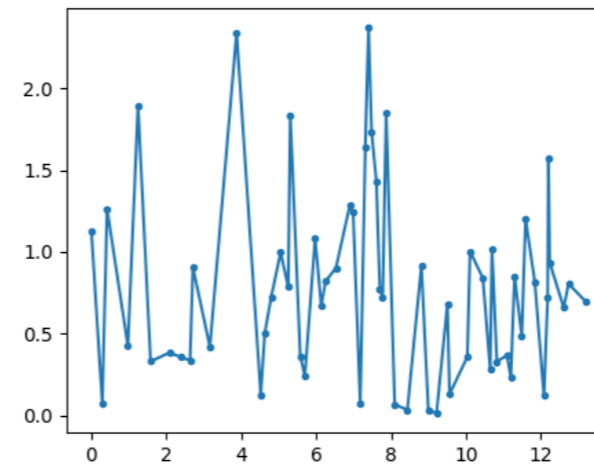
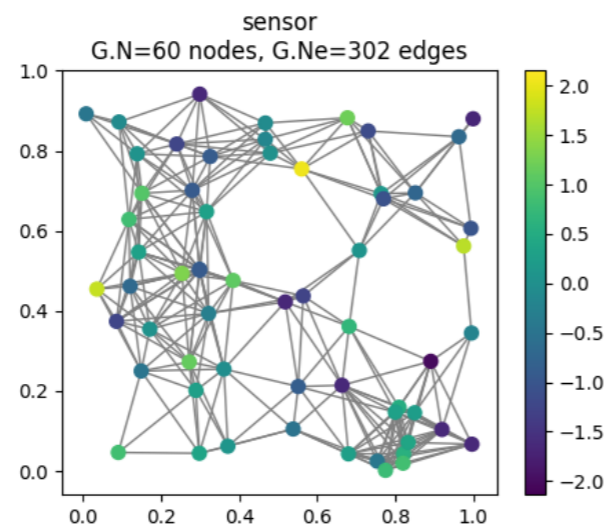
- Fourier-like analysis



Recap: Graph signal processing

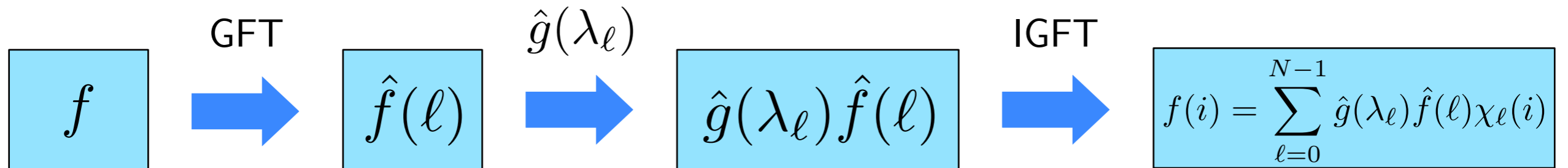
- Graph Fourier transform

$$\hat{f}(\ell) = \langle \chi_\ell, f \rangle : \begin{bmatrix} | & & | \\ \chi_0 & \cdots & \chi_{N-1} \\ | & & | \end{bmatrix}^T \begin{matrix} | \\ f \\ | \end{matrix}$$



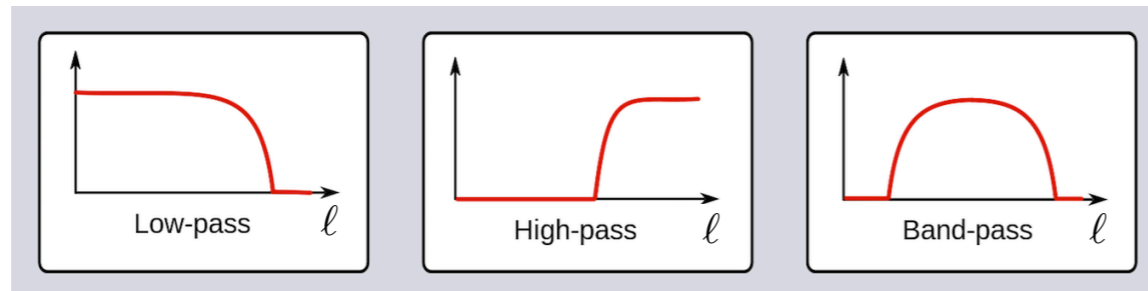
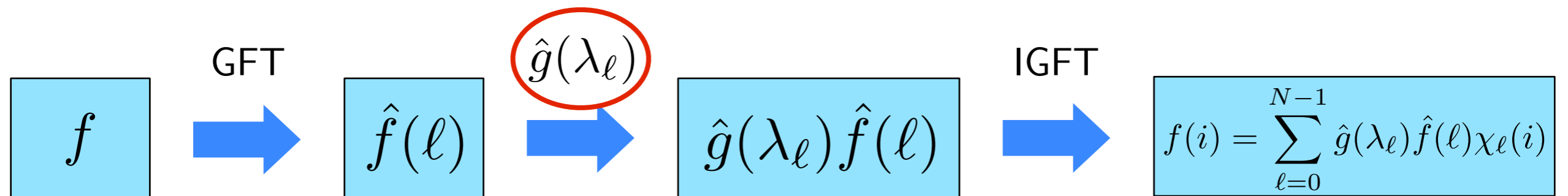
Recap: Graph signal processing

