

First Summer School on Geometric Deep Learning – Pescara 2022

GEOMETRIC DEEP LEARNING THE ERLANGEN PROGRAMME OF ML

Michael Bronstein



UNIVERSITY OF
OXFORD

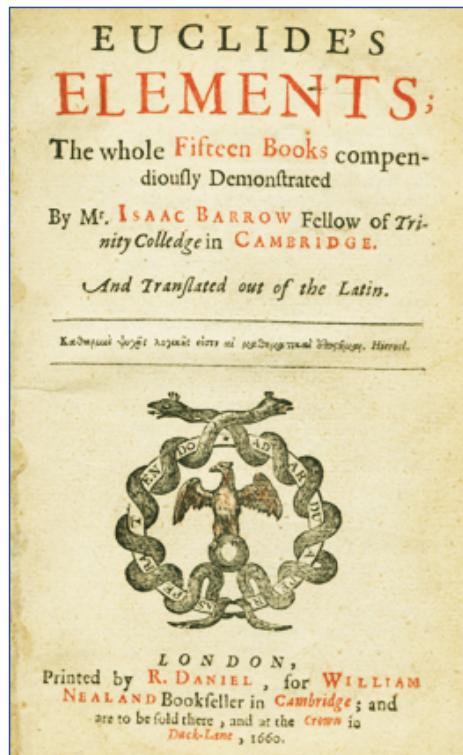


“Symmetry, as wide or as narrow as you may define its meaning, is one idea by which man through the ages has tried to comprehend and create order, beauty, and perfection”



Hermann Weyl

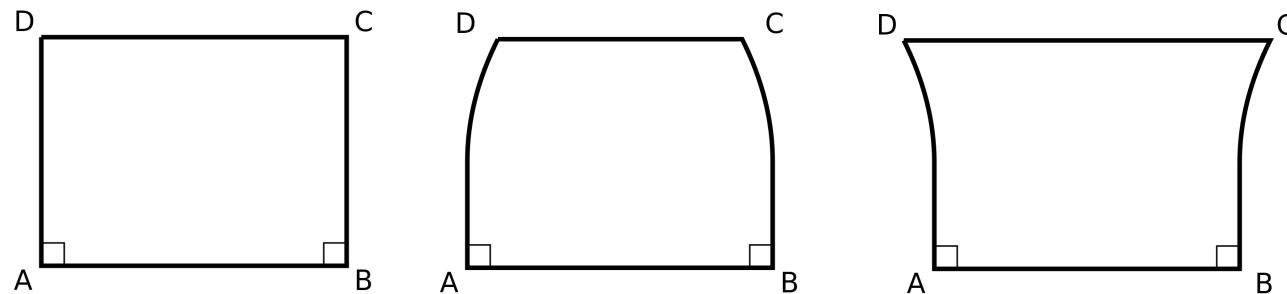
The Origins



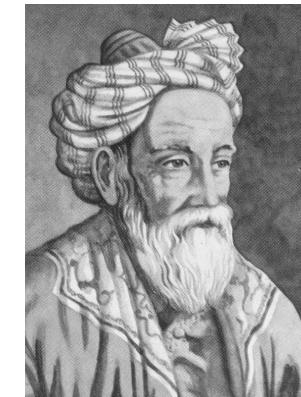
Euclid

~300 BC

Early attempts



Khayyam-Saccheri quadrilateral



Omar Khayyam

“Three cases of angles in a quadrilateral: Fifth Postulate follows from the right-angle assumption”

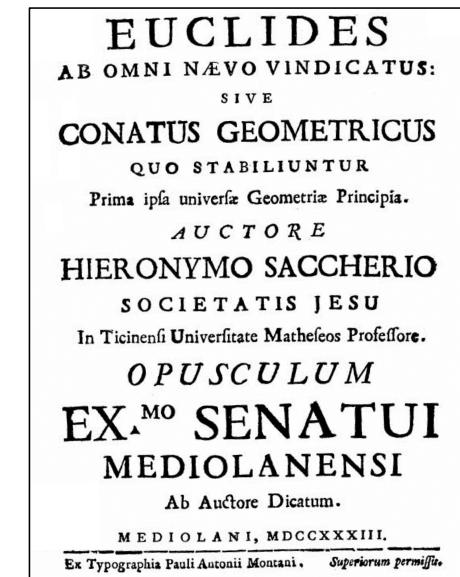
1077

Early attempts

Et hujus quidem (post multa, ne dicam omnia, conditionatè expensa) absolutam falsitatem in XXXIII. tandem ostendo, quia repugnantis naturæ lineæ rectæ, circa quam multa ibi interfero necessaria Lemmata . Tandem verò in præcedente Propositione absolutè demonstro sibi ipsi repugnantem hypothesin anguli acuti .

"repugnant to the nature of straight lines"

— Giovanni Saccheri



1736

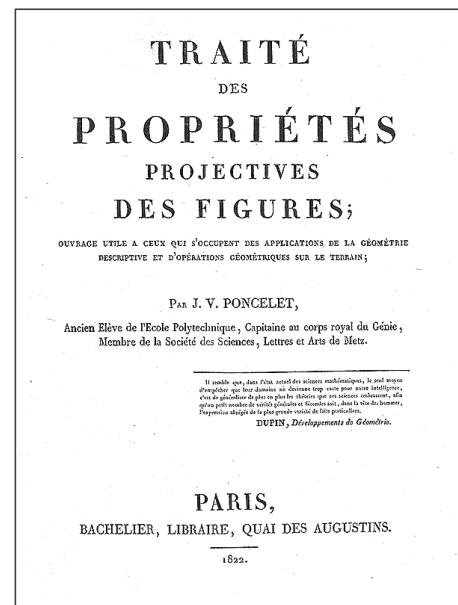
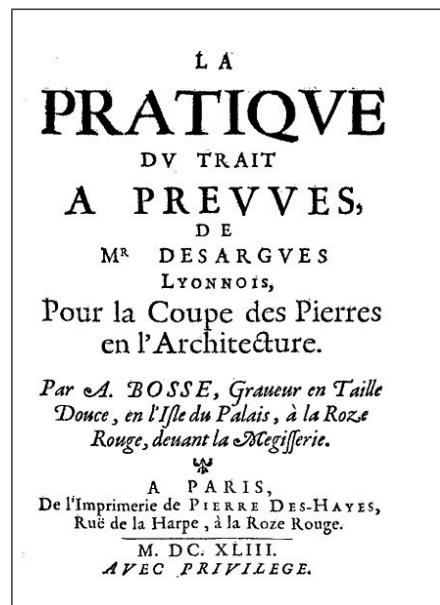


End of Euclid's Monopoly



G. Desargues

1643



"Projective geometry"



J. V. Poncelet

1822

End of Euclid's Monopoly



C. F. Gauss

“I have discovered such wonderful things that I was amazed...out of nothing I have created a strange new world.” — Jánus Bolyai to his father

“To praise it would amount to praising myself. For the entire content of the work...coincides almost exactly with my own meditations [in the] past thirty or thirty-five years.”

— Gauss to Farkas Bolyai



J. Bolyai

1823

Gauss ~1800; Bolyai (1823) 1832

End of Euclid's Monopoly

"In geometry I find certain imperfections which I hold to be the reason why this science [...] can as yet make no advance from that state in which it came to us from Euclid. I consider [...] the momentous gap in the theory of parallels, to fill which all efforts of mathematicians have so far been in vain."

Lobachevsky (1826) 1829

178

О НАЧАЛАХЪ ГЕОМЕТРИИ (*).

(Г. Лобачевского.)

Кажется, трудность понятий увеличивается по мѣрѣ ихъ приближенія къ начальными истинамъ въ природѣ; также какъ она возрастаетъ въ другомъ направлѣніи, къ той границѣ, куда стремится умъ за новыми познаніями. Вотъ почему трудности въ Геометрии должны принадлежать, впервыхъ, самому предмету. Даѣвъ средства, къ которымъ надобно прибѣгнуть чтобы достичнуть здѣсь послѣдней строгости, едва ли могутъ отвѣтить цѣли и простотѣ сего ученія. Тѣ, которые хотѣли удовлетворить симъ требованіямъ, заключили себя въ такой тѣсной кругѣ, что всѣ усилия ихъ не могли быть вознаграждены успѣхомъ. Наконецъ скажемъ и то, что со временеми Ньютона и Декарта, вся Математика, сдѣлавшись Аналитикой, пошла столь быстрыми шагами впередъ, что оставила далеко за собой то ученіе, безъ котораго могла уже об-

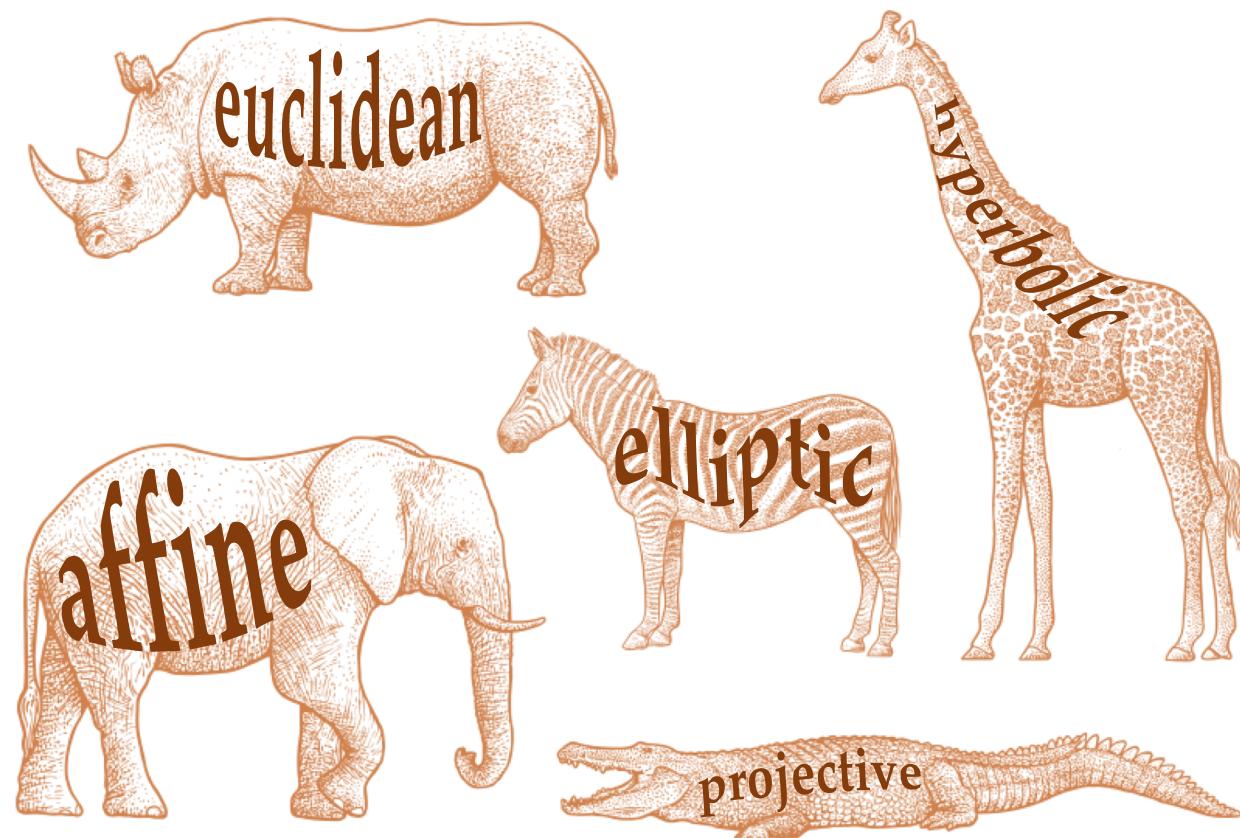
(*.) Извлечено самимъ Сочинителемъ изъ разсужденія, подъ названіемъ: *Exposition succincte des principes de la Géométrie etc.*, читаннаго имъ въ засѣданіи Отдѣленія Физико-Математическихъ наукъ, 12 Февраля 1826 года.



N. Lobachevsky

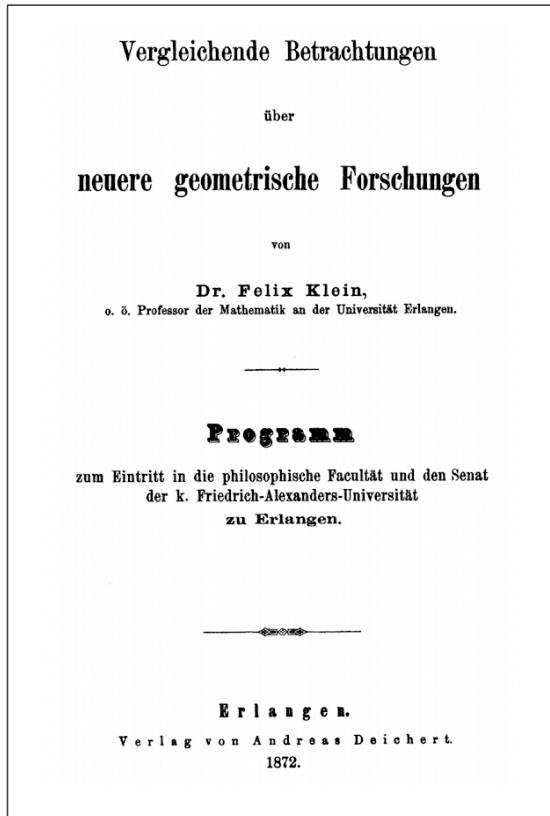
1829

Nineteenth Century Zoo of Geometries



The Erlangen Programme

“Given a [homogeneous] manifold and a transformation group acting [transitively] on it, to investigate those properties of figures on that manifold which are invariant under transformations of that group”



F. Klein

1872

Klein 1872

The Erlangen Programme

	Euclidean	Affine	Projective
<i>angle</i>	+	—	—
<i>distance</i>	+	—	—
<i>area</i>	+	—	—
<i>parallelism</i>	+	+	—
<i>intersection</i>	+	+	+

Noether's Theorem

“Every [differentiable] symmetry of the action of a physical system [with conservative forces] has a corresponding conservation law”

Invariante Variationsprobleme.
(F. Klein zum fünfzigjährigen Doktorjubiläum.)
Von
Emmy Noether in Göttingen.
Vorgelegt von F. Klein in der Sitzung vom 26. Juli 1918¹⁾.
Es handelt sich um Variationsprobleme, die eine kontinuierliche Gruppe (im Lieschen Sinne) gestatten; die daraus sich ergebenden Folgerungen für die zugehörigen Differentialgleichungen finden ihren allgemeinsten Ausdruck in den in § 1 formulierten, in den folgenden Paragraphen bewiesenen Sätzen. Über diese aus Variationsproblemen entspringenden Differentialgleichungen lassen sich viel präzisere Aussagen machen als über beliebige, eine Gruppe gestattende Differentialgleichungen, die den Gegenstand der Lieschen Untersuchungen bilden. Das folgende beruht also auf einer Verbindung der Methoden der formalen Variationsrechnung mit denen der Lieschen Gruppentheorie. Für spezielle Gruppen und Variationsprobleme ist diese Verbindung der Methoden nicht neu; ich erwähne Hamel und Herglotz für spezielle endliche, Lorentz und seine Schüler (z. B. Fokker), Weyl und Klein für spezielle unendliche Gruppen²⁾. Insbesondere sind die zweite Kleinsche Note und die vorliegenden Ausführungen gegenseitig durch einander beeinflusst.

1) Die endgültige Fassung des Manuskriptes wurde erst Ende September eingereicht.
2) Hamel: Math. Ann. Bd. 59 und Zeitschrift f. Math. u. Phys. Bd. 50. Herglotz: Ann. d. Phys. (4) Bd. 36, bes. § 9, S. 511. Fokker, Verslag d. Amsterdamer Akad., 27/I. 1917. Für die weitere Literatur vergl. die zweite Note von Klein: Göttinger Nachrichten 19. Juli 1918.
In einer eben erschienenen Arbeit von Kneser (Math. Zeitschrift Bd. 2) handelt es sich um Aufstellung von Invarianten nach ähnlicher Methode.
Kgl. Ges. d. Wiss. Nachrichten, Math.-phys. Klasse, 1918, Heft 2. 17



E. Noether

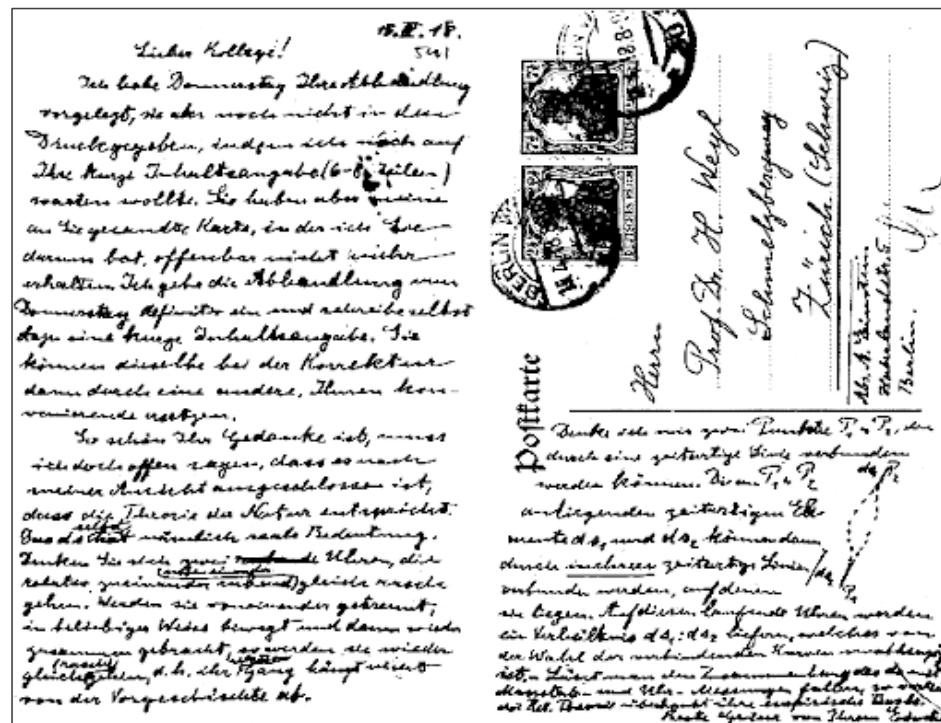
1918

Noether 1918

Gauge invariance

Lieber Kollege! —

Postcard dated 15 April
1918 from Einstein to
Weyl arguing with his
initially proposed
gauge theory

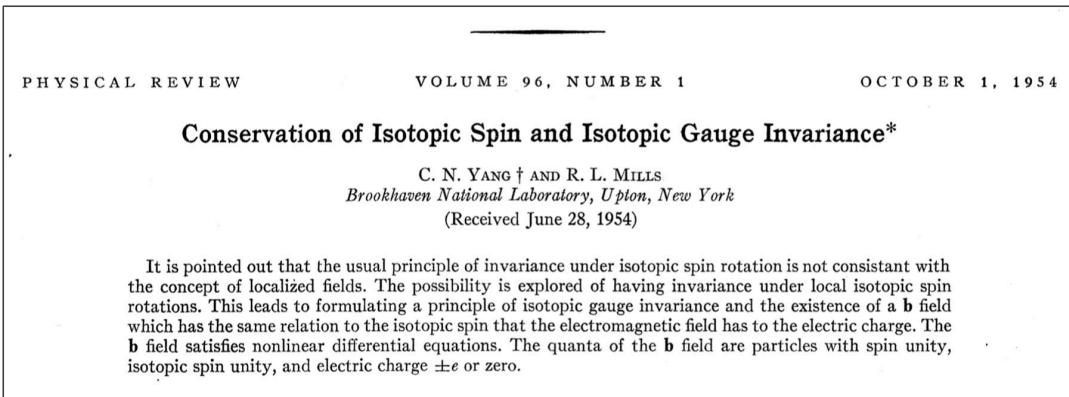


H. Weyl

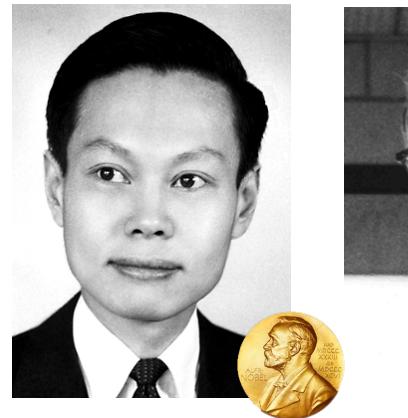
1929

Weyl 1919; 1929 (see Straumann 1987)

Unification of forces

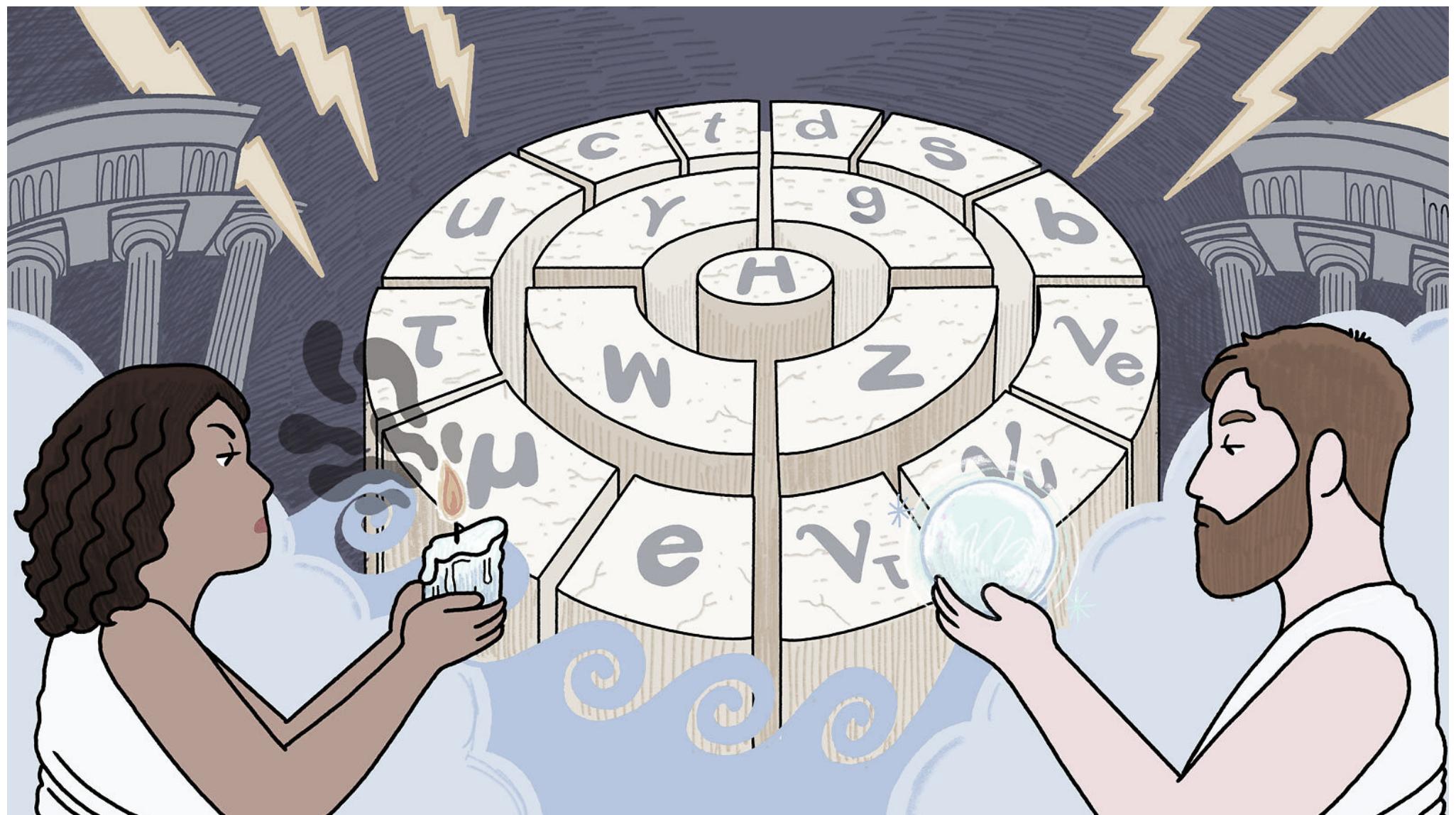


Unification of electromagnetic and weak forces (modelled with the groups $U(1) \times SU(2)$) and the strong force (based on the group $SU(3)$)



1954

Yang, Mills 1954



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



P. Anderson

Anderson 1972

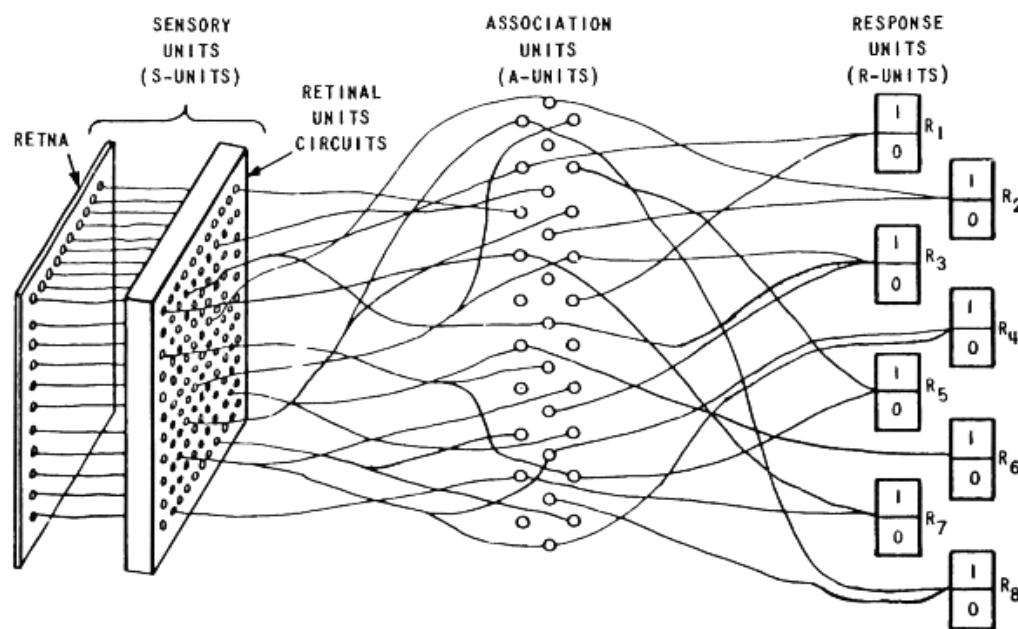
?

EARLY NEURAL NETWORKS & THE AI WINTER



Dartmouth AI Conference 1956

Early neural networks



F. Rosenblatt

1957

Perceptron, one of the first neural network architectures

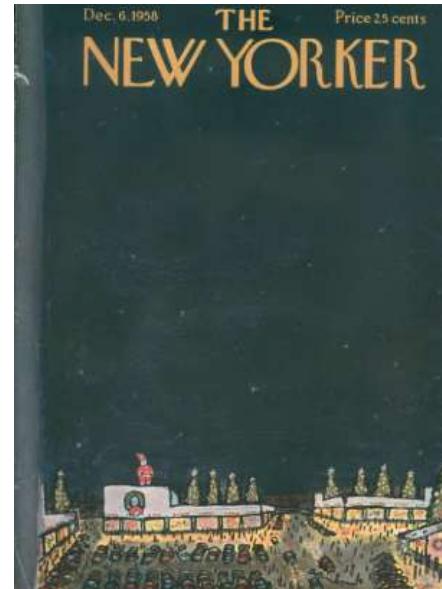
Rosenblatt 1957

Early hype

“First serious rival to the
human brain even devised.”

“Remarkable machine
capable of what amounts to
thought”

— The New Yorker



Manson, Stewart, Gill 1958

Early hype

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
PROJECT MAC

Artificial Intelligence Group
Vision Memo. No. 100.

July 7, 1966

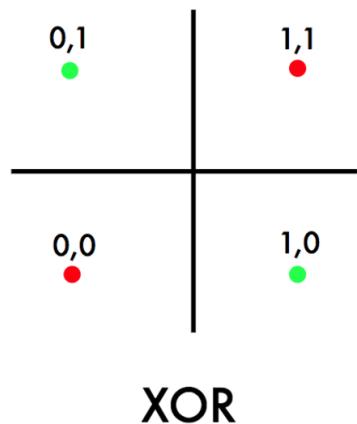
THE SUMMER VISION PROJECT

Seymour Papert

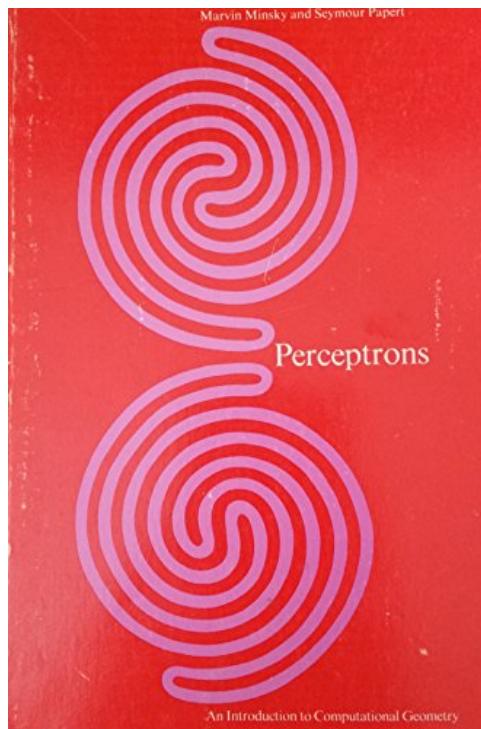
The summer vision project is an attempt to use our summer workers effectively in the construction of a significant part of a visual system. The particular task was chosen partly because it can be segmented into sub-problems which will allow individuals to work independently and yet participate in the construction of a system complex enough to be a real landmark in the development of "pattern recognition".

Papert 1966

The “XOR Affair”



“[simple] perceptron
cannot represent even
the XOR function”



M. Minsky

S. Papert

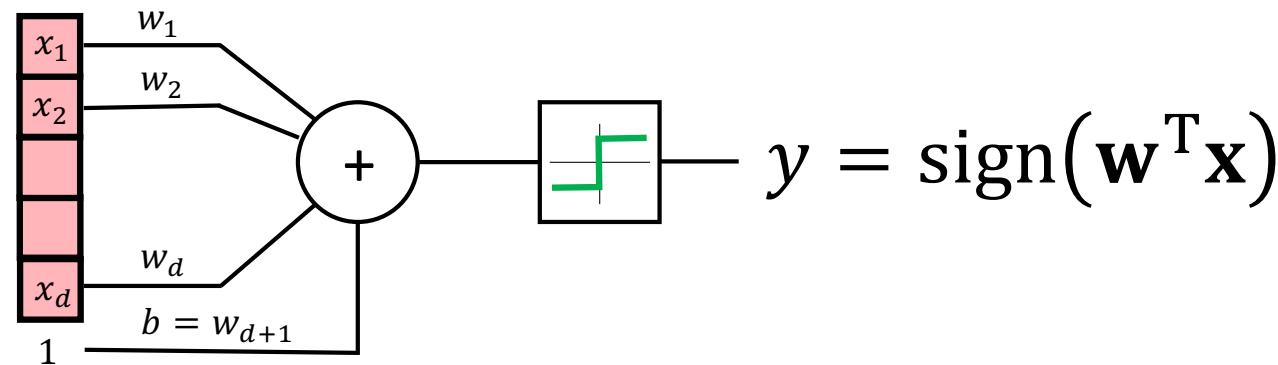
1969

Minsky, Papert 1969



“AI WINTER”

“Simple perceptron”



First “geometric” machine learning

Group Invariance Theorem: “if a neural network is invariant to a group, then its output can be expressed as functions of the orbits of the group”



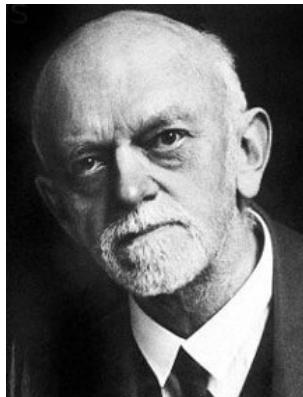
M. Minsky

S. Papert

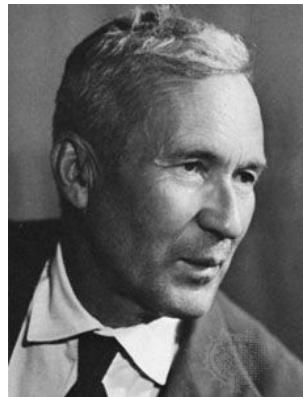
1969

Minsky, Papert 1969

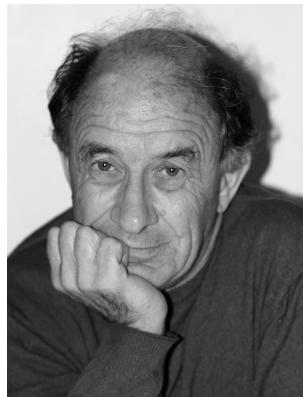
Universal approximation



D. Hilbert



A. Kolmogorov



V. Arnold



G. Cybenko



K. Hornik

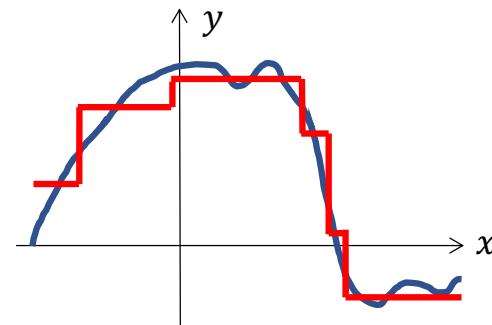
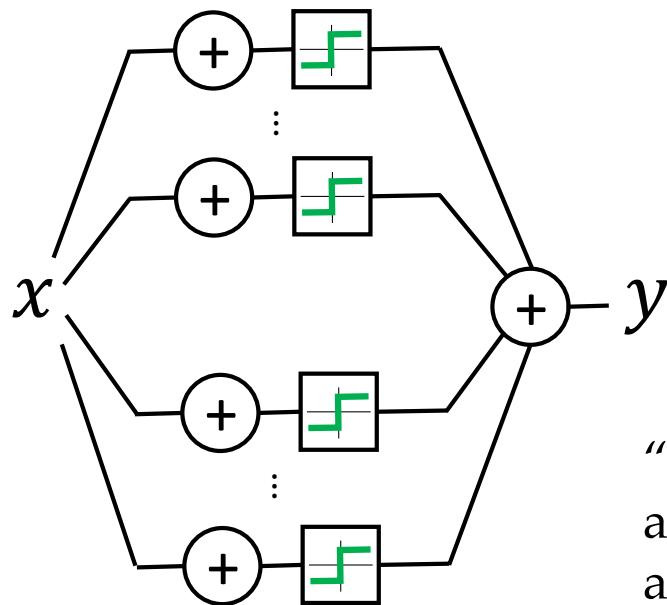
13th Problem

$$f(x_1, \dots, x_n) = \sum_{q=0}^{2n} \Phi_q \left(\sum_{p=1}^n \phi_{q,p}(x_p) \right)$$

Results specific to multilayer
neural networks

Hilbert 1900; Arnold 1956; Kolmogorov 1957; Cybenko 1989; Hornik 1991

Universal approximation



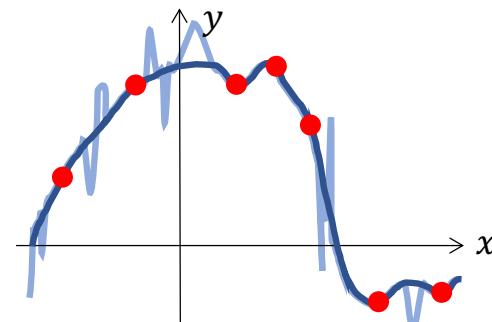
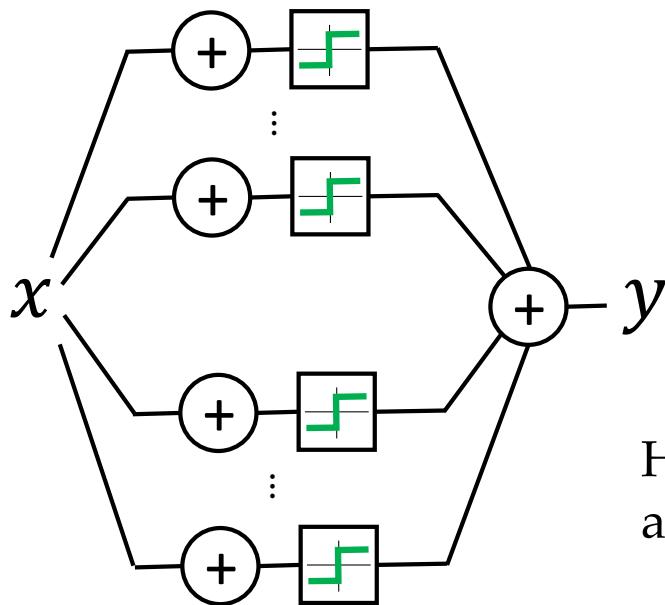
“A 2-layer perceptron can approximate a continuous function to any desired accuracy”

Cybenko 1989; Hornik 1991; Barron 1993; Leshno et al 1993; Maiorov 1999; Pinkus 1999

Deep learning = glorified curve fitting



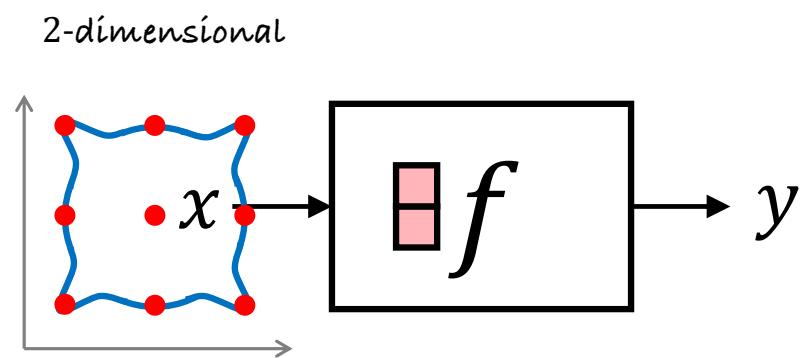
Universal approximation



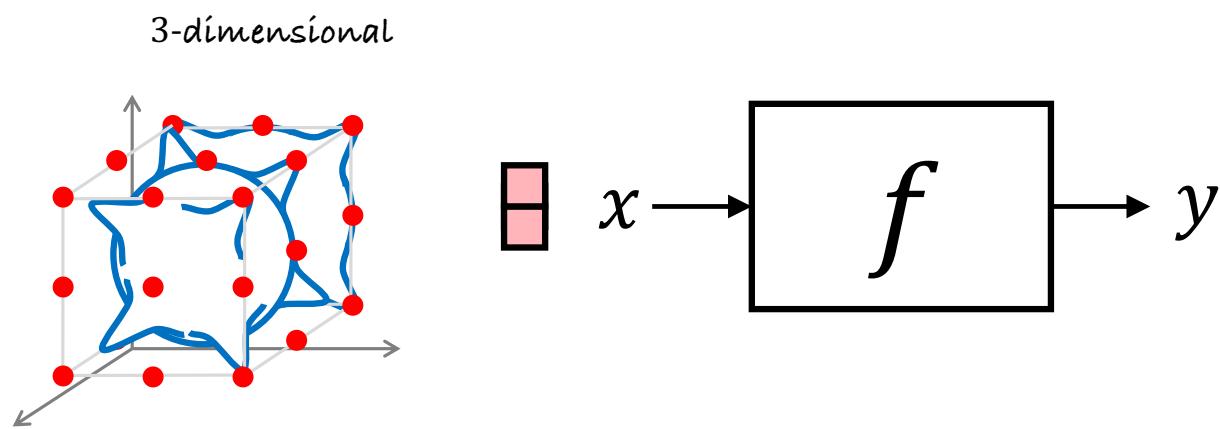
How many samples are needed to approximate to accuracy ε ?

Cybenko 1989; Hornik 1991; Barron 1993; Leshno et al 1993; Maiorov 1999; Pinkus 1999

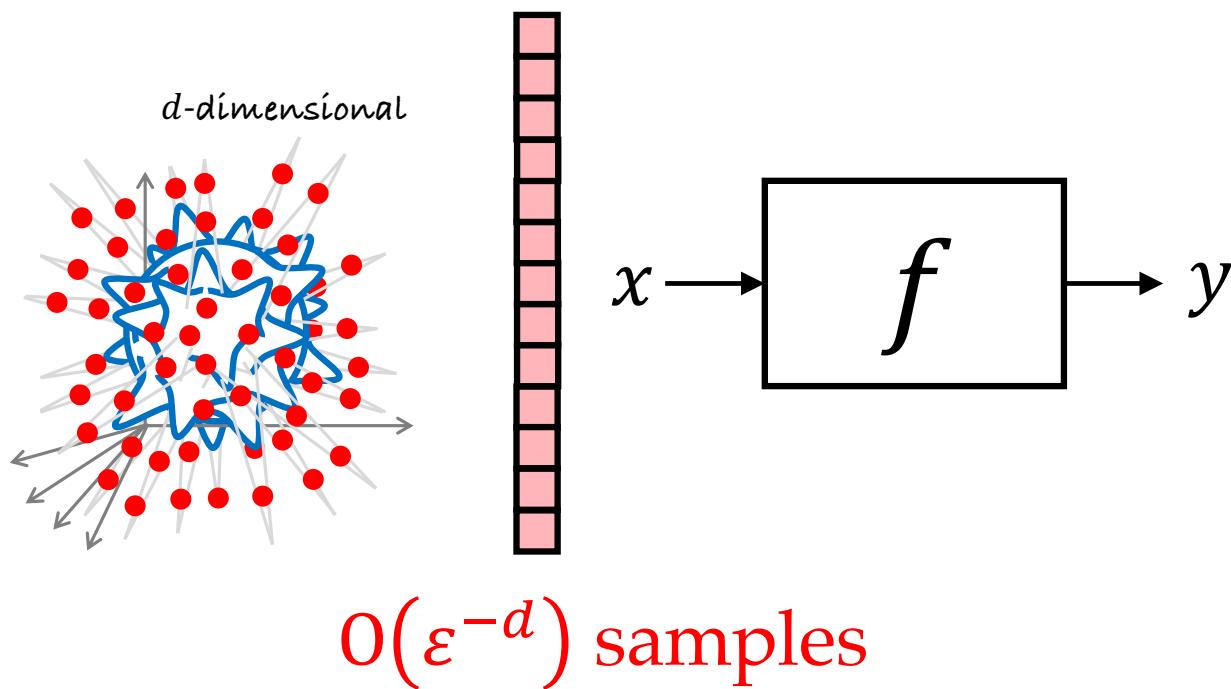
The Curse of Dimensionality



The Curse of Dimensionality



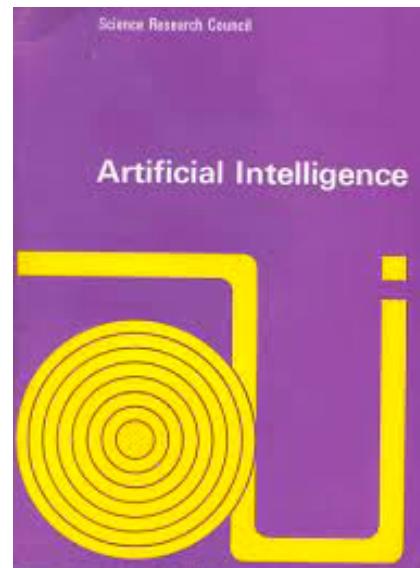
The Curse of Dimensionality





The Lighthill Report

“Most workers in AI research and in related fields confess to a pronounced feeling of disappointment in what has been achieved in the past twenty-five years. [...] In no part of the field have the discoveries made so far produced the major impact that was then promised.”

