"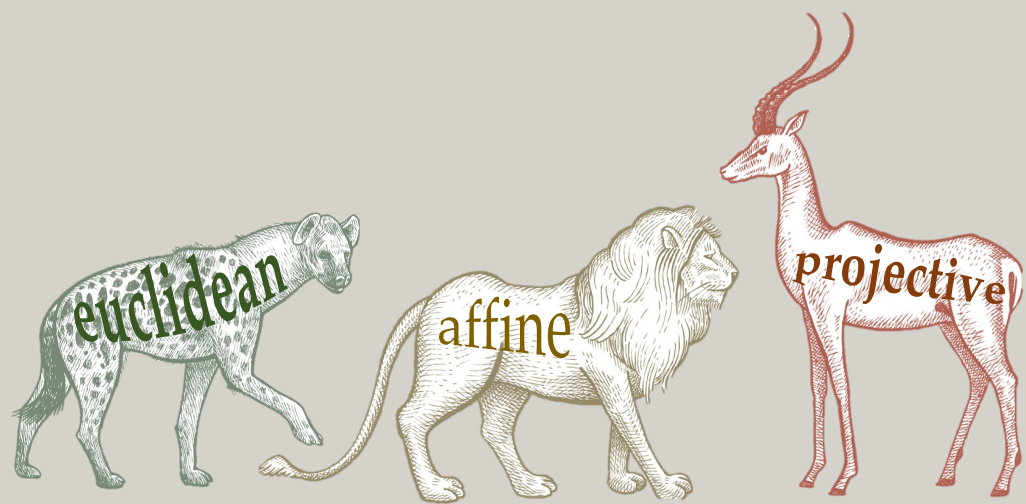Symmetry, as wide or as narrow as you may define its meaning, is one idea by which man through the ages has tried to comprehend and create order, beauty, and perfection"
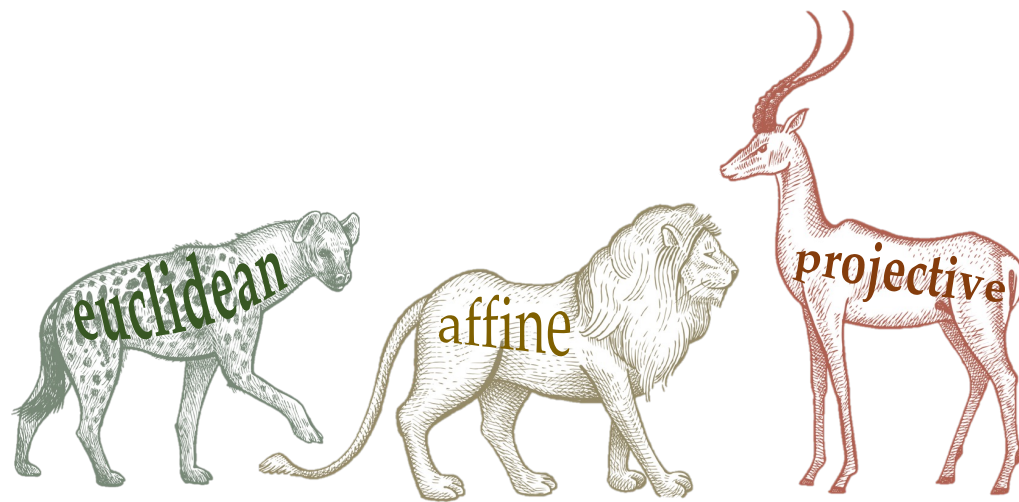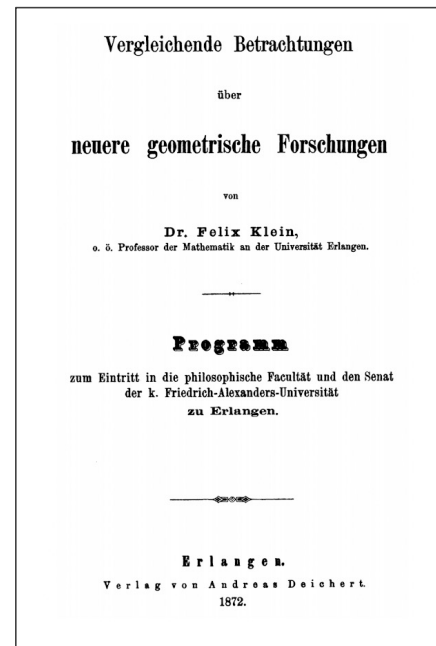
**Hermann Weyl**

XIX century

euclidean

affine

projective

# *The Erlangen Programme*



**Geometry = space + transformation group**

Klein 1872

Felix Klein

# Cultural Impact in Mathematics
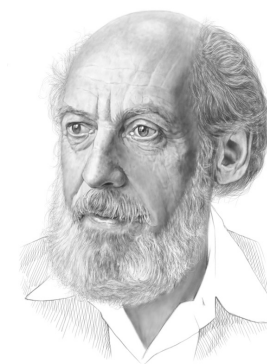
**E. Beltrami**

**E. Cartan**

1868

1920s

**GENERAL THEORY OF NATURAL EQUIVALENCES**

BY

SAMUEL EILENBERG AND SAUNDERS MacLANE

CONTENTS

**Introduction.** The subject matter of this paper is best explained by an example, such as that of the relation between a vector space $L$ and its "dual"

Presented to the Society, September 8, 1942; received by the editors May 15, 1945.

**Category Theory**

**S. Eilenberg**

**S. Mac Lane**

1945

Portraits: Ihor Gorskyi

# *New Physics*



| H. Poincaré | H. Minkowski | E. Noether | H. Weyl | C. N. Yang | R. L. Mills |
|:-----------:|:------------:|:----------:|:-------:|:----------:|:-----------:|
| 1904 | 1907 | 1918 | 1929 | 1954 | |

$R^{1,3}$

$O(1,3)$

**External symmetry**

$U(1)$

$SU(2)$

$SU(3)$

**Internal symmetry**

"It is only slightly overstating the case to say that Physics is the study of symmetry"

— "More is different"

**P. Anderson**

# Geometric Deep Learning Blueprint

**domain** $\Omega$  |  **signals** $\mathcal{X}(\Omega)$  |  **functions** $\mathcal{F}\big(\mathcal{X}(\Omega)\big)$



**symmetry group** $\mathfrak{G}$  |  **group representation** $\rho(\mathfrak{G})$  |  **Invariance / Equivariance**

$$f(\rho(g)x) = f(x)$$
$$f(\rho(g)x) = \rho(g)f(x)$$

B et al. 2021

# *Example: Convolutional Neural Networks*

**Plane** $\mathbb{R}^2$  | **images** $\mathcal{X}(\mathbb{R}^2)$ | **functions** $\mathcal{F}(\mathcal{X}(\Omega))$



**Translation group** $T(2)$

**Shift operator** $S$

$$S_v x(u) = x(u - v)$$

**Convolutional layer**

$$(Sx \star y) = S(x \star y)$$

# *Example: Graph Neural Networks*



**Graph** $G = (V, E)$

**Node features** $\mathcal{X}(G)$

**functions** $\mathcal{F}\big(\mathcal{X}(\Omega)\big)$

**F**

**Permutation group** $\Sigma_n$

**Permutation matrix P**

$$\mathbf{PX} = \big(x_{\pi^{-1}(i),j}\big)$$

**Message passing**

$$\mathbf{F}\big(\mathbf{PX}, \mathbf{PAP}^\top\big) = \mathbf{PF}(\mathbf{X}, \mathbf{A})$$

# *Graphs = Systems of Relations and Interactions*



**Molecules**          **Interactomes**          **Social networks**

Graph Neural Networks

*Graphs: The Basics*



**node feature** $\mathbf{x}_i$

**node** $i$

**edge** $i \sim j$

graph

# Key Structural Properties of Graphs



**node feature** $\mathbf{x}_i$

**node** $i$

**edge** $i \sim j$

arbitrary ordering of nodes

graph

# Key Structural Properties of Graphs



**Feature matrix** $n \times d$

**X**

arbitrary ordering of nodes

# Key Structural Properties of Graphs

**Adjacency matrix $n \times n$**

**Feature matrix $n \times d$**

**A**

**X**

arbitrary ordering of nodes

# Key Structural Properties of Graphs



**Adjacency matrix $n \times n$**

$\mathbf{PAP}^{\mathrm{T}}$

**Feature matrix $n \times d$**

$\mathbf{PX}$

arbitrary ordering of nodes

# Key Structural Properties of Graphs



**Adjacency matrix $n \times n$**

**Feature matrix $n \times d$**

$\mathbf{PAP}^{\mathrm{T}}$

$\mathbf{PX}$

arbitrary ordering of nodes

$n!$ permutations

*Invariant vs Equivariant tasks*

water solubility?

who is a spammer?

*Invariant Graph Functions*

graph function $f(\mathbf{X}, \mathbf{A})$

*Invariant Graph Functions*

graph function $f(\mathbf{X}, \mathbf{A})$

*Invariant Graph Functions*

permutation-invariant

$$f\left(\mathbf{P}\mathbf{X}, \mathbf{P}\mathbf{A}\mathbf{P}^{\top}\right) = f(\mathbf{X}, \mathbf{A})$$

*Invariant vs Equivariant tasks*



water solubility?

who is a spammer?