- Graph spectral filtering



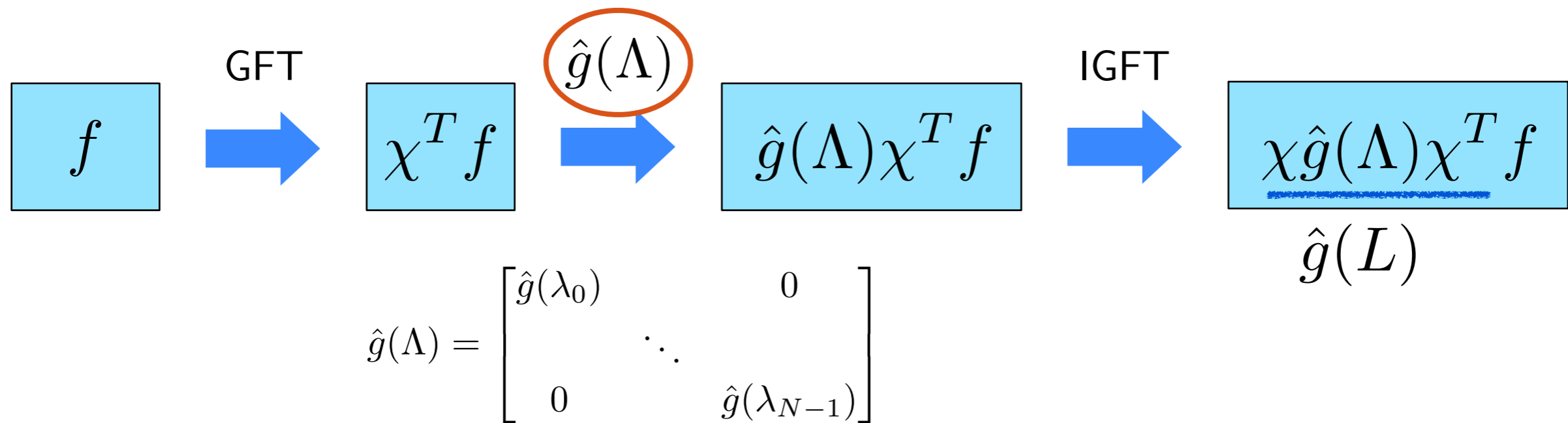
Recap: Graph signal processing

- Graph spectral filtering

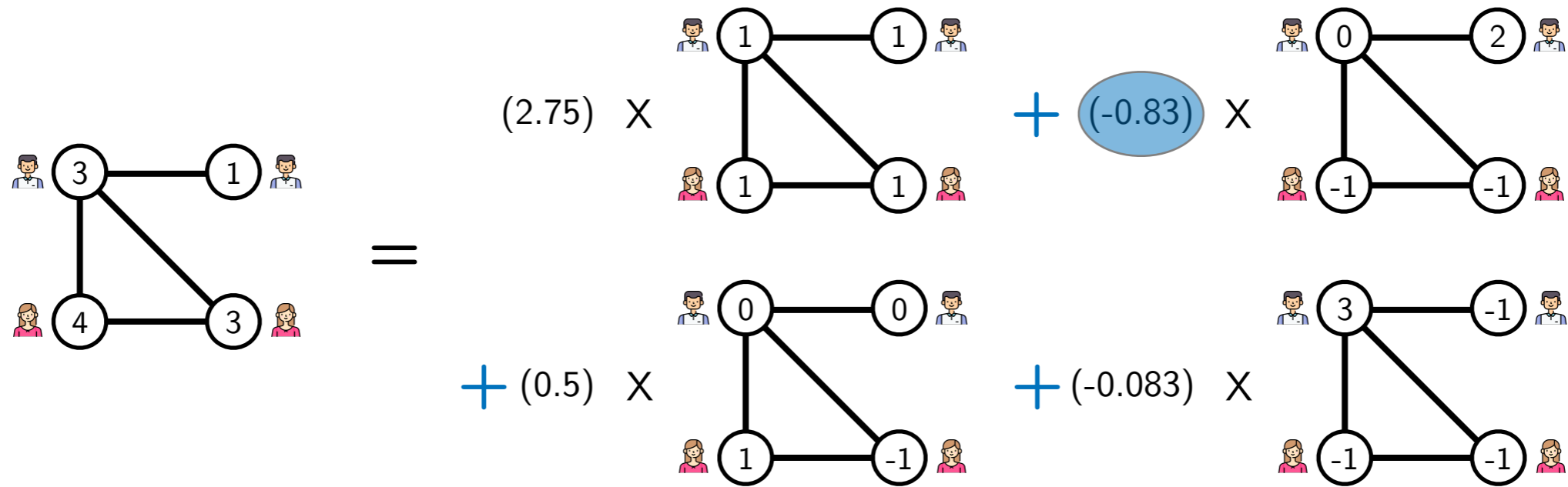


Recap: Graph signal processing

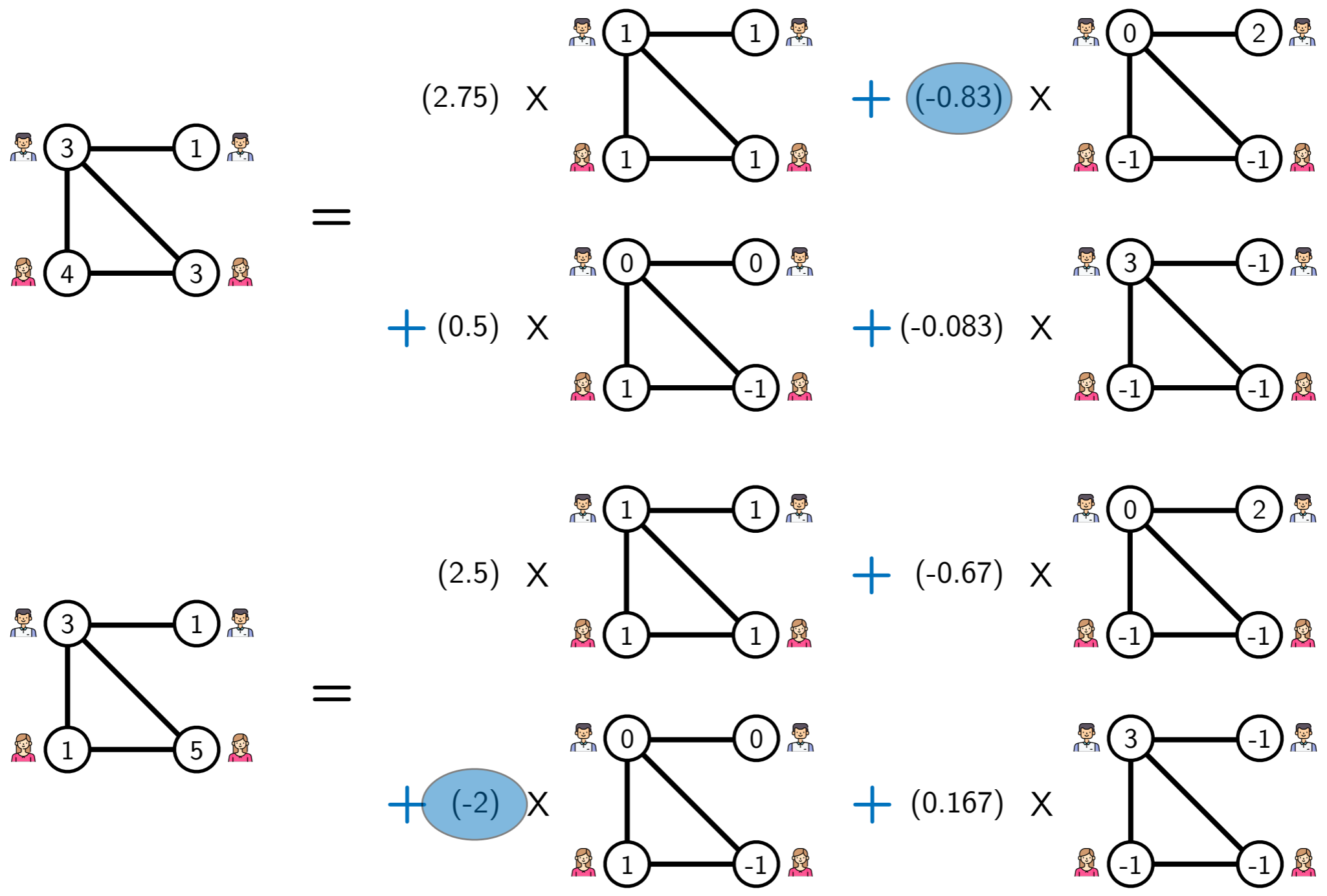
- Graph spectral filtering



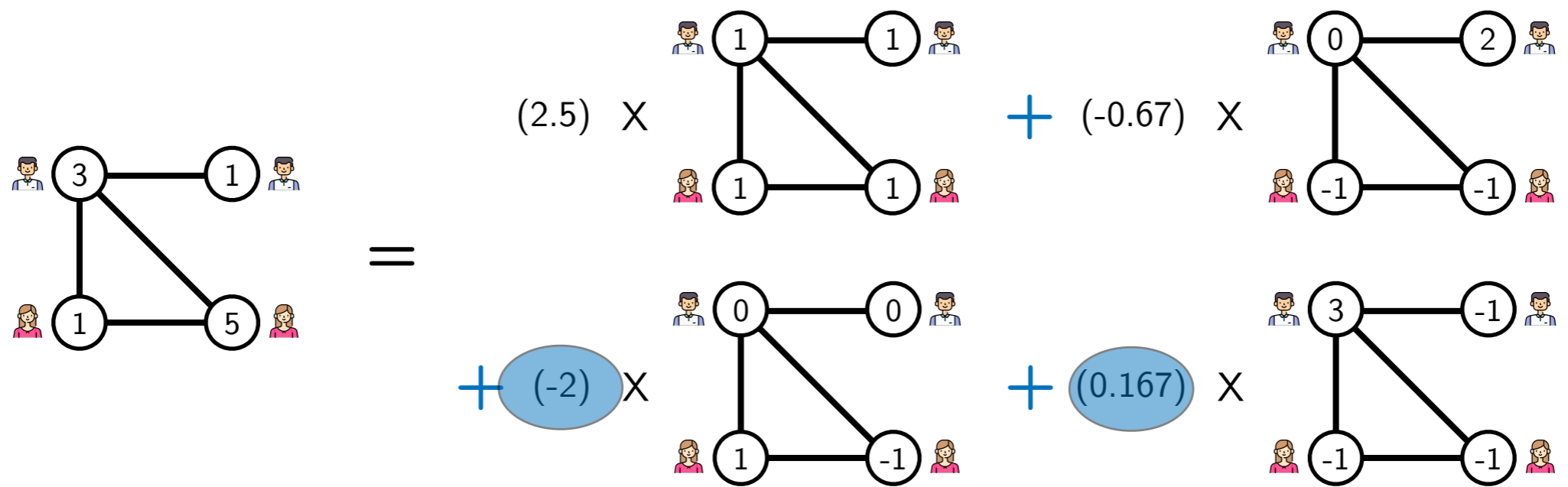
An example on movie rating



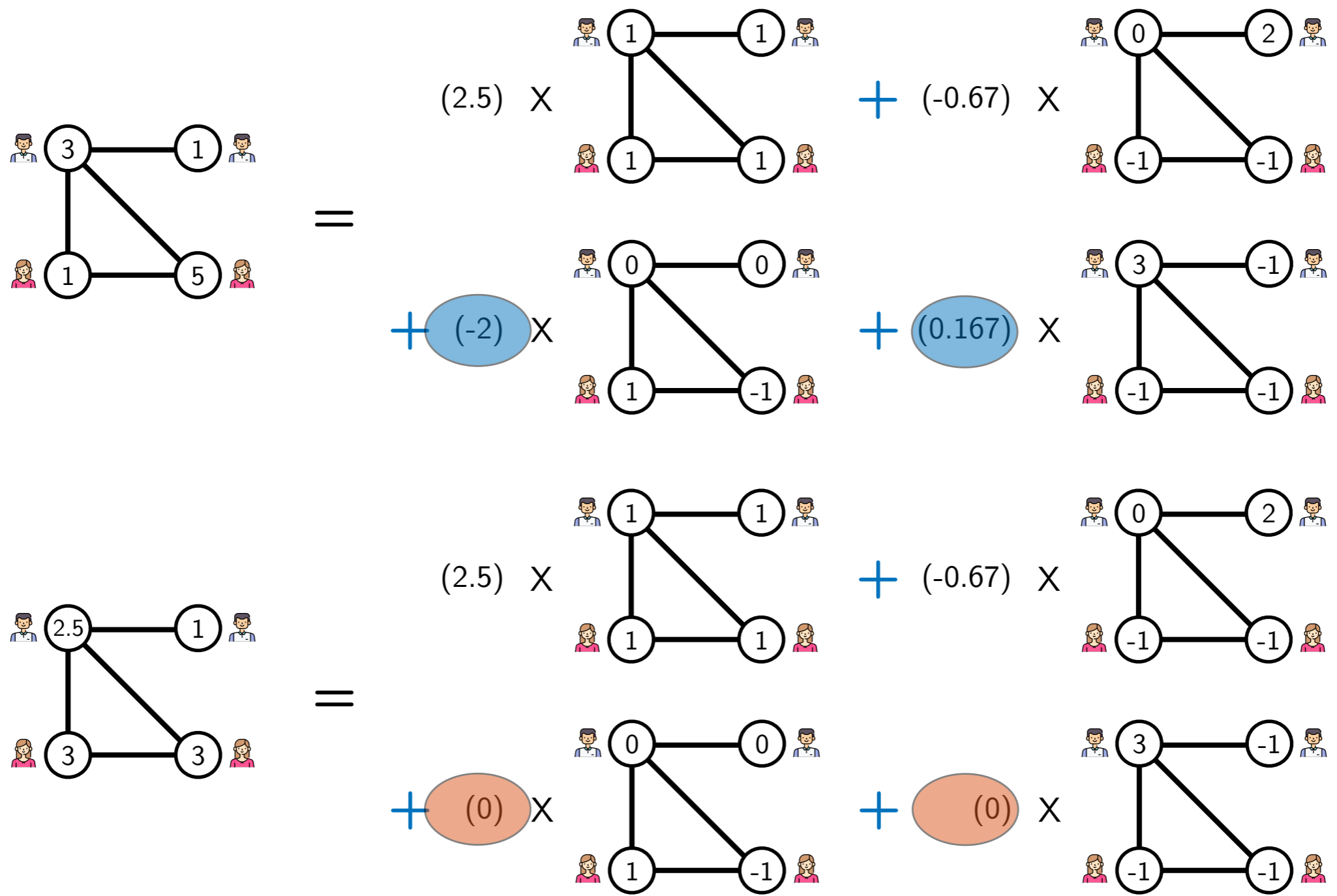
An example on movie rating



An example on movie rating

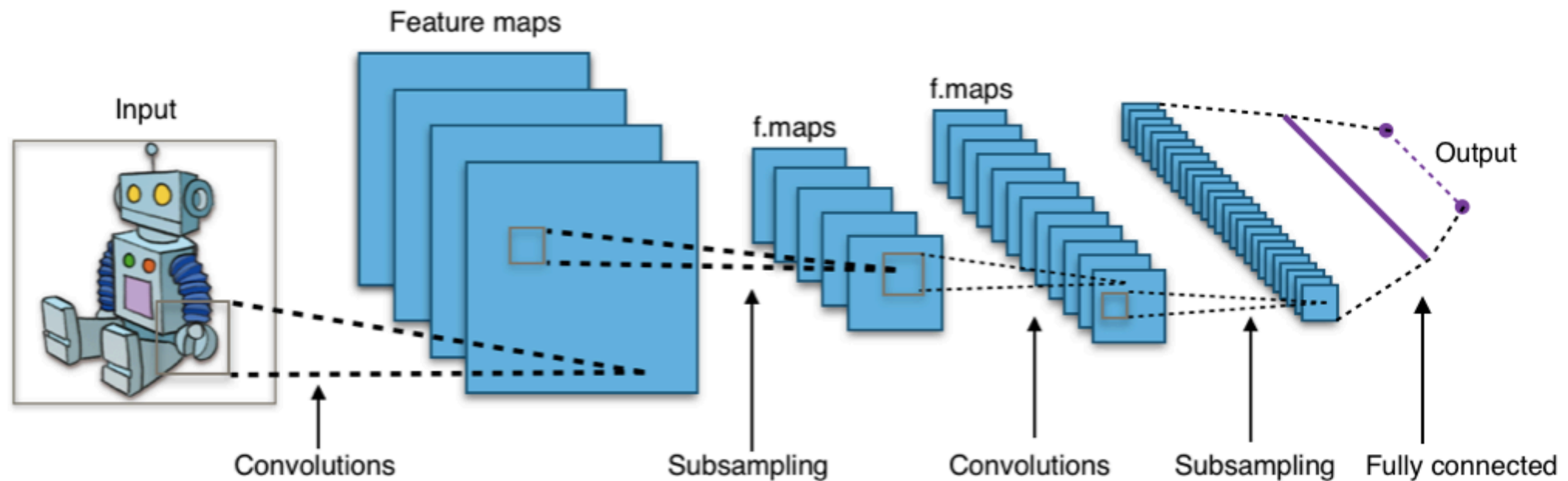


An example on movie rating



Graph Neural Networks

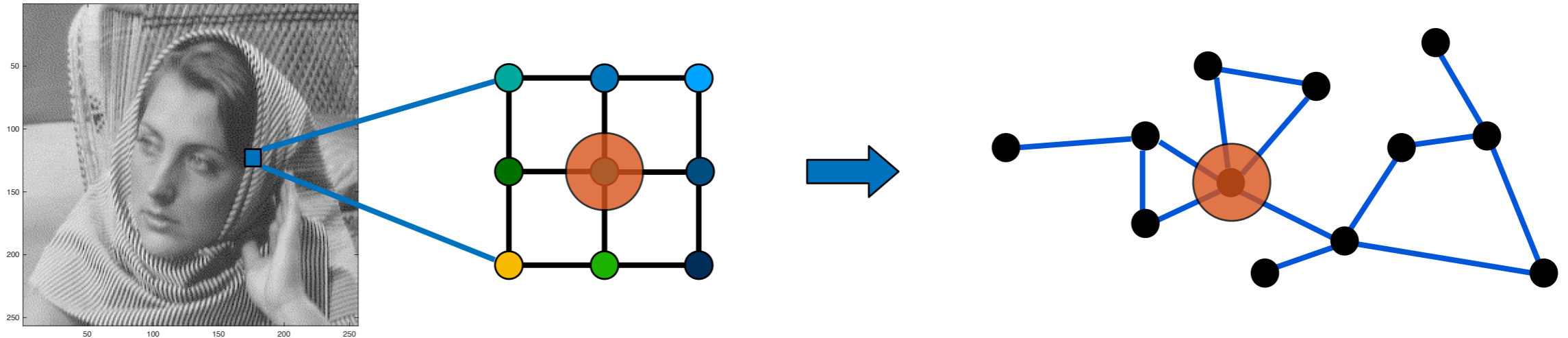
CNNs exploit structure within data



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 do it on graphs?
- **localisation:** what's the notion of locality?
- **multi-scale:** how to down-sample on graphs?
- **efficiency:** how to keep the computational complexity low?

Convolution on graphs

classical convolution

time domain

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

3 ₀	3 ₁	2 ₂	1	0
0 ₂	0 ₂	1 ₀	3	1
3 ₀	1 ₁	2 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

Convolution on graphs

classical convolution

time domain

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

3 ₀	3 ₁	2 ₂	1	0
0 ₂	0 ₂	1 ₀	3	1
3 ₀	1 ₁	2 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

Convolution on graphs

classical convolution

time domain

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



frequency domain

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

3 ₀	3 ₁	2 ₂	1	0
0 ₂	0 ₂	1 ₀	3	1
3 ₀	1 ₁	2 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

Convolution on graphs

classical convolution

time domain

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



frequency domain

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

convolution on graphs

graph spectral domain

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

Convolution on graphs

classical convolution

time domain

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



frequency domain

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

convolution on graphs

spatial (node) domain

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



graph spectral domain

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



Convolution on graphs

classical convolution

time domain

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



frequency domain

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

convolution on graphs

spatial (node) domain

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



graph spectral domain

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

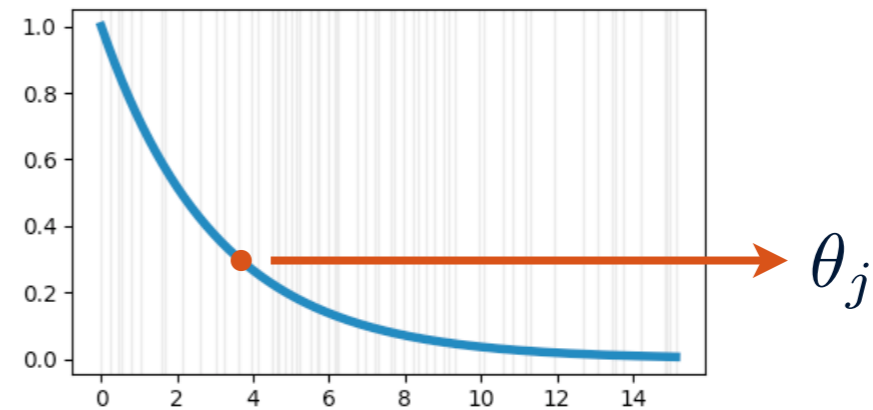
A non-parametric filter

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



learning a non-parametric filter:

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



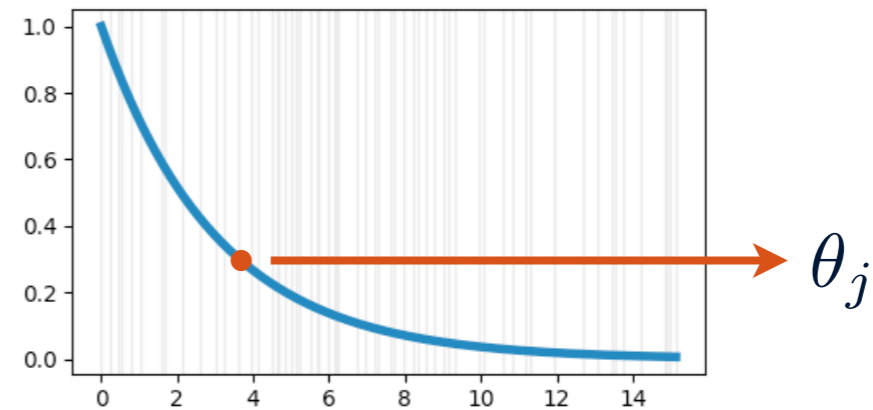
A non-parametric filter

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



learning a non-parametric filter:

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



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

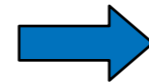
A parametric filter

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



parametric filter as polynomial of Laplacian

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



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

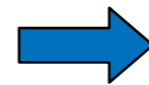
A parametric filter

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



parametric filter as polynomial of Laplacian

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

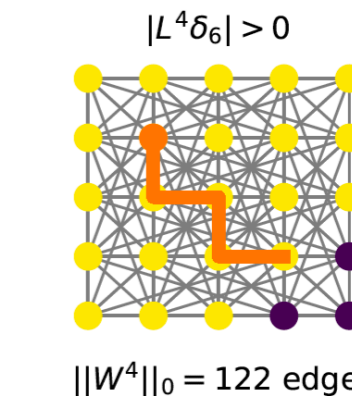
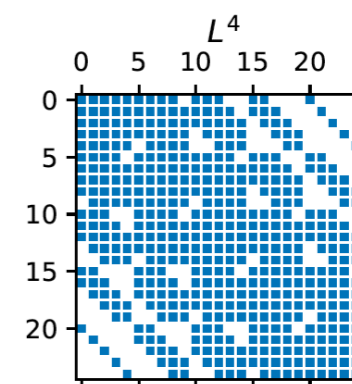
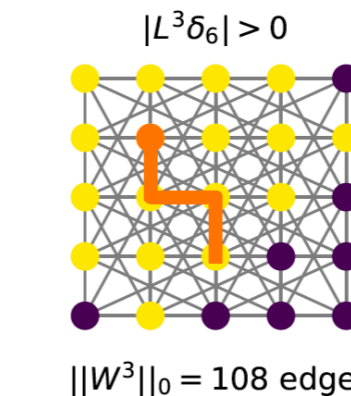
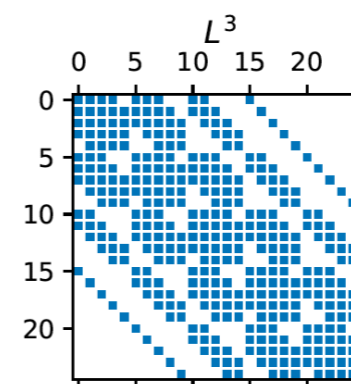
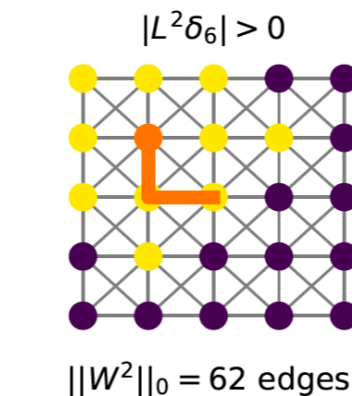
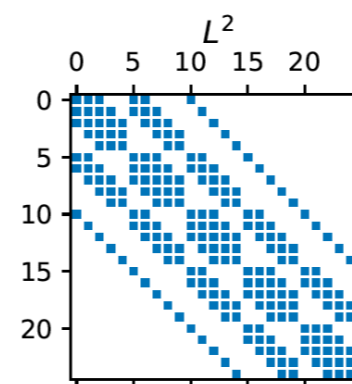
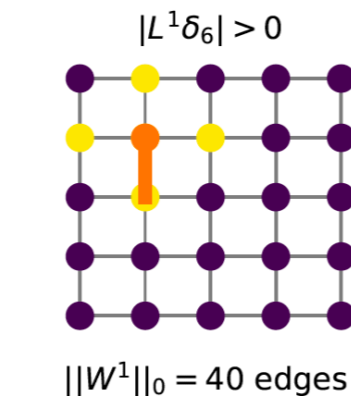
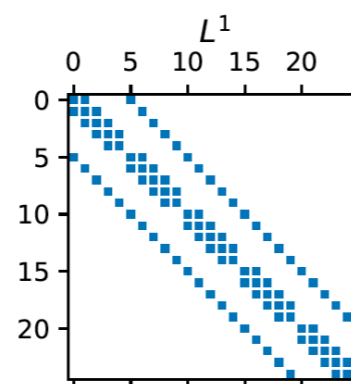
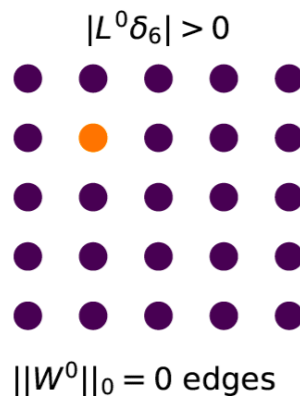
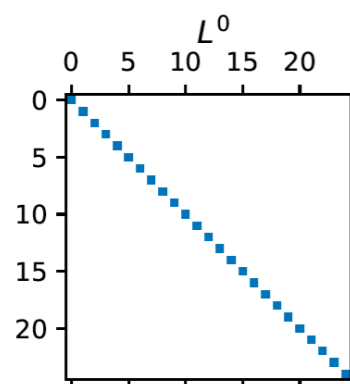


$$\hat{g}_\theta(L) = \sum_{j=0}^K \theta_j L^j$$

what do powers of graph Laplacian capture?