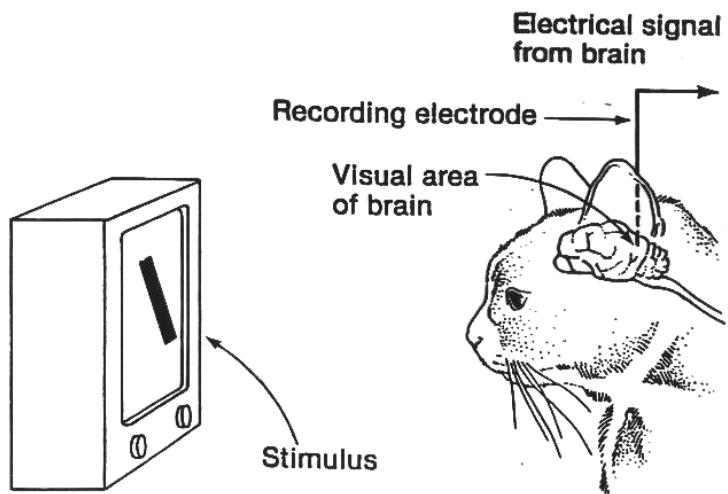


J. Lighthill

1972

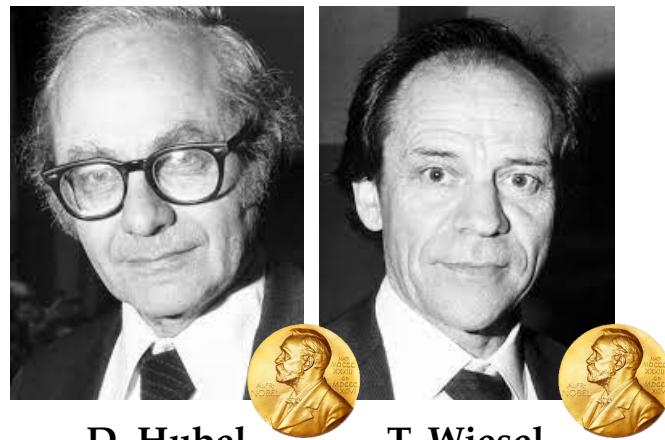
THE EMERGENCE OF GEOMETRIC ARCHITECTURES

Secrets of the visual cortex



Experiments of Hubel and Wiesel that established the structure of the visual cortex

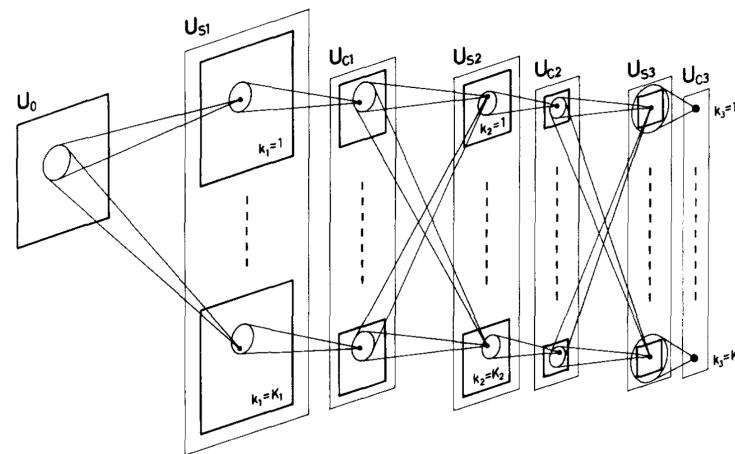
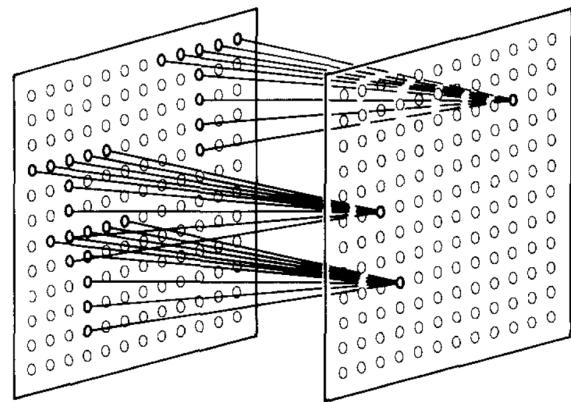
Hubel, Wiesel 1959; 1962



D. Hubel T. Wiesel

1959

Neocognitron

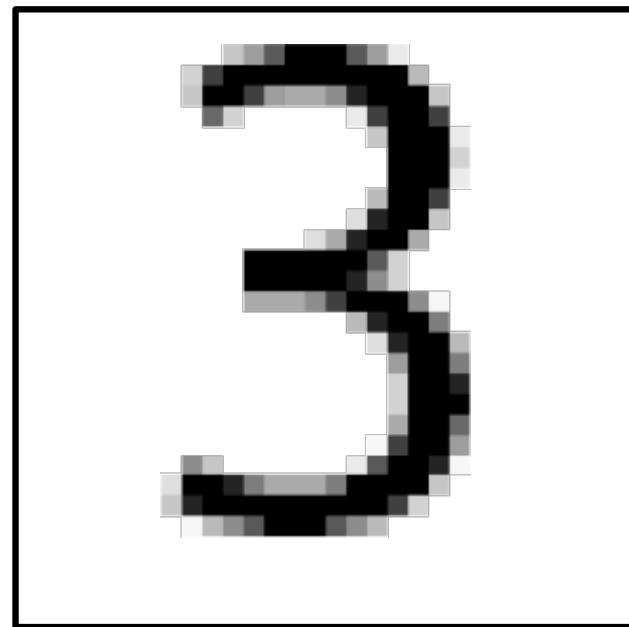


Neocognitron, an early geometric neural network

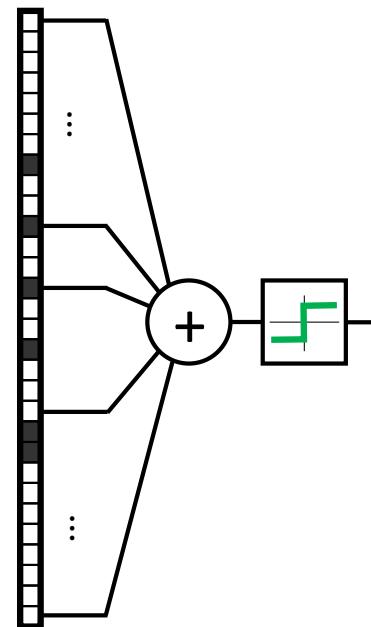


K. Fukushima

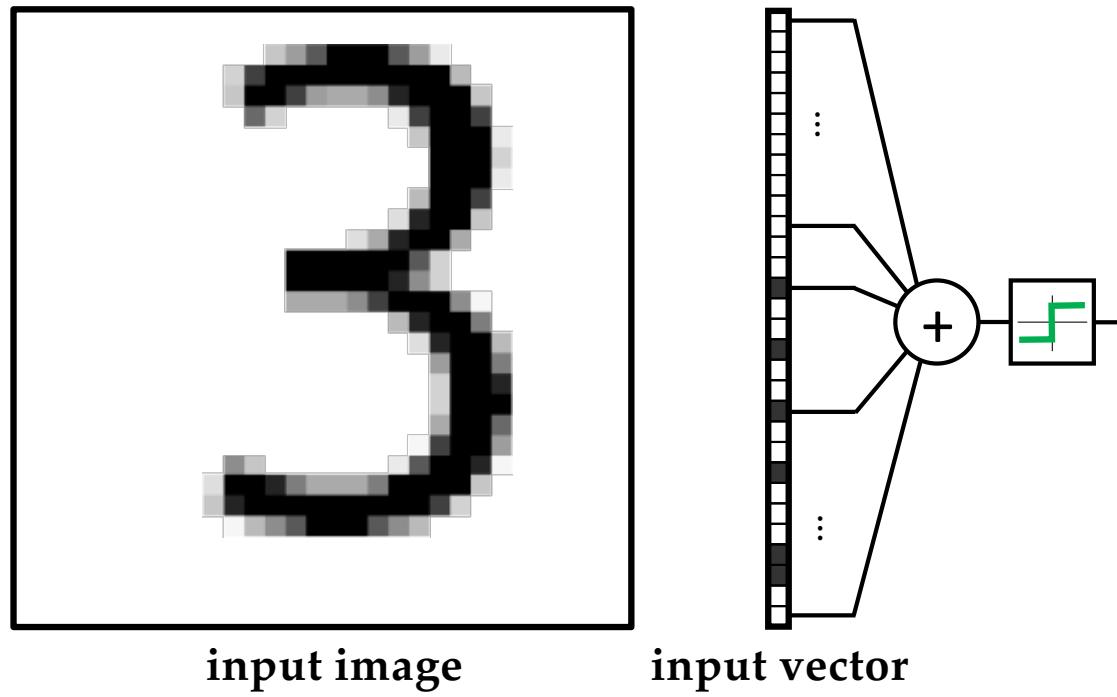
1980



input image

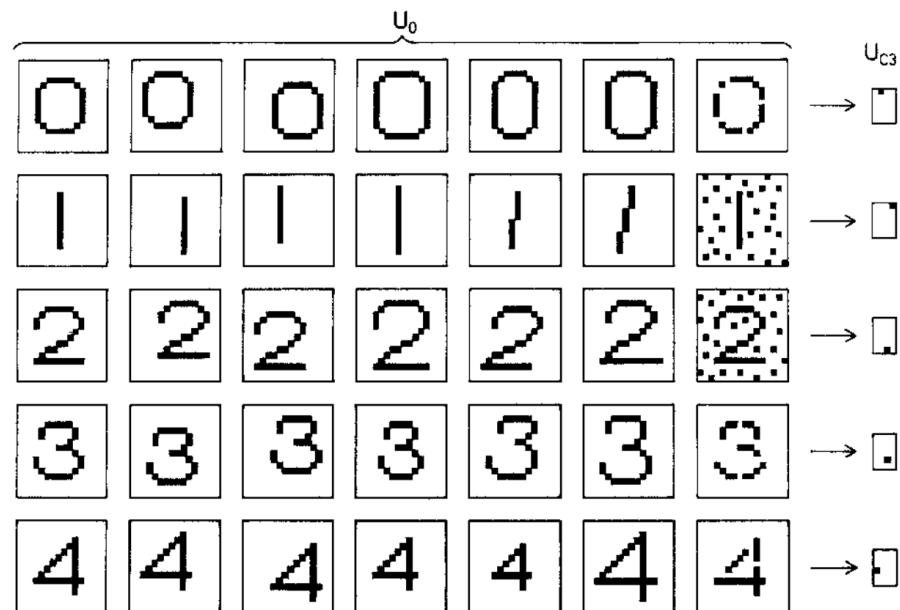


input vector



“The response of [Perceptrons] was severely affected by the shift in position [...] of the input patterns. Hence, their ability for pattern recognition was not so high.” — Fukushima

Neocognitron



Experimental evaluation of the Neocognitron



K. Fukushima

1980

Neocognitron

- Deep neural network (7 layers tested)
- Local connectivity (“receptive fields”)
- Nonlinear filters with shared weights (S-layers)
- Average pooling (C-layers)
- ReLU activation function
- “Self-organised” (unsupervised) – **no backprop yet!**



K. Fukushima

1980

How to train your neural network?



F. Rosenblatt

Perceptron
learning rule
(1 layer)



A. Ivakhnenko

Group method of
data handling



S. Linnainmaa



P. Werbos

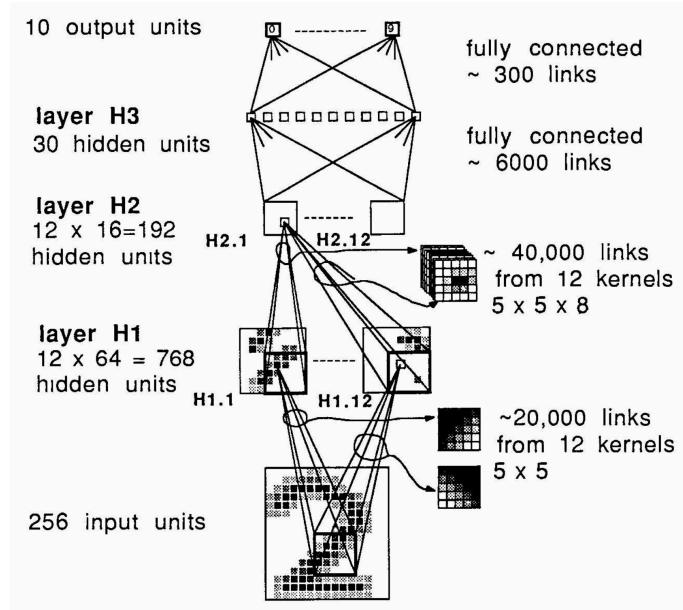
Backpropagation



D. Rumelhart

Rosenblatt 1957; Ivakhnenko, Lapa 1966; Linnainmaa 1970; Werbos 1982; Rumelhart et al. 1986

Convolutional neural networks



First version of a CNN

LeCun et al. 1989

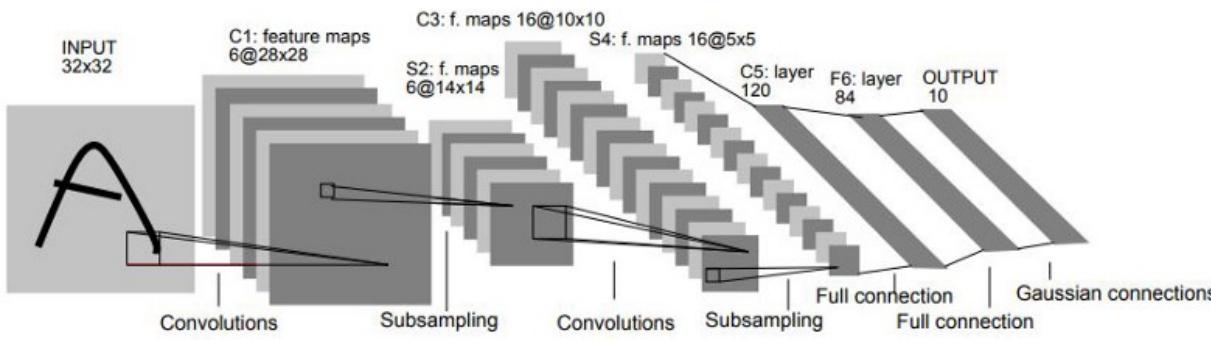


AT&T DSP-32C
capable of 125m floating
point multiply-accumulate
operations/sec

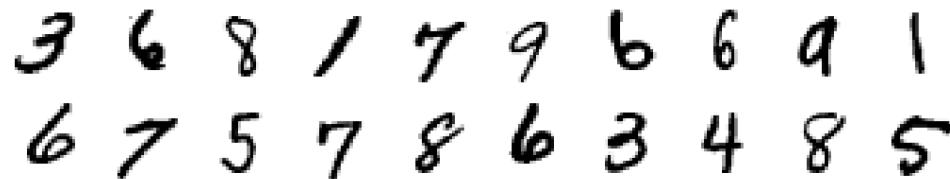


Y. LeCun

LeNet-5



LeNet-5 classical CNN architecture



MNIST digits dataset



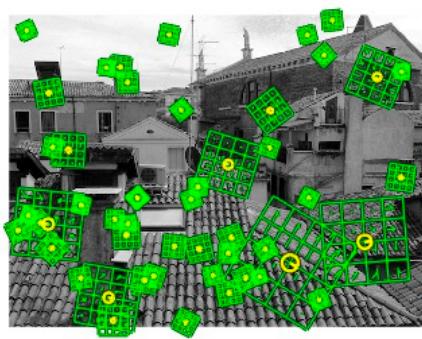
Y. LeCun

LeCun et al. 1998

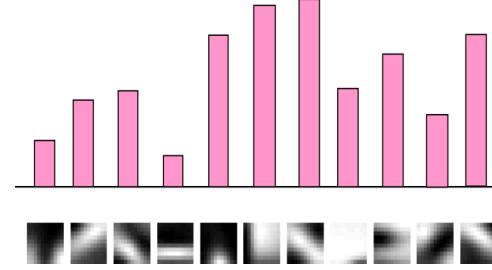
Computer vision in the 2000s



Feature detection



Feature description



Feature aggregation

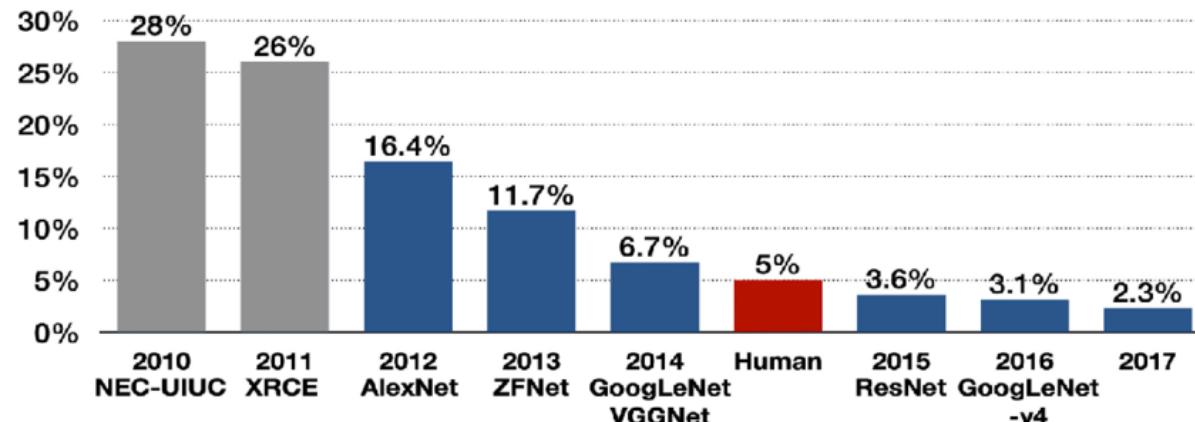


Classification

A typical image classification pipeline from the 2000s

ImageNet

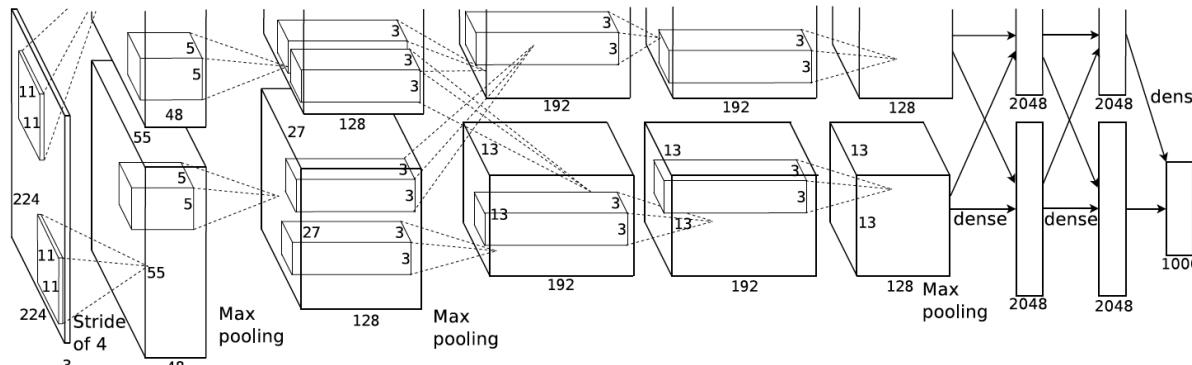
Top-5 error



L. Fei-Fei

AlexNet beating all “handcrafted” approaches on ImageNet benchmark—the moment of truth for computer vision

AlexNet



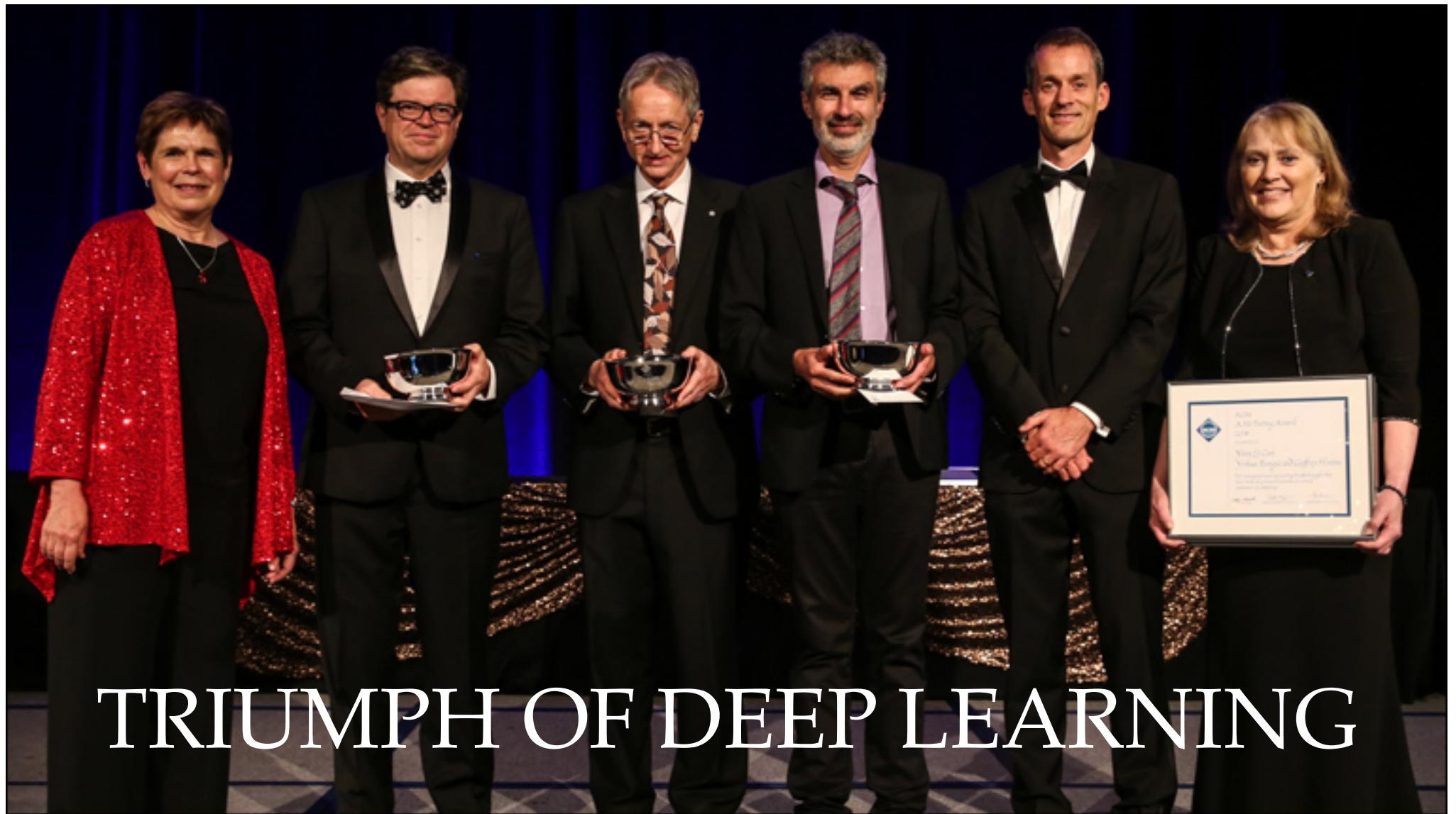
AlexNet architecture

Nvidia GTX 580 GPU capable of
~200G FLOP / sec



A. Krizhevsky

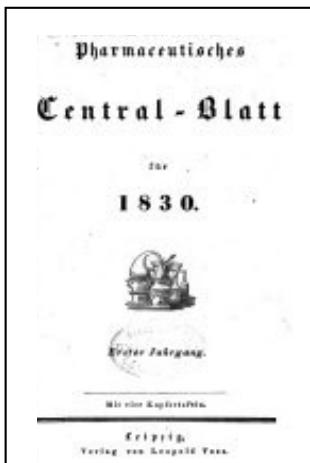
Krizhevsky et al. 2012



TRIUMPH OF DEEP LEARNING

GRAPH NEURAL NETWORKS & THEIR CHEMICAL PRECURSORS

Early chemoinformatics



First chemical abstracts journal
Chemisches Zentralblatt 1830–1969

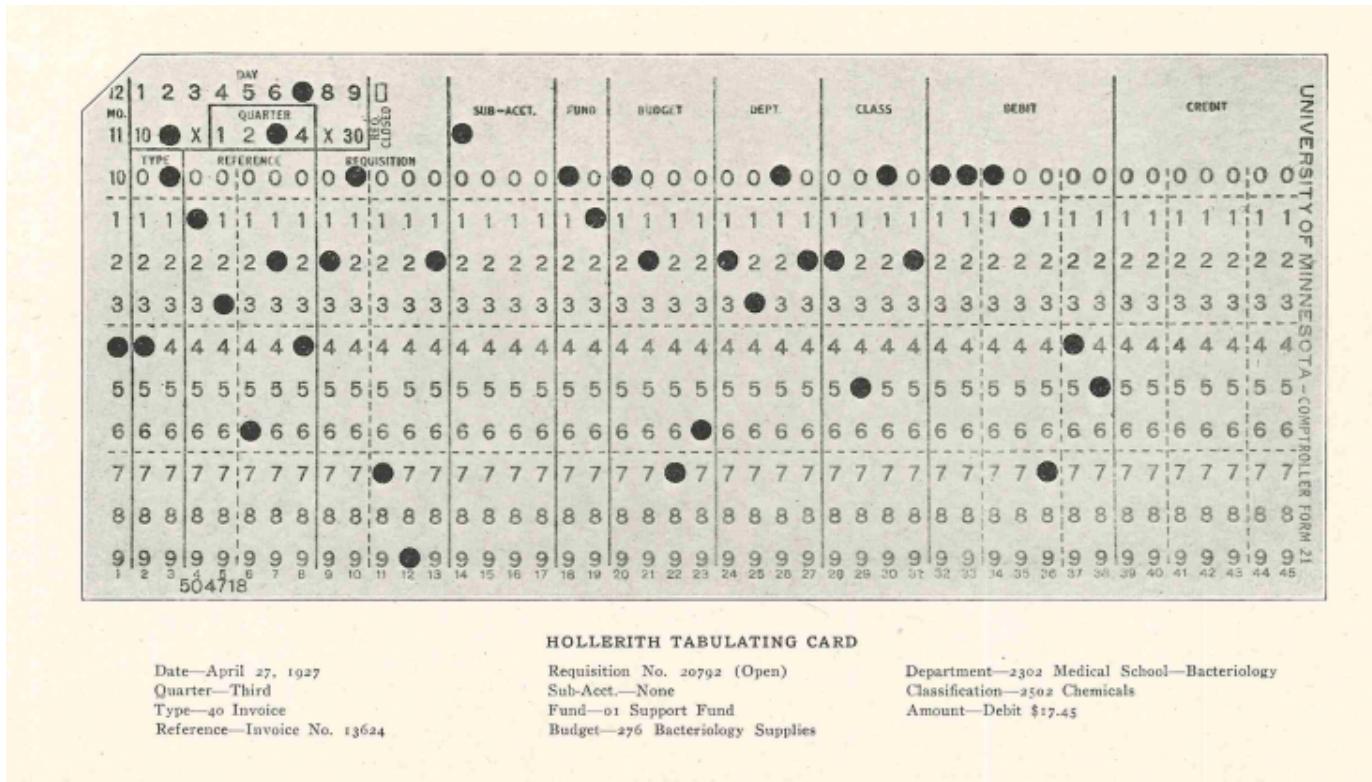


Beilstein Handbuch
~500 volumes, 500k pages



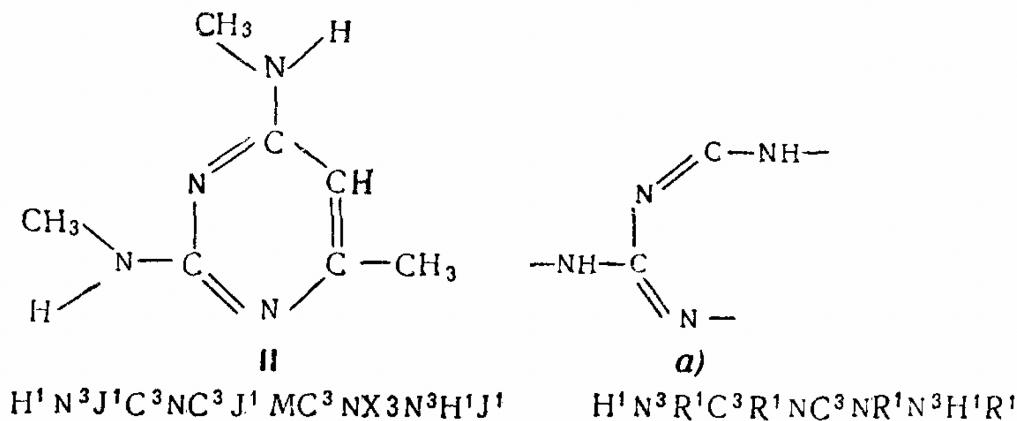
Chemical Abstracts Service
as of today ~200m compounds

Early chemoinformatics



Punch card for early computer

Structural similarity of molecules



Early “chemical ciphers” used for molecule representations
fail to capture structural similarity



G. Vlăduț

1959

Graph theory & Chemistry

CHEMISTRY AND ALGEBRA

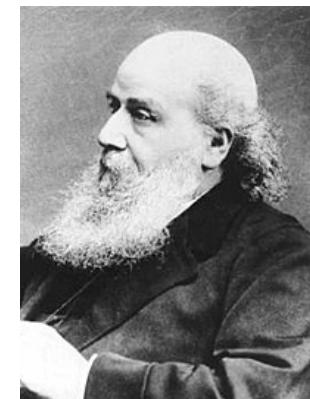
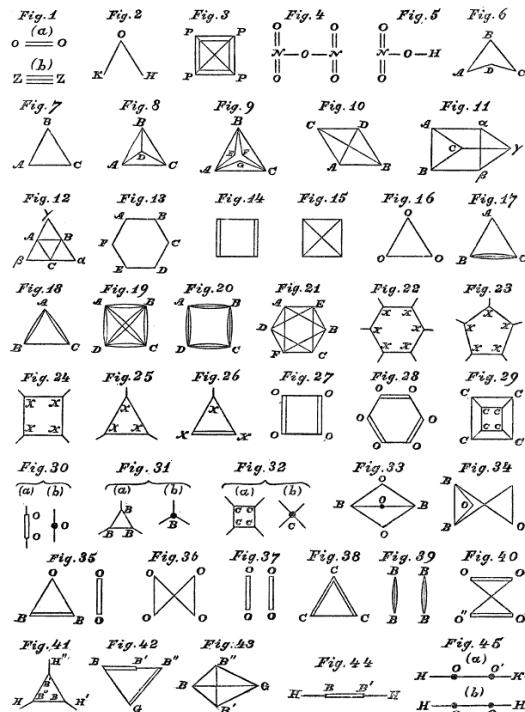
IT may not be wholly without interest to some of the readers of NATURE to be made acquainted with an analogy that has recently forcibly impressed me between branches of human knowledge apparently so dissimilar as modern chemistry and modern algebra.

The weight of an invariant is identical with the number of the bonds in the chemicograph of the analogous chemical substance, and the weight of the leading term (or basic differentiant) of a co-variant is the same as the number of bonds in the chemicograph of the analogous compound radical. Every invariant and covariant thus becomes expressible by a *graph* precisely identical with a Kekuléan diagram or chemicograph.

Baltimore, January 1

J. J. SYLVESTER

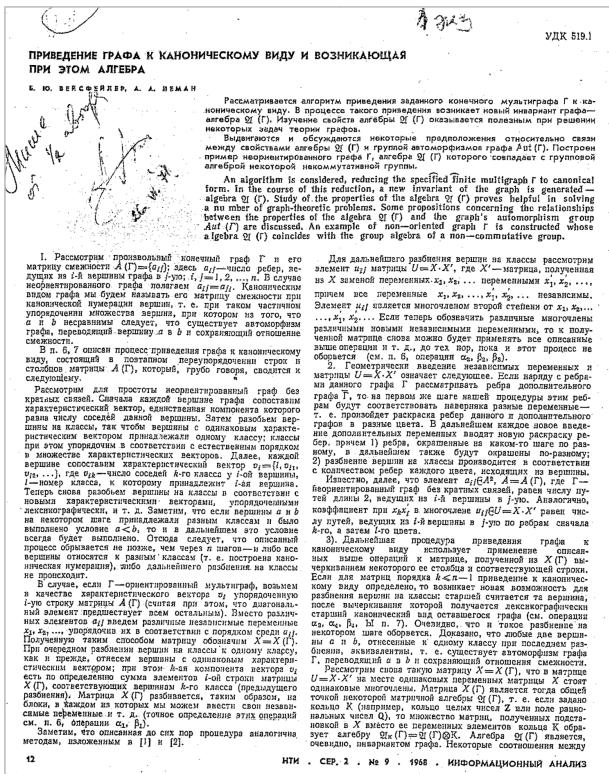
The term “graph” appeared first in the chemical context



J. Sylvester

1878

Weisfeiler-Lehman test



A. Lehman



B. Weisfeiler

1968

Weisfeiler, Lehman 1968; Portraits: Ihor Gorskiy

First Graph Neural Networks



A. Sperduti

Labeling RAAM

1994



C. Goller

Backprop through structure

1996



A. Küchler



M. Gori

"Graph Neural Networks"

2005, 2008



F. Scarselli

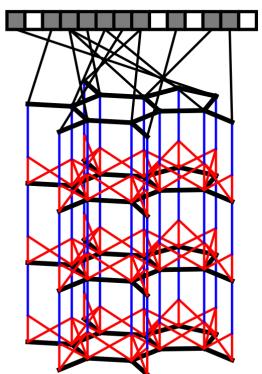


Y. Li

Gated GNN

2015

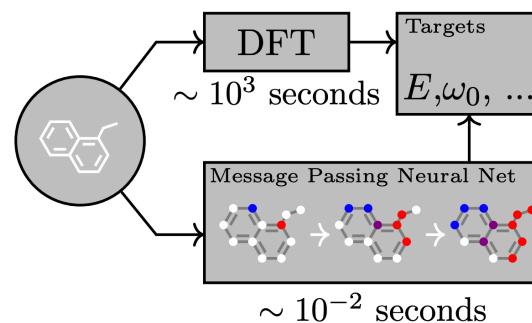
Back to the chemical roots



GNN-based
molecular fingerprints



D. Duvenaud



Chemical property prediction
using message passing GNNs



J. Gilmer

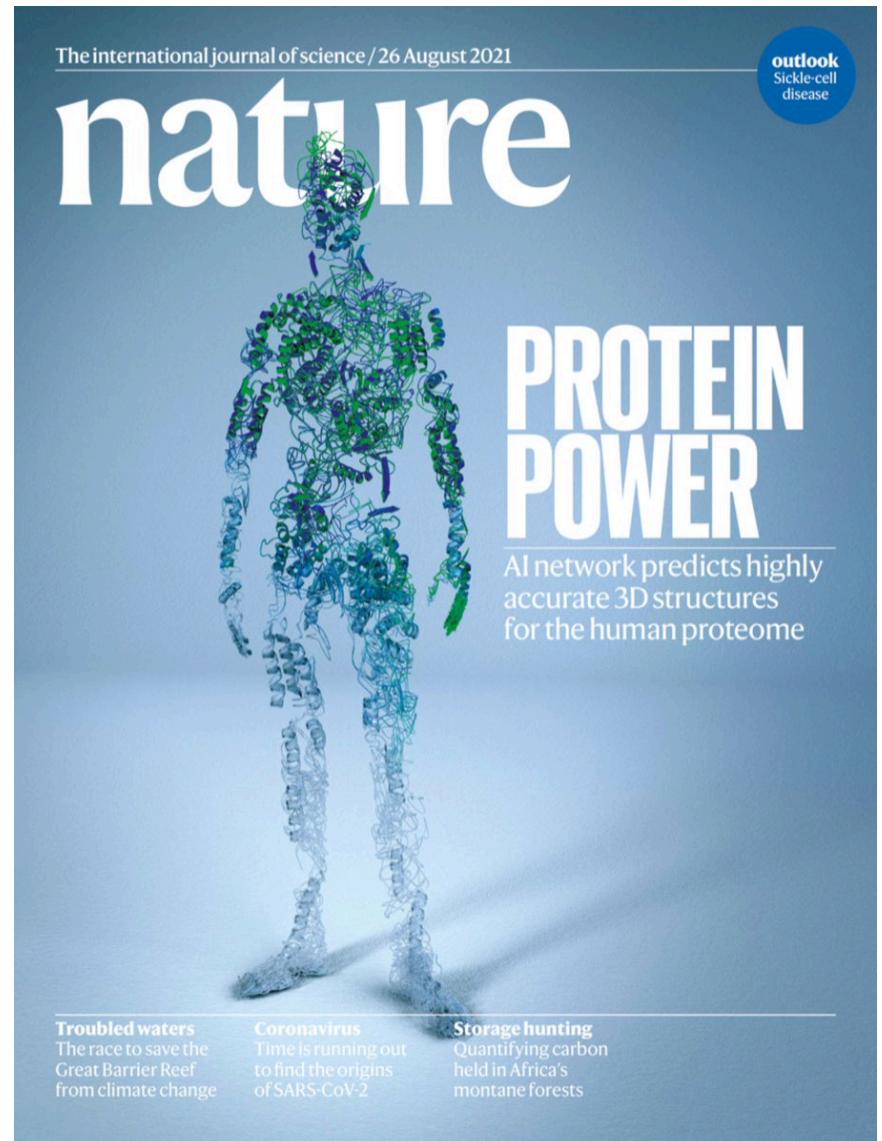
Duvenaud et al. 2015; Gilmer et al. 2017

Back to the chemical roots



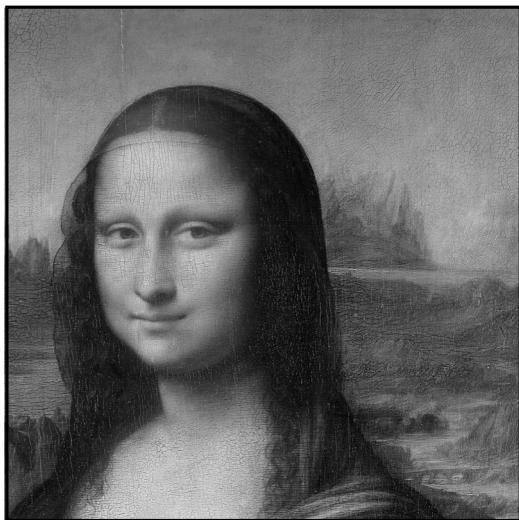
An “ImageNet” moment of structural biology

Jumper et al. 2021

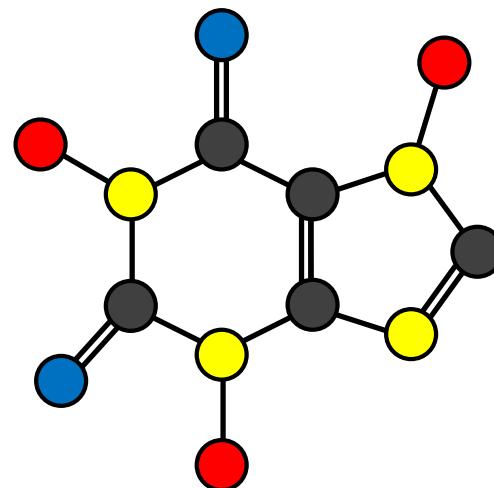


THE BLUEPRINT

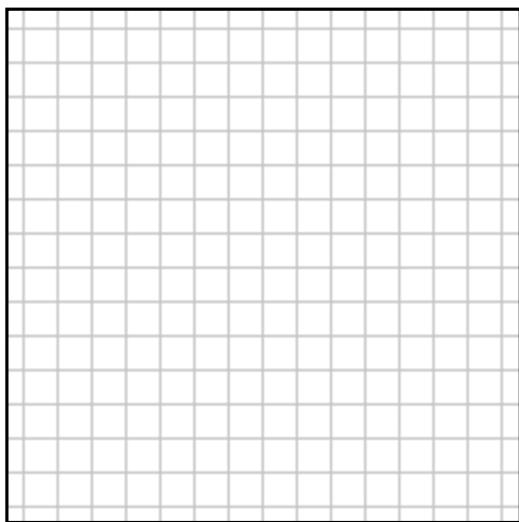
Convolutional Neural Network



Graph Neural Network

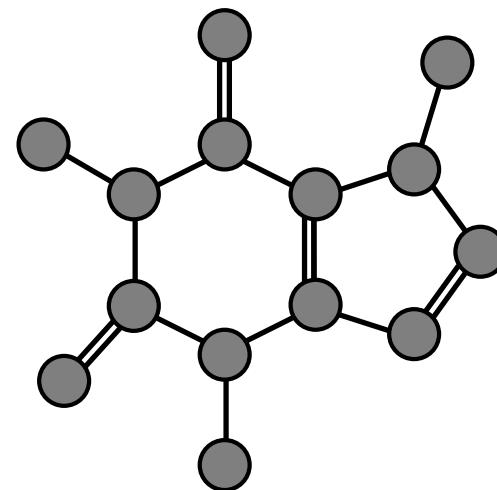


Convolutional Neural Network



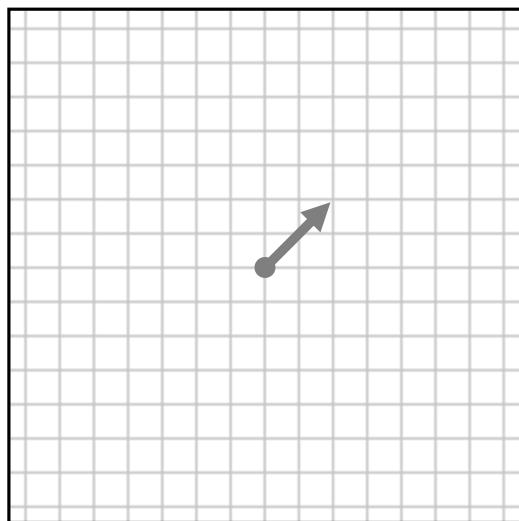
Underlying domain:
grid

Graph Neural Network



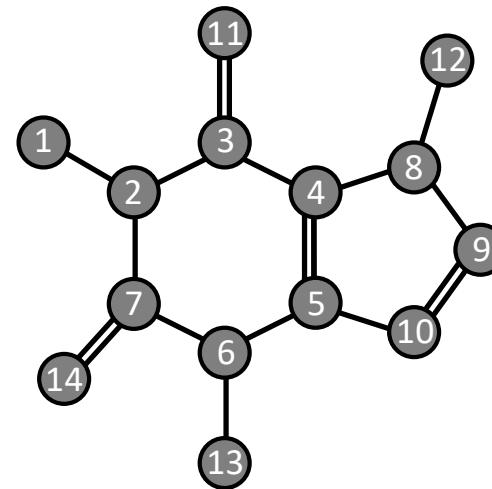
Underlying domain:
graph

Convolutional Neural Network



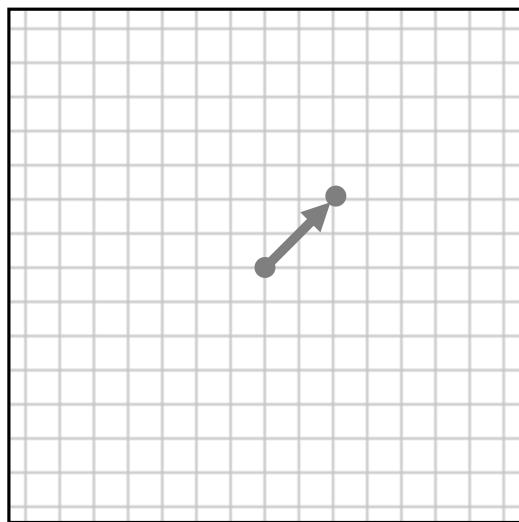
Symmetry:
Translation

Graph Neural Network



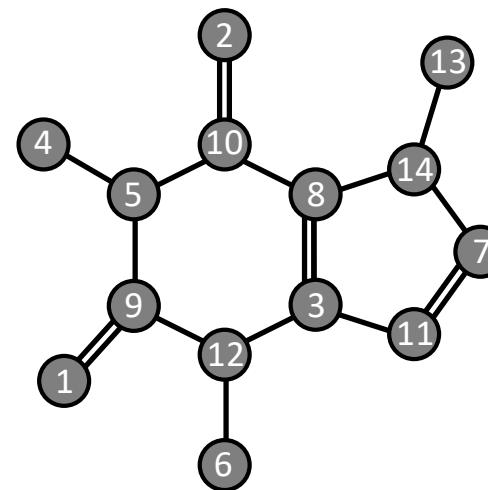
Symmetry:
Permutation

Convolutional Neural Network



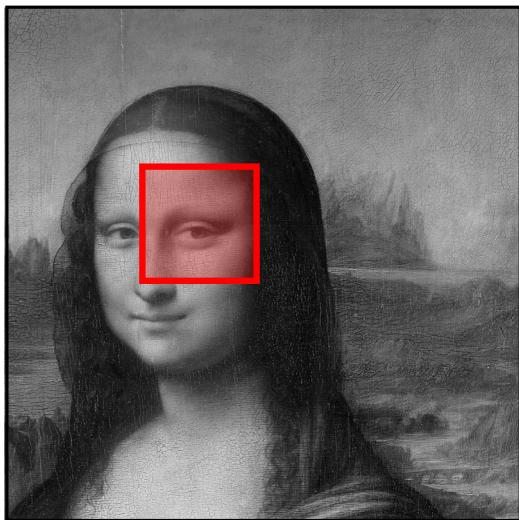
Symmetry:
Translation

Graph Neural Network



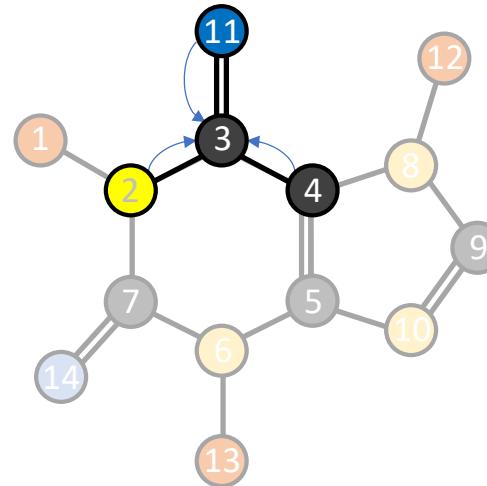
Symmetry:
Permutation

Convolutional Neural Network



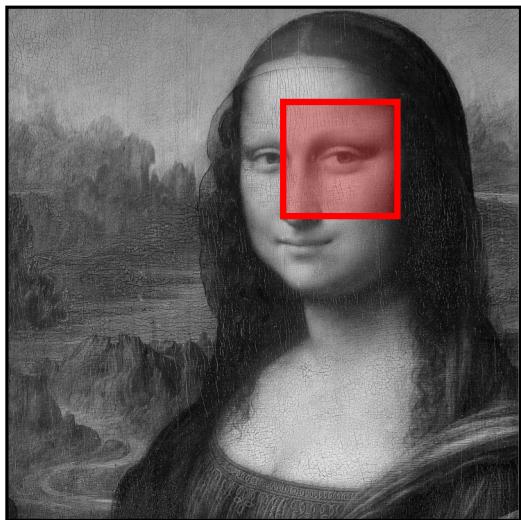
Convolution:
translation equivariant

Graph Neural Network



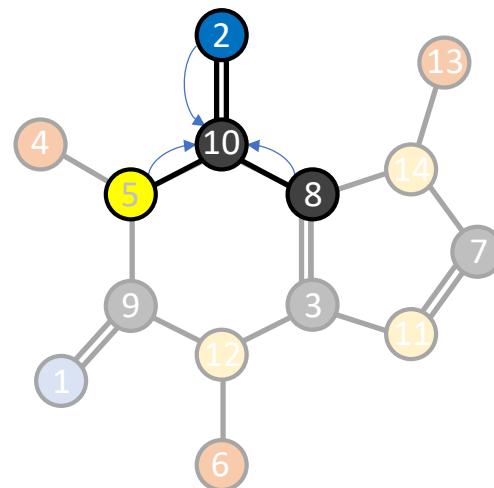
Message passing:
permutation equivariant

Convolutional Neural Network



Convolution:
translation equivariant

Graph Neural Network



Message passing:
permutation equivariant

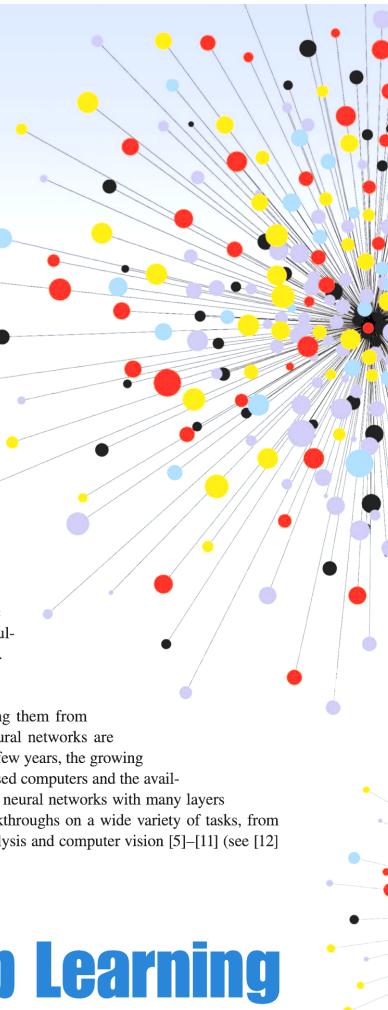
Michael M. Bronstein, Joan Bruna, Yann LeCun,
Arthur Szlam, and Pierre Vandergheynst

Many scientific fields study data with an underlying structure that is non-Euclidean. Some examples include social networks in computational social sciences, sensor networks in communications, functional networks in brain imaging, regulatory networks in genetics, and meshed surfaces in computer graphics. In many applications, such geometric data are large and complex (in the case of social networks, on the scale of billions) and are natural targets for machine-learning techniques. In particular, we would like to use deep neural networks, which have recently proven to be powerful tools for a broad range of problems from computer vision, natural-language processing, and audio analysis. However, these tools have been most successful on data with an underlying Euclidean or grid-like structure and in cases where the invariances of these structures are built into networks used to model them.

