

Geometric Deep Learning

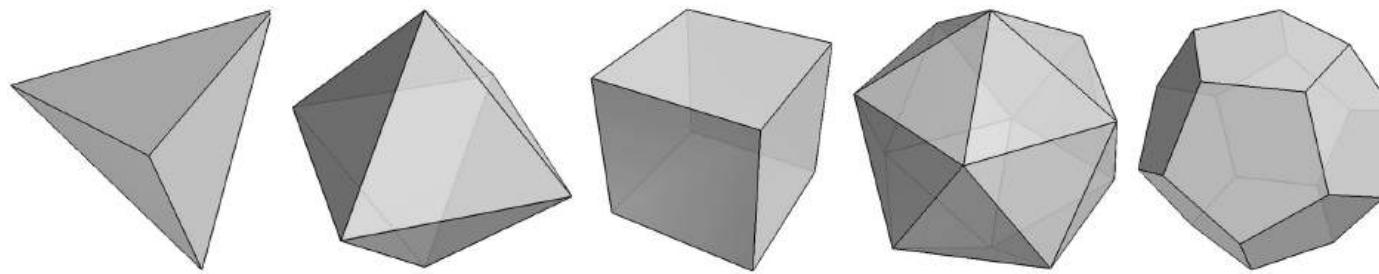
Michael Bronstein

University of Oxford & AITHYRA

“**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**”



H. Weyl

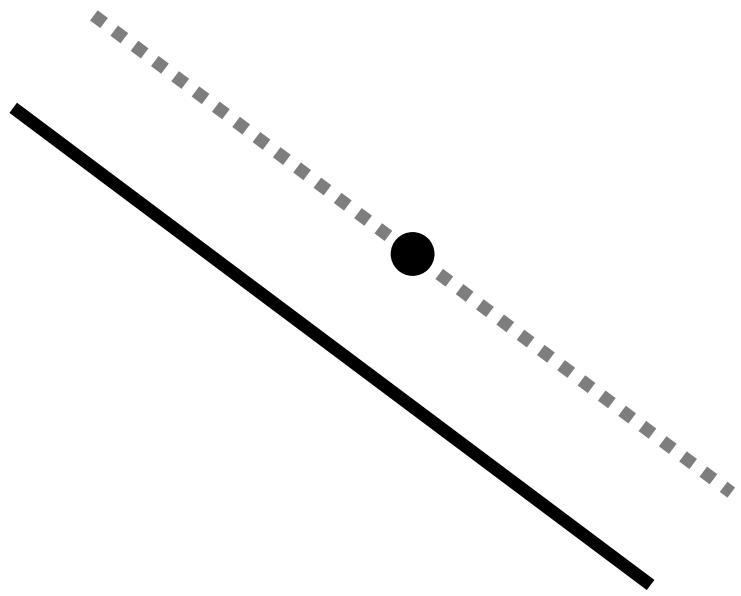


Plato

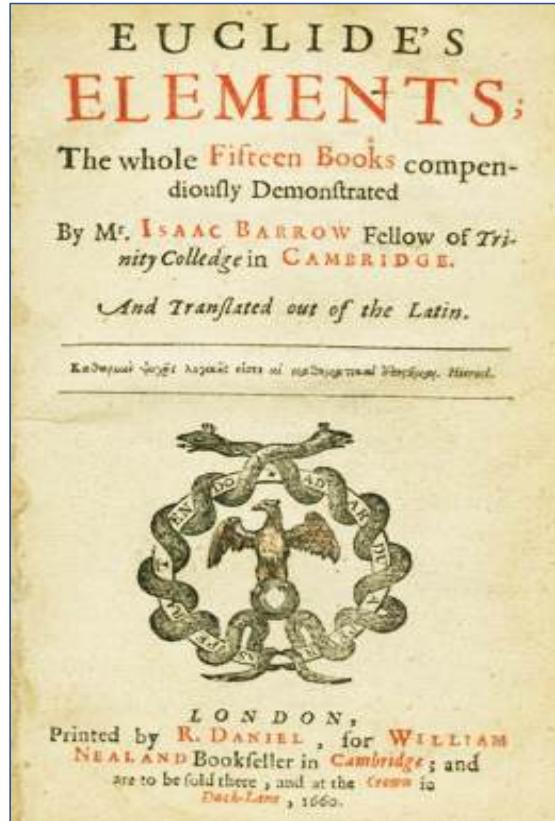
~370 BC

Portrait: Ihor Gorskyi

ΑΓΕΩΜΕΤΡΗΤΟΣ
ΜΗΔΕΙΣ, ΕΙΣΙΤΩ



Fifth Postulate



Euclid

~300 BC

Portrait: Ihor Gorskyi

XIX century





The Erlangen Programme



Geometry = space + transformation group

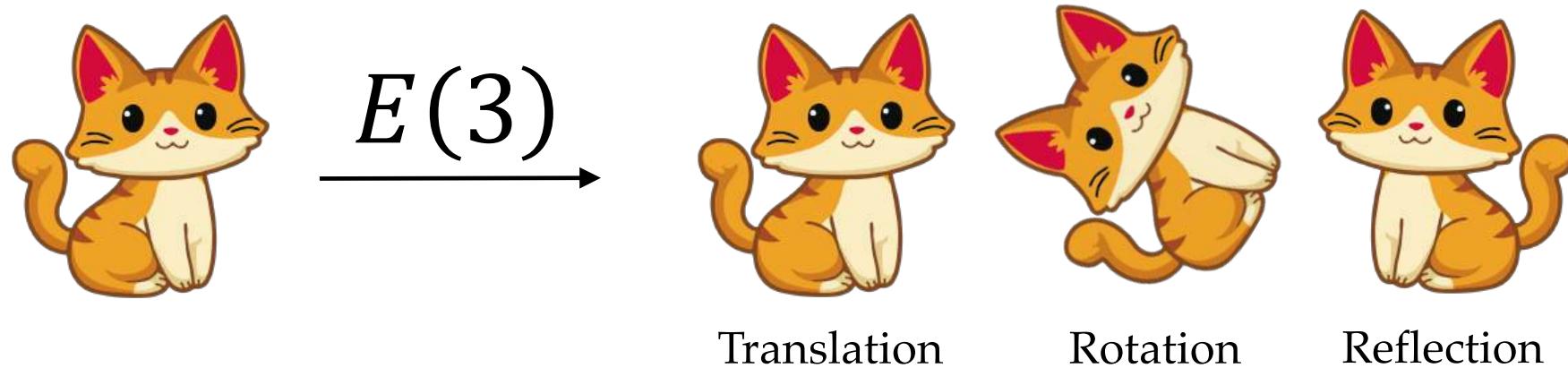
Klein 1872



F. Klein

1872

Euclidean geometry





H. Poincaré

1904



H. Minkowski

1907



E. Noether

1918



H. Weyl

1929



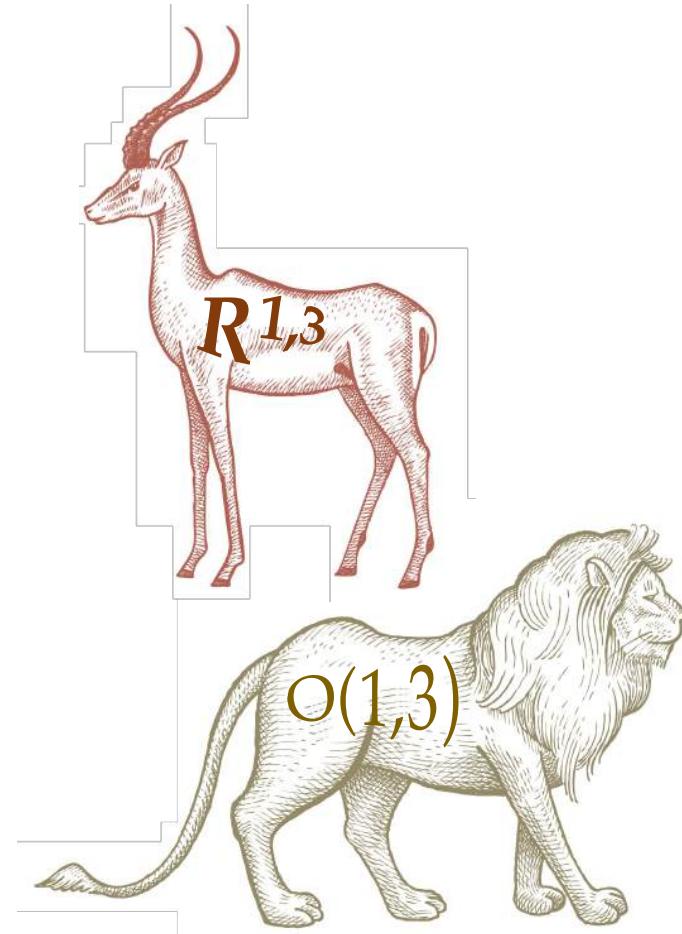
C. N. Yang

1954

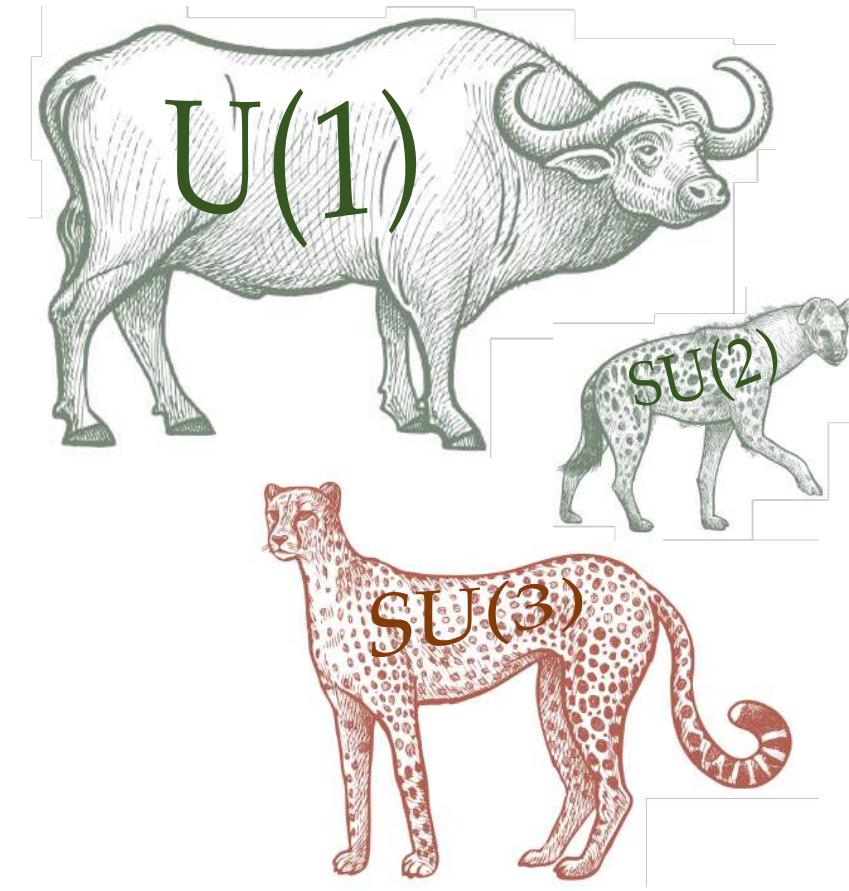


R. L. Mills

Poincaré 1904; Noether 1918; Weyl 1929; Yang & Mills 1954; Portraits: Ihor Gorskyy



External symmetry



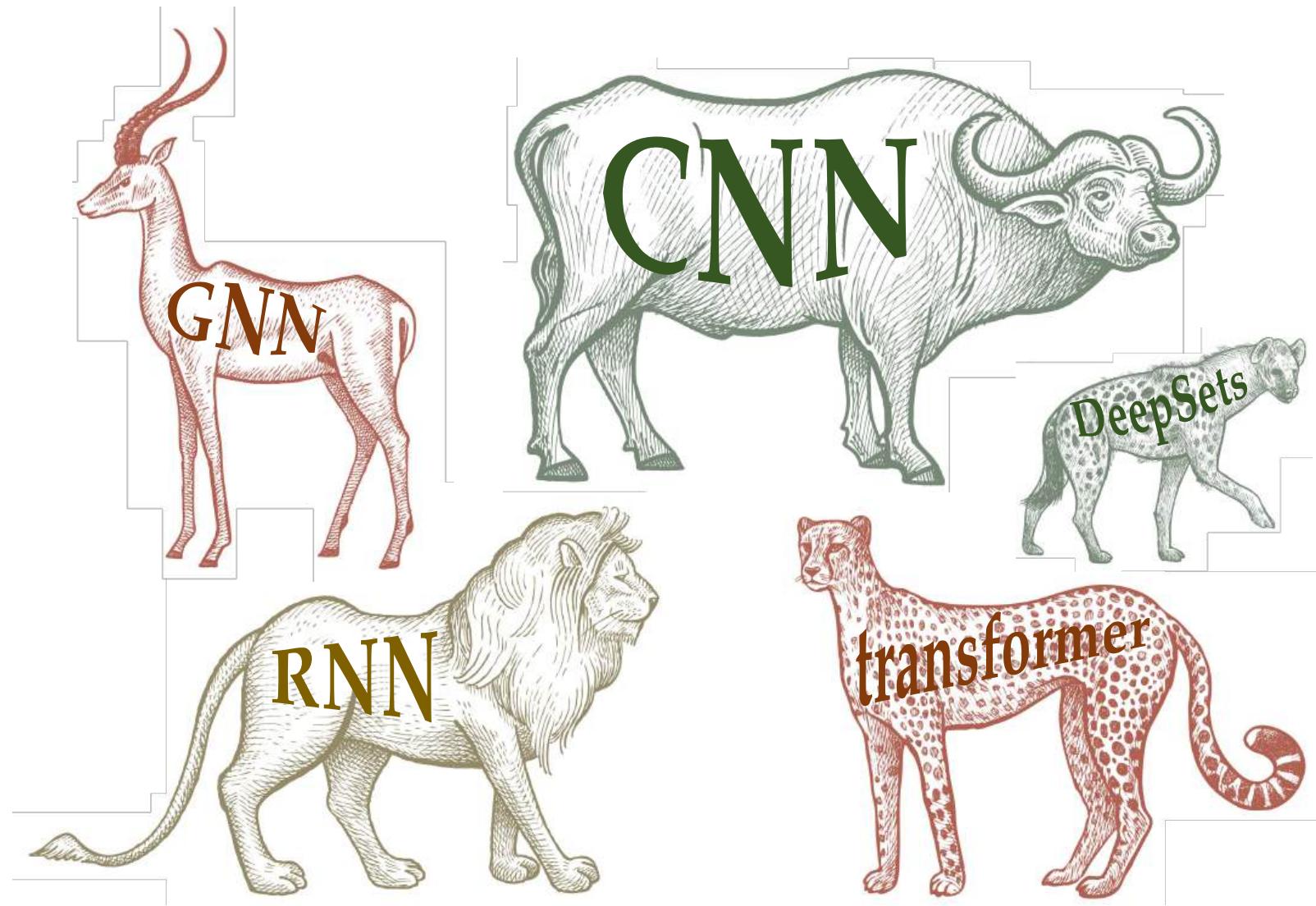
Internal symmetry

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

— More is different



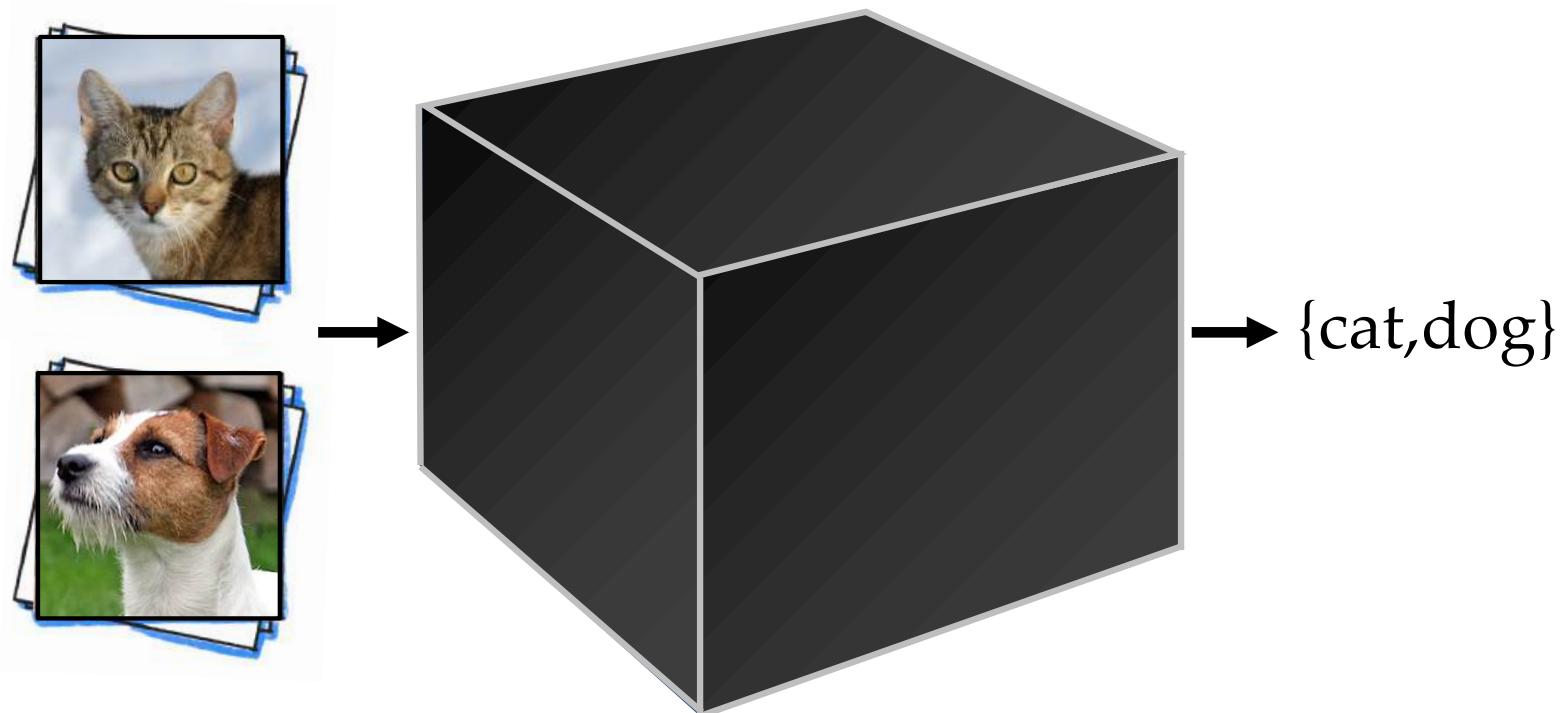
P. Anderson



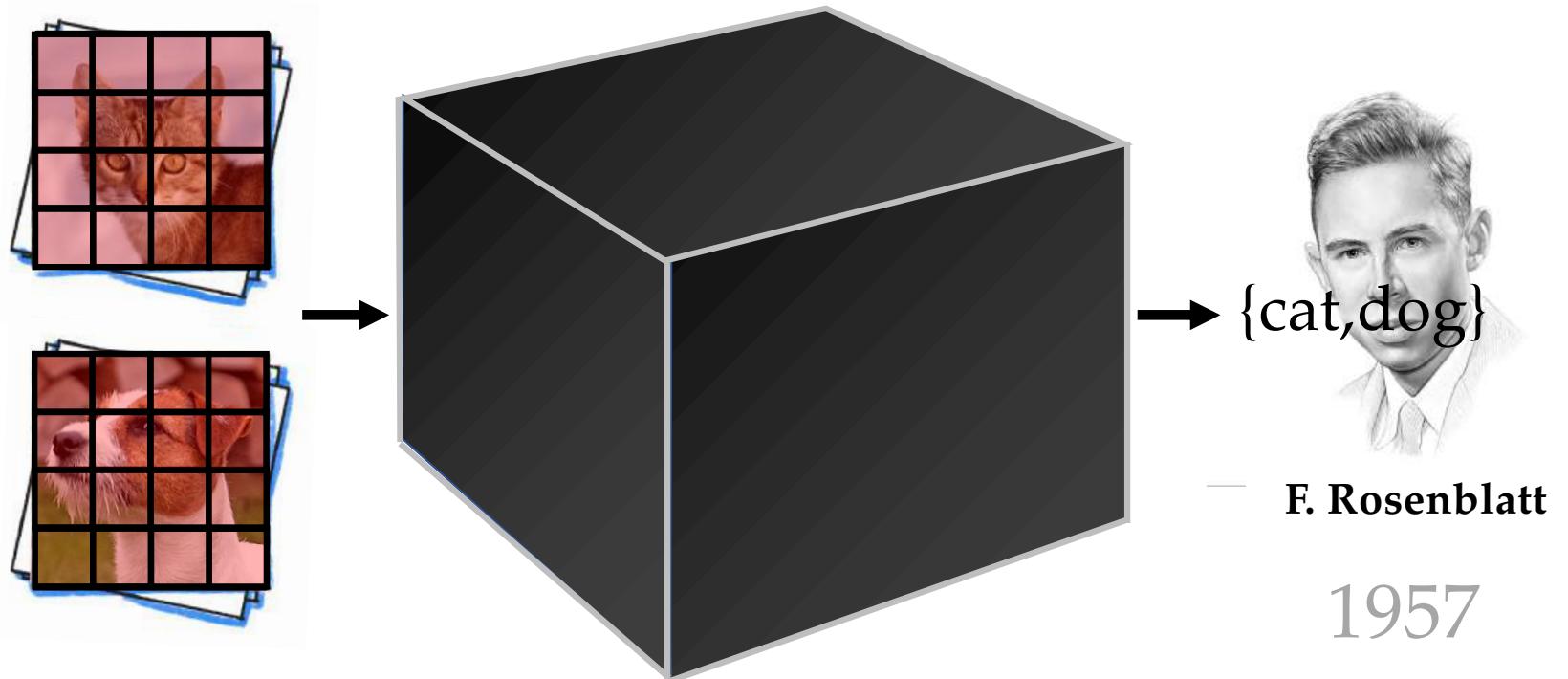


Geometric Deep Learning

Supervised ML = Function Approximation

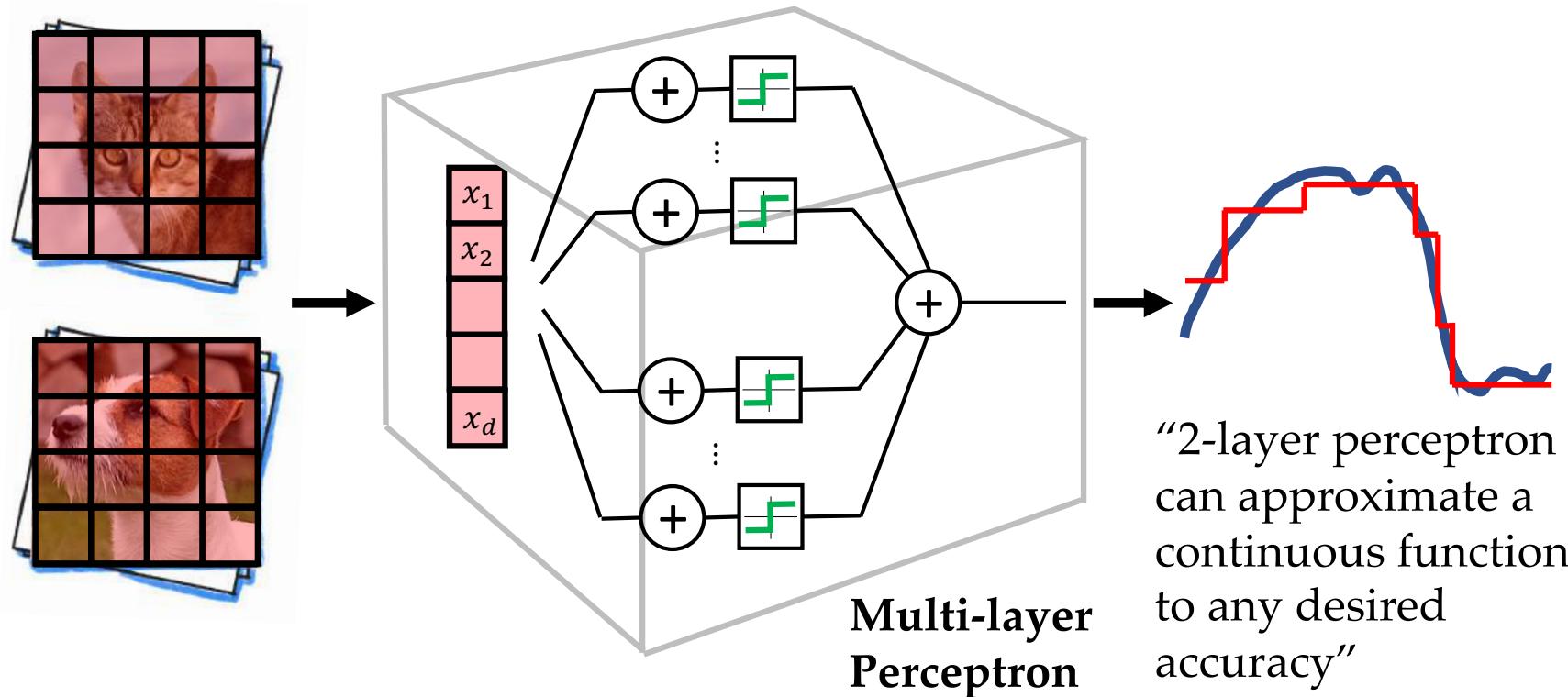


Supervised ML = Function Approximation



Rosenblatt 1957; Portrait: Ihor Gorskyi

Universal Approximation

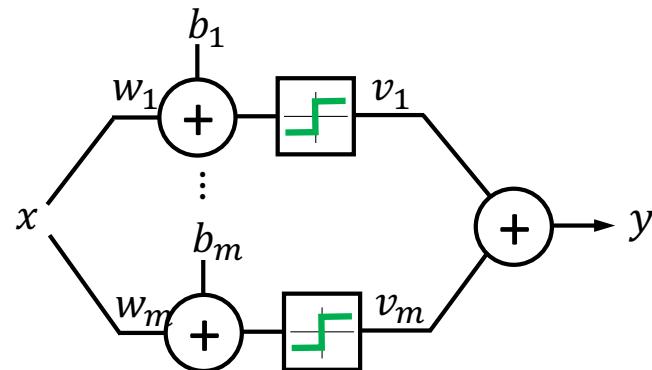


Universal Approximation: Hilbert’s 13th problem 1900; Kolmogorov 1956; Arnold 1957; Cybenko 1989; Hornik 1991; Barron 1993; Leshno et al 1993; Maiorov 1999; Pinkus 1999

