
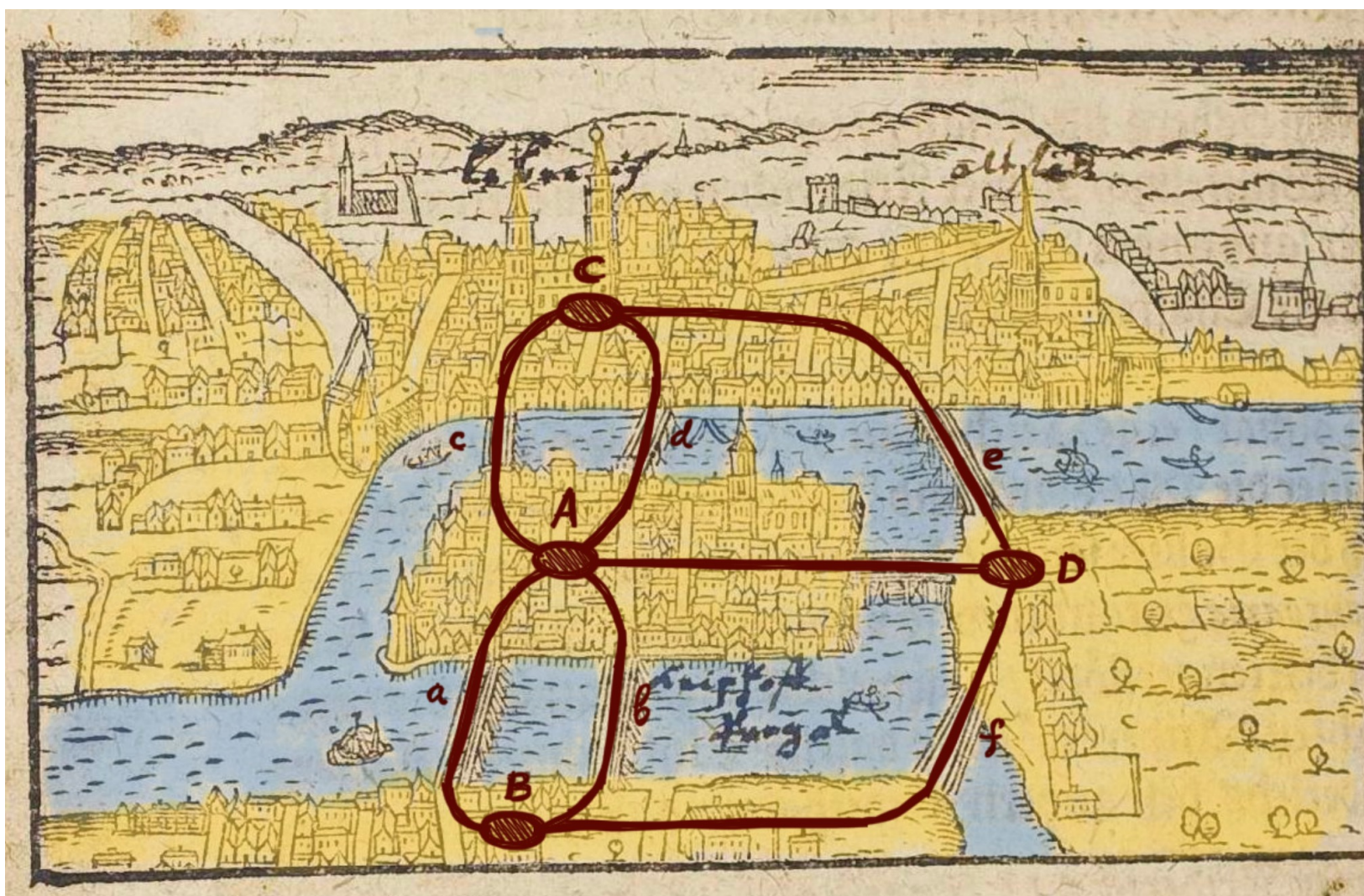


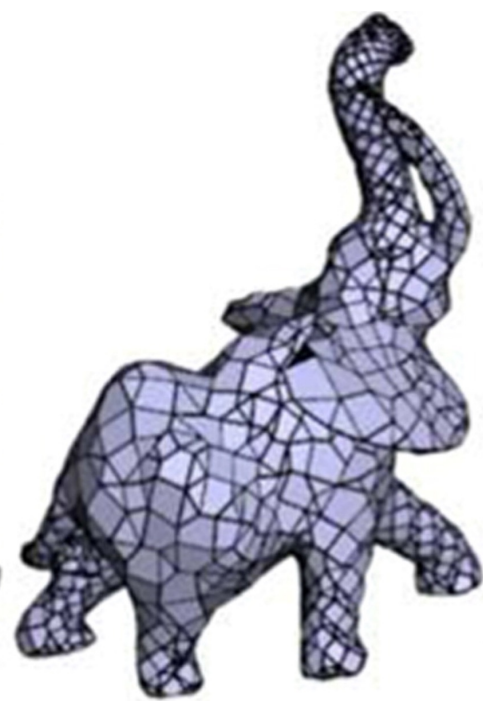
# A Continuous Perspective on Graph Neural Networks

Ben Chamberlain

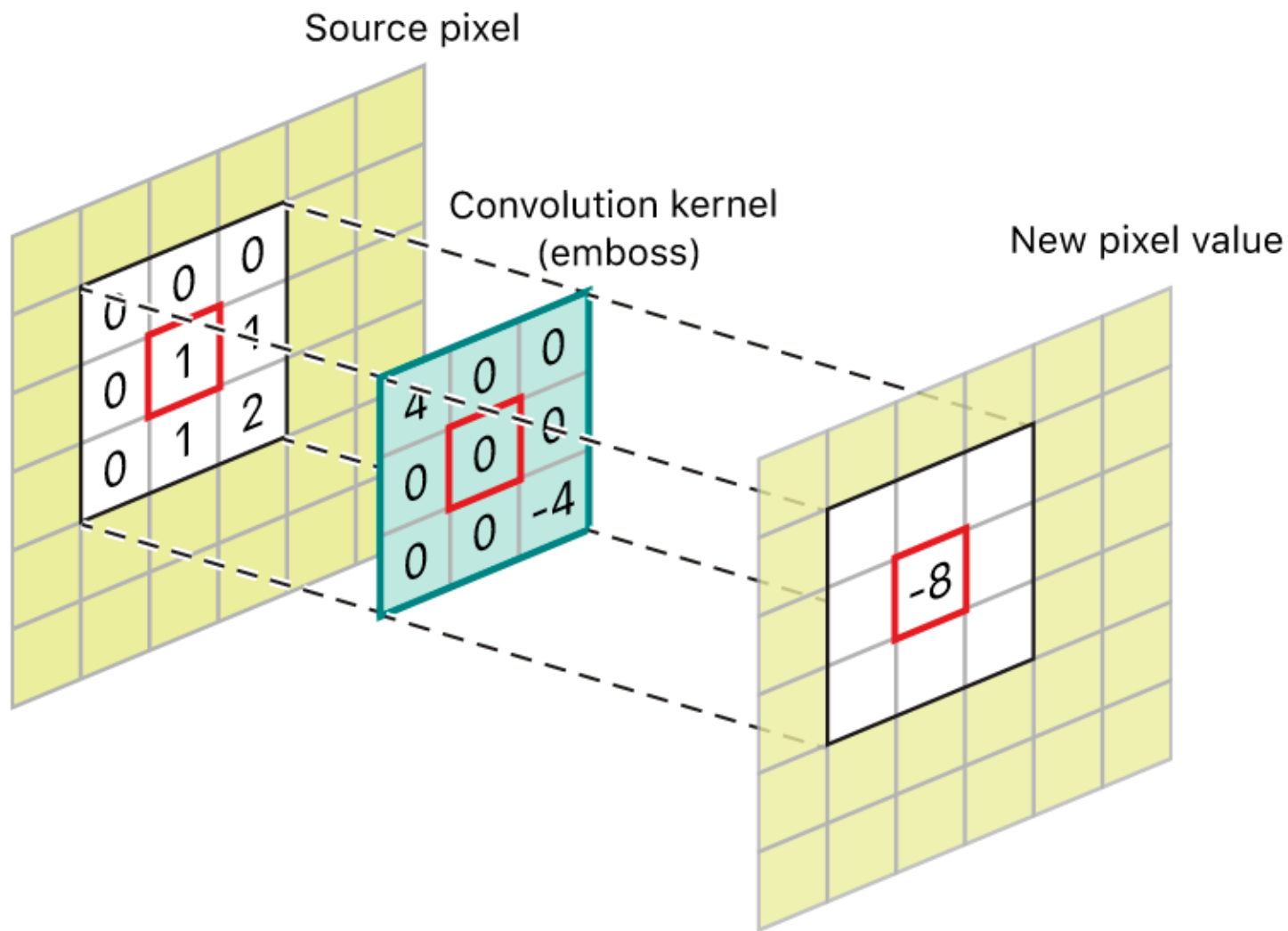
 @b\_p\_chamberlain





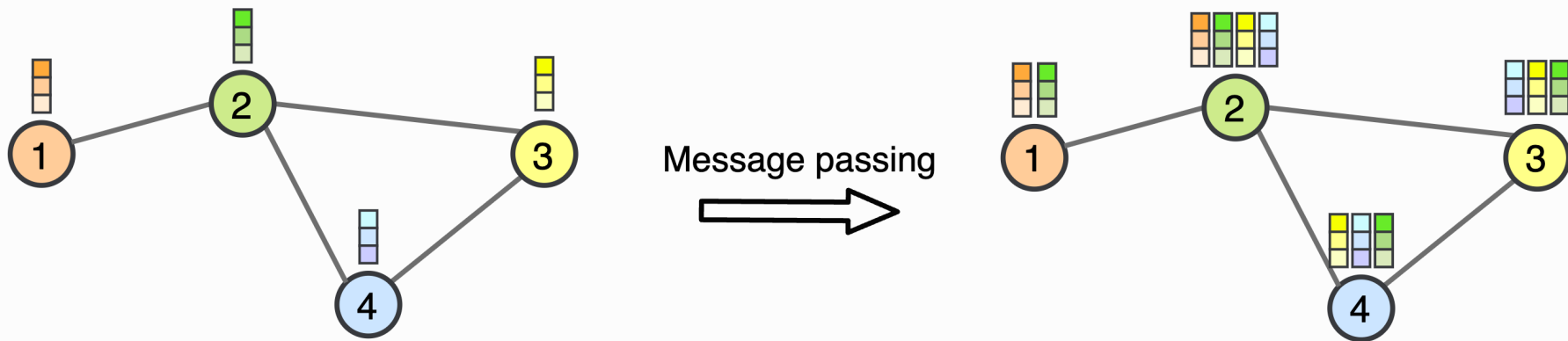








# Message Passing GNNs



- Most GNNs are Message Passing GNNs
- Message Passing GNNs interleave two steps
  1. Propagation along edge (message passing)
  2. Shared feature transformation by a neural network

$$\mathbf{H}^{(k+1)} = \sigma \left( \underbrace{\mathbf{D}^{-\frac{1}{2}} \tilde{\mathbf{A}} \mathbf{D}^{-\frac{1}{2}}}_{\text{Message passing}} \mathbf{H}^{(k)} \underbrace{\mathbf{W}_k}_{\text{Feature transformation}} \right)$$