Geometric deep learning is an umbrella term for emerging techniques attempting to generalize (structured) deep neural models to non-Euclidean domains, such as graphs and manifolds. The purpose of this article is to overview different examples of geometric deep-learning problems and present available solutions, key difficulties, applications, and future research directions in this nascent field.

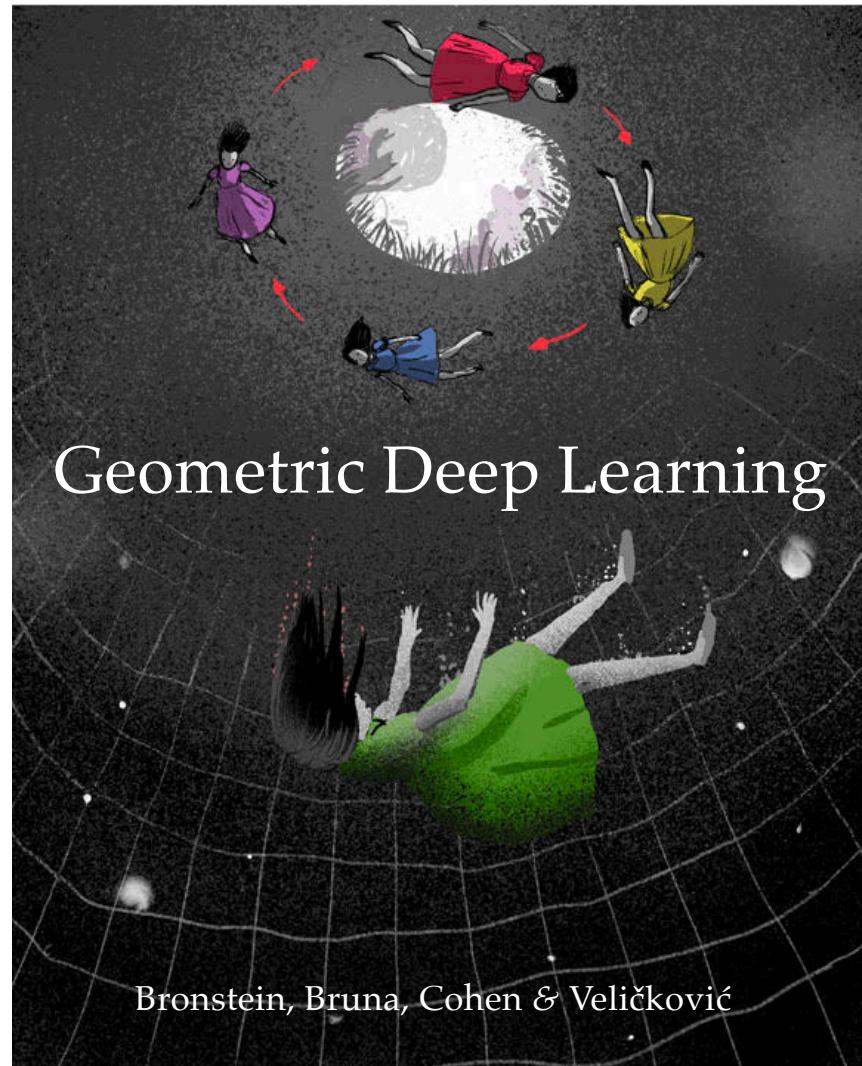
Overview of deep learning

Deep learning refers to learning complicated concepts by building them from simpler ones in a hierarchical or multilayer manner. Artificial neural networks are popular realizations of such deep multilayer hierarchies. In the past few years, the growing computational power of modern graphics processing unit (GPU)-based computers and the availability of large training data sets have allowed successfully training neural networks with many layers and degrees of freedom (DoF) [1]. This has led to qualitative breakthroughs on a wide variety of tasks, from speech recognition [2], [3] and machine translation [4] to image analysis and computer vision [5]–[11] (see [12]

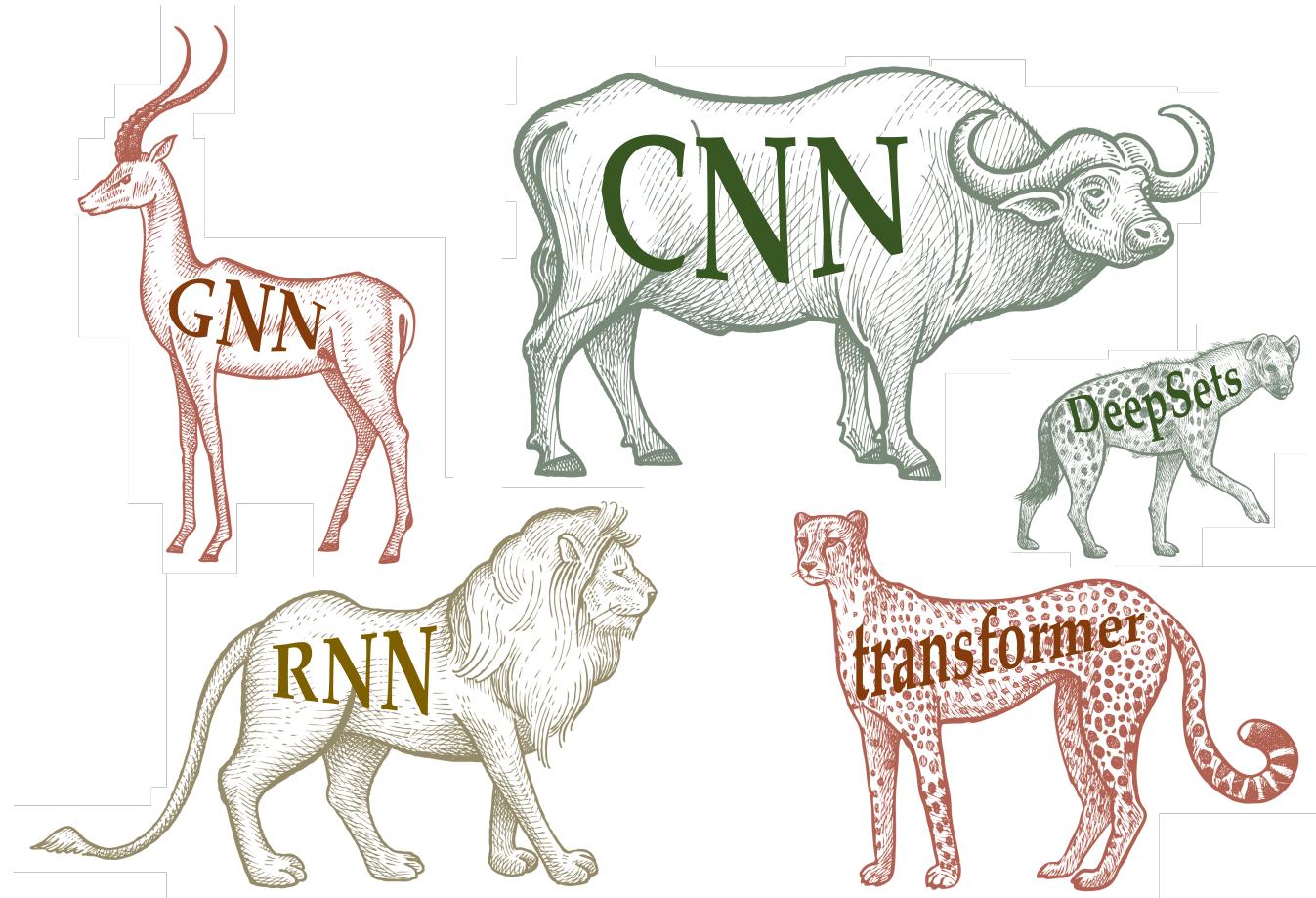


Geometric Deep Learning

Going beyond Euclidean data



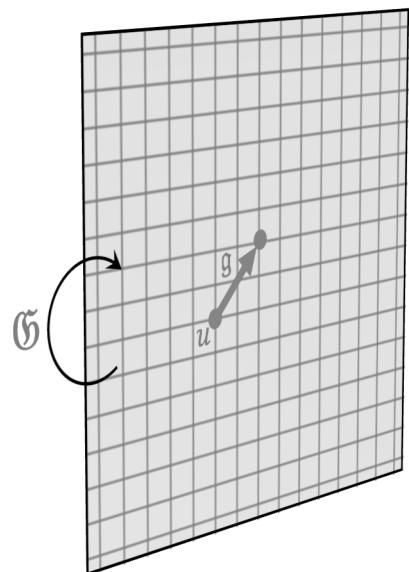
Twentieth Century Zoo of Neural Network Architectures



The Erlangen Programme of ML
Geometric Deep Learning

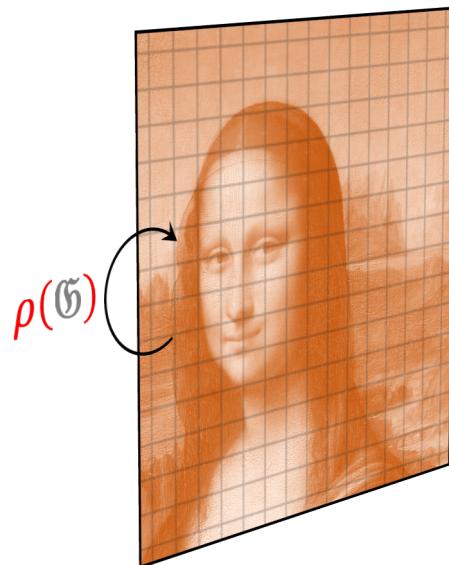
Geometric Deep Learning Blueprint

domain Ω



symmetry group \mathfrak{G}

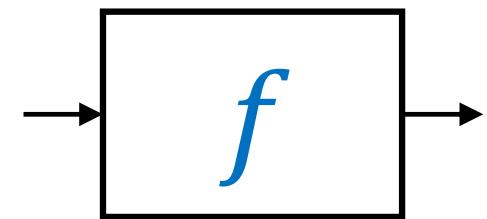
signals $\mathcal{X}(\Omega)$



group representation $\rho(\mathfrak{G})$

$$\rho(\mathfrak{g})x(u) = x(\mathfrak{g}^{-1}u)$$

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



equivariance

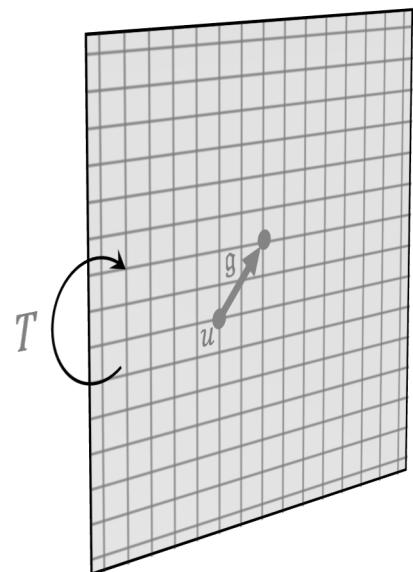
$$f(\rho(\mathfrak{g})x) = \rho(\mathfrak{g})f(x)$$

invariance

$$f(\rho(\mathfrak{g})x) = f(x)$$

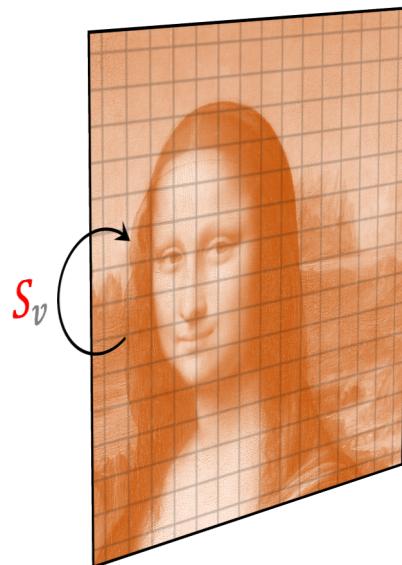
Example: Convolutional Neural Networks

Plane \mathbb{R}^2



Translation group $T(2)$

images $\mathcal{X}(\mathbb{R}^2)$



Shift operator S

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

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

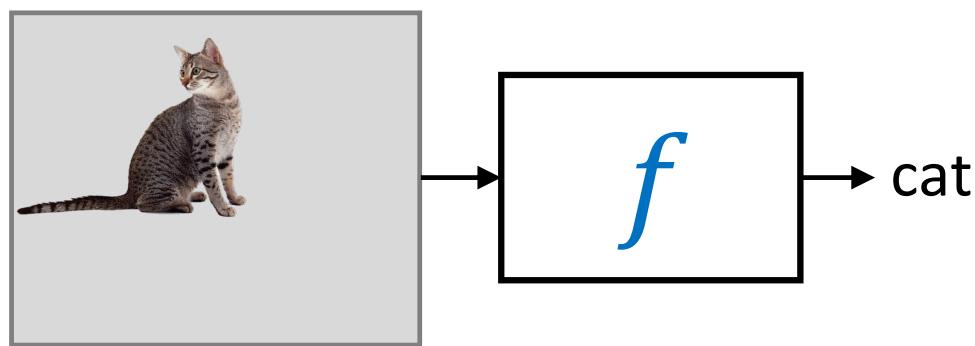


Convolutional layer

$$(Sx * y) = S(x * y)$$

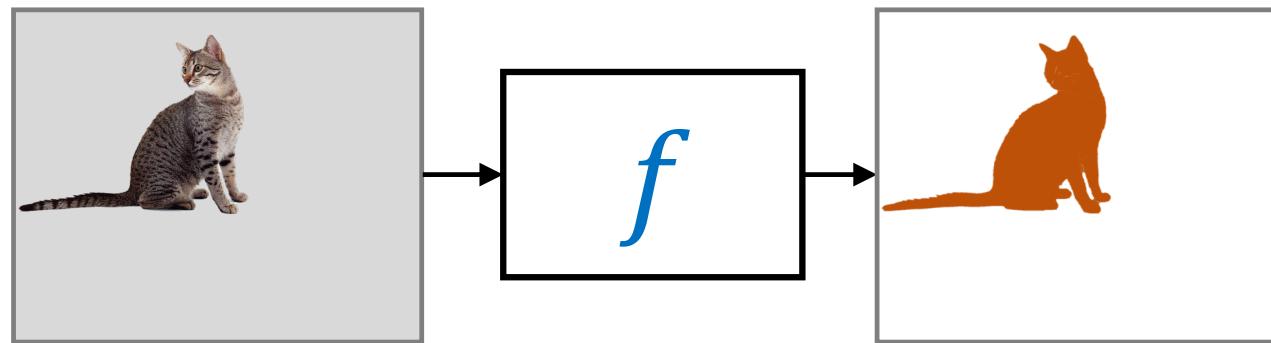
Invariant functions: Image Classification

\mathfrak{G} -invariance $f(\rho(g)x) = f(x)$



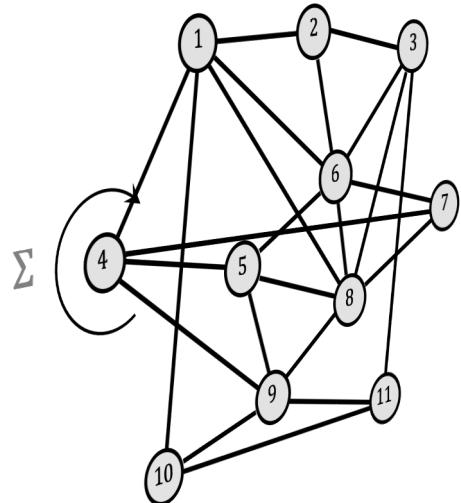
Equivariant functions: Image Segmentation

$$\mathfrak{G}\text{-equivariance } f(\rho(g)x) = \rho(g)f(x)$$



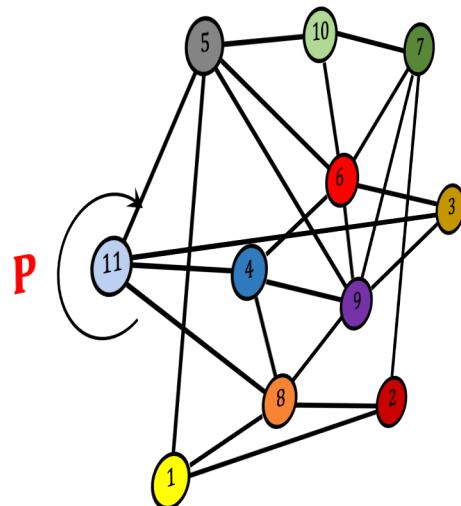
Example: Graph Neural Networks

Graph $G = (V, E)$



Permutation group Σ_n

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



Permutation matrix P

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

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

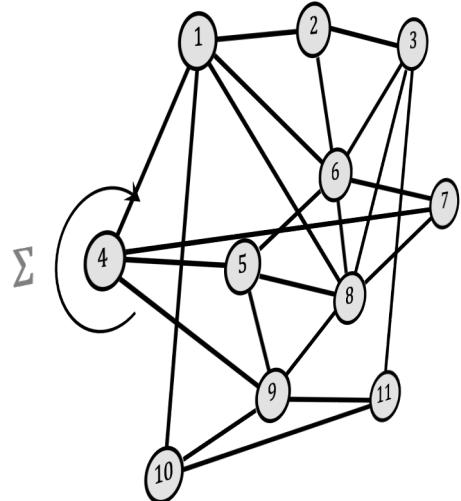


Message passing

$$\mathcal{F}(PX, PAP^T) = PF(X, A)$$

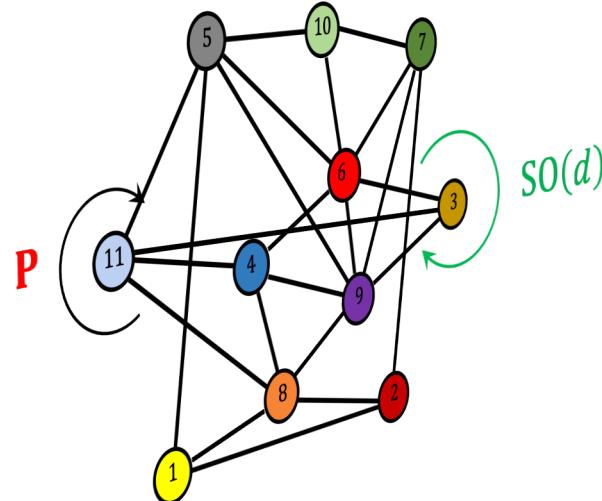
Example: Equivariant Graph Neural Networks

Graph $G = (V, E)$



Permutation group Σ_n

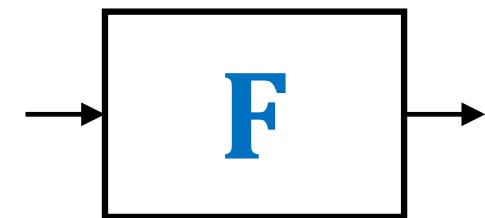
Node features $\mathcal{X}(G)$



Permutation matrix P

Rotation R

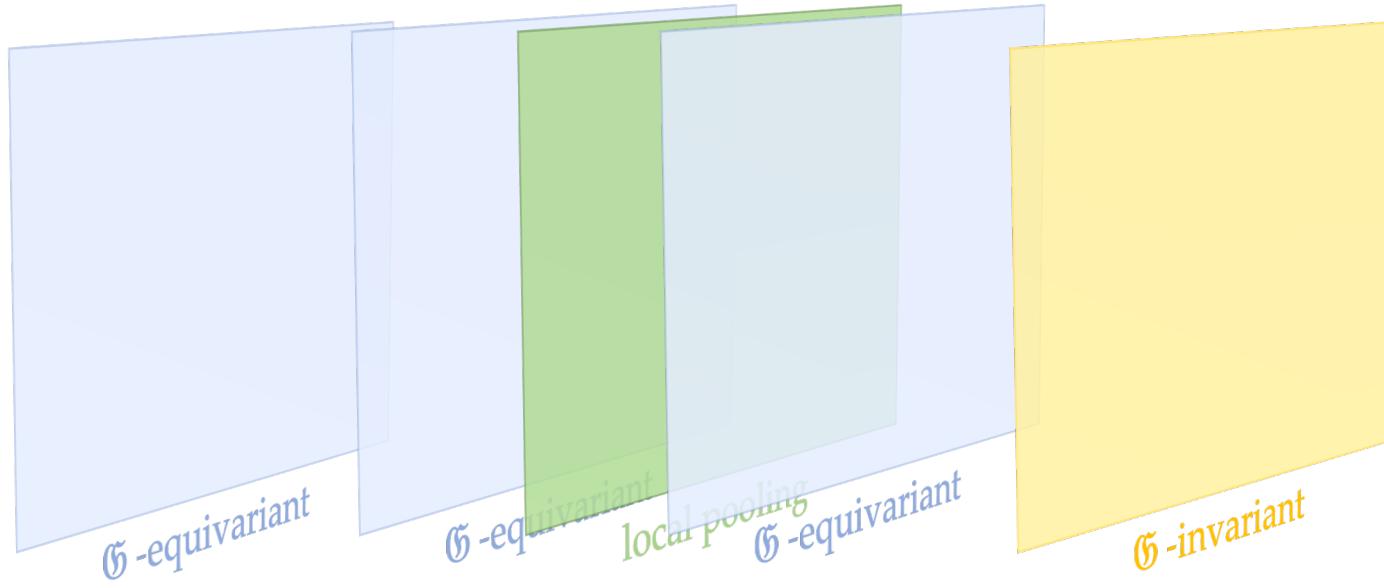
functions $\mathcal{F}(\mathcal{X}(\Omega))$



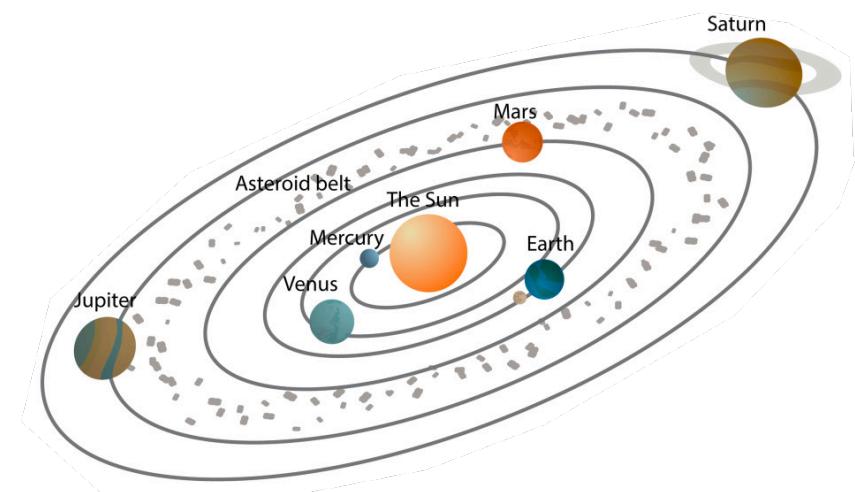
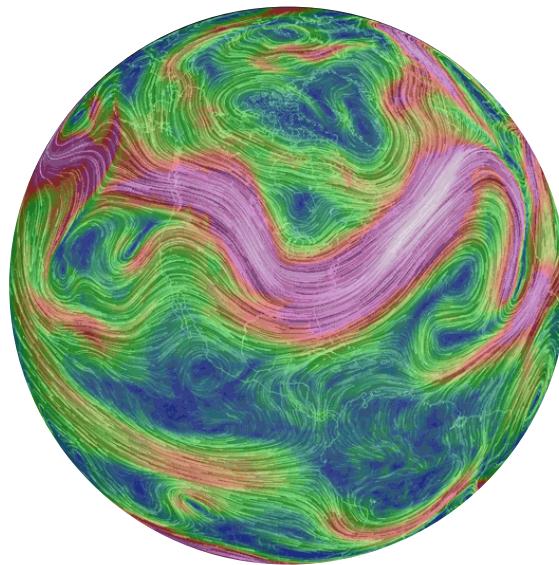
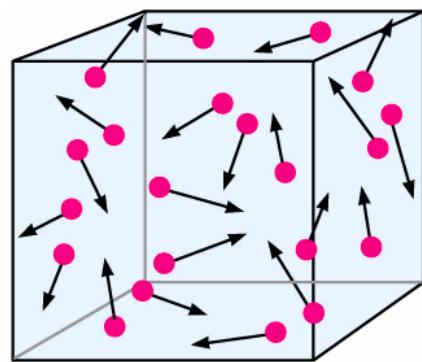
Equivariant message passing

$$\mathcal{F}(PXR, PAP^T) = PF(X, A)R$$

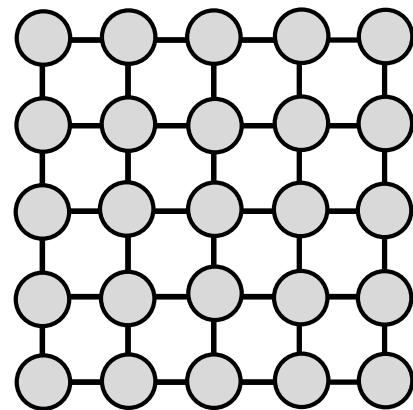
Geometric Deep Learning Blueprint



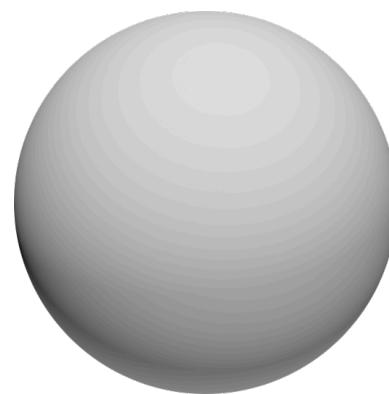
Scale Separation in Physics



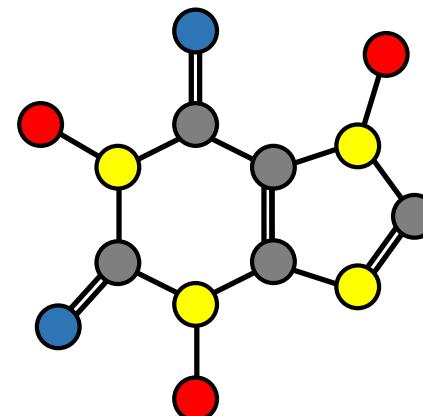
The “5G” of Geometric Deep Learning



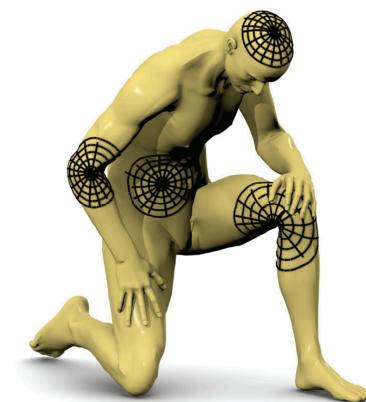
Images &
Sequences



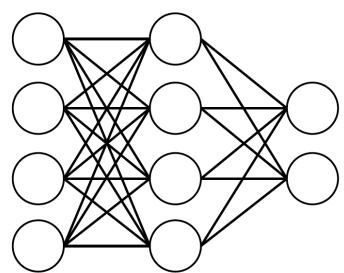
Homogeneous
spaces



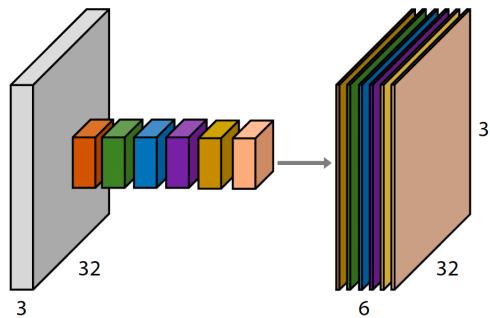
Graphs & Sets



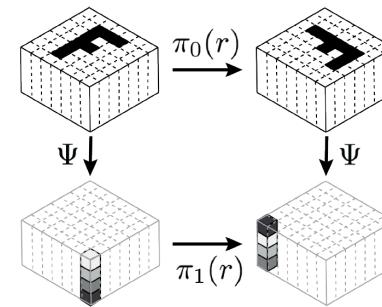
Manifolds, Meshes &
Geometric graphs



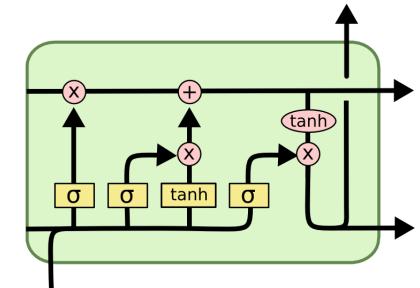
Perceptrons
Function regularity



CNNs
Translation



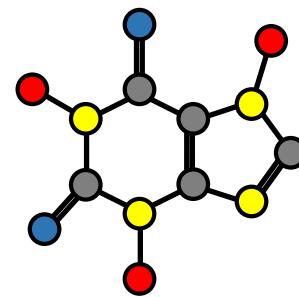
Group-CNNs
Translation+Rotation,
Global groups



LSTMs
Time warping



DeepSets / Transformers
Permutation



GNNs
Permutation

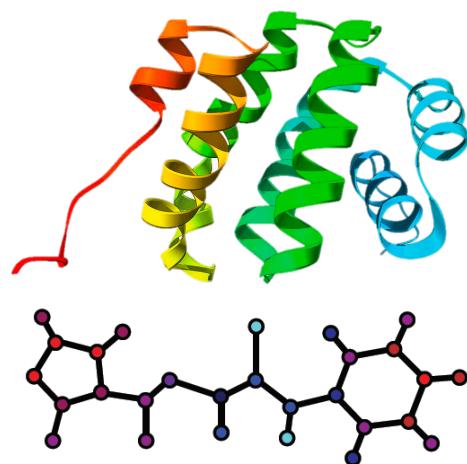


Intrinsic CNNs
Isometry / Gauge choice

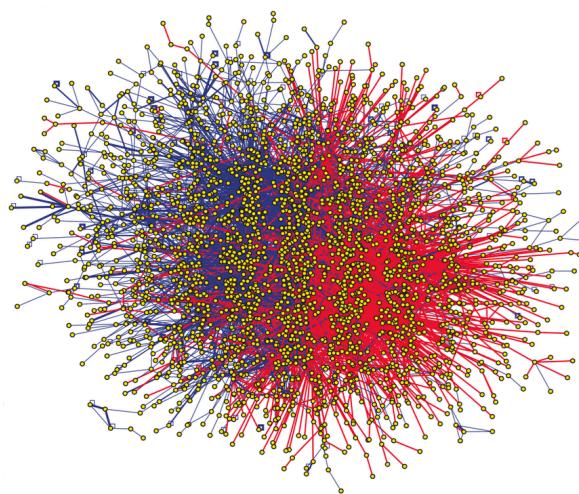
geometricdeeplearning.com

GRAPHS

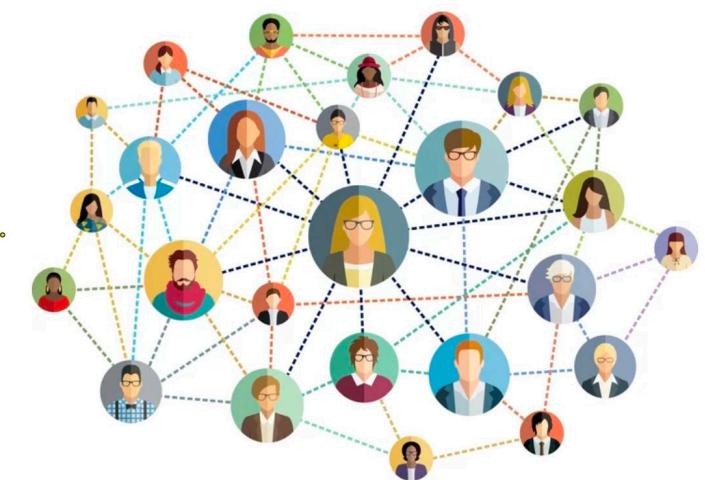
Graphs = Systems of Relations and Interactions



Molecules

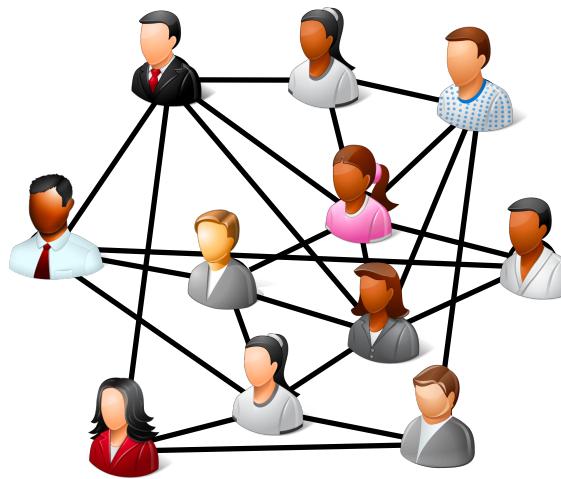


Interactomes



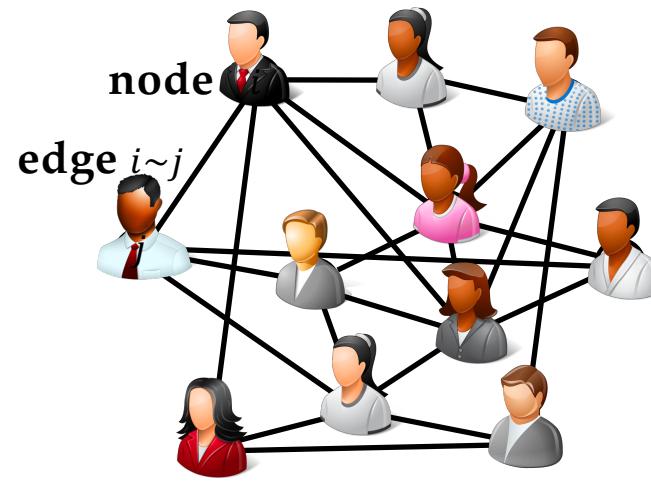
Social networks

Graphs



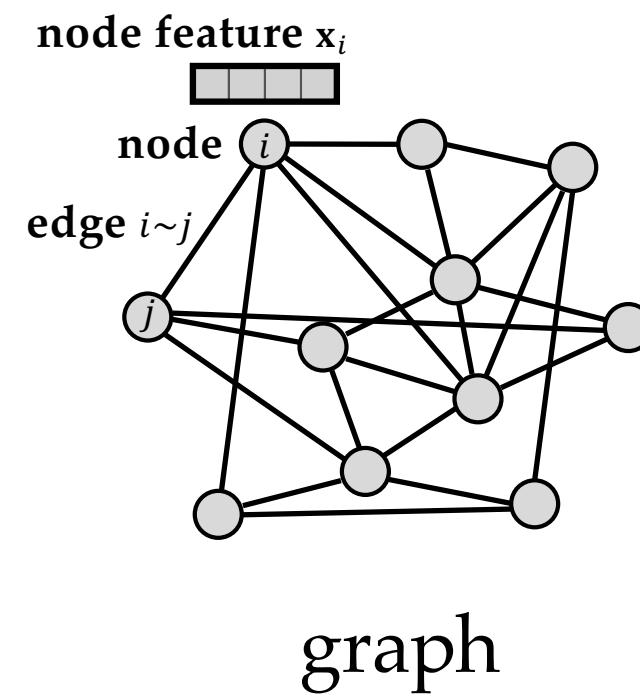
social network

Graphs

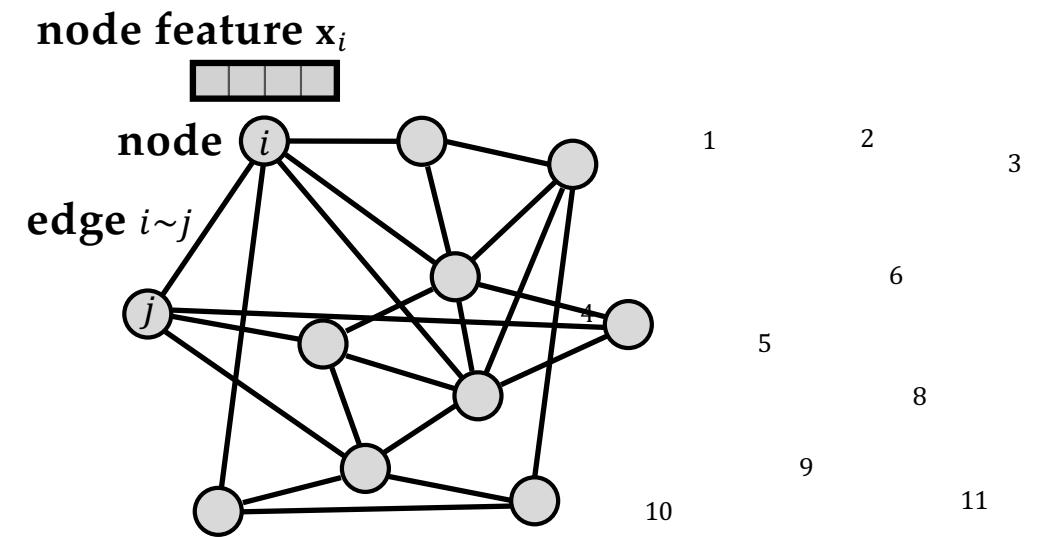


social network

Graphs



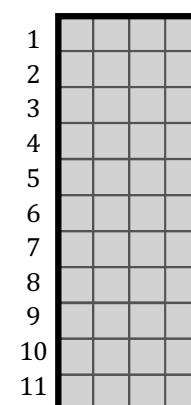
Key Structural Properties of Graphs



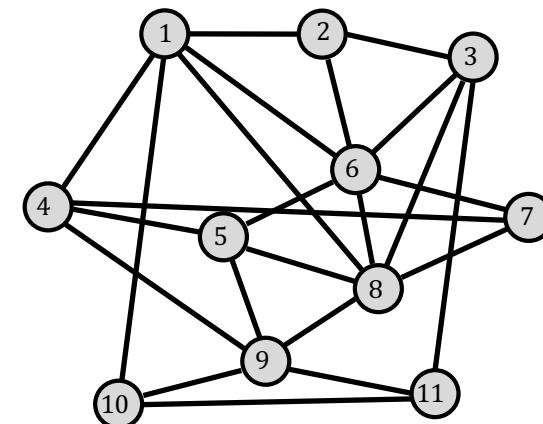
arbitrary graph ordering of nodes

Key Structural Properties of Graphs

Feature
matrix $n \times d$

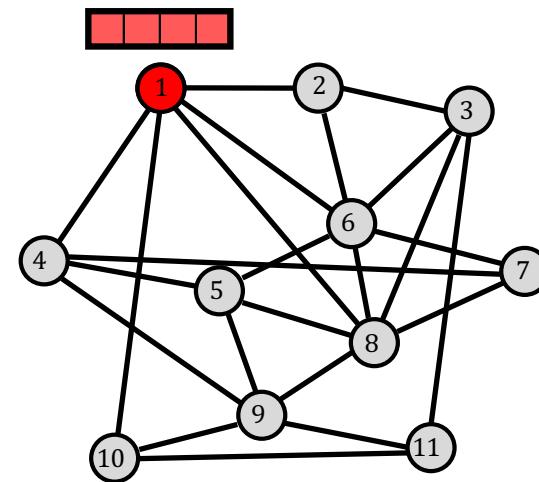
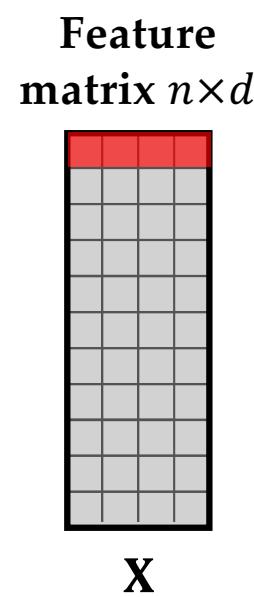
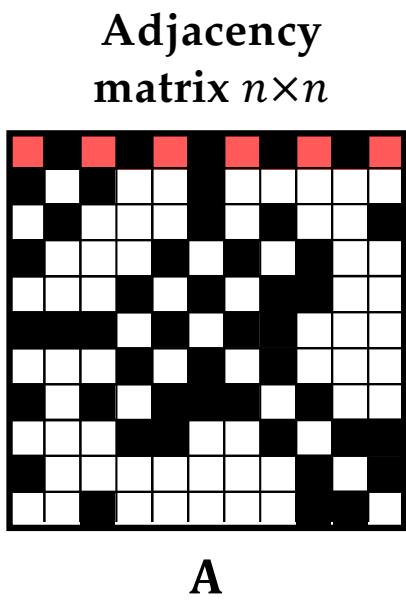


X



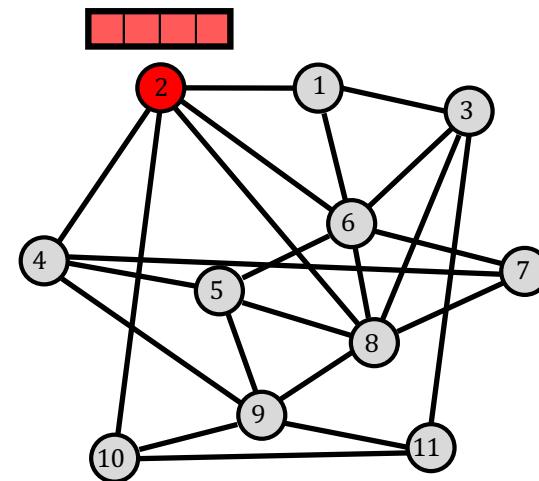
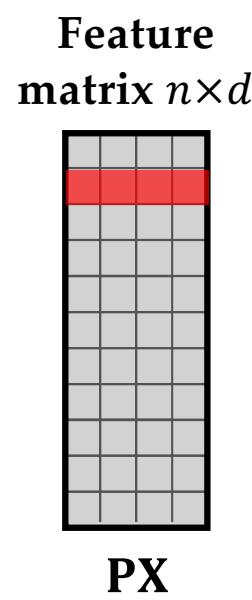
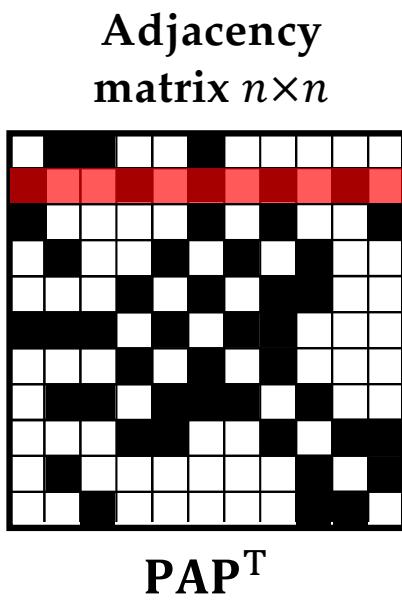
arbitrary ordering of nodes