Powers of graph Laplacian

L^k defines the k -neighborhood



Localization: $d_G(v_i, v_j) > K$ implies $(L^K)_{ij} = 0$

(slide by Michaël Deferrard)

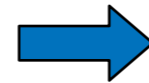
A parametric filter

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



parametric filter as polynomial of Laplacian

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



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

A parametric filter

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



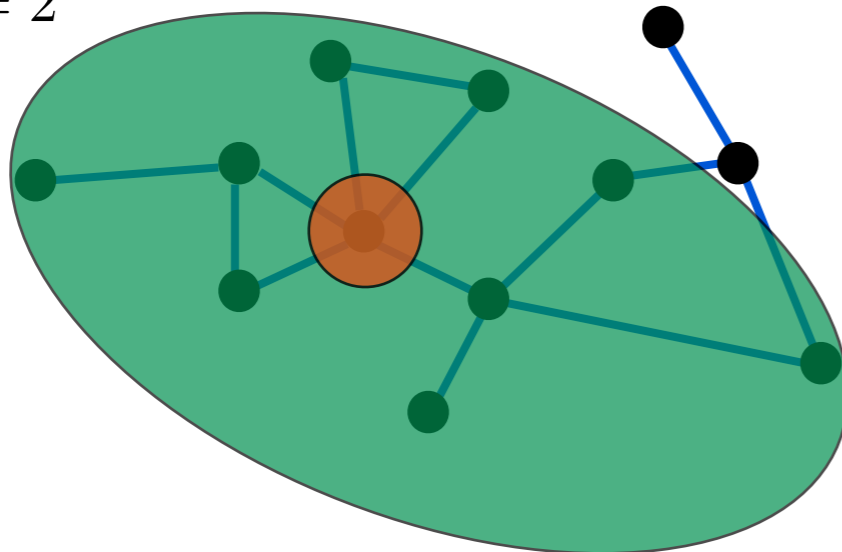
parametric filter as polynomial of Laplacian

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



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

$K = 2$



- localisation within K -hop neighbourhood
- efficient computation via recursive multiplication with L

A simplified parametric filter

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

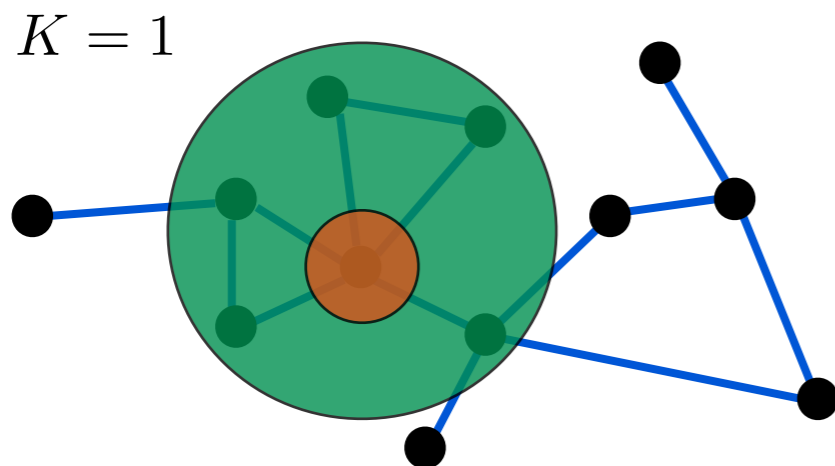


simplified parametric filter

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

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

(localisation within **1-hop** neighbourhood)



A simplified parametric filter

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

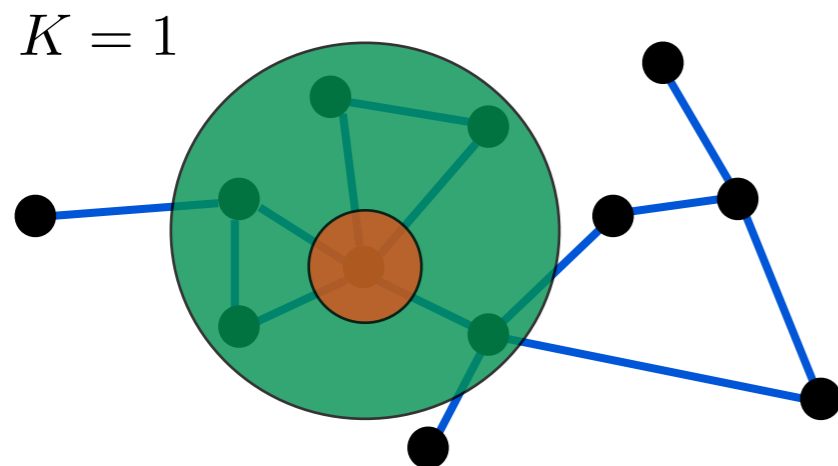


simplified parametric filter

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

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

(localisation within **1-hop** neighbourhood)



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

A simplified parametric filter

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



simplified parametric filter

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

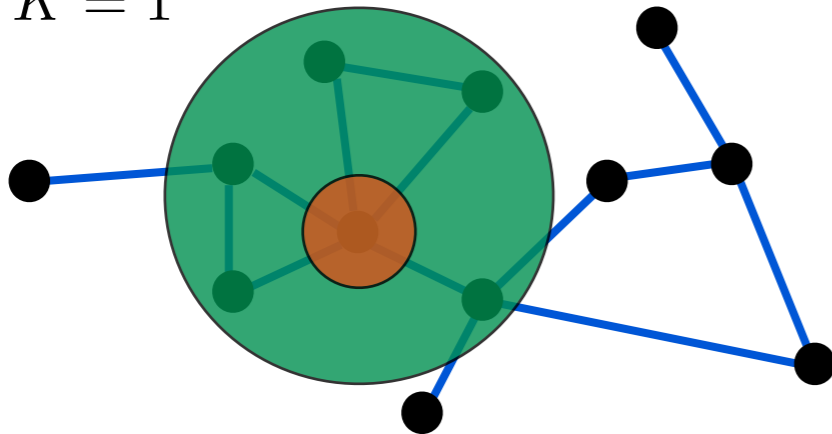
$K = 1$



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

(localisation within **1-hop** neighbourhood)

$K = 1$



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



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

renormalisation



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

A simplified parametric filter

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



simplified parametric filter

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

A simplified parametric filter

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

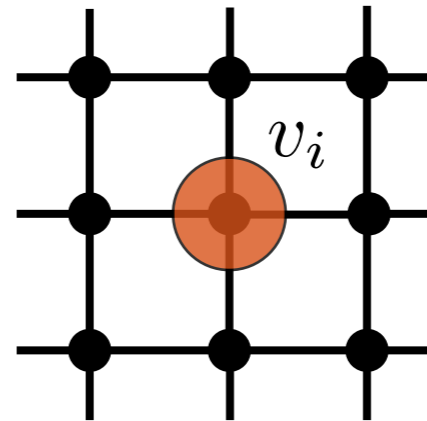


simplified parametric filter

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



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



A simplified parametric filter

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



simplified parametric filter

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

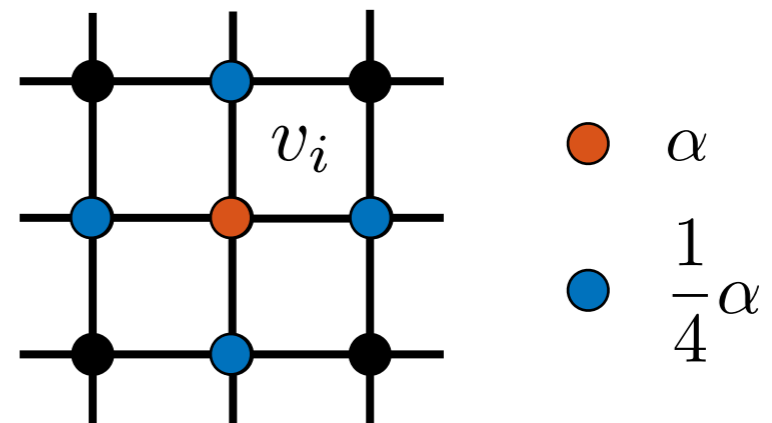


$$y_i = \alpha f_i + \alpha \frac{1}{\sqrt{d_i}} \sum_{j:(i,j) \in \mathcal{E}} w_{ij} \frac{1}{\sqrt{d_j}} f_j$$



unitary edge weights

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



A simplified parametric filter

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



simplified parametric filter

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



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

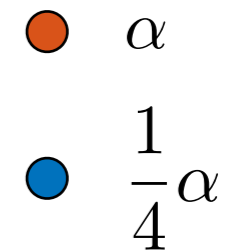
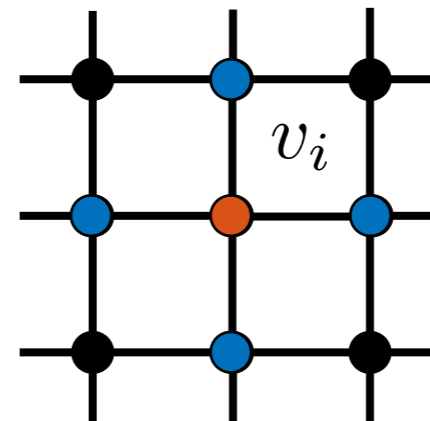


unitary edge weights

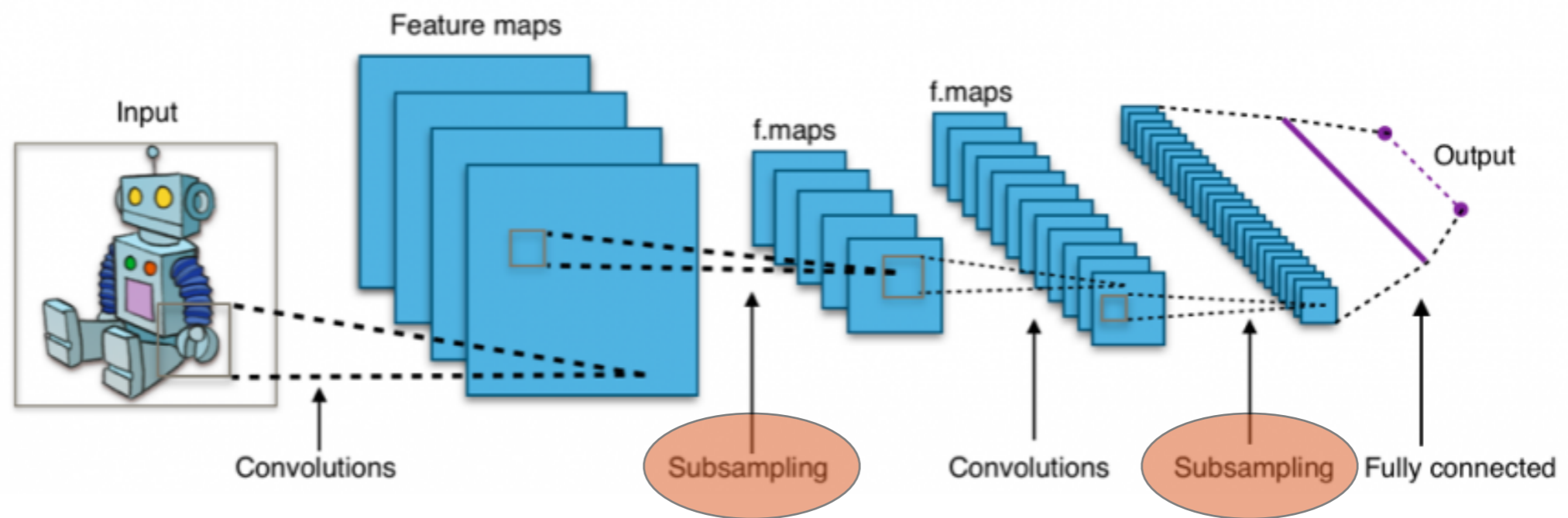
$$y_i = \alpha f_i + \frac{1}{4} \alpha \sum_{j:(i,j) \in \mathcal{E}} f_j$$

3 ₀	3 ₁	2 ₂	1	0
0 ₂	0 ₂	1 ₀	3	1
3 ₀	1 ₁	2 ₂	2	3
2	0	0	2	2
2	0	0	0	1

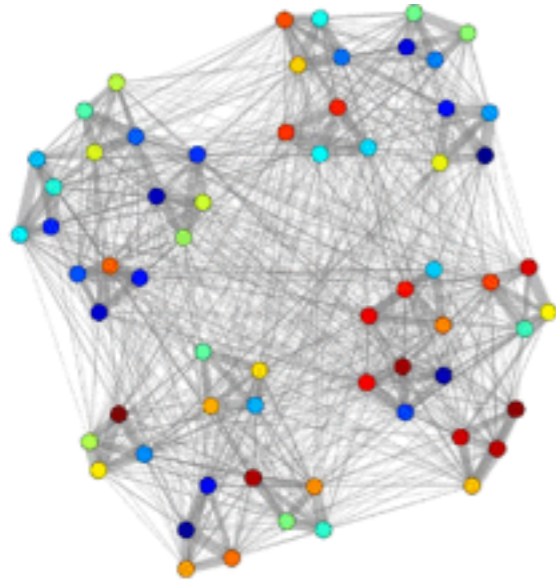
12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0



Pooling on graphs

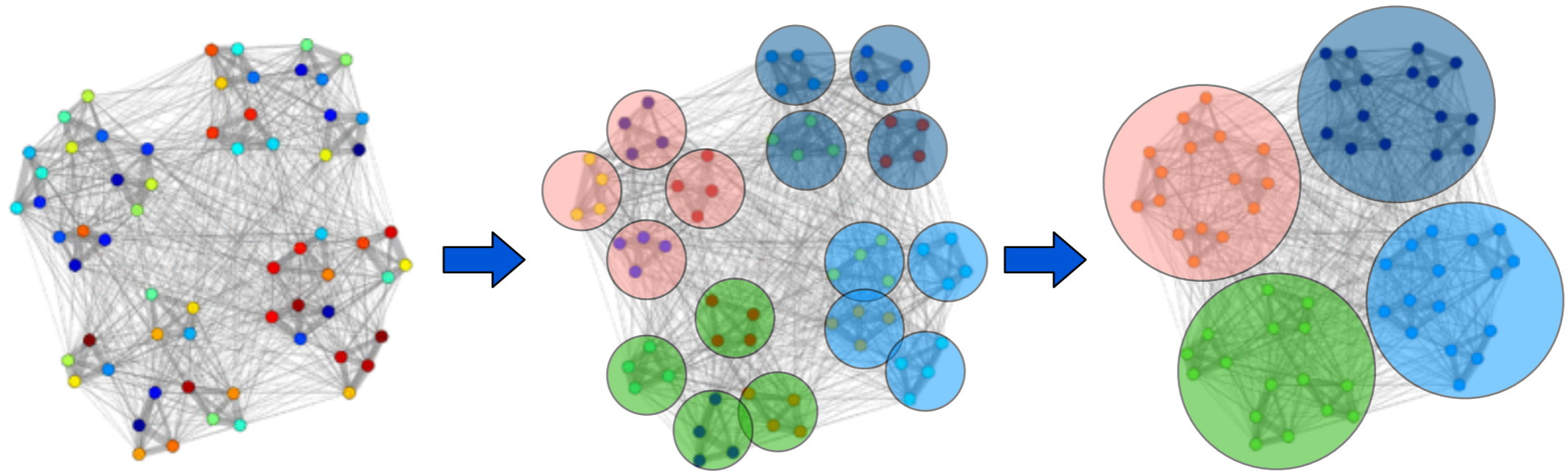


Pooling on graphs



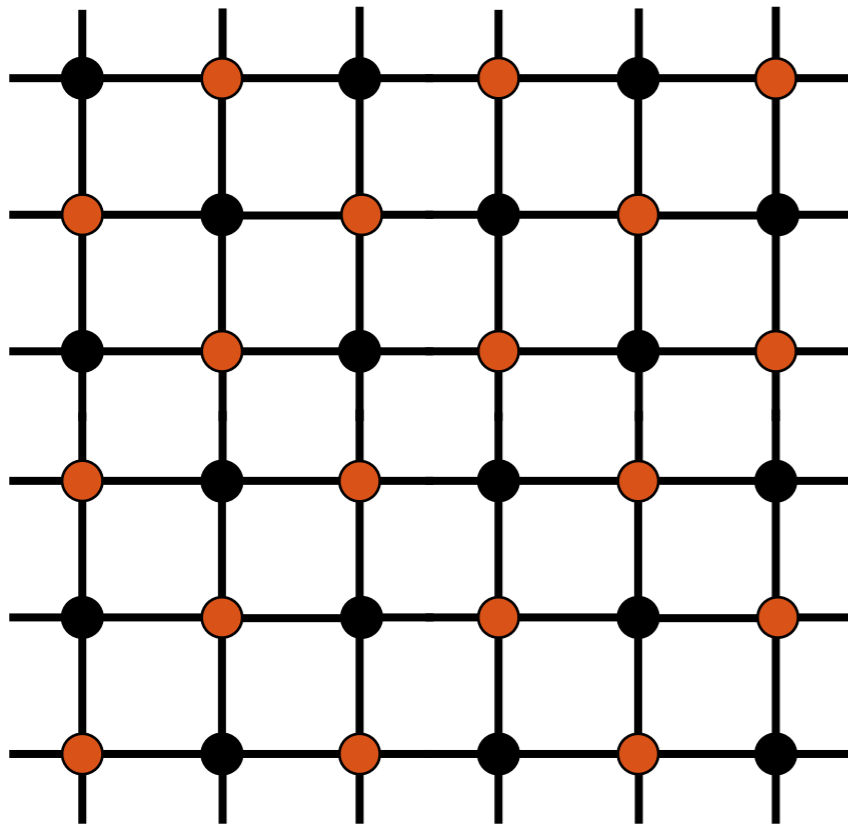
- pooling = downsampling on graphs, but how?