Shallow Perceptrons are universal approximators

Universal Approximation Theorem: σ is not polynomial iff for every continuous function $f: K \subset \mathbb{R}^d \rightarrow \mathbb{R}$ defined on a *compact set* K and $\varepsilon > 0$, there exists a two-layer Perceptron with m neurons and weights $\mathbf{W}, \mathbf{b}, \mathbf{v}$ s.t.

$$\max_{x \in K} \left| f(x) - \sum_{j \leq m} v_j \sigma(w_j^T x + b_j) \right| < \varepsilon$$



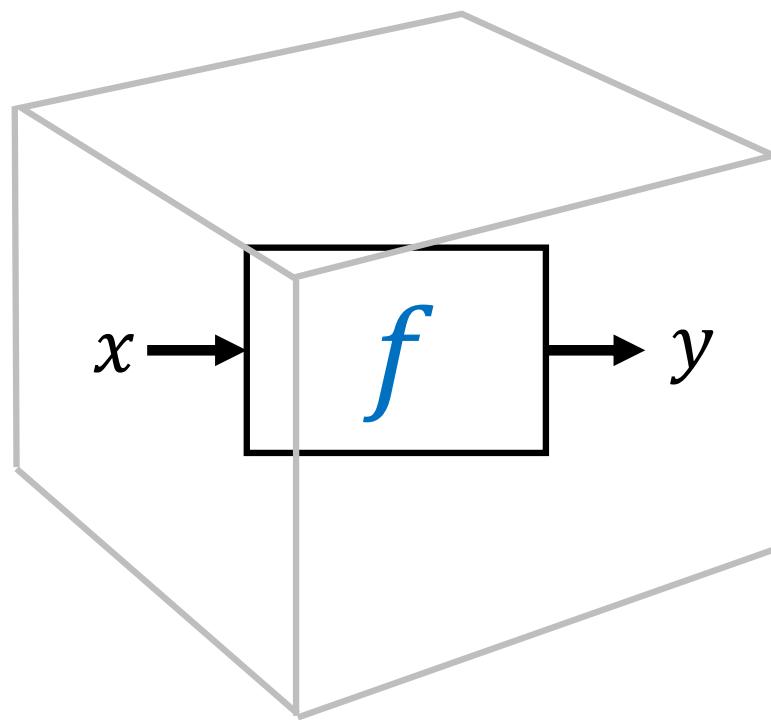
Cybenko 1989; Hornik 1991; Pinkus 1999

Shallow Perceptrons are universal approximators

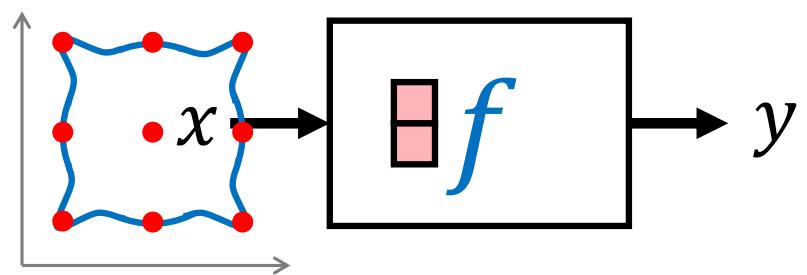
Universal Approximation Theorem: σ is not polynomial iff for every continuous function $f: K \subset \mathbb{R}^d \rightarrow \mathbb{R}$ defined on a *compact set* K and $\varepsilon > 0$, there exists a two-layer Perceptron with m neurons and weights $\mathbf{W}, \mathbf{b}, \mathbf{v}$ s.t.

$$\max_{x \in K} \left| f(x) - \sum_{j \leq m} v_j \sigma(w_j^T x + b_j) \right| < \varepsilon$$

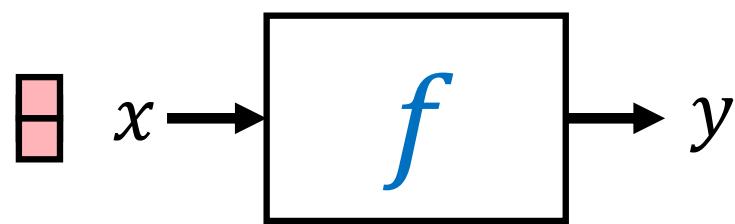
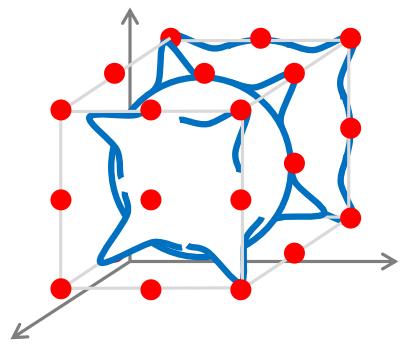
**What is the relation between dimension d ,
number of neurons m , and the error ε ?**

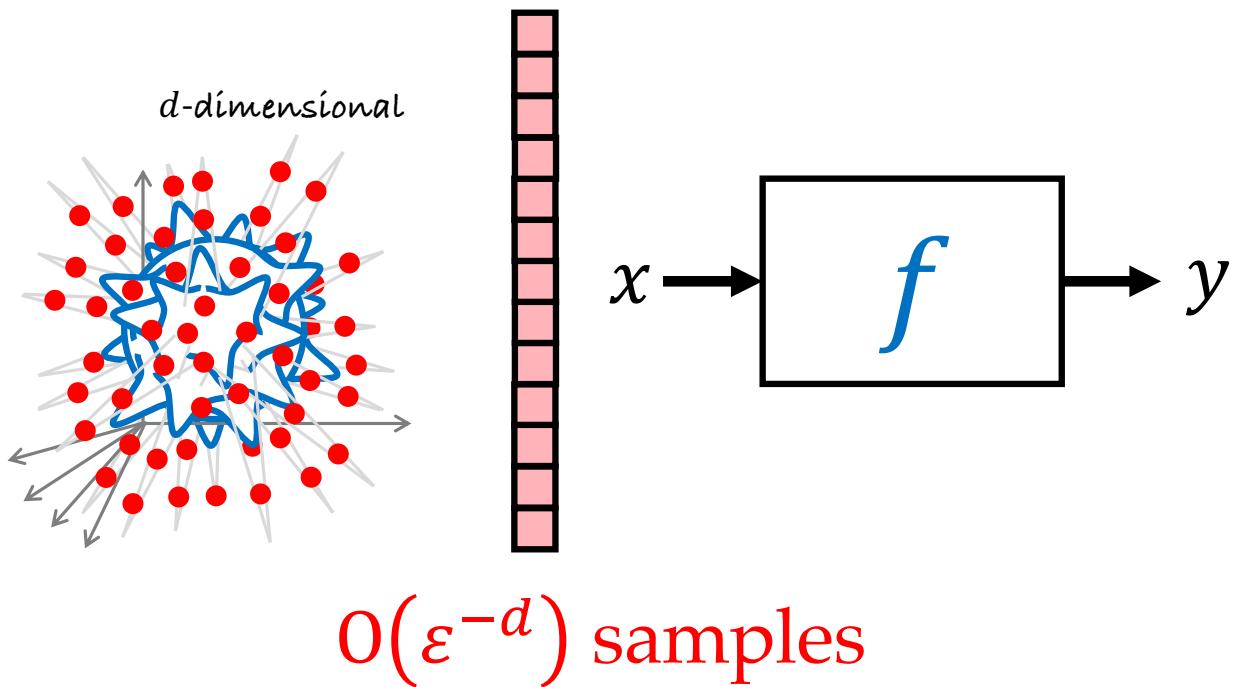


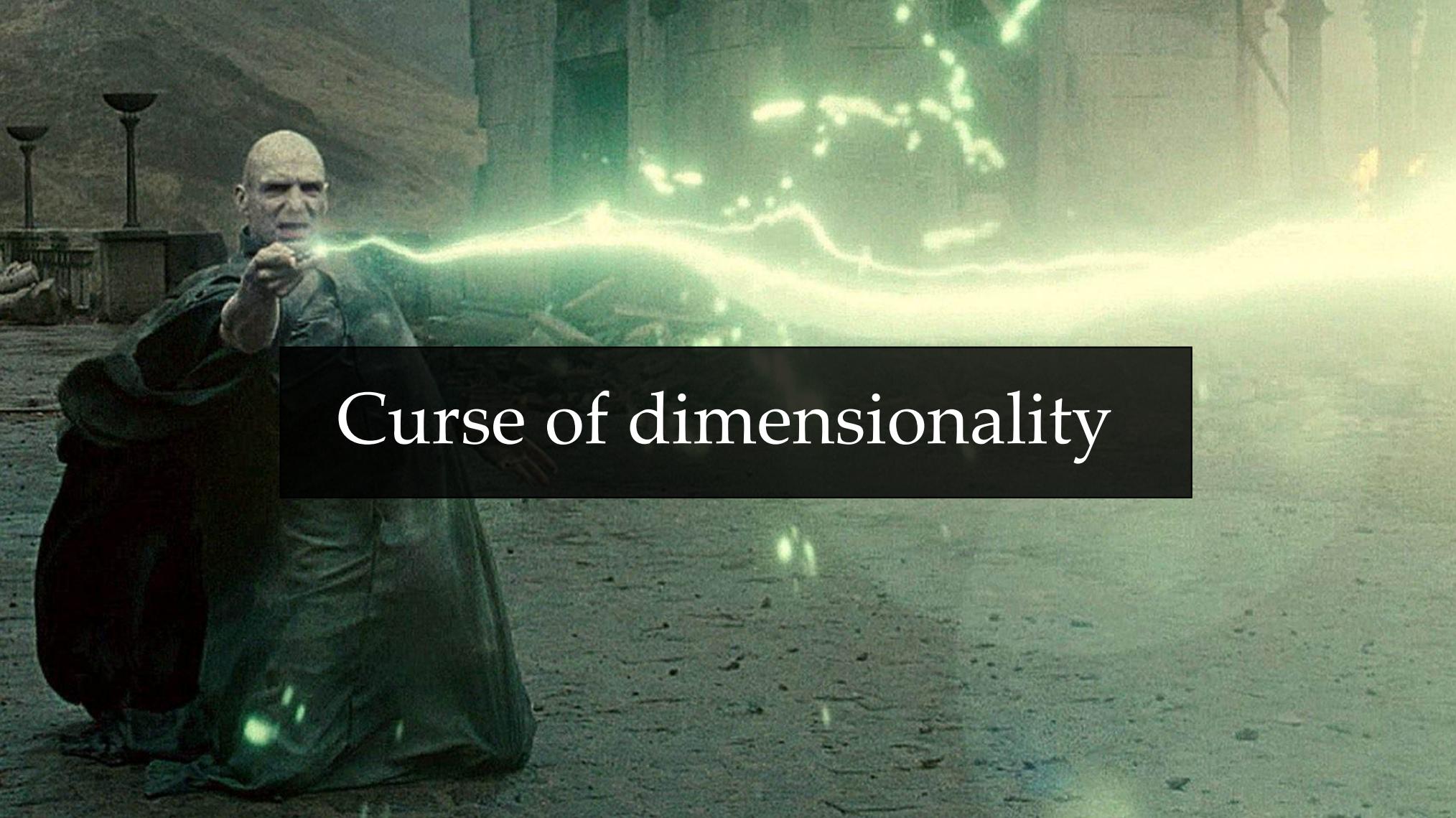
2-dimensional



3-dimensional







Curse of dimensionality

Approximation rates

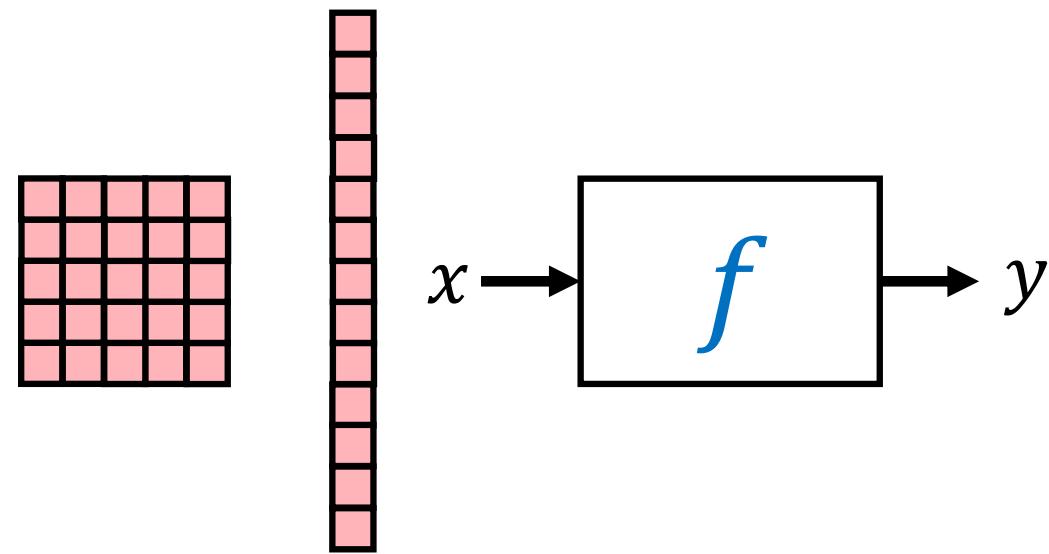
- Bound on the approximation error

$$\varepsilon = \inf_{g \in \mathcal{F}} \sup_{x \in K \subset \mathbb{R}^d} |f(x) - g(x)|$$

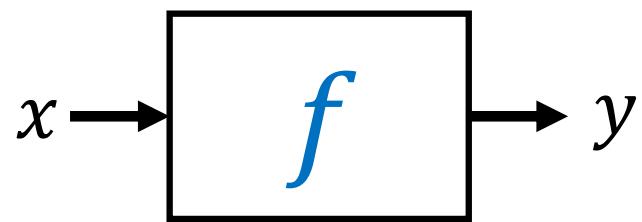
w.r.t. d (input dimension) and m (model complexity) for different classes of functions

- *Sobolev class* $f \in H^s(\mathbb{R}^d) = \left\{ f \in L_2(\mathbb{R}^d) : \int_{\mathbb{R}^d} (1 + \|\omega\|^2)^s |\hat{f}(\omega)|^2 d\omega < \infty \right\}$
error is exponential $\varepsilon = \mathcal{O}(m^{-s/d})$ **dimensionality-cursed!**
- *Barron class* $f \in \left\{ f \in L_2(\mathbb{R}^d) : \int_{\mathbb{R}^d} \|\omega\|^2 |\hat{f}(\omega)|^2 d\omega < \infty \right\}$
error is $\varepsilon = \mathcal{O}(m^{-1})$ **too strong assumption in practice!**

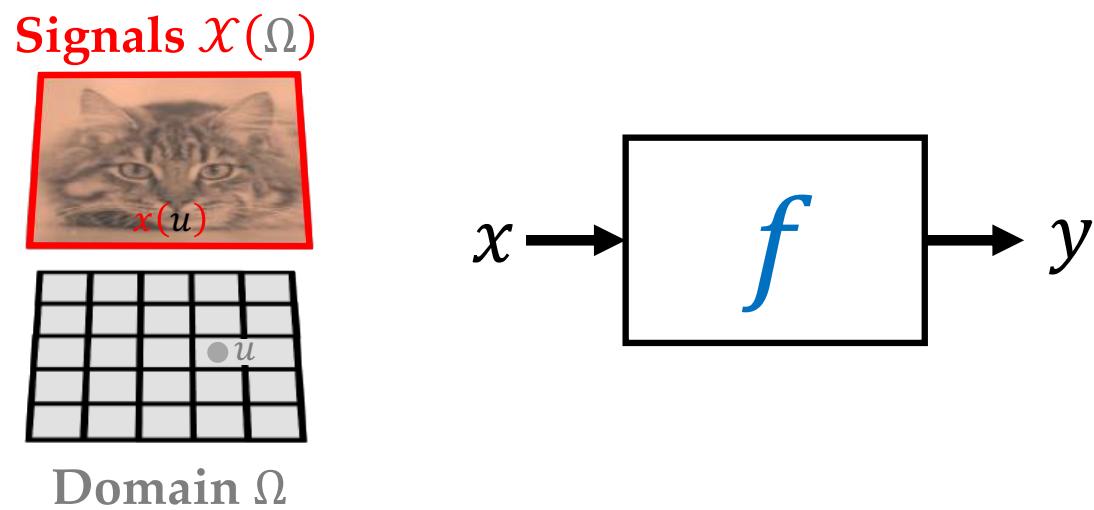
Geometric priors



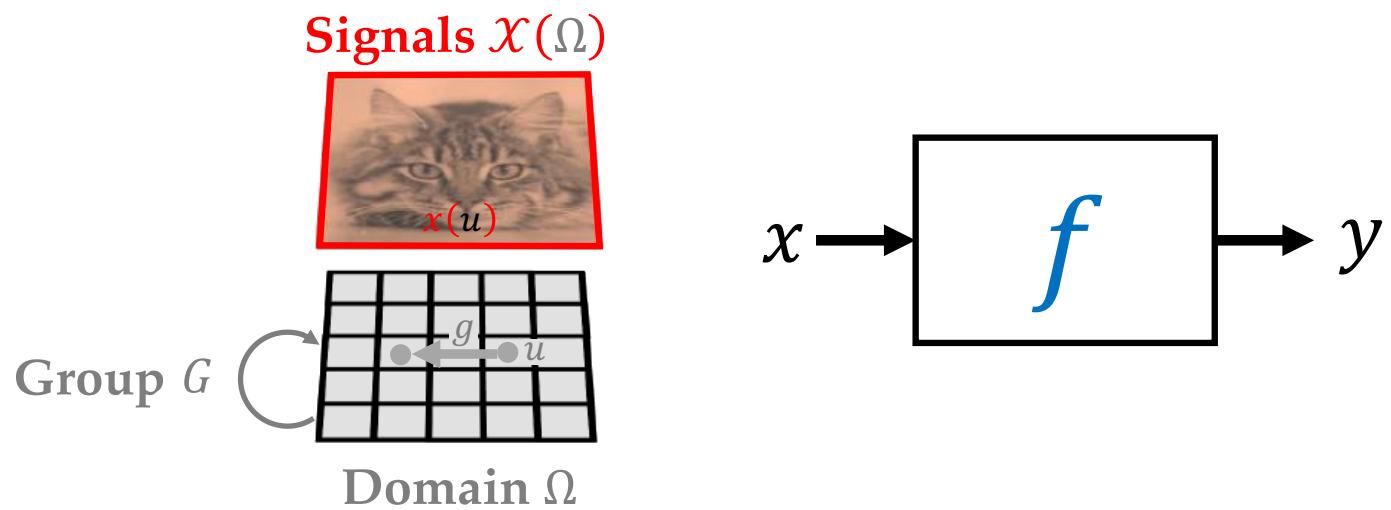
Geometric priors



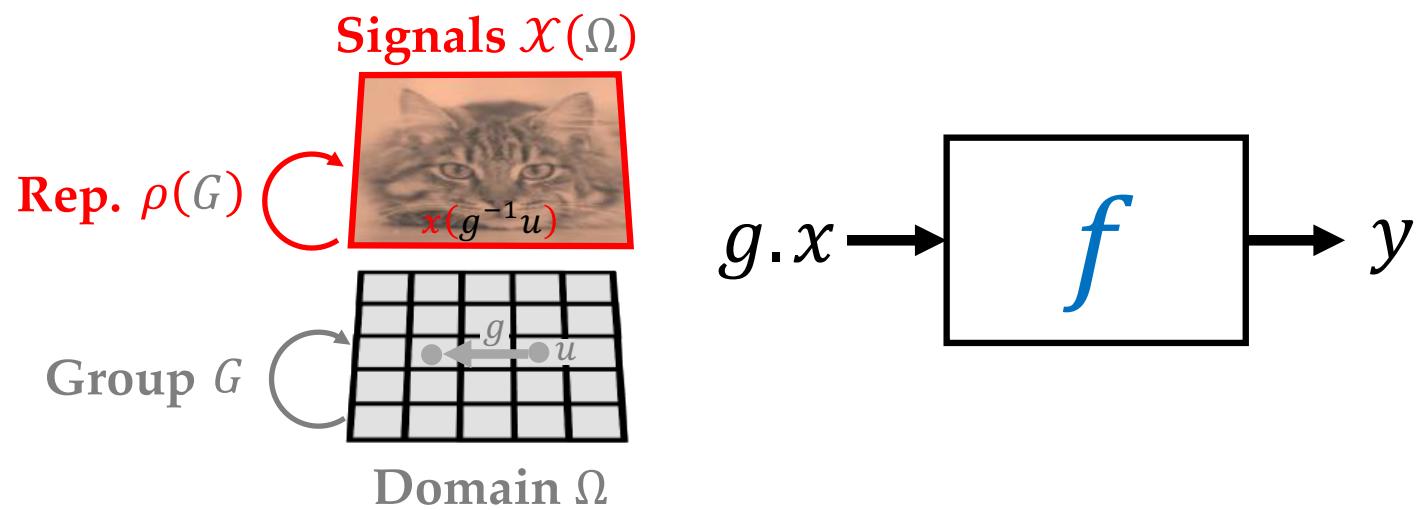
Geometric priors



Geometric priors

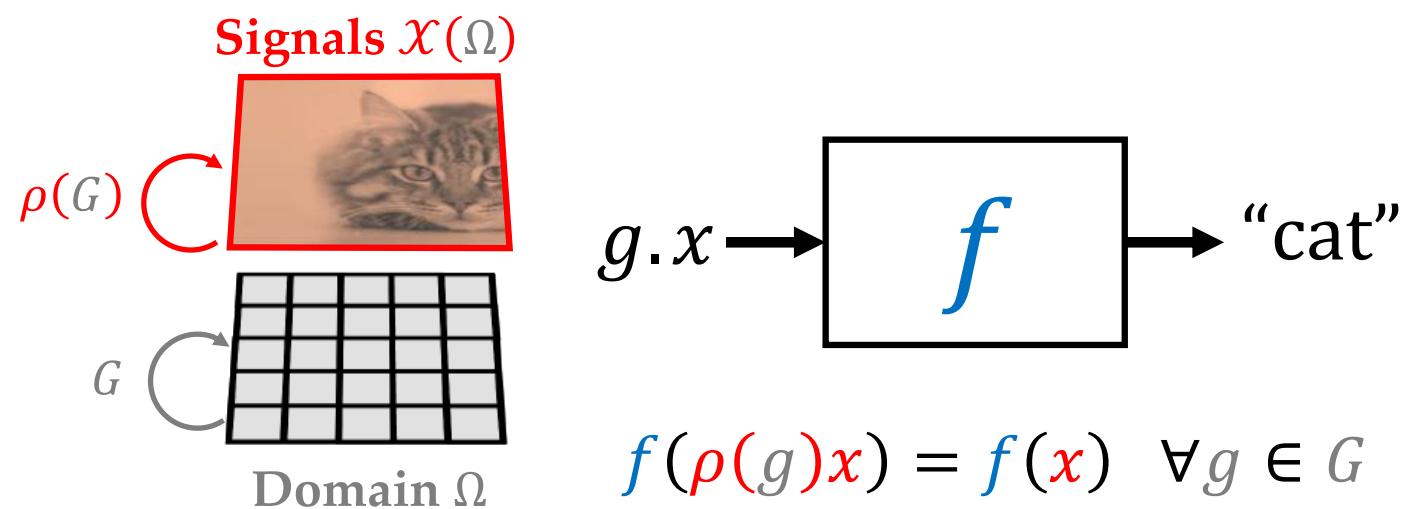


Geometric priors

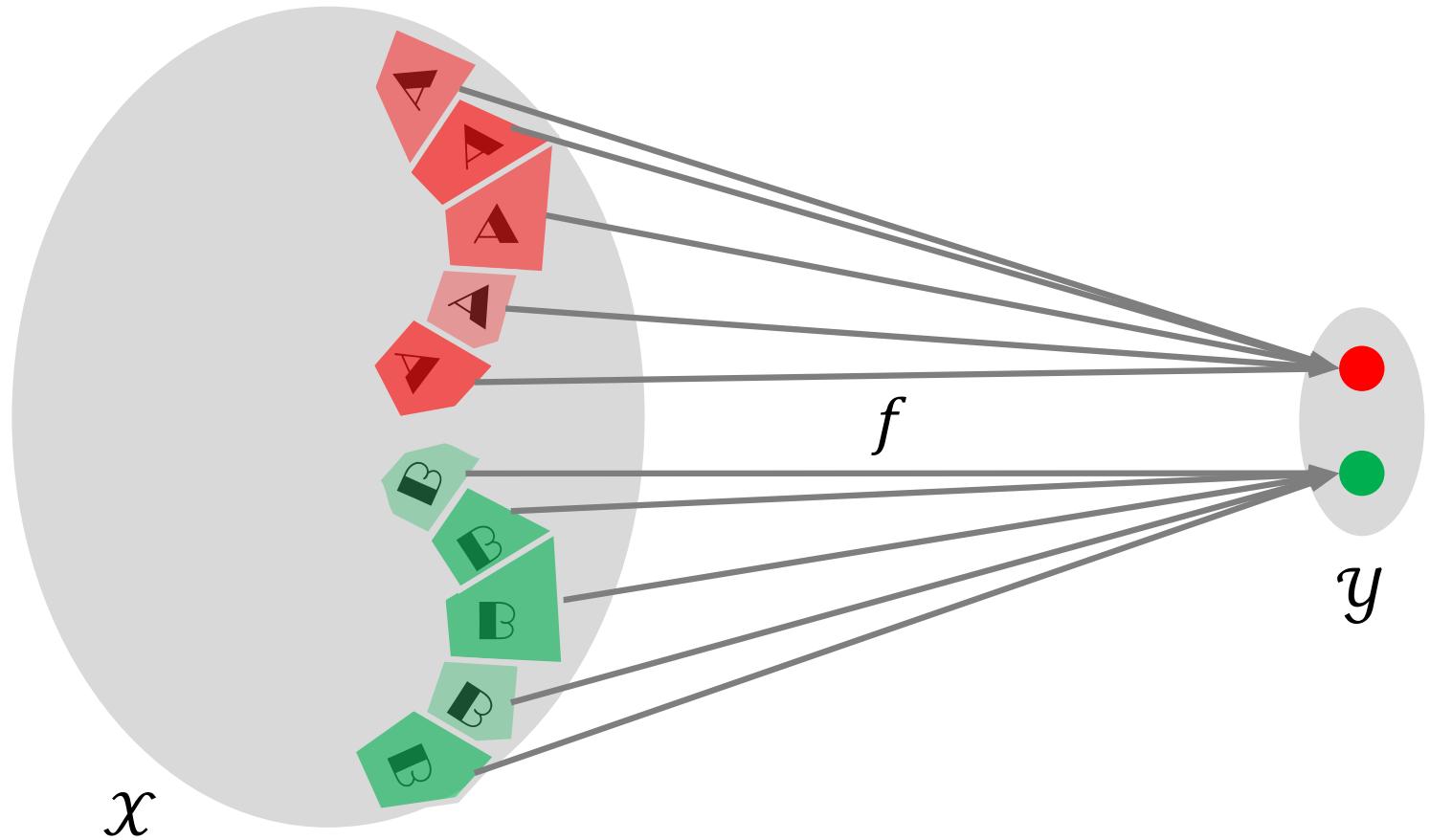


“How f interacts with the group G ? ”