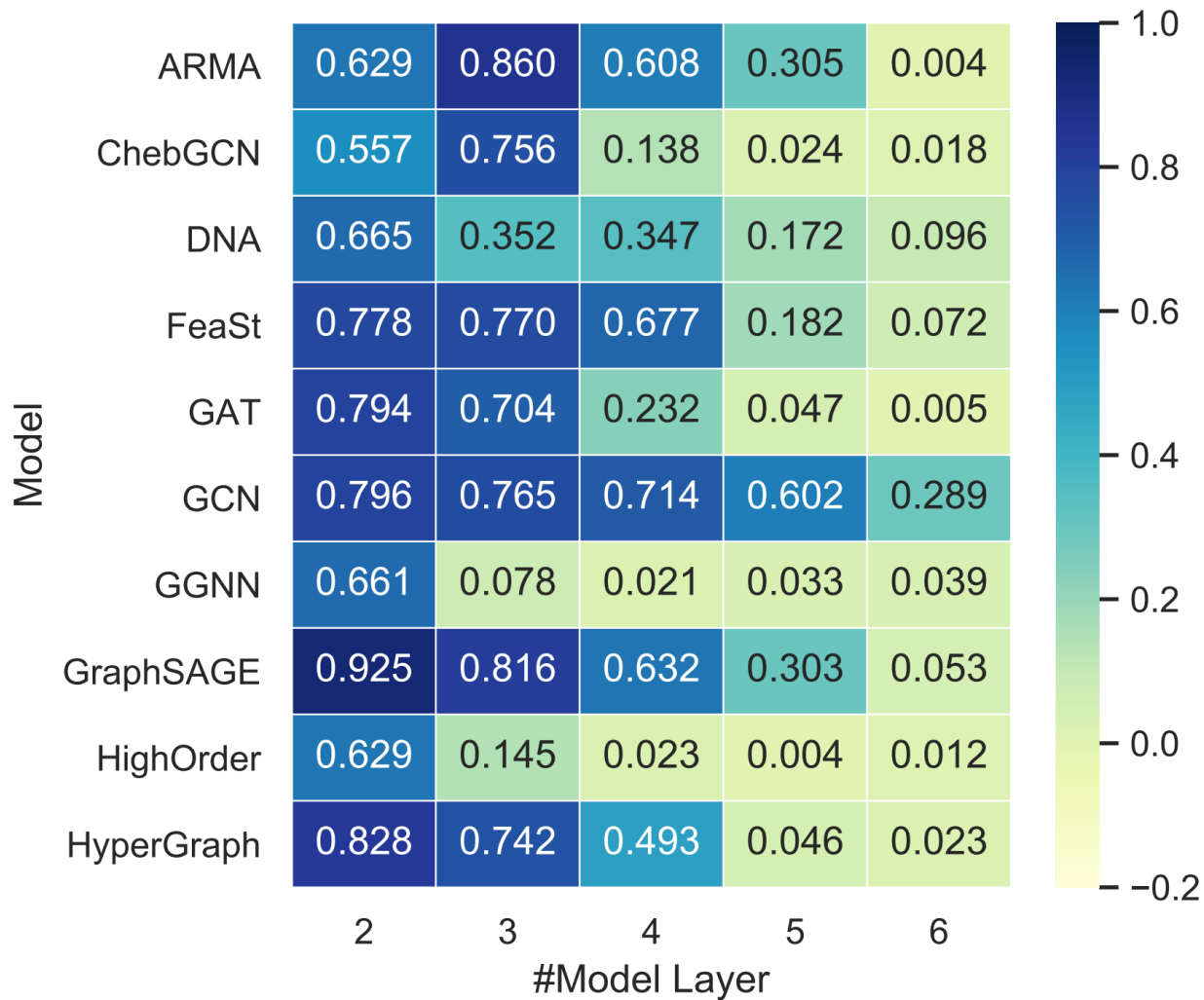
## *What's Wrong with (Message Passing) GNNs?*

- **Deep GNNs  $\Rightarrow$  over-smoothing**
- **Input graph = computational graph  $\Rightarrow$  bottlenecks & “over-squashing”**
- Graph & features incompatible structure (“heterophily”)  $\Rightarrow$  poor performance
- Limited expressive power  $\Rightarrow$  cannot detect simple structures such as cycles or cliques





# Oversmoothing

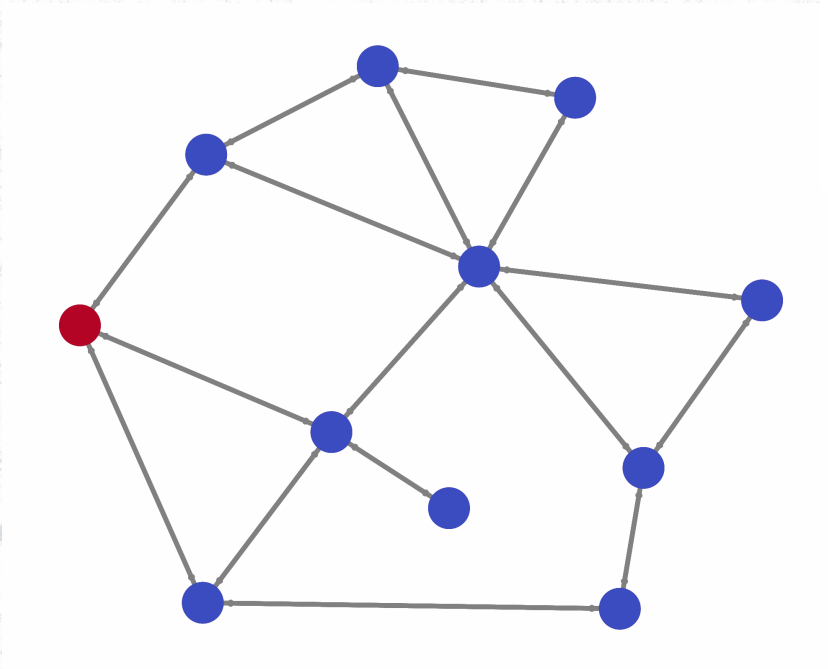


# Message Passing is Information Diffusion



Laplacian diffusion on graphs:

$$X_t = LX_{t-1}$$



$$\mathbf{E}(\mathbf{X}) = \frac{1}{v} \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} \|\mathbf{X}_i - \mathbf{X}_j\|^2$$

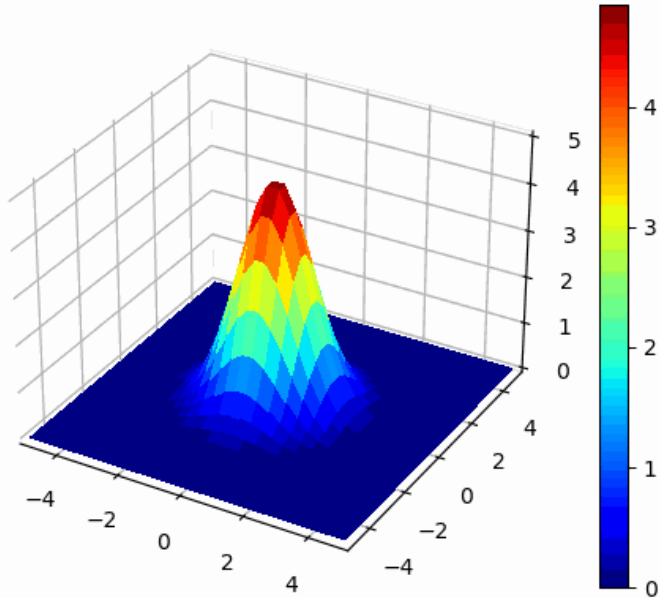


Pierre-Simon Laplace

# The continuous analogy is heat diffusion



Heat diffusion:



Joseph Fourier

$$\frac{\partial x}{\partial t} = \Delta x$$

$$\mathcal{E}(f, h) = \frac{1}{2} \int_{\mathbb{R}^n} \|\nabla f\|_h^2 dx$$

# Diffusion in Image Processing



$$\frac{\partial}{\partial t}x = c\Delta x$$

$$\frac{\partial}{\partial t}x = \operatorname{div} \left[ \frac{1}{1 + (|\nabla x|/\lambda)^2} \nabla x \right]$$