Pooling on graphs



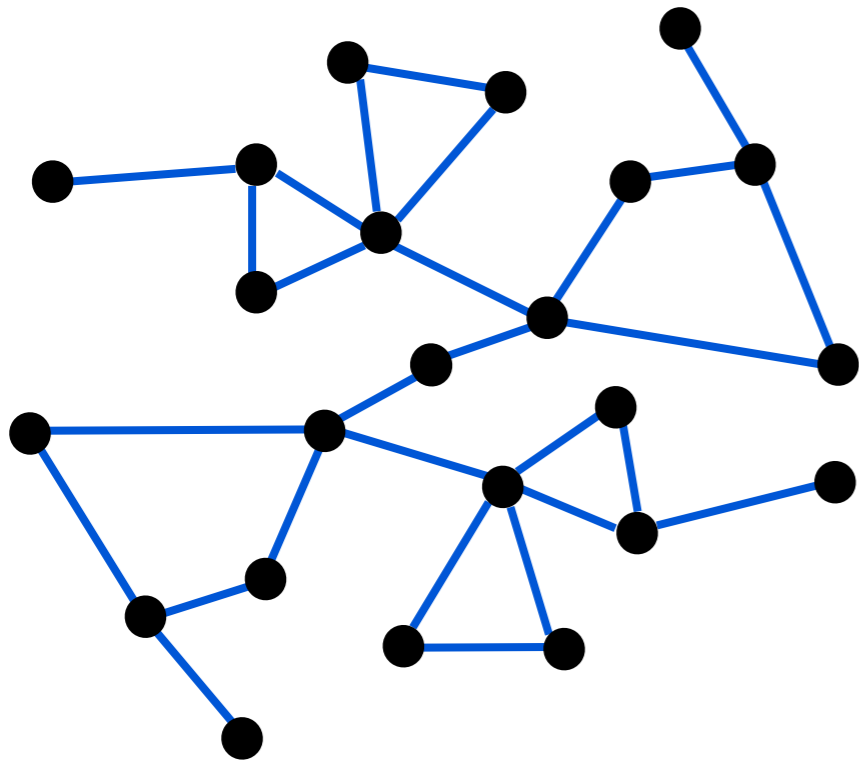
- pooling = downsampling on graphs, but how?
- natural idea: graph coarsening

Pooling on graphs



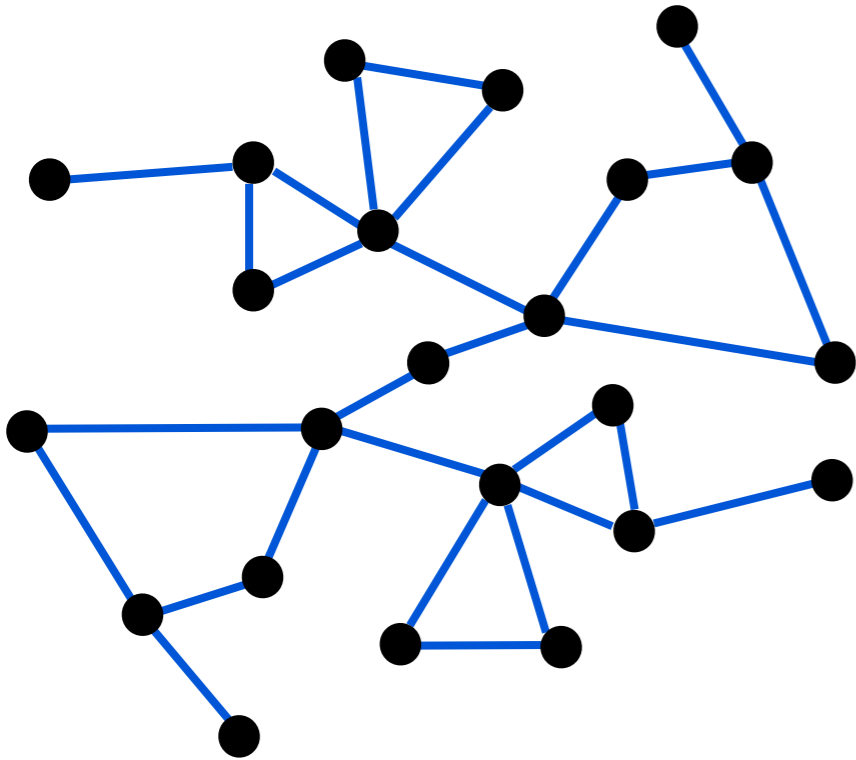
- coarsening is straightforward on regular grids

Pooling on graphs



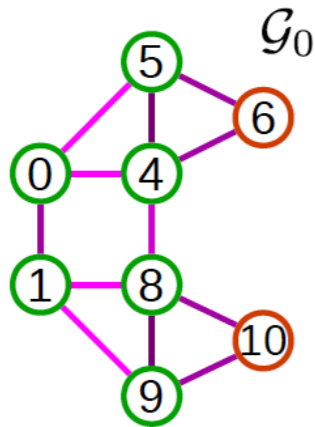
- coarsening is straightforward on regular grids
- not so much on irregular graphs

Pooling on graphs



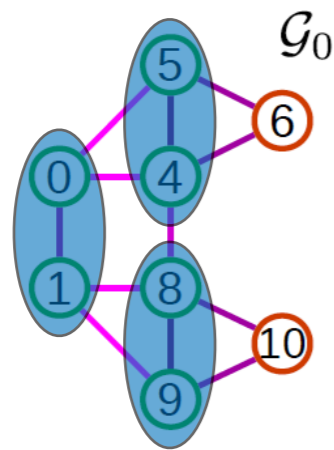
- coarsening is straightforward on regular grids
- not so much on irregular graphs
- can be achieved via node clustering
 - multi-level partitioning
 - roughly fixed downsampling factor (e.g., 2)
 - need for efficiency

Pooling on graphs



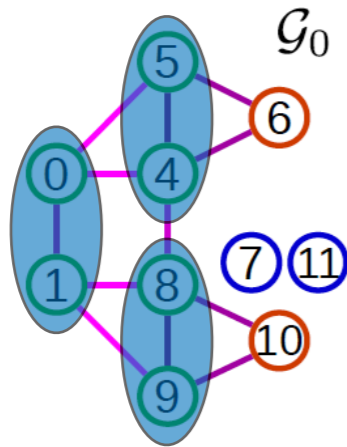
- pooling based on Graclus algorithm (Dhillon et al. 2007)
 - local greedy way of merging vertices: maximising $w_{ij}(1/d_i + 1/d_j)$
 - adding artificial vertices to ensure two children for each node

Pooling on graphs



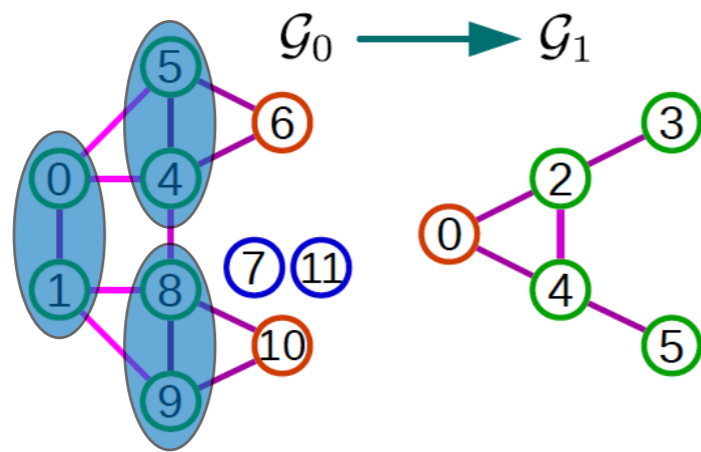
- pooling based on Graclus algorithm (Dhillon et al. 2007)
 - local greedy way of merging vertices: maximising $w_{ij}(1/d_i + 1/d_j)$
 - adding artificial vertices to ensure two children for each node

Pooling on graphs



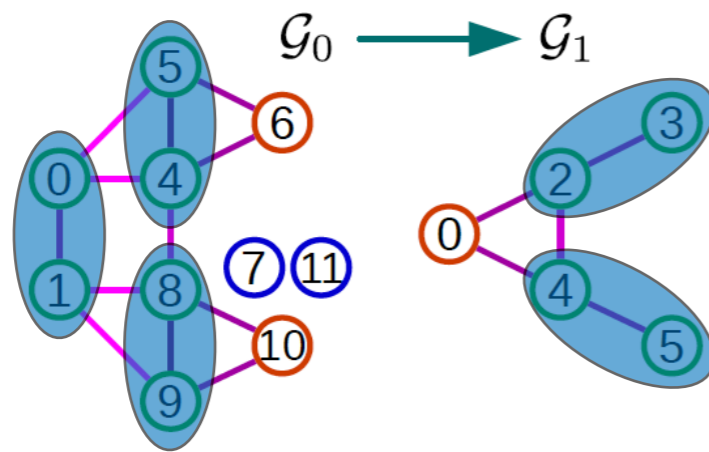
- pooling based on Graclus algorithm (Dhillon et al. 2007)
 - local greedy way of merging vertices: maximising $w_{ij}(1/d_i + 1/d_j)$
 - adding artificial vertices to ensure two children for each node

Pooling on graphs



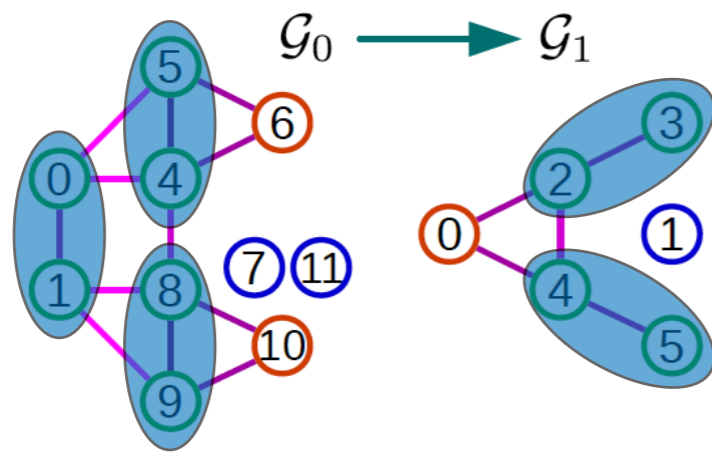
- pooling based on Graclus algorithm (Dhillon et al. 2007)
 - local greedy way of merging vertices: maximising $w_{ij}(1/d_i + 1/d_j)$
 - adding artificial vertices to ensure two children for each node

Pooling on graphs



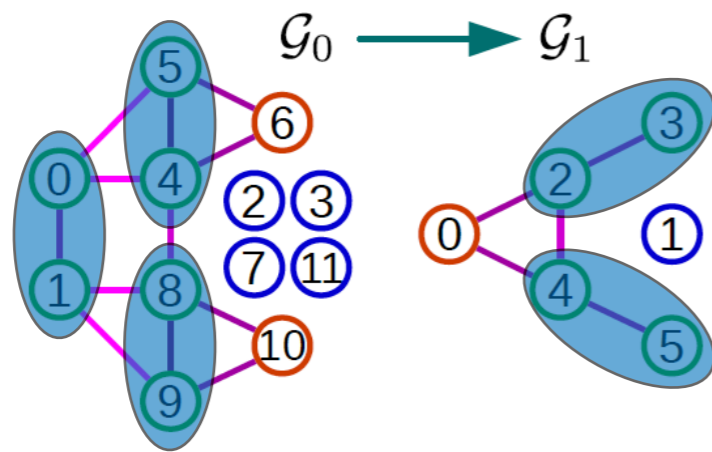
- pooling based on Graclus algorithm (Dhillon et al. 2007)
 - local greedy way of merging vertices: maximising $w_{ij}(1/d_i + 1/d_j)$
 - adding artificial vertices to ensure two children for each node

Pooling on graphs



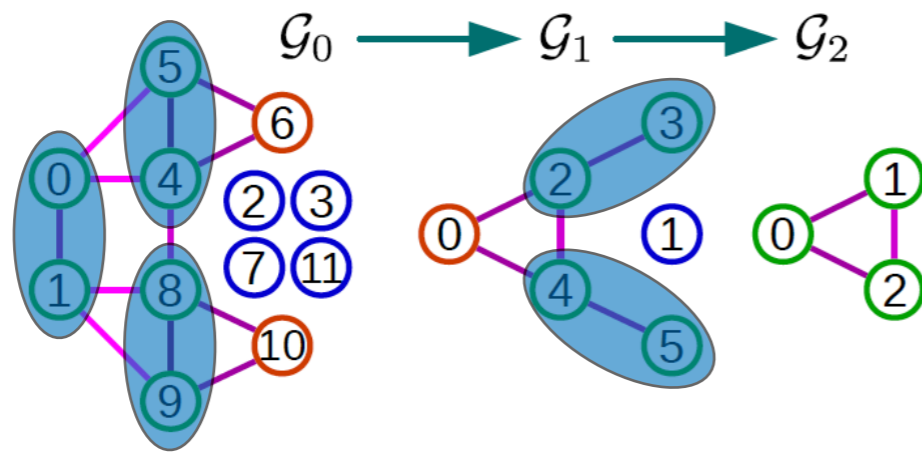
- pooling based on Graclus algorithm (Dhillon et al. 2007)
 - local greedy way of merging vertices: maximising $w_{ij}(1/d_i + 1/d_j)$
 - adding artificial vertices to ensure two children for each node