*Equivariant Graph Functions*

node function $\mathbf{F}(\mathbf{X}, \mathbf{A})$

*Equivariant Graph Functions*

node function $\mathbf{F}(\mathbf{X}, \mathbf{A})$

# Equivariant Graph Functions

permutation-equivariant

$$\mathbf{F}\left(\mathbf{PX}, \mathbf{PAP}^{\top}\right) = \mathbf{PF}(\mathbf{X}, \mathbf{A})$$

# Graph Neural Networks: Node tasks



Permutation-equivariant

Permutation-equivariant

Permutation-equivariant

Node-wise predictions

# Graph Neural Networks: Graph tasks

# Neighbour Aggregation



neighbourhood
$$\mathcal{N}_i = \{j : i \sim j\}$$

*Neighbour Aggregation*

multiset of
neighbour features

$$\mathbf{X}_{\mathcal{N}_i} = \left\{\!\!\left\{ \mathbf{x}_{j \in \mathcal{N}_i} \right\}\!\!\right\}$$

neighbourhood
$\mathcal{N}_i = \{j : i \sim j\}$

*Neighbour Aggregation*

local function

$$\phi \left( \begin{array}{c} \mathbf{x}_i \\ \mathbf{X}_{\mathcal{N}_i} \end{array} \right)$$

*Neighbour Aggregation*



local function

$$\phi \left( \begin{array}{c} \mathbf{x}_i \\ \mathbf{X}_{\mathcal{N}_i} \end{array} \right)$$

permutation invariant

*GNN Layer*

$$\mathbf{F}(\mathbf{X}, \mathbf{A}) = \begin{pmatrix} - \phi(\mathbf{x}_1, \mathbf{X}_{\mathcal{N}_1}) - \\ \vdots \\ - \phi(\mathbf{x}_i, \mathbf{X}_{\mathcal{N}_i}) - \\ \vdots \\ - \phi(\mathbf{x}_n, \mathbf{X}_{\mathcal{N}_n}) - \end{pmatrix}$$

permutation equivariant

# Convolutional GNNs



$$\mathbf{x}_i \leftarrow \sigma\left(\sum_{j \in \mathcal{N}_i \cup \{i\}} a_{ij}\psi(\mathbf{x}_j)\right)$$

Defferard et al. 2016; Kipf, Welling 2016 (GCN)

# Convolutional GNNs



$$\mathbf{x}_i \leftarrow \sigma\left( \sum_{j \in \mathcal{N}_i \cup \{i\}} a_{ij} \psi(\mathbf{x}_j) \right)$$

graph adjacency

Defferard et al. 2016; Kipf, Welling 2016 (GCN)

# Convolutional GNNs



$$\mathbf{x}_i \leftarrow \sigma\left(\sum_{j\in\mathcal{N}_i\cup\{i\}} a_{ij}\psi(\mathbf{x}_j)\right)$$

nonlinear activation

graph adjacency

node-wise transformation

Defferard et al. 2016; Kipf, Welling 2016 (GCN)

# Convolutional GNNs



$$\mathbf{x}_i \leftarrow \sigma\left( \sum_{j \in \mathcal{N}_i \cup \{i\}} a_{ij} \mathbf{W} \mathbf{x}_j \right)$$

nonlinear
activation

graph
adjacency

node-wise linear
transformation

Defferard et al. 2016; Kipf, Welling 2016 (GCN)

# *Convolutional GNNs*

- Simplest GNN

- Highly scalable

- Industrial use cases

- Folklore: works only on ==homophilic graphs==



$$\mathbf{X} \leftarrow \sigma(\mathbf{A}\mathbf{X}\mathbf{W})$$

diffusion $n \times n$

channel mixing $d \times d$

Defferard et al. 2016; Kipf, Welling 2016 (GCN)
Rossi, Frasca et B 2020 (SIGN); Ying et al. 2018 (PinSAGE)

# Attentional GNNs



$$\mathbf{x}_i \leftarrow \sigma \left( \sum_{j \in \mathcal{N}_i \cup \{i\}} \alpha_{ij}(\mathbf{x}_i, \mathbf{x}_j) \psi(\mathbf{x}_j) \right)$$

learnable attention weights

Monti et al. 2017; Veličković et al. 2018 (GAT)

# Message-Passing GNNs



$$\mathbf{x}_i \leftarrow \sigma\left(\sum_{j \in \mathcal{N}_i \cup \{i\}} \psi(\mathbf{x}_i, \mathbf{x}_j)\right)$$

message from
node $j$ to node $i$

Gilmer et al. 2017 (MPNN); Battaglia et al 2018 (Graph Networks)
Wang et B, Solomon 2018 (edgeconv)

# *Message-Passing GNNs*



Message Passing GNNs with <mark>injective aggregation</mark> are equivalent to Weisfeiler-Lehman graph isomorphism test

Xu et al. 2019; Morris et al. 2019

# *Weisfeiler-Lehman Test & Chemical precursors of GNNs*



**George Vlăduţ**

**Andrey Lehman**  **Boris Weisfeiler**

Vlăduţ et al. 1959; Weisfeiler, Lehman 1968

# *Weisfeiler-Lehman Test*

$$\phi(\ ,\{\!\{\ ,\ ,\ \}\!\})$$

# *Weisfeiler-Lehman Test*

$$\phi(\bullet,\{\bullet,\bullet,\bullet\})$$



$$\phi(\ ,\{\ ,\ \})$$

Weisfeiler, Lehman 1968

# Weisfeiler-Lehman Test



$\phi(\bullet, \{\!\!\{\bullet, \bullet, \bullet\}\!\!\})$

$\phi(\bullet, \{\!\!\{\bullet, \bullet\}\!\!\})$

# Weisfeiler-Lehman Test

# Weisfeiler-Lehman Test



$\phi(\bigcirc,\{\!\!\{\bigcirc,\bigcirc,\bigcirc\}\!\!\})$

$\phi(\bigcirc,\{\!\!\{\bigcirc,\bigcirc\}\!\!\})$

$\phi(\bigcirc,\{\!\!\{\bigcirc,\bigcirc\}\!\!\})$

Weisfeiler, Lehman 1968

# Weisfeiler-Lehman Test

non-isomorphic graphs that are WL-equivalent

**Necessary but insufficient condition!**

*Message-Passing GNNs have limited expressive power!*

decalin

bicyclopnetyl

# *Towards More Expressive GNNs*



| Higher-order WL tests | Positional & Structural encoding | Subgraph GNNs | Topological message passing |
|---|---|---|---|

Maron et al. 2019
Morris et al. 2019

Monti, Otness et B 2018
Sato 2020
Dwivedi et al. 2020
Bouritsas, Frasca et B 2020
…many more

Papp et al. 2021
Cotta et al. 2021
Zhao et al. 2021
Bevilacqua, Frasca et B, Maron 2021
Frasca et B, Maron 2022

Bodnar, Frasca et B 2021

**CNN**
canonical node ordering

**Equivariant GNN**
+data symmetry group

**Subgraph GNN**
product symmetry group

**Cellular GNN**
high-order complex

*Positional encoding*
+extra features

**GNN**
input graph

MORE STRUCTURE

**CNN**
canonical node
ordering

**Equivariant GNN**
+data symmetry
group

**Subgraph GNN**
product symmetry
group

**Cellular GNN**
high-order
complex

*Positional encoding*
+extra features

**DeepSet/PointNet**
no graph

LESS INTERACTION

**GNN**
input graph

MORE INTERACTION

**Transformer**
learnable graph

MORE STRUCTURE

**CNN**
canonical node
ordering

**Equivariant GNN**
+data symmetry
group

**Subgraph GNN**
product symmetry
group

**Cellular GNN**
high-order
complex

*Positional encoding*
+extra features

LESS INTERACTION

**DeepSet/PointNet**
no graph

**GNN**
input graph

*Graph rewiring*
precomputed graph

**Transformer**
learnable graph

# Weisfeiler-Lehman hierarchy



increasingly expressive test

$k$-WL

3-WL

*CFI graphs*

2-WL

*strongly regular*

1-WL

*d-regular*

## GNN expressive power

Weisfeiler, Lehman 1968 (2-WL); Babai, Mathon 1979 (*k*-WL);
Cai, Fürer, Immerman 1992 (CFI graphs)

# Graphs may be ==unfriendly== for message passing resulting in "bottlenecks"



# Gap between Theory & Practice

## Graph rewiring

Alon, Yahav 2020 (bottlenecks); Hamilton et al. 2017 (neighbour sampling); Klicpera et al. 2019 (diffusion); Topping, Di Giovanni et B 2022 (Ricci flow); Deac et al. 2022 (expanders)

# Graphs vs Meshes vs Grids



**Grid**                    **Mesh**                    **Graph**

# *Graphs vs Meshes vs Grids*



**Grid**
Fixed

**Mesh**

**Graph**

*Graphs vs Meshes vs Grids*

**Grid**
Fixed

**Mesh**
Rotation

**Graph**

# Graphs vs Meshes vs Grids



**Grid**
Fixed

**Mesh**
Rotation

**Graph**
Permutation

# Graphs have the least structure

# Graphs vs Meshes vs Grids



**Grid**                **Mesh**                **Graph**

*Graphs vs Meshes vs Grids*



**Grid**   **Mesh**   **Graph**

Continuous models for GNNs?

Physics-inspired GNNs

# *Physical metaphor of Graph ML*

GNN = dynamic system



$$\dot{\mathbf{X}}(t) = \mathbf{F}_{\boldsymbol{\theta}(t)}(\mathbf{X}(t), \mathcal{G})$$

Haber, Ruthotto 2017; Chen et al. 2019 (Neural ODEs); Xhonneau et al. 2020 (CGNN); Chamberlain, Rowbottom, et B. 2021 (GRAND, BLEND)
Eliasof, Haber 2021 (PDE-GCN); Di Giovanni, Rowbottom et B 2022 (GRAFF), Rusch et B 2022 (GraphCON)

# Physical metaphor of Graph ML

*Input*  *Layer 1*  *Layer 2*

feature space

$t$   $t + \tau$   *time*

$$\mathbf{X}(t + \tau) = \mathbf{X}(t) + \tau \mathbf{F}_{\boldsymbol{\theta}(t)}(\mathbf{X}(t), \mathcal{G})$$

GNN = dynamic system

layers = discretisation of time

graph = coupling function
(discretisation of space)

Haber, Ruthotto 2017; Chen et al. 2019 (Neural ODEs); Xhonneau et al. 2020 (CGNN); Chamberlain, Rowbottom, et B. 2021 (GRAND, BLEND)
Eliasof, Haber 2021 (PDE-GCN); Di Giovanni, Rowbottom et B 2022 (GRAFF), Rusch et B 2022 (GraphCON)

# *Heat Diffusion*

**Newton Law of Cooling:** "the [temperature] a hot body loses in a given time is proportional to the temperature difference between the object and the environment"



( 824 )

with a little preffing, I took a drop thereof, and in it difcover'd a mighty number of living Creatures. I repeated my obfervation the fame evening with the fame fuccefs, but the next day I could find none of them alive; and whereas I had laid that drop upon a fmall Copper Plate, i fancied to my felf that the exhalation of the moifture might be the caufe of their death, and not the cold weather, which at that time was very moderate.