Key Structural Properties of Graphs

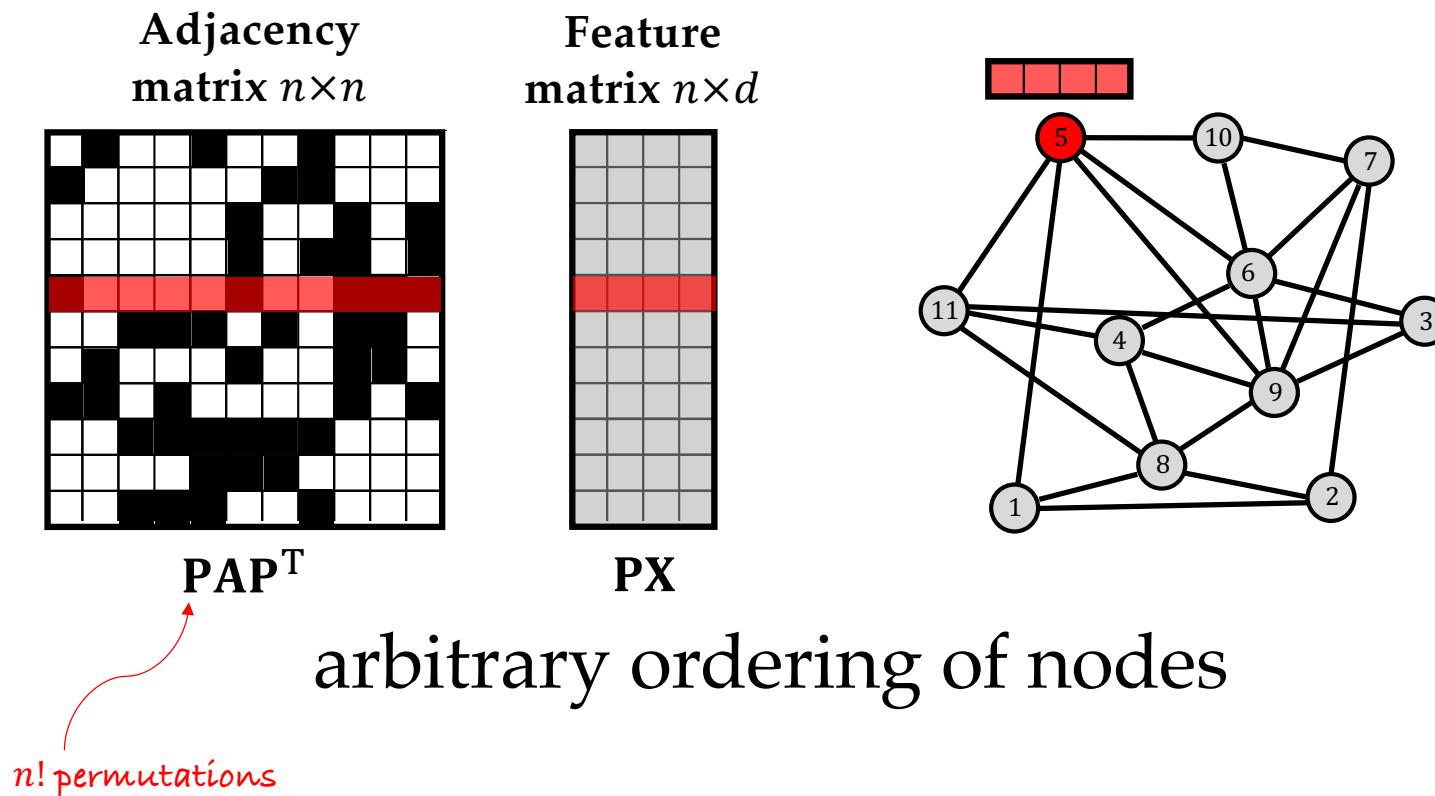


arbitrary ordering of nodes

Key Structural Properties of Graphs

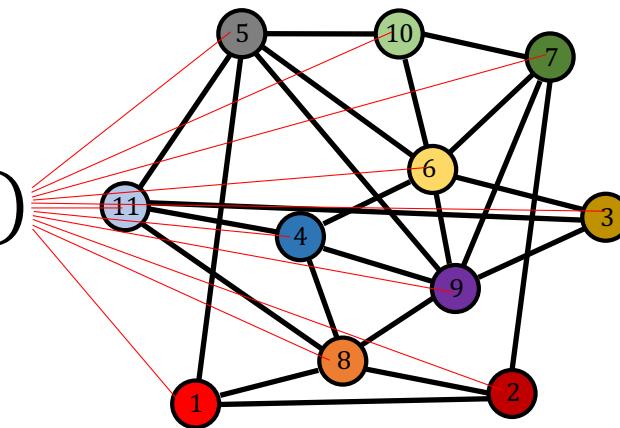


Key Structural Properties of Graphs



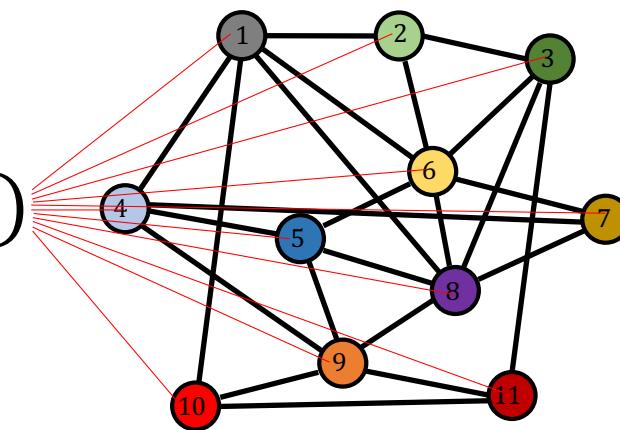
Invariant Graph Functions

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



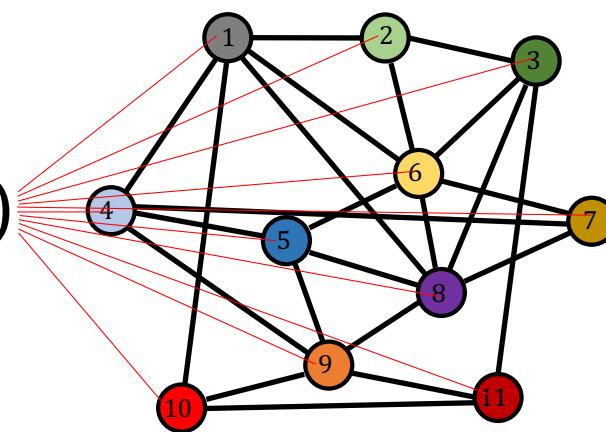
Invariant Graph Functions

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



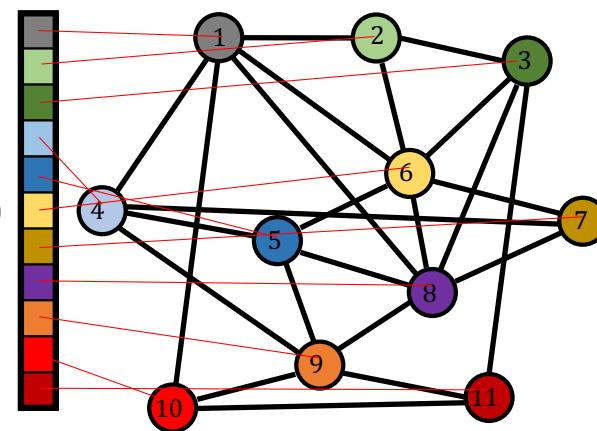
Invariant Graph Functions

permutation-invariant
 $f(\mathbf{P}\mathbf{X}, \mathbf{P}\mathbf{A}\mathbf{P}^T) = f(\mathbf{X}, \mathbf{A})$



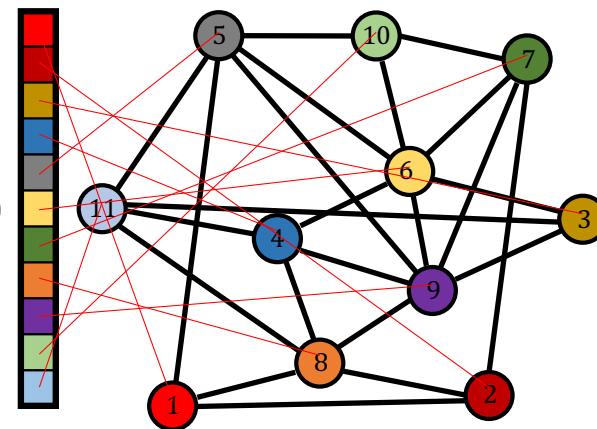
Equivariant Graph Functions

node function $F(X, A)$



Equivariant Graph Functions

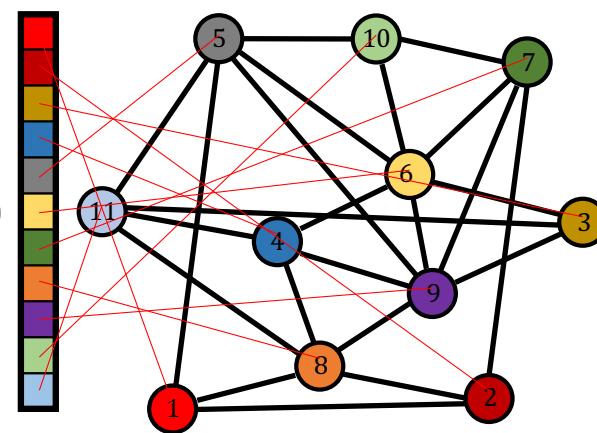
node function $F(X, A)$



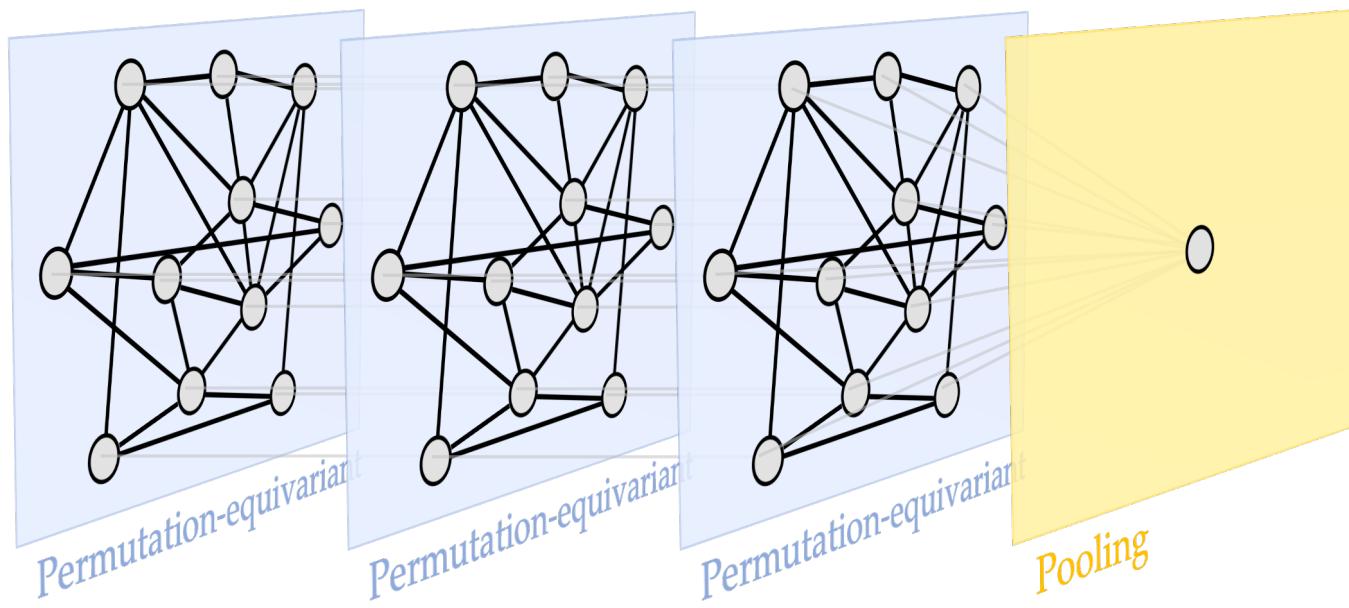
Equivariant Graph Functions

permutation-equivariant

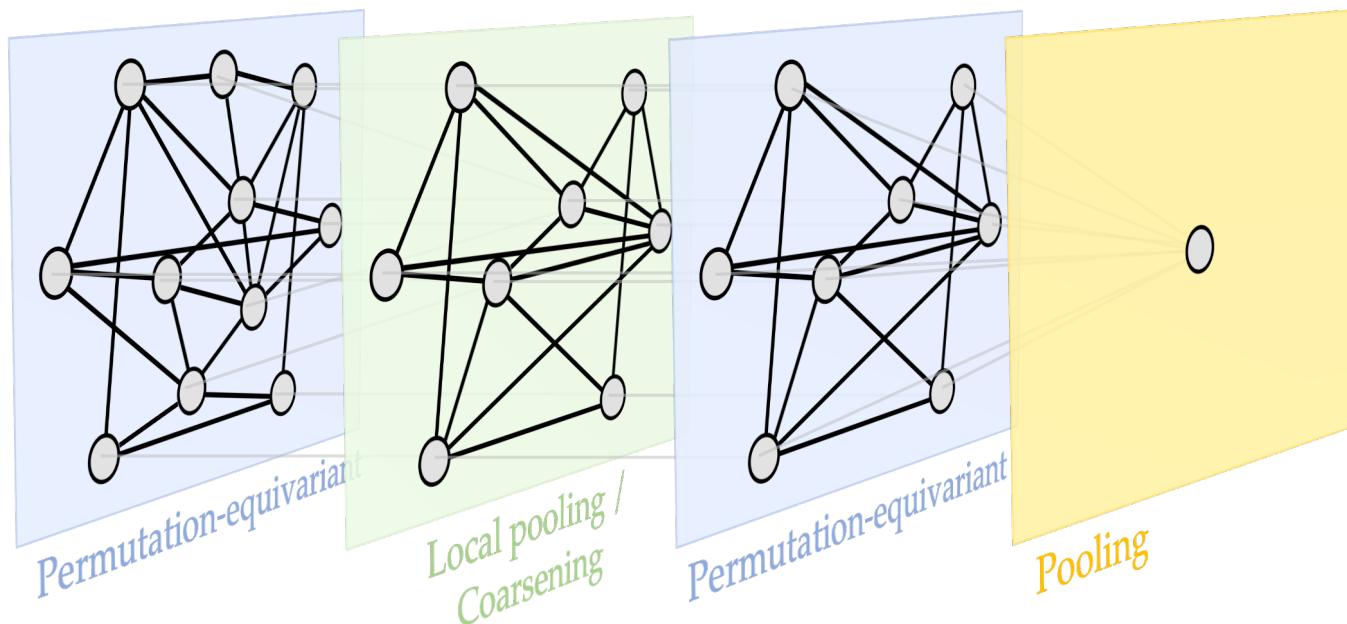
$$F(\mathbf{P}\mathbf{X}, \mathbf{P}\mathbf{A}\mathbf{P}^\top) = \mathbf{P}F(\mathbf{X}, \mathbf{A})$$



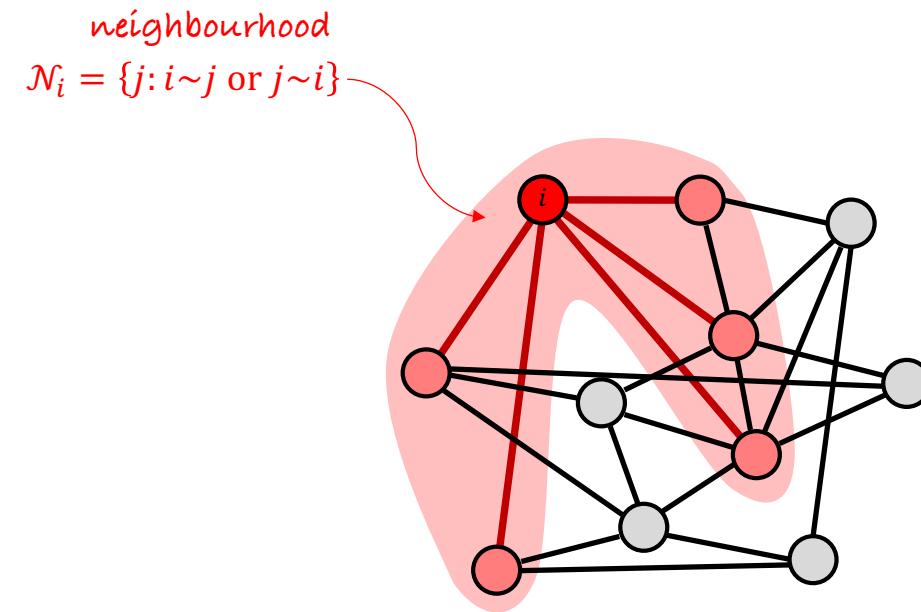
Graph Neural Networks



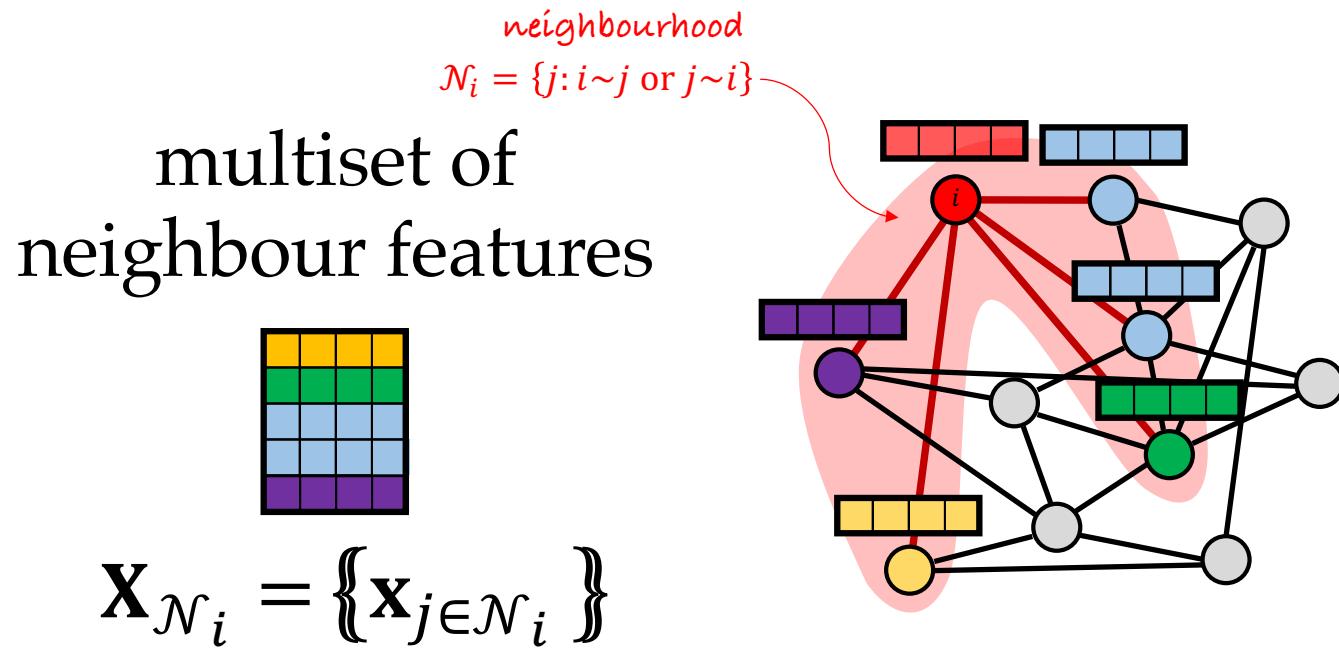
Graph Neural Networks



A General Blueprint for Constructing Graph Functions

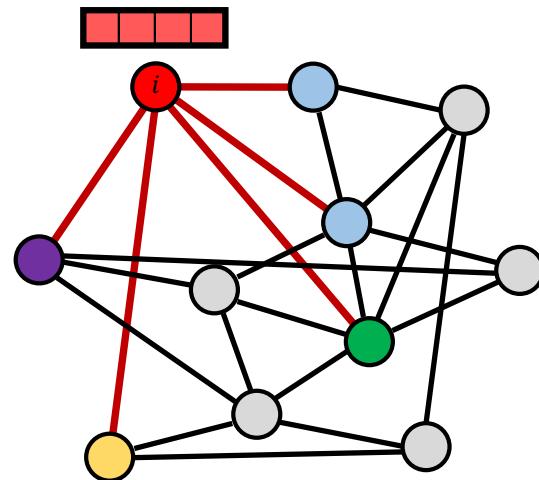


A General Blueprint for Constructing Graph Functions



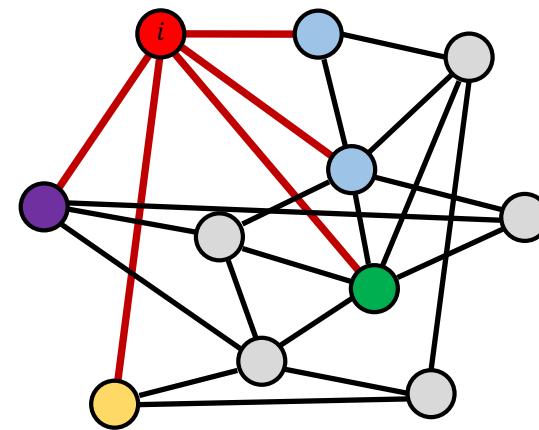
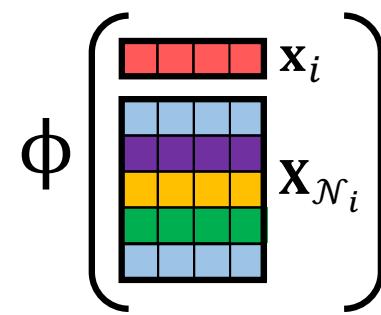
A General Blueprint for Constructing Graph Functions

multiset of
local function
neighbour features

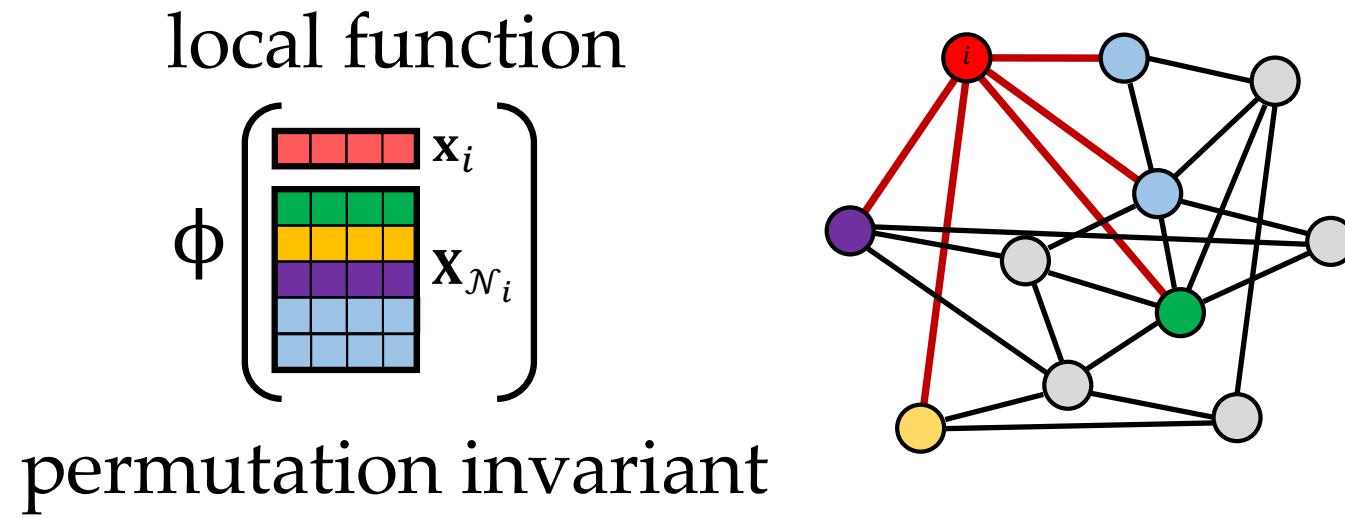
$$\phi \left(\begin{array}{c} \text{matrix} \\ \text{of} \\ \text{neighbour} \\ \text{features} \end{array} \right) \quad \mathbf{x}_i$$
$$\mathbf{X}_{\mathcal{N}_i} = \{\mathbf{x}_{j \in \mathcal{N}_i}\}$$


A General Blueprint for Constructing Graph Functions

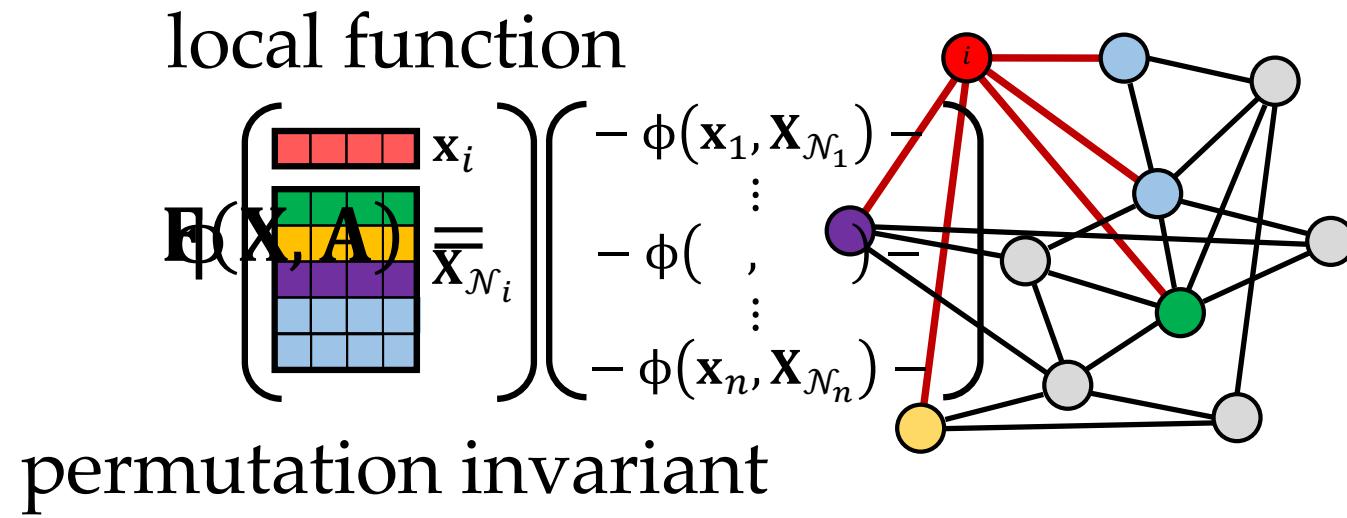
local function



A General Blueprint for Constructing Graph Functions



A General Blueprint for Constructing Graph Functions



A General Blueprint for Constructing Graph Functions

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

permutation equivariant

"Flavours" of Graph Neural Networks

$$f(\mathbf{x}_i) = \phi \left(\mathbf{x}_i, \bigcup_{j \in \mathcal{N}_i} \psi(\mathbf{x}_j) \right)$$

permutation-invariant
aggregation operator, e.g. sum

new feature of
node i

learnable
functions

The diagram illustrates the update rule for a node i in a graph. The new feature vector $f(\mathbf{x}_i)$ is computed as a function ϕ of the node's own initial feature \mathbf{x}_i and the aggregated features from its neighborhood. The neighborhood features are obtained by applying a learnable function ψ to the features of each node j in the neighborhood \mathcal{N}_i . The aggregation is performed using a permutation-invariant operator, such as a sum.

"Flavours" of Graph Neural Networks

$$f(\mathbf{x}_i) = \phi \left(\mathbf{x}_i, \bigcup_{j \in \mathcal{N}_i} c_{ij} \psi(\mathbf{x}_j) \right)$$

“convolutional”

importance of node j to
the representation of i

"Flavours" of Graph Neural Networks

$$f(\mathbf{x}_i) = \phi \left(\mathbf{x}_i, \bigcup_{j \in \mathcal{N}_i} a(\mathbf{x}_i, \mathbf{x}_j) \psi(\mathbf{x}_j) \right)$$

“attentional”

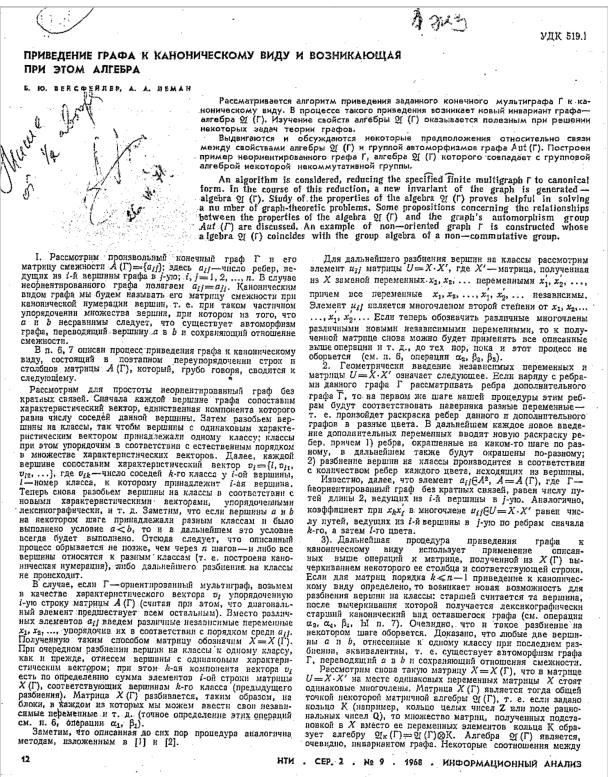
Monti et al. 2017; Veličković et al. 2018 (GAT)

Message Passing

$$f(\mathbf{x}_i) = \phi \left(\mathbf{x}_i, \bigcup_{j \in \mathcal{N}_i} \psi(\mathbf{x}_i, \mathbf{x}_j) \right)$$

“message passing”

Weisfeiler-Lehman Test



A. Lehman

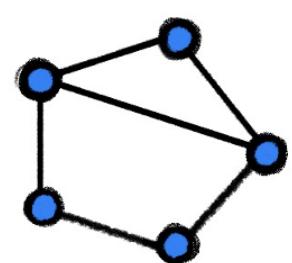


B. Weisfeiler

1968

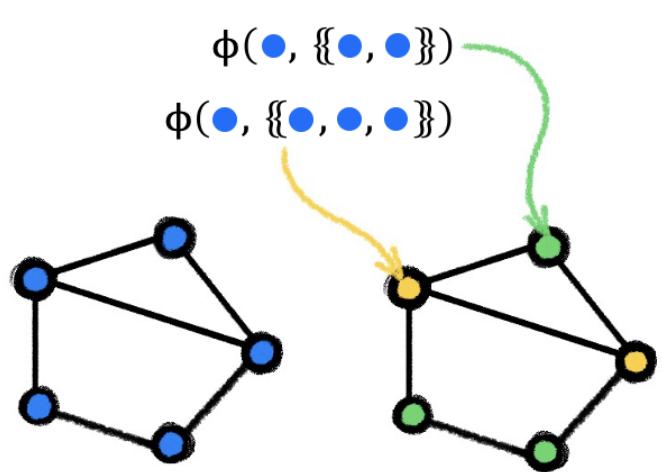
Weisfeiler, Lehman 1968; Portraits: Ihor Gorskiy

Weisfeiler-Lehman Test



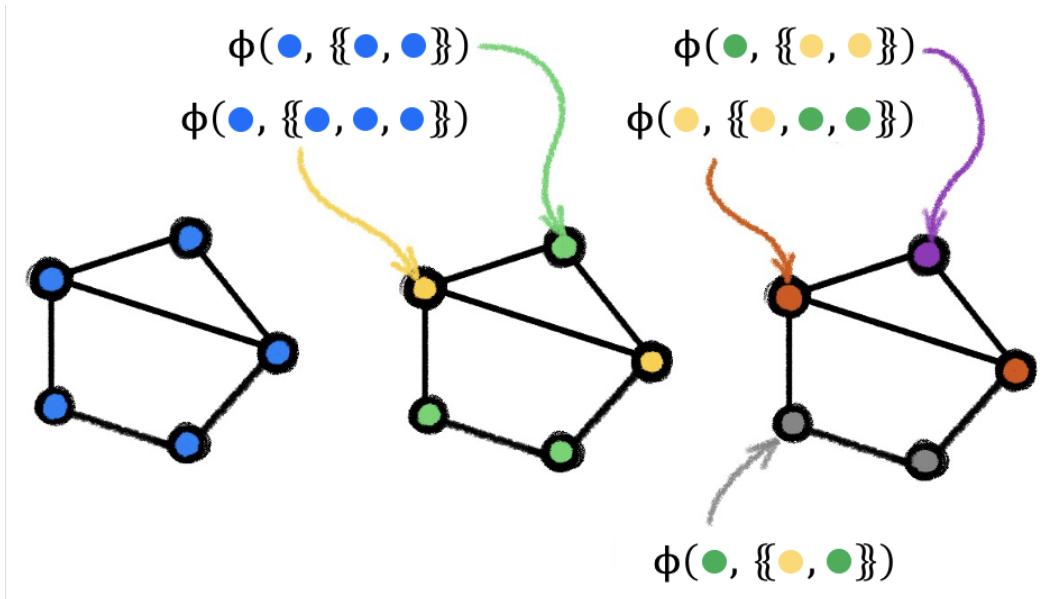
Weisfeiler, Lehman 1968

Weisfeiler-Lehman Test



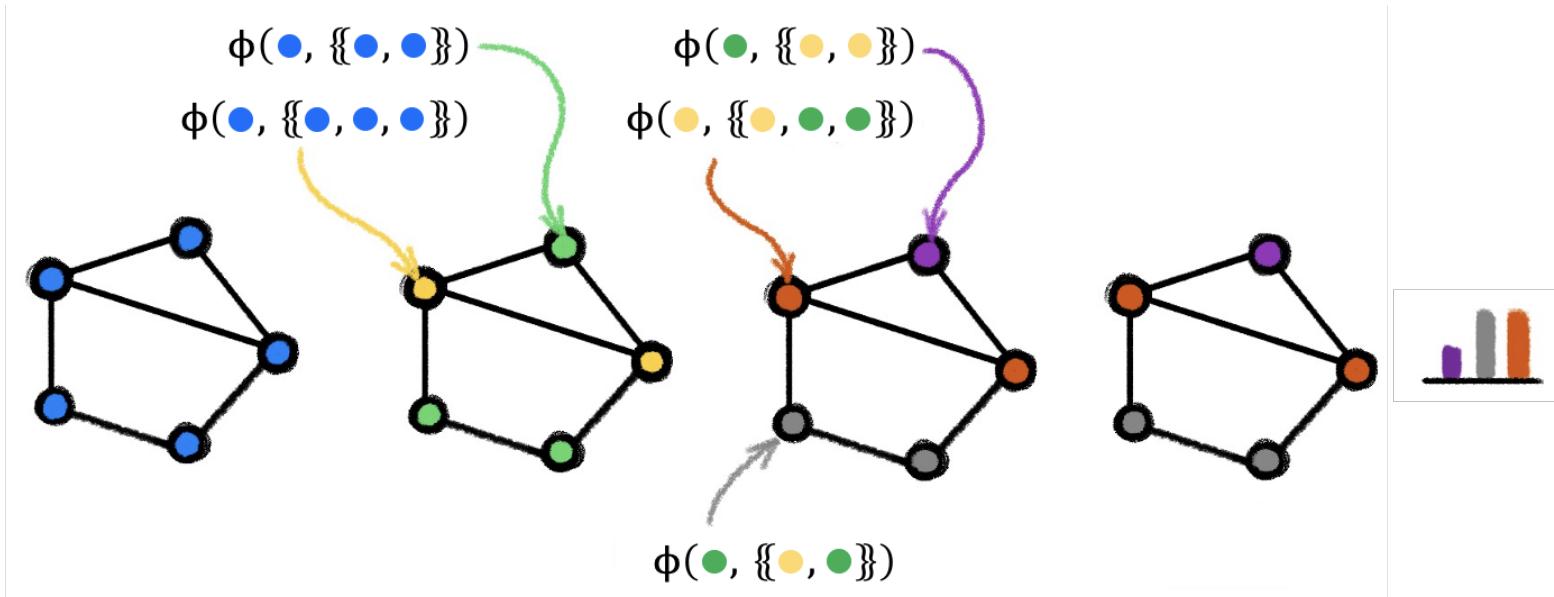
Weisfeiler, Lehman 1968

Weisfeiler-Lehman Test



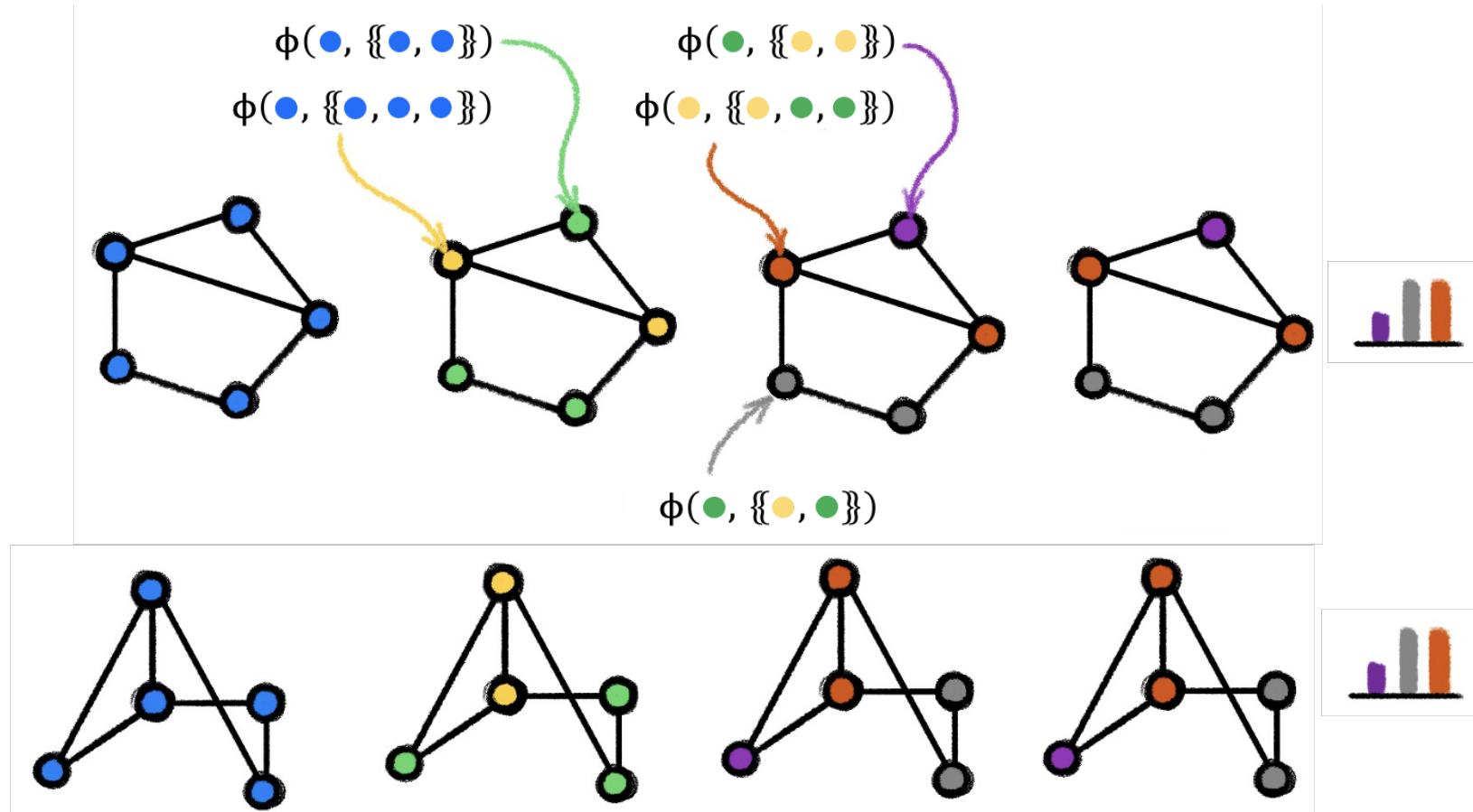
Weisfeiler, Lehman 1968

Weisfeiler-Lehman Test

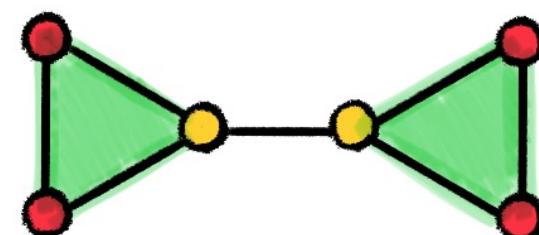
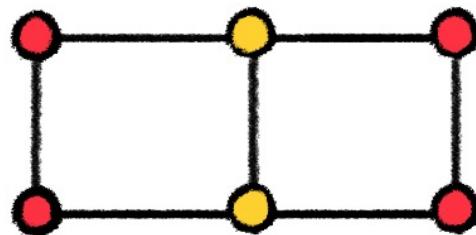


Weisfeiler, Lehman 1968

Weisfeiler-Lehman Test

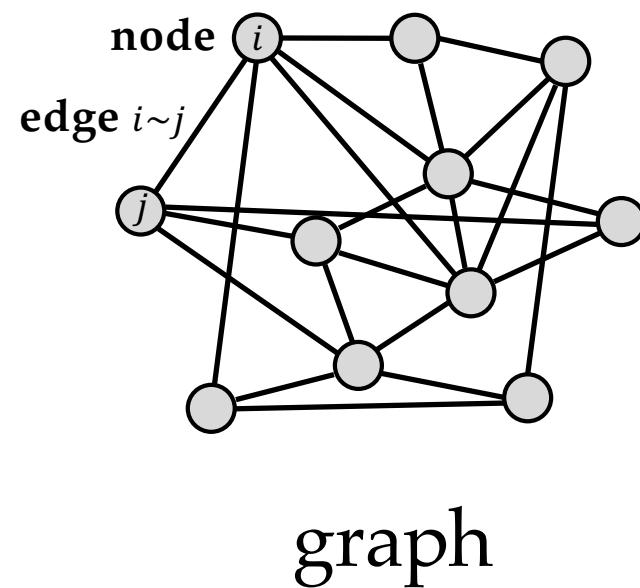


necessary but insufficient condition!

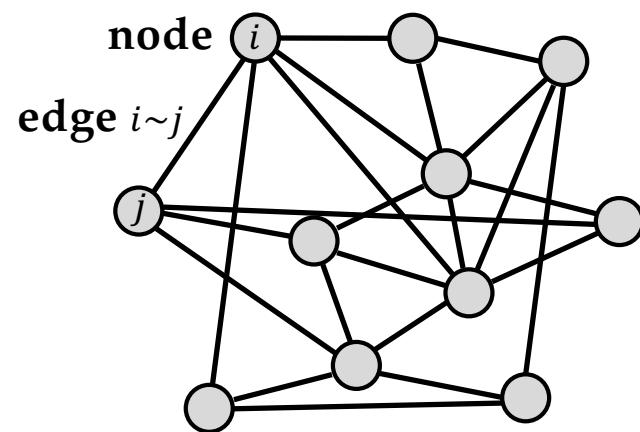


non-isomorphic graphs that are WL-equivalent

Special Cases of GNNs

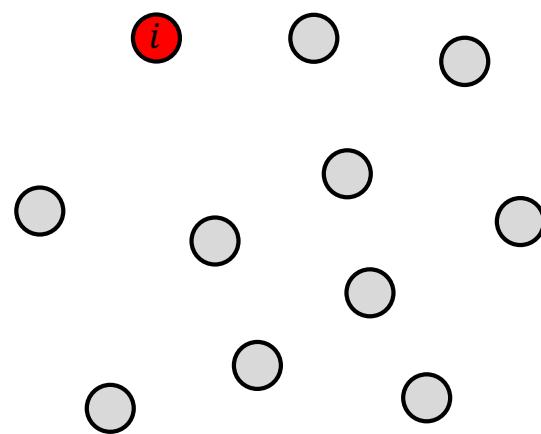


Special Cases of GNNs



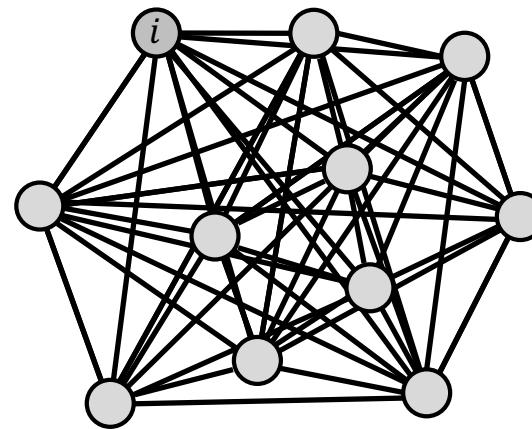
graph

DeepSets



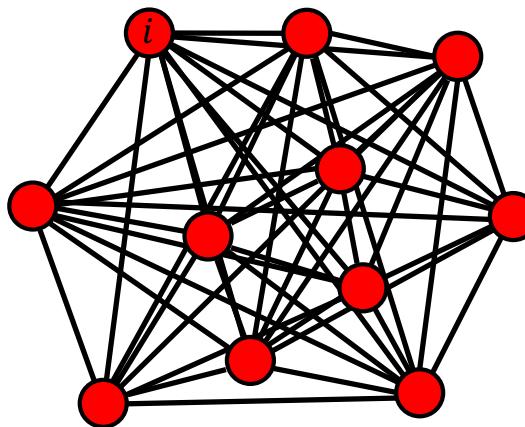
Zaheer et al 2017; Qi et al 2017

Transformers



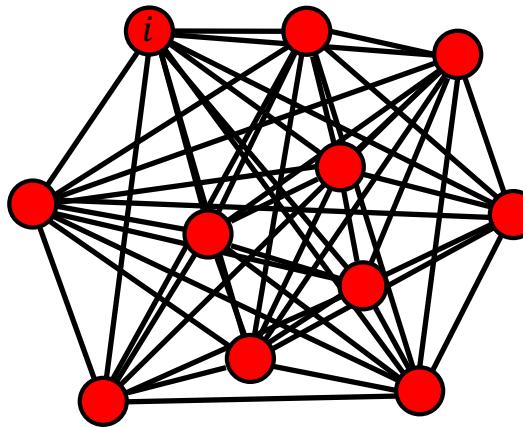
complete graph

Transformers



$$\phi \left(\mathbf{x}_i, \sum_{j=1}^n c_{ij} \psi(\mathbf{x}_j) \right)$$

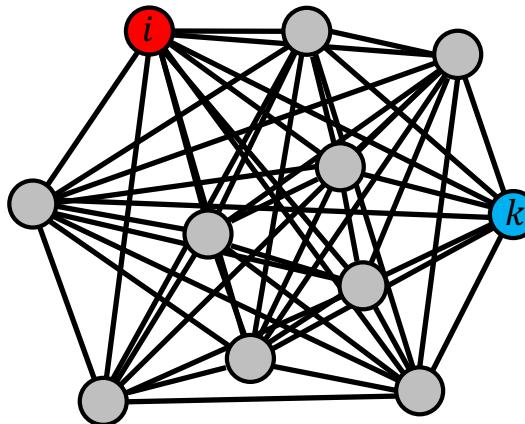
Transformers



$$\phi \left(\mathbf{x}_i, \bigcup_{j=1}^n a(\mathbf{x}_i, \mathbf{x}_j) \psi(\mathbf{x}_j) \right)$$

Vaswani et al. 2017

Transformers

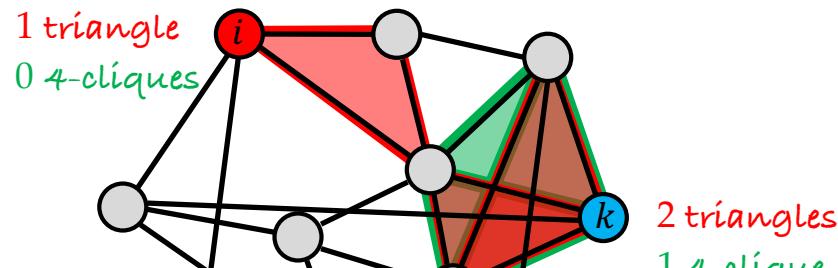


$$\phi \left(\mathbf{x}_i, \bigcup_{j=1}^n a(\mathbf{x}_i, \mathbf{x}_j, \mathbf{p}_i, \mathbf{p}_j) \psi(\mathbf{x}_j) \right)$$

↑
positional encoding

Vaswani et al. 2017

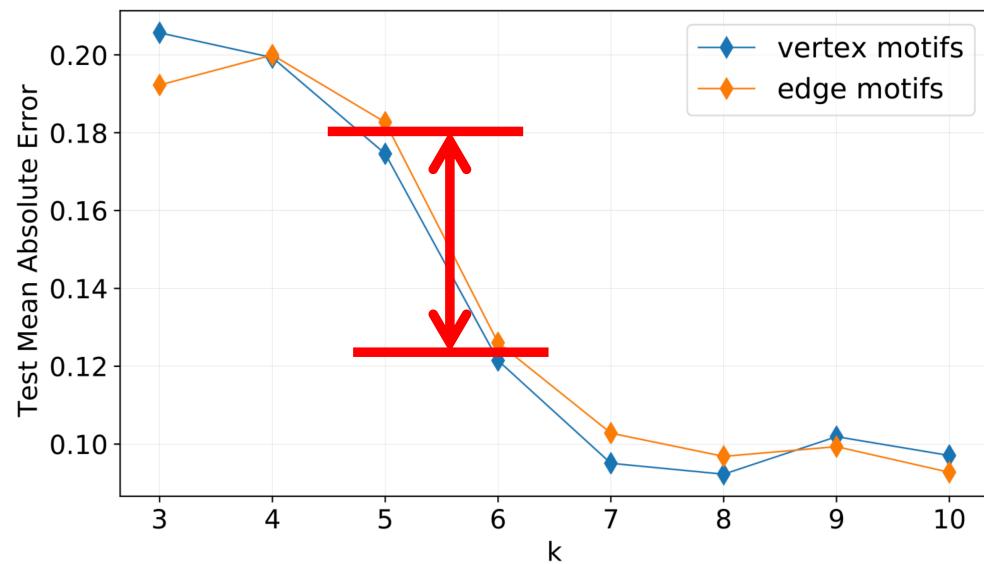
Graph Substructure Networks



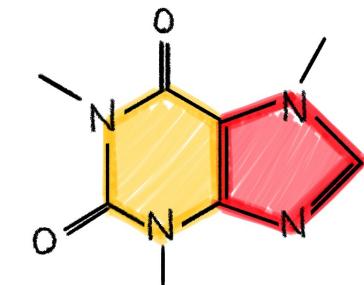
$$\phi \left(\mathbf{x}_i, \bigwedge_{j \in \mathcal{N}_i} \psi(\mathbf{x}_i, \mathbf{x}_j, \mathbf{p}_i) \right)$$