Pooling on graphs



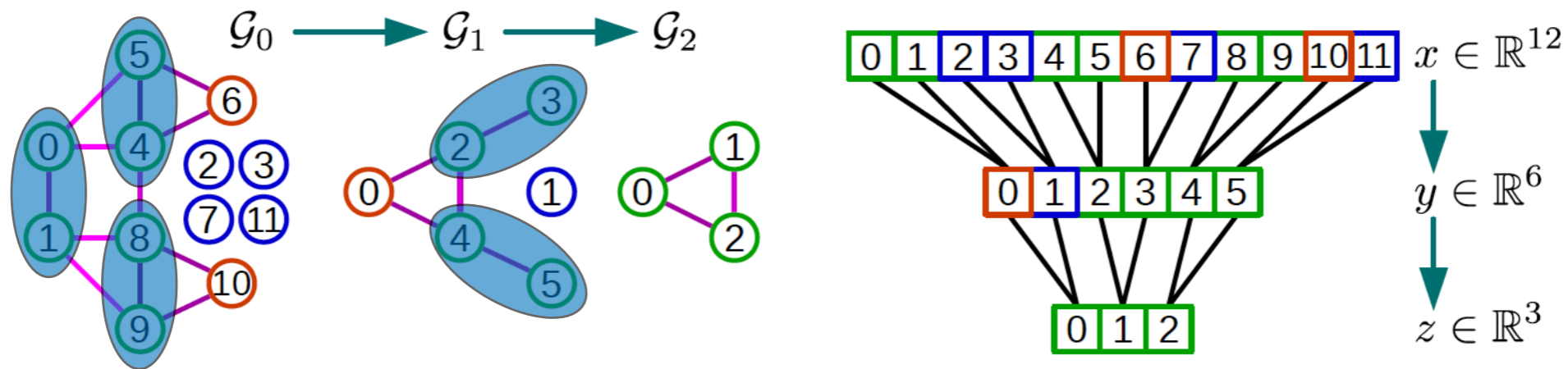
- pooling based on Graclus algorithm (Dhillon et al. 2007)
 - local greedy way of merging vertices: maximising $w_{ij}(1/d_i + 1/d_j)$
 - adding artificial vertices to ensure two children for each node

Pooling on graphs



- pooling based on Graclus algorithm (Dhillon et al. 2007)
 - 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

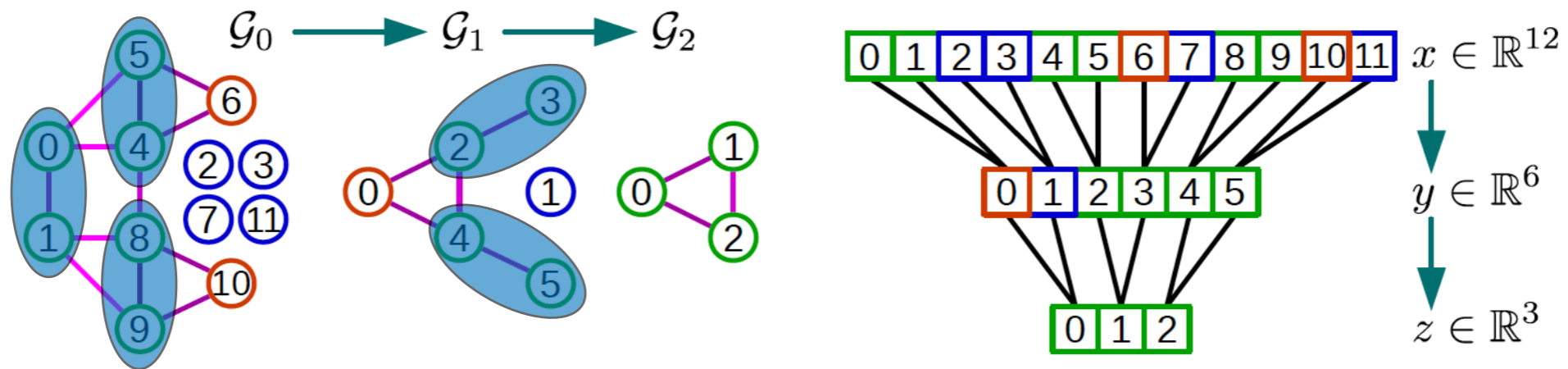
Pooling on graphs



Defferrard et al. 2016

- pooling based on Graclus algorithm (Dhillon et al. 2007)
 - 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

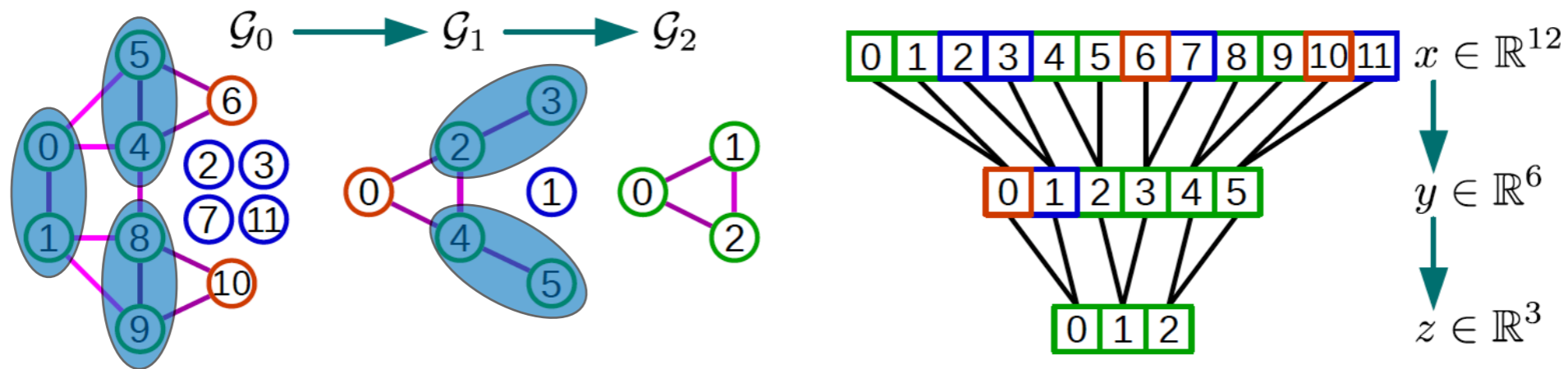
Pooling on graphs



Defferrard et al. 2016

- pooling based on Graclus algorithm (Dhillon et al. 2007)
 - 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: [$\max(0,1)$ $\max(4,5,6)$ $\max(8,9,10)$]

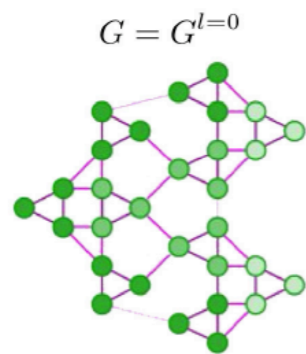
Pooling on graphs



Defferrard et al. 2016

- pooling based on Graclus algorithm (Dhillon et al. 2007)
 - 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

CNN on graphs: Graph classification



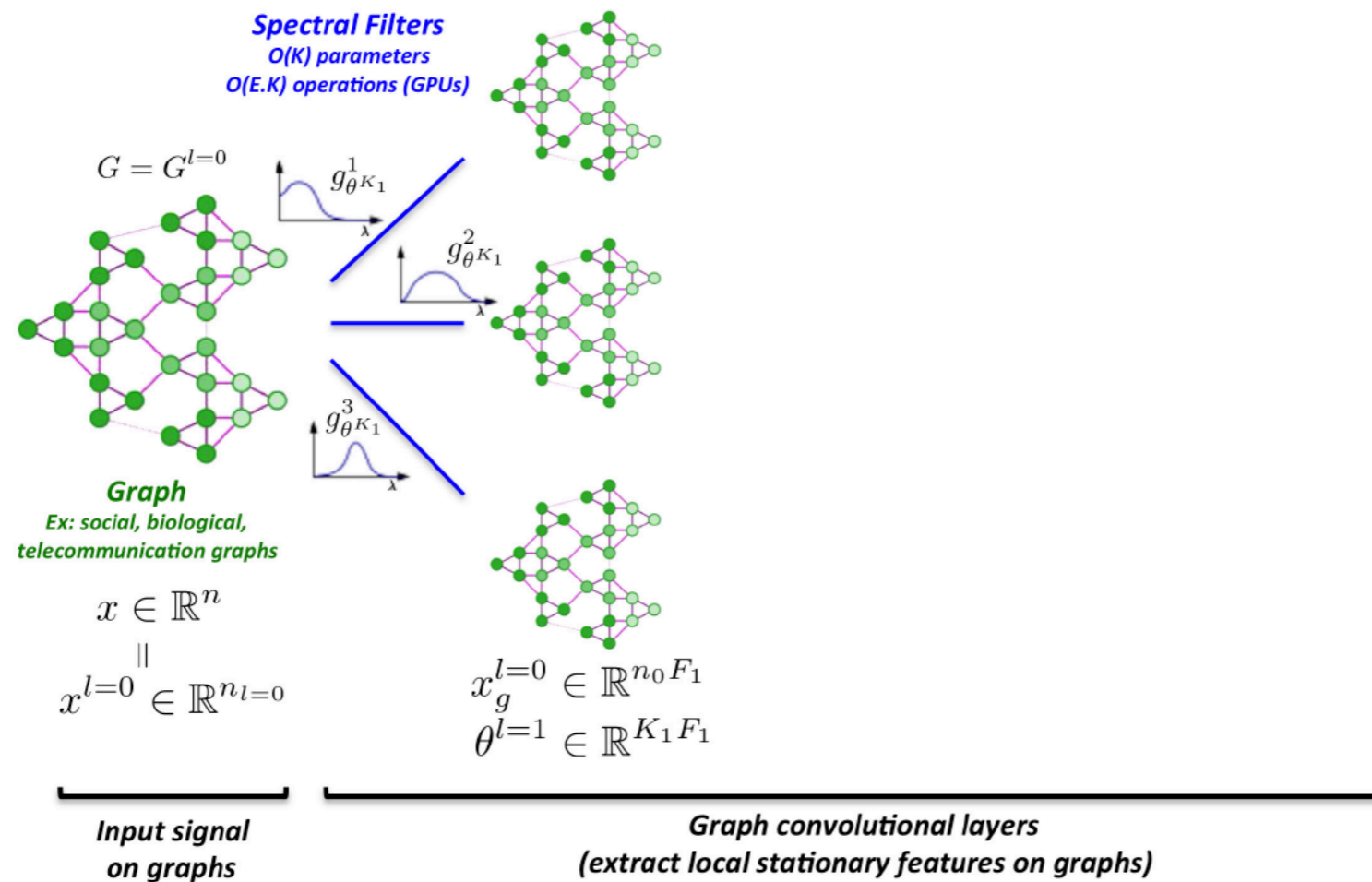
Graph

*Ex: social, biological,
telecommunication graphs*

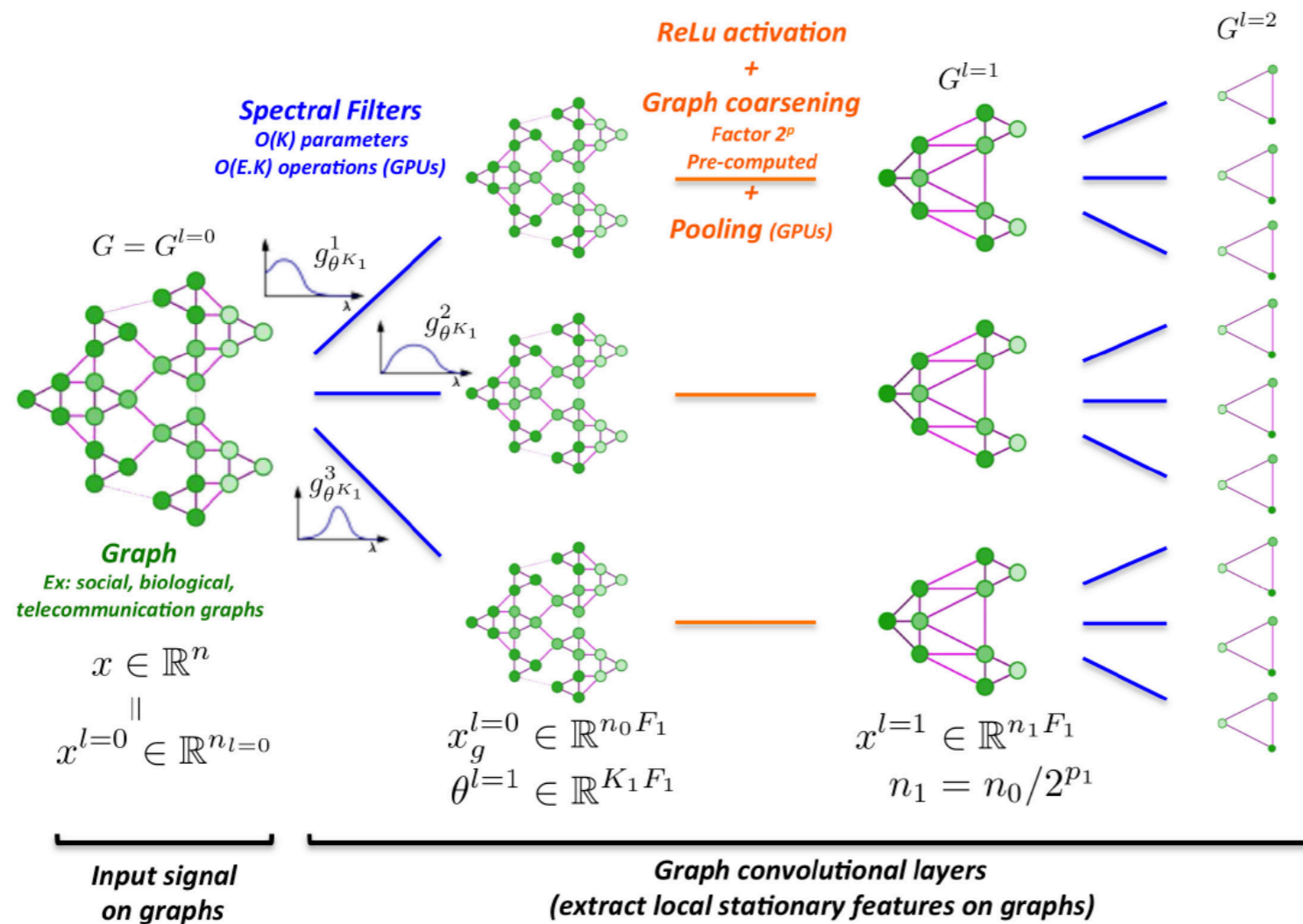
$$\begin{aligned} x &\in \mathbb{R}^n \\ &\parallel \\ x^{l=0} &\in \mathbb{R}^{n_{l=0}} \end{aligned}$$