In the beginning of *April* I took the Male feed of a Jack or Pike, but could difcover nothing more than in that of a Cod-fifh, but having added about four times as much Water in quantity as the matter itfelf was, and then making my remarks, I could perceive that the *Animalcula* did not only wax ftronger and fwitter, but, to my great amazement, I faw them move with that celerity, that I could compare it to nothing more than what we have feen with our naked Eye, a River Fifh chafed by its powerful Enemy, which is juft ready to devour it: You muft obferve that this whole Courfe was not longer than the Diameter of a fingle Hair of ones Head.

VII. *Scala graduum Caloris.*

*Calorum Defcriptiones & figna.*

Alor aeris hyberni ubi aqua incipit gelu rigefcere. Innotefcit hic calor accurate locando Thermometrum in nive compreffa quo tempore gelu folvitur.

0,1,2. Calores aeris hyberni.
2,3,4. Calores aeris verni & autumnalis.
4,5,6. Calores aeris æftivi.
6 Calor aeris meridiani circa menfem Julium.
12 1 Calor maximus quem Thermometer ad contactum

**Isaac Newton**

# Heat Diffusion Equation on Graphs



$$\dot{\mathbf{x}}_i(t) = \mathbf{x}_i(t) - \frac{1}{d_i} \sum_{j \in \mathcal{N}_i} a_{ij} \mathbf{x}_j(t)$$

# Heat Diffusion Equation on Graphs



$$\dot{\mathbf{x}}_i(t) = \mathbf{x}_i(t) - \frac{1}{d_i} \sum_{j \in \mathcal{N}_i} a_{ij} \mathbf{x}_j(t)$$

rate of temperature change     self temperature     temperature of the environment

# Heat Diffusion Equation on Graphs



$$\dot{\mathbf{x}}_i(t) \; = \; \underbrace{\frac{1}{d_i} \sum_{j \in \mathcal{N}_i}}_{\substack{\text{divergence} \\ \text{div}}} a_{ij} \underbrace{\Big( \mathbf{x}_i(t) - \mathbf{x}_j(t) \Big)}_{\substack{\text{gradient} \\ -(\nabla \mathbf{X})_{ij}}}$$

# Heat Diffusion Equation on Graphs



$$\dot{\mathbf{X}}(t) = -\mathrm{div}\big(\nabla \mathbf{X}(t)\big)$$

# Heat Diffusion Equation on Graphs



$$\dot{\mathbf{X}}(t) \ = \ \Delta\mathbf{X}(t)$$

# Heat Equation as a prototypical Gradient Flow

$$\dot{\mathbf{X}}(t) = -\nabla \mathcal{E}(\mathbf{X}(t))$$

$$\mathcal{E}_{\mathrm{DIR}}(\mathbf{X}) = \frac{1}{2} \sum_{j \in \mathcal{N}_i} \left\| (\nabla \mathbf{X})_{ij} \right\|^2 = \frac{1}{2} \mathrm{trace}\left( \mathbf{X}^{\mathrm{T}} \Delta \mathbf{X} \right)$$

**G. Dirichlet**

- Heat diffusion equation is the gradient flow of the Dirichlet energy

- "Smoothness" of the node features

- Dirichlet energy <mark>decreases along the flow</mark>

- In the limit $t \to \infty$ results in "oversmoothing"

- Not very expressive: works only in homophilic graphs ("similar neighbours")

Zhou, Schölkopf 2005 (label propagation); Rossi et B 2021 (feature propagation)

# *Gradient Flow Framework (GRAFF)*



**Traditional GNNs**

$$\dot{\mathbf{X}}(t) = \mathbf{F}_{\boldsymbol{\theta}(t)}(\mathbf{X}(t), \mathcal{G})$$

- Parametrize <mark>evolution equations</mark>

Di Giovanni, Rowbottom et B 2022

**GRAFF**

$$\dot{\mathbf{X}}(t) = -\nabla \mathcal{E}_{\boldsymbol{\theta}(t)}(\mathbf{X}(t), \mathcal{G})$$

- Parametrize <mark>energy</mark>

- Derive evolution equation as GF

- Better interpretability

# Generalised parametric Dirichlet energy

$$\mathcal{E}_{\boldsymbol{\theta}}(\mathbf{X}) = \frac{1}{2}\sum_{i=1}^{n}\langle \mathbf{x}_i, \boldsymbol{\Omega}\mathbf{x}_i\rangle - \frac{1}{2}\sum_{j\in\mathcal{N}_i}\bar{a}_{ij}\langle \mathbf{x}_i, \mathbf{W}\mathbf{x}_j\rangle$$

- Energy of a system of particles (nodes) parameterised by $d\times d$ matrices $\boldsymbol{\Omega}$ and $\mathbf{W}$

Di Giovanni, Rowbottom et B 2022

# Generalised parametric Dirichlet energy

$$\mathcal{E}_\theta(\mathbf{X}) \;=\; \frac{1}{2}\sum_{i=1}^{n}\langle \mathbf{x}_i, \mathbf{\Omega}\mathbf{x}_i\rangle \;-\; \frac{1}{2}\sum_{j\in\mathcal{N}_i}\bar{a}_{ij}\langle \mathbf{x}_i, \mathbf{W}\mathbf{x}_j\rangle$$

external energy

- Energy of a system of particles (nodes) parameterised by $d\times d$ matrices $\mathbf{\Omega}$ and $\mathbf{W}$

- External energy term acting on all particles

Di Giovanni, Rowbottom et B 2022

# *Generalised parametric Dirichlet energy*

$$\mathcal{E}_{\boldsymbol\theta}(\mathbf{X}) \;=\; \frac{1}{2}\sum_{i=1}^{n}\langle \mathbf{x}_i, \boldsymbol\Omega\mathbf{x}_i\rangle \;-\; \frac{1}{2}\sum_{j\in\mathcal{N}_i}\bar{a}_{ij}\langle \mathbf{x}_i, \mathbf{W}\mathbf{x}_j\rangle$$

external energy        internal energy
(pair-wise interactions)

- Energy of a system of particles (nodes) parameterised by $d{\times}d$ matrices $\boldsymbol\Omega$ and $\mathbf{W}$

- External energy term acting on all particles

- Internal energy term: interactions between nodes along edges of the graph

  - **Attractive** interactions along positive eigenvectors of $\mathbf{W}$

  - **Repulsive** interactions along negative eigenvectors of $\mathbf{W}$

Di Giovanni, Rowbottom et B 2022

time

attraction

repulsion

Di Giovanni, Rowbottom et B 2022

# Gradient Flow of $\mathcal{E}_{\boldsymbol{\theta}}$

$$\dot{\mathbf{X}}(t) = -\nabla \mathcal{E}_{\boldsymbol{\theta}}\big(\mathbf{X}(t)\big) = -\mathbf{X}(t)\frac{\boldsymbol{\Omega} + \boldsymbol{\Omega}^{\mathrm{T}}}{2} + \bar{\mathbf{A}}\mathbf{X}(t)\frac{\mathbf{W} + \mathbf{W}^{\mathrm{T}}}{2}$$

Di Giovanni, Rowbottom et B 2022

# *Gradient Flow of $\mathcal{E}_{\boldsymbol{\theta}}$*

$$\dot{\mathbf{X}}(t) = -\nabla\mathcal{E}_{\boldsymbol{\theta}}\big(\mathbf{X}(t)\big) = -\mathbf{X}(t)\frac{\boldsymbol{\Omega} + \boldsymbol{\Omega}^{\mathbf{T}}}{2} + \bar{\mathbf{A}}\mathbf{X}(t)\frac{\mathbf{W} + \mathbf{W}^{\mathbf{T}}}{2}$$

*matrices appear only in symmetrized form*

Di Giovanni, Rowbottom et B 2022

# Gradient Flow of $\mathcal{E}_\theta$

$$\dot{\mathbf{X}}(t) = -\mathbf{X}(t)\mathbf{\Omega} + \overline{\mathbf{A}}\mathbf{X}(t)\mathbf{W}$$

- Symmetric matrices $\mathbf{\Omega}$ and $\mathbf{W}$

- Time-independent parameters: $\mathbf{\Omega}(t) = \mathbf{\Omega}$, $\mathbf{W}(t) = \mathbf{W}$

Di Giovanni, Rowbottom et B 2022

# Discretised Gradient Flow of $\mathcal{E}_\theta$

$$\mathbf{X}(t + \tau) = \mathbf{X}(t) + \tau(-\mathbf{X}(t)\mathbf{\Omega} + \overline{\mathbf{A}}\mathbf{X}(t)\mathbf{W})$$

- Residual convolutional-type GNN

- Symmetric weights

- Symmetry constraint does not diminish expressive power

- Shared weights across layers

Di Giovanni, Rowbottom et B 2022; Xu, Zagoruyko, Komodakis 2019 (universal approximation in NNs with symmetric weights)

# GRAFF

$$\mathbf{X}(t + \tau) = \mathbf{X}(t) + \tau\overline{\mathbf{A}}\mathbf{X}(t)\mathbf{W}$$

- Residual
- Symmetric weights
- Shared weights
- No nonlinear activation

- Gradient flow (interpretability)
- $d^2/2$ weights
- Can induce both low- and high-frequency dominated dynamics
- Works with heterophilic graphs

Di Giovanni, Rowbottom et B 2022

# GCN

$$\mathbf{X}(t + \tau) = \tau\sigma\big(\overline{\mathbf{A}}\mathbf{X}(t)\mathbf{W}(t)\big)$$

- Non-residual
- Non-symmetric weights
- Different weights per layer
- Nonlinear activation