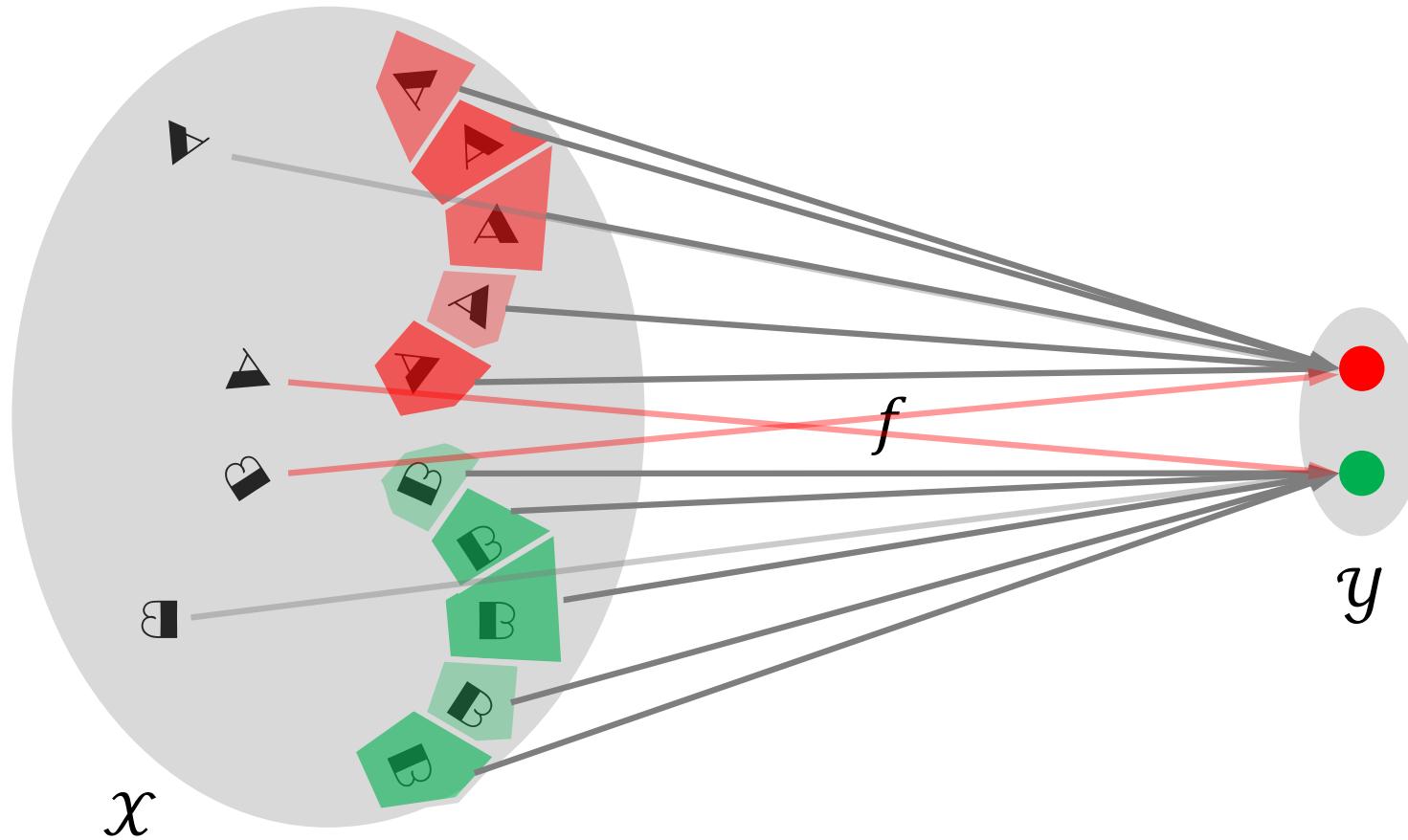
Geometric priors: Invariance



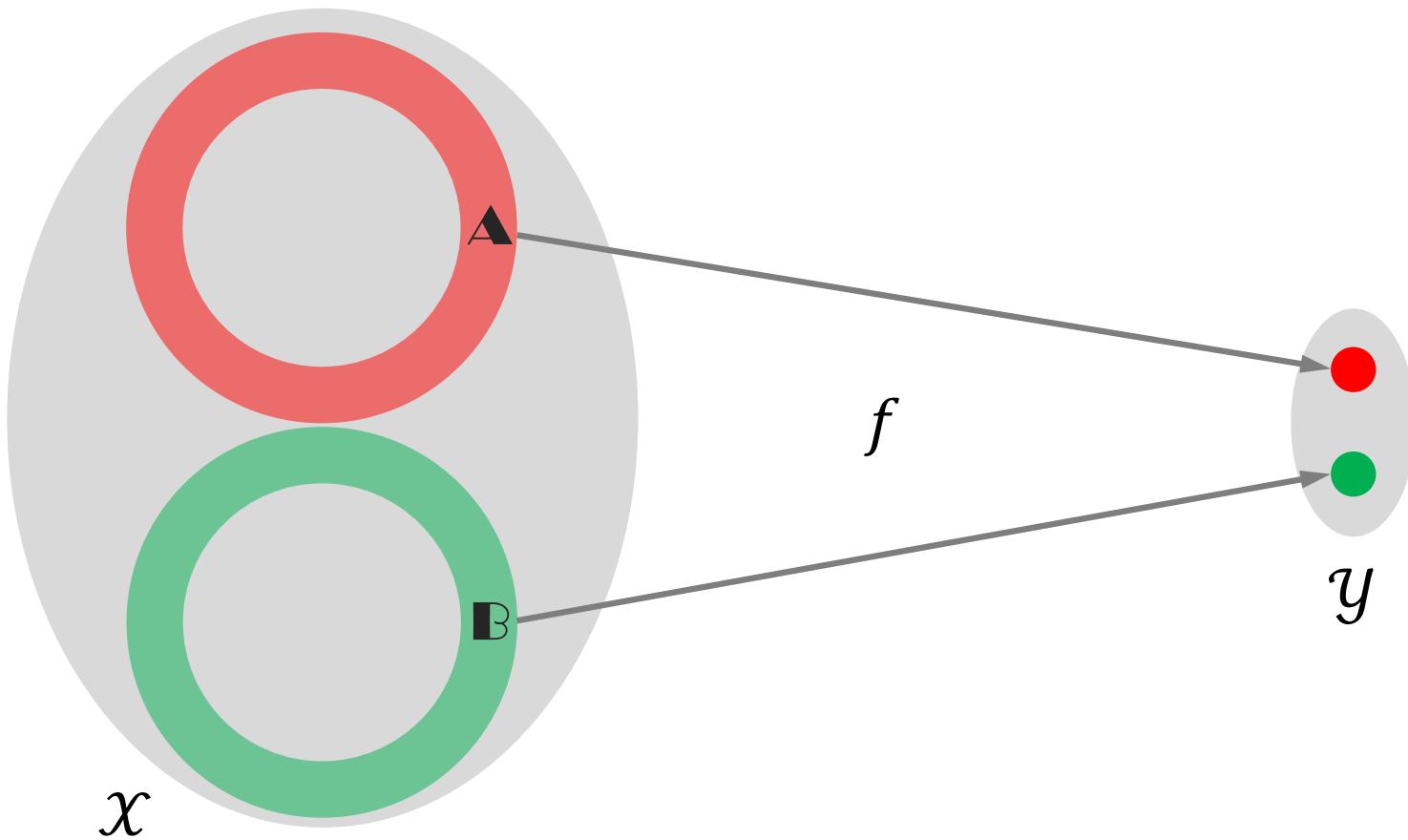
Symmetries of the Label Function



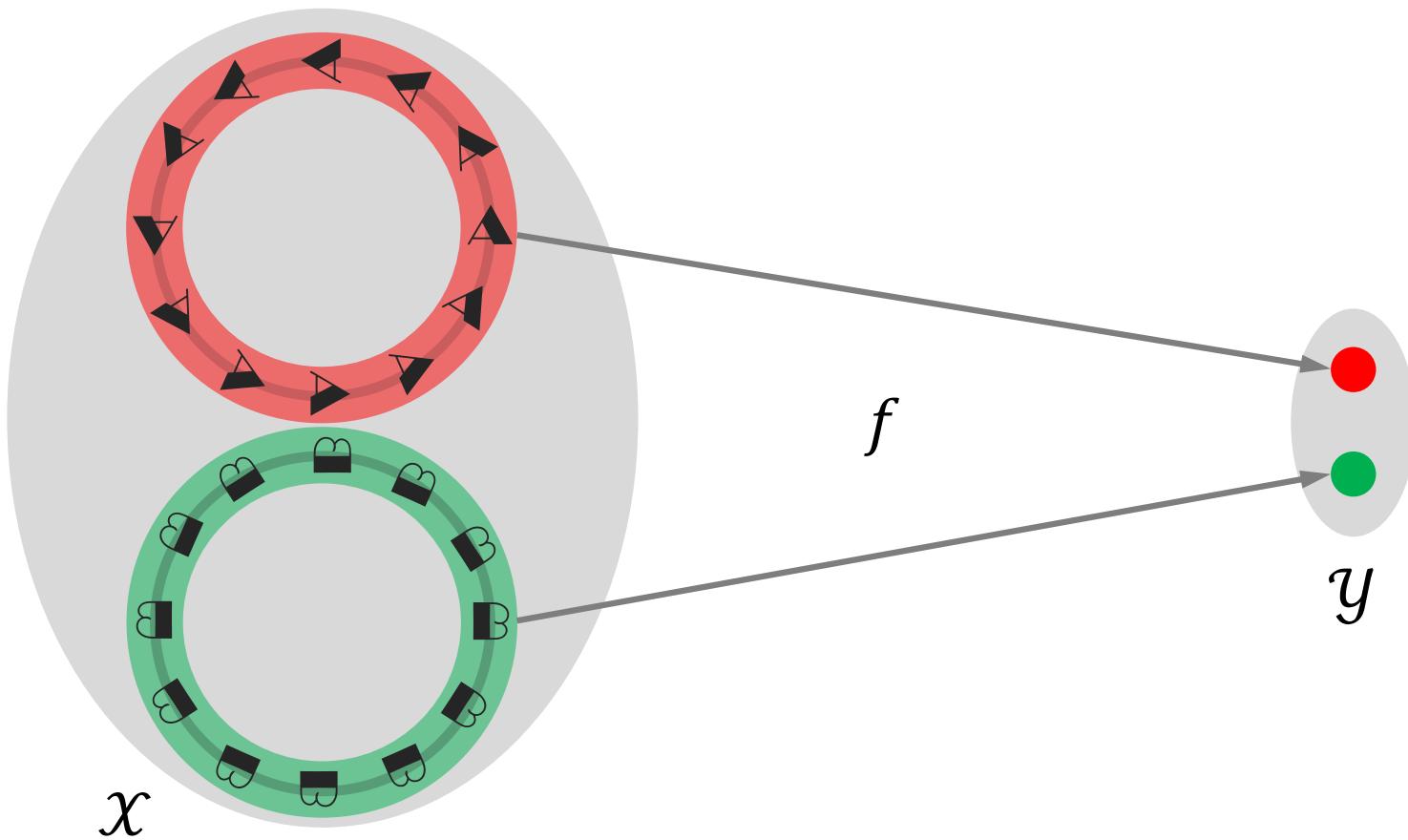
Symmetries of the Label Function



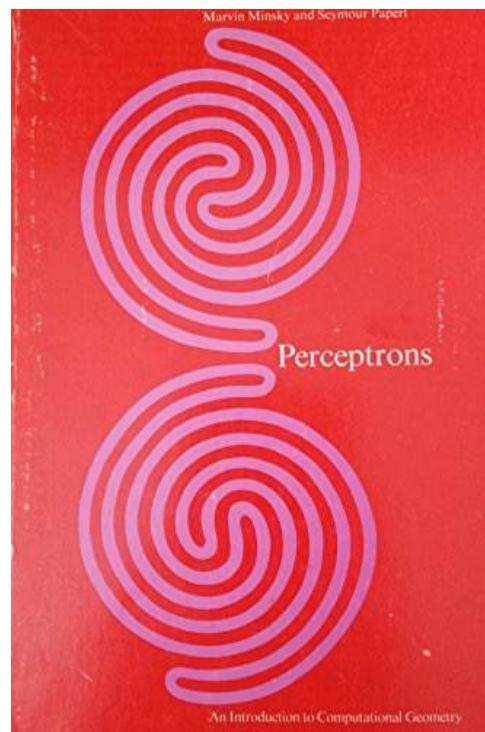
Symmetries of the Label Function



Symmetries of the Label Function



First “geometric” machine learning



M. Minsky S. Papert

1969

Minsky, Papert 1969

First “geometric” machine learning

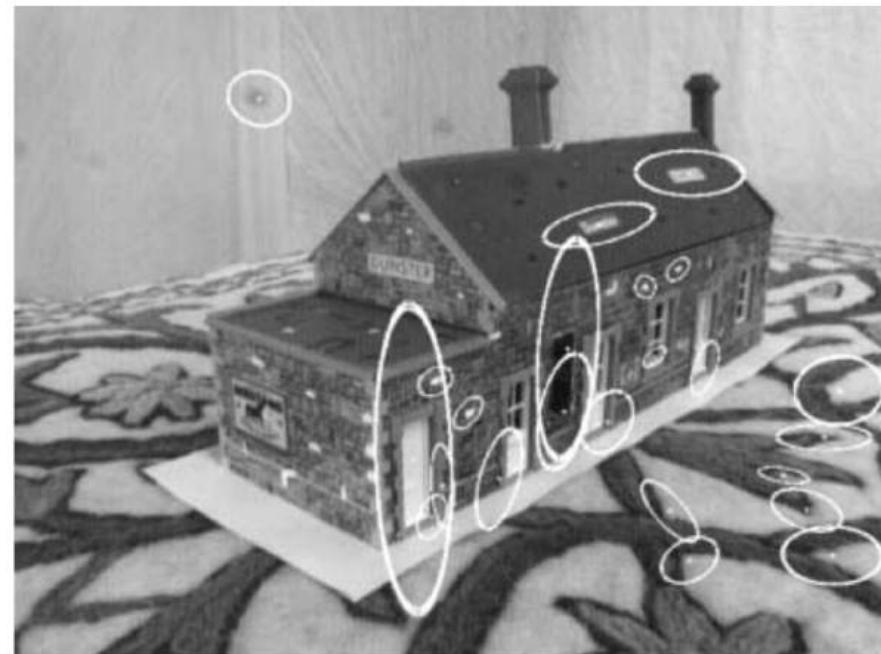
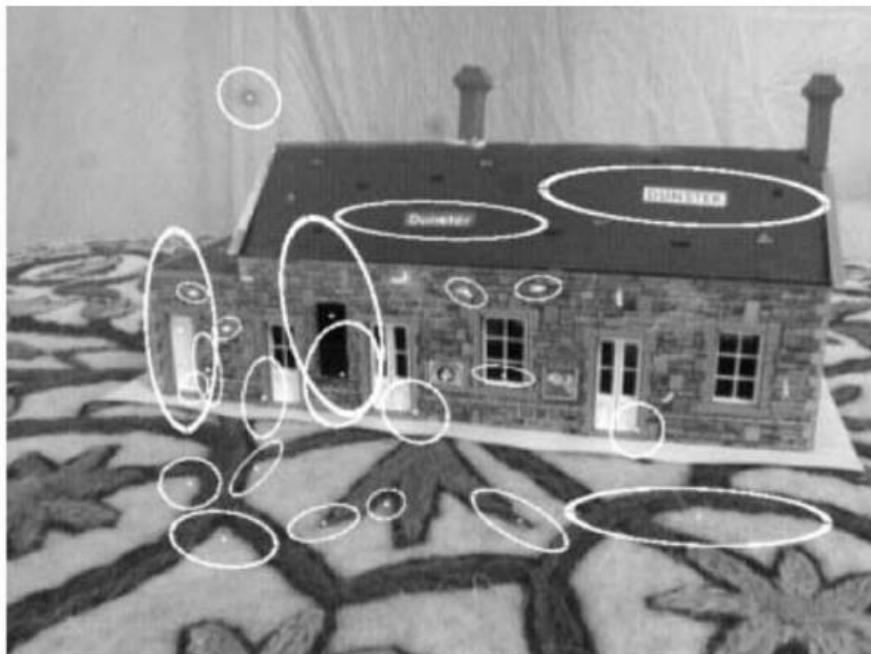
Group Invariance Theorem: “if a neural network is invariant to a group, then its output can be expressed as functions of the orbits of the group”



M. Minsky S. Papert

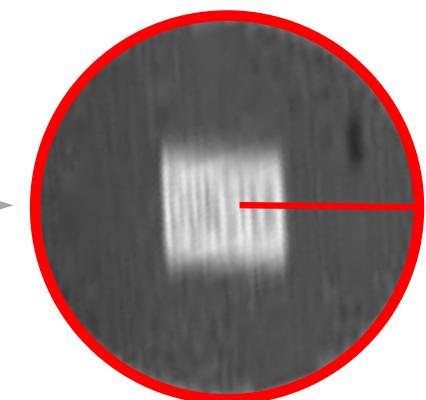
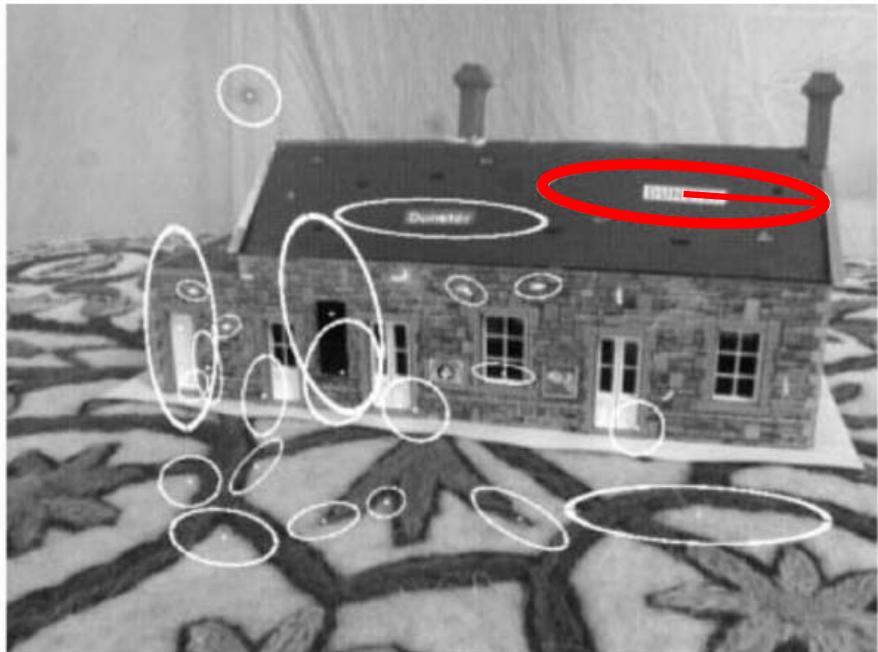
1969