Homogeneous  
diffusion



Non-homogeneous  
diffusion

# Variational methods in image processing – 50 frame diffusion



Gaussian:

$$g(x) = c$$

Perona & Malik:

$$g(|\nabla x|) = \frac{1}{1 + (|\nabla x|/\lambda)^2}$$

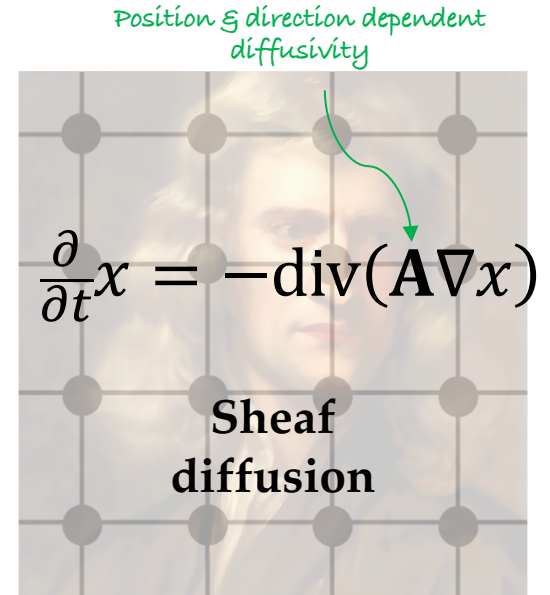
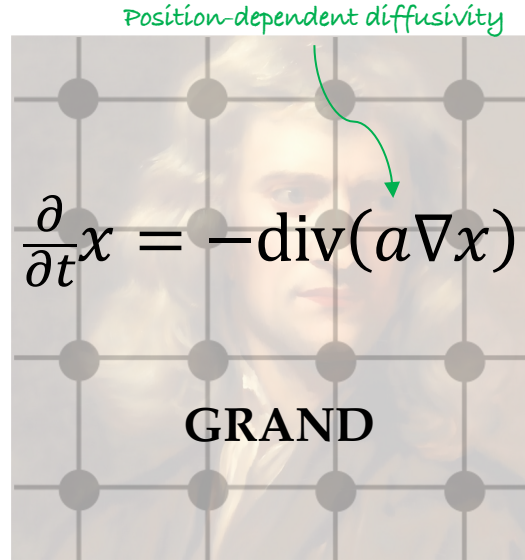
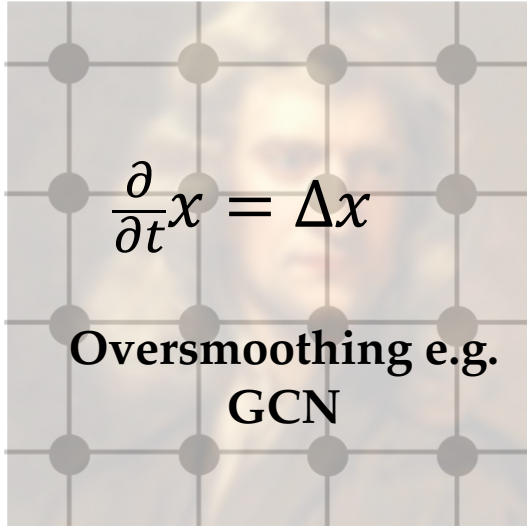
Diffusion Iteration 1



Diffusion Iteration 1



## Attentional diffusivity for GNNs



Kipf 2017, Chamberlain 2021a, Bodner 2022

For more on our work on sheaves see: Aleksa Gordic's  
YouTube channel The AI Epiphany  
<https://www.youtube.com/c/TheAIEpiphany>

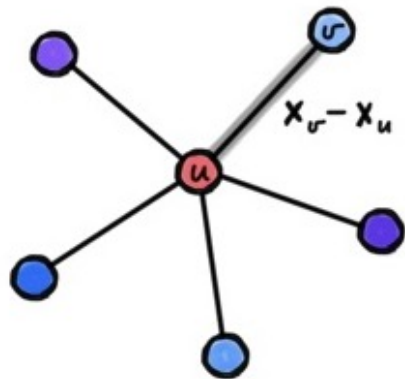


## Spatial discretisation: Continuous to Discrete

$$\frac{\partial x(u,t)}{\partial t} = \text{div}[\nabla x(u,t)] \quad (1)$$

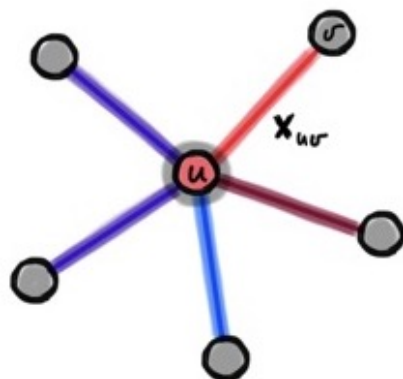
gradient - flow along edges

$$(\nabla X)_{uv} = x_u - x_v$$



divergence - aggregation of edges

$$(\text{div}(X))_u = \sum_v w_{uv} x_{uv}$$



$$\dot{X}(t) = (A(X(t)) - I)X(t) \quad (2)$$





# Temporal discretisation

Neural ODEs:

ResNet:

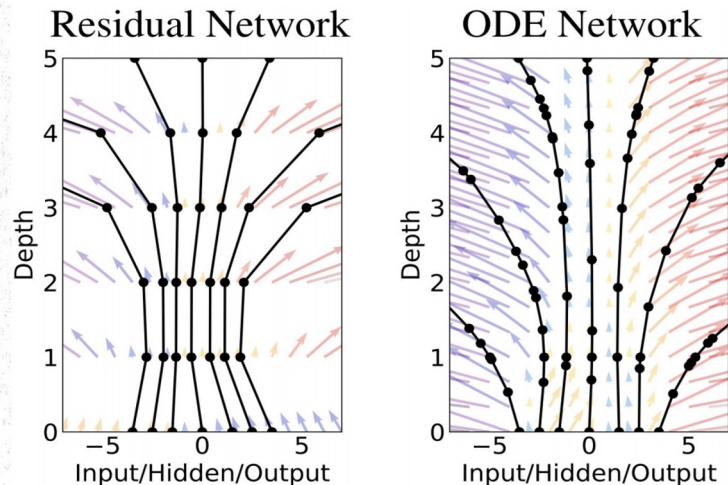
$$\mathbf{h}_{t+1} = \mathbf{h}_t + f(\mathbf{h}_t, \theta_t)$$

Dynamics:

$$\frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}_t, \theta_t, t)$$

ODE Solvers:

$$\mathbf{z}(t_1) = \mathbf{z}(t_0) + \int_{t_0}^{t_1} f(\mathbf{z}(t), t, \theta) dt = \text{ODESolve}(\mathbf{z}(t_0), f, t_0, t_1, \theta)$$





## Simple GNN by discretising time in the graph diffusion equation:

$$\dot{\mathbf{X}}(t) = (\mathbf{A}(\mathbf{X}(t)) - \mathbf{I}) \mathbf{X}(t) \quad (2)$$

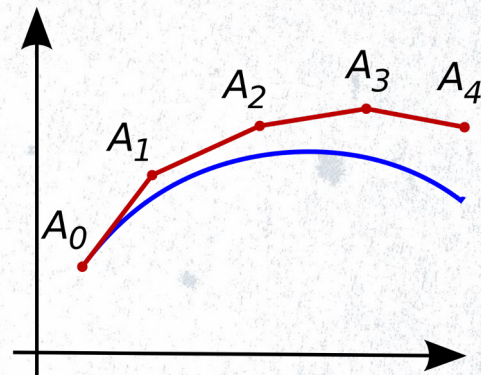
### Explicit Euler temporal discretisation

$$\mathbf{X}(k + 1) = \mathbf{X}(k) + \tau [\mathbf{A}(\mathbf{X}(k)) - \mathbf{I}] \mathbf{X}(k)$$

Set time step  $\tau = 1$  get simplified GCN

$$\mathbf{X}(k + 1) = \mathbf{A}(\mathbf{X}(k)) \mathbf{X}(k)$$

without feature channel mixing and non-linearities



# Better ODE Solvers

Eg: Runge-Kutta 4



## Multi-step methods:

$$x_{n+1} + \sum_{i=1}^s \alpha_i x_{n+1-i} = \tau \sum_{i=0}^s \beta_i f_{n+1-i}$$

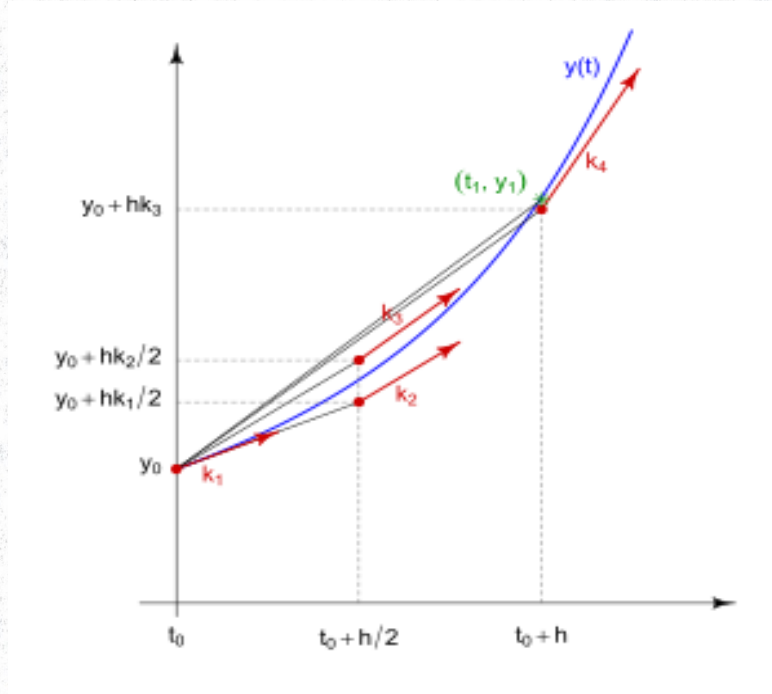
## Adaptive step-size solvers:

order p:  $x_p(\tau)$

order p-1:  $x_{p-1}(\tau)$

error:  $\varepsilon = x_p(\tau) - x_{p-1}(\tau)$

tolerance:  $etol = atol + rtol * \max(|x_0|, |x_1|)$



If  $\varepsilon > etol$

then decrease step-size  $\tau$

## GRAND depth analysis:

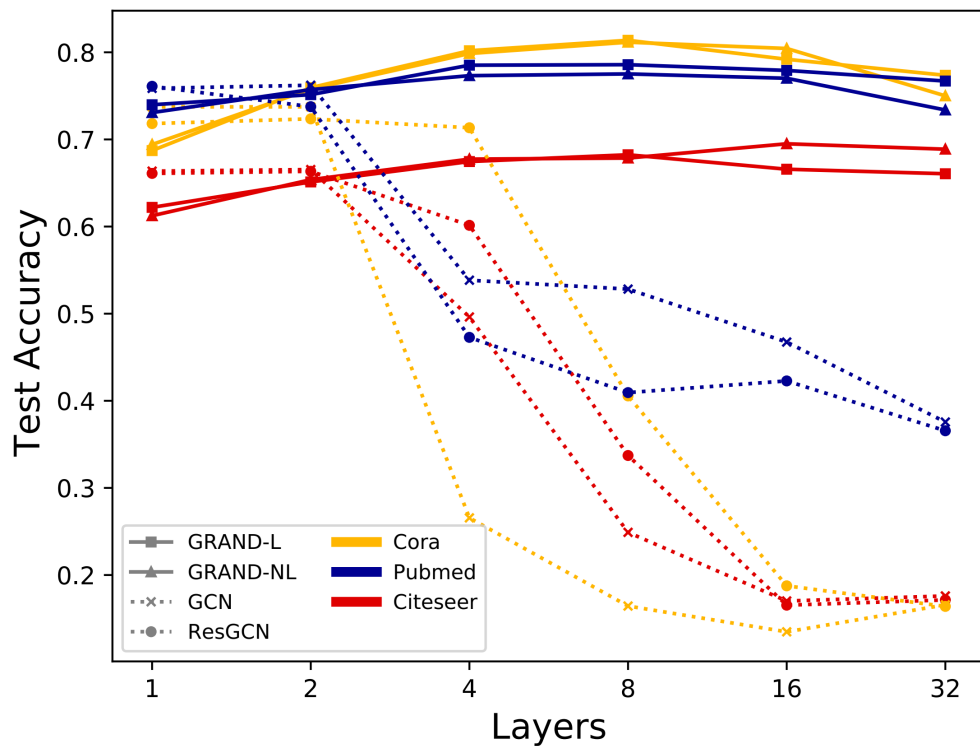


Figure 2. Performance of architectures of different depth.