- Not a gradient flow
- $Ld^2$ weights
- Only low-frequency dominated dynamics (oversmoothing)
- Only homophilic graphs

Kipf, Welling 2017

# *Homophily vs Heterophily*



Synthetic Cora node classification task

Di Giovanni, Rowbottom et B 2022

*Scalability*



Graph convolution

$\mathbf{X}_{\text{in}} \rightarrow$ ENC $\rightarrow \mathbf{X}$

$\mathbf{X} \mathbf{W}$

$L\mathbf{A}\mathbf{X} \mathbf{W}$

$(L-1)\mathbf{A}^2\mathbf{X} \mathbf{W}^2$

$\vdots$

$\mathbf{A}^L\mathbf{X} \mathbf{W}^L$

$+ \rightarrow$ DEC $\rightarrow \mathbf{X}_{\text{out}}$

Rossi, Frasca et B 2020 (SIGN)

# *Scalability*

Graph convolution



$$\mathbf{X}_{\text{in}} \rightarrow \boxed{\text{ENC}} \rightarrow \mathbf{X}$$

$$\mathbf{X}\,\mathbf{W}$$
$$L\mathbf{A}\mathbf{X}\,\mathbf{W}$$
$$(L-1)\mathbf{A}^2\mathbf{X}\,\mathbf{W}^2$$
$$\vdots$$
$$\mathbf{A}^L\mathbf{X}\,\mathbf{W}^L$$

**Pre-compute**

$$+ \rightarrow \boxed{\text{DEC}} \rightarrow \mathbf{X}_{\text{out}}$$

Rossi, Frasca et B 2020 (SIGN)

# *Spectral analysis of GRAFF*

- $\Delta = \Phi \Lambda \Phi^{\mathrm{T}}$ orthogonal eigendecomposition of the graph Laplacian

- $\mathbf{W} = \Psi \mathbf{M} \Psi^{\mathrm{T}}$ orthogonal eigendecomposition of channel mixing weights

- Output of $L$ layers GRAFF

$$\mathbf{X}(L\tau) = \sum_{k=1}^{d} \sum_{l=0}^{n-1} \left(1 + \tau \mu_k (1 - \lambda_l)\right)^L \langle \mathbf{X}(0), \boldsymbol{\psi}_l \otimes \boldsymbol{\phi}_l \rangle \boldsymbol{\psi}_l \otimes \boldsymbol{\phi}_l$$

- **Low frequencies** ($\lambda_l < 1$) magnified by **positive** eigenvalues of $\mathbf{W}$ ($\mu_k > 0$)

- **High frequencies** ($\lambda_l > 1$) magnified by **negative** eigenvalues of $\mathbf{W}$ ($\mu_k < 0$)

# Spectral analysis of GRAFF

- $\boldsymbol{\Delta} = \boldsymbol{\Phi}\boldsymbol{\Lambda}\boldsymbol{\Phi}^{\mathrm{T}}$ orthogonal eigendecomposition of the graph Laplacian

- $\mathbf{W} = \boldsymbol{\Psi}\mathbf{M}\boldsymbol{\Psi}^{\mathrm{T}}$ orthogonal eigendecomposition of channel mixing weights

- Output of $L$ layers GRAFF

$$\mathbf{X}(L\tau) = \sum_{k=1}^{d}\sum_{l=0}^{n-1}\left(1 + \tau\mu_k(1-\lambda_l)\right)^L \langle\mathbf{X}(0), \boldsymbol{\psi}_l \otimes \boldsymbol{\phi}_l\rangle\boldsymbol{\psi}_l \otimes \boldsymbol{\phi}_l$$

If eigenvalues of **W** are sufficiently negative, then GRAFF dynamics is high-frequency dominant: $\mathcal{E}_{\mathrm{DIR}}(\mathbf{X}(t)/\|\mathbf{X}(t)\|) \to \lambda_{\max}/2$

**No oversmoothing!**

# Non-homogeneous Diffusion in Image Processing

$$\dot{\mathbf{X}}(t) = -\mathrm{div}\left(\frac{\nabla\mathbf{X}(t)}{1 + c\|\nabla\mathbf{X}(t)\|^2}\right)$$

<span style="color:green">*edge indicator*</span>



"Do not diffuse across edges"

Perona, Malik 1990

# Non-homogeneous Diffusion in Image Processing



Homogeneous
diffusion

Non-homogeneous
diffusion

Perona, Malik 1990; Kimmel et al. 1997; Sochen et al. 1998; Tomasi, Manduchi 1998; Weickert 1998; Buades et al. 2005

# Diffusion in Image Processing



Perona, Malik 1990; Kimmel et al. 1997; Sochen et al. 1998; Tomasi, Manduchi 1998; Weickert 1998; Buades et al. 2005

# Non-homogeneous Diffusion Equation on Graphs

$$\dot{\mathbf{x}}_i(t) = \sum_{j \in \mathcal{N}_i} a\left(\mathbf{x}_i(t), \mathbf{x}_j(t)\right)\left(\mathbf{x}_j(t) - \mathbf{x}_i(t)\right)$$

*Explicit (Forward Euler)* discretization with timestep $\tau$:

$$\mathbf{x}_i(t + \tau) = \mathbf{x}_i(t) + \tau \sum_{j \in \mathcal{N}_i} a\left(\mathbf{x}_i(t), \mathbf{x}_j(t)\right)\left(\mathbf{x}_j(t) - \mathbf{x}_i(t)\right)$$

Chamberlain, Rowbottom, et B. 2021

# Non-homogeneous Diffusion Equation on Graphs

$$\dot{\mathbf{x}}_i(t) = \sum_{j \in \mathcal{N}_i} a\Big(\mathbf{x}_i(t), \mathbf{x}_j(t)\Big)\Big(\mathbf{x}_j(t) - \mathbf{x}_i(t)\Big)$$

*Explicit (Forward Euler)* discretization with timestep $\tau$:

$$\mathbf{x}_i(t+\tau) = \sum_{j \in \mathcal{N}_i} a\Big(\mathbf{x}_i(t), \mathbf{x}_j(t)\Big)\mathbf{x}_j(t)$$

normalised $\sum_j a_{ij} = 1$
Unit step $\tau = 1$

## GAT is a particular discretisation of graph diffusion

Chamberlain, Rowbottom, et B. 2021

Geometric Flows & Graph Rewiring

# Spatial Derivative: Graph Rewiring?



Different discretisations of 2D Laplacian

# *Images as embedded manifolds*



$$\dot{\mathbf{X}} = -\mathrm{div}(a(\mathbf{X})\nabla\mathbf{X})$$

Non-linear diffusion

$$\dot{\mathbf{Z}} = \Delta_{\mathbf{G}}\mathbf{Z}$$

Non-Euclidean diffusion

Kimmel et al. 1997; Sochen et al. 1998

# Beltrami flow

- Consider image as embedded 2-*manifold*

$$\mathbf{Z}(\mathbf{u}) = \big(\mathbf{u}, \alpha \mathbf{X}(\mathbf{u})\big)$$

- *Pullback metric:* 2×2 matrix

$$\mathbf{G} = \mathbf{I} + \alpha^2 \big(\nabla_{\mathbf{u}} \mathbf{X}(\mathbf{u})\big)^{\mathsf{T}} \nabla_{\mathbf{u}} \mathbf{X}(\mathbf{u})$$

- *Beltrami flow* = gradient flow of the *Polyakov energy* (harmonic energy of the embedding used in string theory)



feature coordinates

$x$

$u_1$

$u_2$

positional coordinates

$$\dot{\mathbf{Z}} = \Delta_{\mathbf{G}} \mathbf{Z}$$

**Eugenio Beltrami**

Kimmel et al. 1997; Sochen et al. 1998

# *Graph Beltrami flow*



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

- **Graph Beltrami flow**

$$\dot{\mathbf{z}}_i(t) = \sum_{j \in \mathcal{N}_i} a\left(\mathbf{z}_i(t), \mathbf{z}_j(t)\right)\left(\mathbf{z}_j(t) - \mathbf{z}_i(t)\right)$$

Chamberlain, Rowbottom, et B. 2021

# *Graph Beltrami flow*

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

- **Graph Beltrami flow**

$$\dot{\mathbf{z}}_i(t) = \sum_{j \in \mathcal{N}_i} a\Big(\mathbf{z}_i(t), \mathbf{z}_j(t)\Big)\Big(\mathbf{z}_j(t) - \mathbf{z}_i(t)\Big)$$

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



Chamberlain, Rowbottom, et B. 2021

# Graph Beltrami flow

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

- **Graph Beltrami flow**

$$\dot{\mathbf{z}}_i(t) = \sum_{j \in \mathcal{N}_i} a\left(\mathbf{z}_i(t), \mathbf{z}_j(t)\right)\left(\mathbf{z}_j(t) - \mathbf{z}_i(t)\right)$$