Canonisation

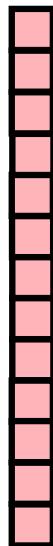


Mikolajczyk, Schmid 2004

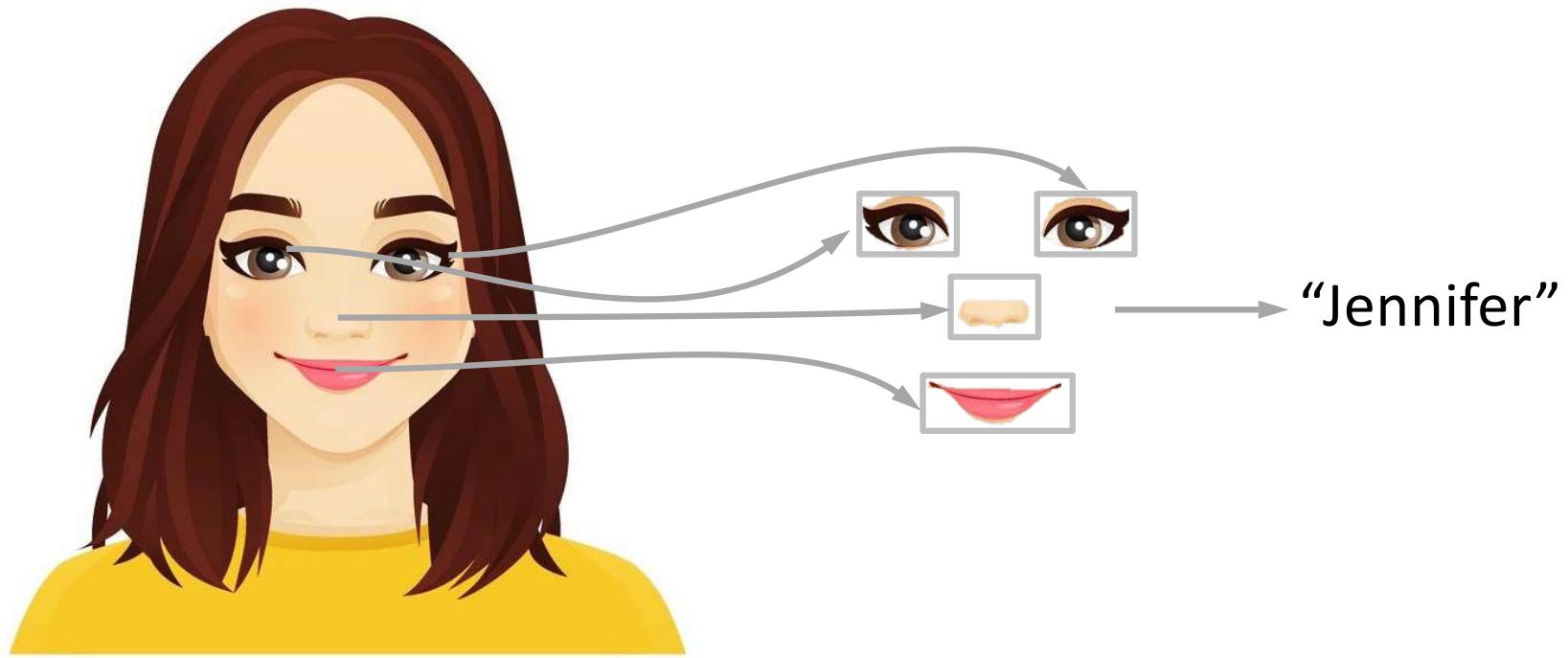
Canonisation



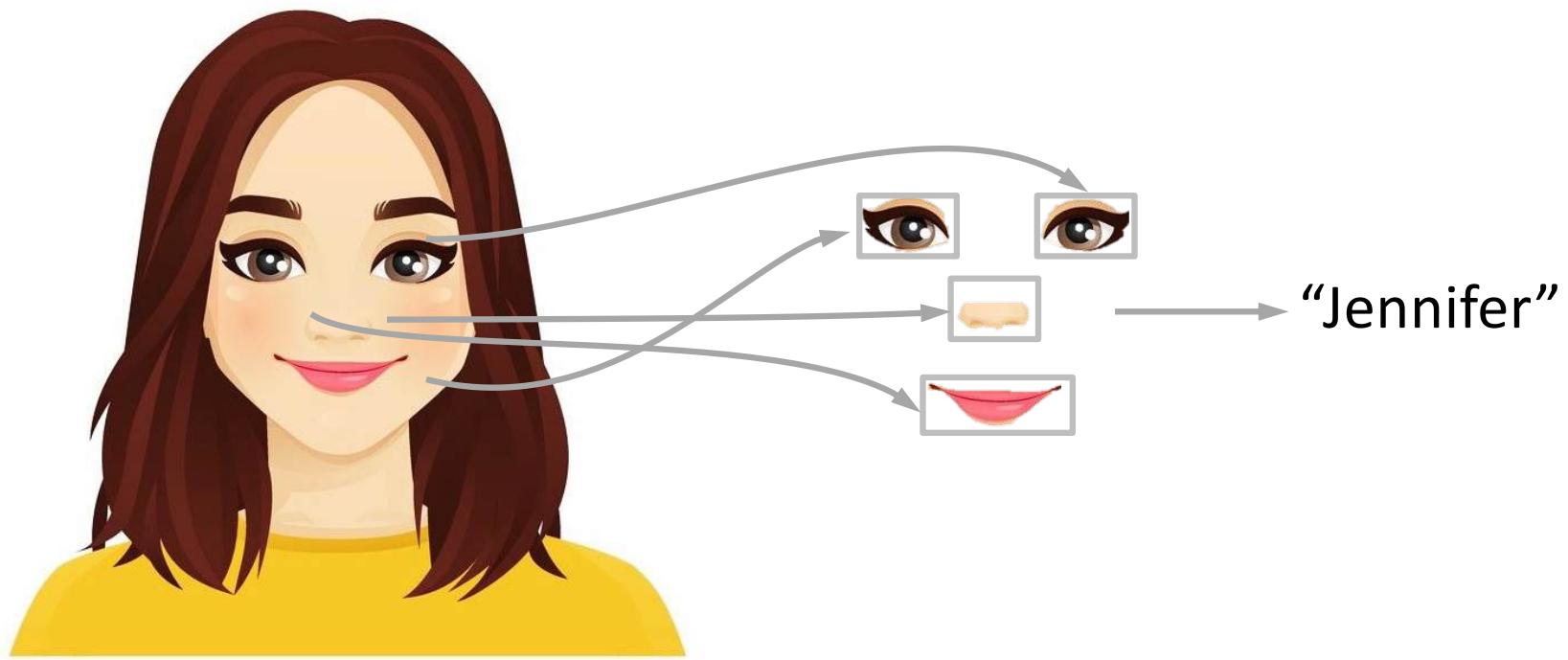
canonisation



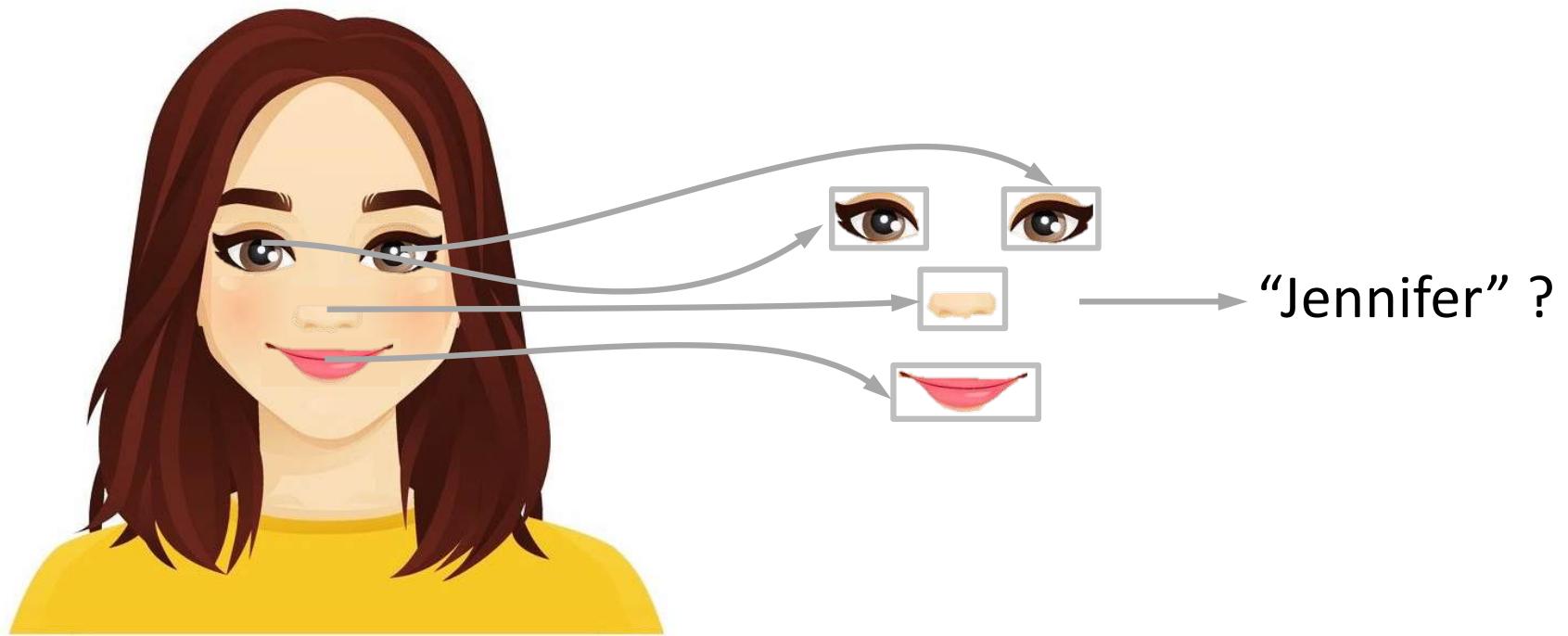
Canonisation

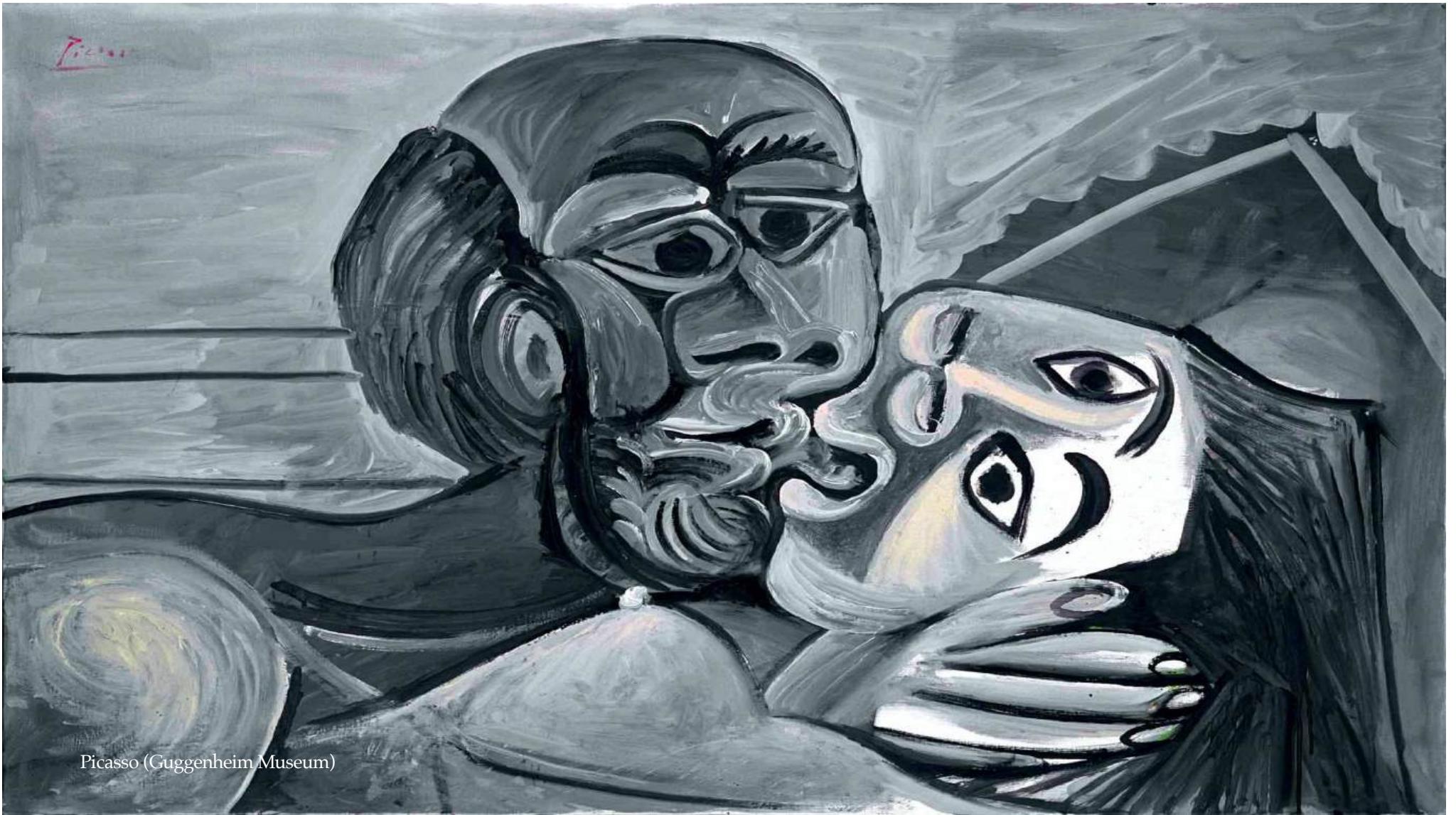


Canonisation



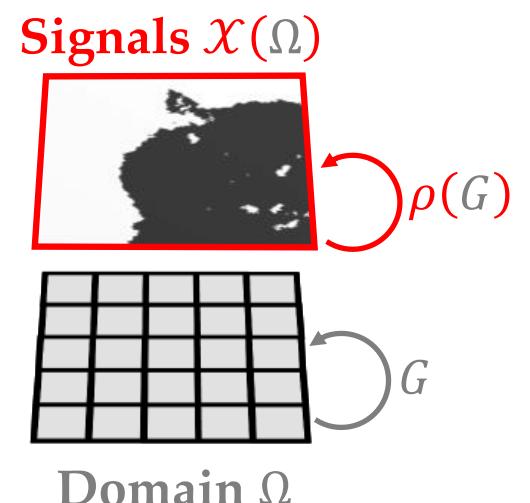
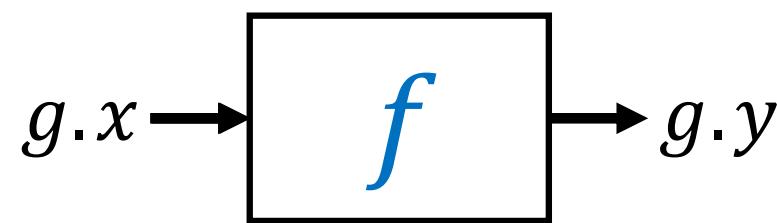
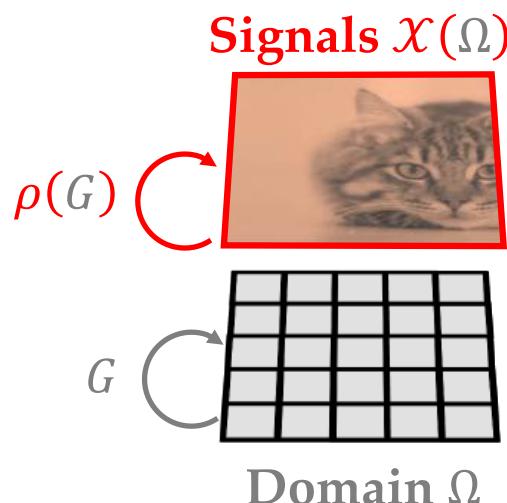
Canonisation



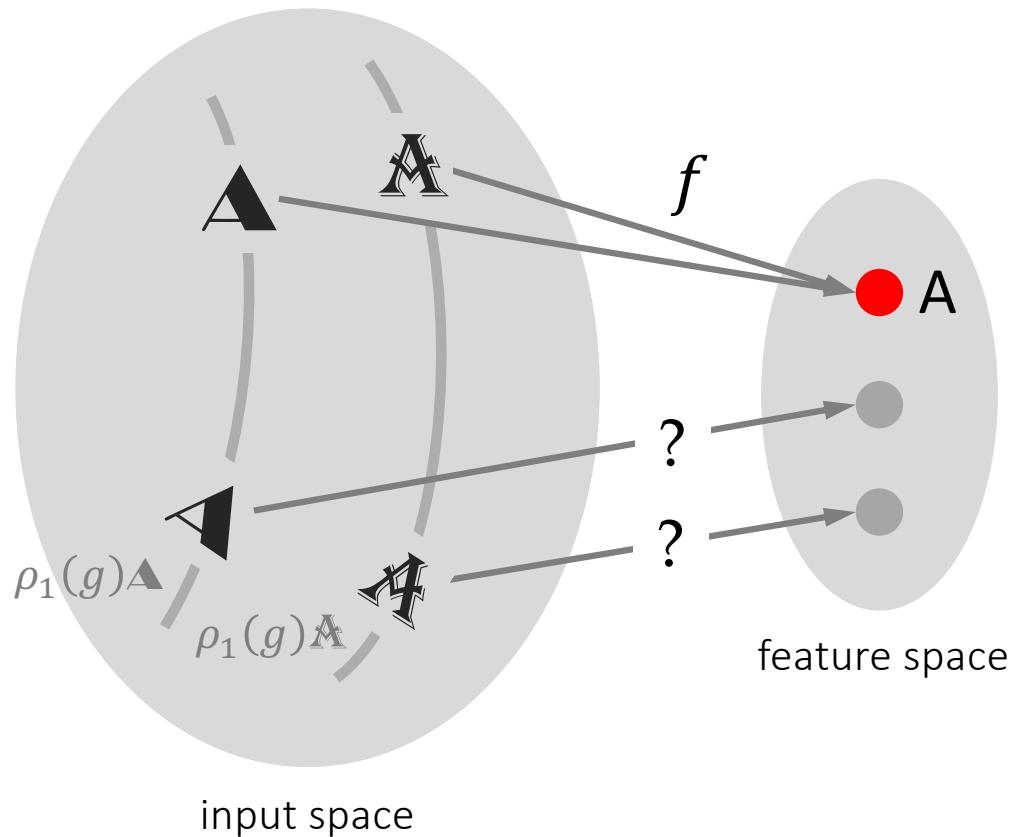


Picasso (Guggenheim Museum)

Geometric priors: Equivariance



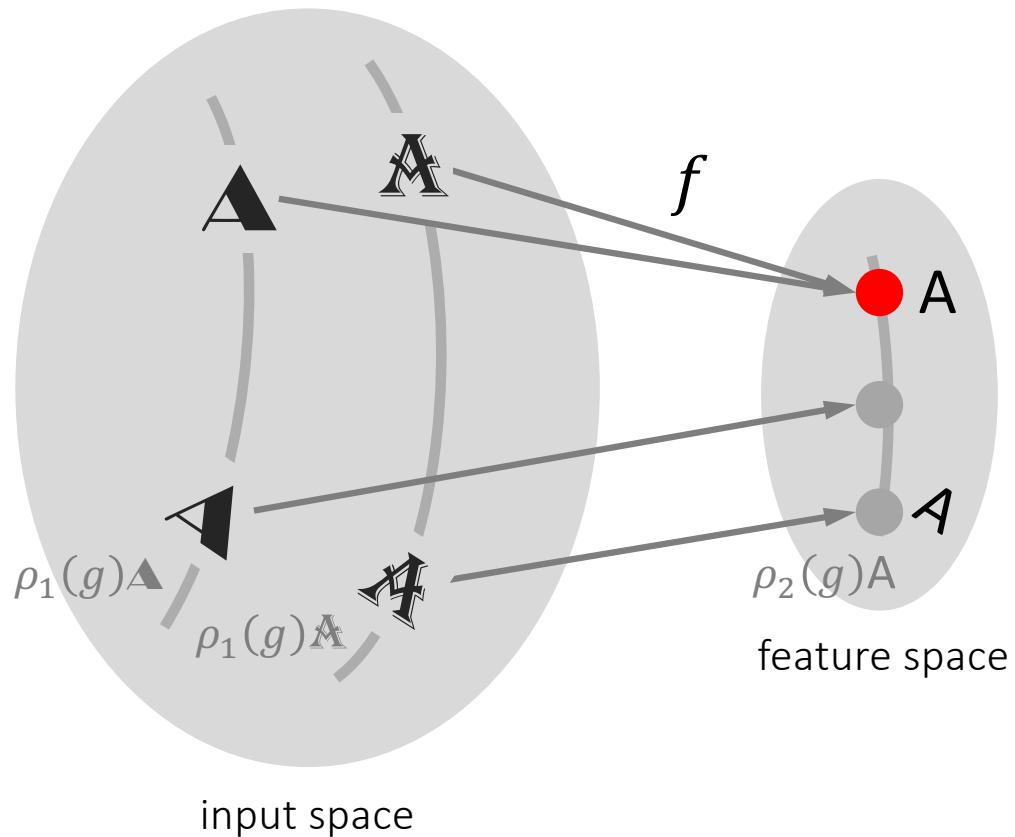
Equivariance = Symmetry-consistent generalisation



$$f(\rho_1(g)\mathbf{A}) = \rho_2(g) f(\mathbf{A})$$

$$f(\rho_1(g)\mathbf{A}) = \rho_2(g) f(\mathbf{A})$$

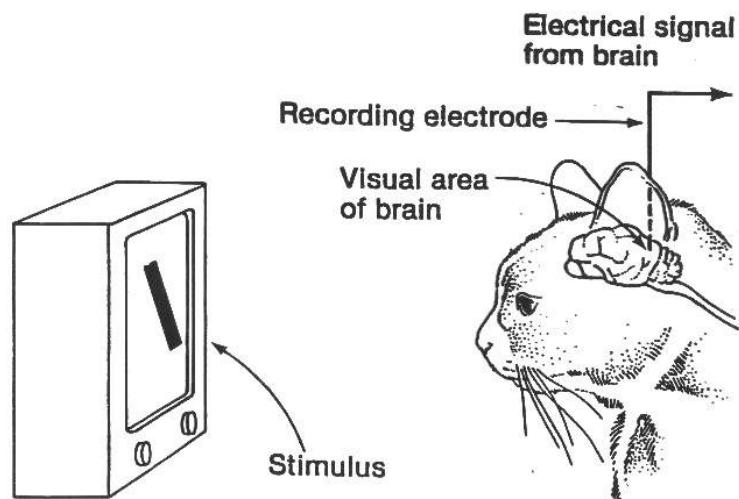
Equivariance = Symmetry-consistent generalisation



$$f(\rho_1(g)\mathbb{A}) = \rho_2(g) f(\mathbb{A})$$

$$f(\rho_1(g)\mathbb{A}) = \rho_2(g) f(\mathbb{A})$$

Early Geometric Architectures

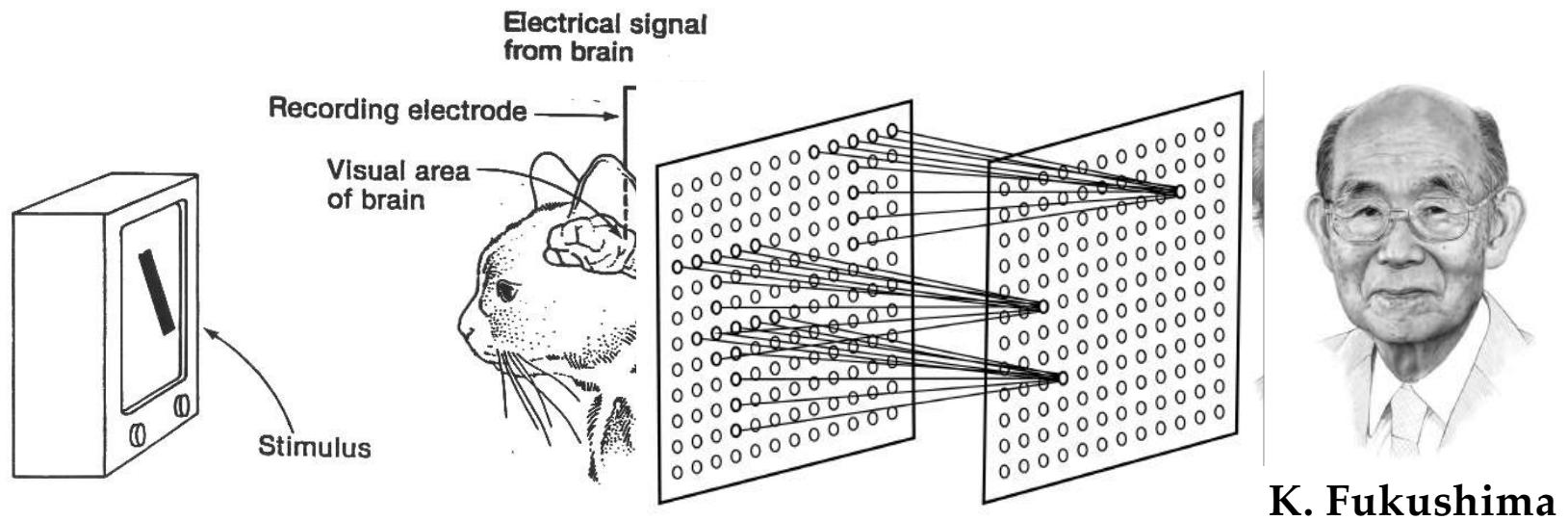


D. Hubel T. Wiesel

1959

Hubel, Wiesel 1959, 1962; Portraits: Ihor Gorskyi

Early Geometric Architectures



1959 1980

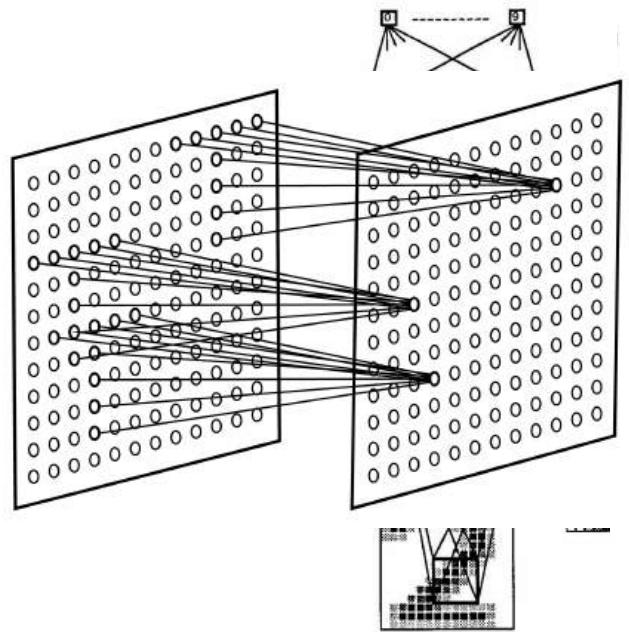
Hubel, Wiesel 1959, 1962; Fukushima 1980; Portraits: Ihor Gorskyi

Early Geometric Architectures



D. Hubel T. Wiesel

1959

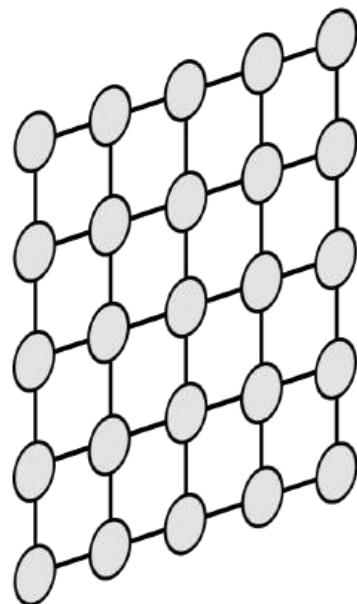


K. Fukushima

1980

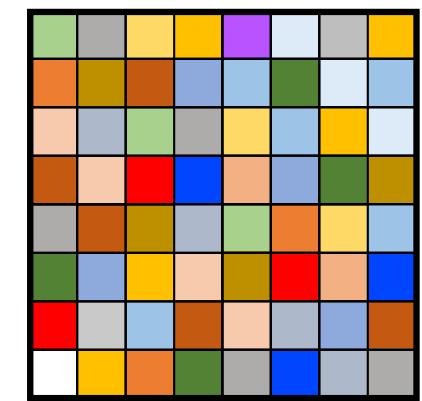
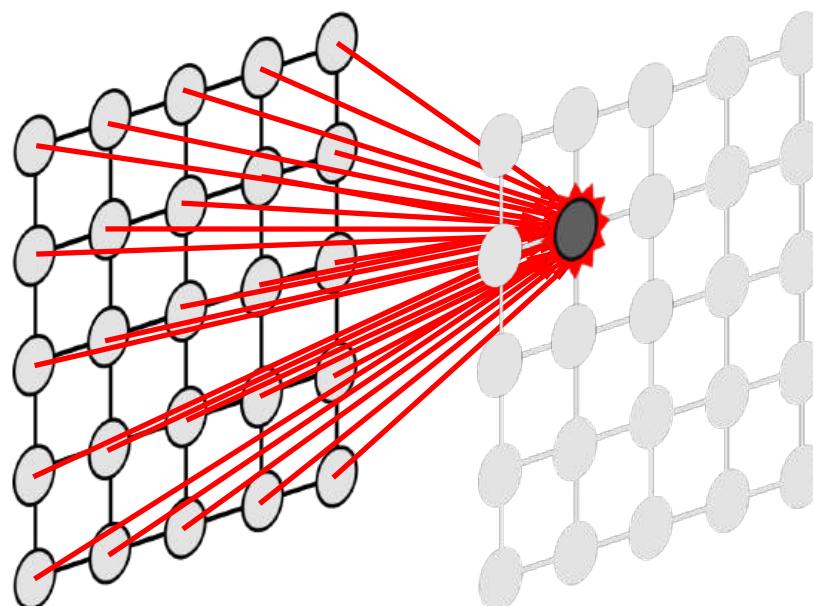
Hubel, Wiesel 1959, 1962; Fukushima 1980; LeCun et al. 1989; Portraits: Ihor Gorskyy