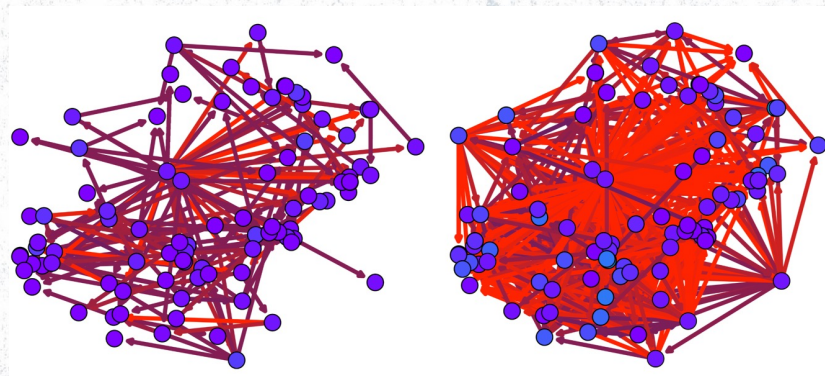
# Decoupling the Input graph and the computational graph

# Spatial discretisation: Graph Rewiring



Decouple **input graph** from **information propagation graph**

- **Neighbourhood sampling (GraphSAGE)**<sup>1</sup>
- **Multi-hop filters (SIGN)**<sup>2</sup>
- **Complete graph**<sup>3</sup>
- **Topology diffusion (DIGL)**<sup>4</sup>



<sup>1</sup>Hamilton et al. 2017; <sup>2</sup>Rossi, Frasca, et B. 2020; <sup>3</sup>Alon, Yahav 2020; <sup>4</sup>Klicpera et al. 2019

## Implicit Versus Explicit Euler's Methods

- **Explicit Euler's Method:**

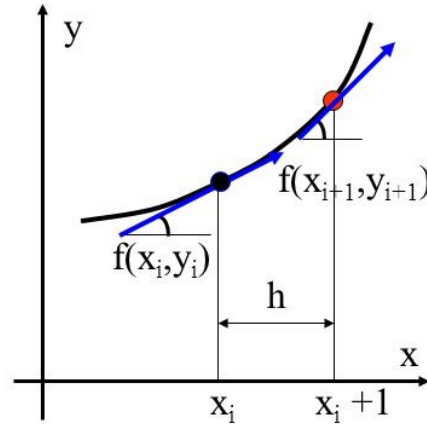
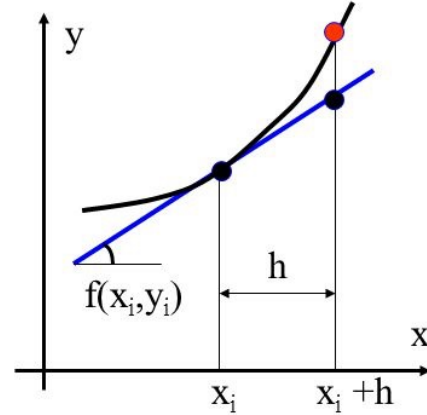
$$\frac{dy}{dx} = f(x, y)$$

$$y_{i+1} = y_i + f(x_i, y_i) \cdot h$$

- **Implicit Euler's Method:**

$$\frac{dy}{dx} = f(x, y)$$

$$y_{i+1} = y_i + f(x_{i+1}, y_{i+1}) \cdot h$$



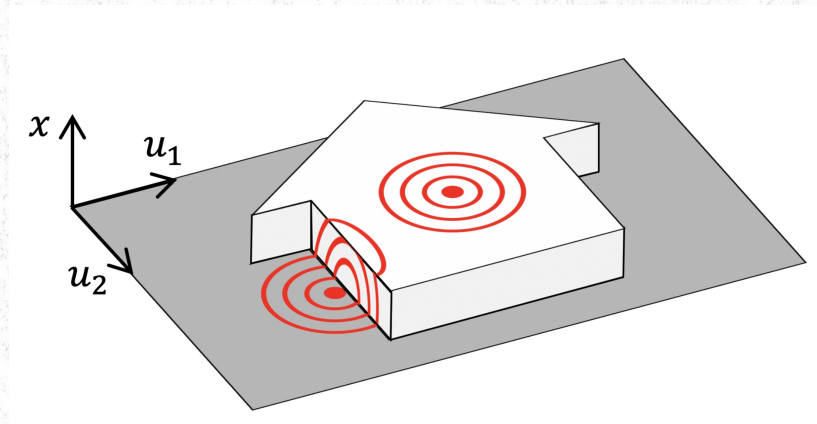


# Images as embedded manifolds



$$\frac{\partial}{\partial t} \mathbf{x} = -\text{div}(a(\mathbf{x}) \nabla \mathbf{x})$$

Non-linear diffusion



$$\frac{\partial}{\partial t} \mathbf{z} = \Delta_{\mathbf{G}} \mathbf{z}$$

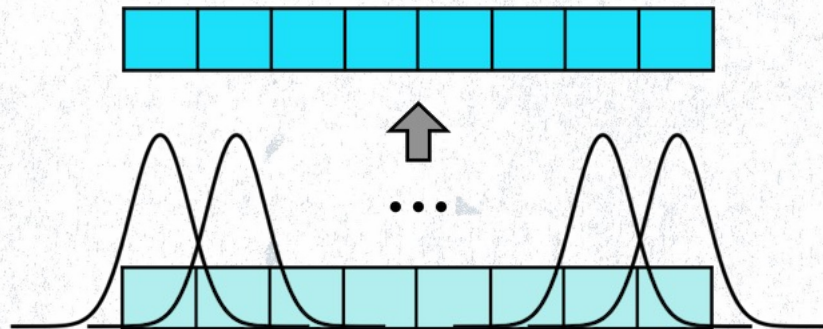
Non-Euclidean diffusion



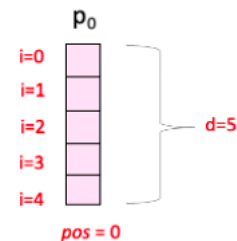


# Intro to positional encodings

- Popularized in the transformer
- Now widely used in GNNs
- Many options
  - Random node features<sup>1</sup>
  - Graph Laplacian eigenvectors<sup>2</sup>
  - Graph substructure counts<sup>3</sup>
  - Bags of subgraphs<sup>4</sup>



$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$





# *What to do with the positional encodings*

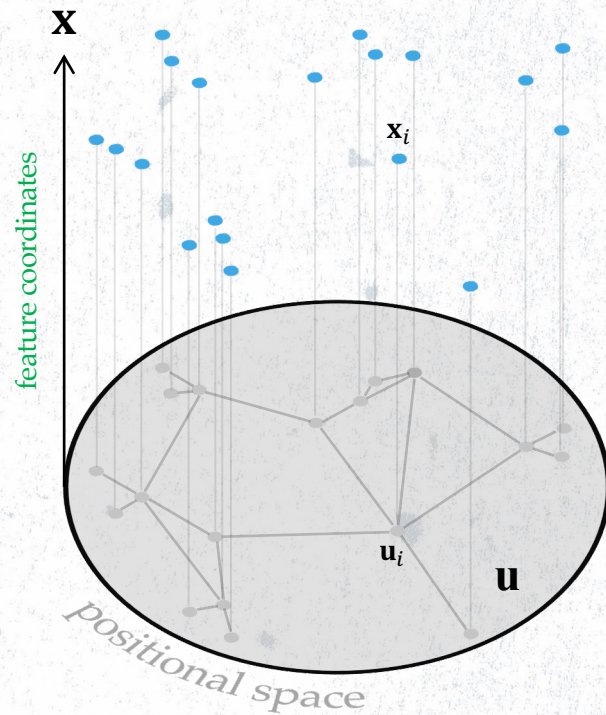
- In computer vision they are just for edge detection and then thrown away
- The graph approach is more elegant as the evolved positional encodings can be used to rewire the graph?
- Why rewire the graph?
  - Decouples the given and computational graph can remove bottlenecks or increase performance and scalability



# Graph Beltrami flow

- Graph with positional and feature node coordinates  $\mathbf{z}_i = (\mathbf{u}_i, \mathbf{x}_i)$
- **Graph Beltrami flow**

$$\frac{\partial}{\partial t} \mathbf{z}_i = \sum_{j:(i,j) \in E} a(\mathbf{z}_i, \mathbf{z}_j) (\mathbf{z}_j - \mathbf{z}_i)$$



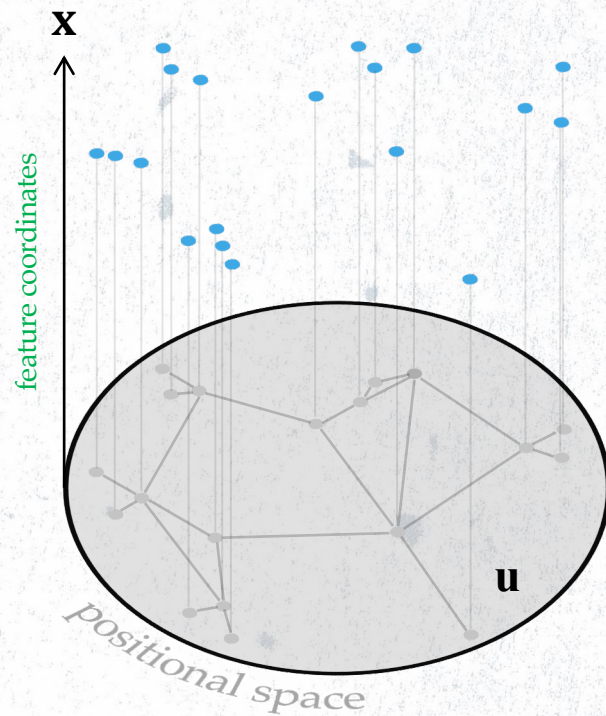


# Graph Beltrami flow

- Graph with positional and feature node coordinates  $\mathbf{z}_i = (\mathbf{u}_i, \mathbf{x}_i)$
- **Graph Beltrami flow**

$$\frac{\partial}{\partial t} \mathbf{z}_i = \sum_{j:(i,j) \in E} a(\mathbf{z}_i, \mathbf{z}_j) (\mathbf{z}_j - \mathbf{z}_i)$$