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

  - Evolution of $\mathbf{u}$ = graph rewiring

Chamberlain, Rowbottom, et B. 2021
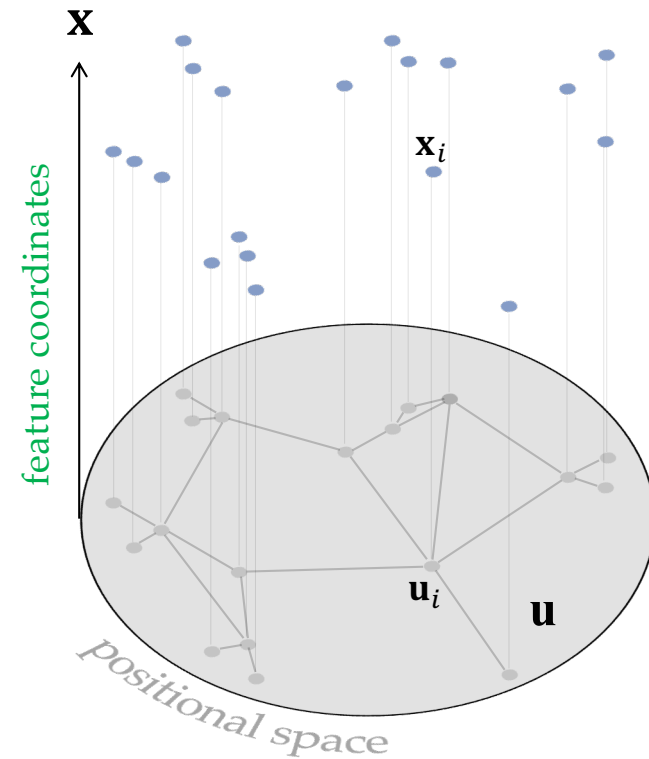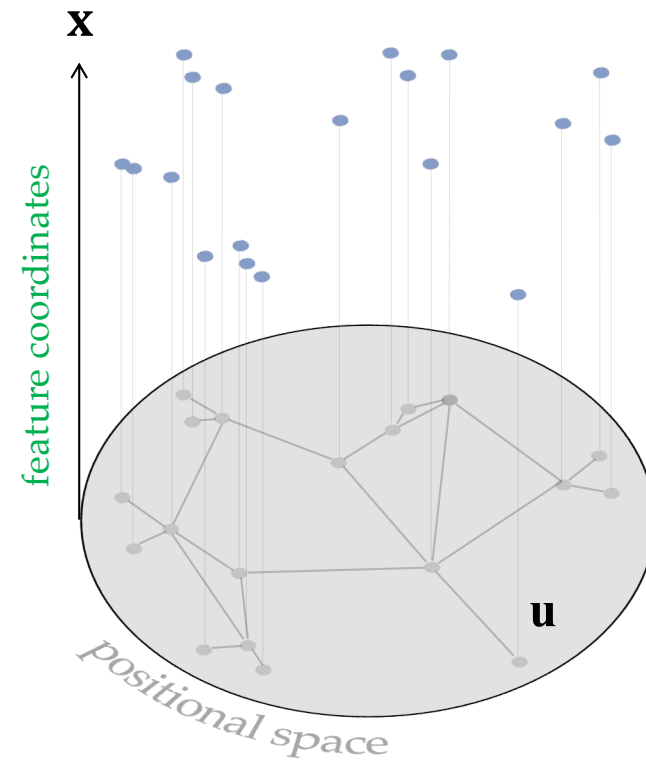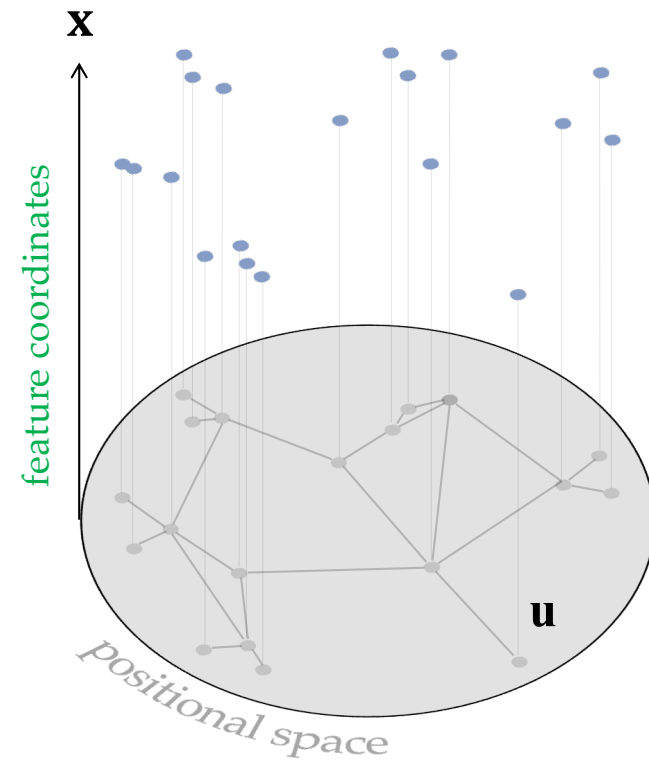
# *Graph Beltrami flow*

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

- **Graph Beltrami flow**

$$\dot{\mathbf{z}}_i(t) = \sum_{j \in \mathcal{N}'_i} a\Big(\mathbf{z}_i(t), \mathbf{z}_j(t)\Big)\Big(\mathbf{z}_j(t) - \mathbf{z}_i(t)\Big)$$

*rewired graph*

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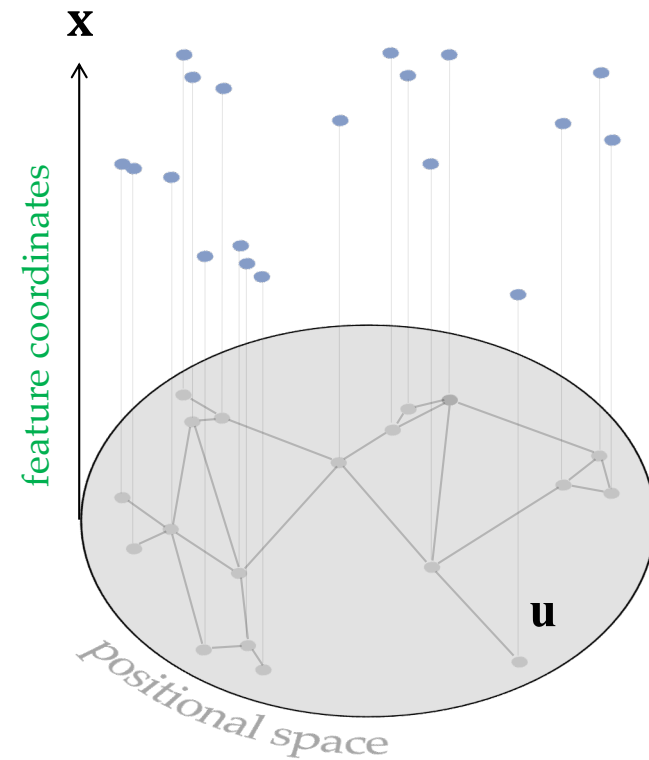- Evolution of $\mathbf{u}$ = graph rewiring



Chamberlain, Rowbottom, et B. 2021

# Graph Beltrami flow



Evolution of positional/feature components + rewiring of the Cora graph

# *Ricci flow*

- **Ricci flow:** "diffusion of the Riemannian metric"

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



Evolution of a manifold under Ricci flow

**G. Ricci-Curbastro**     **R. Hamilton**

Ricci 1903; Hamilton 1988;

Science

22 December 2006 | $10

Breakthrough
of the Year

The
Poincaré
Conjecture
PROVED

Ricci 1903; Hamilton 1988; Perelman 2003

G. Perelman        G. Ricci-Curbastro        R. Hamilton

# *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

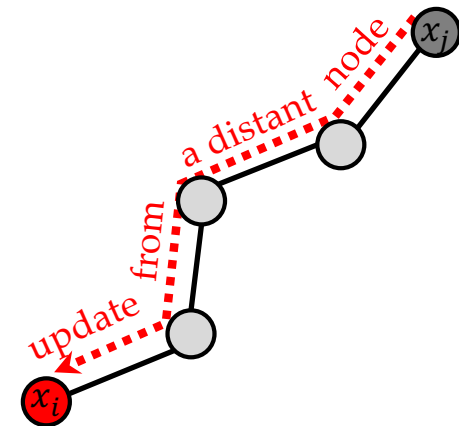# *Over-squashing*

- Consider an MPNN of the form

$$\mathbf{x}_i^{(k+1)} = \sigma \left( \mathbf{W}_1 \mathbf{x}_i^{(k)} + \sum_j a_{ij} \, \mathbf{W}_2 \mathbf{x}_j^{(k)} \right)$$

- $L = depth$ (number of layers)
- $p = width$ (hidden dimension)
- Nonlinearity $\sigma$ is $c_\sigma$-*Lipschitz-continuous*
- $w$ = maximum element of weight matrices $\mathbf{W}_1$, $\mathbf{W}_2$

> **Theorem (Sensitivity bound):** For any $i, j \in V$
> $$\left\| \frac{\partial \mathbf{x}_i^{(L)}}{\partial \mathbf{x}_j^{(0)}} \right\|_1 \leq (c_\sigma w p)^L (\mathbf{I} + \mathbf{A})_{ij}^L$$

Topping, Di Giovanni et B 2021; Di Giovanni et B 2023



**Over-squashing:** small Jacobian $\left\| \partial \mathbf{x}_i^{(L)} / \partial \mathbf{x}_j^{(0)} \right\|$ indicates poor information propagation from input node
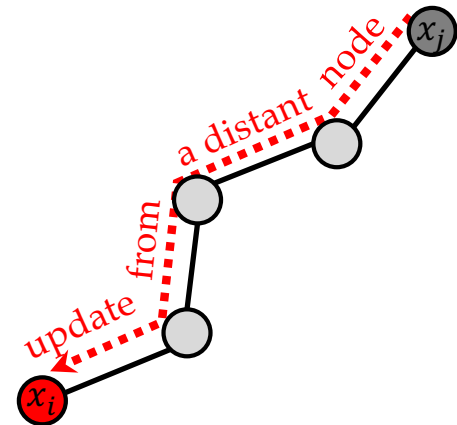
# *Over-squashing*

- Consider an MPNN of the form

$$\mathbf{x}_i^{(k+1)} = \sigma\left(\mathbf{W}_1\mathbf{x}_i^{(k)} + \sum_j a_{ij}\,\mathbf{W}_2\mathbf{x}_j^{(k)}\right)$$

- $L = depth$ (number of layers)
- $p = width$ (hidden dimension)
- Nonlinearity $\sigma$ is $c_\sigma$-*Lipschitz-continuous*
- $w$ = maximum element of weight matrices $\mathbf{W}_1$, $\mathbf{W}_2$