Multi-Layer Perceptron



LeCun et al. 1989

Multi-Layer Perceptron

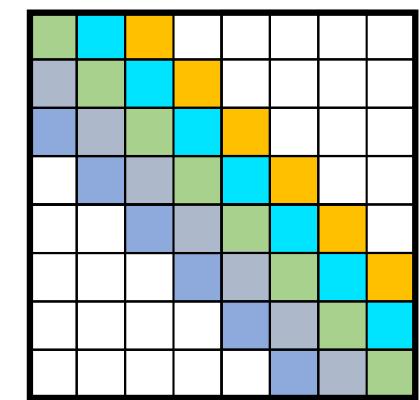
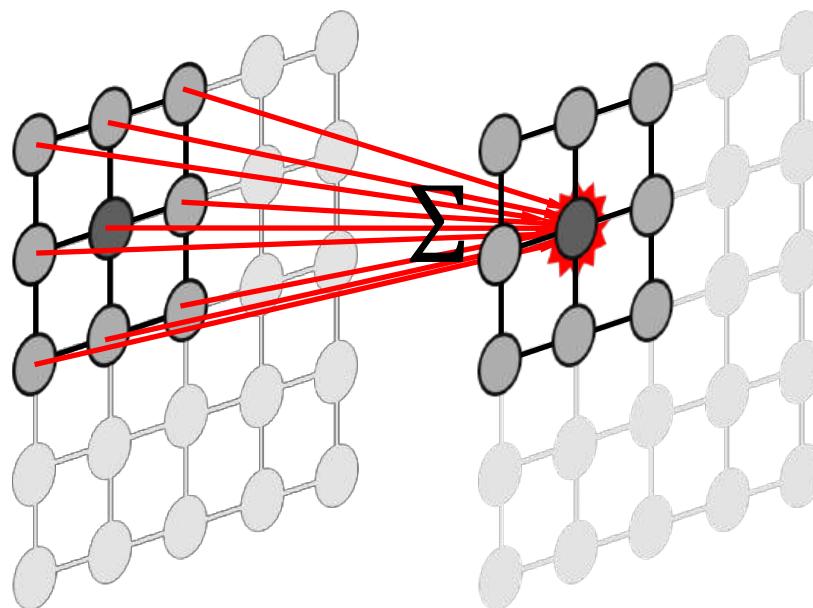


weight matrix
 $\mathcal{O}(n^2)$ DOF

Fully connected layer

LeCun et al. 1989

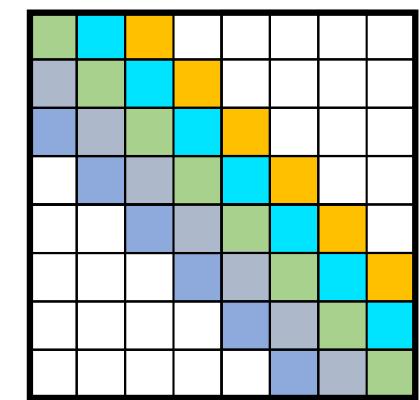
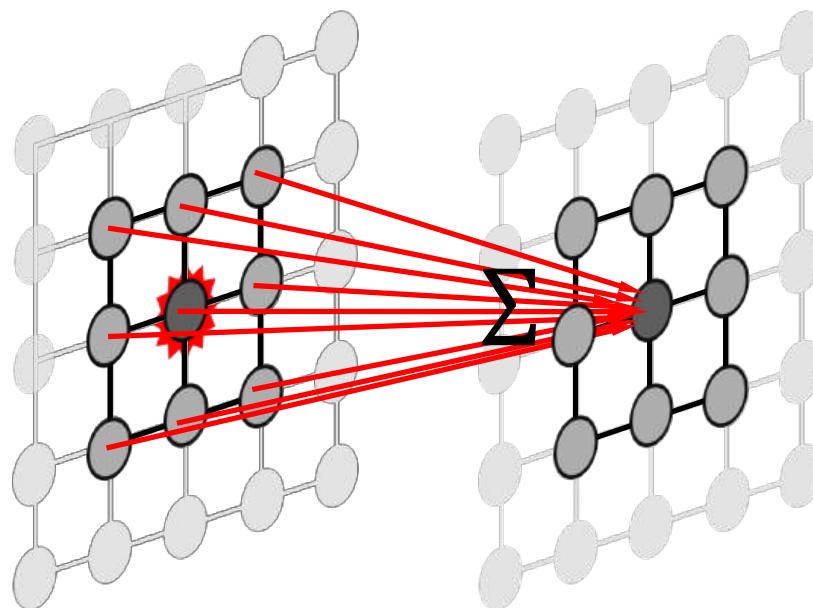
Convolutional Neural Networks



Locality + Shared parameters

LeCun et al. 1989

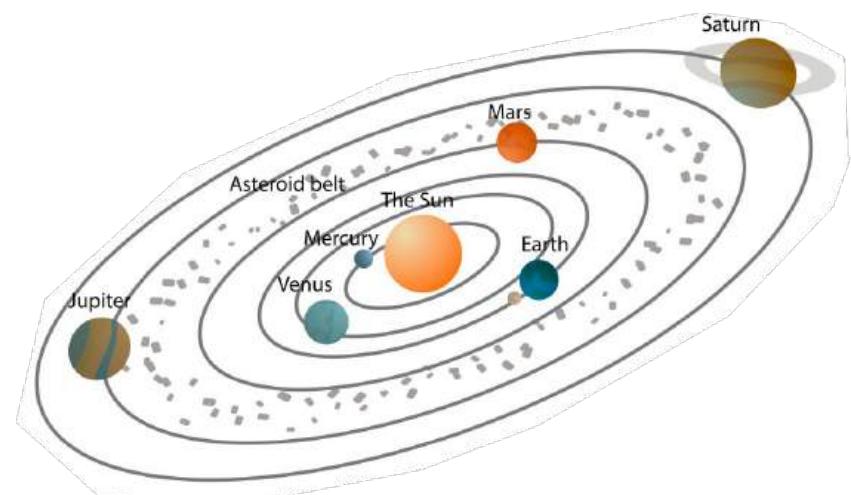
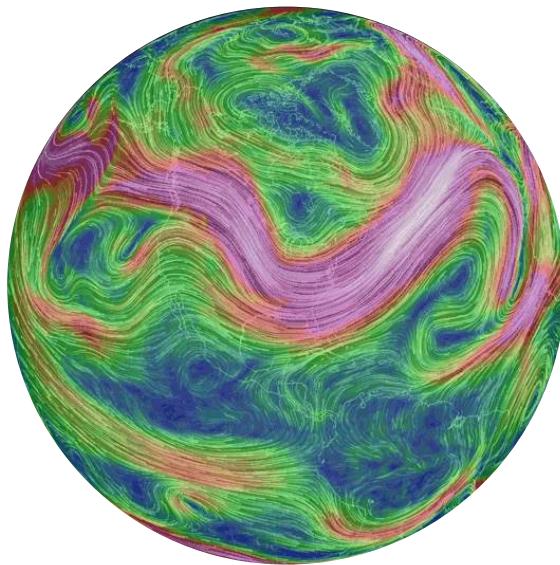
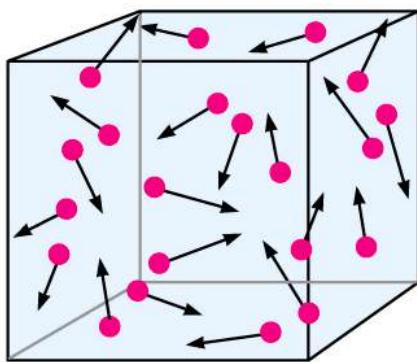
Convolutional Neural Networks



Locality + Shared parameters

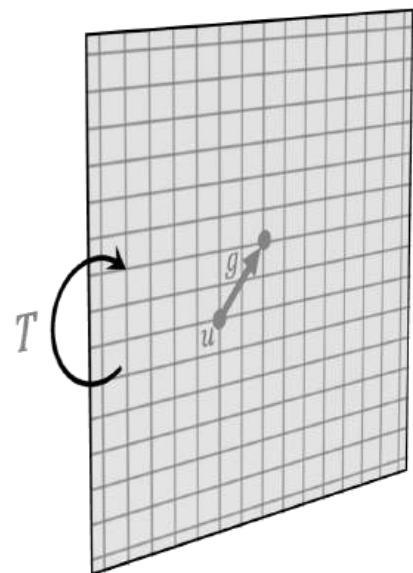
LeCun et al. 1989

Scale Separation



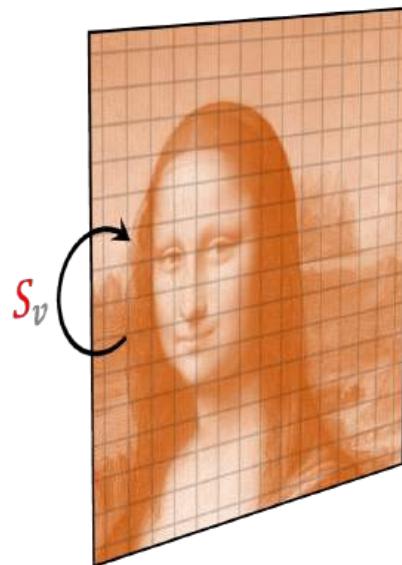
Convolutional Neural Networks

Plane \mathbb{R}^2



Translation group $T(2)$

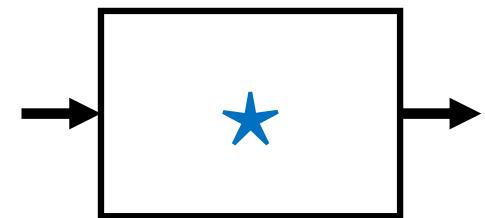
Images $\mathcal{X}(\mathbb{R}^2)$



Shift operator S

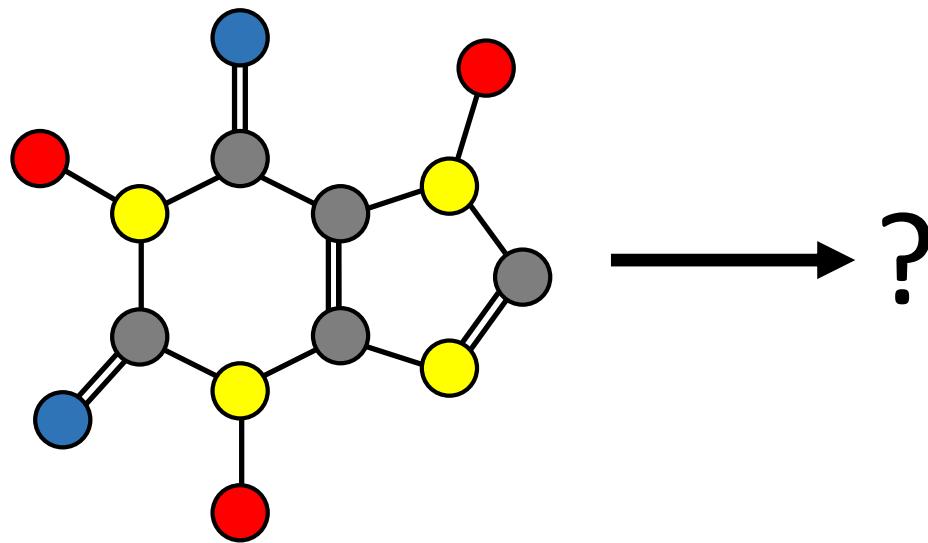
$$S_v x(u) = x(u - v)$$

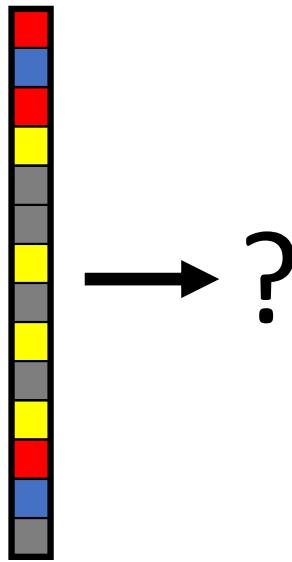
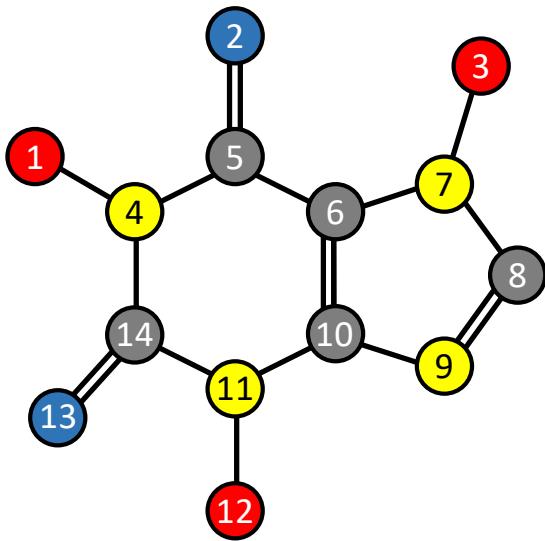
Functions $\mathcal{F}(\mathcal{X}(\mathbb{R}^2))$

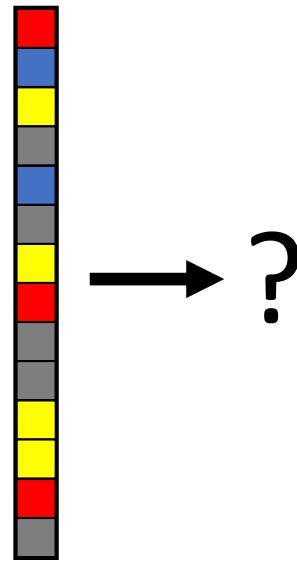
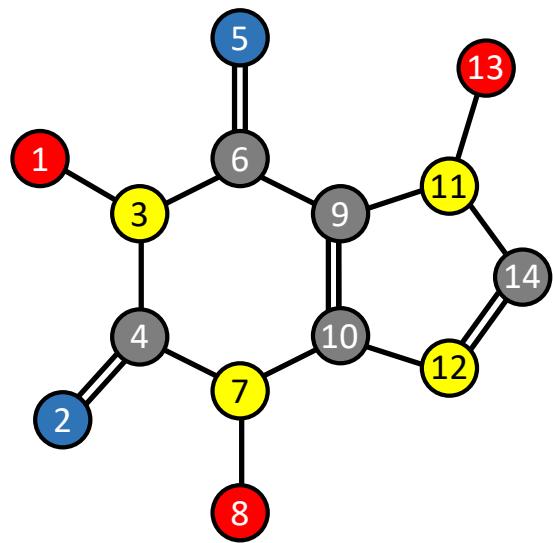


Convolution

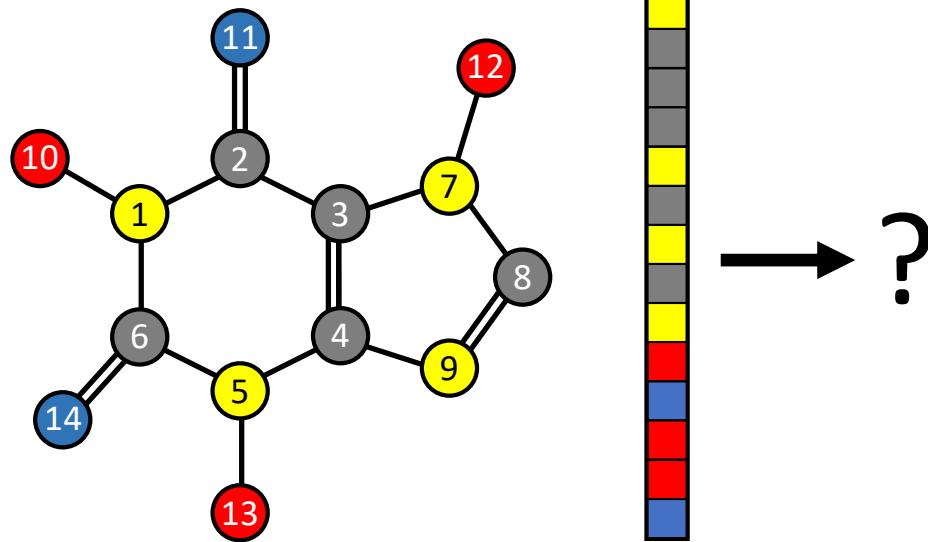
$$(Sx \star y) = S(x \star y)$$





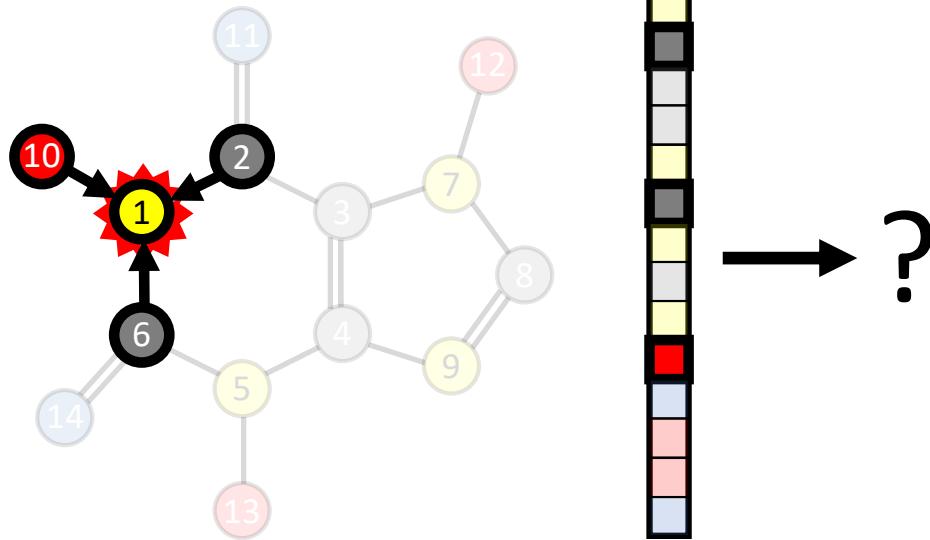


Permutation Invariance



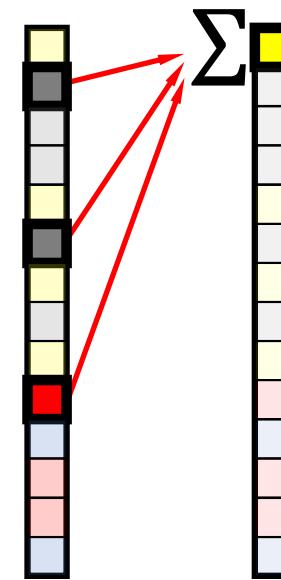
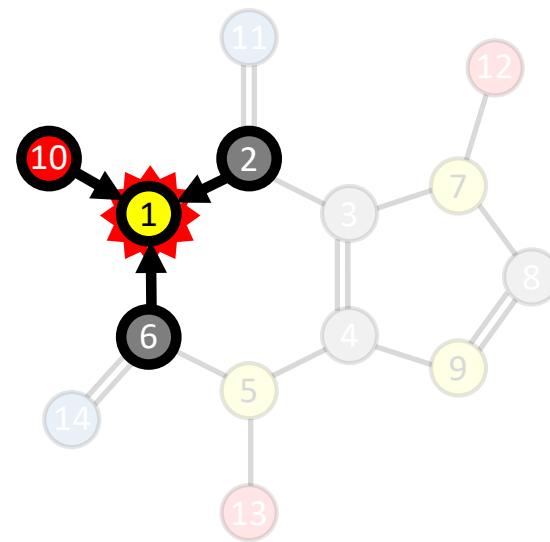
“properties of a molecule do not change if we reorder the atoms”

Graph Neural Networks



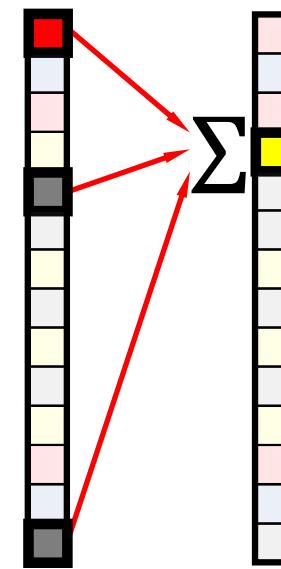
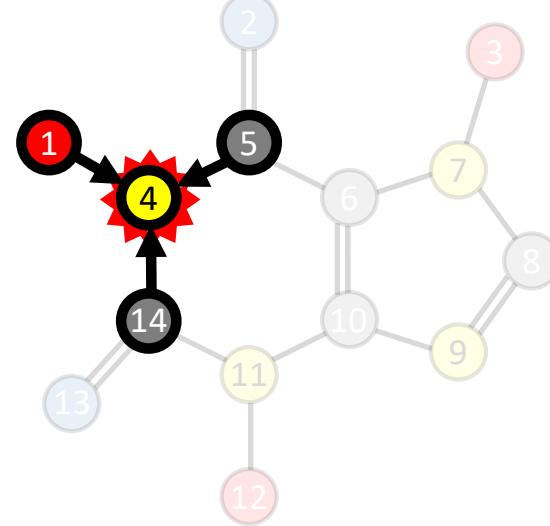
Locality + Shared parameters

Graph Neural Networks



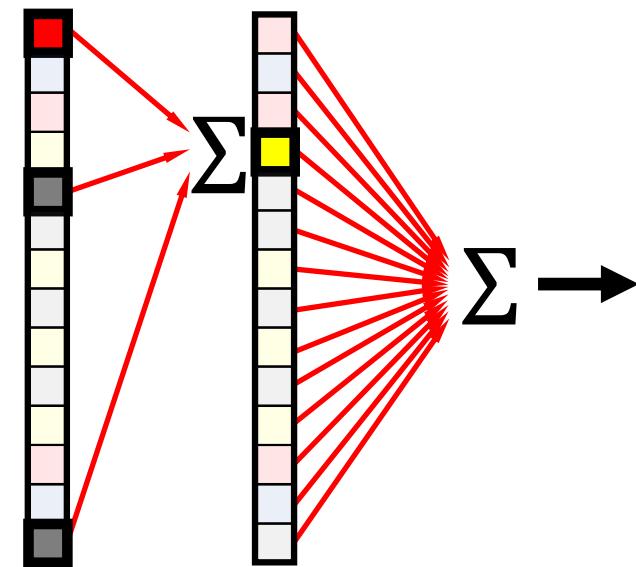
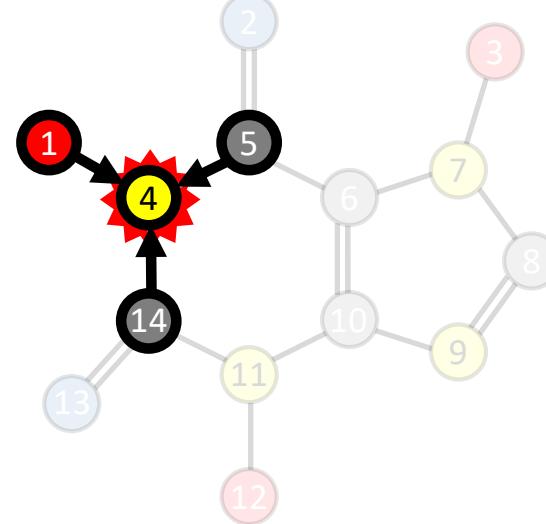
Permutation-
equivariant layer

Graph Neural Networks



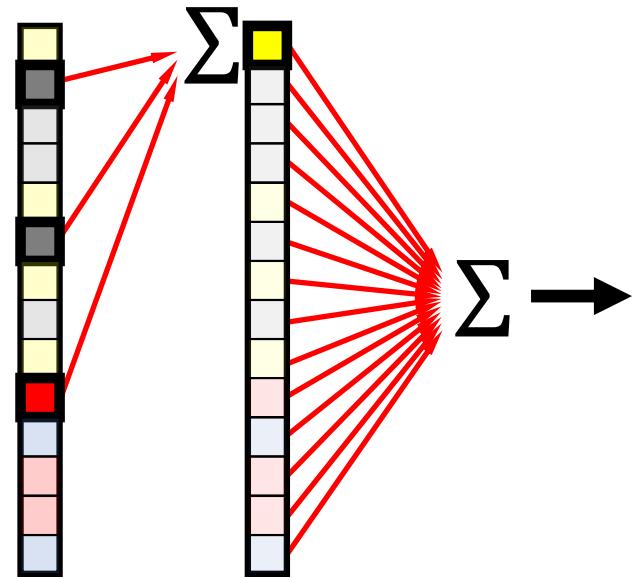
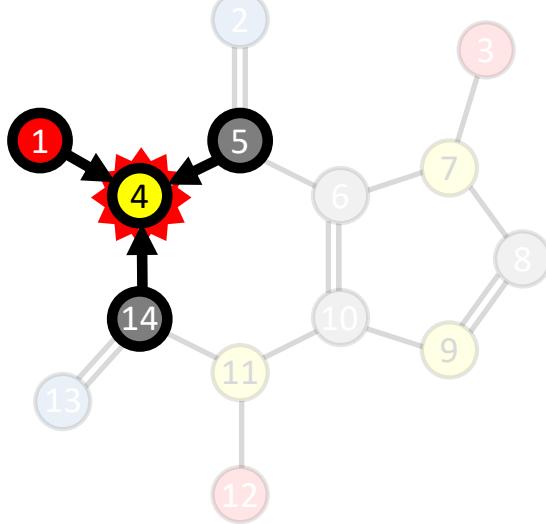
Permutation-
equivariant layer