↑
structural encoding

Graph Substructure Networks

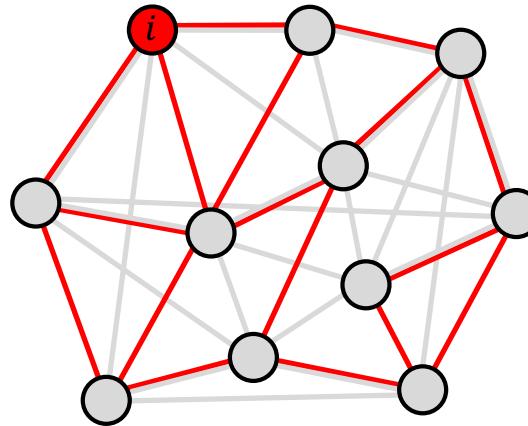


Molecule property prediction on ZINC
using GSN with k -cycles



Molecule of caffeine

decouple computational graph from the input graph



sampling

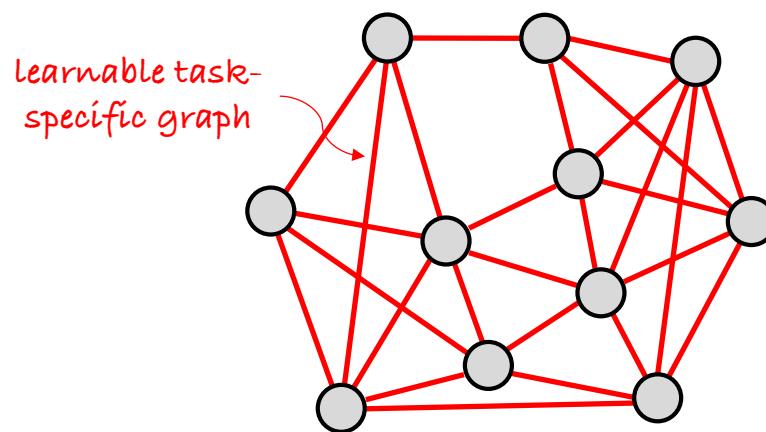
rewiring

multi-hop filters

$$\phi \left(\mathbf{x}_i, \bigwedge_{j \in \mathcal{N}'_i} \psi(\mathbf{x}_i, \mathbf{x}_j) \right)$$

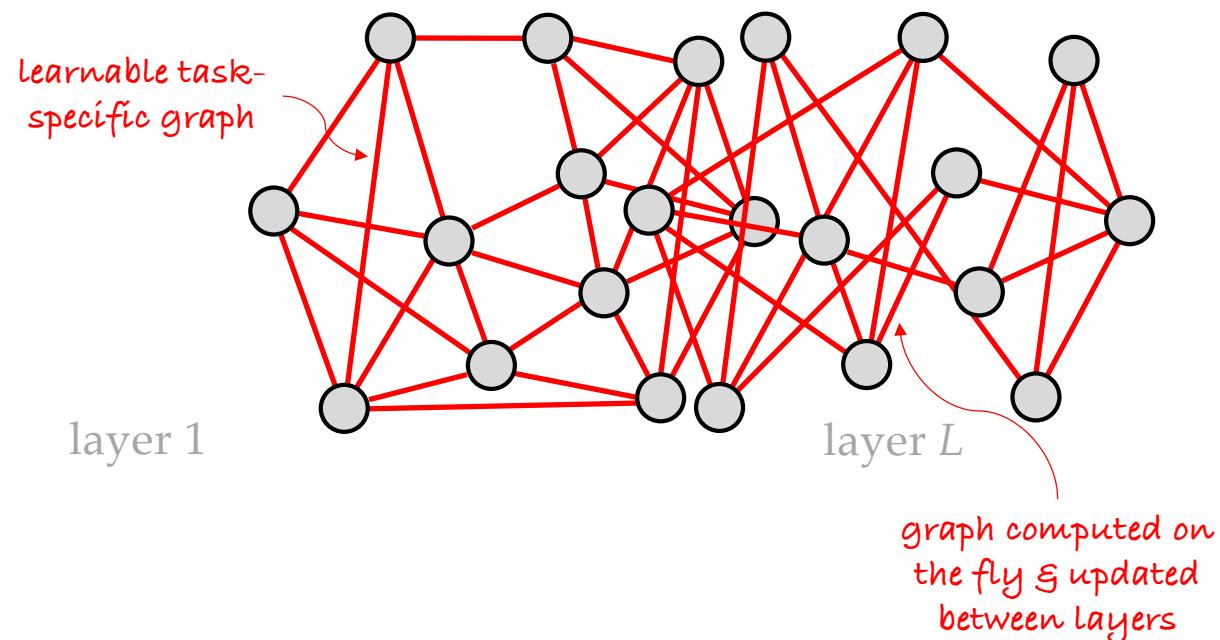
Hamilton et al 2017; Klicpera et al. 2019; Alon & Yahav 2020; Frasca, Rossi et B 2020

Latent Graph Learning

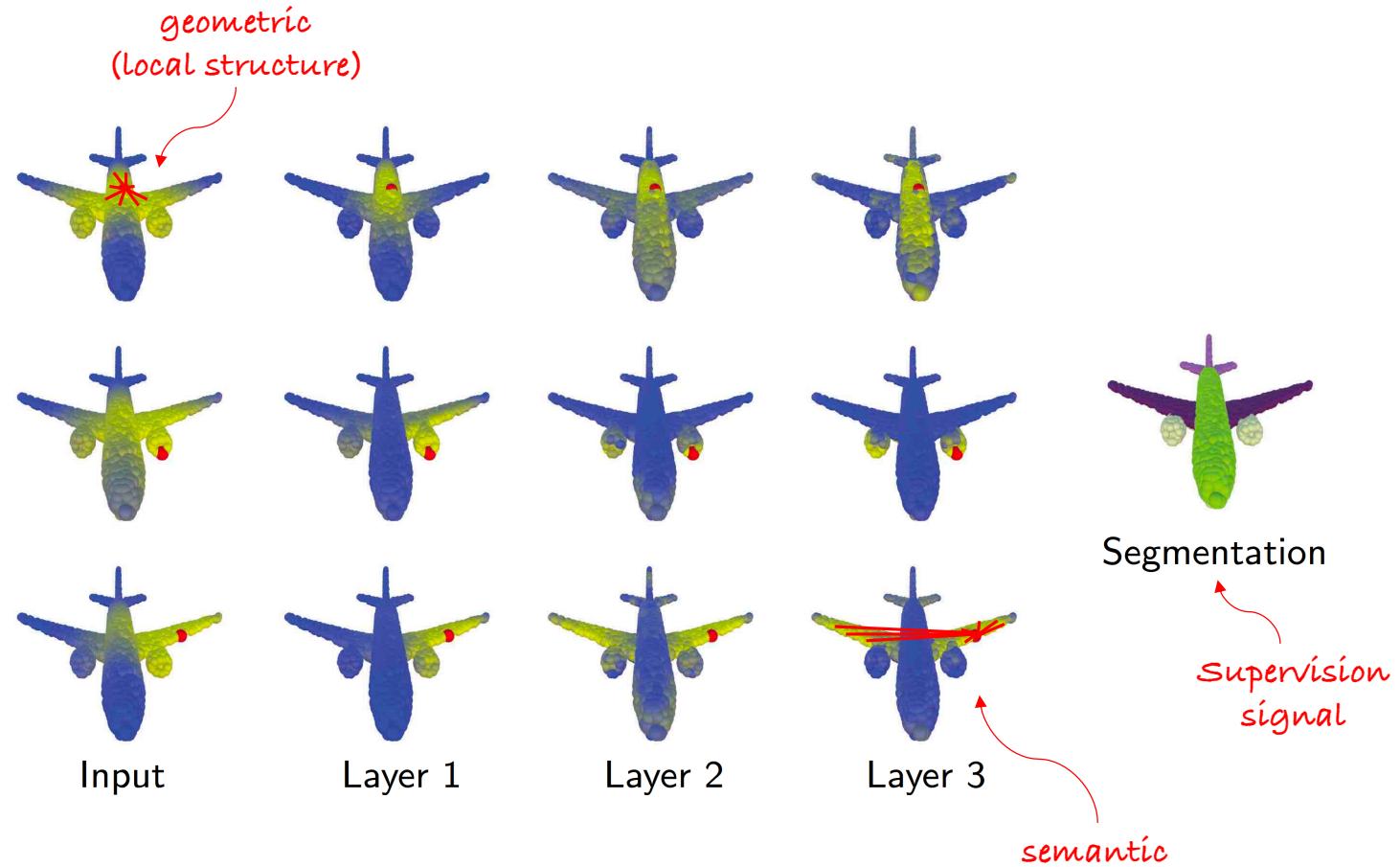


Wang et al 2018; Franceschi et al. 2019; Kipf et al. 2020; Kazi, Cosmo et al. 2020; Cranmer et al. 2020

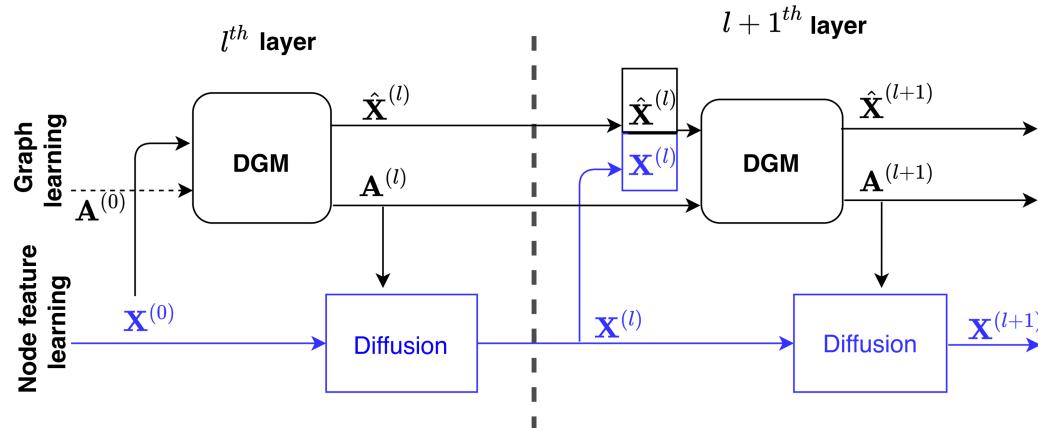
Dynamic Graph CNN



Wang et al 2018; Franceschi et al. 2019; Kipf et al. 2020; Kazi, Cosmo et al. 2020; Cranmer et al. 2020



Differentiable Graph Module

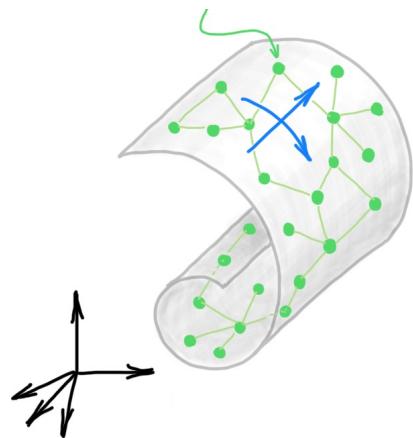


Differentiable Graph Module (DGM) allowing to construct the graph from the data and use it for feature learning

Method	TADPOLE		UK Biobank	
	Transductive	Inductive	Transductive	Inductive
DGCNN	84.59±4.33	82.99±4.91	58.35±0.91	51.84±8.16
LDS	87.06±3.67	†	OOM	†
cDGM	92.91±2.50	91.85±2.62	61.32±1.51	55.91±3.49
dDGM	94.10±2.12	92.17±3.65	63.22±1.12	57.34±5.32

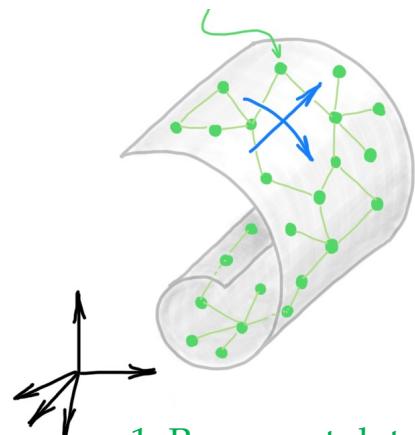
Manifold Learning

intrinsically low-dimensional
data in a high-dimensional space

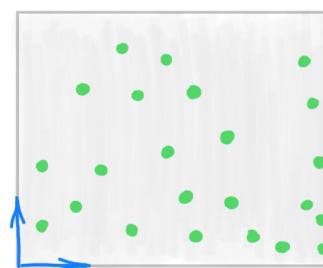


Manifold Learning

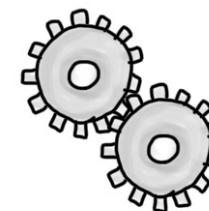
intrinsically low-dimensional
data in a high-dimensional space



1. Represent data
structure as a graph



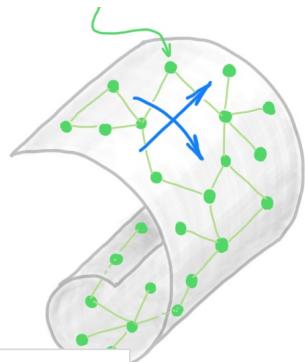
2. Compute low-
dimensional embedding



3. Apply ML

Manifold Learning 2.0

intrinsically low-dimensional
data in a high-dimensional space



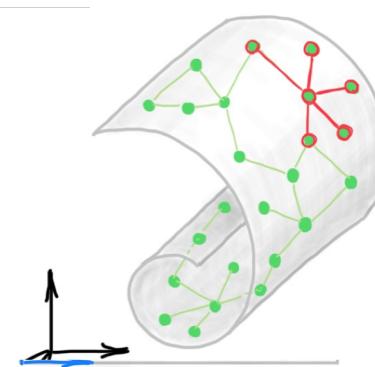
Latent graph neural networks: Manifold learning 2.0?

Can we use graph neural networks when
the graph is unknown?

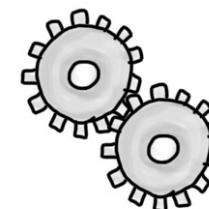


Michael Bronstein
Sep 10, 2020 · 12 min read ★

resent data
re as a graph



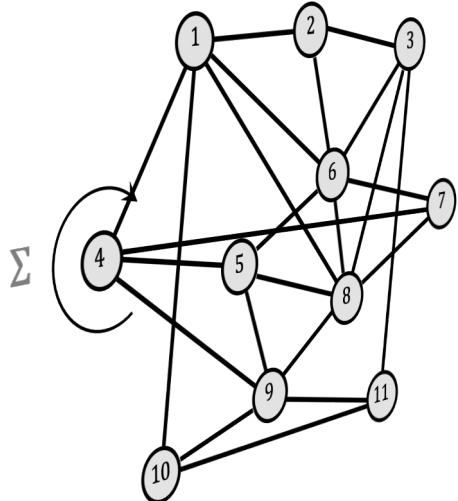
1. Represent data structure
as a graph
2. apply ML
directly on the graph
3. Apply ML



GNNs BEYOND NODES & EDGES

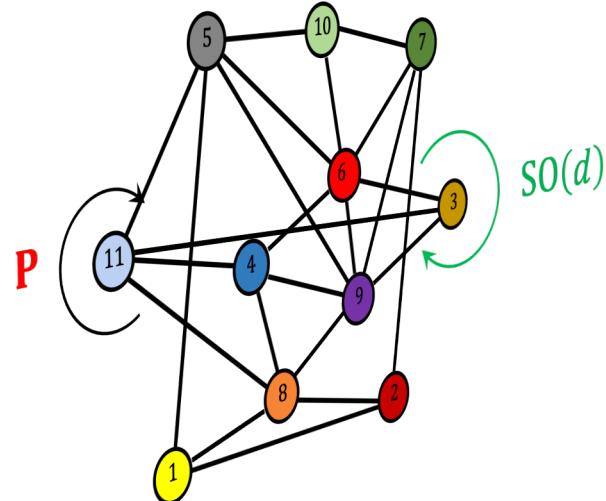
Data Symmetry in Geometric Graphs

Graph $G = (V, E)$



Permutation group Σ_n

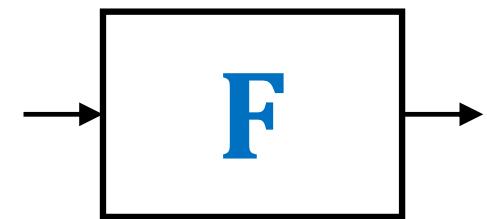
Node features $\chi(G)$



Permutation matrix P

Rotation R

functions $\mathcal{F}(\chi(\Omega))$



Equivariant message passing

$$F(PXR, PAP^T) = PF(X, A)R$$

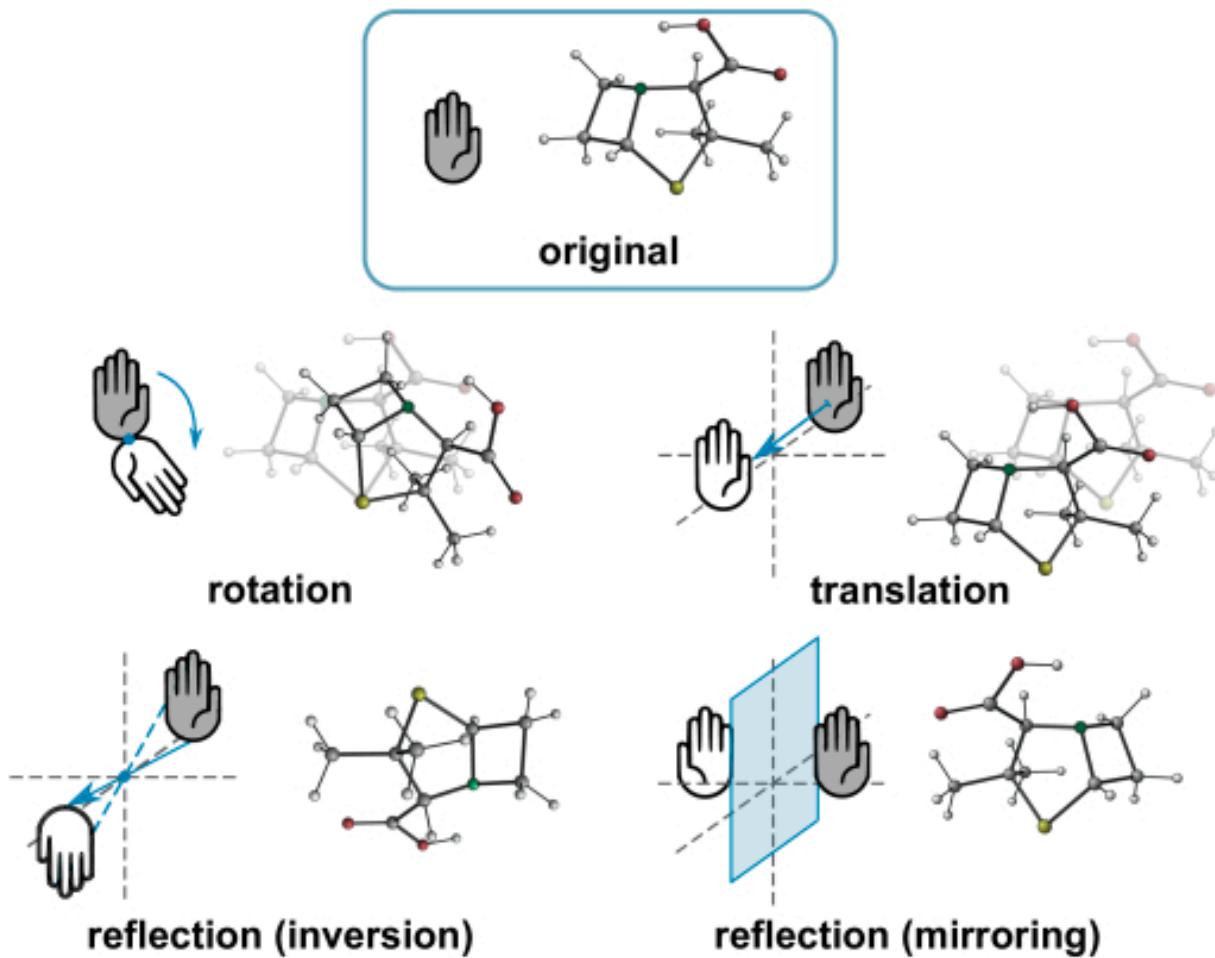
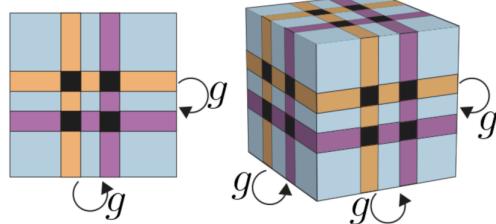


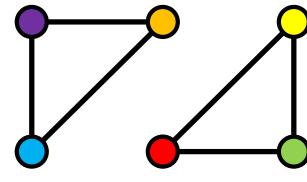
Figure: Atz et al. 2021

Towards More Expressive GNNs



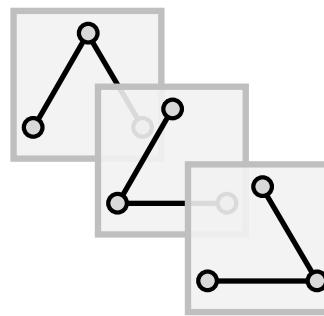
k-WL tests

Maron et al. 2019
Morris et al. 2019



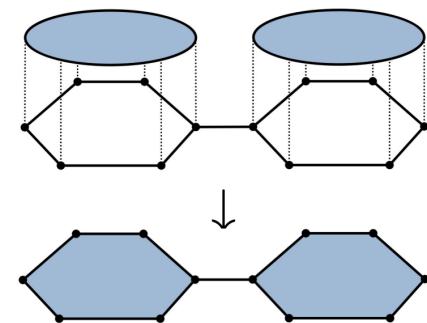
Positional &
Structural
encoding

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



Subgraph
GNNs

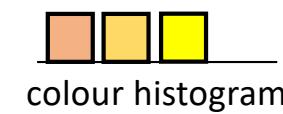
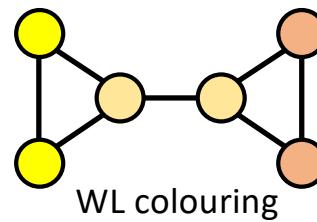
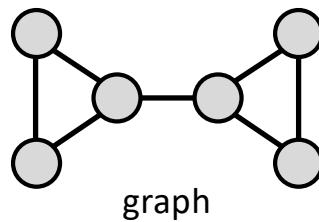
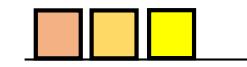
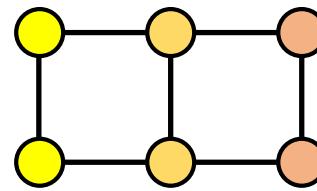
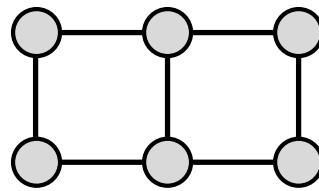
Papp et al. 2021
Cotta et al. 2021
Zhao et al. 2021
Bevilacqua, Frasca et B 2021



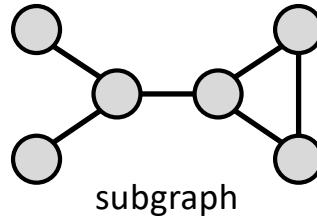
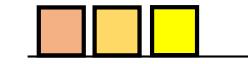
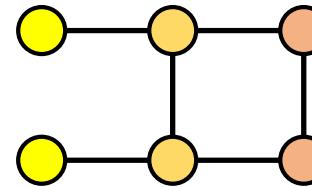
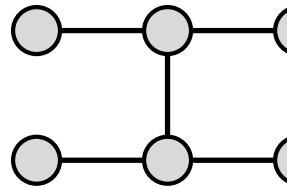
Topological
message passing

Bodnar, Frasca et B 2021

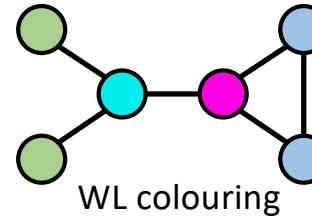
Subgraph GNNs



Subgraph GNNs



subgraph



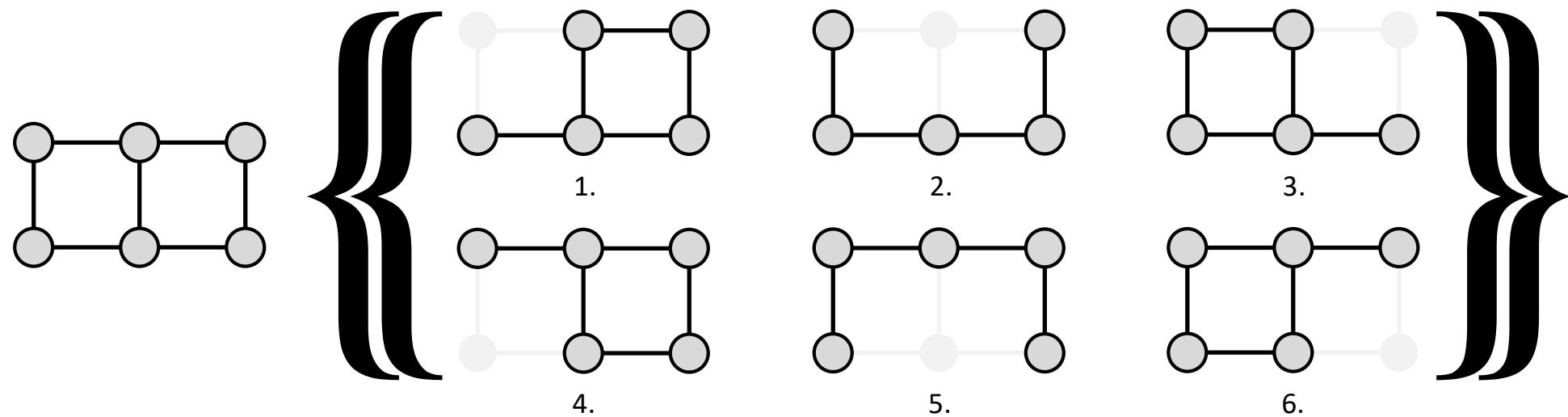
WL colouring



colour histogram

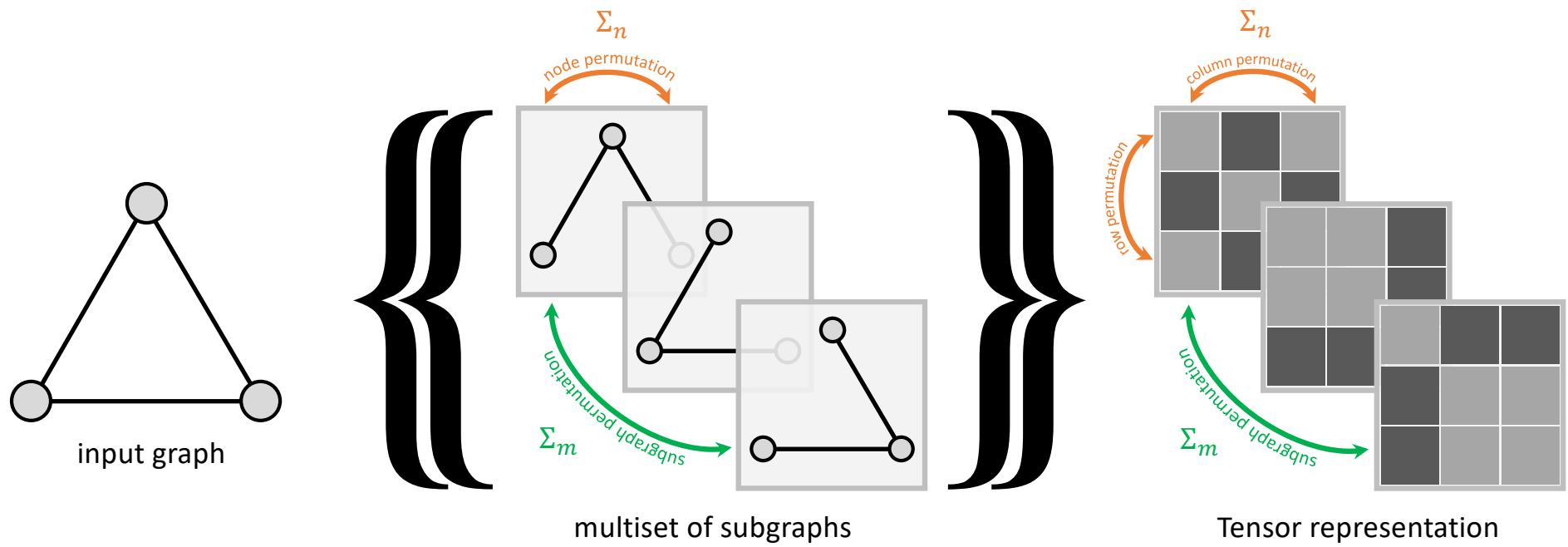
Graph perturbation allows to distinguish between structures
otherwise indistinguishable by Weisfeiler-Lehman

Collection of Subgraphs

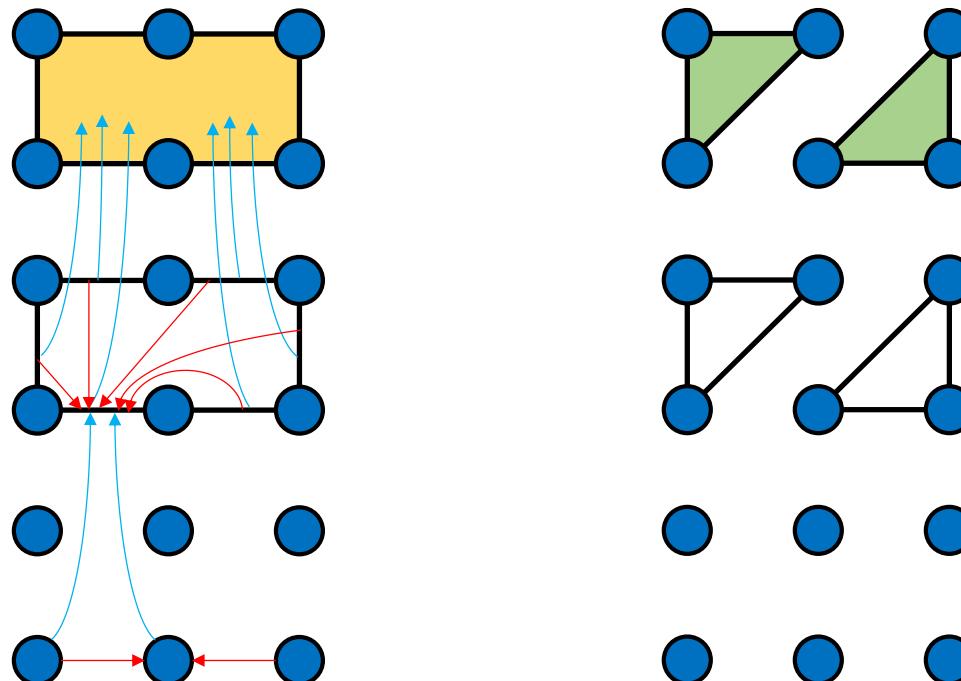


A multiset of subgraphs obtained by node deletion

Equivariant Subgraph Aggregation Networks

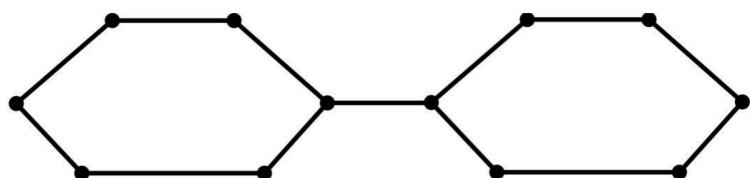


Topological Message Passing

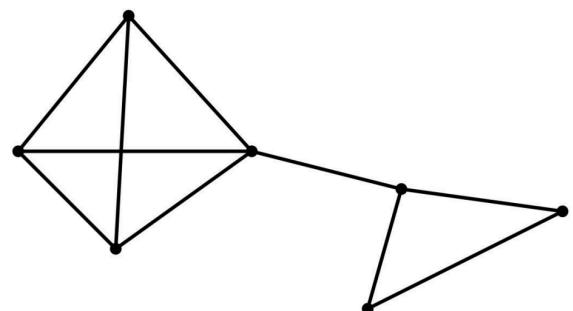
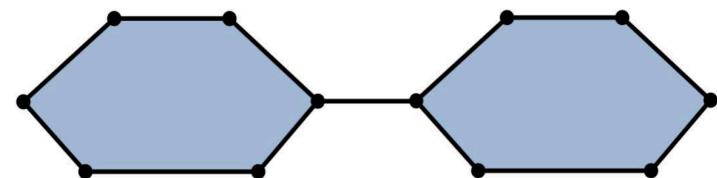


Lift the graph into a cell complex

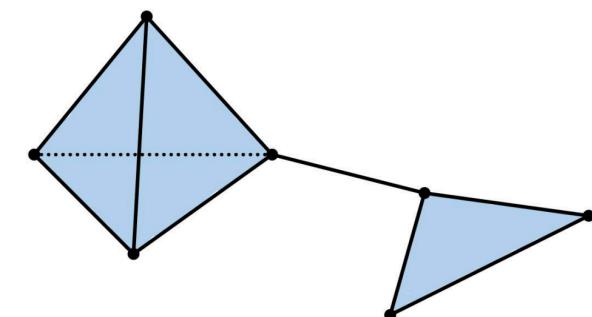
Cellular Lifting Maps



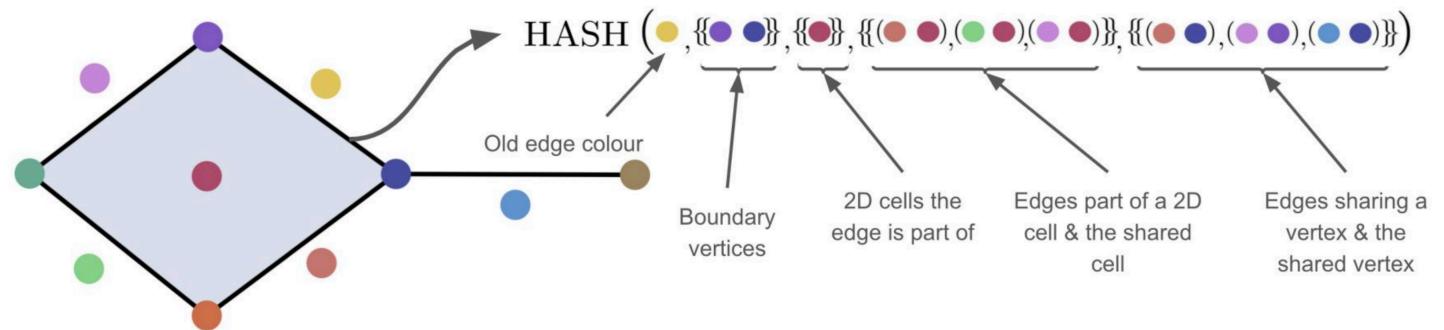
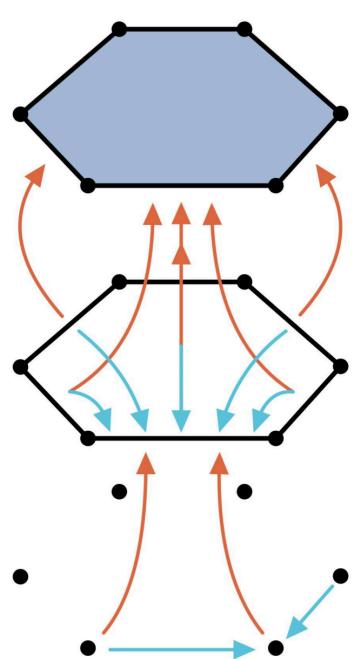
cycles



cliques



Cellular Weisfeiler-Lehman



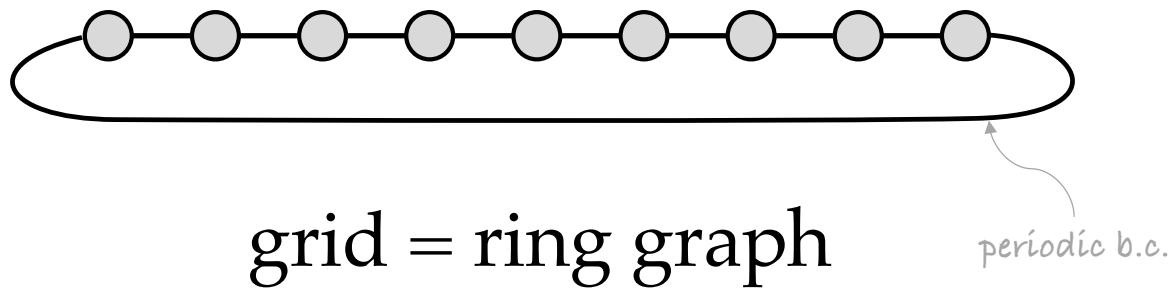
Theorem 15. For some finite k , there exists a pair of graphs indistinguishable by 3-WL but distinguishable by CWL with the lifting maps from Corollary 14. For the clique complex and induced cycle liftings, the statement holds for $k \geq 4$. For the simple cycle based lifting, it holds for $k \geq 8$.

Corollary 14. For all $k \geq 3$, the following lifting transformations make CWL strictly more powerful than the WL test. (1) The clique complex lifting considering cliques of size at most k . (2) The map that attaches 2-cells to all the simple cycles of size at most k . (3) The map that attaches 2-cells to all the induced cycles of size at most k . (4) The union of all the transformations above.

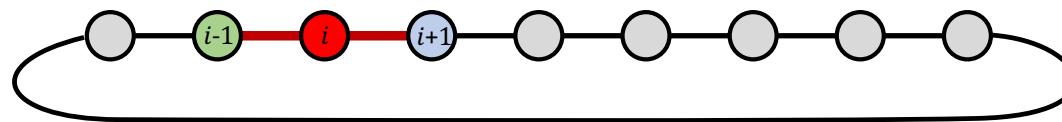
Cellular WL is strictly more powerful than WL with appropriate lifting transformation

GRIDS

Grids vs Graphs

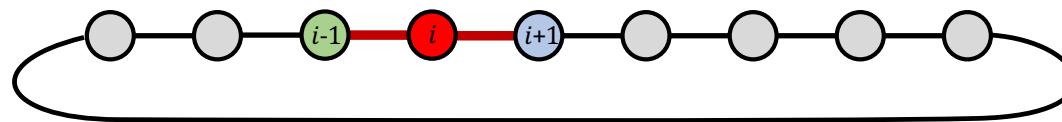


Grids vs Graphs



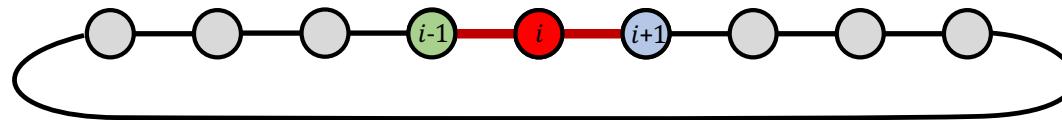
fixed neighbourhood structure

Grids vs Graphs



fixed neighbourhood structure

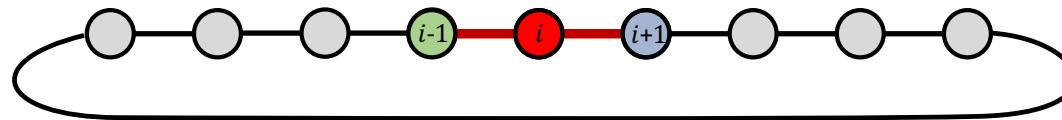
Grids vs Graphs



local aggregation function

$$f(\mathbf{x}_i) = \phi(\mathbf{x}_i, \{\mathbf{x}_{i-1}, \mathbf{x}_{i+1}\})$$

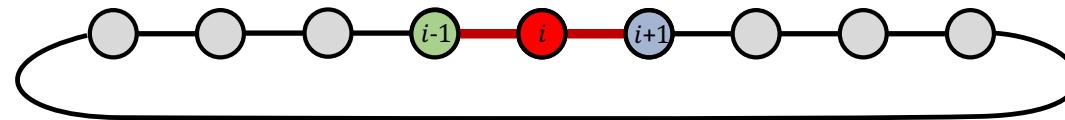
Grids vs Graphs



local aggregation function

$$f(\mathbf{x}_i) = \phi(\mathbf{x}_{i-1}, \mathbf{x}_i, \mathbf{x}_{i+1})$$

Grids vs Graphs



linear local aggregation function

$$f(\mathbf{x}_i) = a\mathbf{x}_{i-1} + b\mathbf{x}_i + c\mathbf{x}_{i+1}$$

Convolution

$$f(\mathbf{X}) = \begin{bmatrix} b & c \\ a & b & c \\ & a & b & c \\ c & a & b \end{bmatrix} \mathbf{X}$$

circulant matrix = convolution

Convolution

vector of parameters θ

$$f(\mathbf{X}) = \begin{pmatrix} b & c & & \\ a & b & c & \\ & a & b & c \\ c & & a & b \end{pmatrix} \mathbf{X}$$

circulant matrix $\mathbf{C}(\theta)$

Deriving Convolution from Symmetry

vector of parameters θ

$$\begin{bmatrix} b & c & a \\ a & b & c \\ c & a & b \end{bmatrix} f(\mathbf{X}) = \mathbf{C} \begin{bmatrix} b & c & a \\ a & b & c \\ c & a & b \end{bmatrix} = \mathbf{C} \begin{bmatrix} a & x & y & z \\ x & a & y & z \\ y & z & a & x \\ z & x & y & a \end{bmatrix} \mathbf{X}$$

circulant matrices a commute
circulant matrix $\mathbf{C}(\theta)$

Deriving Convolution from Symmetry

$$\begin{bmatrix} b & c & a \\ a & b & c \\ c & a & b \end{bmatrix} \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix} = \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix} \begin{bmatrix} b & c & a \\ a & b & c \\ c & a & b \end{bmatrix}$$

circulant matrix \Rightarrow commutes with shift

Deriving Convolution from Symmetry

$$\begin{bmatrix} b & c & a \\ a & b & c \\ c & a & b \end{bmatrix} \left[\begin{array}{c} 1 \\ 1 \\ 1 \end{array} \right] = \left[\begin{array}{c} 1 \\ 1 \\ 1 \end{array} \right] \begin{bmatrix} b & c & a \\ a & b & c \\ c & a & b \end{bmatrix}$$

shift \mathbf{S}

convolution \Rightarrow shift-equivariant

$$\mathbf{CS} = \mathbf{SC}$$

Deriving Convolution from Symmetry

$$\begin{bmatrix} b & c & a \\ a & b & c \\ c & a & b \end{bmatrix} \left[\begin{array}{c} 1 \\ 1 \\ 1 \end{array} \right] \xrightarrow{\text{shift } S} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} b & c & a \\ a & b & c \\ c & a & b \\ a & b & c \\ b & a & b \end{bmatrix}$$

convolution \Leftrightarrow shift-equivariant

convolution emerges from translation symmetry

Deriving Fourier Transform from Symmetry

same eigenbasis for all convolutions

different eigenvalues for each convolution

$$\begin{bmatrix} b & c & & \\ a & b & c & \\ & a & b & c \\ c & & a & b \end{bmatrix} = \begin{bmatrix} | & & | & & | & & | \\ u_1 & \dots & u_n & & & & u_1^* \\ | & & | & & | & & | \end{bmatrix} \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix} \begin{bmatrix} | & & | & & | & & | \\ & \vdots & & & & & \\ | & & | & & | & & | \end{bmatrix} \begin{bmatrix} u_1^* \\ \vdots \\ u_n^* \end{bmatrix}$$

commuting matrices are jointly diagonalizable
vectors of \mathbf{S}

Deriving Fourier Transform from Symmetry

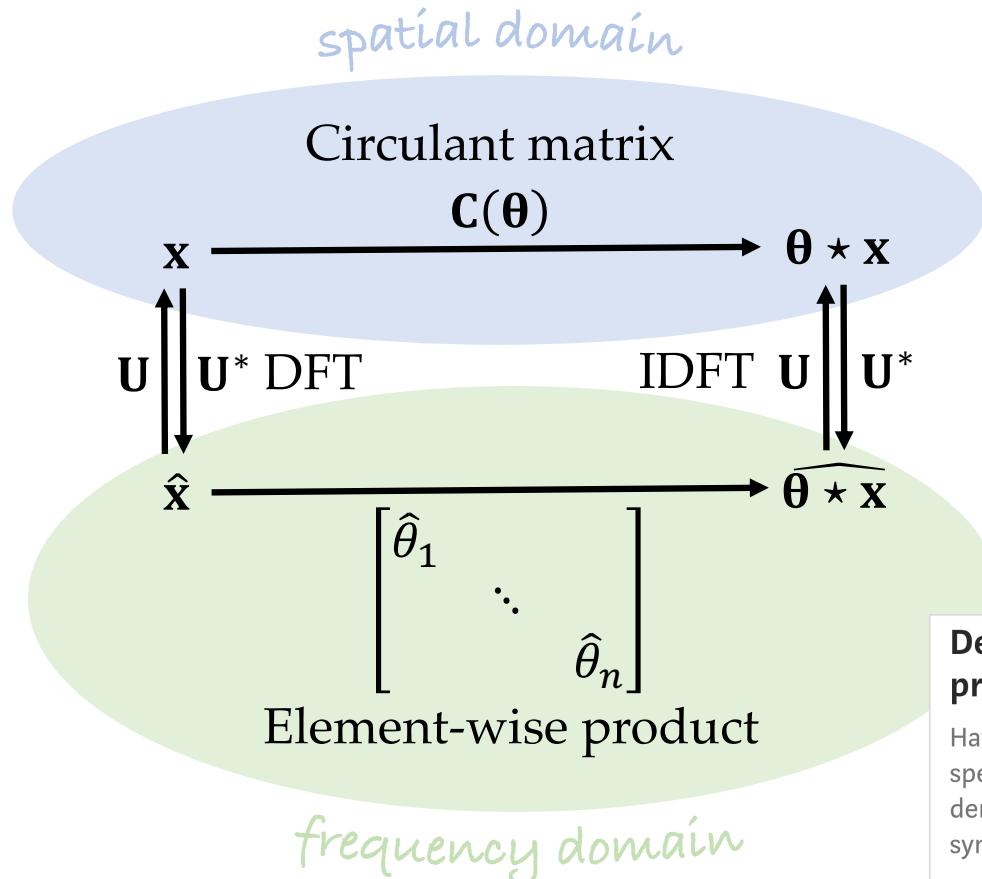
Fourier basis

$$\mathbf{u}_k = \frac{1}{\sqrt{n}} \begin{bmatrix} 1 \\ e^{i\frac{2\pi}{n}k} \\ \vdots \\ e^{i\frac{2\pi}{n}(n-1)k} \end{bmatrix}$$

$$\begin{bmatrix} b \\ a \\ a \\ c \end{bmatrix} = \begin{bmatrix} | & & & | \\ \mathbf{u}_1 & \dots & & \mathbf{u}_n \\ | & & & | \end{bmatrix} \begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \vdots \\ \hat{\theta}_n \end{bmatrix}$$