**Theorem (Sensitivity bound):** For any $i, j \in V$

$$\left\|\frac{\partial\mathbf{x}_i^{(L)}}{\partial\mathbf{x}_j^{(0)}}\right\|_1 \leq (c_\sigma w p)^L (\mathbf{I} + \mathbf{A})_{ij}^L$$

model    topology

Topping, Di Giovanni et B 2021; Di Giovanni et B 2023



**Over-squashing:** small Jacobian $\left\|\partial\mathbf{x}_i^{(L)}/\partial\mathbf{x}_j^{(0)}\right\|$ indicates poor information propagation from input node

# *Preventing over-squashing*

$$\left\|\frac{\partial \mathbf{x}_i^{(L)}}{\partial \mathbf{x}_j^{(0)}}\right\|_1 \leq (c_\sigma w p)^L (\mathbf{I} + \mathbf{A})_{ij}^L$$

model  topology

- **Width** $p$ helps mitigate over-squashing (potentially at the risk of worse generalization)

- **Depth** $L$ does not help
  - If $L \sim \mathrm{diam}(G)$, over-squashing occurs between distant nodes
  - If $L \gg 1$, we transition from over-squashing to vanishing gradients

Di Giovanni et B 2023

# *Preventing over-squashing*

$$\left\| \frac{\partial \mathbf{x}_i^{(L)}}{\partial \mathbf{x}_j^{(0)}} \right\|_1 \leq (c_\sigma w p)^L (\mathbf{I} + \mathbf{A})_{ij}^L$$

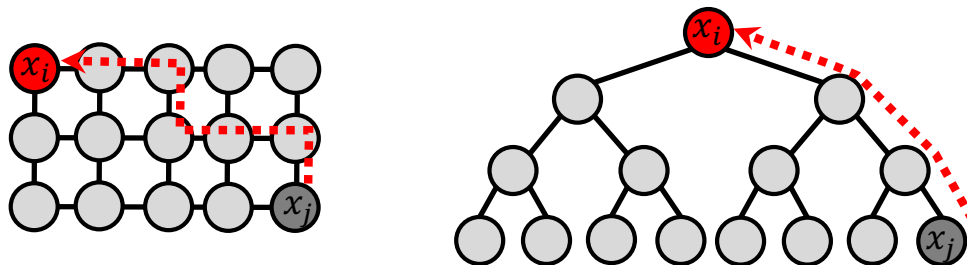$$\underset{\text{model}}{} \qquad \underset{\text{topology}}{}$$

- **Width** $p$ helps mitigate over-squashing (potentially at the risk of worse generalization)

- **Depth** $L$ does not help
  - If $L \sim \text{diam}(G)$, over-squashing occurs between distant nodes
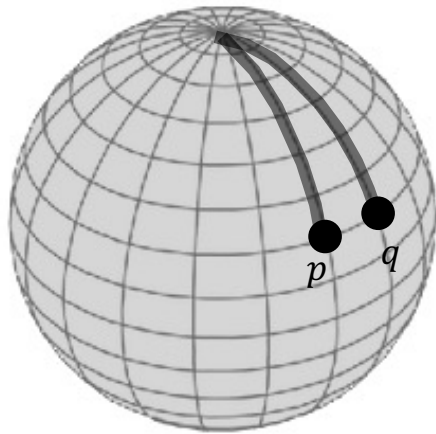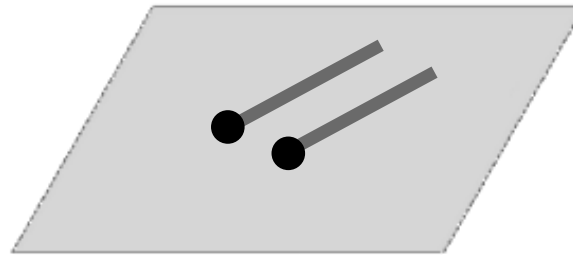  - If $L \gg 1$, we transition from over-squashing to vanishing gradients

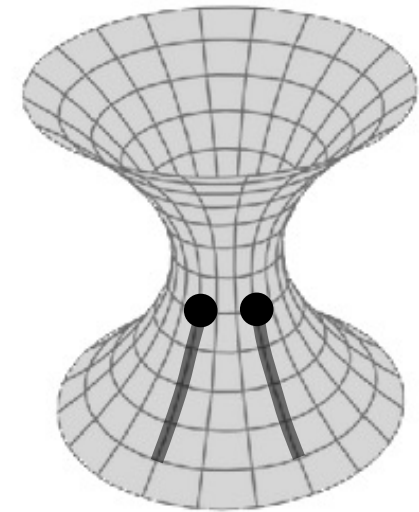- **Topology** of $G$ has the largest effect on over-squashing
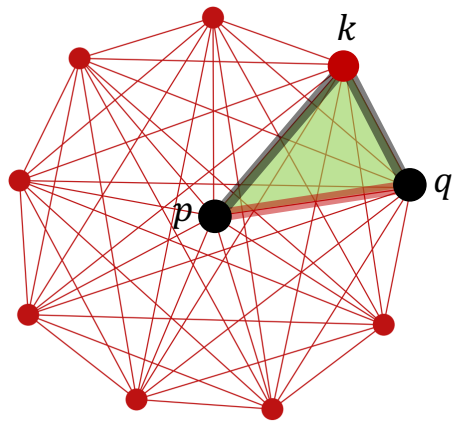
# Ricci Curvature on Manifolds



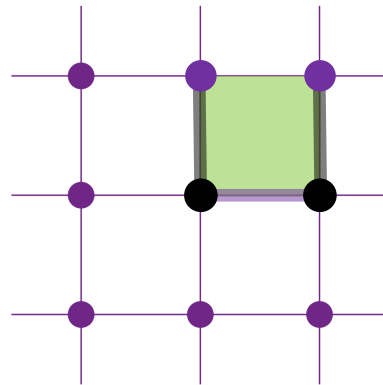Spherical (>0)                Euclidean (=0)                Hyperbolic (<0)

**"geodesic dispersion"**

Ricci 1903

# Ricci Curvature on Graphs



Clique (>0)          Grid (=0)          Tree (<0)

Forman 2003; Ollivier 2007; Topping, di Giovanni, et B. 2021

# Balanced Forman Curvature

**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

Triangles based at i~j

Max number of 4-cycle based at i~j
traversing the same node

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

Degree of i

Neighbours of i forming a 4-cycle
based at i~j (w/o diagonals)

Forman 2003; Topping, di Giovanni, et B. 2021

# Balanced Forman Curvature

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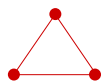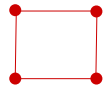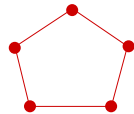
$$\text{Ric}(i,j) = \frac{2}{d_i} + \frac{2}{d_j} + 2\frac{|\sharp_\Delta(i,j)|}{\max\{d_i, d_j\}} + \frac{|\sharp_\Delta(i,j)|}{\min\{d_i, d_j\}} + \frac{\gamma_{\max}^{-1}}{\max\{d_i, d_j\}}\left(\left|\sharp_\square^i(i,j)\right| + \left|\sharp_\square^j(i,j)\right|\right) - 2$$
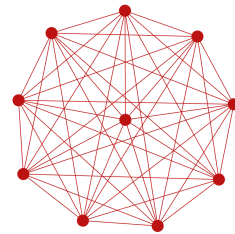


Cycle $C_3$: $\frac{3}{2}$     $C_4$: 1     $C_{n \geq 5}$: 0     Clique $K_n$: $\frac{n}{n-1}$     Grid $G_n$: 0     Tree $T_r$: $\frac{4}{r+1} - 2$

Forman 2003; Topping, di Giovanni, et B. 2021

*Over-squashing & Bottleneck via Curvature*

**Theorem 1 (main result):** Consider an MPNN with $L \geq 2$ layers and $|\nabla \phi_\ell| \leq \alpha$ and $|\nabla \psi_\ell| \leq \beta$. Let $i \sim j$ with $d_i \leq d_j$ and assume $\exists \delta$ s.t. $0 < \delta < \max\{d_i, d_j\}^{1/2}$, $\delta < \gamma_{\max}^{-1}$ and $\mathrm{Ric}(i,j) \leq -2 + \delta$. Then, there exist nodes $Q \subset \{s : d_G(i,s) = 2\}$ of size $|Q| > 1/\delta$ s.t.

small $\delta$ =
negative curvature

more nodes

$$\frac{1}{|Q|} \sum_{k \in Q} \left| \frac{\partial x_k^{(\ell+2)}}{\partial x_i^{(\ell)}} \right| < (\alpha \beta)^2 \delta^{1/4}$$
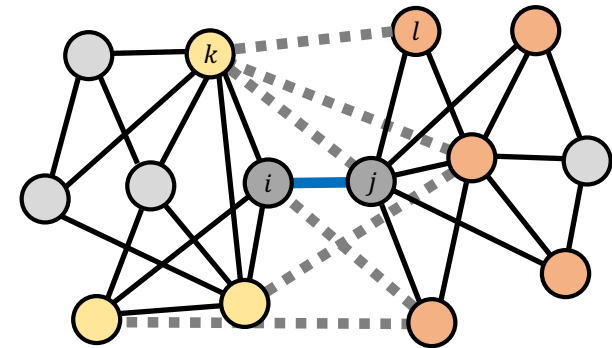
stronger over-squashing