Graph Neural Networks



Permutation-
invariant readout

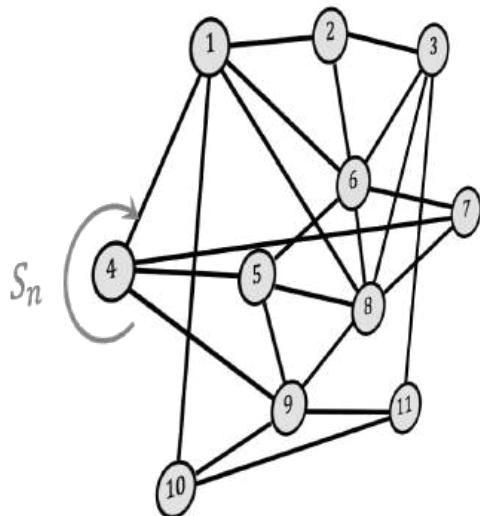
Graph Neural Networks



Permutation-
invariant readout

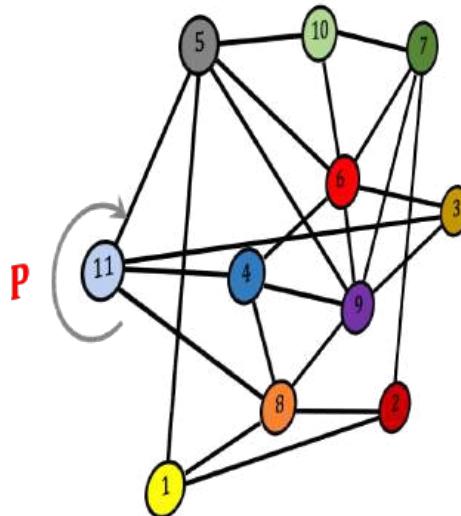
Graph Neural Networks

Graph $G = (V, E)$



Permutation group S_n

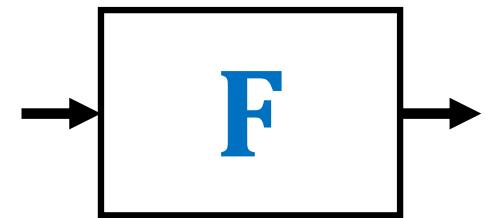
Node features $\mathcal{X}(G)$



Permutation matrix P

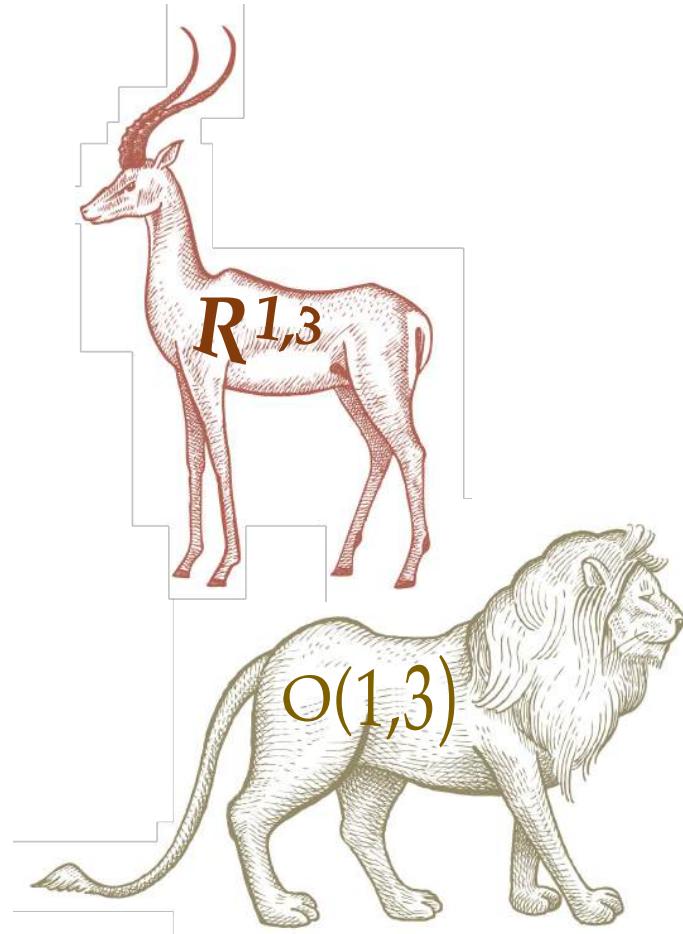
$$\mathbf{P}\mathbf{X} = (x_{\pi^{-1}(i),j})$$

Functions $\mathcal{F}(\mathcal{X}(G))$

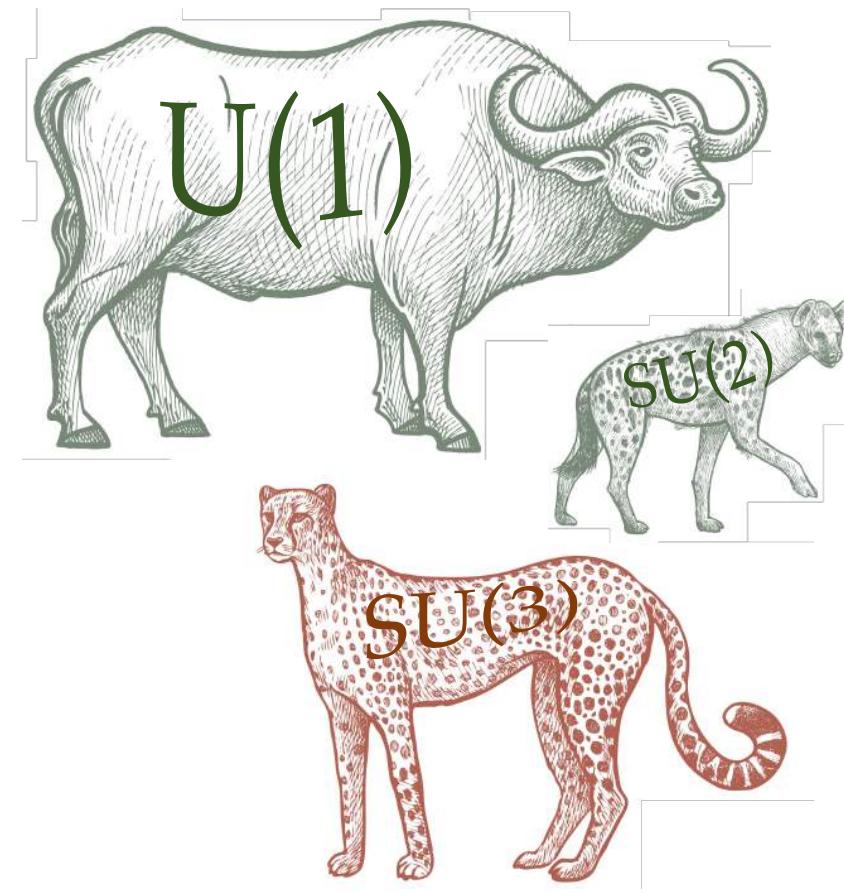


Message passing

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

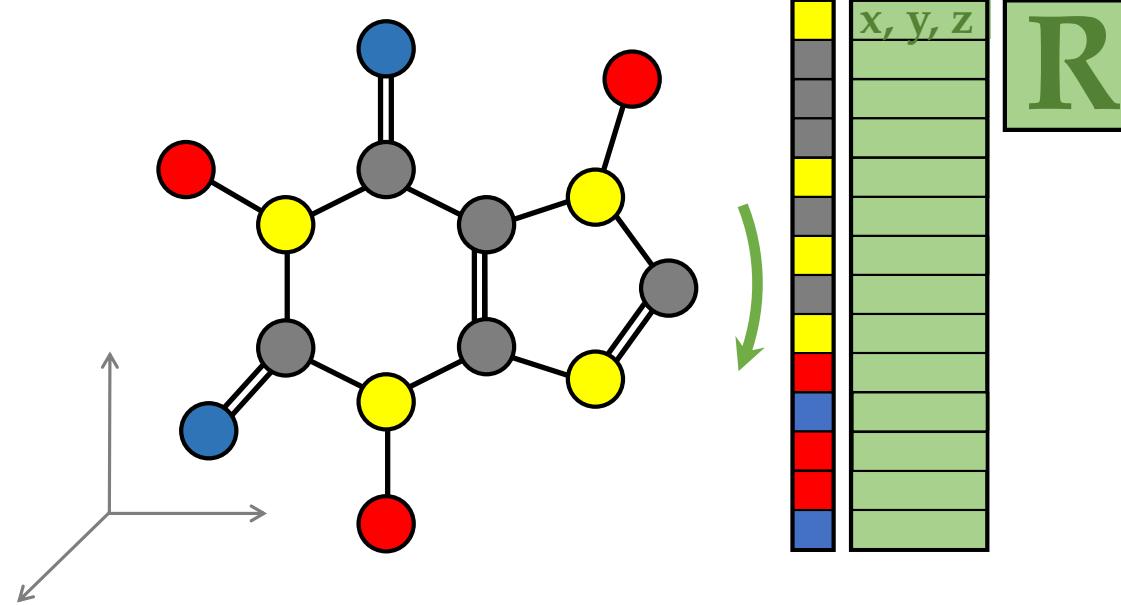


External symmetry



Internal symmetry

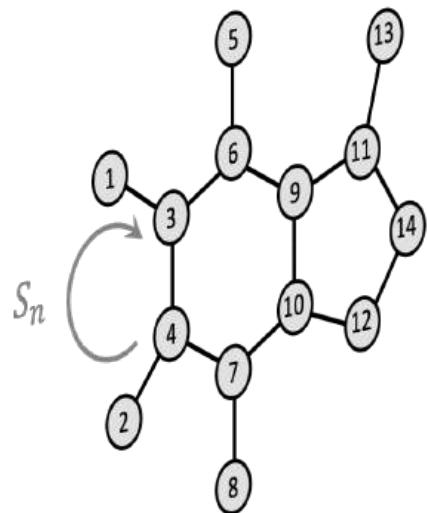
$\text{SO}(3)$ -invariance



"properties of a molecule do not change if we rotate it"

Geometric (“Equivariant”) Graph Neural Networks

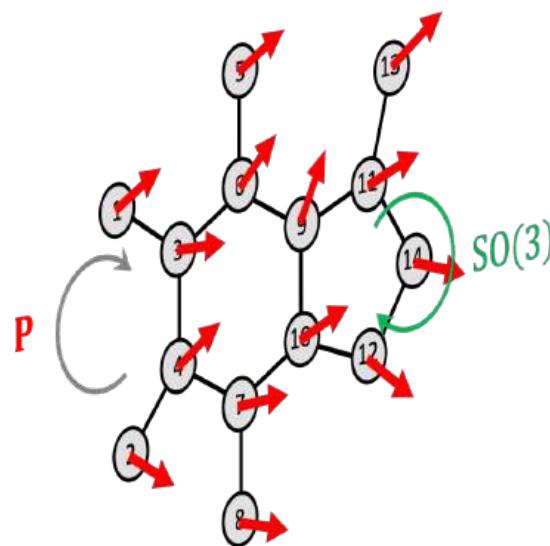
Geometric Graph G



Permutation group S_n

“domain symmetry”

Node features $\chi(G)$



Functions $\mathcal{F}(\chi(G))$



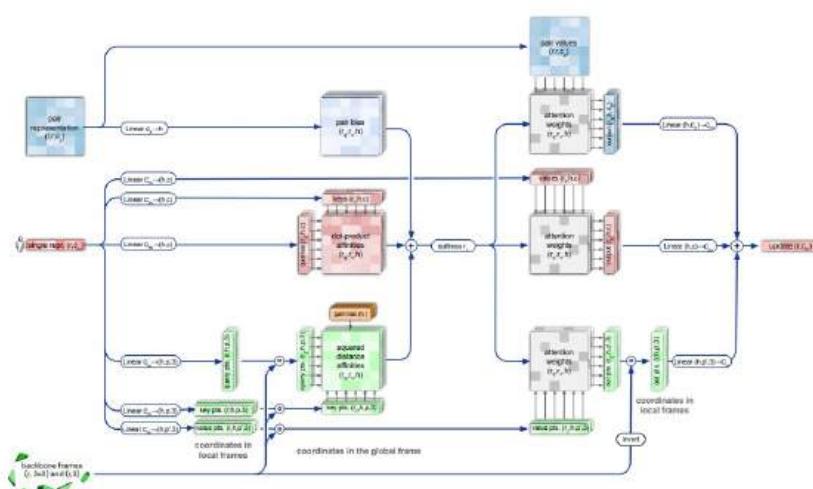
Permutation matrix P

Rotation R
“data symmetry”

Geometric message passing

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

Revolution in Structural Biology



Jumper et al. 2021

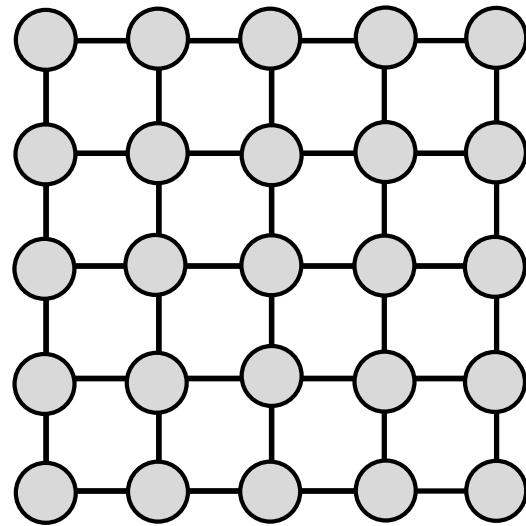
AlphaFold 2
“Invariant point
attention”



Baek et al. 2021

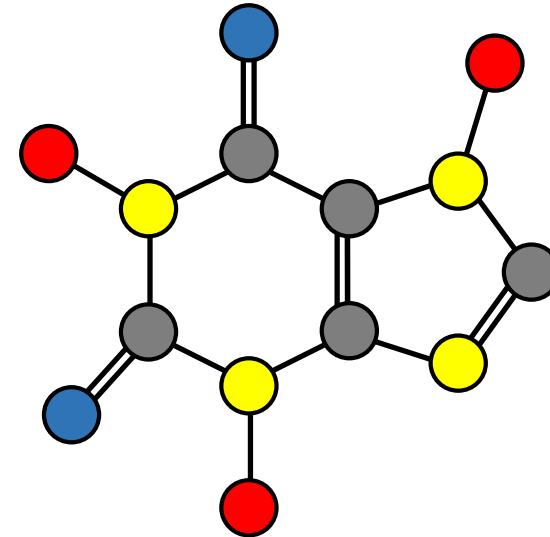
RosettaFold
SE(3)-equivariant
Transformer

Grids



Translation

Graphs



Permutation

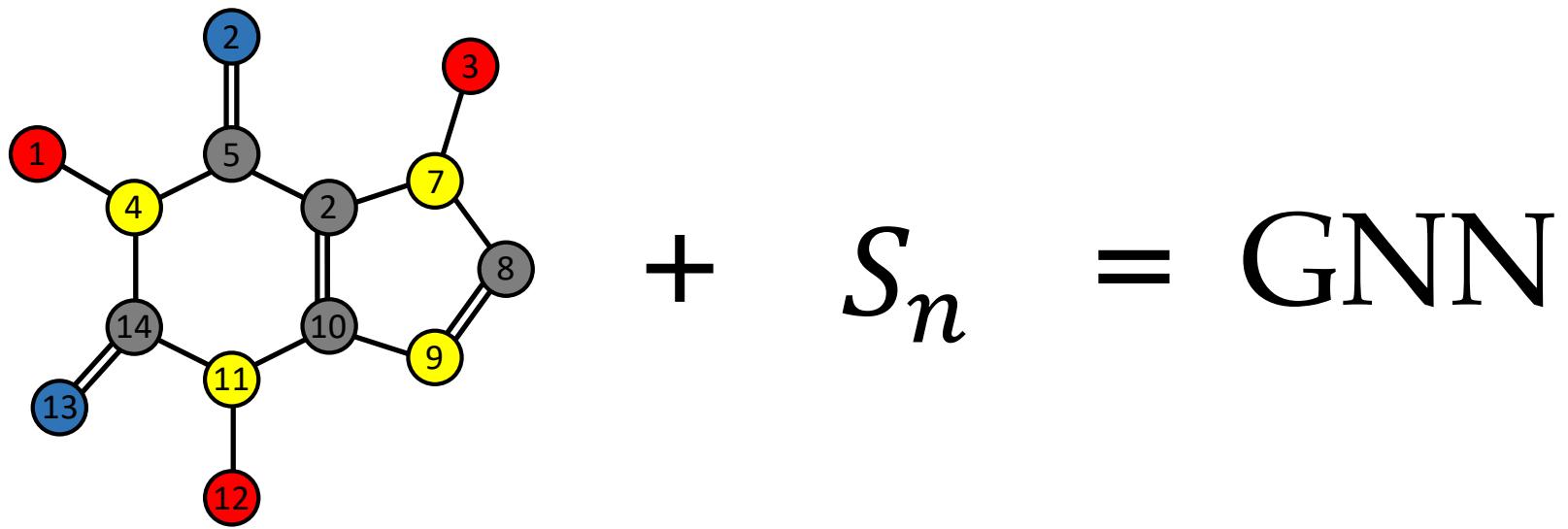
Meshes



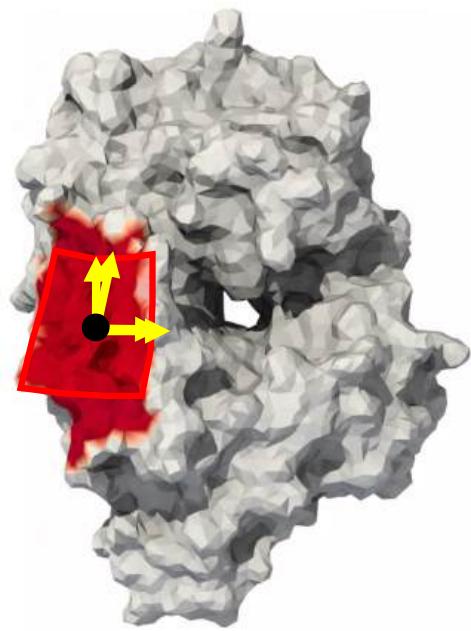
Local Rotation



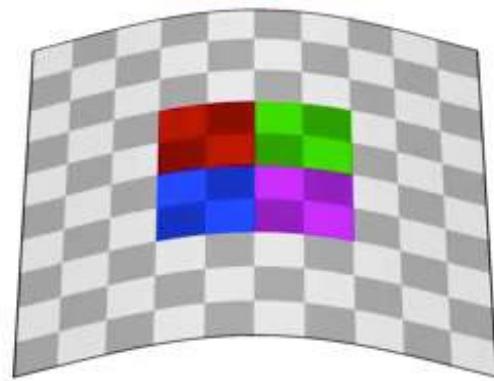
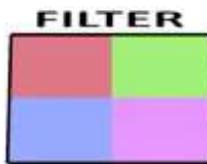
+ $T(2)$ = CNN



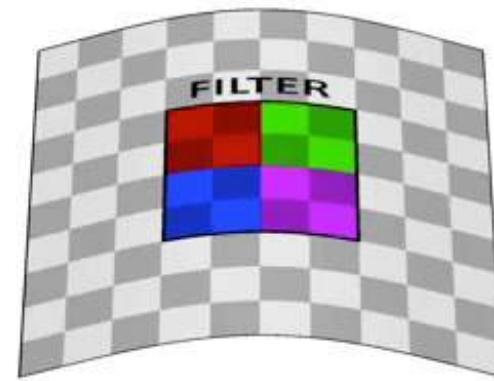
B et al 2017; 2021



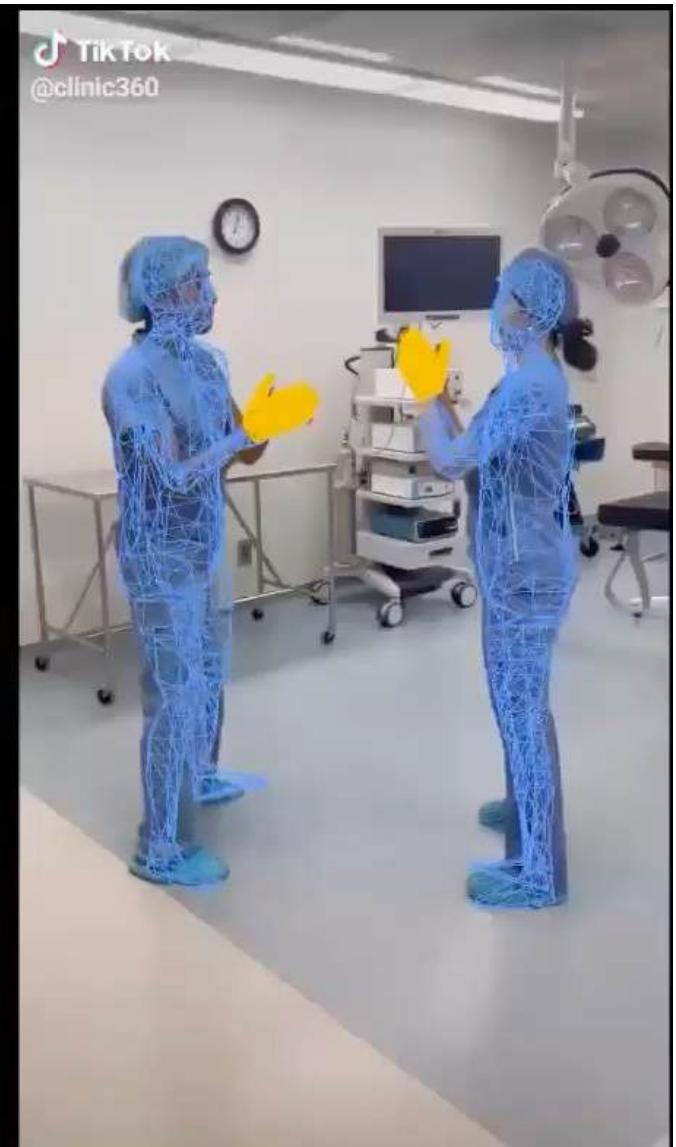
+ $\text{SO}(2)$ = MeshCNN



**Euclidean (extrinsic)
convolution**

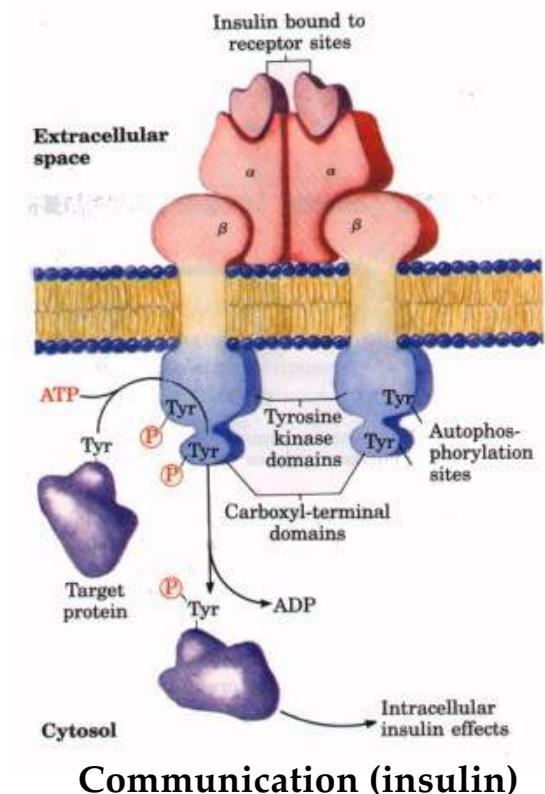
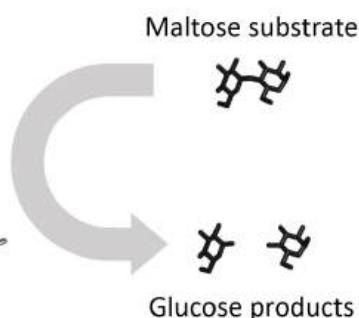
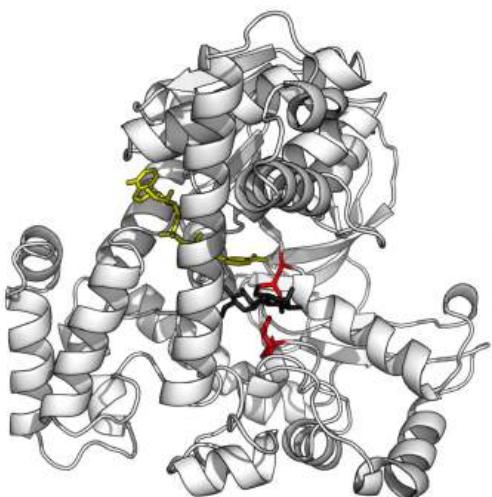
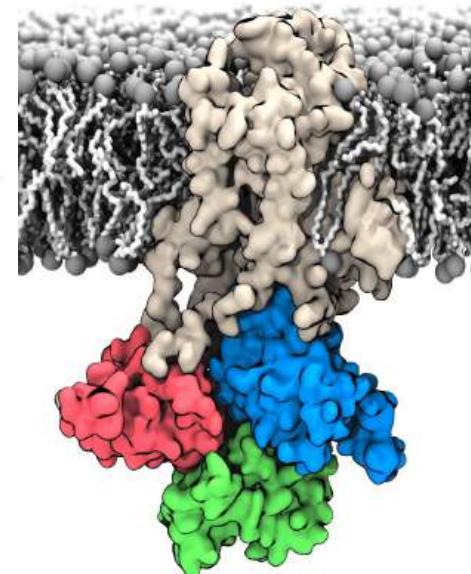
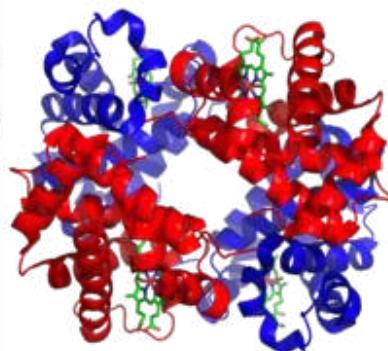
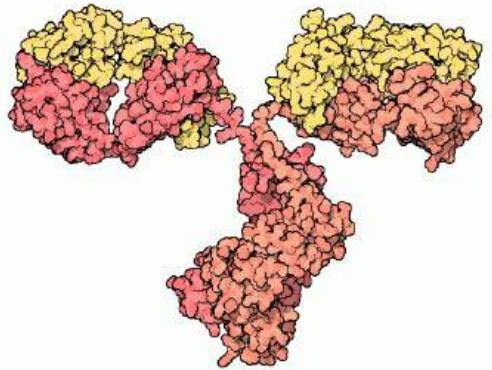