Fourier transform
 $\hat{\theta} = \mathbf{U}^* \theta$

commuting matrices are jointly
diagonalisable by Fourier Transform



Deriving convolution from first principles

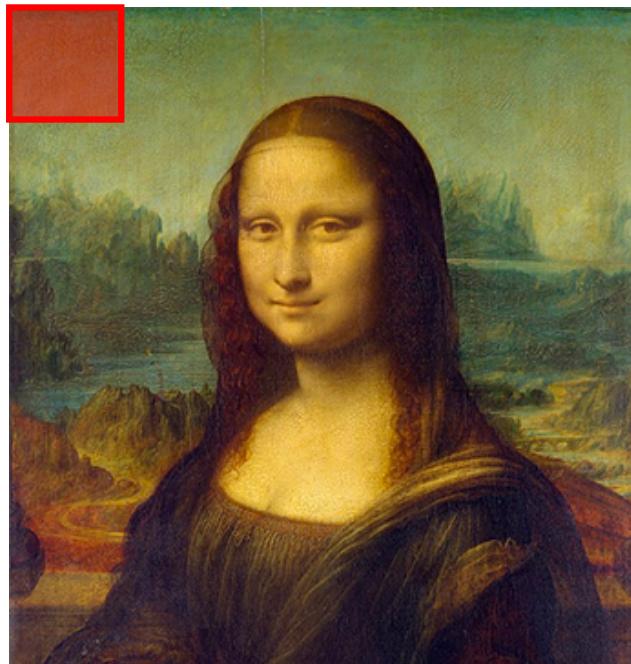
Have you even wondered what is so special about convolution? I show how to derive the convolution from translational symmetry



Michael Bronstein
Jul 26, 2020 · 9 min read ★

GROUPS

Convolution, revisited



Convolution, revisited

convolution = matching shifted filter

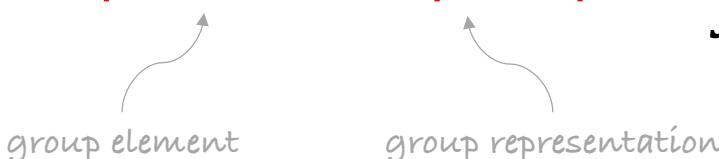
$$(x \star \psi)(u) = \langle x, T_u \psi \rangle = \int_{-\infty}^{+\infty} x(v) \psi(u - v) dv$$


A diagram illustrating the components of the convolution formula. Two curved arrows point from the labels "shift vector" and "shift operator" to the corresponding terms in the equation. The "shift vector" arrow points to the variable u in the term $T_u \psi$. The "shift operator" arrow points to the term $u - v$ in the integral expression.

domain Ω = symmetry group \mathfrak{G}

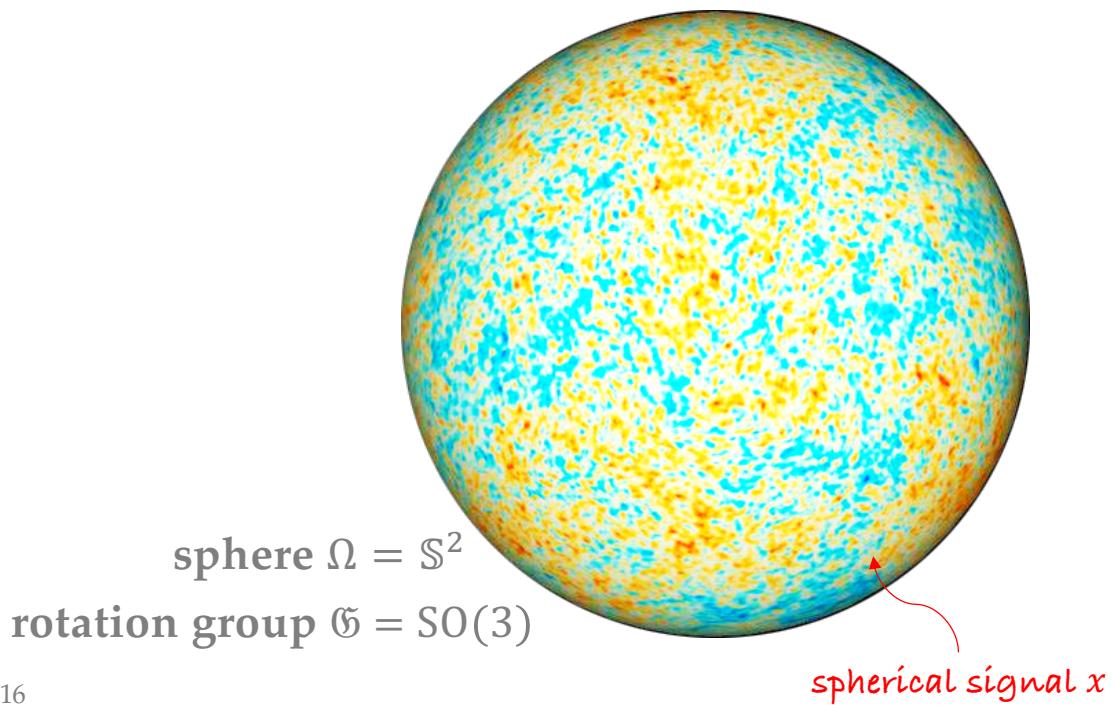
Group Convolution

convolution = matching transformed filter

$$(x \star \psi)(g) = \langle x, \rho(g)\psi \rangle = \int_{\Omega} x(v)\psi(g^{-1}v)dv$$


group element group representation

Convolution on the Sphere



Cohen, Welling 2016

Convolution on the Sphere

$$(x \star \psi)(R) = \int_{\mathbb{S}^2} x(u)\psi(R^{-1}u)du$$

signal on $SO(3)$

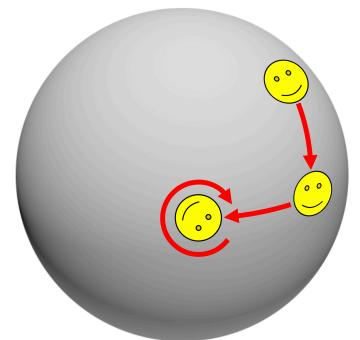
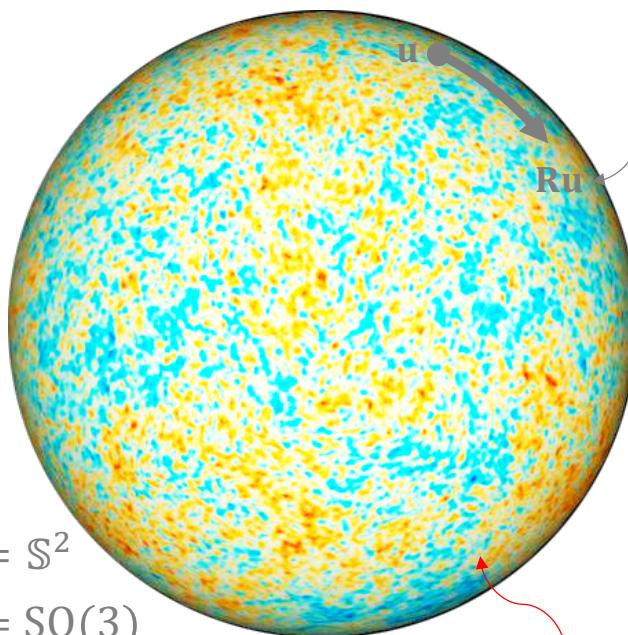
sphere $\Omega = \mathbb{S}^2$

rotation group $G = SO(3)$

spherical signal x

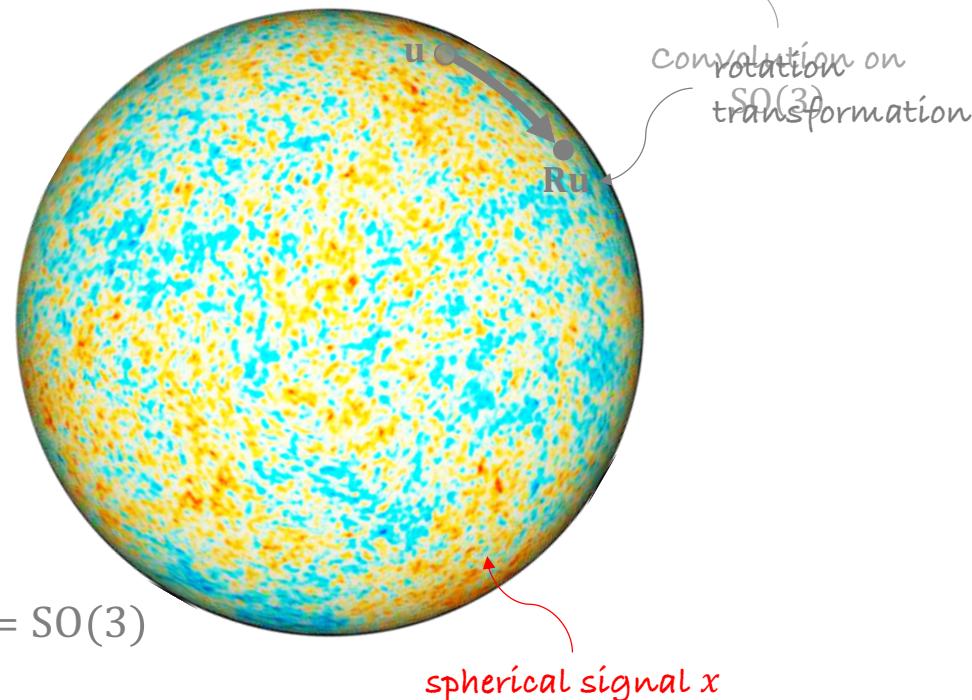
rotation transformation

Cohen, Welling 2016



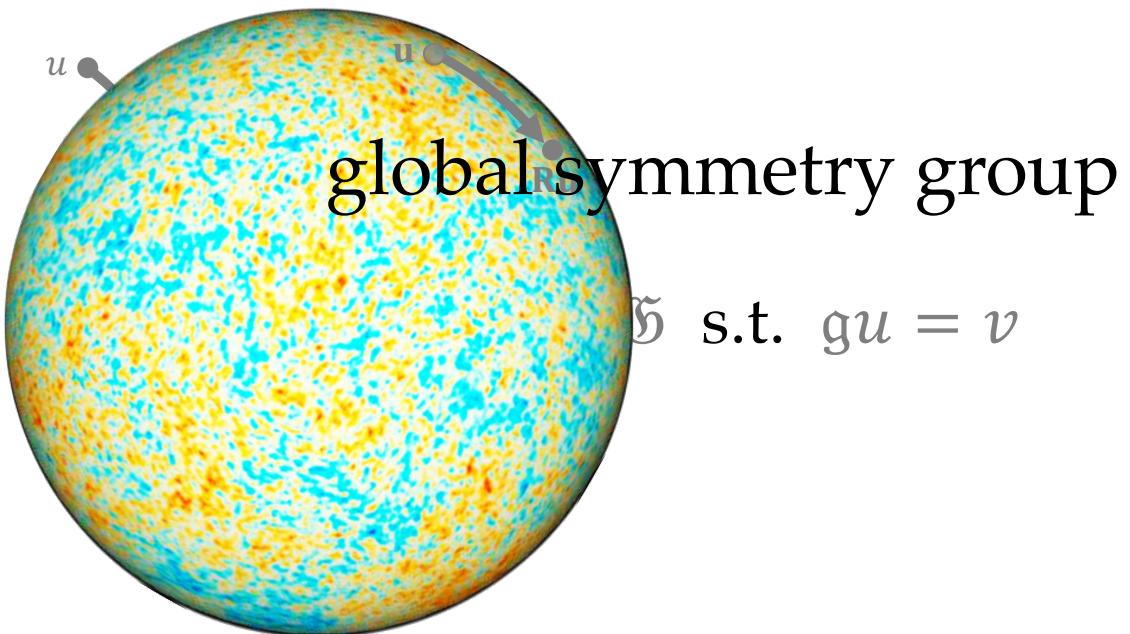
Convolution on the Sphere

$$((x \star \psi) \star \phi)(R) = \int_{SO(3)} (x \star \psi)(Q) \phi(R^{-1}Q)dQ$$



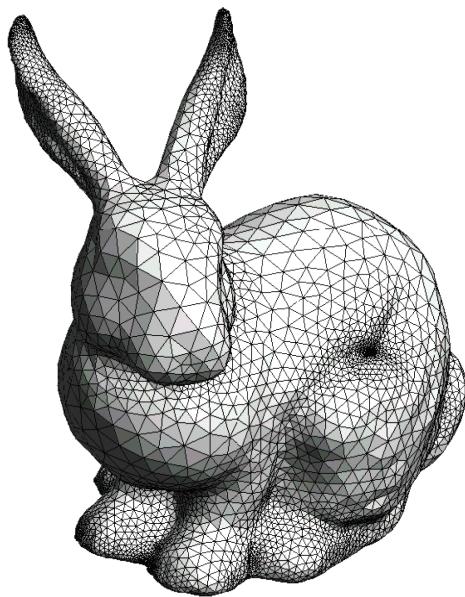
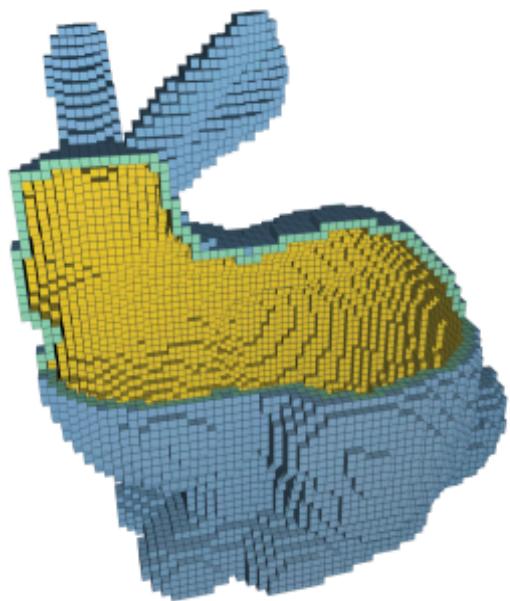
Cohen, Welling 2016

Homogeneous Spaces



MANIFOLDS & MESHES

Why Manifolds?



More efficient representation: no “waste”
for internal structures

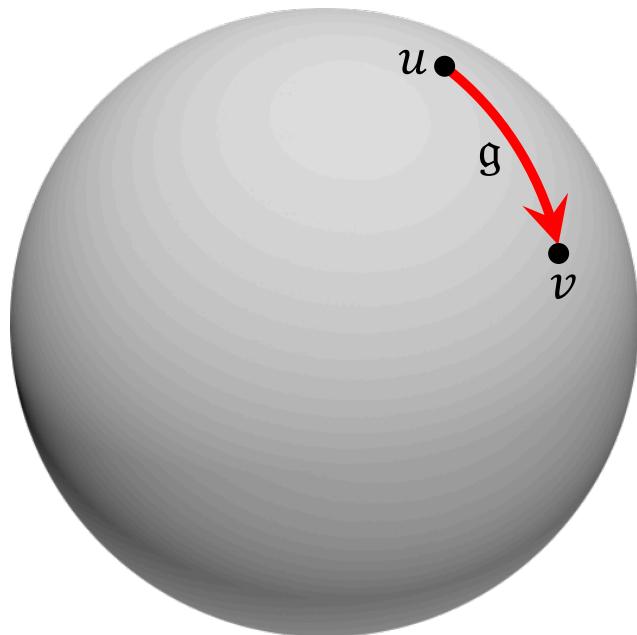
Natural model for
deformable shapes

Why Manifolds?

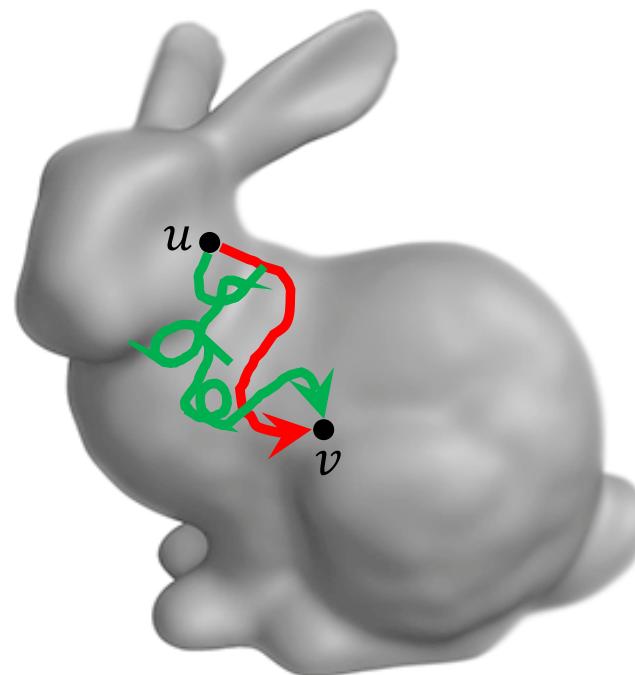
In protein modeling,
abstract out internal
structure that is irrelevant
for interactions + allow
some conformation
changes



Homogeneous Spaces

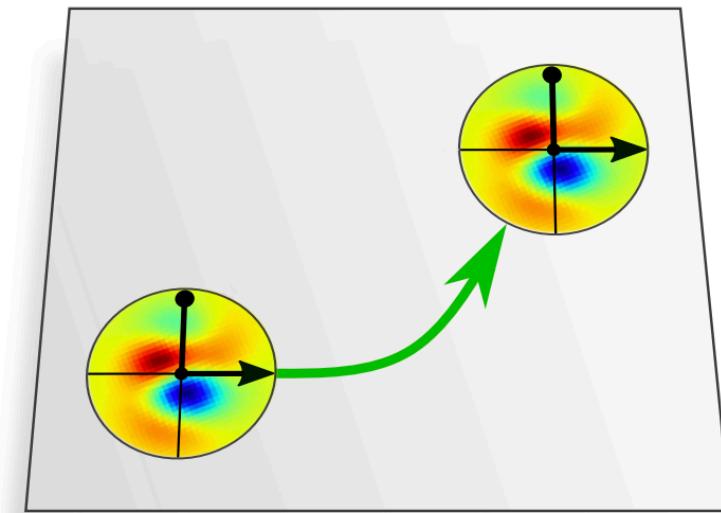
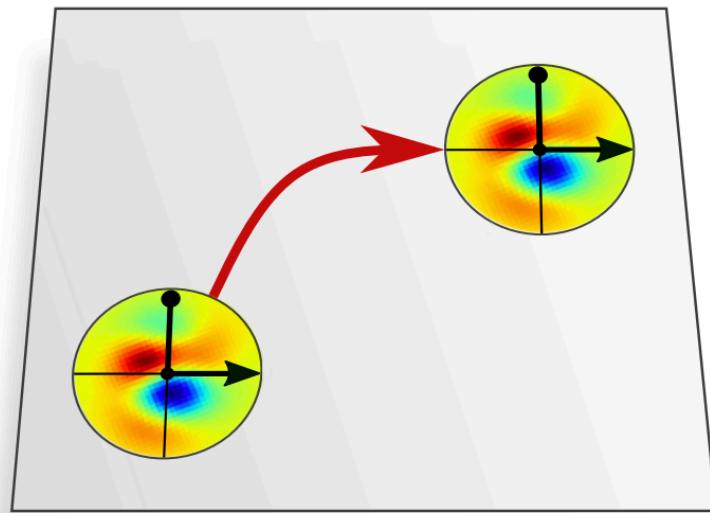


global symmetry group
 $\exists g \in \mathfrak{G} \text{ s.t. } gu = v$



no useful global
symmetry group

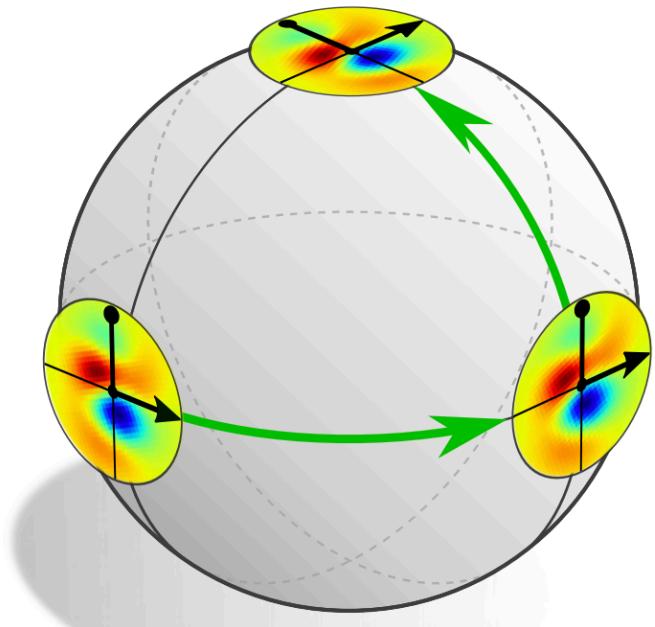
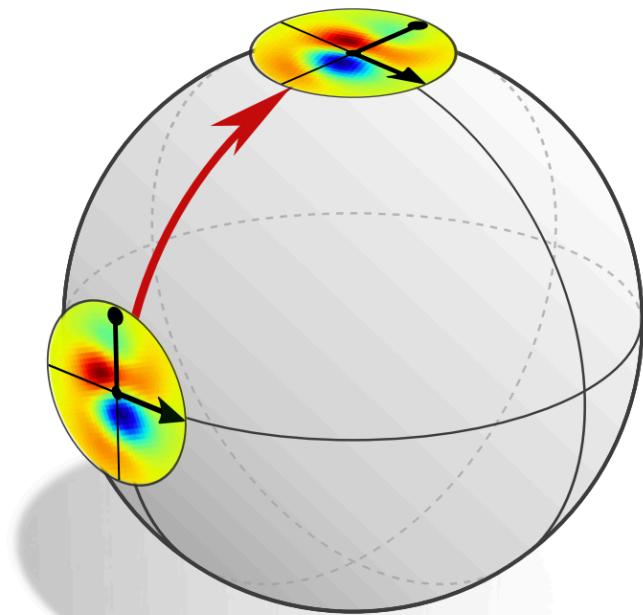
Euclidean Convolution



Euclidean space: Transport the filter around the domain

Figure: Weiler et al. 2021

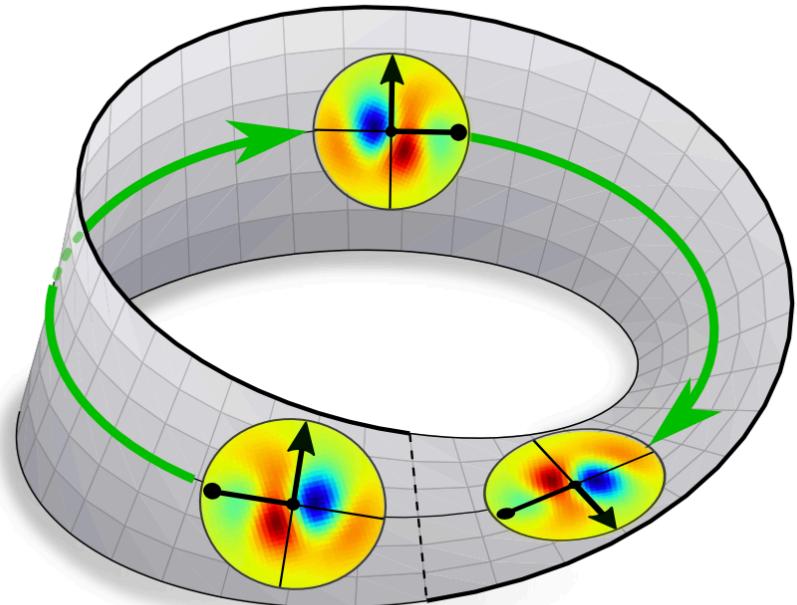
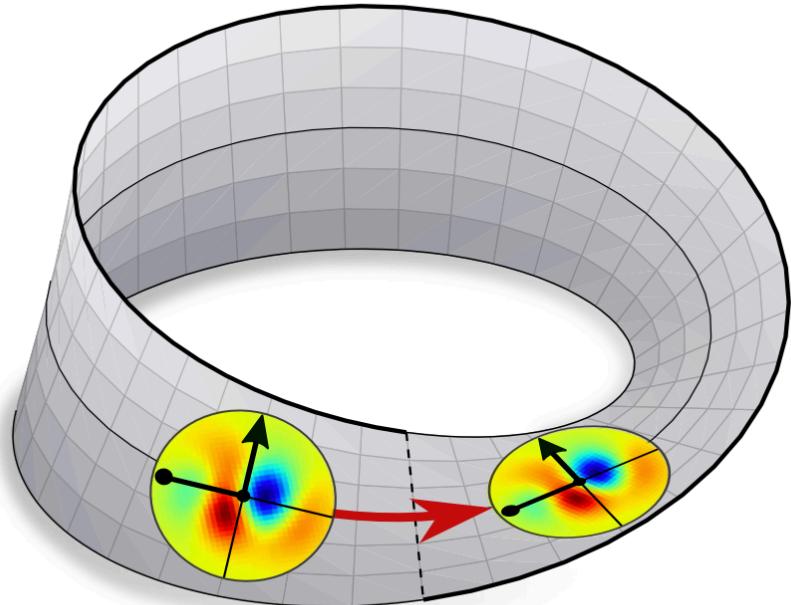
Non-Euclidean Convolution



Manifold: Result of transport is *path dependent*

Figure: Weiler et al. 2021

Non-Euclidean Convolution



Manifold: Result of transport is *path dependent*

Figure: Weiler et al. 2021

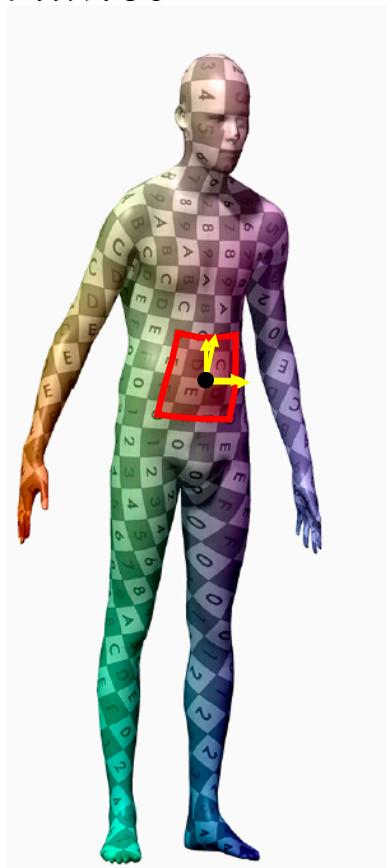
Two Types of Invariance

**Local gauge
transformation**

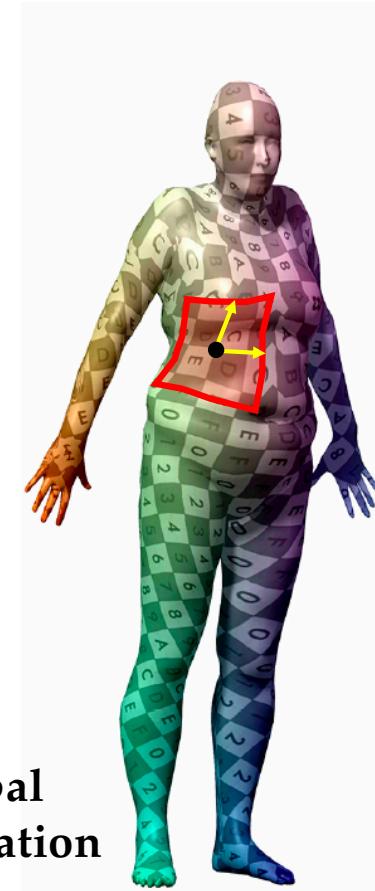


Two Types of Invariance

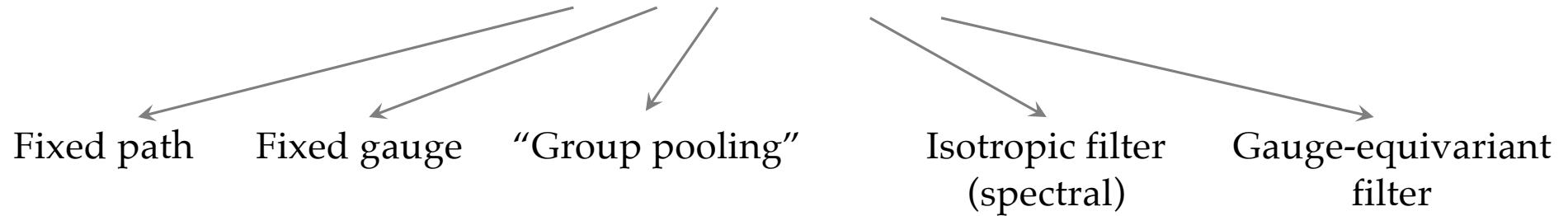
**Local gauge
transformation**



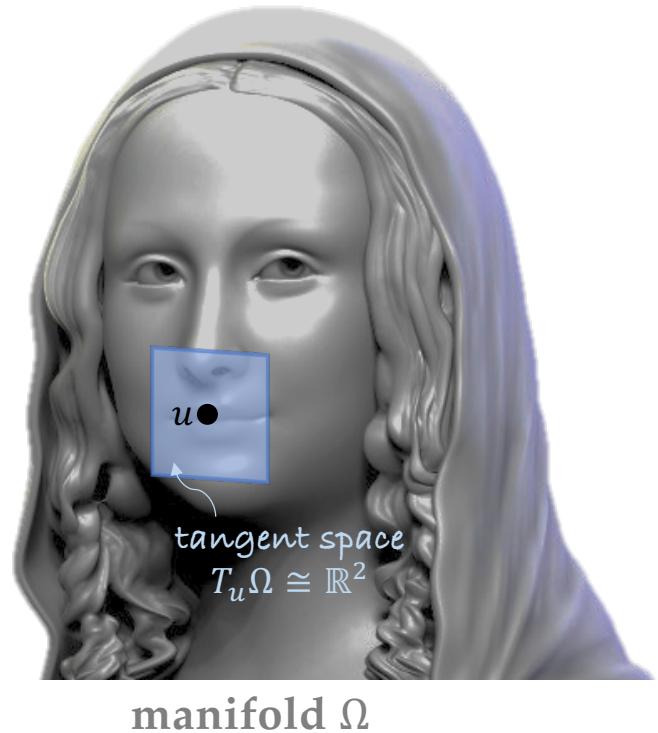
**Global
deformation**



Non-Euclidean Convolution Recipes



Manifolds



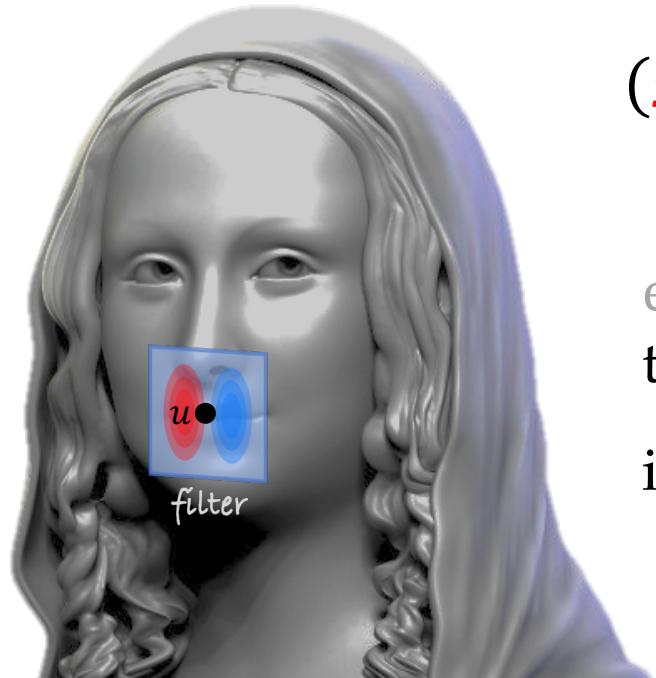
manifold = locally Euclidean space

Riemannian metric = local length/ direction

Intrinsic quantity = expressed solely in terms of the Riemannian metric

Isometry = metric-preserving deformation

Geodesic CNNs



manifold Ω
isometry group $\text{Iso}(\Omega)$

Masci et al. 2015; Boscaini et al. 2016; Monti et al. 2017

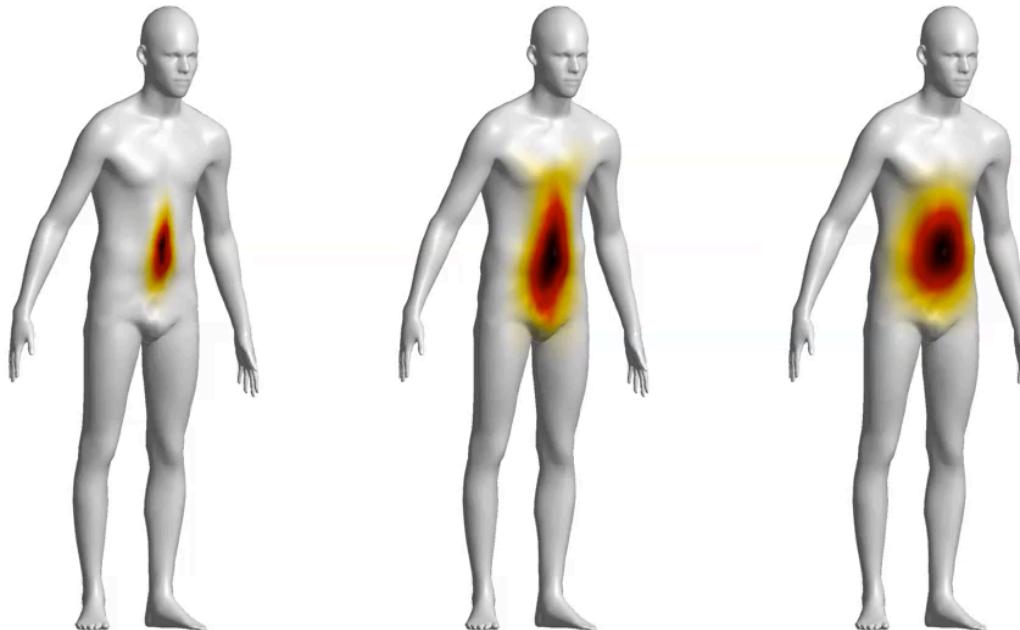
$$(\mathbf{x} \star \psi)(u) = \int_{T_u \Omega} \psi(v) \mathbf{x}(\exp_u v) dv$$

Exponential map
 $\exp_u: T_u \Omega \rightarrow \Omega$

\exp_u is an intrinsic map allowing to express the signal \mathbf{x} locally in the tangent space $T_u \Omega$

intrinsic filter = invariant to isometries

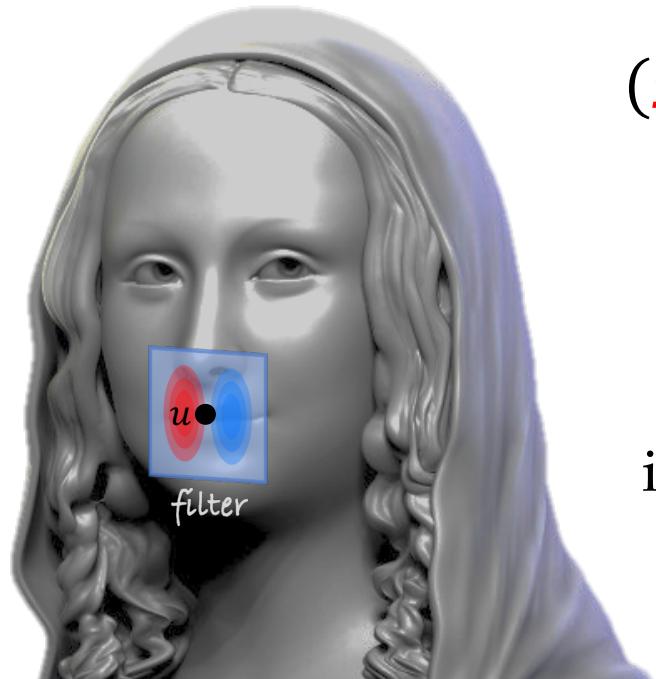
Geodesic CNNs



Anisotropic intrinsic filters on a manifold

Masci et al. 2015; Boscaini et al. 2016; Monti et al. 2017

Geodesic CNNs



manifold Ω
isometry group $\text{Iso}(\Omega)$

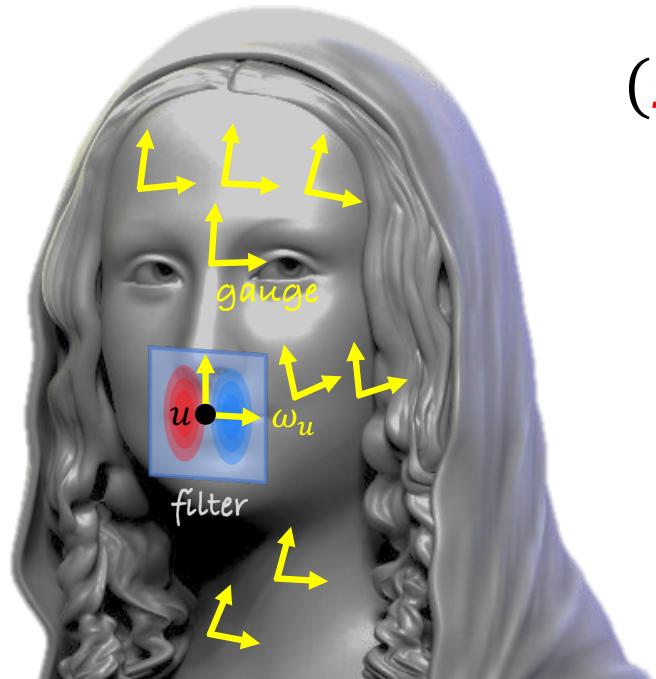
$$(x \star \psi)(u) = \int_{T_u \Omega} \psi(v) x(\exp_u v) dv$$

Problem: these are abstract vectors!

intrinsic filter = invariant to isometries

Masci et al. 2015; Boscaini et al. 2016; Monti et al. 2017

Geodesic CNNs

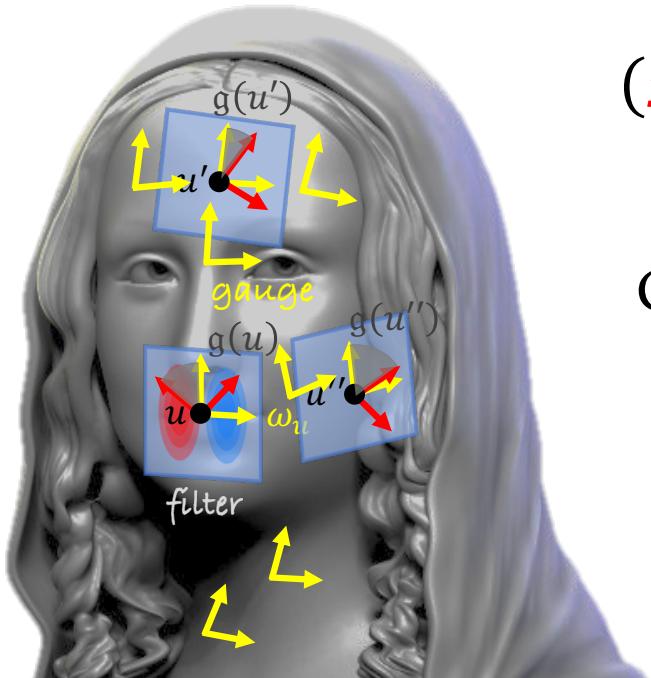


manifold Ω
isometry group $\text{Iso}(\Omega)$

$$(x \star \psi)(u) = \int_{\mathbb{R}^2} \psi(v) x(\exp_u \omega_u v) dv$$

local reference frame
 $\omega_u: \mathbb{R}^2 \rightarrow T_u \Omega$

Gauge Transformations



manifold Ω
structure group \mathfrak{G}

$$(x \star \psi)(u) = \int_{\mathbb{R}^2} \psi(v) x(\exp_u \omega_u v) dv$$

Gauge defined up to gauge transformation

$$g: \Omega \rightarrow \mathfrak{G}$$

Cohen et al. 2019; Weiler et al. 2021

Structure Group

A gauge is defined up to a *gauge transformation* $g: \Omega \rightarrow \mathfrak{G}$

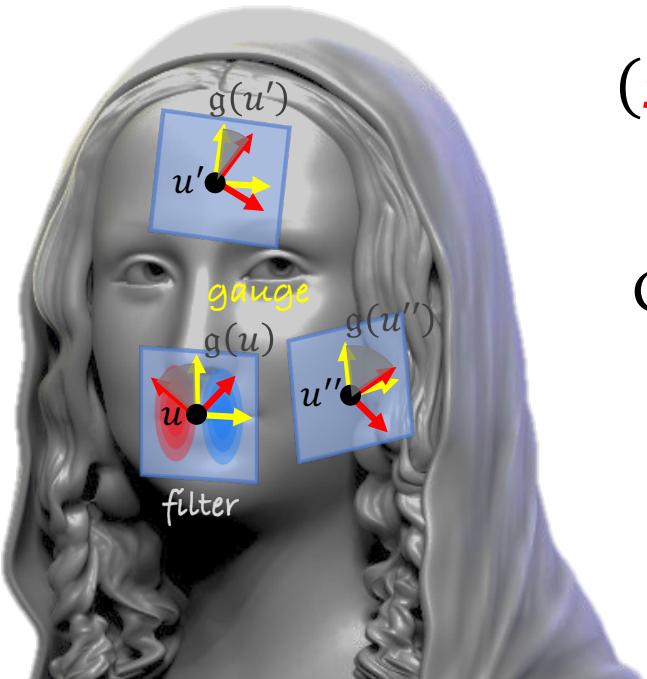
- | | | |
|-------------------------------|-----------------|-------------------------------------|
| • “Naked” manifold | $GL(s)$ | invertible matrices |
| • Manifold+orientation | $GL^+(s)$ | invertible matrices with $\det > 0$ |
| • Manifold+volume | $SL(s)$ | matrices with $\det = 1$ |
| • Manifold+metric | $O(s)$ | orthogonal matrices |
| • Manifold+metric+orientation | $SO(s)$ | orthogonal matrices with $\det = 1$ |
| • Manifold+frame field | $\{\text{id}\}$ | identity (no ambiguity) |

Structure Group

A gauge is defined up to a *gauge transformation* $g: \Omega \rightarrow \mathfrak{G}$

• “Naked” manifold	$GL(s)$	invertible matrices
• Manifold+orientation	$GL^+(s)$	invertible matrices with $\det > 0$
• Manifold+volume	$SL(s)$	matrices with $\det = 1$
• Manifold+metric	$O(s)$	orthogonal matrices
• Manifold+metric+orientation	$SO(s)$	orthogonal matrices with $\det = 1$
• Manifold+frame field	$\{\text{id}\}$	identity (no ambiguity)