**<span style="color:red">Over-squashing is caused by strongly negatively-curved edges!</span>**

<span style="color:gray">Topping, di Giovanni, et B. 2021</span>
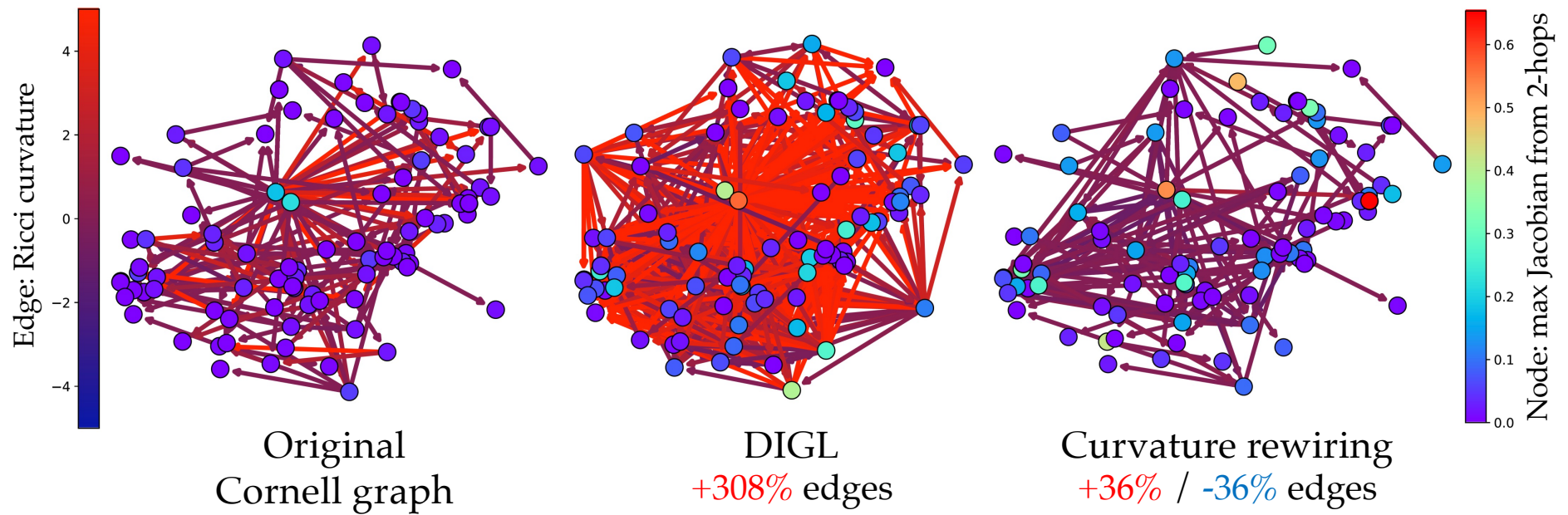
# Stochastic Discrete Ricci Flow (SDRF)

**Input:** graph $G = (V, E)$, temperature $\tau > 0$, (optional $C$)

- For edge $i \sim j$ with smallest $\text{Ric}(i, j)$

  - Calculate the improvement $\delta_{kl} = \text{Ric}_{G'}(i, j) - \text{Ric}(i, j)$ from adding edge $k \sim l$ with $k \in B_1(i)$ and $l \in B_1(j)$

  - Sample index $k, l$ with probability $\text{Softmax}(\tau \delta_{kl})$ and add edge $k \sim l$ to $E'$

- (optional) Remove edge $i \sim j$ with largest $\text{Ric}(i, j) > C$
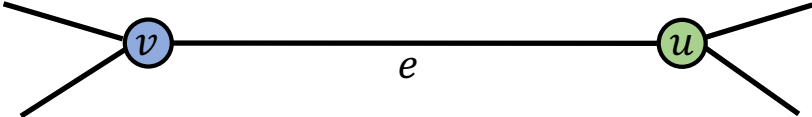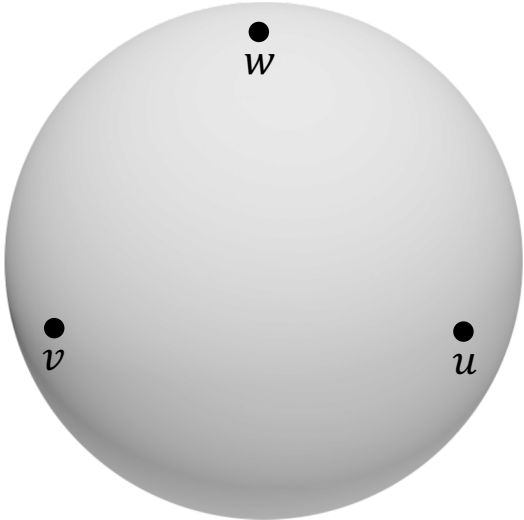
**Output:** new graph $G' = (V, E')$



Topping, di Giovanni, et B. 2021

# *Curvature- vs Diffusion-based Rewiring*



Original
Cornell graph

DIGL
+308% edges

Curvature rewiring
+36% / -36% edges

Edge: Ricci curvature

Node: max Jacobian from 2-hops

Topping, di Giovanni, et B. 2021; Klicpera et al. 2019 (DIGL)

Even more Exotic Stuff

# Cellular Sheaves



Graph

Manifold

Bodnar, di Giovanni, et B. 2022

# Cellular Sheaves



Cellular sheaf $\mathcal{F}$
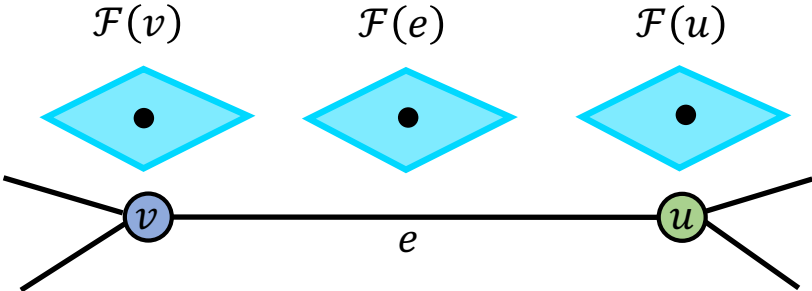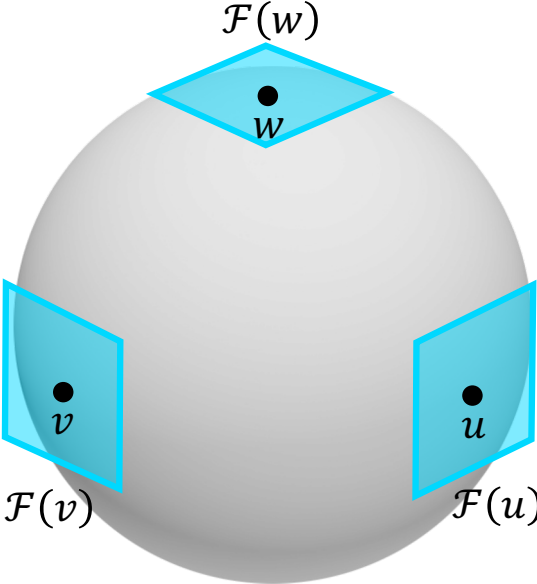
Manifold + Connection

Bodnar, di Giovanni, et B. 2022
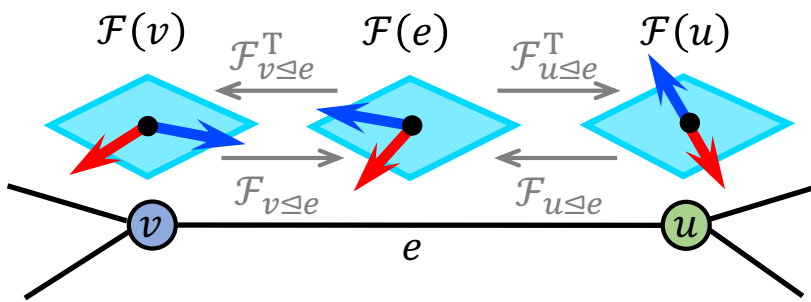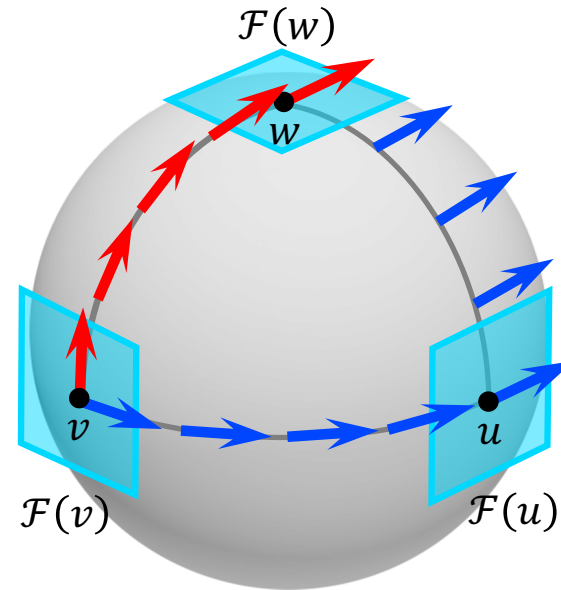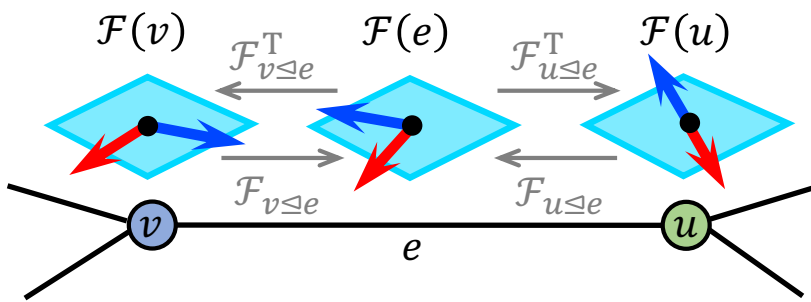
# *Cellular Sheaves*
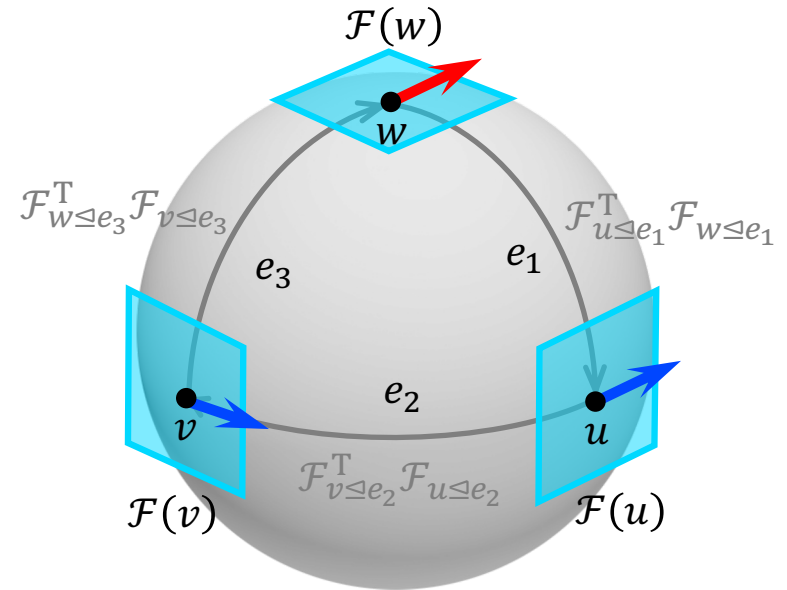