- Evolution of  $\mathbf{x}$  = feature diffusion



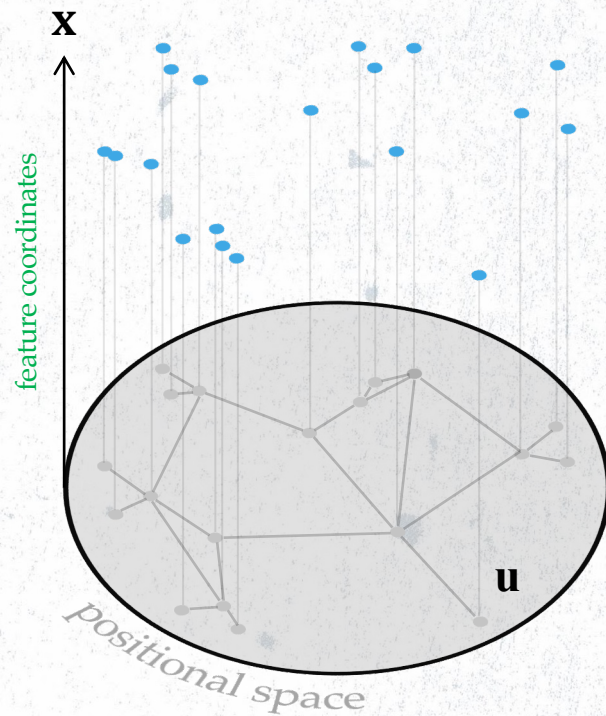


# Graph Beltrami flow

- Graph with positional and feature node coordinates  $\mathbf{z}_i = (\mathbf{u}_i, \mathbf{x}_i)$
- **Graph Beltrami flow**

$$\frac{\partial}{\partial t} \mathbf{z}_i = \sum_{j:(i,j) \in E} a(\mathbf{z}_i, \mathbf{z}_j) (\mathbf{z}_j - \mathbf{z}_i)$$

- Evolution of  $\mathbf{x}$  = feature diffusion



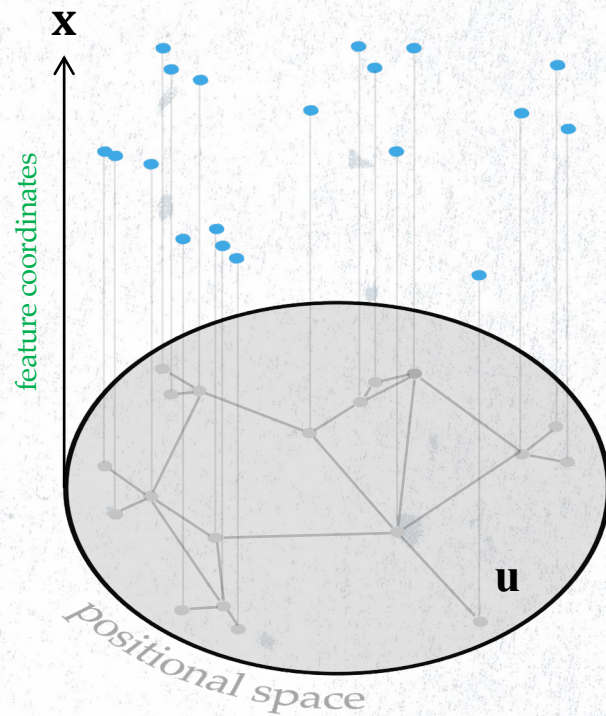


# Graph Beltrami flow

- Graph with positional and feature node coordinates  $\mathbf{z}_i = (\mathbf{u}_i, \mathbf{x}_i)$
- **Graph Beltrami flow**

$$\frac{\partial}{\partial t} \mathbf{z}_i = \sum_{j:(i,j) \in E} a(\mathbf{z}_i, \mathbf{z}_j) (\mathbf{z}_j - \mathbf{z}_i)$$

- Evolution of  $\mathbf{x}$  = feature diffusion
- Evolution of  $\mathbf{z}$  = graph rewiring





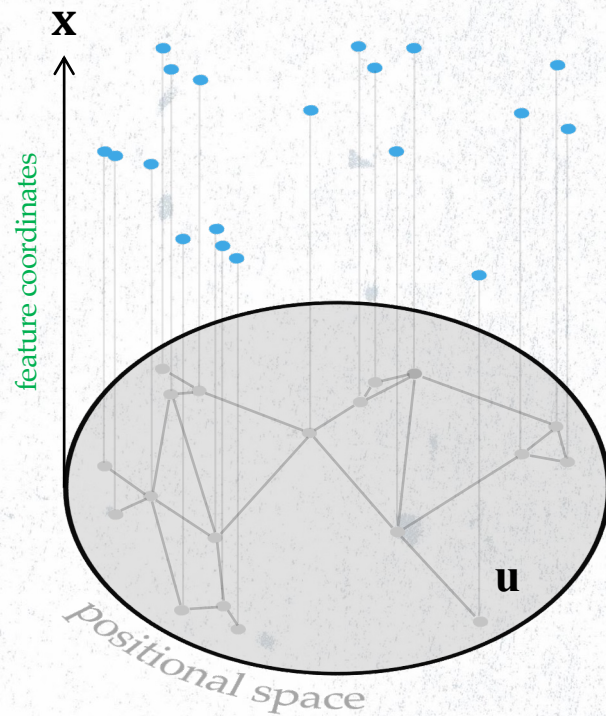
# Graph Beltrami flow

- Graph with positional and feature node coordinates  $\mathbf{z}_i = (\mathbf{u}_i, \mathbf{x}_i)$
- **Graph Beltrami flow**

$$\frac{\partial}{\partial t} \mathbf{z}_i = \sum_{j:(i,j) \in E'} a(\mathbf{z}_i, \mathbf{z}_j) (\mathbf{z}_j - \mathbf{z}_i)$$

*rewired graph*

- Evolution of  $\mathbf{x}$  = feature diffusion
- Evolution of  $\mathbf{z}$  = graph rewiring





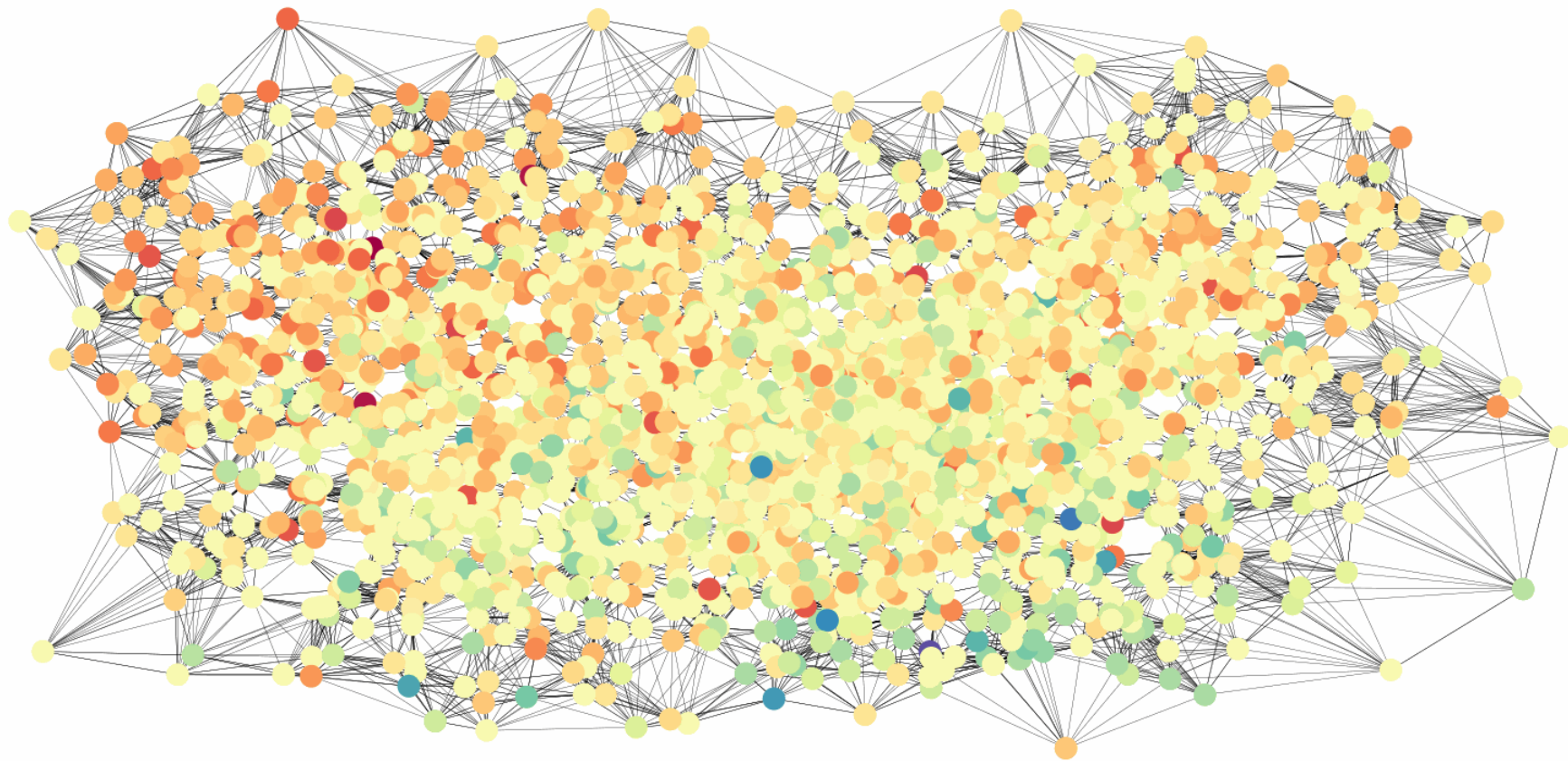
# GNN Architectures as instances of BLEND

$$\mathbf{z}^{(k+1)} = \mathbf{Q}(\mathbf{z}^{(k)})\mathbf{z}^{(k)}$$

<b>Method</b>	<b>Evolution</b>	<b>Diffusivity</b>	<b>Graph</b>	<b><u>Discretisation</u></b>
<u>ChebNet</u>	<b>X</b>	Fixed <b>A</b>	Fixed $E$	Explicit Fixed step
GAT	<b>X</b>	<b>A(X)</b>	Fixed $E$	Explicit Fixed step
<u>MoNet</u>	<b>X</b>	<b>A(U)</b>	Fixed $E$	Explicit Fixed step
Transformer	<b>X</b>	<b>A(U, X)</b>	Fixed $E = V^2$	Explicit Fixed step
DIGL	<b>X</b>	<b>A(X)</b>	Fixed $E(\mathbf{U})$	Explicit Fixed step
DGCNN	<b>X</b>	<b>A(X)</b>	Adaptive $E(\mathbf{X})$	Explicit Fixed step
<b>BLEND</b>	<b>U, X</b>	<b>A(U, X)</b>	<b>Adaptive <math>E(\mathbf{U})</math></b>	<b>Explicit Adaptive step / Implicit</b>