**Input signal
on graphs**

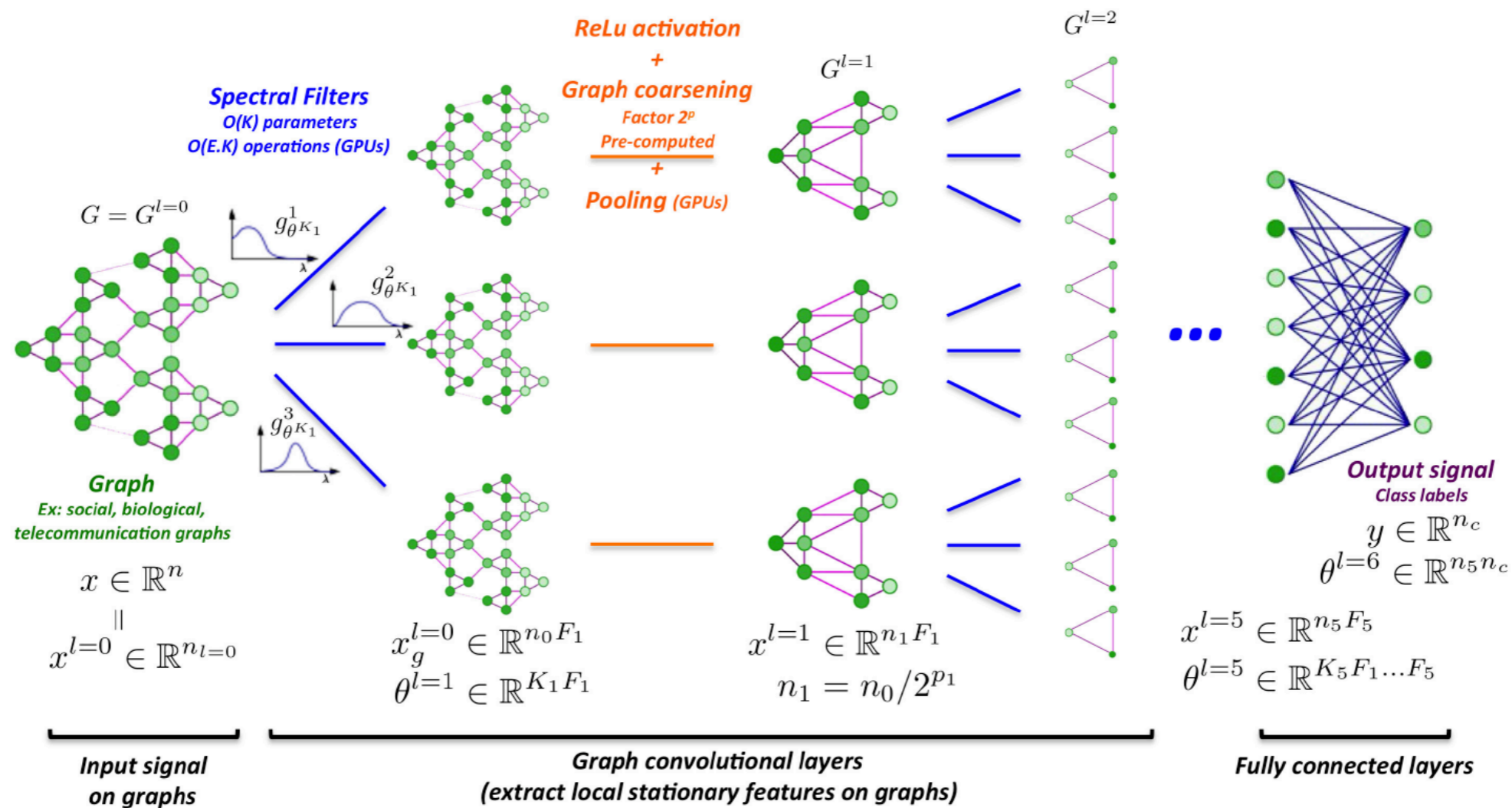
CNN on graphs: Graph classification



CNN on graphs: Graph classification

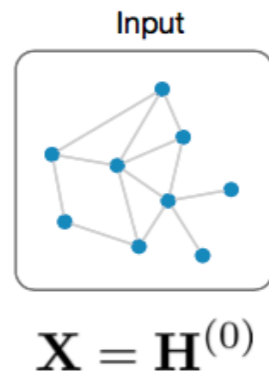


CNN on graphs: Graph classification



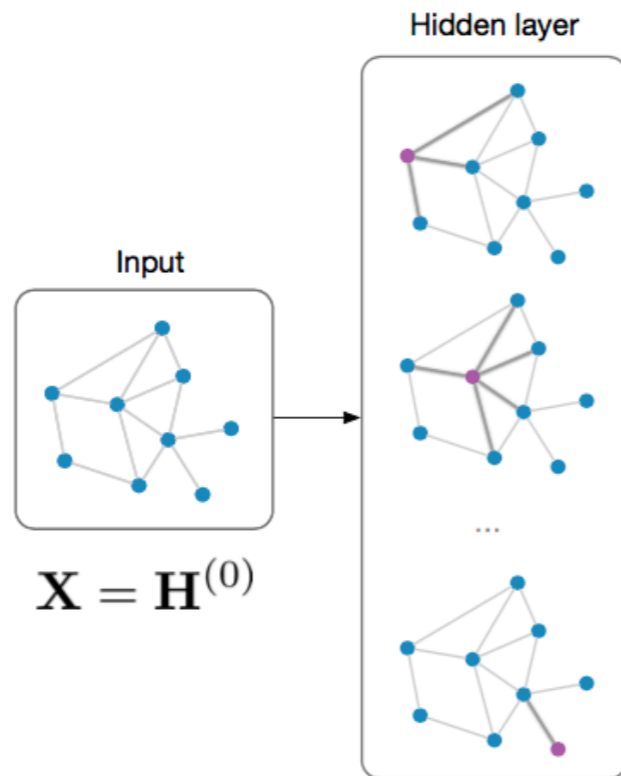
CNN on graphs: Node classification

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



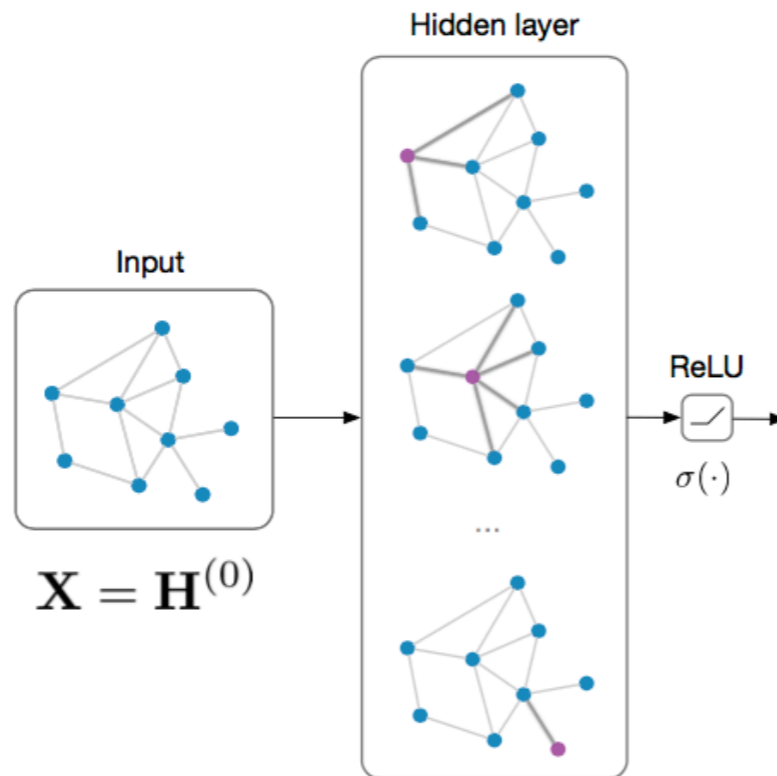
CNN on graphs: Node classification

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



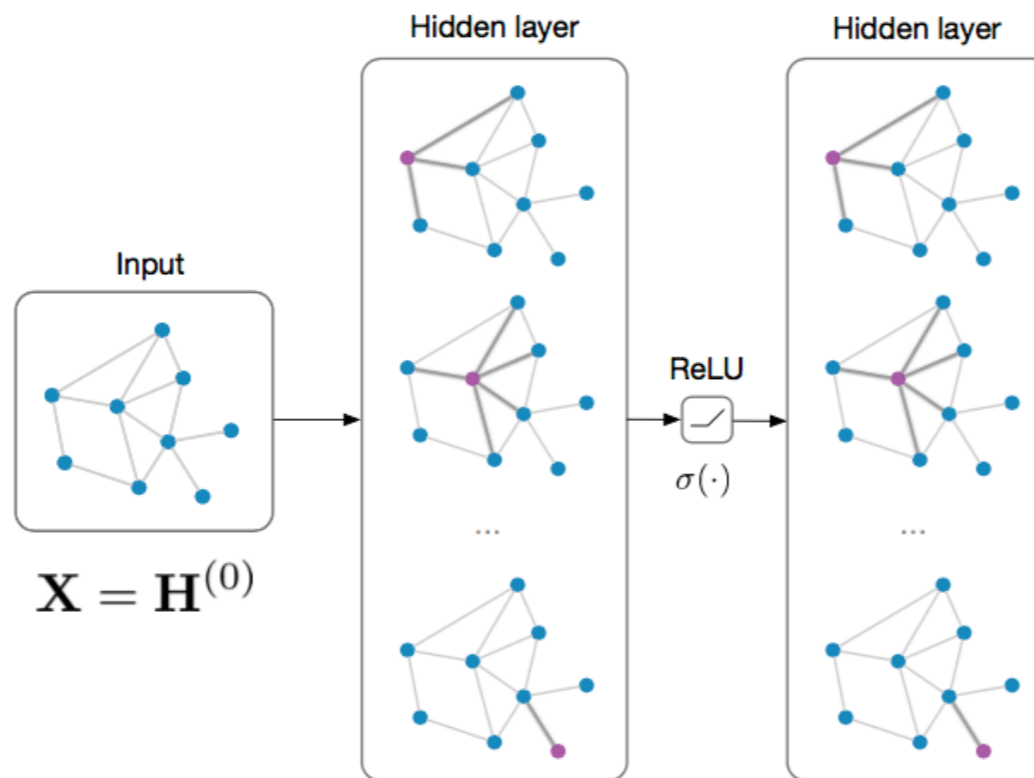
CNN on graphs: Node classification

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



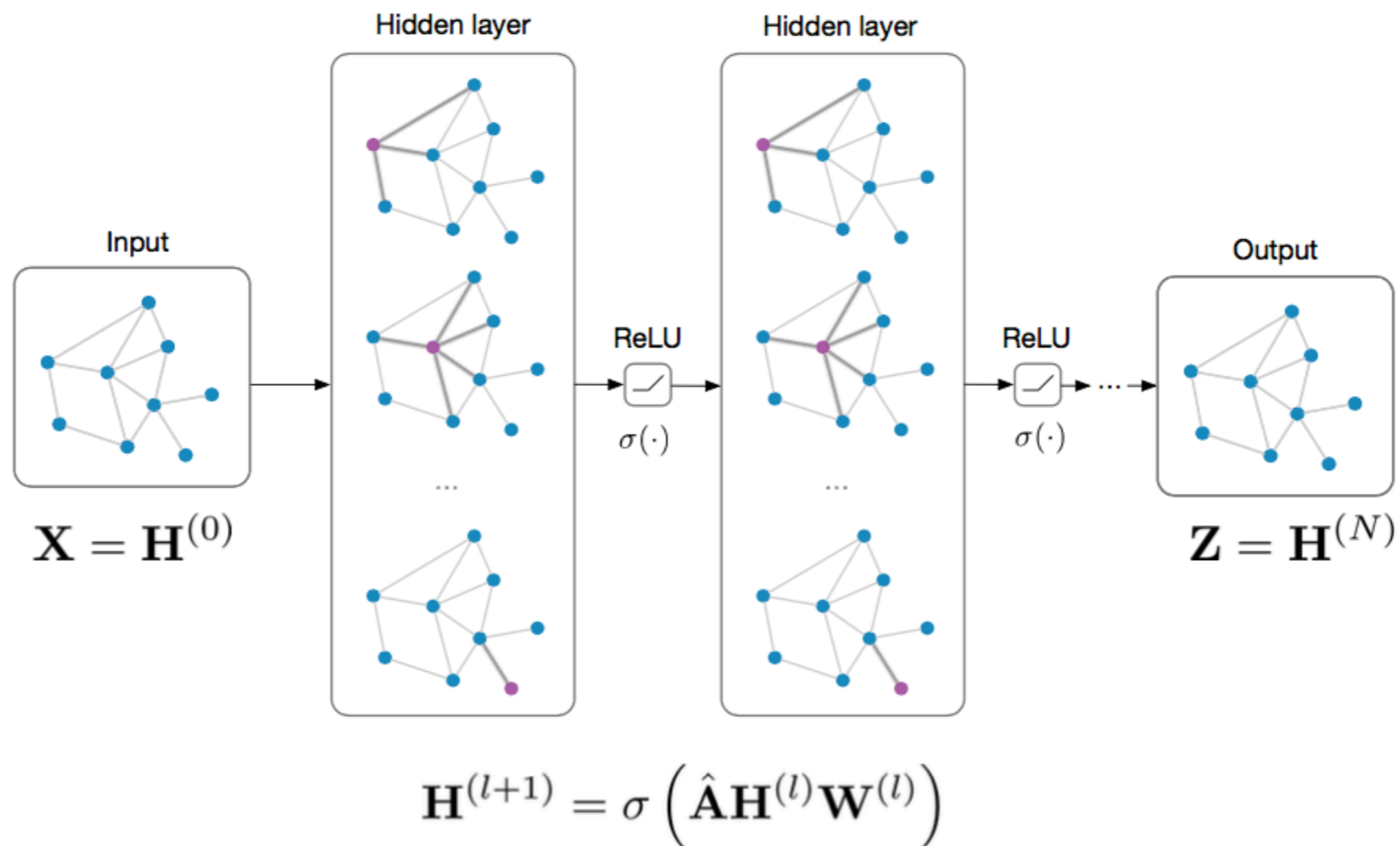
CNN on graphs: Node classification

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



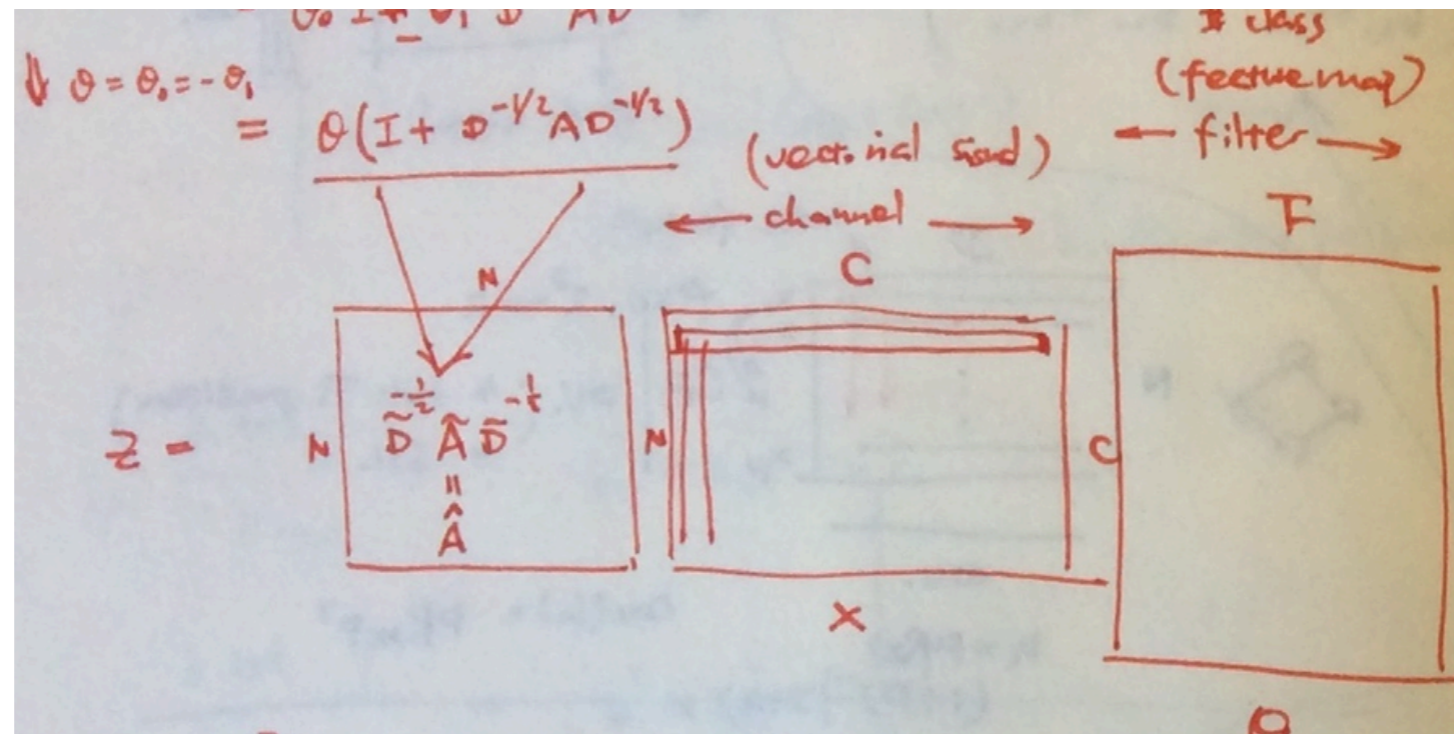
CNN on graphs: Node classification

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



CNN on graphs: Node classification

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