Cellular sheaf $\mathcal{F}$

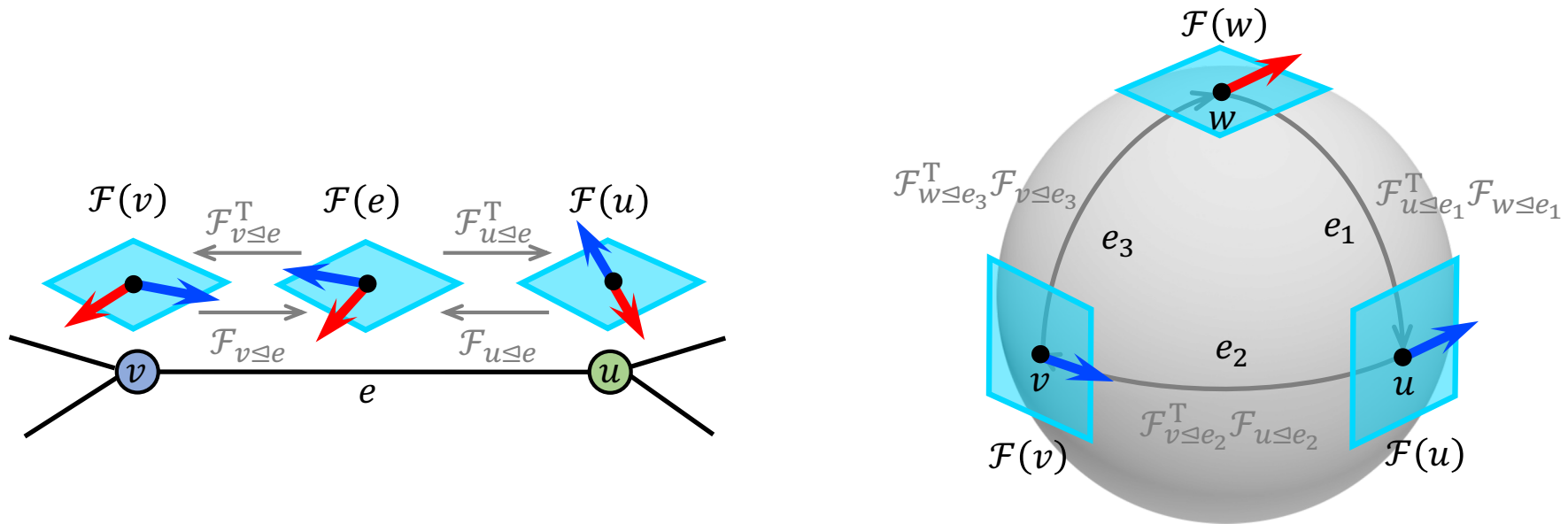Manifold + Connection

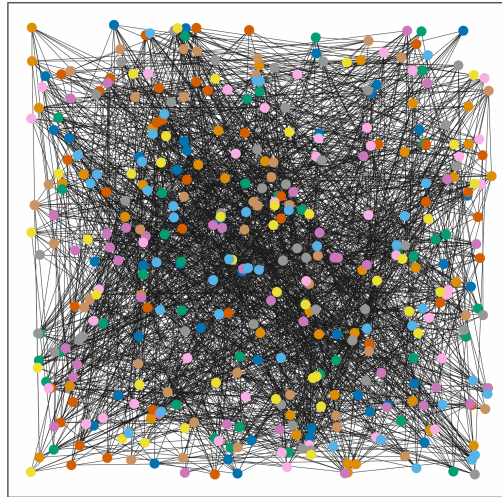# Cellular Sheaves



Cellular sheaf $\mathcal{F}$

Analogy to parallel transport on manifolds

*Cellular Sheaves*



Endow graph with "geometry" leading to richer diffusion with better separation, ability to cope with heterophily, and no oversmoothing

Bodnar, di Giovanni, et B. 2022

# *Diffusion on Cellular Sheaves*



$$\dot{\mathbf{X}}(t) = \Delta_{\mathcal{F}}\mathbf{X}(t) \ \text{ with i.c. } \mathbf{X}(0) = \mathbf{X}$$

Node classification = limit of sheaf diffusion equation
with an appropriate sheaf

Bodnar, di Giovanni, et B. 2022

# *Alternative to Weisfeiler-Lehman for expressive power?*

| Graph type | #Node classes | Sheaf class $\mathcal{F}$, dim=$d$ | |
|---|---|---|---|
| Homophilic | 2 | Symmetric $d$=1 | ✓ |
| Heterophilic | 2 | Symmetric $d$=1 | ✗ |
| | 2 | Non-symmetric $d$=1 | ✓ |
| | ≥3 | Non-symmetric $d$=1 | ✗ |
| | ≤2$d$ | Orthogonal, $d$={2,4} | ✓ |

The capability of sheaf diffusion to solve node classification problem in the limit
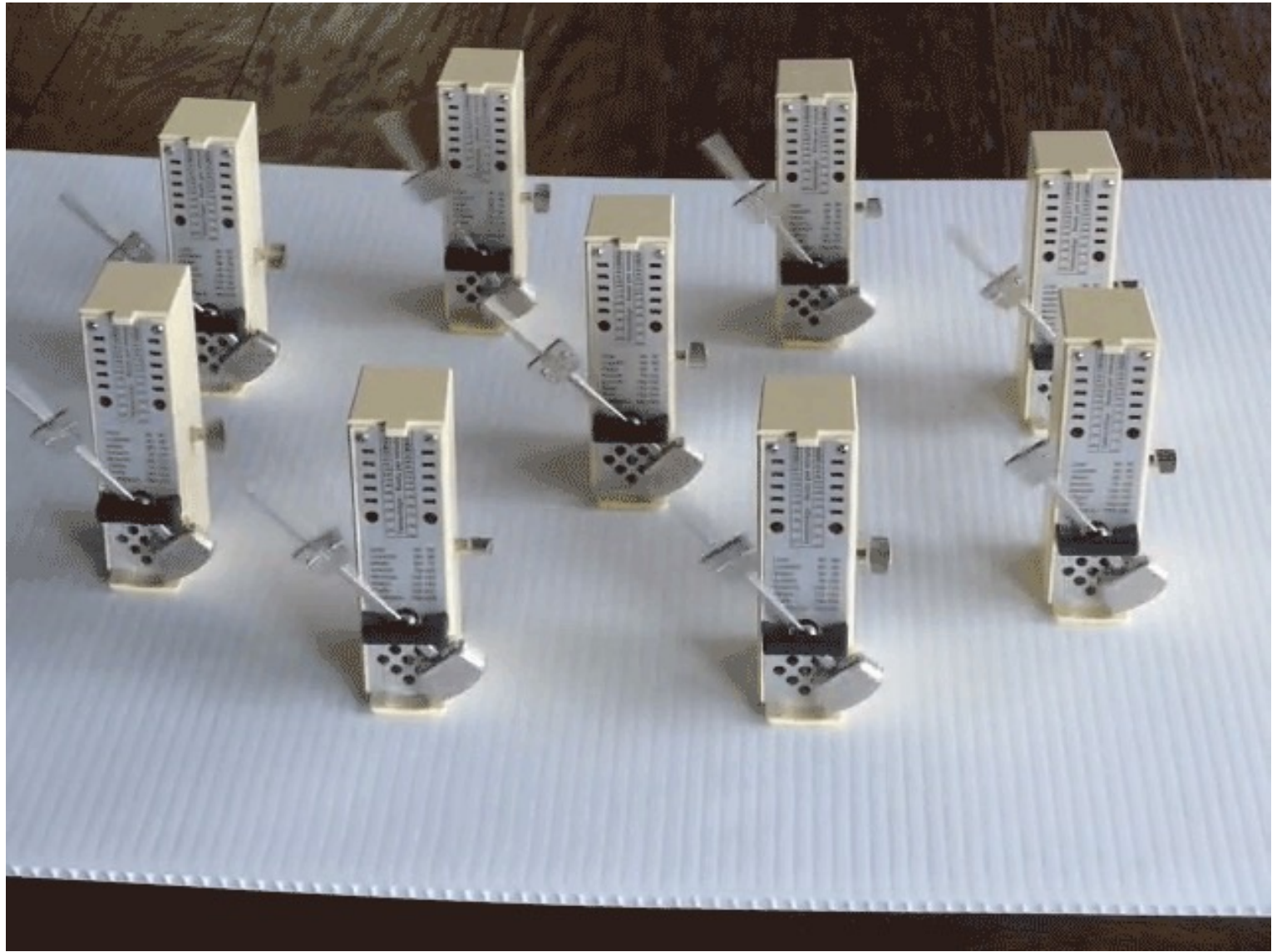
Bodnar, di Giovanni, et B. 2022

## What do we gain from physics-inspired GNNs?

- New perspectives on old problems (e.g. oversmoothing, bottlenecks, etc.)
- New architectures
  - Many GNNs can be formalised as a discretised Graph Diffusion equation
  - More efficient solvers (multistep, adaptive, implicit, multigrid, etc.)
  - Implicit schemes = multi-hop filters
- Principled architectural choices (residual connection, shared symmetric weights)
- Theoretical guarantees (e.g. stability, convergence, expressive power, etc.)
- Deep links to other fields less known in GNN literature (e.g. differential geometry and algebraic topology)
- Other physical models

# *Graph-Coupled Oscillators*



Dynamics of a system of coupled oscillators on a molecular graph

Thank you!