Beltrami Flow, diffusion time=0





# GNNs as Interacting Particle Systems

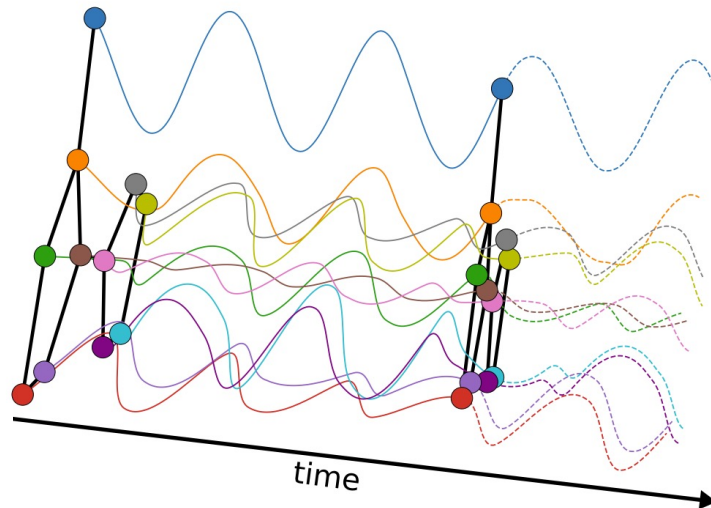
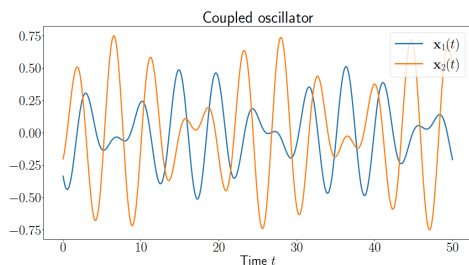
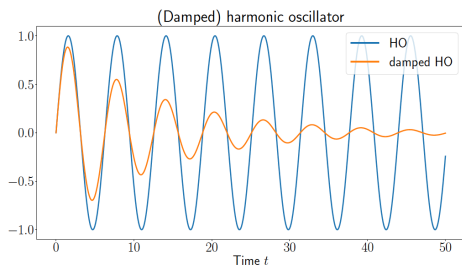
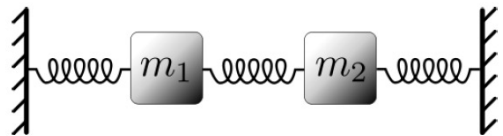
## Interacting Particle Approach



- **Associate positions with node embeddings**
- **Learning generates node trajectories**
- **Trajectories are constrained by an energy and described by a differential equation (the gradient flow of the energy)**

# Graph Coupled Oscillator Networks (GraphCON)

Rusch et al at ICML 22



$$\mathbf{X}'' = \underbrace{\sigma(\mathbf{F}_\theta(\mathbf{X}, t) + \mathbf{X}\bar{\mathbf{W}} + \bar{\mathbf{b}})}_{\text{Graph coupling}} - \underbrace{\gamma\mathbf{X} - \alpha\mathbf{X}'}_{\text{Damped Oscillator}}$$

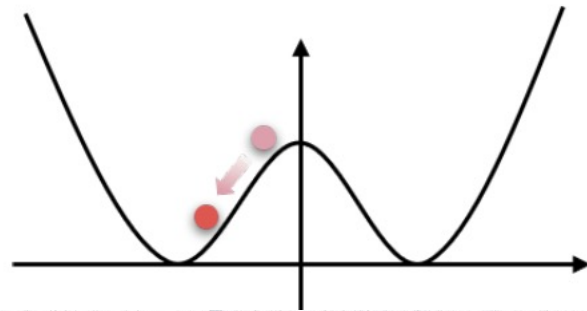
Graph coupling

Damped Oscillator

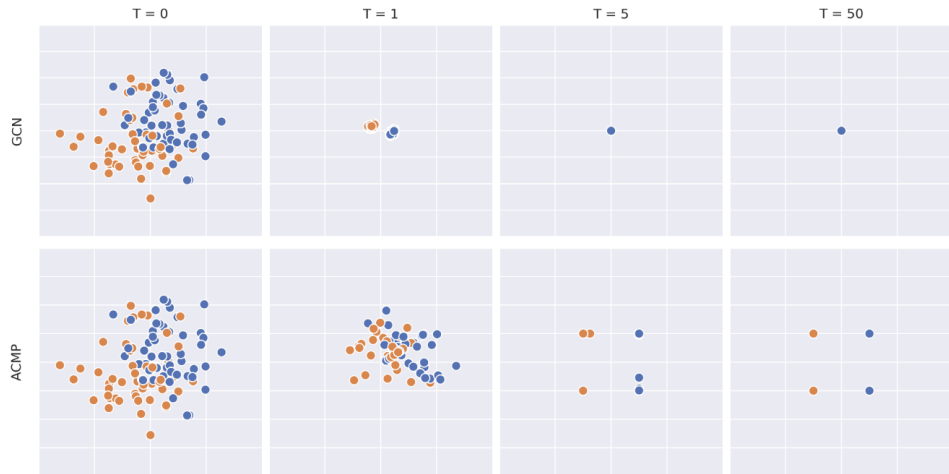
# Allen-Cahn Message Passing (ACMP)

Wang et al. 22

The Allen-Cahn potential introduces repulsive forces that prevent oversmoothing



$$\frac{\partial \mathbf{x}}{\partial t} = -\nabla_{\mathbf{x}} \Phi, \quad \frac{\partial \mathbf{x}_i}{\partial t} = -\frac{\partial \Phi}{\partial \mathbf{x}_i} = \alpha \sum_{j \in \mathcal{N}_i} a_{i,j} (\mathbf{x}_j - \mathbf{x}_i) + \delta \mathbf{x}_i (1 - \|\mathbf{x}_i\|^2)$$



Graph Diffusion

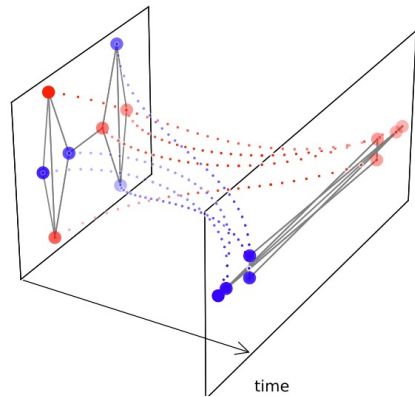
Allen-Cahn double well

# Gradient Flow Framework (GRAFF)

Di Giovanni, Rowbottom et al. 22



The gradient flow of a learnable energy can induce repulsive forces that prevent oversmoothing



$$\mathcal{E}^{\text{tot}}(\mathbf{F}) := \frac{1}{2} \sum_i \langle \mathbf{f}_i, \mathbf{\Omega} \mathbf{f}_i \rangle - \frac{1}{2} \sum_{i,j} \bar{a}_{ij} \langle \mathbf{f}_i, \mathbf{W} \mathbf{f}_j \rangle \equiv \mathcal{E}_{\mathbf{\Omega}}^{\text{ext}}(\mathbf{F}) + \mathcal{E}_{\mathbf{W}}^{\text{pair}}(\mathbf{F})$$

Gradient flow gives diffusion with channel mixing and attractive / repulsive forces

$$\dot{\mathbf{F}}(t) = -\nabla_{\mathbf{F}} \mathcal{E}^{\text{tot}}(\mathbf{F}(t)) = \underbrace{-\mathbf{F}(t)\mathbf{\Omega}}_{\text{Damping}} + \underbrace{\bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W}}_{\text{Onsager diffusion with symmetric W}}$$

Damping

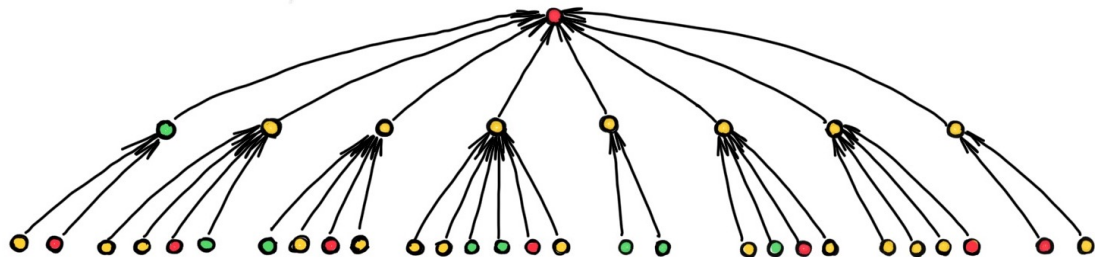
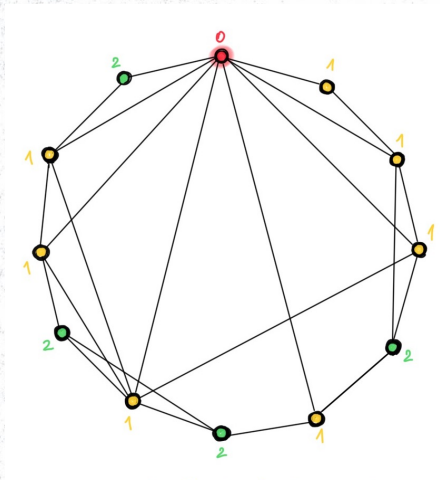
Onsager diffusion with symmetric W



# Bottlenecks and Oversquashing



# Over-squashing & Bottlenecks



In small-world graphs metric ball volume  $\text{vol}(B_k) = \sum_{j \in B_k} d_j$   
grows exponentially with ball radius  $k$

Long-distance dependency + Fast volume growth  
= Over-squashing





# Characterisation of Over-squashing in GNNs

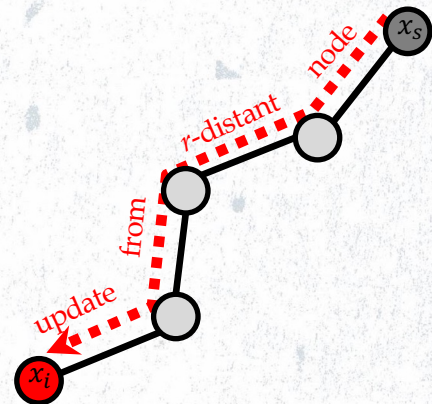
- Multilayer MPNN-type GNN of the form

$$x_i^{(\ell+1)} = \phi_\ell \left( x_i^{(\ell)}, \sum_{j=1}^n \hat{a}_{ij} \psi_\ell \left( x_i^{(\ell)}, x_j^{(\ell)} \right) \right)$$

- $|\nabla \phi_\ell| \leq \alpha$  and  $|\nabla \psi_\ell| \leq \beta$  for  $\ell = 0, 1, \dots, L$ .