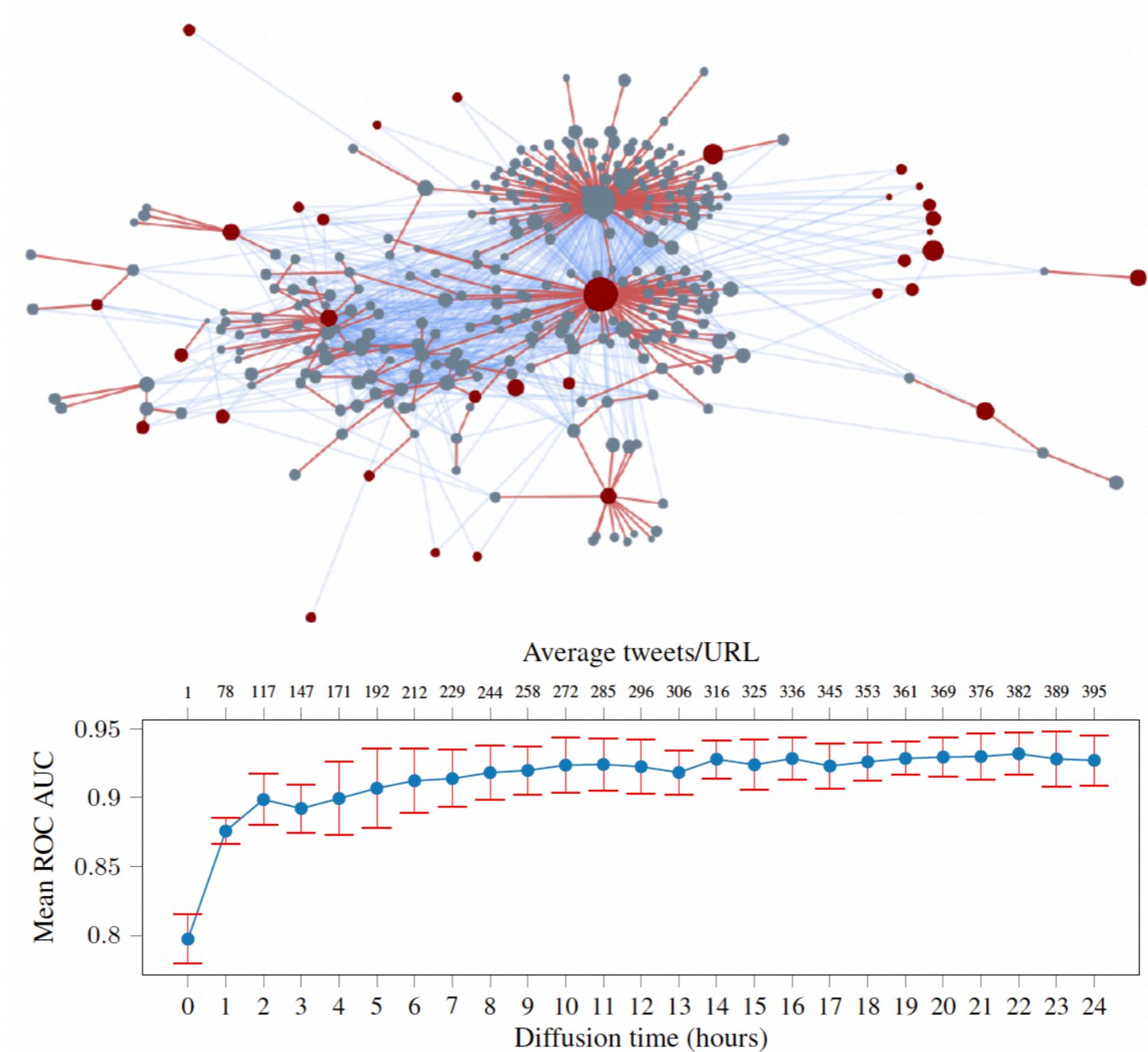
$$\mathbf{H}^{(l+1)} = \sigma(\hat{\mathbf{A}}\mathbf{H}^{(l)}\mathbf{W}^{(l)})$$

Application I: Document classification

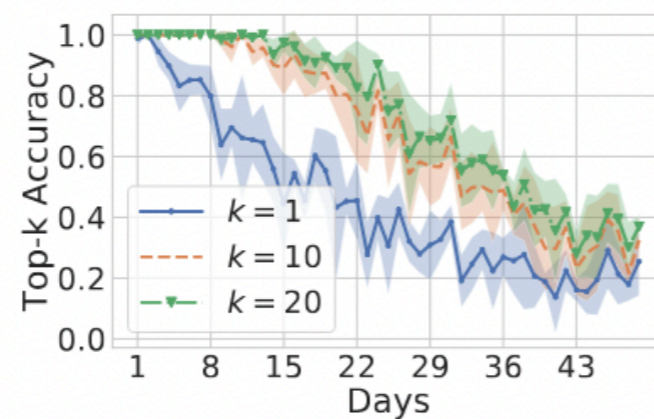
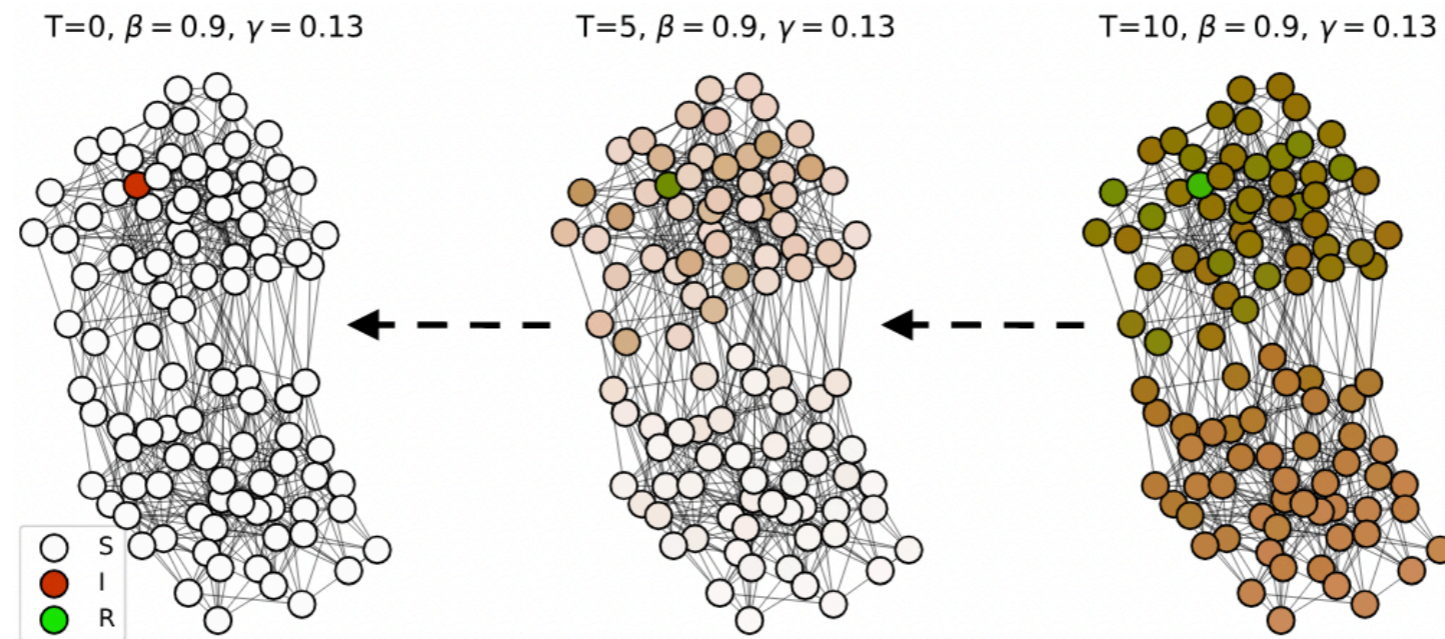


Method	Citeseer	Cora	Pubmed	NELL
ManiReg [3]	60.1	59.5	70.7	21.8
SemiEmb [28]	59.6	59.0	71.1	26.7
LP [32]	45.3	68.0	63.0	26.5
DeepWalk [22]	43.2	67.2	65.3	58.1
ICA [18]	69.1	75.1	73.9	23.1
Planetoid* [29]	64.7 (26s)	75.7 (13s)	77.2 (25s)	61.9 (185s)
GCN (this paper)	70.3 (7s)	81.5 (4s)	79.0 (38s)	66.0 (48s)
GCN (rand. splits)	67.9 ± 0.5	80.1 ± 0.5	78.9 ± 0.7	58.4 ± 1.7

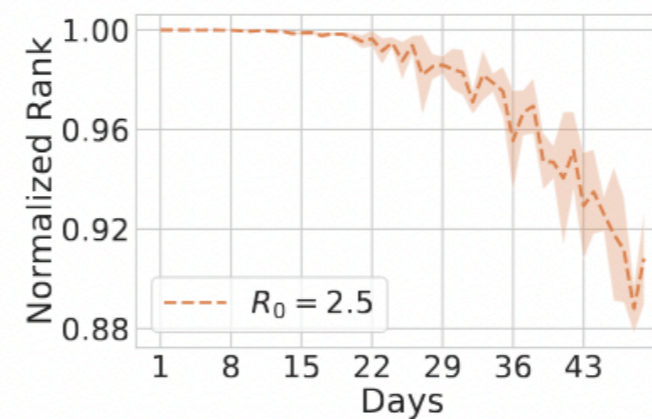
Application II: Fake news detection



Application III: Finding patient zero

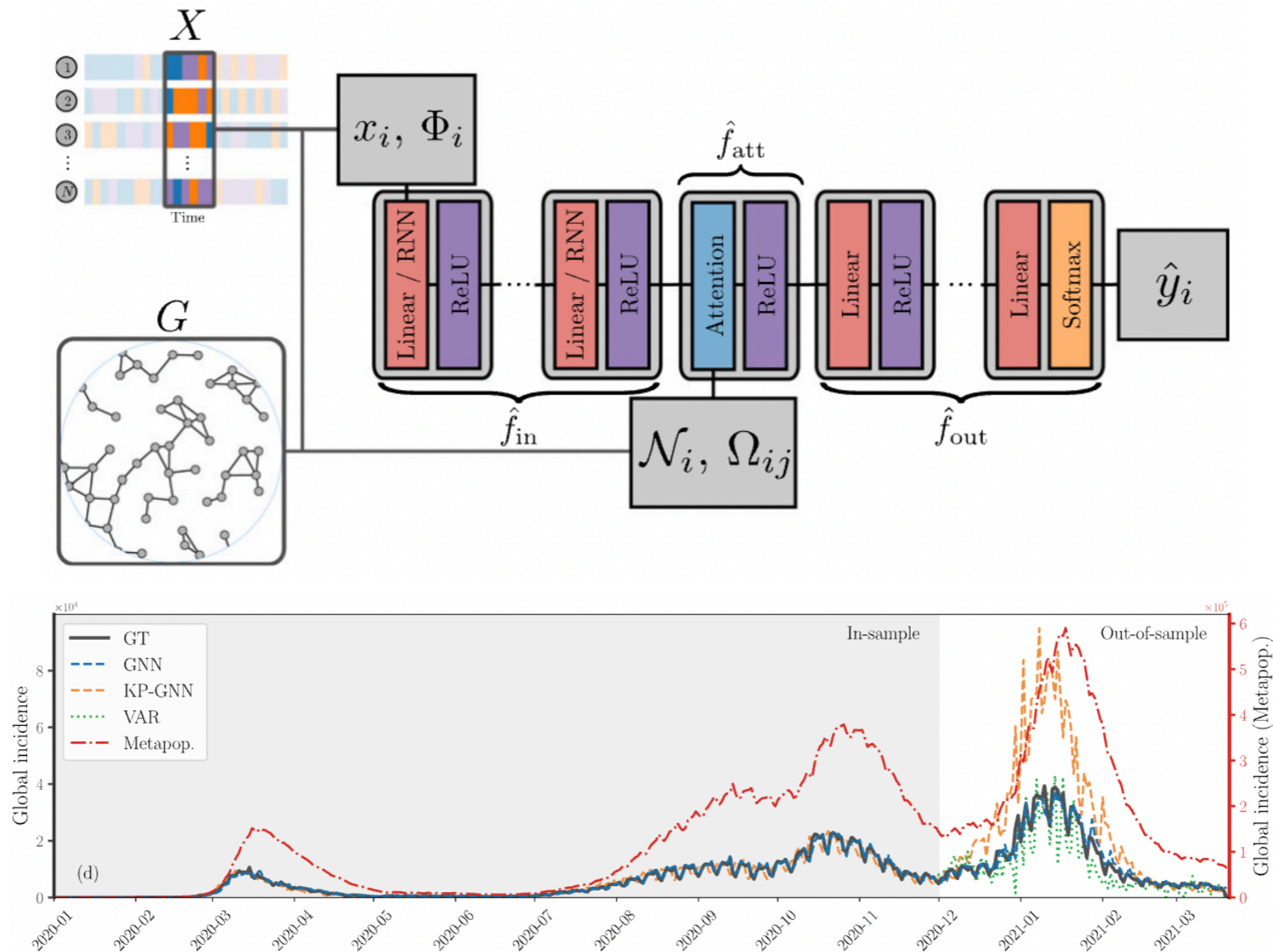


(a) Top-k accuracy

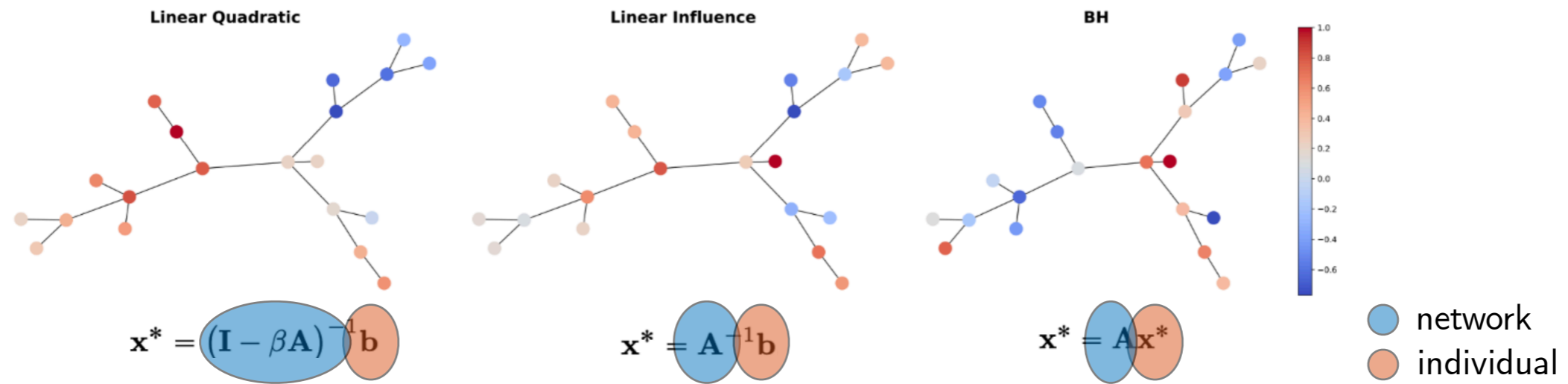


(b) Normalized rank

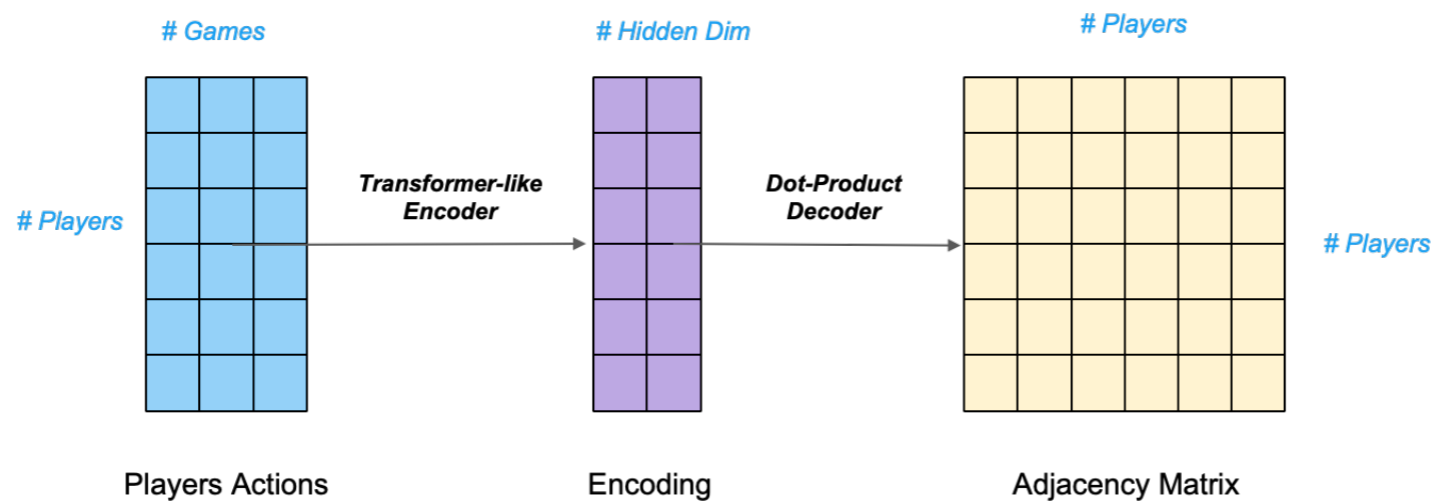
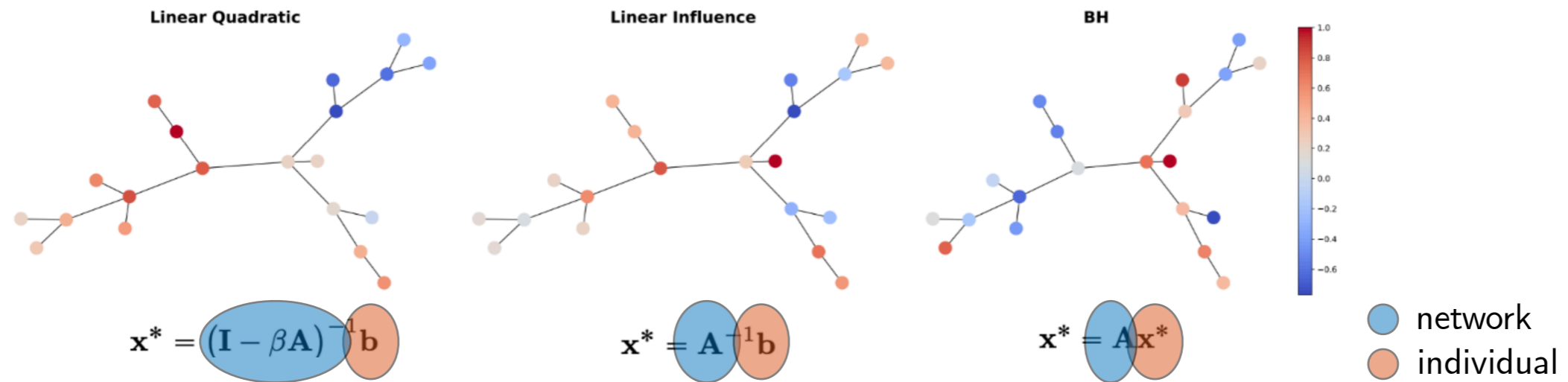
Application IV: Learning contagion dynamics