**Geometric (intrinsic)
convolution**

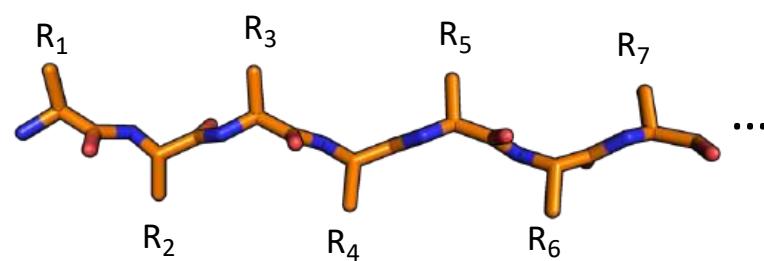


Snap Acquires Ariel AI To Enhance AR Features



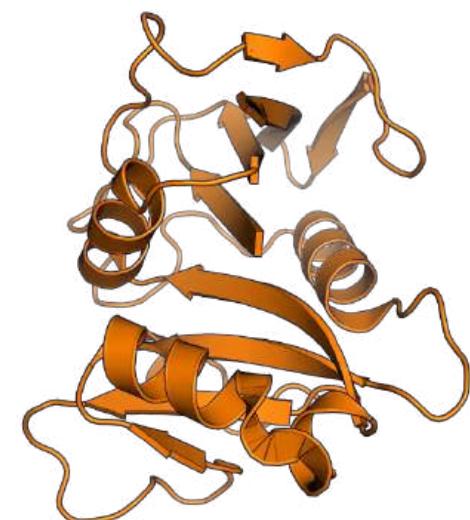
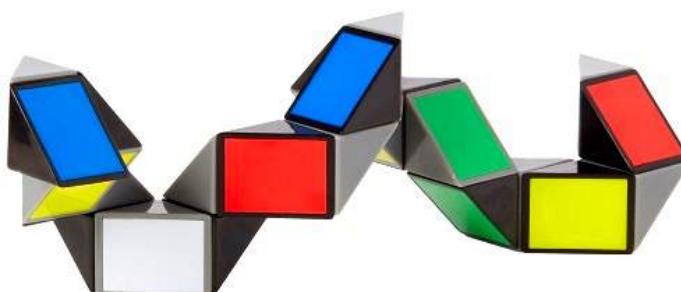


Protein Folding



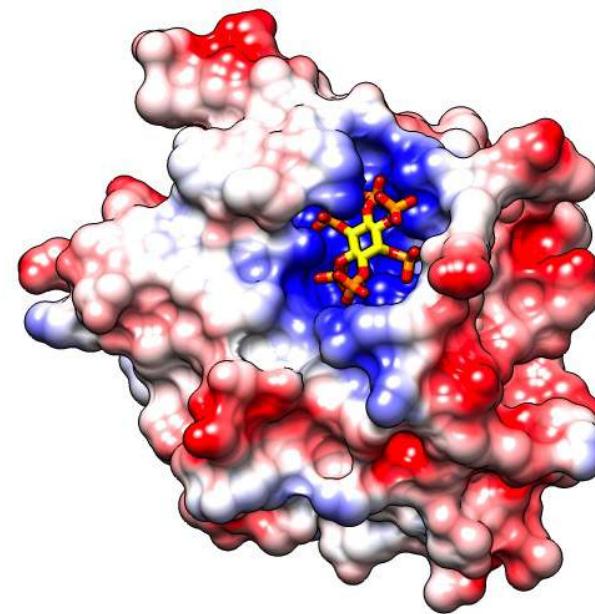
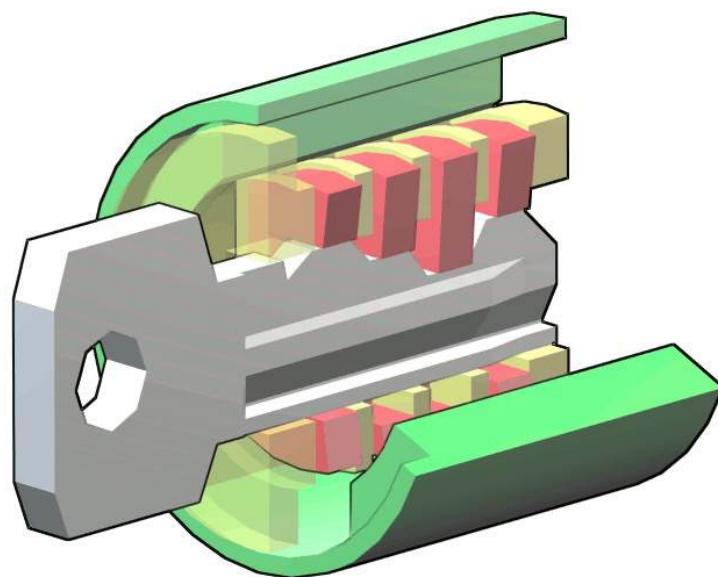
Sequence of
amino acids

Protein folding



3D structure

Lock-Key Metaphor



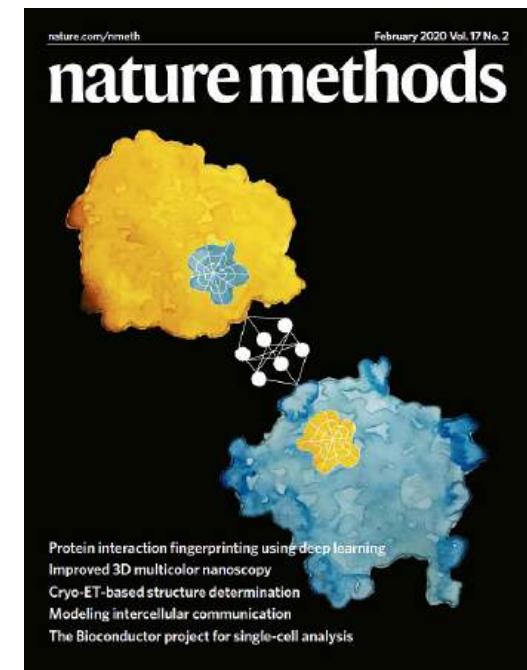
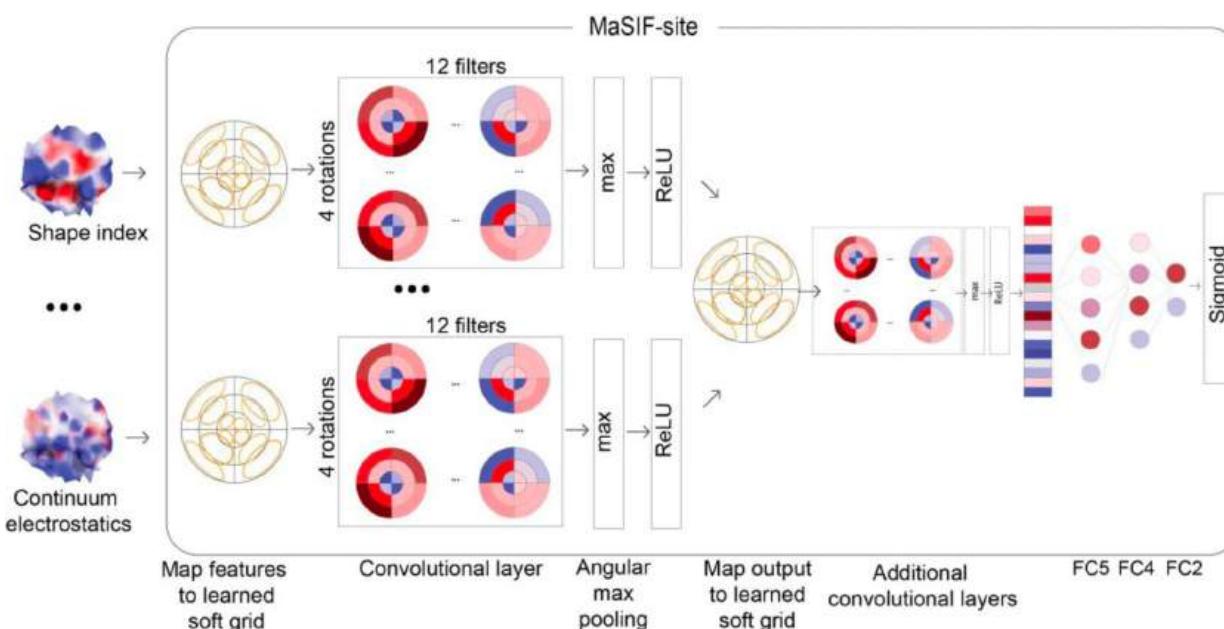
Emil Fischer "Schlüssel-Schloss-Prinzip" 1894

Surface-based Protein Representation

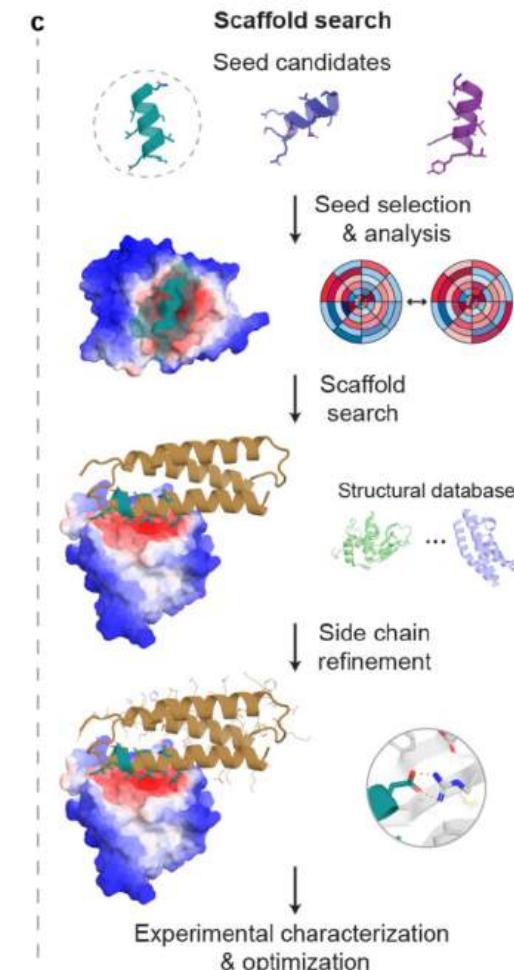
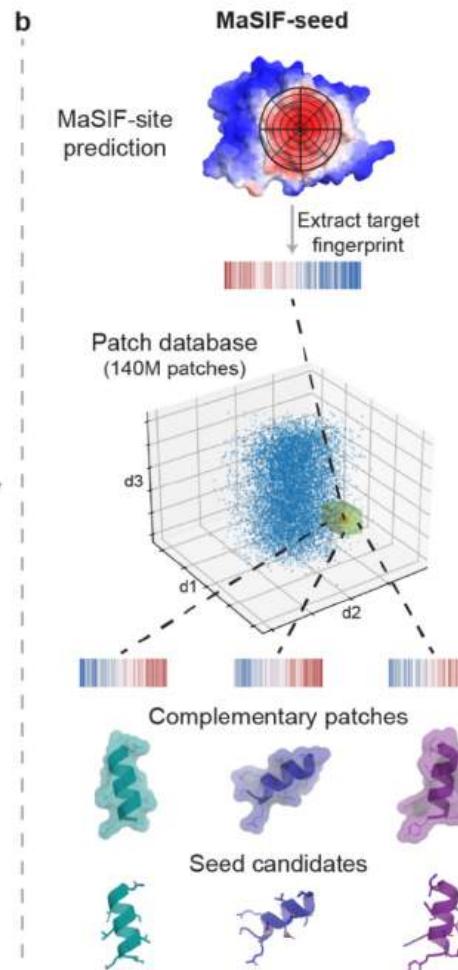
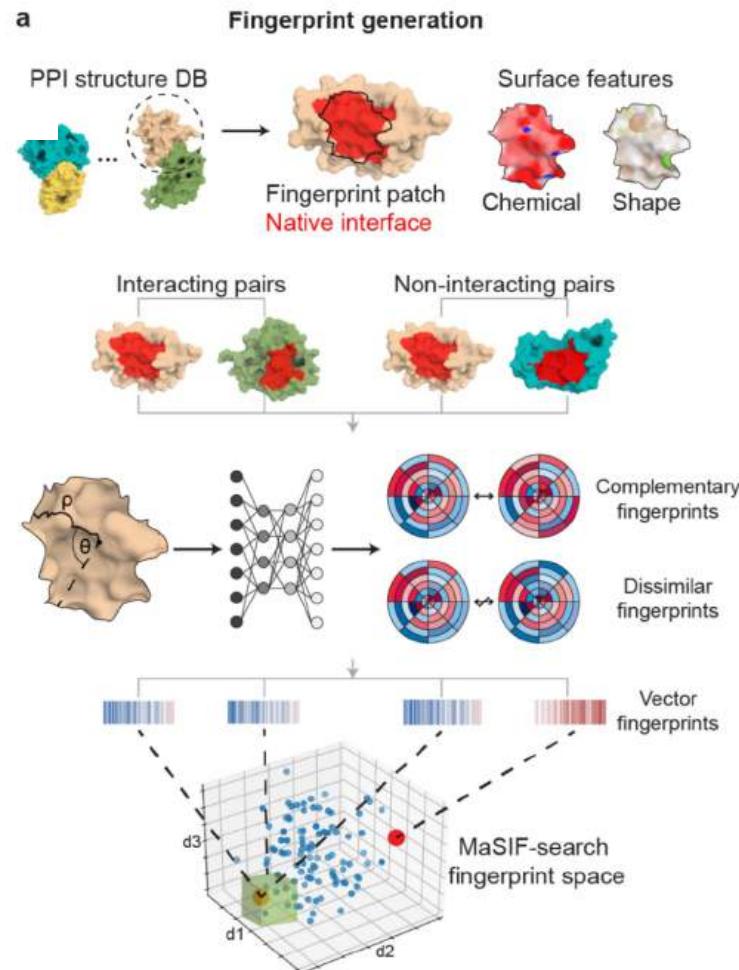
**Abstract out internal
structure** that is irrelevant
for interactions + allow
some **conformation
changes**



MaSIF: Geometric ML for Protein Function Prediction & Design

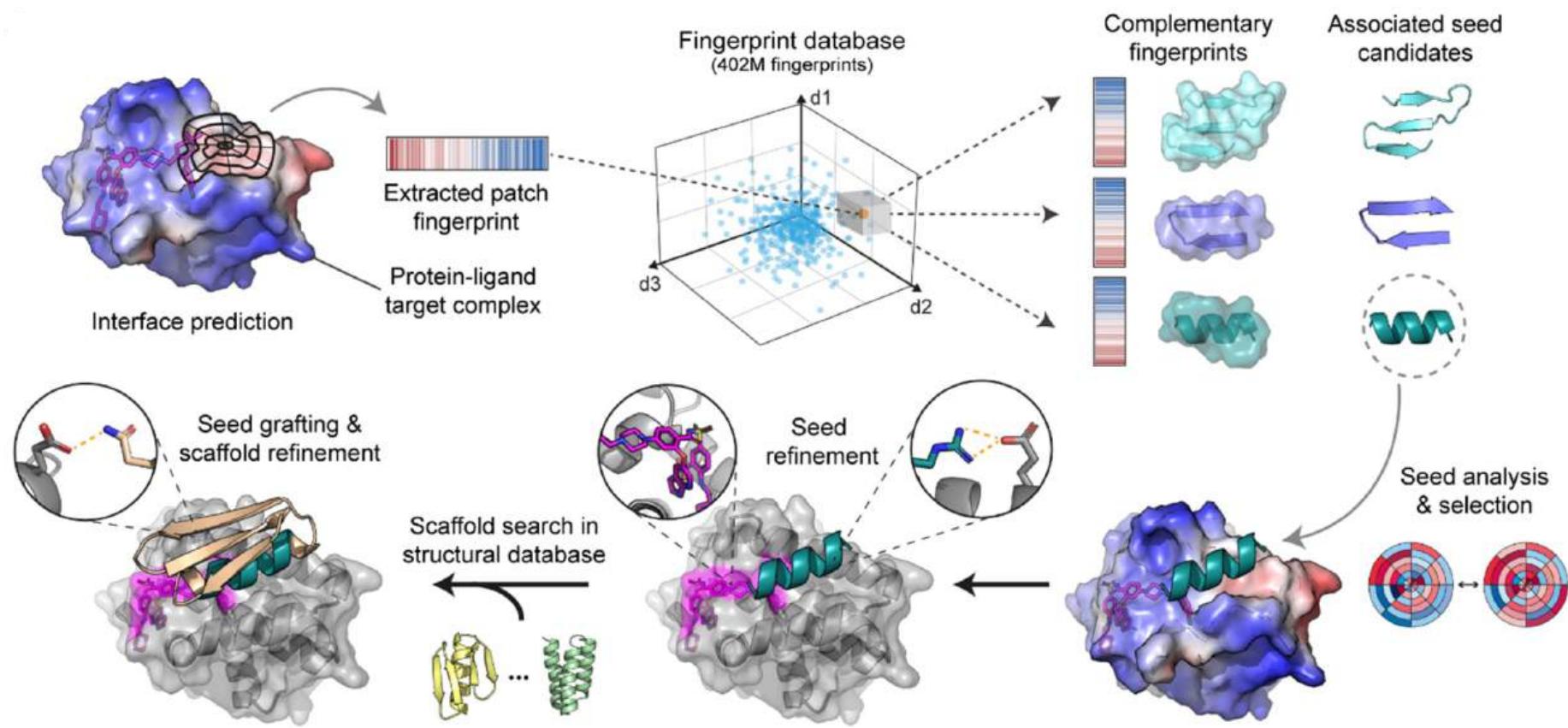


Gainza et B, Correia 2020



Gainza et al., Correia 2022

Switchable interactions



Marchand et B, Correia 2024

Article

De novo design of protein interactions learned surface fingerprints

<https://doi.org/10.1038/s41586-023-05993-x>

Received: 16 June 2022
Accepted: 21 March 2023
Published online: 26 April 2023
Open access

Pablo Gainza^{1,2*}, Sarah Wehrle^{1,2†}, Alessandra Van Hall-Belvaux^{1,2‡}, Andreas Schreck^{1,2}, Zander Hartwig^{1,2}, Stephen Buckley¹, Dongchun Ni¹, Frey Swirnsson¹, Casper Goverde¹, Priscilla Turrell¹, Charlène Racine¹, Alexandra Teslenko¹, Martin Pacea¹, Stéphanie Rosset¹, Sandrine George¹, Jane Marsden^{1,3}, Aparna Petruzzelli¹, Kefang Liu¹, Zepeng Xu¹, Yan Chai¹, Pengfei Guo¹, Elisa Oricio¹, Beate Fierz¹, Didier Tron¹, Henning Blath¹, Michael Broutin¹, and Elena E. Correa^{1,2§}

Physical interactions between proteins are essential for most biological governing life! However, the molecular determinants of such interactions challenging to understand, even as genomic, proteomic and structural. This knowledge gap has been a major obstacle for the comprehensive of cellular protein–protein interaction networks and for the *de novo* binders that are crucial for synthetic biology and translational applications we use a geometric deep-learning framework operating on protein sequences generates fingerprints to describe geometric and chemical features that drive protein–protein interactions^{3,4}. We hypothesized that these fingerprints the key aspects of molecular recognition that represent a new paradigm computational design of novel protein interactions. As a proof of principle, computationally designed several *de novo* protein binders to engage targets: SARS-CoV-2 spike, PD-1, PD-L1 and CTLA-4. Several designs were experimentally optimized, whereas others were generated purely in nanomolar affinity with structural and mutational characterizations accurate predictions. Overall, our surface-centric approach captures chemical determinants of molecular recognition, enabling an approach *de novo* design of protein interactions and, more broadly, of artificial function.

Designing novel protein–protein interactions (PPIs) remains a fundamental challenge in computational protein design, with broad basic and translational applications in biology. The challenge consists of generating amino acid sequences that engage a target site and form a quaternary complex with a given protein. This represents a stringent test of our understanding of the physicochemical determinants that drive biomolecular interactions¹. Robust computational methods to design de novo PPIs could be used to rapidly engineer protein-based therapeutics such as antibodies and protein inhibitors or vaccines among others, and are therefore of considerable interest for biomedical and translational applications^{2,3}.

Despite recent advances in rational PPI design^{2,6} and prediction¹⁵, designing novel protein binders against specific targets is very challenging, particularly when no structural elements from pre-existing ligands are available.



protein–ligand neosurfaces with a able deep learning tool

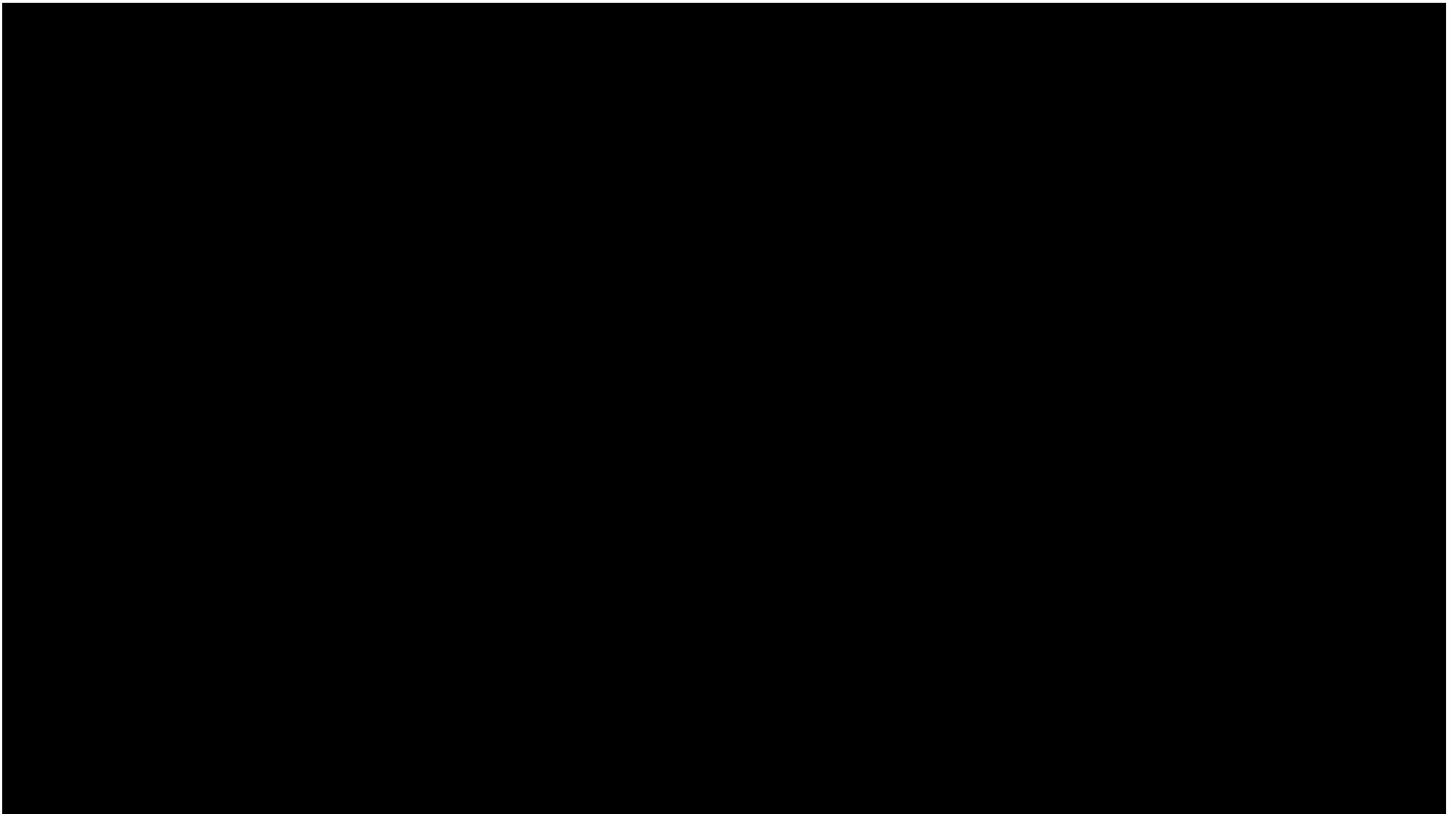
-024-08435-4 Anthony Marchand^a, Stephen Buckley^a, Anne Schnurrung^a, Martin Pacasa^a, Maddalena Elia^a, Pablo Gaitan^a, Evgenia Elizarova^a, Rebecca M. Neeser^b, Pao-Wan Lee^b, Luc Reymond^c, Yangyang Miao^c, Leo Scheller^c, Sandrine Georgeon^c, Joseph Schmidt^c, Philipp Schwaller^c, Sebastian J. Maerkli^c, Michael Brönstrom^d & Bruno E. Correia^d

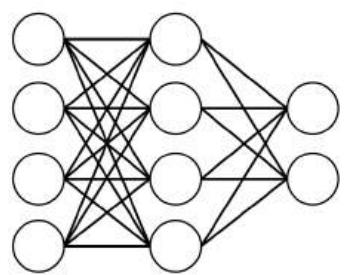
Molecular recognition events between proteins drive biological processes in living systems⁵. However, higher levels of mechanistic regulation have emerged, in which protein–protein interactions are conditioned to small molecules^{3,5}. Despite recent advances, computational tools for the design of new chemically induced protein interactions have remained a challenging task for the field⁶. Here we present a computational strategy for the design of proteins that target neosurfaces, that is, surfaces arising from protein–ligand complexes. To develop this strategy, we leveraged a geometric deep learning approach based on learned molecular surface representations^{7,8} and experimentally validated binders against three drug bound protein complexes: Bcl2–venetoclax, DB3–progesterone and PDF1–actinin. All binders demonstrated high affinity and accurate specificities, as assessed by mutational and structural characterization. Remarkably, surface fingerprints previously trained only on proteins could be applied to neosurfaces induced by interactions with small molecules, providing a powerful demonstration of generalizability that is uncommon in other deep learning approaches. We anticipate that such designed chemically induced protein interactions will have the potential to expand the sensing repertoire and the assembly of new synthetic pathways in engineered cells for innovative drug-controlled cell-based therapies⁹.