**Lemma 1 (sensitivity):** Let node  $s$  be geodesically  $d_G(i, s) = r + 1$  away from node  $i$ . Then

$$\left| \frac{\partial h_i^{(r+1)}}{\partial x_s} \right| \leq (\alpha\beta)^{r+1} (\hat{\mathbf{A}}^{r+1})_{is}$$



Over-squashing: small Jacobian  $\left| \frac{\partial x_i^{(r+1)}}{\partial x_s} \right|$  leads to poor information propagation



# Characterisation of Over-squashing in GNNs

- Multilayer MPNN-type GNN of the form

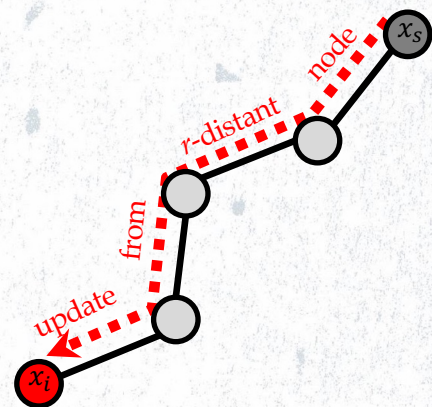
$$x_i^{(\ell+1)} = \phi_\ell \left( x_i^{(\ell)}, \sum_{j=1}^n \hat{a}_{ij} \psi_\ell \left( x_i^{(\ell)}, x_j^{(\ell)} \right) \right)$$

- $|\nabla \phi_\ell| \leq \alpha$  and  $|\nabla \psi_\ell| \leq \beta$  for  $\ell = 0, 1, \dots, L$ .

**Lemma 1 (sensitivity):** Let node  $s$  be geodesically  $d_G(i, s) = r + 1$  away from node  $i$ . Then

$$\left| \frac{\partial x_i^{(r+1)}}{\partial x_s} \right| \leq (\alpha\beta)^{r+1} (\hat{\mathbf{A}}^{r+1})_{is}$$

It's the graph structure ("bottleneck") to blame!



Over-squashing: small Jacobian  $\left| \frac{\partial x_i^{(r+1)}}{\partial x_s} \right|$  leads to poor information propagation



# Characterisation of Over-squashing in GNNs

- Multilayer MPNN-type GNN of the form

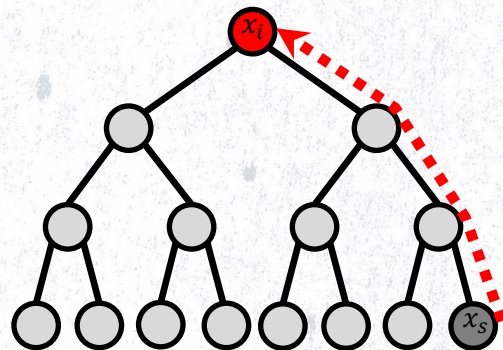
$$x_i^{(\ell+1)} = \phi_\ell \left( x_i^{(\ell)}, \sum_{j=1}^n \hat{a}_{ij} \psi_\ell \left( x_i^{(\ell)}, x_j^{(\ell)} \right) \right)$$

- $|\nabla \phi_\ell| \leq \alpha$  and  $|\nabla \psi_\ell| \leq \beta$  for  $\ell = 0, 1, \dots, L$ .

**Lemma 1 (sensitivity):** Let node  $s$  be geodesically  $d_G(i, s) = r + 1$  away from node  $i$ . Then

$$\left| \frac{\partial x_i^{(r+1)}}{\partial x_s} \right| \leq (\alpha\beta)^{r+1} (\hat{\mathbf{A}}^{r+1})_{is}$$

It's the graph structure ("bottleneck") to blame!

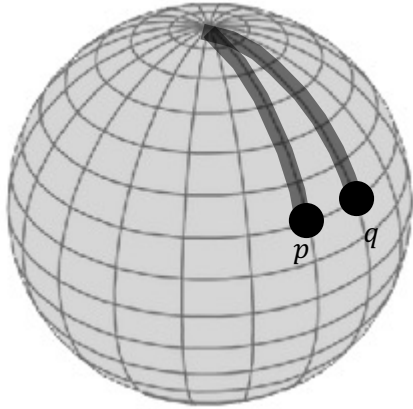


Pathological example: binary tree  
 $(\hat{\mathbf{A}}^{r+1})_{is} = \frac{1}{2} \cdot 3^{-r}$

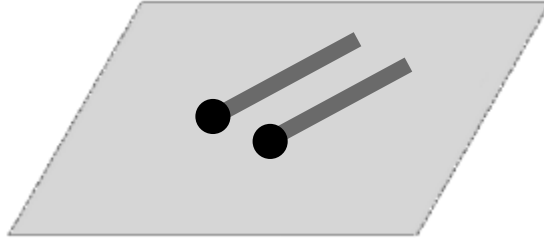


# Graph Curvature

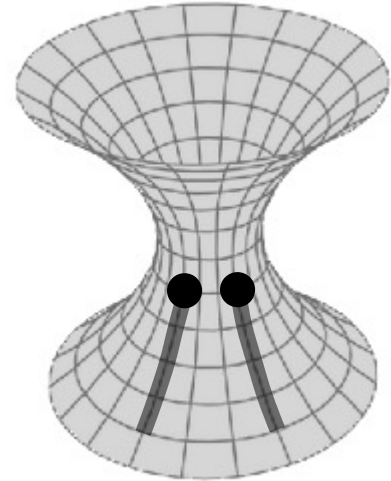
# Ricci Curvature on Manifolds



Spherical ( $>0$ )



Euclidean ( $=0$ )



Hyperbolic ( $<0$ )

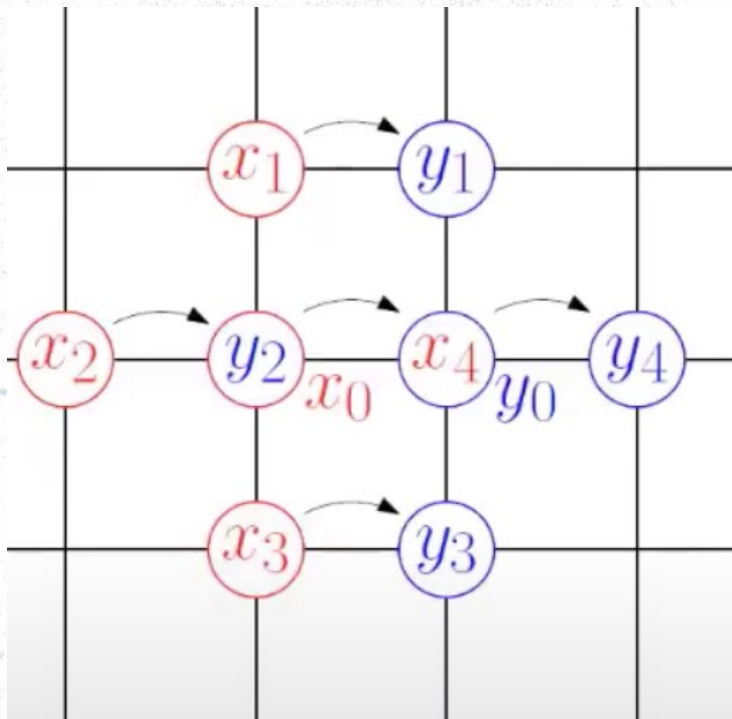
**“geodesic dispersion”**



# Ricci Curvature on Graphs

Sectional curvature defined on edges. Ollivier curvature:

$$\kappa(x, y) := 1 - \frac{W_1(m_x, m_y)}{d(x, y)},$$

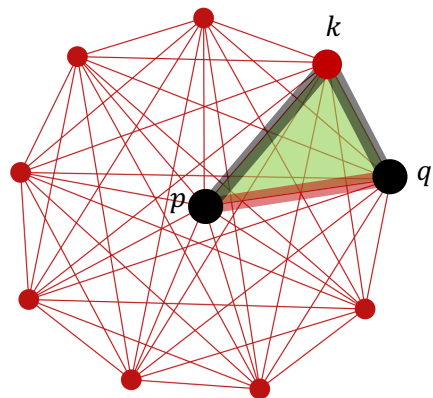




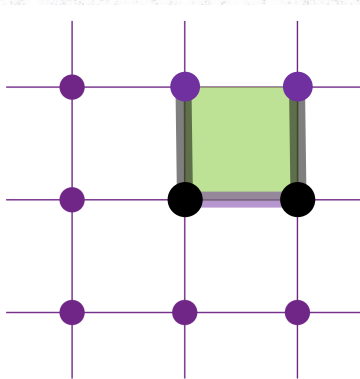
# Ricci Curvature on Graphs

Sectional curvature defined on edges. Ollivier curvature:

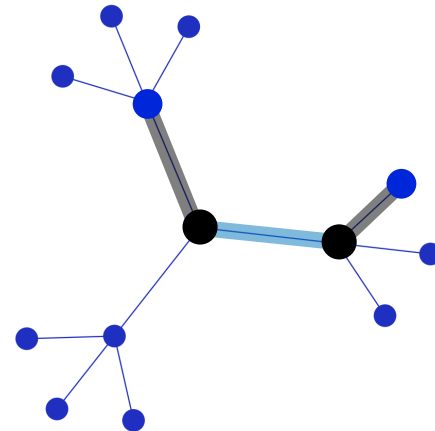
$$\kappa(x, y) := 1 - \frac{W_1(m_x, m_y)}{d(x, y)},$$



Clique ( $>0$ )