Application V: Learning social interactions

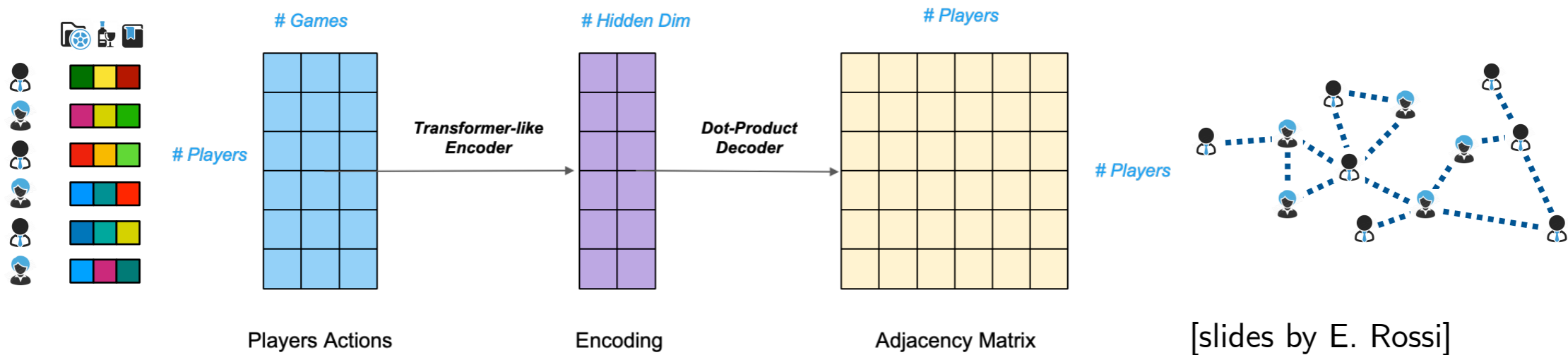
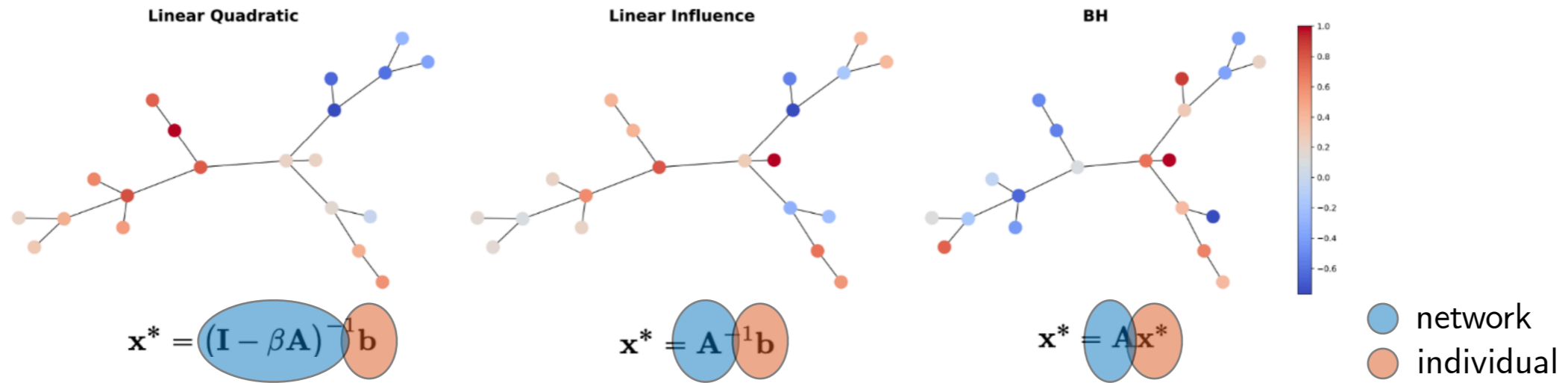


Application V: Learning social interactions

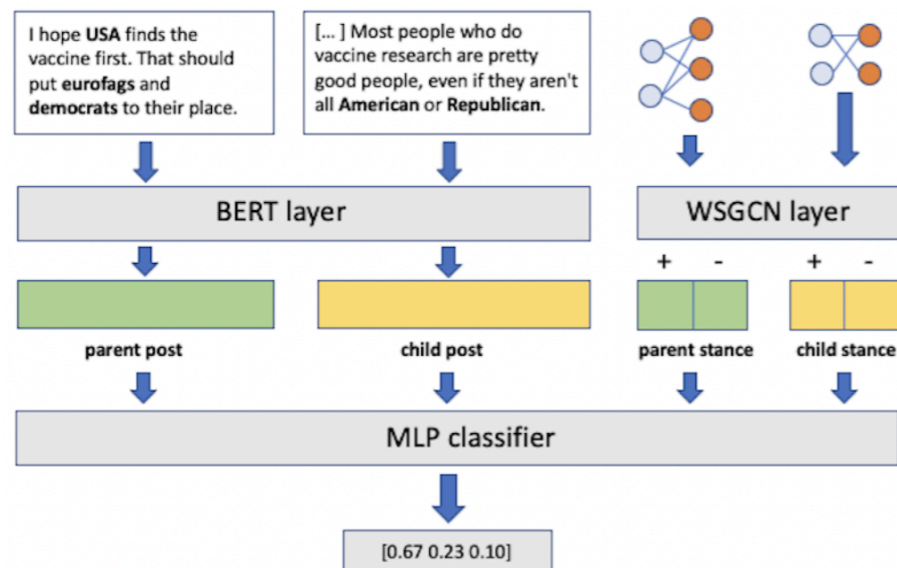


[slides by E. Rossi]

Application V: Learning social interactions



Application VI: Language and social media analysis



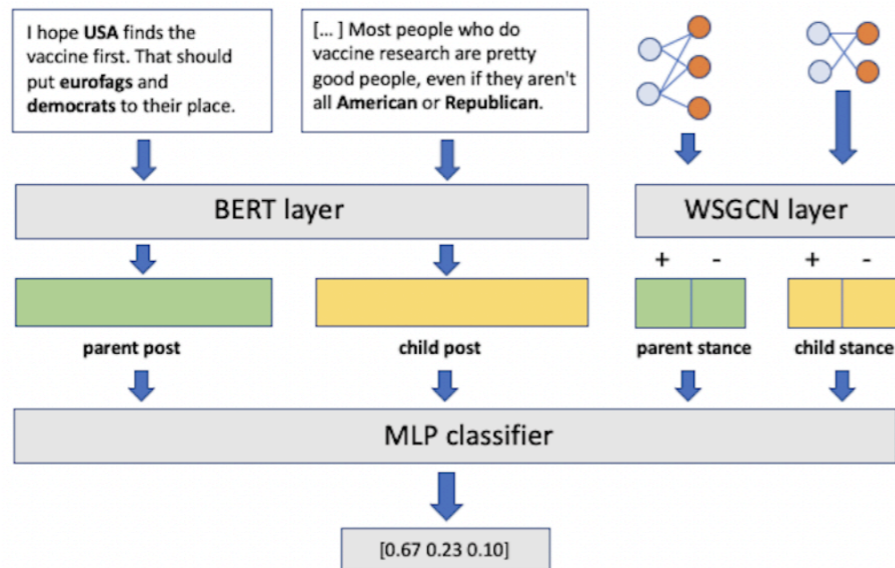
disagreement prediction

Hofmann et al., "Modeling ideological salience and framing in polarized online groups with graph neural networks and structured sparsity," NAACL, 2022.

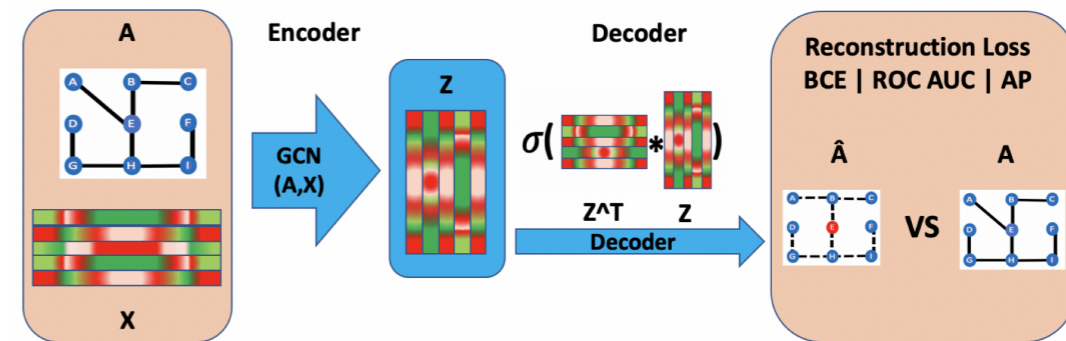
Zhang et al., "Predicting polarisation of dynamic social networks via graph auto-encoders," IC2S2, 2023.

Lorge et al., "STEntConv: Predicting disagreement between reddit users with stance detection and a signed graph convolutional network," LREC-Coling 2024.

Application VI: Language and social media analysis



disagreement prediction



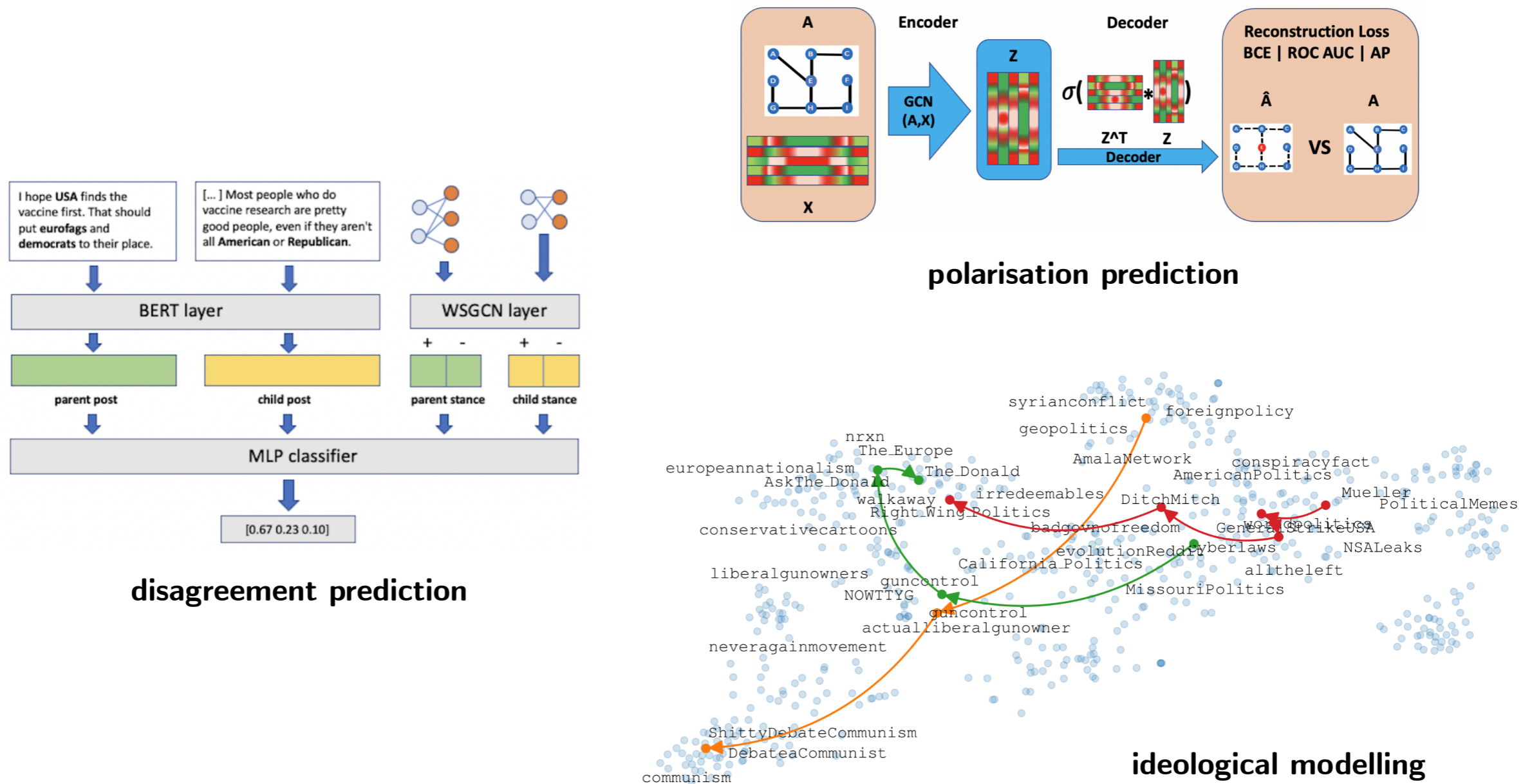
polarisation prediction

Hofmann et al., “Modeling ideological salience and framing in polarized online groups with graph neural networks and structured sparsity,” NAACL, 2022.

Zhang et al., “Predicting polarisation of dynamic social networks via graph auto-encoders,” IC2S2, 2023.

Lorge et al., “STEntConv: Predicting disagreement between reddit users with stance detection and a signed graph convolutional network,” LREC-Coling 2024.

Application VI: Language and social media analysis



Hofmann et al., "Modeling ideological salience and framing in polarized online groups with graph neural networks and structured sparsity," NAACL, 2022.

Zhang et al., "Predicting polarisation of dynamic social networks via graph auto-encoders," IC2S2, 2023.

Lorge et al., "STEntConv: Predicting disagreement between reddit users with stance detection and a signed graph convolutional network," LREC-Coling 2024.

Some Final Thoughts

Summary

- Graph machine learning
 - fast-growing field that extends data analysis to non-Euclidean domain
 - highly interdisciplinary: machine learning, signal processing, harmonic analysis, applies statistics, differential geometry
- Limitations and open challenges
 - models on directed and signed graphs
 - models for temporal dynamics and online/adaptive settings
 - construction/refinement of initial graphs
 - robustness & generalisation & scalability
 - interpretability & causal inference
 - expect more applications in social sciences & economics!