Gauge-equivariant CNNs



manifold Ω
structure group $\mathfrak{G} = \text{SO}(2)$

Cohen et al. 2019

$$(x \star \psi)(u) = \int_{\mathbb{R}^2} \psi(v) \rho(\exp_u \omega_u v) dv$$

Gauge defined up to gauge transformation

$$g: \Omega \rightarrow \text{SO}(2)$$

gauge-equivariant filter

$$\psi(g^{-1}v) = \rho(g^{-1})\psi(v)\rho(g)$$

“Hairy Ball” (a.k.a. Poincaré-Hopf) Theorem

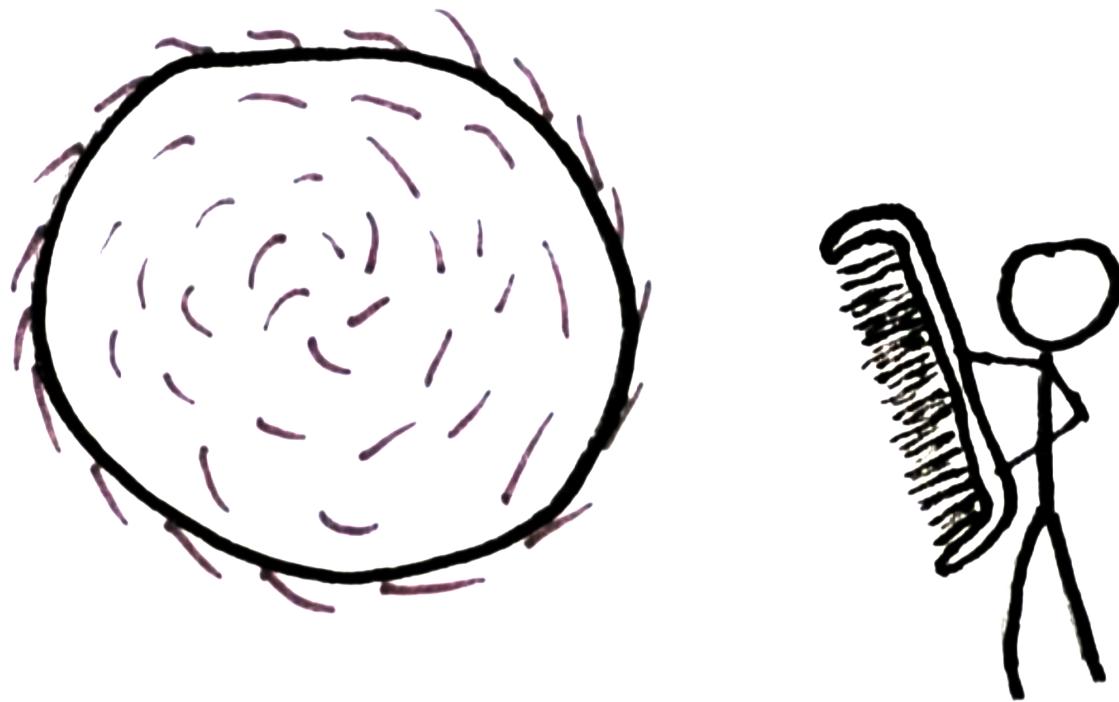
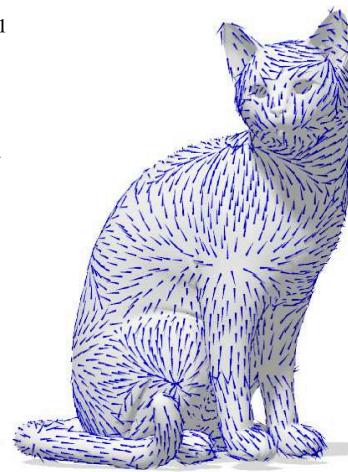
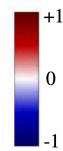
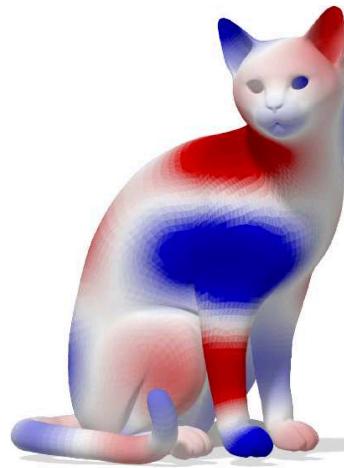
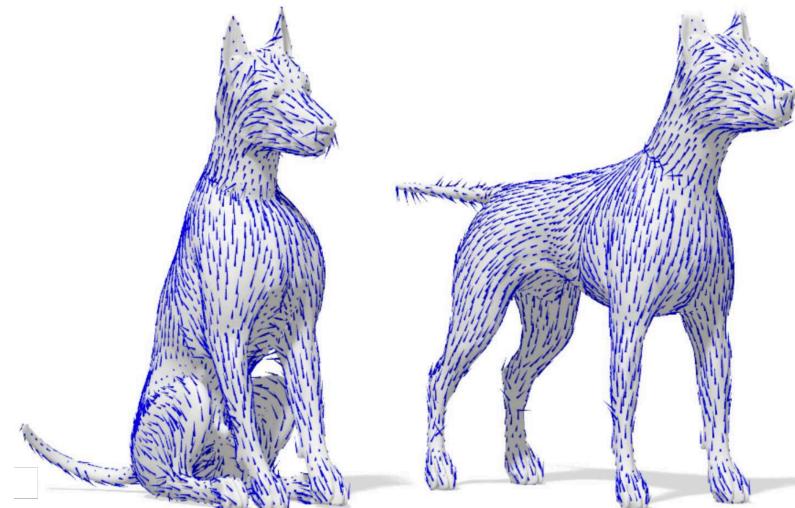


Image: Minutephysics

Theory vs Practice: Stable Gauges



Gradient of intrinsic function

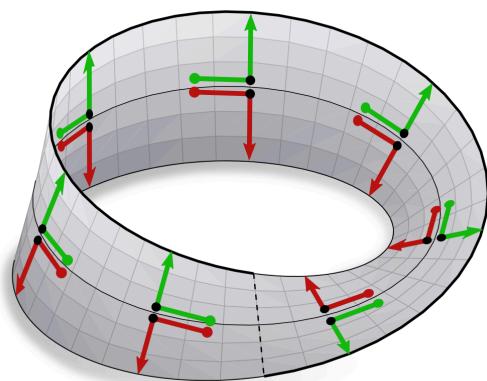


Deformation-invariant
stable gauge

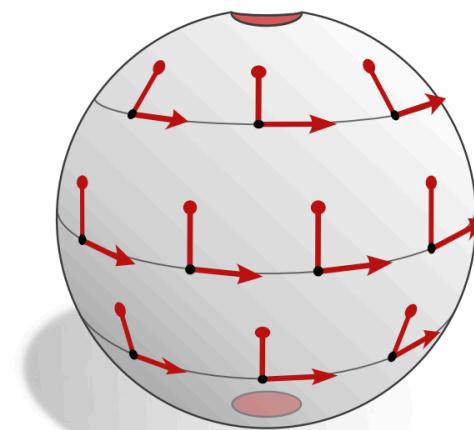
Structure Group



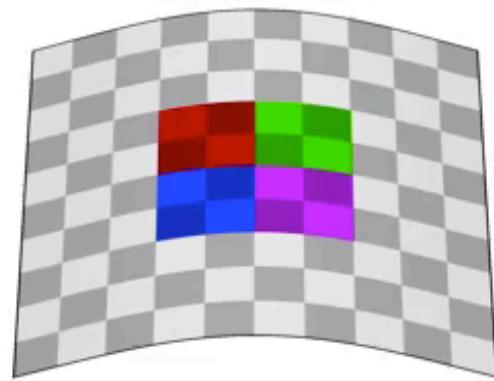
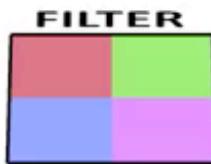
rotation
 $SO(2)$



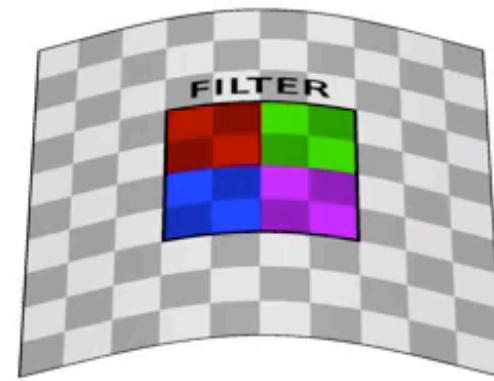
reflection
 R



fixed gauge
 $\{id\}$



**Euclidean (extrinsic)
convolution**



**Geometric (intrinsic)
convolution**

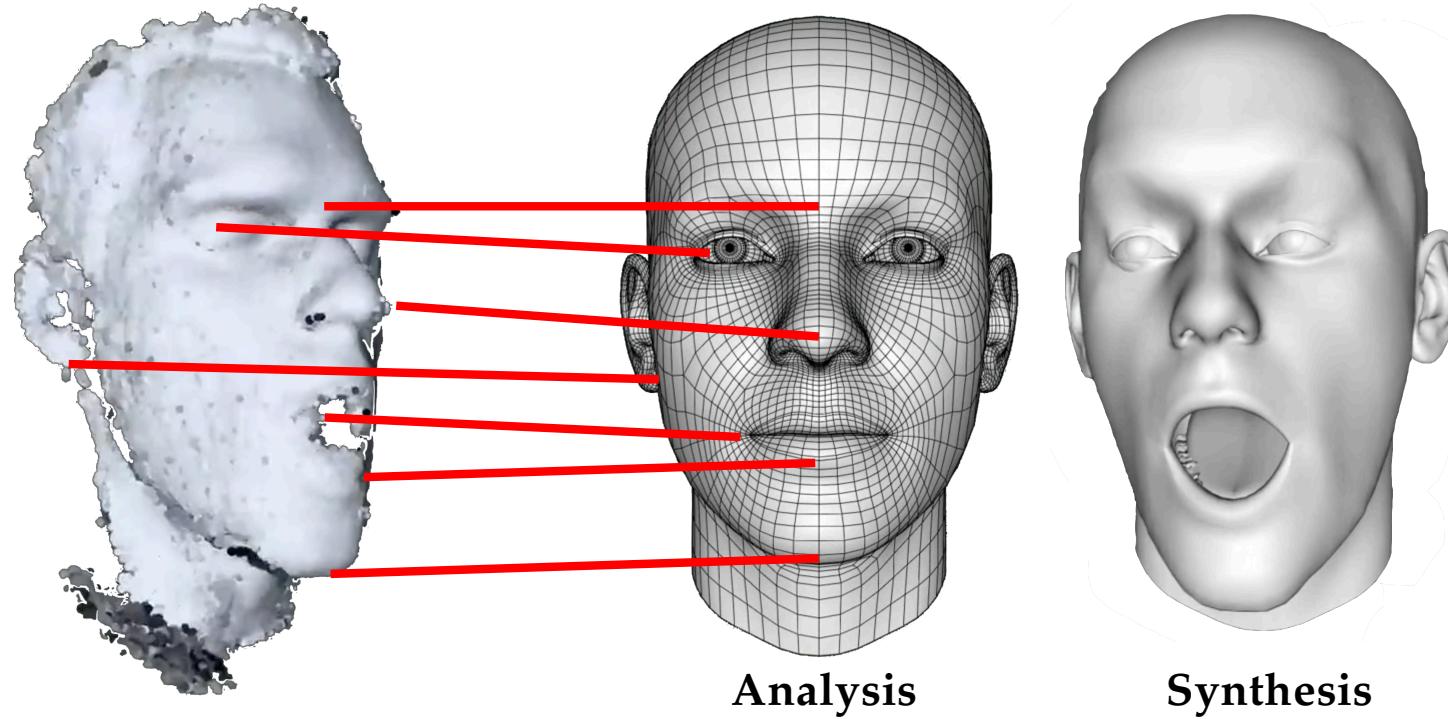
FaceShift 2015



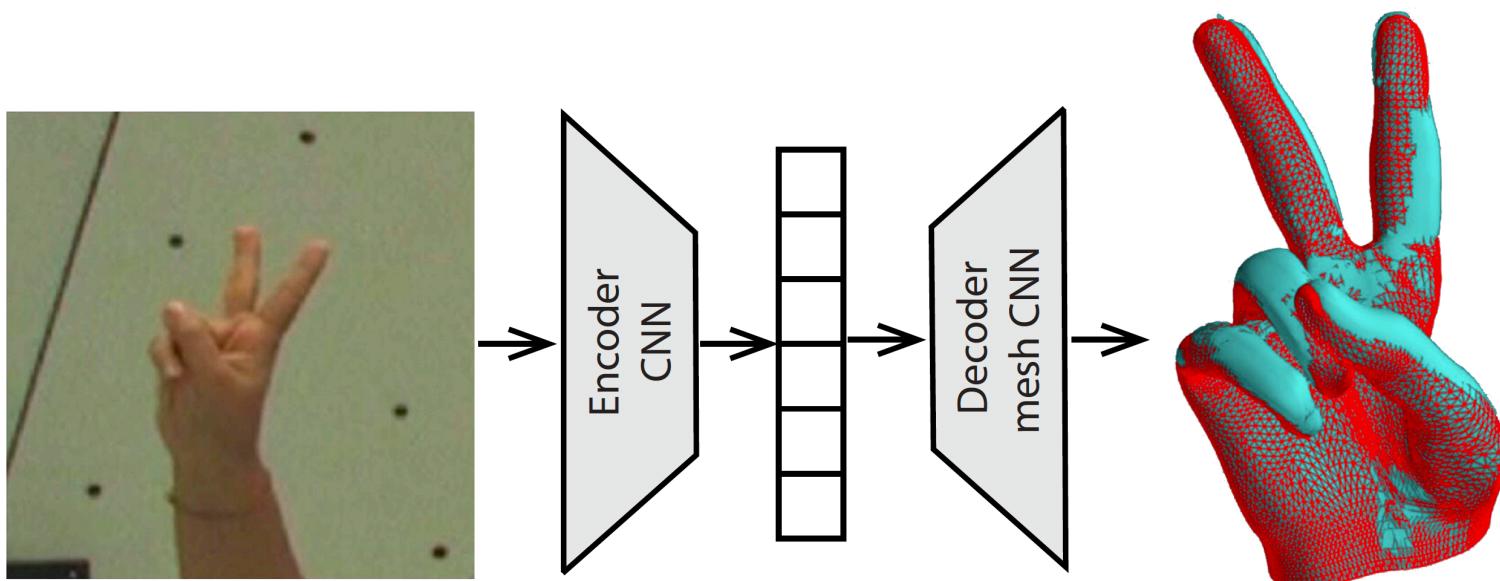
GDC

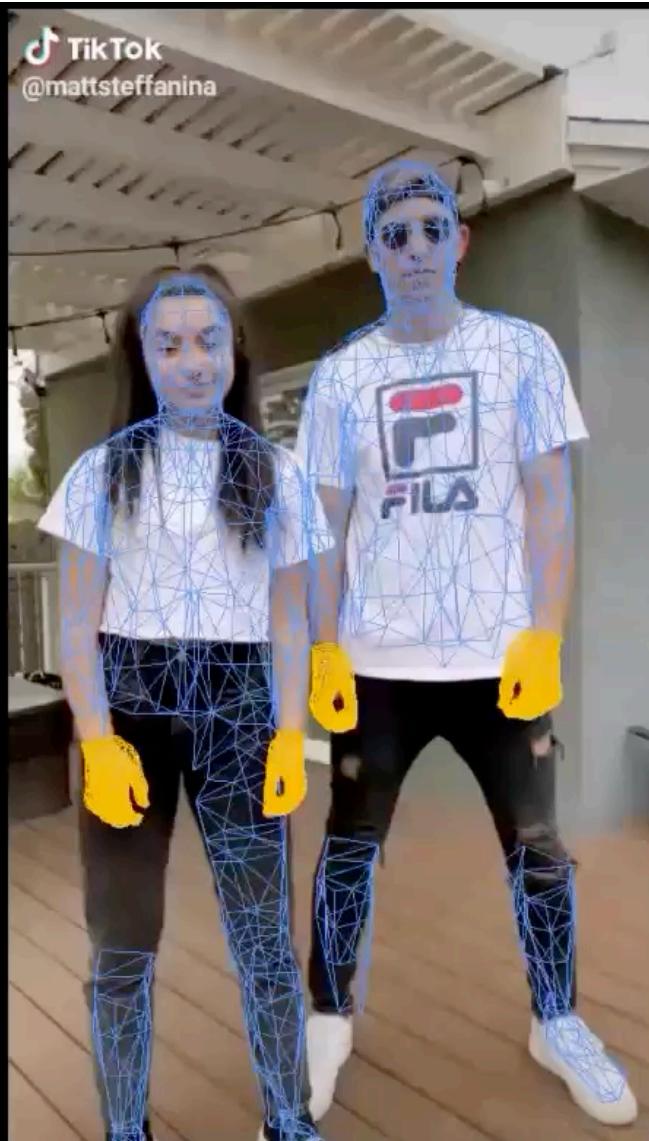


Shape Analysis & Synthesis



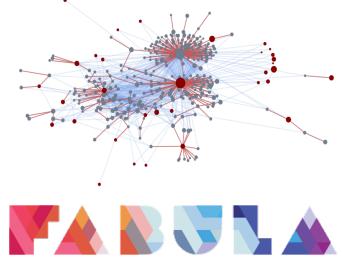
3D Hand Reconstruction



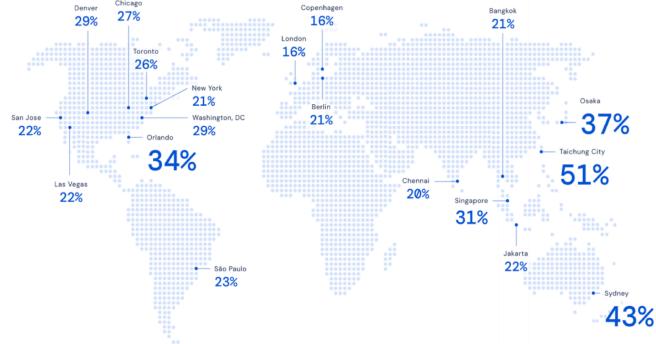


Snap Acquires Ariel AI To Enhance AR Features

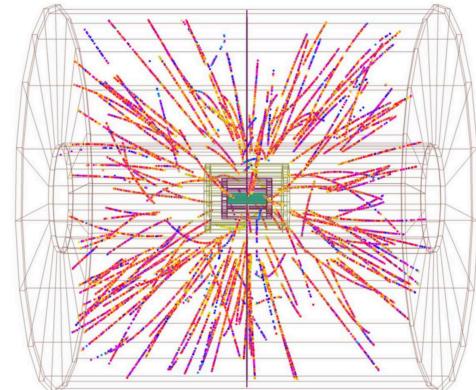




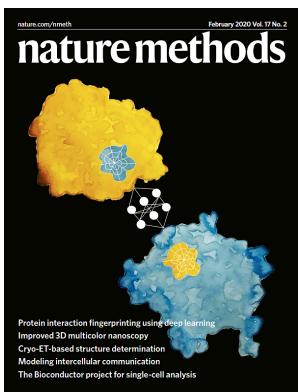
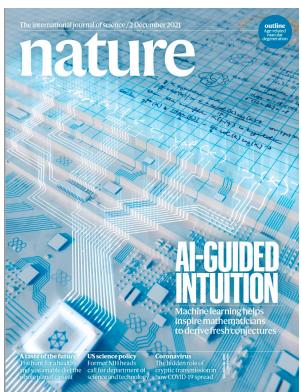
Fake news detection



Navigation



Particle physics



Pure math

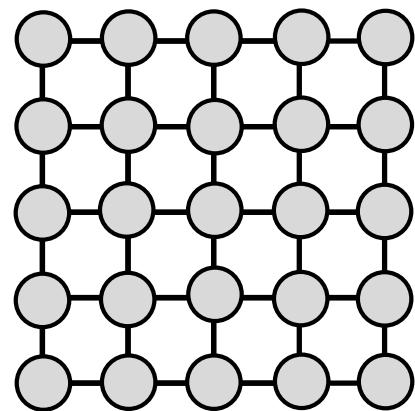
Structural biology



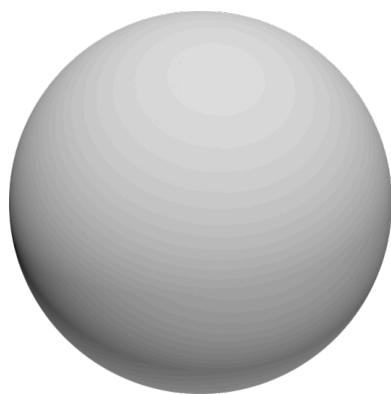
Drug discovery

WRAP UP

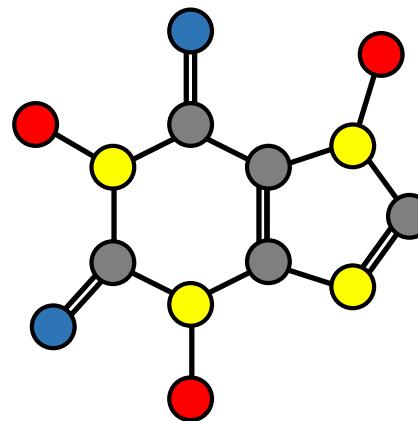
The “5G” of Geometric Deep Learning



Grids



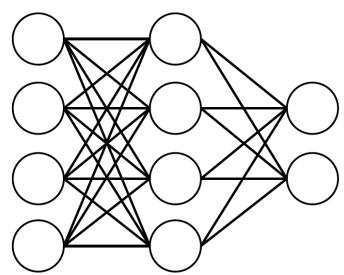
Groups



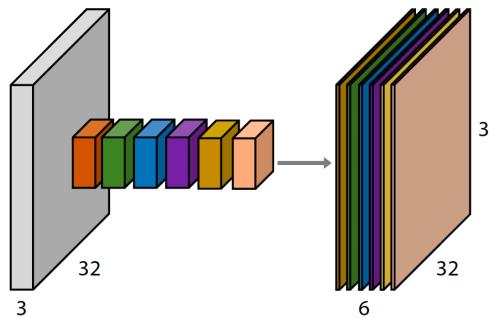
Graphs



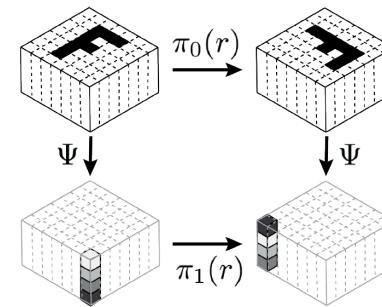
Geodesics &
Gauges



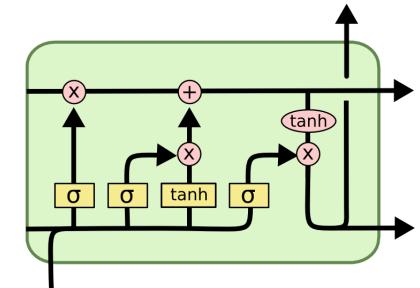
Perceptrons
Function regularity



CNNs
Translation



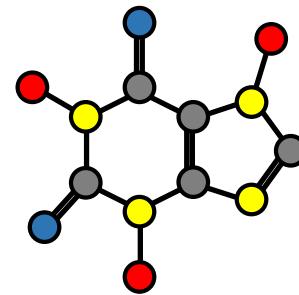
Group-CNNs
Translation+Rotation,
Global groups



LSTMs
Time warping



DeepSets / Transformers
Permutation



GNNs
Permutation



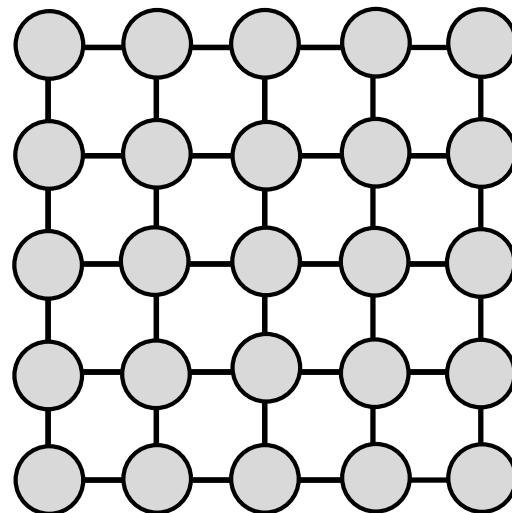
Intrinsic CNNs
Isometry / Gauge choice

“The knowledge of certain principles easily compensates the lack of knowledge of certain facts”

—Claude Adrien Helvétius

PHYSICS-INSPIRED GDL

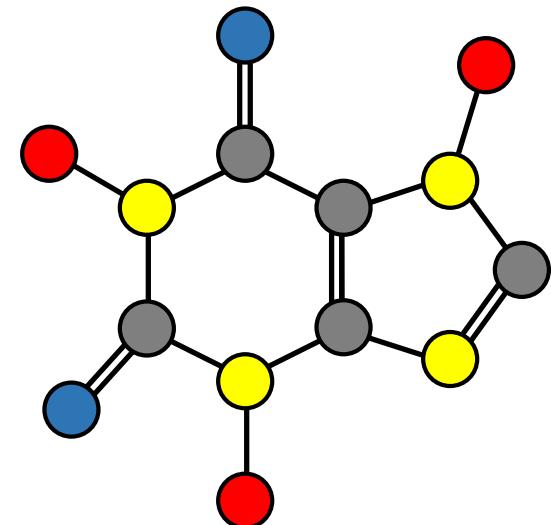
Instances of GDL Blueprint: Different domains



Grid



Mesh



Graph

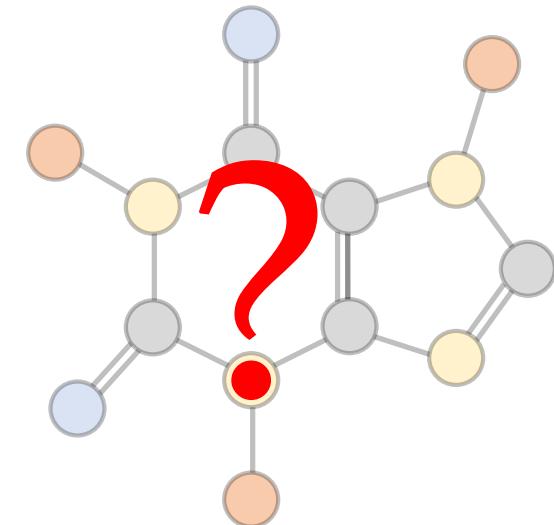
Instances of GDL Blueprint: Different domains



Plane (Homogeneous space)



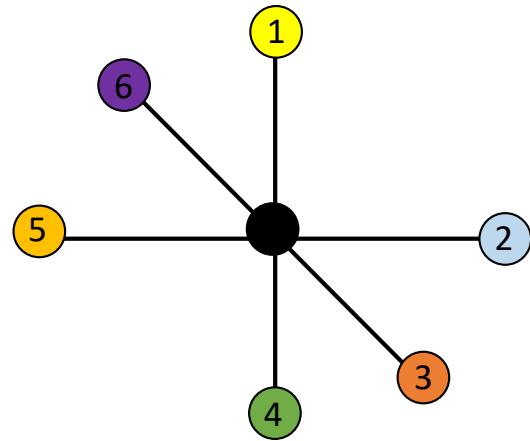
Manifold



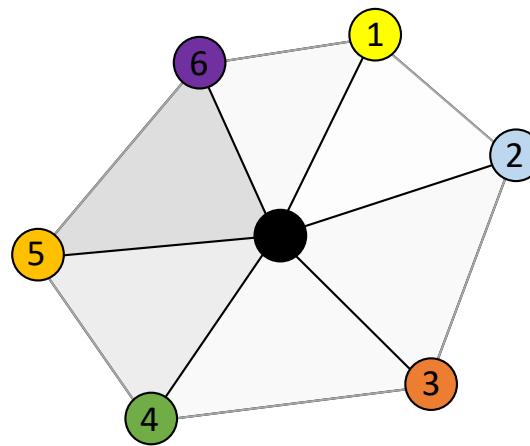
Graph

1. Continuous model for graphs?

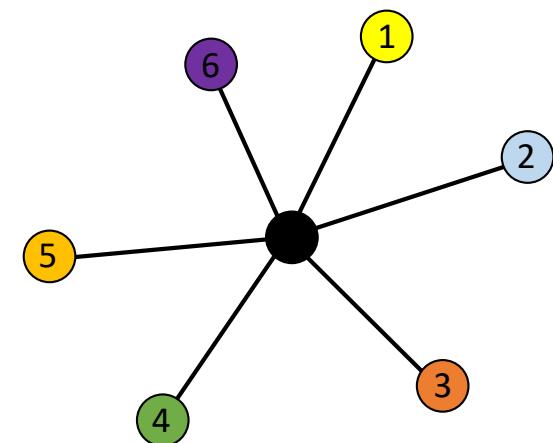
Graphs vs Meshes vs Grids



Grid

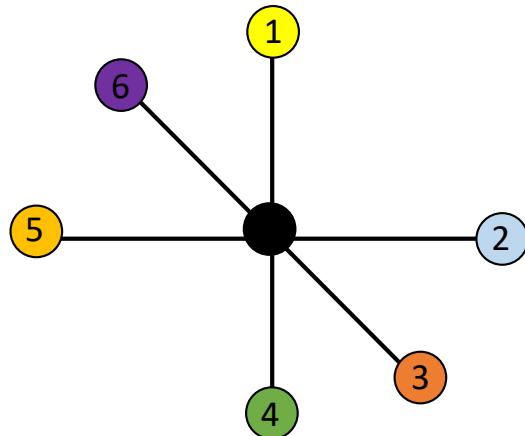


Mesh



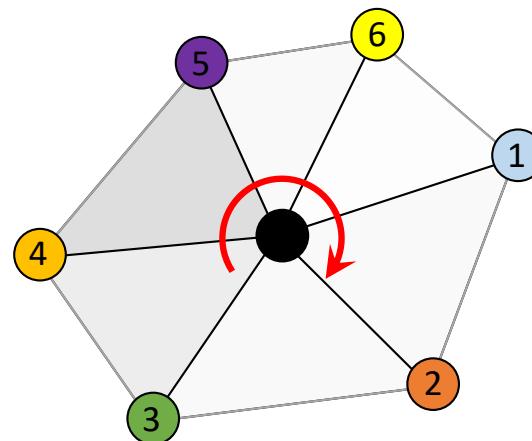
Graph

Graphs vs Meshes vs Grids



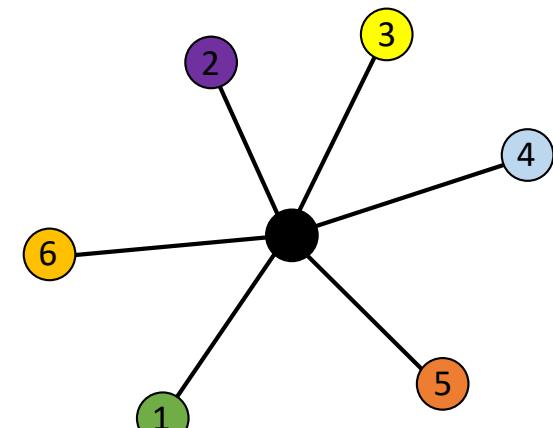
Grid

Fixed



Mesh

Rotation



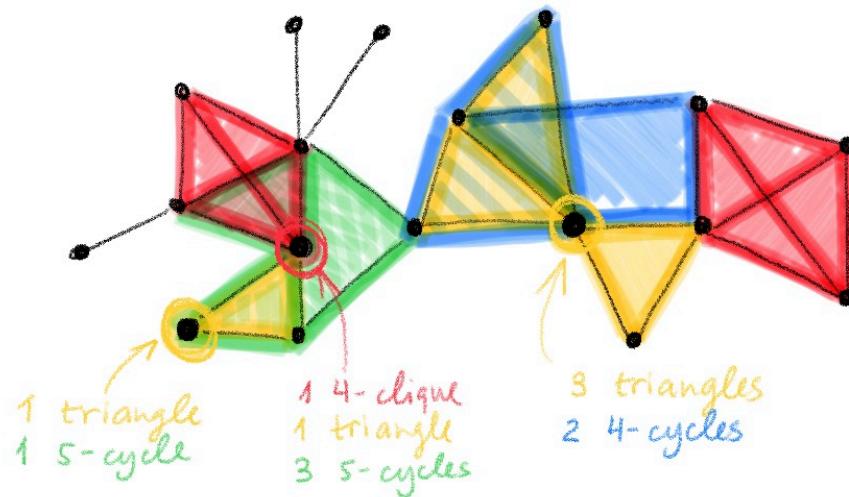
Graph

Permutation

Graphs have the least structure

Positional Encoding Approaches

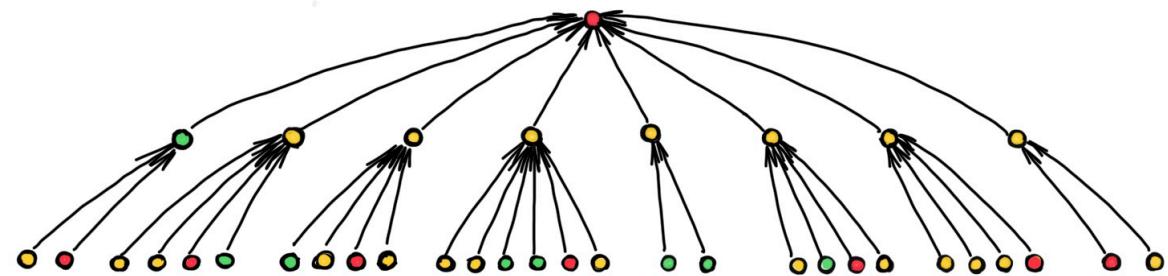
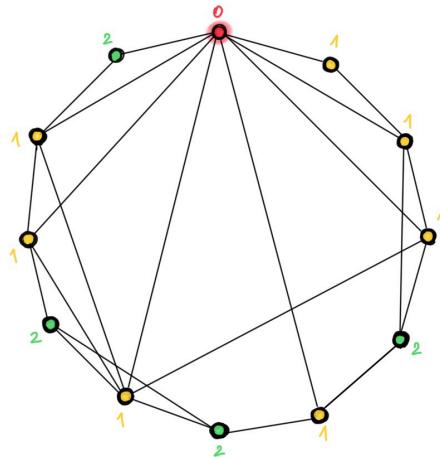
- Random node features¹
- Graph Laplacian eigenvectors²
- Graph substructure counts³
- Bags of subgraphs⁴



2. How to choose positional encoding?

¹Sato et al. 2020; ²Vaswani et al. 2017; Qiu et al. 2020; Dwivedi et al. 2020; ³Bouritsas, Frasca, et B. 2020; ⁴Bevilacqua, Frasca, Lim, et B., Maron 2021

Over-squashing & Bottlenecks



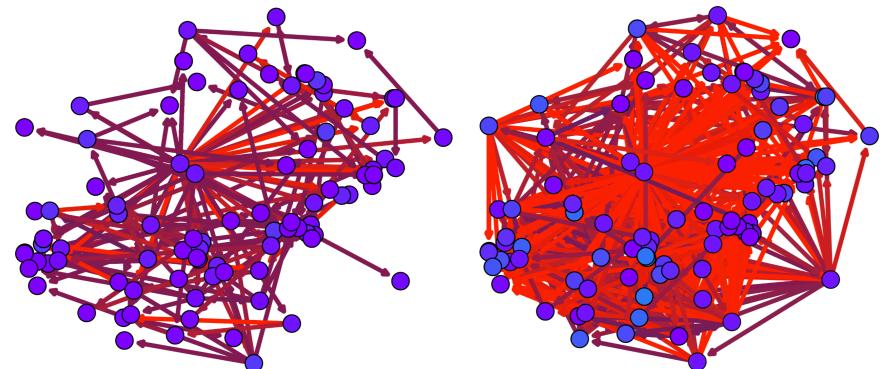
In small-world graphs metric ball volume $\text{vol}(B_k) = \sum_{j \in B_k} d_j$
grows exponentially with ball radius k

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

Graph Rewiring

Decouple **input graph** from **information propagation graph** (at the expense of link to WL)

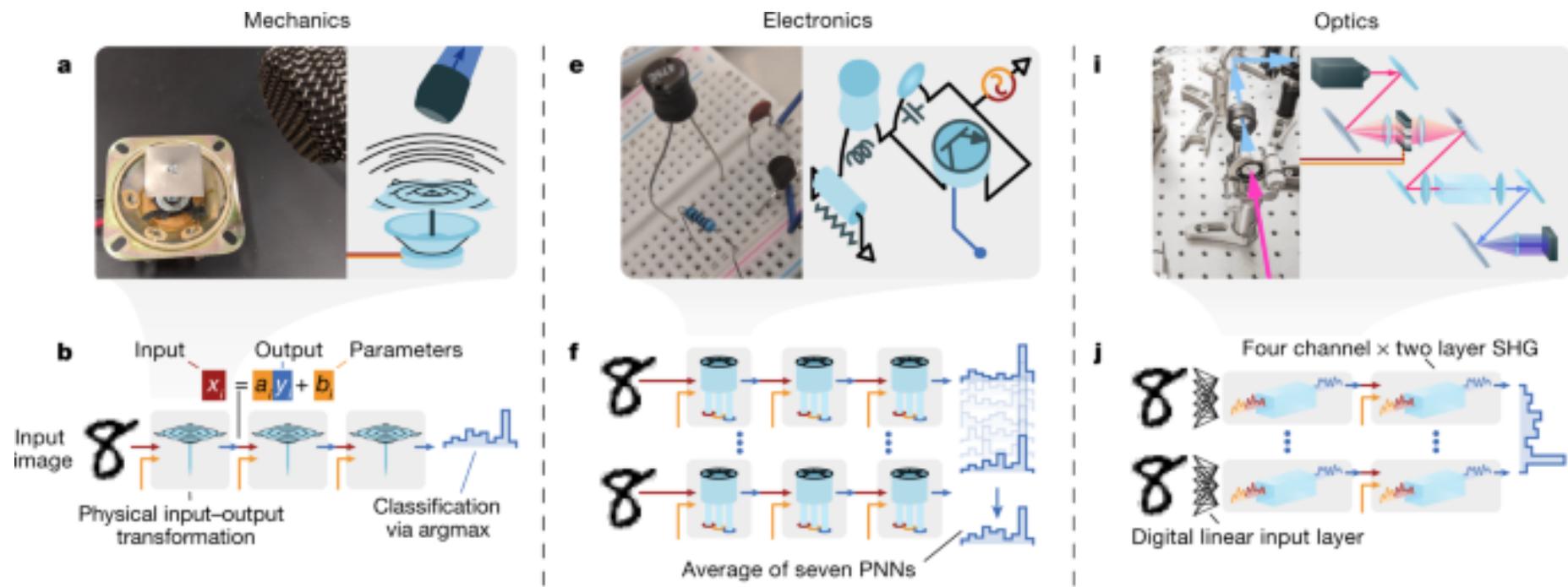
- Neighbourhood sampling (GraphSAGE)¹
- Multi-hop filters (SIGN)²
- Complete graph³
- Topology diffusion (DIGL)⁴
- Learnable graph (Dynamic Graph CNN)⁵



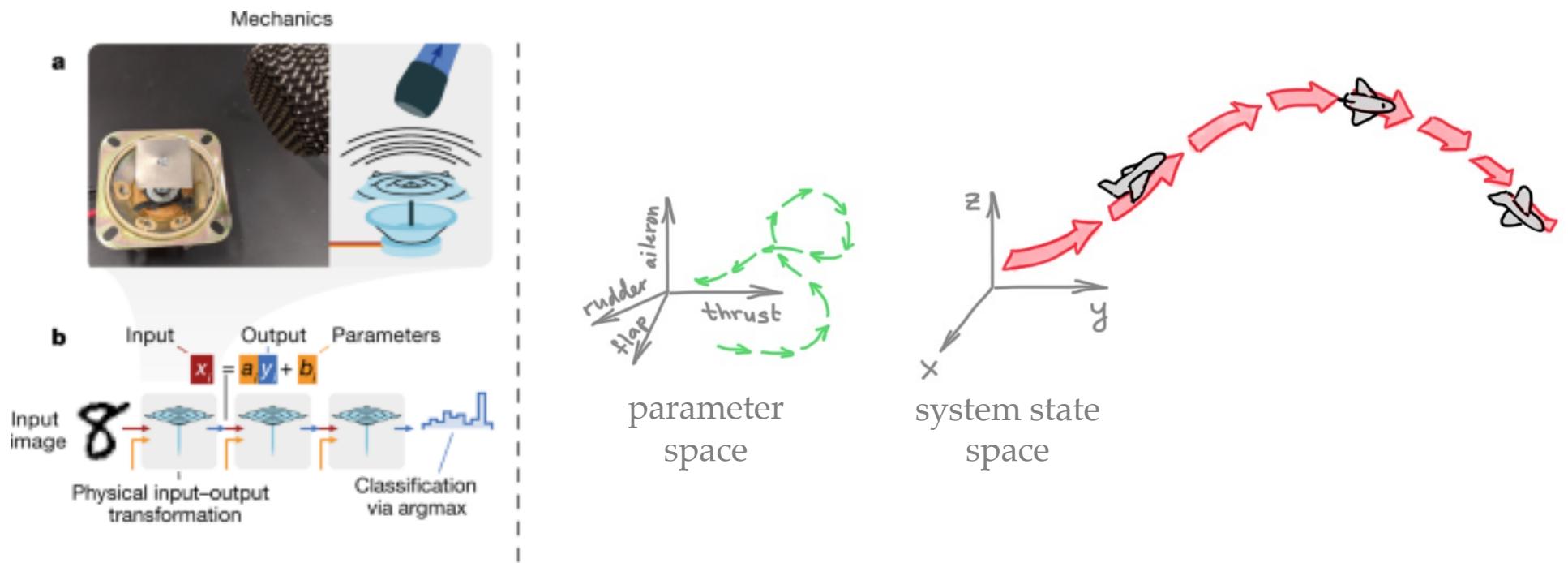
3. How to rewire the graph?

¹Hamilton et al. 2017; ²Rossi, Frasca, et B. 2020; ³Alon, Yahav 2020; ⁴Klicpera et al. 2019; ⁵Wang et B 2018; Kazi, Cosmo, et B. 2020

Physical systems as learning metaphor



Physical systems as learning metaphor



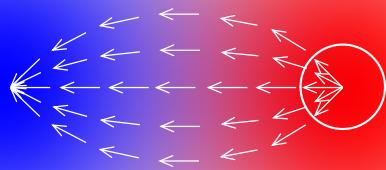
Wright et al. 2022

Diffusion Equation

Fourier 1822; Fick 1855

Diffusion Equation

heat flux $h \propto -\nabla x$



conservation condition: $\frac{\partial}{\partial t}x = -\operatorname{div}(h)$
("no heat created or disappears")

$$\frac{\partial}{\partial t}x(\mathbf{u}, t) = \operatorname{div}(a\nabla x(\mathbf{u}, t))$$

Fourier 1822; Fick 1855

Diffusion Equation

constant diffusivity


$$\frac{\partial}{\partial t}x = \operatorname{div}(c\nabla x)$$

Diffusion Equation

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

Homogeneous
Isotropic

1. Gradient flow of the *Dirichlet energy*

$$\mathcal{E}[x] = \frac{1}{2} \int_{\Omega} \|\nabla x(\mathbf{u})\|^2 d\mathbf{u}$$

Diffusion Equation

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

Homogeneous
Isotropic

1. Gradient flow of the *Dirichlet energy*

$$\mathcal{E}[x] = \frac{1}{2} \int_{\Omega} \|\nabla x(\mathbf{u})\|^2 d\mathbf{u}$$

2. Closed form solution: Gaussian filter

$$x(\mathbf{u}, t) = x(\mathbf{u}, 0) \star \frac{1}{(4\pi t)^{d/2}} e^{-\|\mathbf{u}\|^2/4t}$$

Diffusion Equation

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

Homogeneous
Isotropic

$$\frac{\partial}{\partial t}x = \operatorname{div}(a\nabla x)$$

Non-homogeneous
Isotropic

Position-dependent
diffusivity $a(u)$

Position & direction
dependent diffusivity $\mathbf{A}(u)$

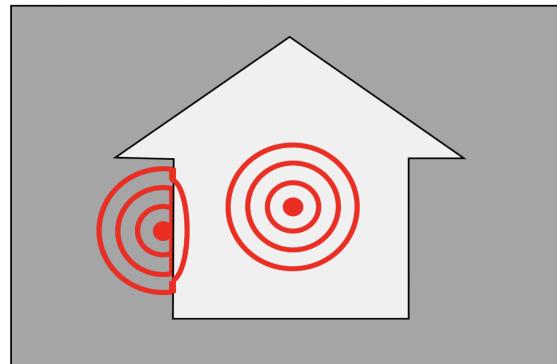
$$\frac{\partial}{\partial t}x = \operatorname{div}(\mathbf{A}\nabla x)$$

Non-homogeneous
Anisotropic

Diffusion Equation in Image Processing

Edge indicator $a(u) \propto \|\nabla x(u)\|^{-1}$

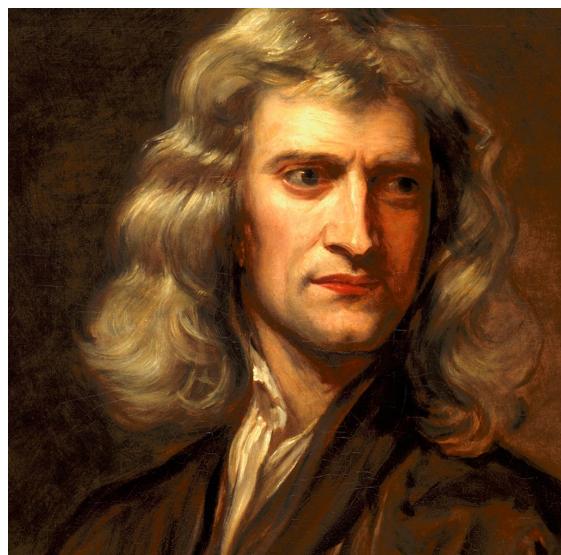
$$\frac{\partial}{\partial t}x = \operatorname{div}(a(x)\nabla x)$$



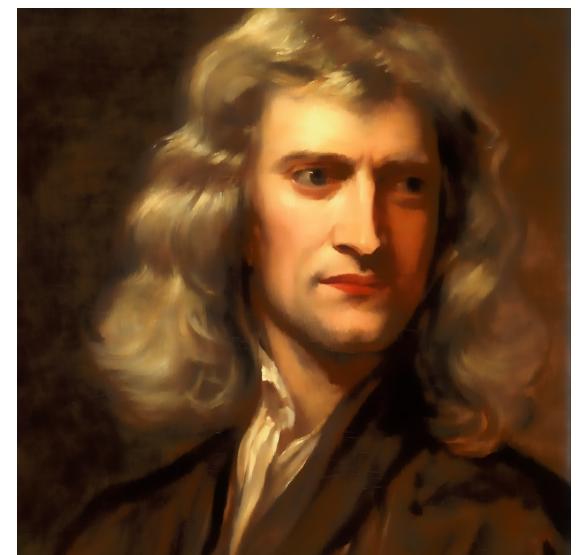
“Do not diffuse across edges”

Perona, Malik 1990

Diffusion in Image Processing



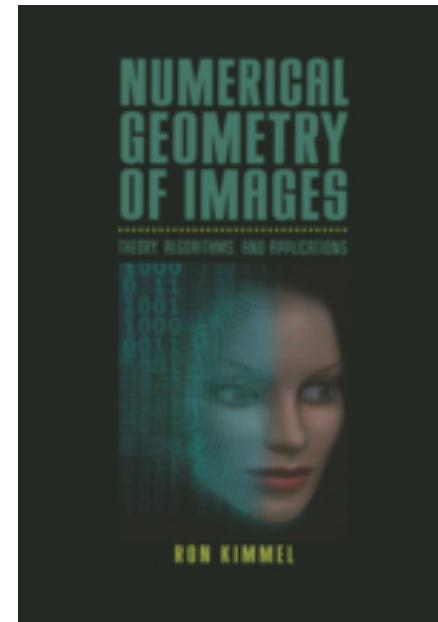
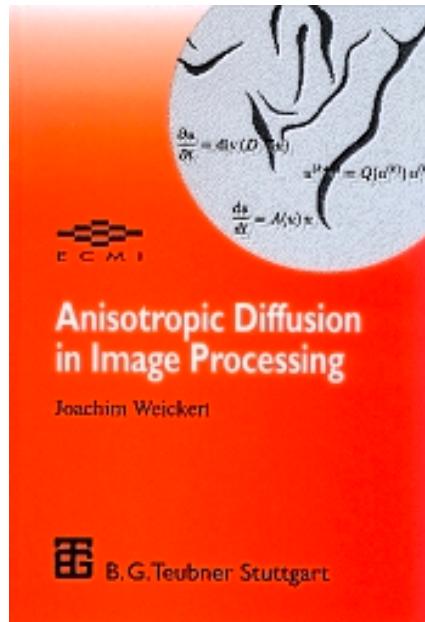
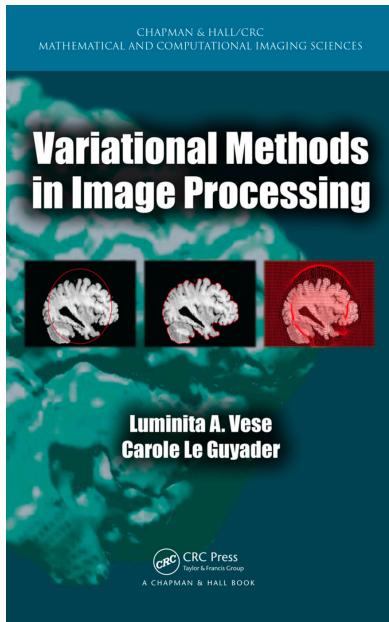
Homogeneous
diffusion



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

Diffusion Equation on Graphs

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

div diffusivity
 $\mathbf{A}(\mathbf{X})$ gradient
 $\nabla \mathbf{X}$

Diffusion Equation on Graphs

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

Explicit (Forward Euler) discretization: $t = k\tau$

$$\frac{\mathbf{x}_i^{(k+1)} - \mathbf{x}_i^{(k)}}{\tau} = \sum_{j:(i,j) \in E} a(\mathbf{x}_i^{(k)}, \mathbf{x}_j^{(k)}) (\mathbf{x}_j^{(k)} - \mathbf{x}_i^{(k)})$$

forward difference

Diffusion Equation on Graphs

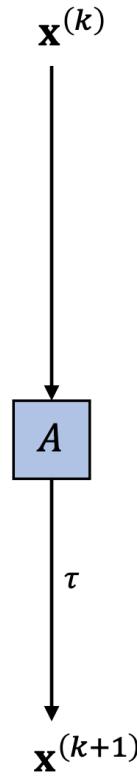
$$\frac{\partial}{\partial t} \mathbf{x}_i(t) = \sum_{j:(i,j) \in E} a(\mathbf{x}_i(t), \mathbf{x}_j(t)) (\mathbf{x}_j(t) - \mathbf{x}_i(t))$$

Explicit (Forward Euler) discretization: $t = k\tau$

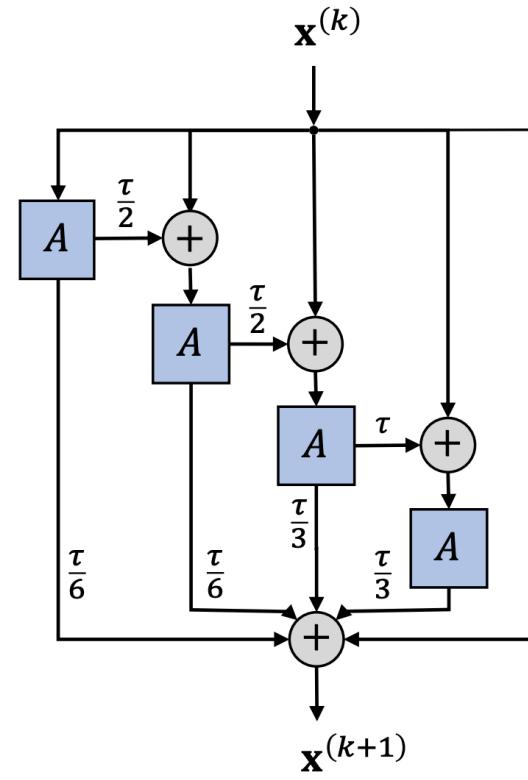
$$\mathbf{x}_i^{(k+1)} = \sum_{j:(i,j) \in E} a(\mathbf{x}_i^{(k)}, \mathbf{x}_j^{(k)}) \mathbf{x}_j^{(k)}$$

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

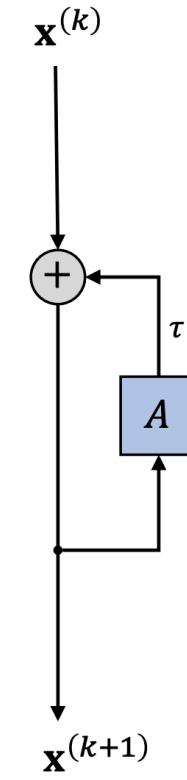
GAT is a particular discretisation of graph diffusion



Explicit
Fixed step



Explicit
Multi-step (Runge-Kutta)



Implicit

Graph Neural Diffusion (GRAND)

Given graph $G = (V, E)$ with input node features \mathbf{X}_{in}

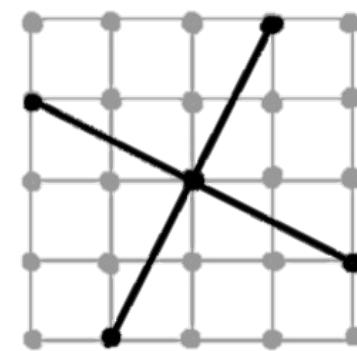
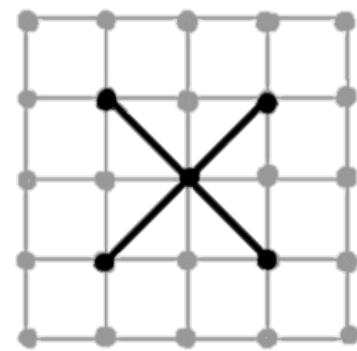
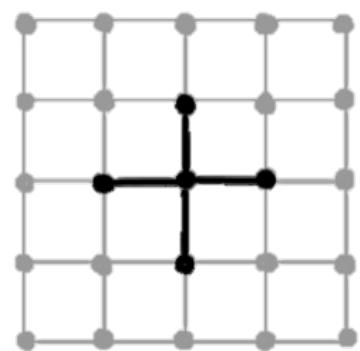
- Set initial condition: $\mathbf{X}(0) = \phi(\mathbf{X}_{\text{in}})$
- Solve graph diffusion eqn: $\mathbf{X}(T) = \mathbf{X}(0) + \int_0^T \text{div} \left(\mathbf{A}(\mathbf{X}(t)) \nabla \mathbf{X}(t) \right) dt$
using an iterative solver
- Output: $\mathbf{Y} = \psi(\mathbf{X}(T))$

where ϕ, ψ and the diffusivity A are learnable functions

What do we gain?