Grid ( $=0$ )



Tree ( $<0$ )



# Balanced Forman Curvature

Ollivier curvature is expensive to calculate due to the optimal transport map

**Balanced Forman Curvature** of edge  $i \sim j$  in simple unweighted graph  $\text{Ric}(i, j) = 0$  if  $\min\{d_i, d_j\} = 1$  and otherwise

$$\text{Ric}(i, j) = \frac{2}{d_i} + \frac{2}{d_j} + 2 \frac{|\#_{\Delta}(i, j)|}{\max\{d_i, d_j\}} + \frac{|\#_{\Delta}(i, j)|}{\min\{d_i, d_j\}} + \frac{\gamma_{\max}^{-1}}{\max\{d_i, d_j\}} (|\#_{\square}^i(i, j)| + |\#_{\square}^j(i, j)|) - 2$$

Degree of  $i$  (points to  $d_i$ )  
 Triangles based at  $i \sim j$  (points to  $|\#_{\Delta}(i, j)|$ )  
 Max number of 4-cycle based at  $i \sim j$  traversing the same node (points to  $\gamma_{\max}^{-1}$ )  
 Neighbours of  $i$  forming a 4-cycle based at  $i \sim j$  (w/o diagonals) (points to  $|\#_{\square}^i(i, j)|$ )



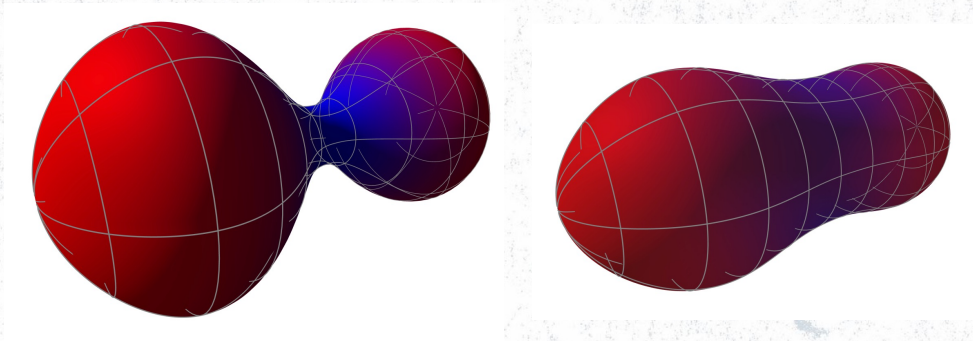


# Ricci Flow

# Ricci flow

- **Ricci flow:** “diffusion of the Riemannian metric”

$$\frac{\partial g_{ij}}{\partial t} = R_{ij}$$



Evolution of a manifold under Ricci flow

Ricci 1903; Hamilton 1988; Perelman 2003

# Science

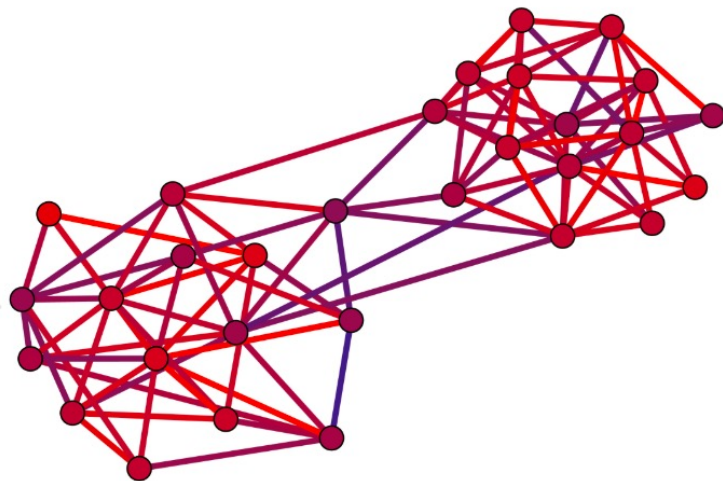
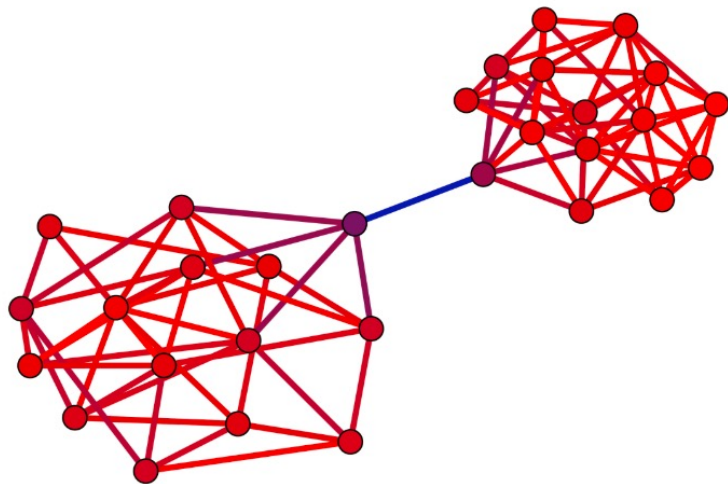
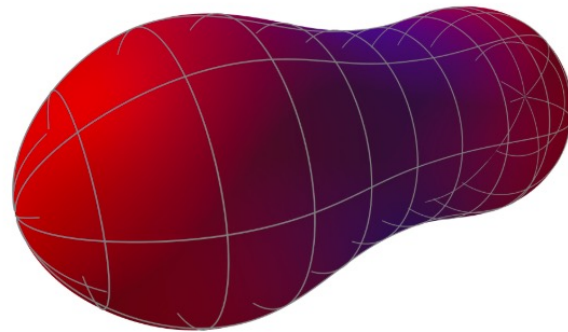
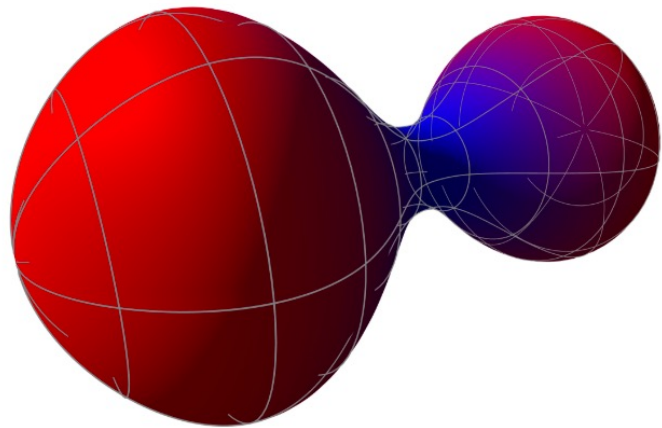
22 December 2006 | \$10

Breakthrough  
of the Year



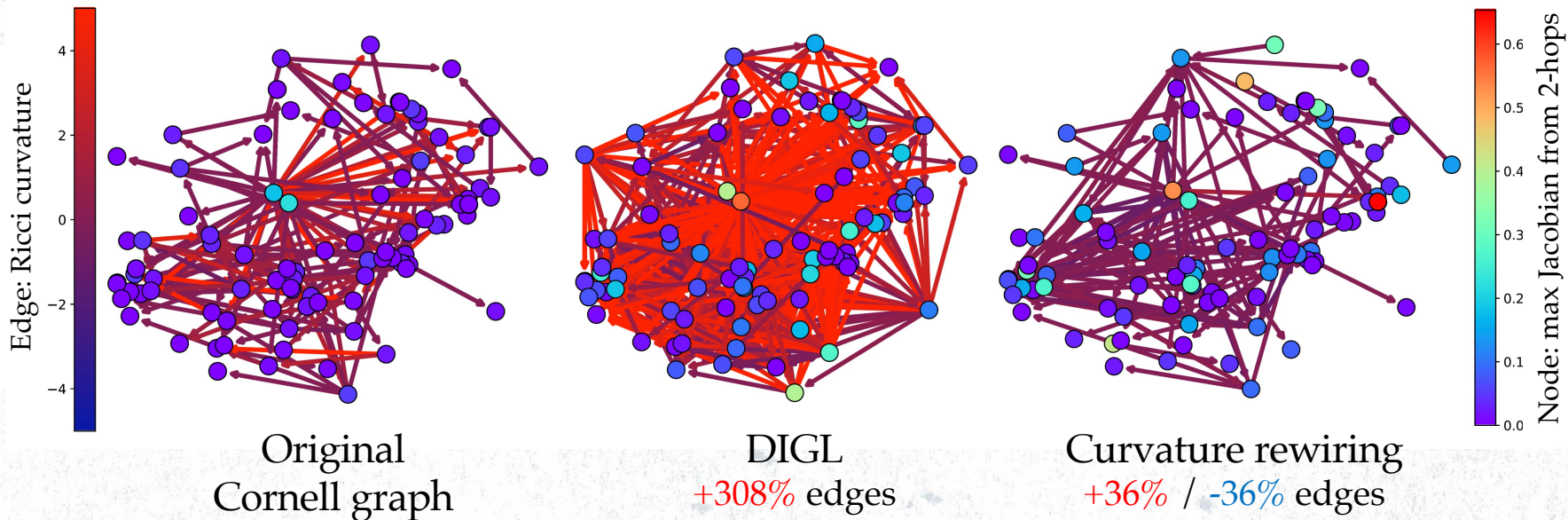
The  
Poincaré  
Conjecture  
PROVED

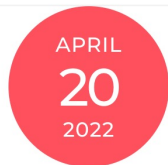
AAAS





# Curvature- vs Diffusion-based Rewiring





# Announcing the ICLR 2022 Outstanding Paper Award Recipients

YEJIN CHOI / 2022 Conference / awards ICLR 2022

*By ICLR 2022 Senior Program Chair Yan Liu and Program Chairs Chelsea Finn, Yejin Choi,  
Marc Deisenroth*

## Outstanding Paper Honorable Mentions

[Understanding over-squashing and bottlenecks on graphs via curvature](#)

*By Jake Topping, Francesco Di Giovanni, Benjamin Paul Chamberlain, Xiaowen Dong, Michael M.  
Bronstein*

For an intro to the paper see: Aleksa Gordic's YouTube channel The AI Epiphany  
<https://www.youtube.com/c/TheAIEpiphany>



# Summary

- **GNNs as differential equations – layers and continuous time**
  - Stability conditions
  - Architectures based on numerical solvers
  - Control smoothing
- **Evolving the underlying space and rewiring**
- **Graph Curvature**
  - Decoupling the input and computational graph
  - Remove bottlenecks

# Thank You and Questions



@b\_p\_chamberlain