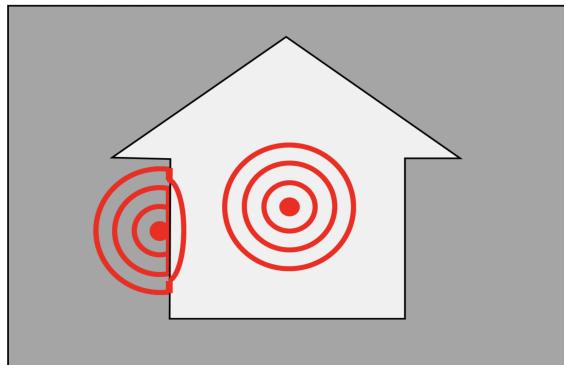
- New perspectives on old problems (e.g. oversmoothing, bottlenecks, etc)
- New architectures
 - Many GNNs can be formalised as a discretised Graph Diffusion equation
 - More efficient solvers (multistep, adaptive, implicit, multigrid, etc.)
 - Implicit schemes = multi-hop filters
- Theoretical guarantees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

Spatial Derivative: Graph Rewiring?



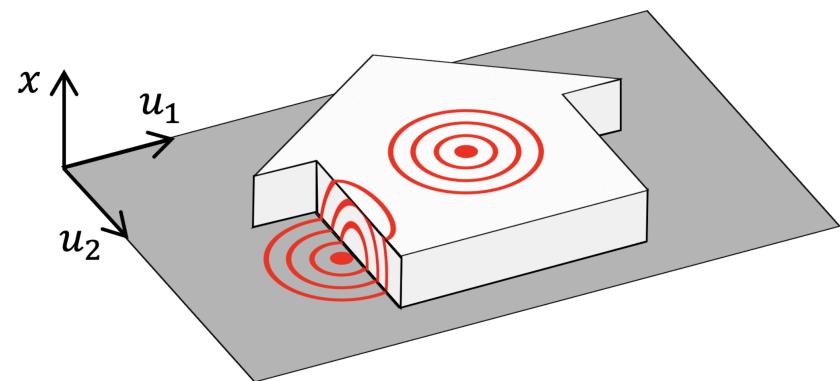
Different discretisations of 2D Laplacian

Images as embedded manifolds



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

Non-linear diffusion



$$\frac{\partial}{\partial t} \mathbf{z} = \Delta_G \mathbf{z}$$

Non-Euclidean diffusion

Beltrami flow

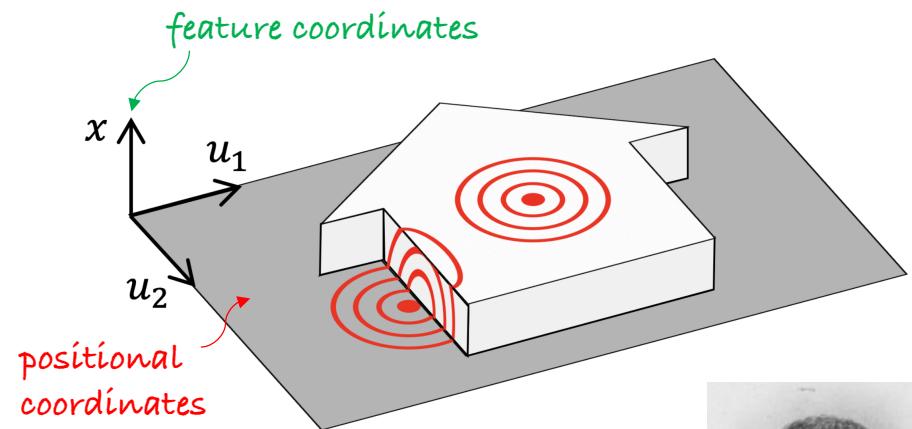
- Consider image as embedded 2-manifold

$$\mathbf{z}(\mathbf{u}) = (\mathbf{u}, \alpha \mathbf{x}(\mathbf{u}))$$

- Pullback metric: 2×2 matrix

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

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



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



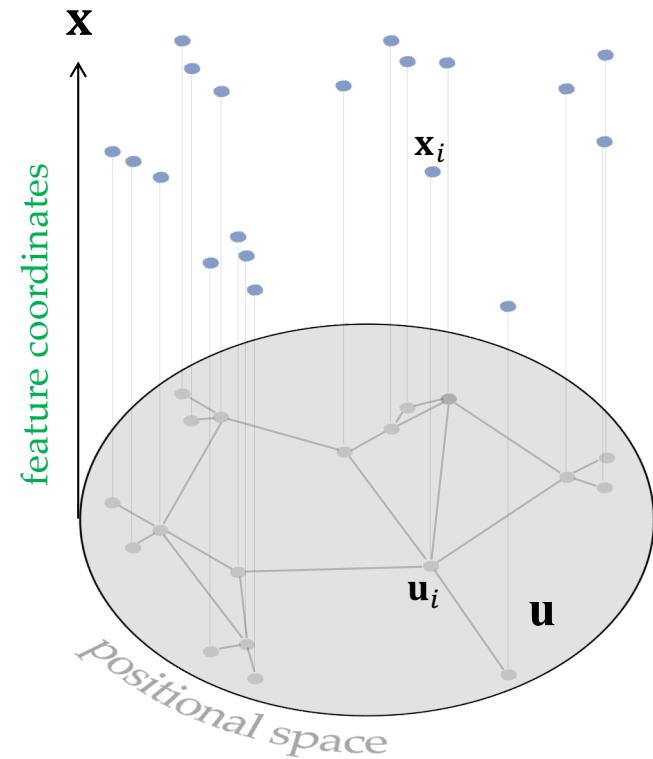
Eugenio Beltrami

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**

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

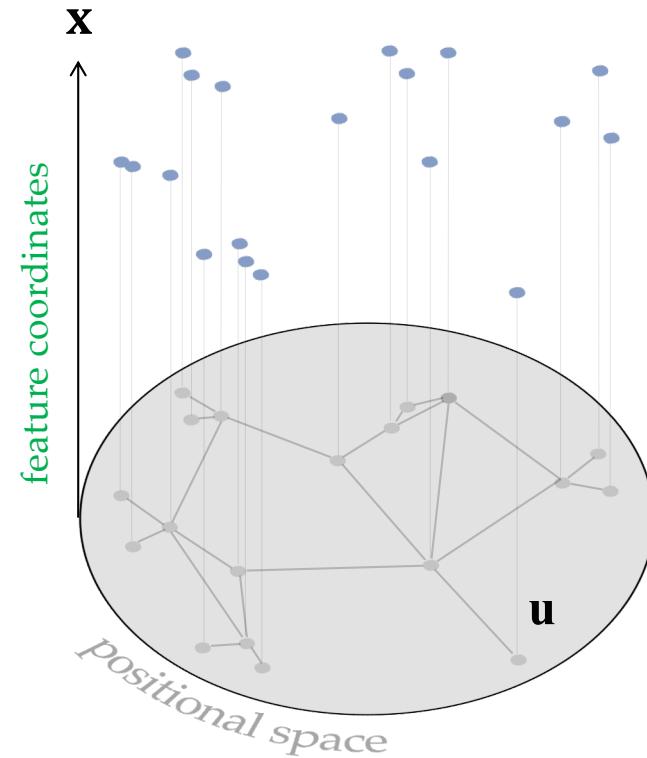


Graph Beltrami flow

- 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

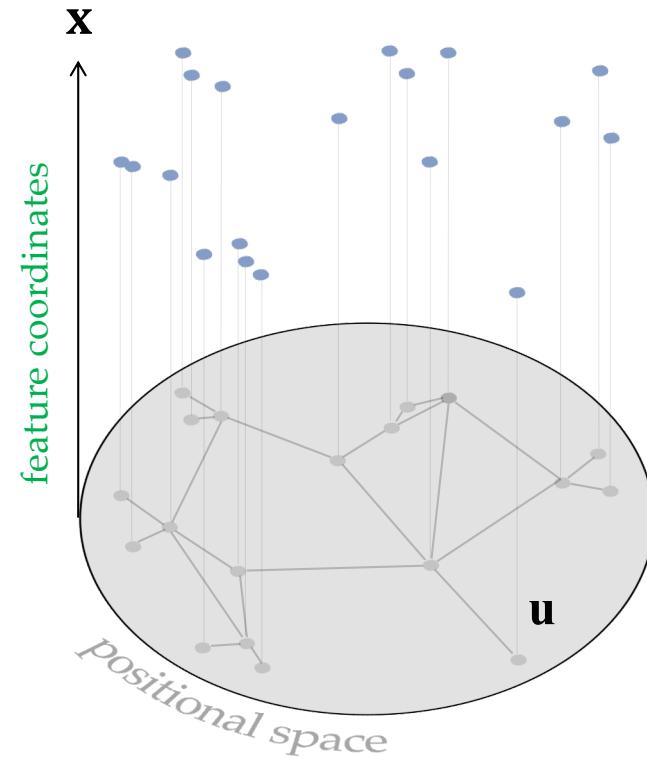


Graph Beltrami flow

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



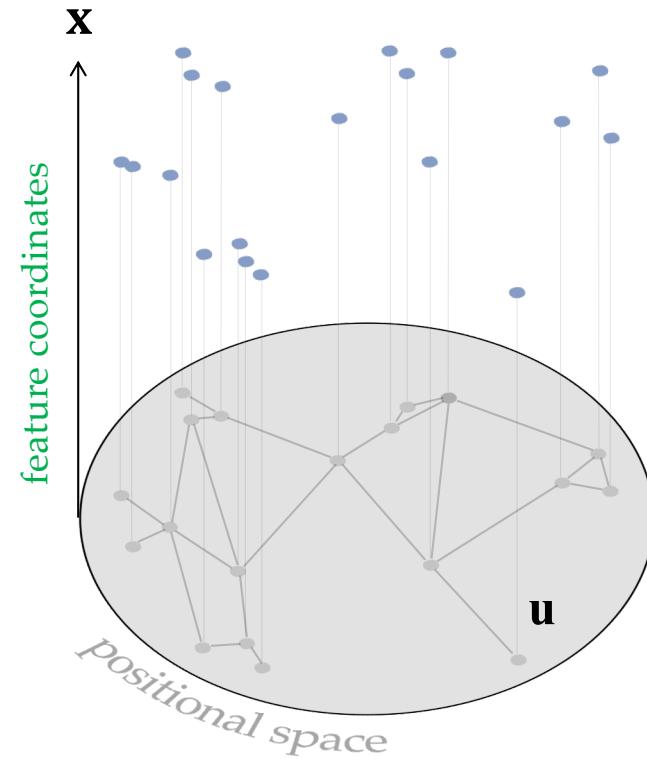
Graph Beltrami flow

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

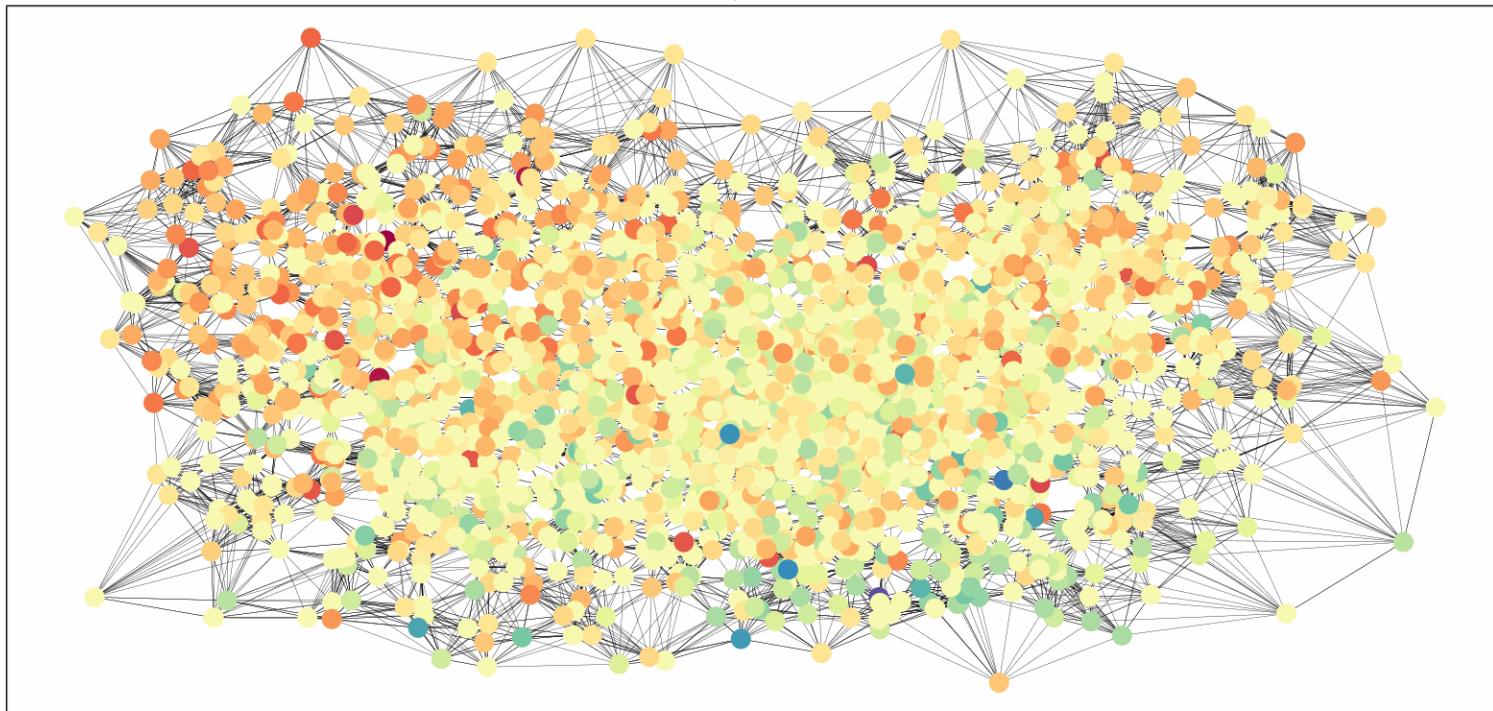
$$\frac{\partial}{\partial t} \mathbf{z}_i = \sum_{j:(i,j) \in E'} a(\mathbf{z}_i, \mathbf{z}_j)(\mathbf{z}_j - \mathbf{z}_i)$$

rewired graph

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



Graph Beltrami flow



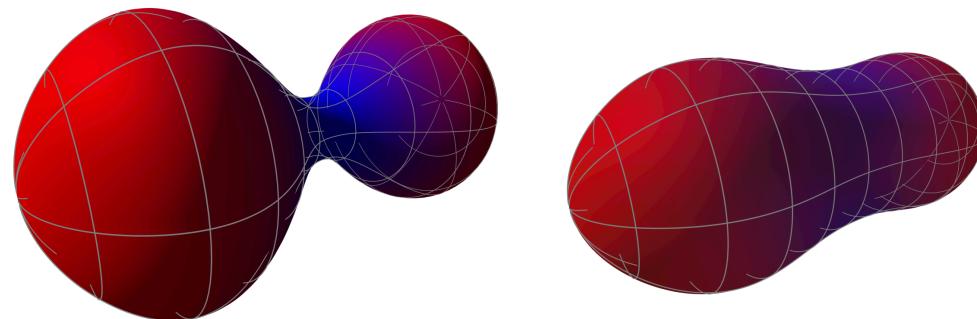
Evolution of positional/feature components + rewiring of the Cora graph

Chamberlain, Rowbottom, et B. 2021

Ricci flow

- Ricci flow: “diffusion of the Riemannian metric”

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



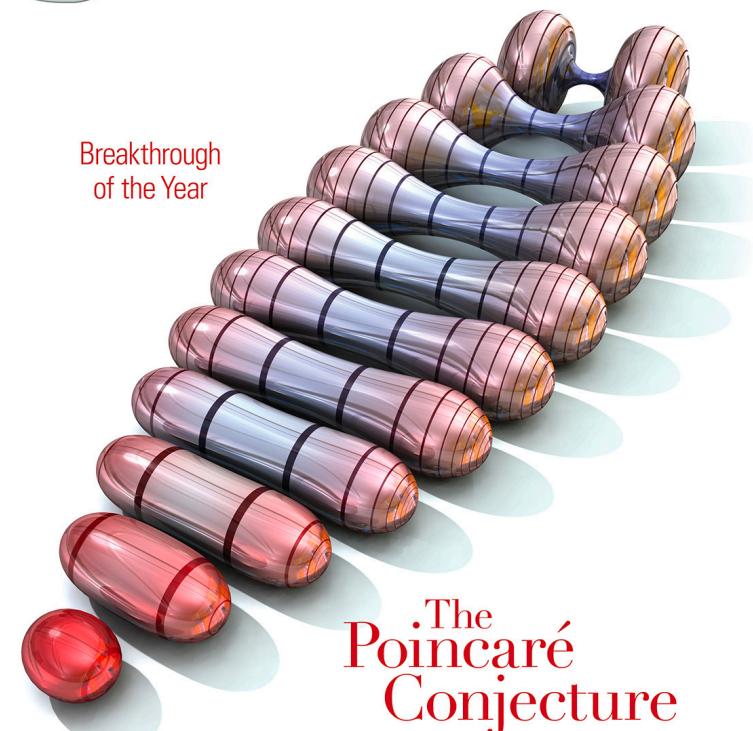
Evolution of a manifold under Ricci flow

Ricci 1903; Hamilton 1988; Perelman 2003

Science

22 December 2006 | \$10

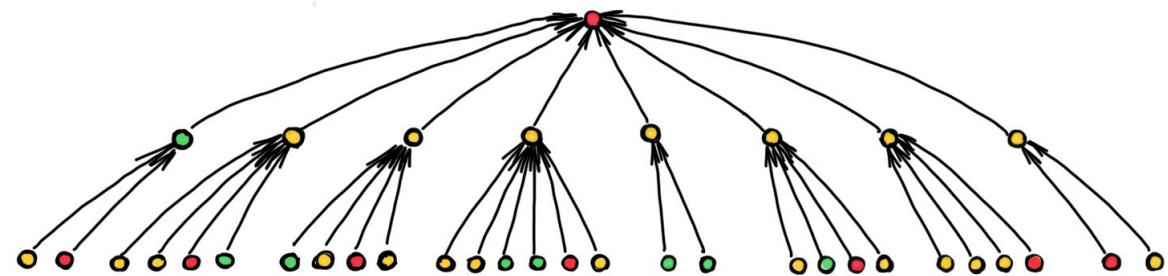
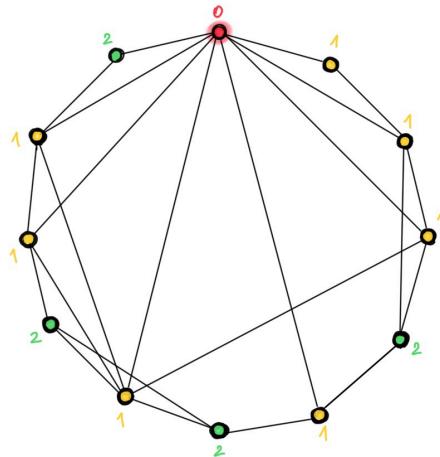
Breakthrough
of the Year



The
Poincaré
Conjecture
PROVED

AAAS

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

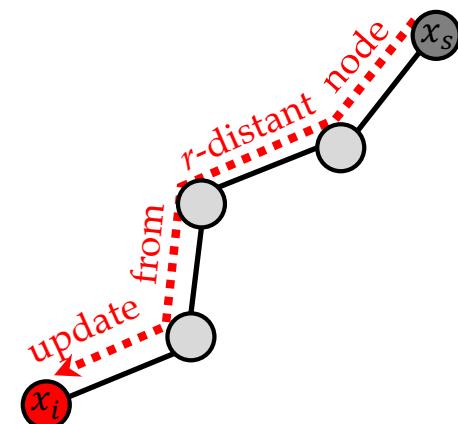
- Multilayer MPNN-type GNN of the form

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

- $|\nabla \phi_\ell| \leq \alpha$ 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

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



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

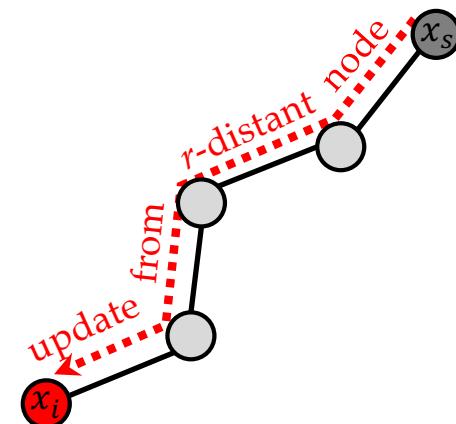
$$x_i^{(\ell+1)} = \phi_\ell \left(x_i^{(\ell)}, \sum_{j=1}^n \hat{a}_{ij} \psi_\ell \left(x_i^{(\ell)}, x_j^{(\ell)} \right) \right)$$

- $|\nabla \phi_\ell| \leq \alpha$ 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

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

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

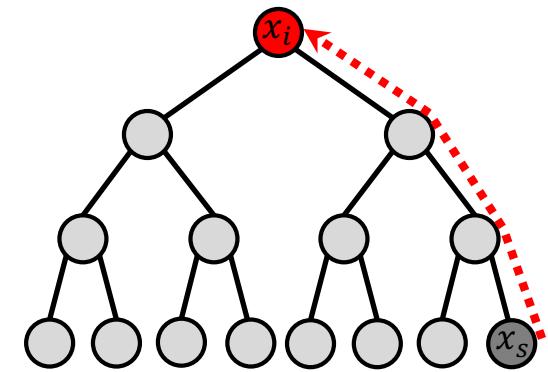
$$x_i^{(\ell+1)} = \phi_\ell \left(x_i^{(\ell)}, \sum_{j=1}^n \hat{a}_{ij} \psi_\ell \left(x_i^{(\ell)}, x_j^{(\ell)} \right) \right)$$

- $|\nabla \phi_\ell| \leq \alpha$ 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

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

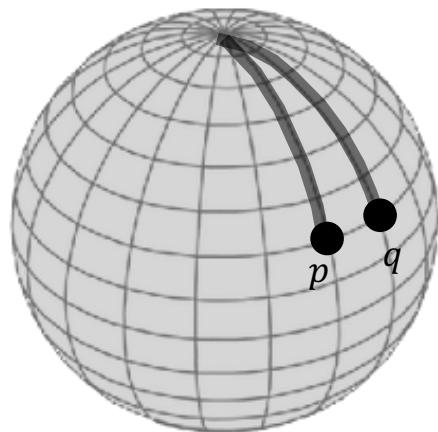
It's the graph structure
("bottleneck") to blame!



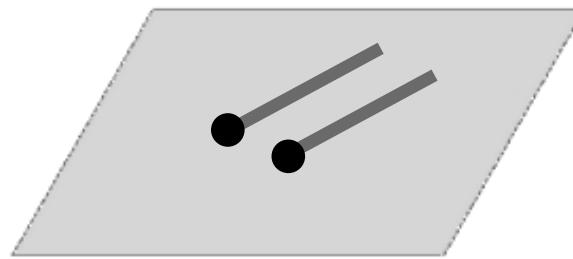
Pathological example: binary tree

$$(\widehat{\mathbf{A}}^{r+1})_{is} = \frac{1}{2} \cdot 3^{-r}$$

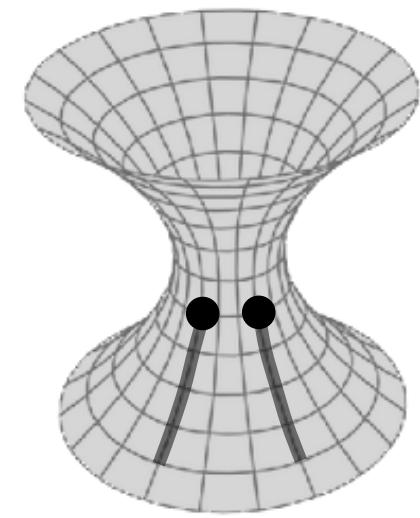
Ricci Curvature on Manifolds



Spherical (>0)



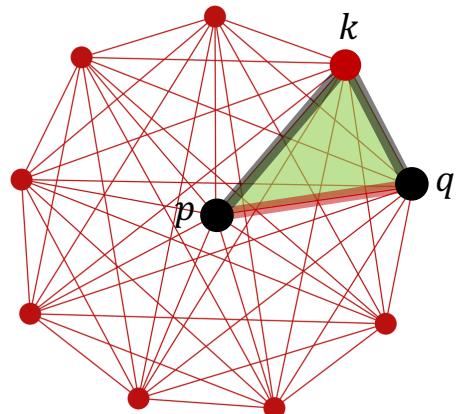
Euclidean ($=0$)



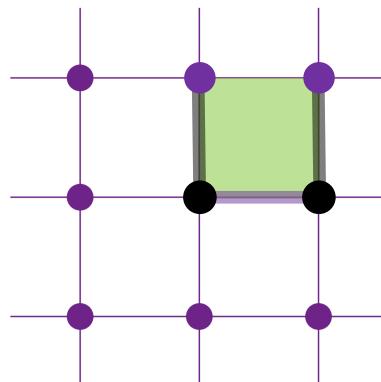
Hyperbolic (<0)

“geodesic dispersion”

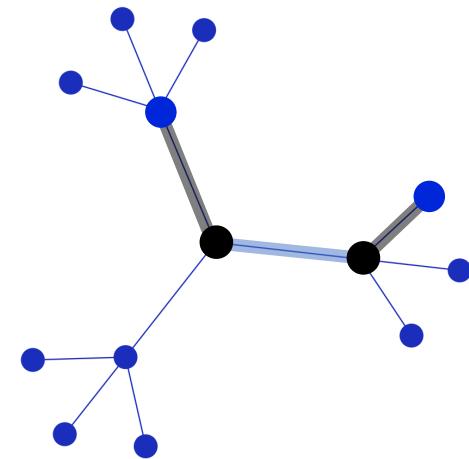
Ricci Curvature on Graphs



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

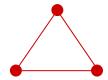
$$\text{Ric}(i, j) = \frac{2}{d_i} + \frac{2}{d_j} + 2 \frac{|\#_{\Delta}(i, j)|}{\max\{d_i, d_j\}} + \frac{|\#_{\Delta}(i, j)|}{\min\{d_i, d_j\}} + \frac{\gamma_{\max}^{-1}}{\max\{d_i, d_j\}} (|\#_{\square}^i(i, j)| + |\#_{\square}^j(i, j)|) - 2$$

Degree of i ↗ Triangles based at $i \sim j$ ↗ Max number of 4-cycle based at $i \sim j$ traversing the same node ↗ Neighbours of i forming a 4-cycle based at $i \sim j$ (w/o diagonals) ↗

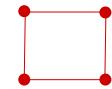
Balanced Forman Curvature

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

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



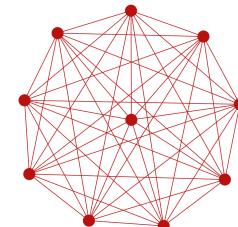
Cycle C_3 : $\frac{3}{2}$



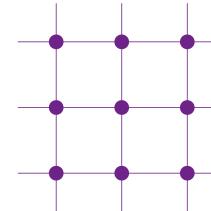
C_4 : 1



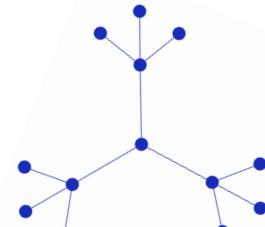
$C_{n \geq 5}$: 0



Clique K_n : $\frac{n}{n-1}$



Grid G_n : 0



Tree T_r : $\frac{4}{r+1} - 2$

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 $\text{Ric}(i, j) \leq -2 + \delta$. Then, there exist nodes $Q \subset \{s: d_G(i, s) = 2\}$ of size $|Q| > 1/\delta$ s.t.

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

Small δ = negative curvature more nodes stronger over-squashing

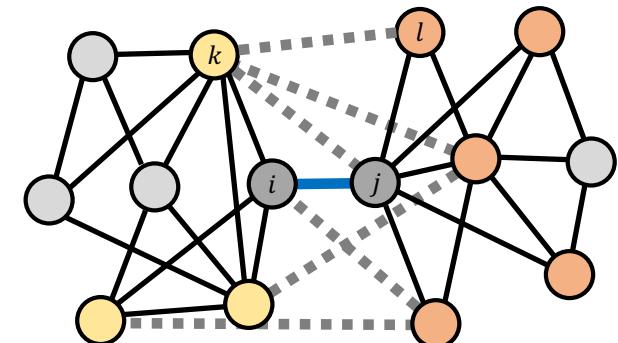
Over-squashing is caused by negatively-curved edges!

Stochastic Discrete Ricci Flow (SDRF)

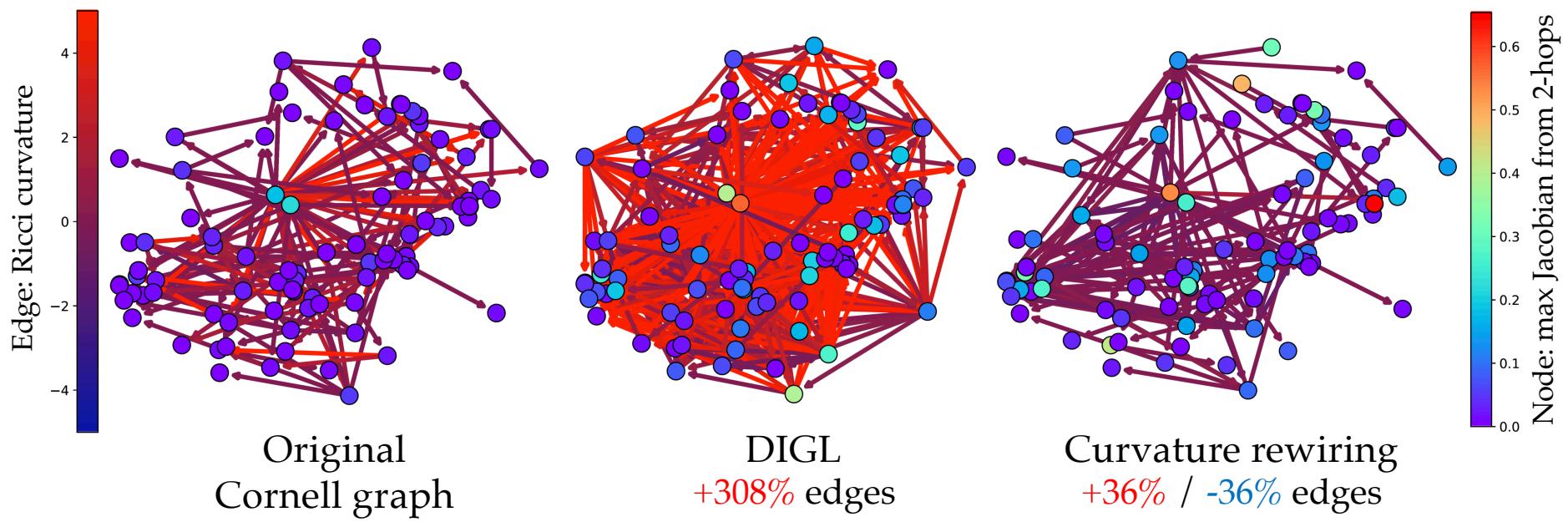
Input: graph $G = (V, E)$, temperature $\tau > 0$, (optional C)

- • 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 $G' = (V, E')$

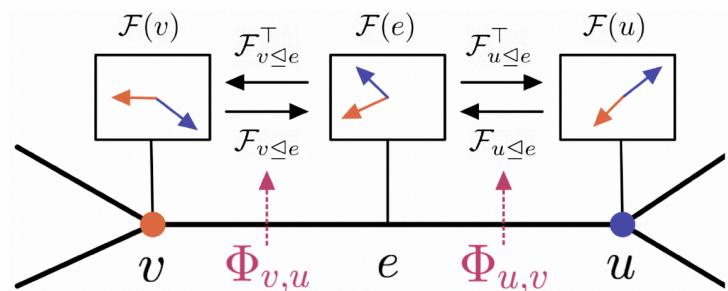


Curvature- vs Diffusion-based Rewiring

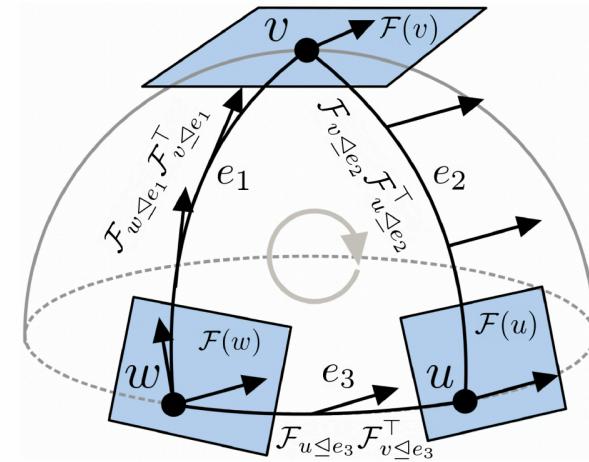


Topping, di Giovanni, et B. 2021; Klicpera et al. 2019 (DIGL)

Cellular Sheaves



Cellular sheaf



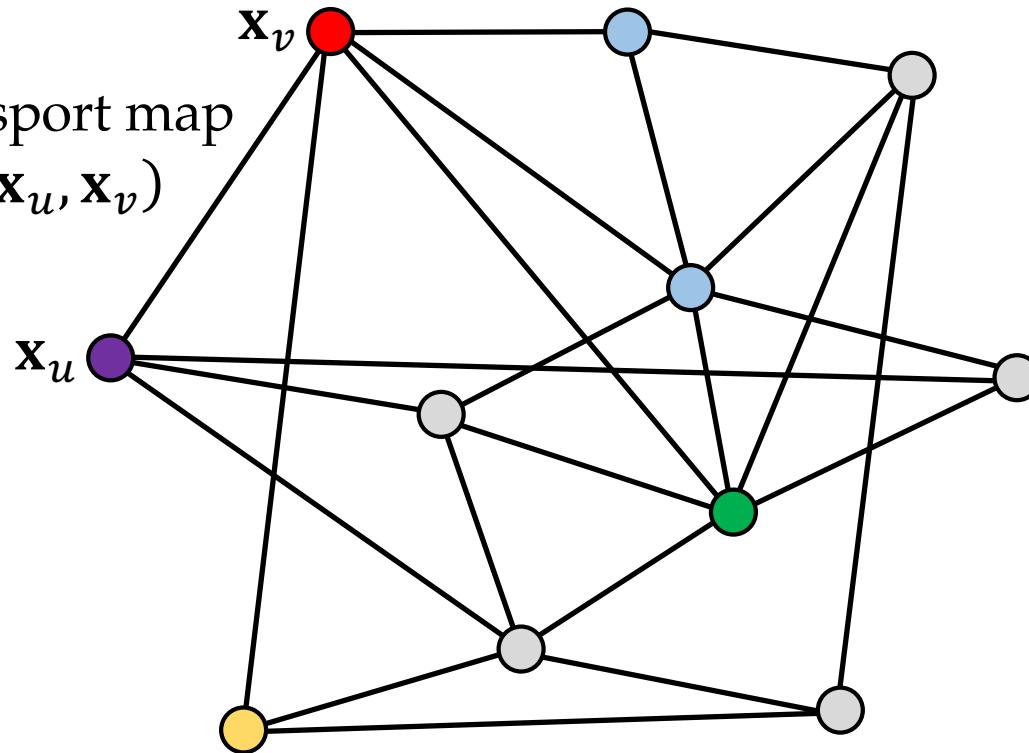
Analogy to parallel transport
on manifolds

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

Cellular Sheaves

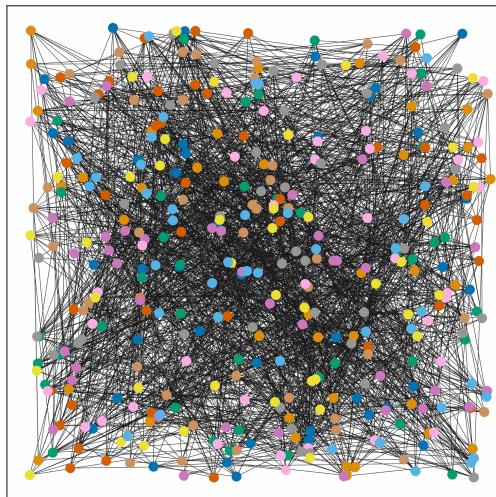
learnable transport map

$$\mathcal{F}_{u \leq v} = \Phi(\mathbf{x}_u, \mathbf{x}_v)$$



$$\text{Sheaf Laplacian: } (\Delta_{\mathcal{F}} \mathbf{x})_v = \sum_{v, u \leq e} \mathcal{F}_{v \leq e}^T (\mathcal{F}_{v \leq e} \mathbf{x}_v - \mathcal{F}_{u \leq e} \mathbf{x}_u)$$

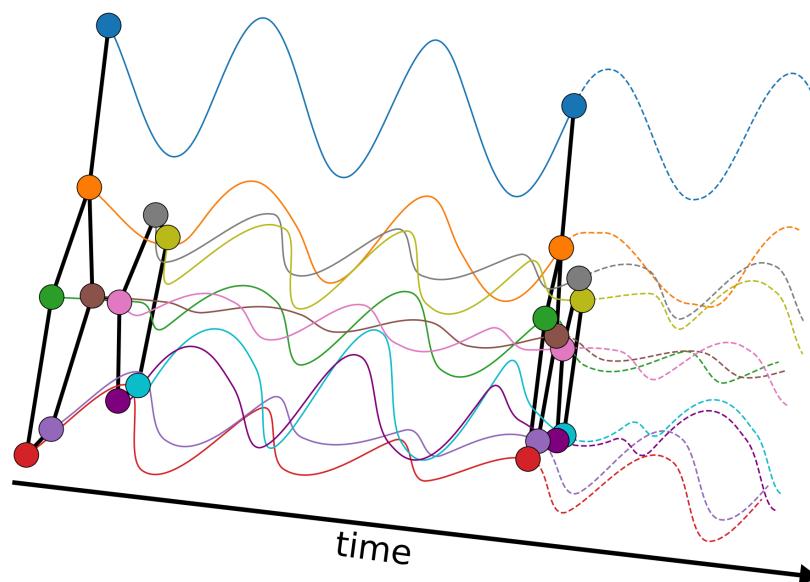
Diffusion on Cellular Sheaves



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

Node classification = limit of sheaf diffusion equation
with an appropriate sheaf, alternative to WL

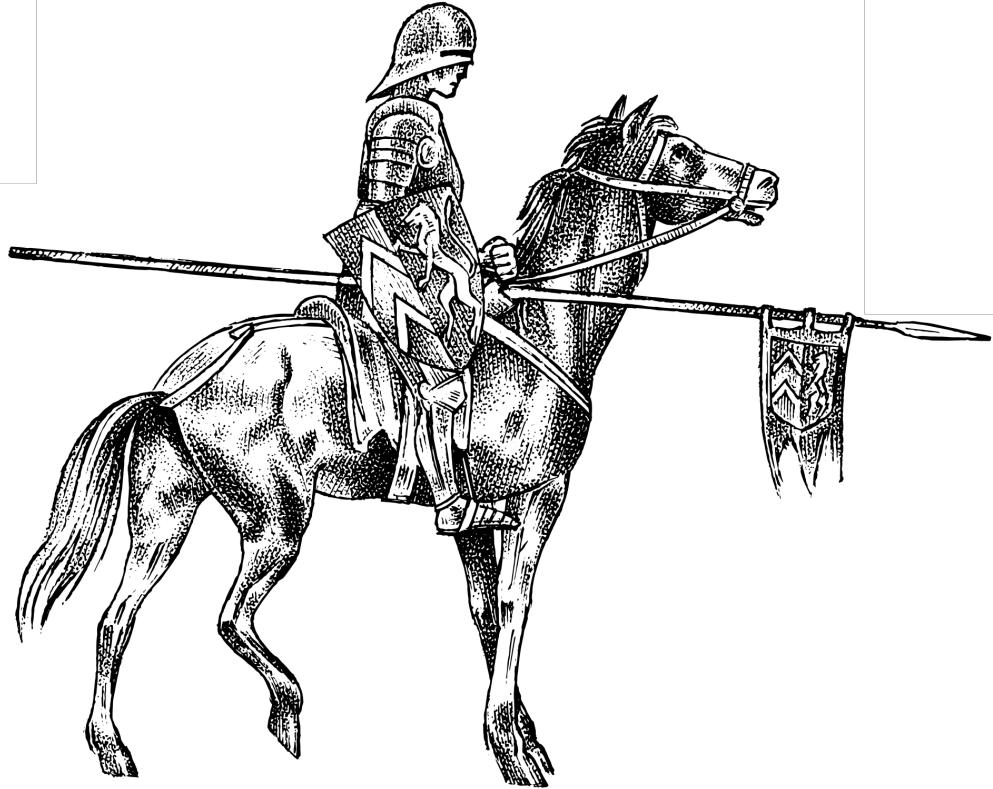
Graph-Coupled Oscillators



Dynamics of a system of coupled oscillators on a molecular graph



Image: Michael Galkin



Are we done with Message Passing?