PPPs) have essential roles in healthy cell division in numerous diseases¹³. For this, PPis have been developed over the computational tools for the design of new entity been proposed¹⁴. The governing quantity of proteins to form interactions is play of several contributions, such as semipermeability, dynamics and solvation, challenging to predict and design new of evolutionary constraints. Native PPis story-lines such as allosteric¹⁵, posttranslational¹⁶ and C1 ligand binding¹⁷. Compound-bound resources, are among the most fascial¹⁸ recognition instances, as relatively binding site can have a large impact on its interaction has been fuelled by the whites, specifically molecular glue that¹⁹ protein interactions for drug development and is represent a promising route for the therapeutics.

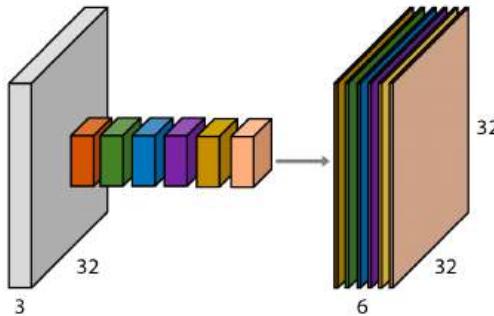
engineering, Institute of Biointerfaces, Ecole polytechnique fédérale de Lausanne, Lausanne, Switzerland; "Laboratory of Chemical Analysis and Engineering, Ecole polytechnique fédérale de Lausanne, Lausanne, Switzerland; "Laboratory of Biological Network Characterization, Institut Voltaire de Lausanne, Lausanne, Switzerland; "Stemcellomics Screening Core Facility, School of Life Sciences, Ecole polytechnique fédérale de Lausanne, Lausanne, Switzerland; "Department of Computer Science, Oxford University, Oxford, UK; "Kittay Research Institute for Biomedicine/Artificial Intelligence, Austin Academy Inc., Monte Rosa Therapeutics, Boston, MA, USA. "These authors contributed equally. A. Kirchhoff, Stephen Buder, Xiang Schneiders

surface feature
using a geometric
approach we developed two ap-
proaches of protein surface
at another protein; at
other partners on the
use search, we extract
patches with complex
fingerprints, whereas in
maturity, Surface finger-
prints free manner
values with fingerprint
are stored in three dimensial
structure to refine the
in its initial concept
as part of the pro-
tein with small molecu-
lular size. MSW-means, w/
molecular surface repr-
esentations and partners on
the search, MSW is
a geometric surface
representing the underlying struc-
ture. It should in principle
be able to represent protein surfaces
of protein-ligand com-
plexes. Shape index
is a geometric feature
which can be comput-
ed from bond/acc-

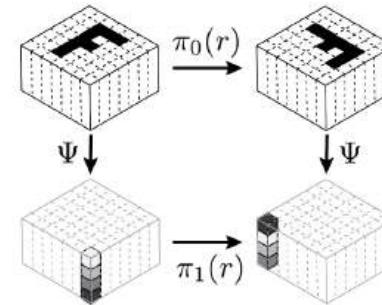




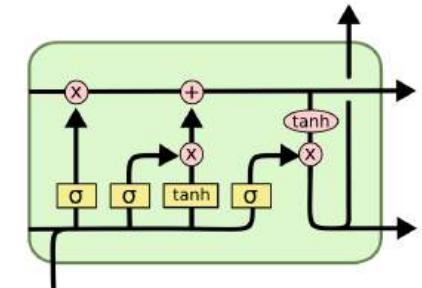
Perceptrons
Function regularity



CNNs
Translation



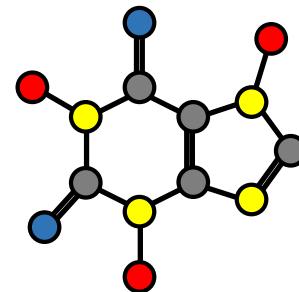
Group-CNNs
Translation+Rotation,
Global groups



LSTMs
Time warping



DeepSets / Transformers
Permutation



GNNs
Permutation



Intrinsic CNNs
Isometry / Gauge choice



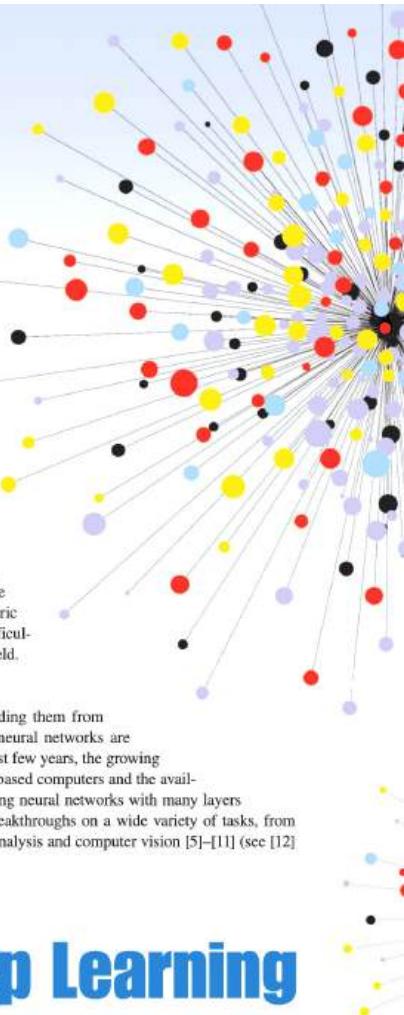
Michael M. Bronstein, Joan Bruna, Yann LeCun,
Arthur Szlam, and Pierre Vandergheynst

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.

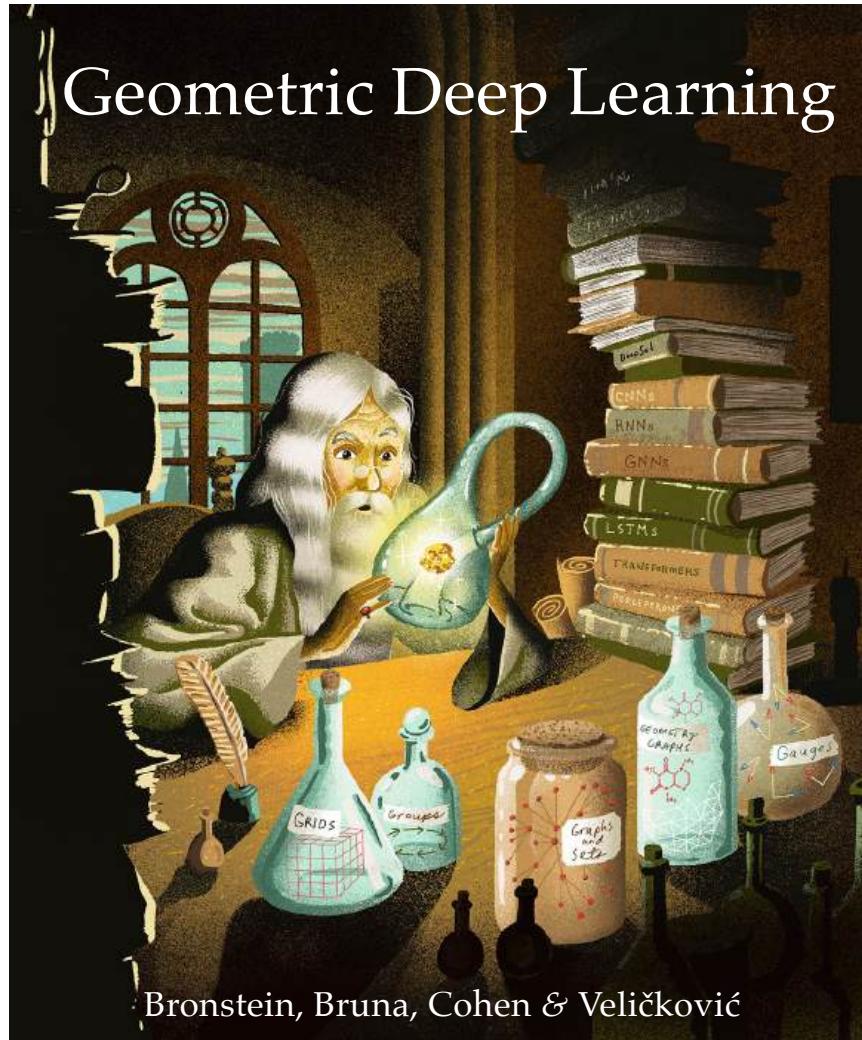
Overview of deep learning

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]

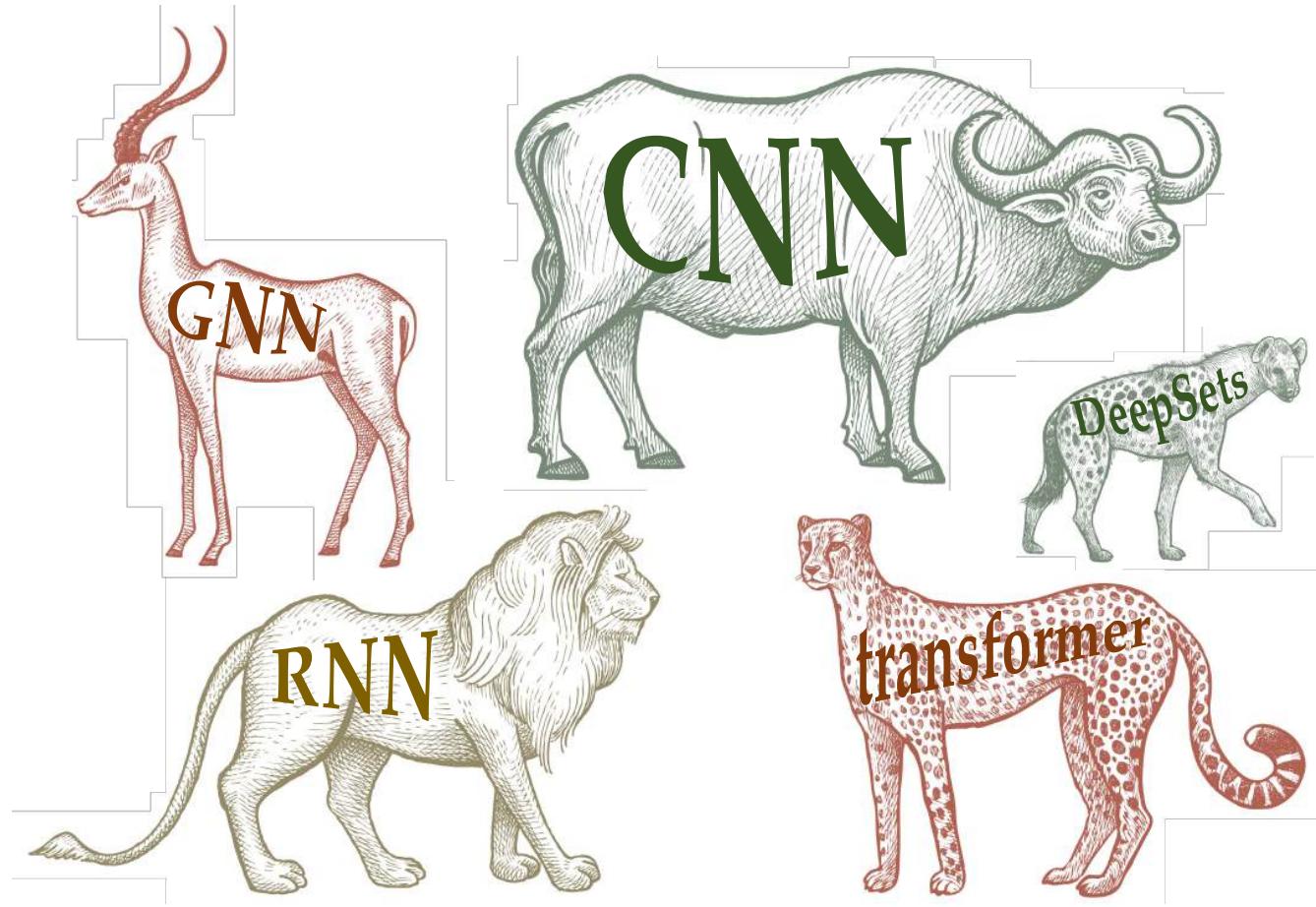


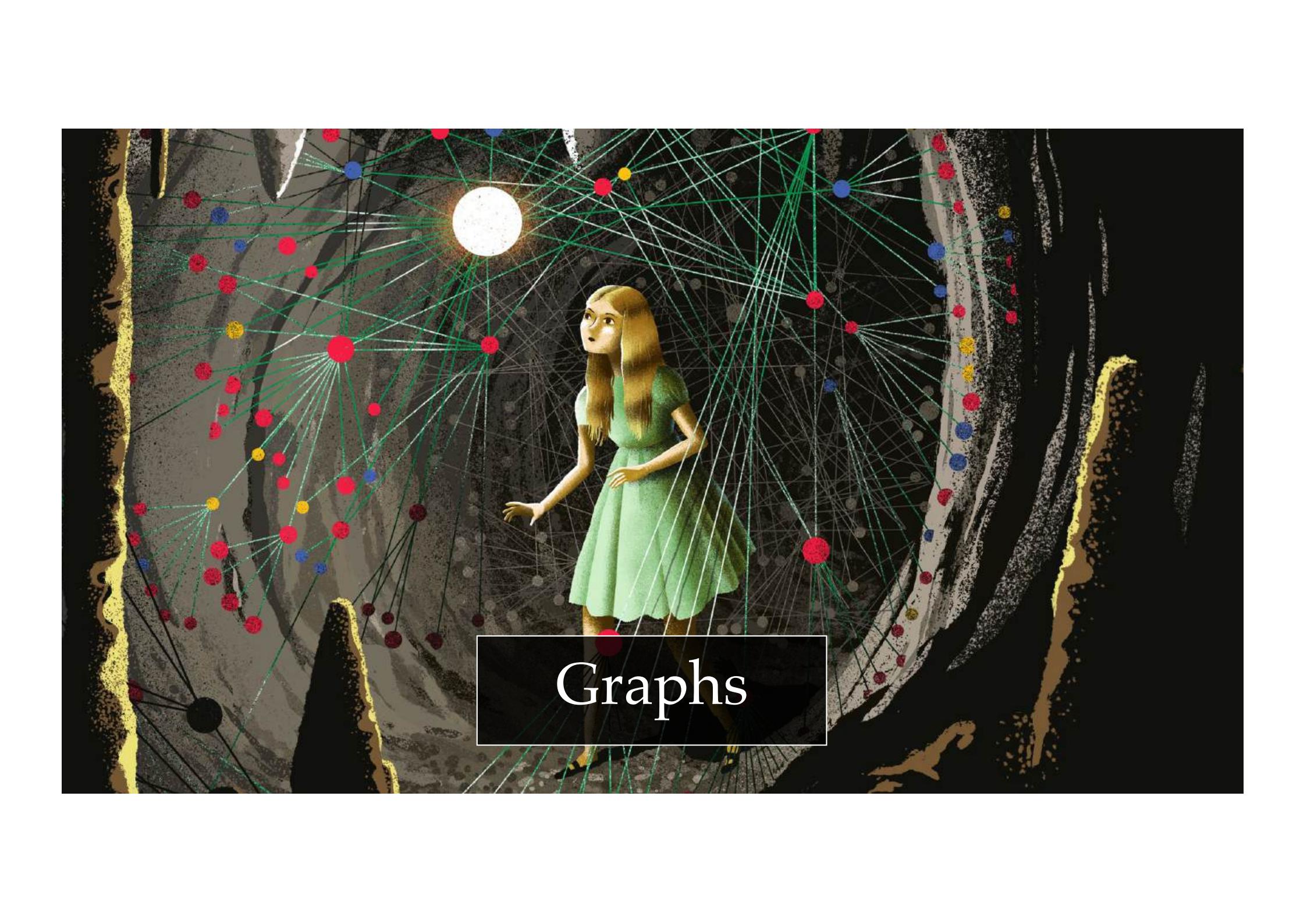
Geometric Deep Learning

Going beyond Euclidean data



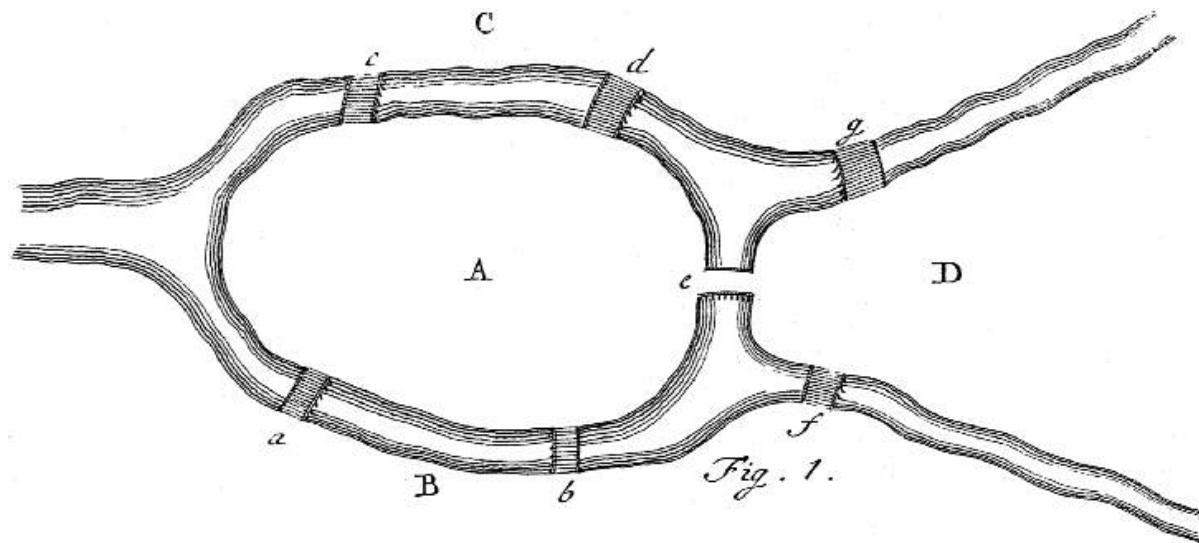
“The Erlangen Programme of ML”
Geometric Deep Learning





Graphs

Origins of Graph Theory



L. Euler

The solution of the classical problem of the “Bridges of Königsberg” by Euler in 1736 first showed the power of graphs to abstract out the geometry (“*geometria situs*”)

1741

Euler 1741

Origins of Topology

Poincaré's "*analysis situs.*"
His famous Conjecture
appeared in a supplement
published in 1904.



H. Poincaré



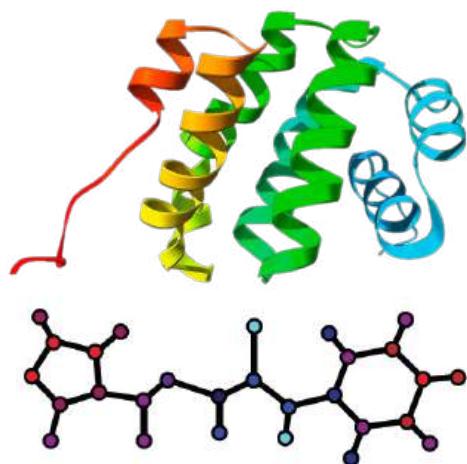
L. Euler

1895

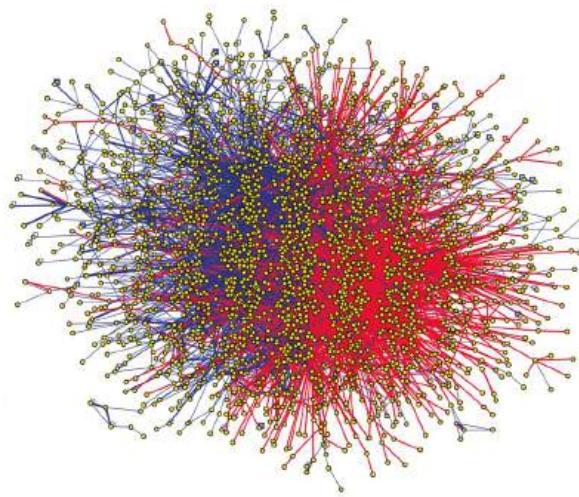
1741

Euler 1741; Poincaré 1895

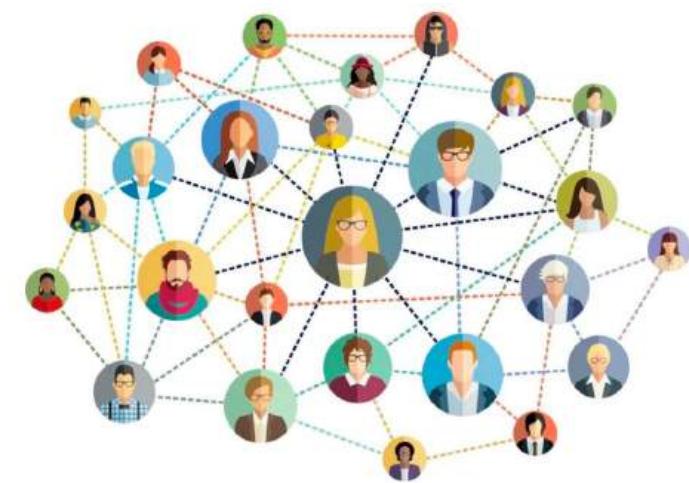
Graphs = Systems of Relations and Interactions



Molecules

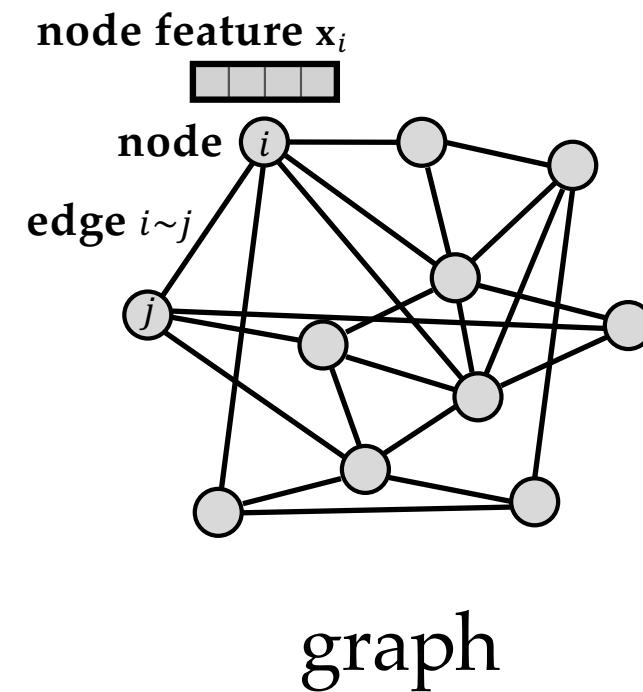


Interactomes

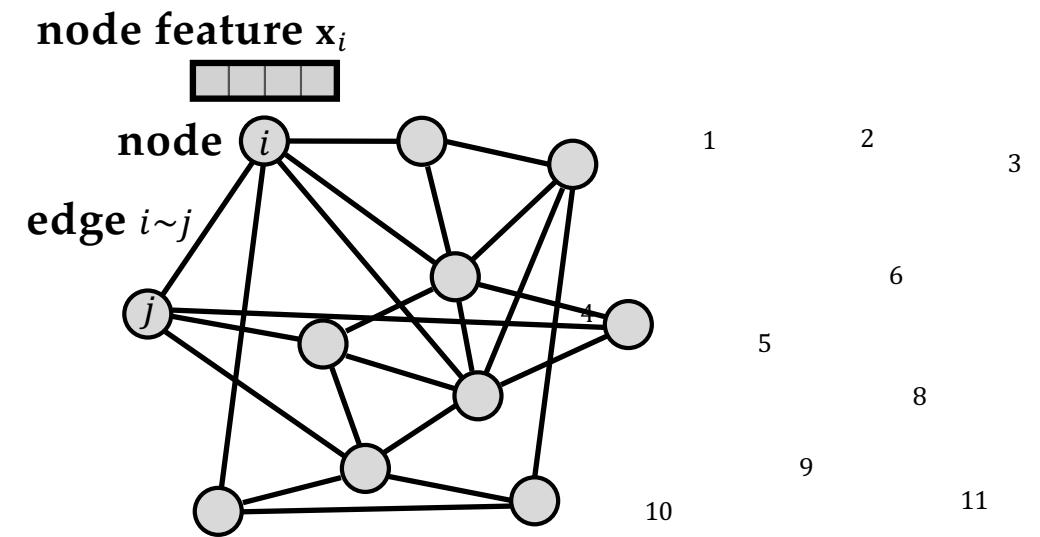


Social networks

Graphs: The Basics



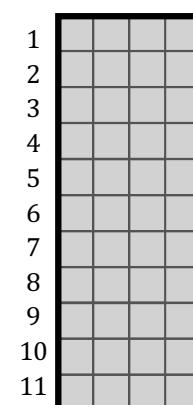
Key Structural Properties of Graphs



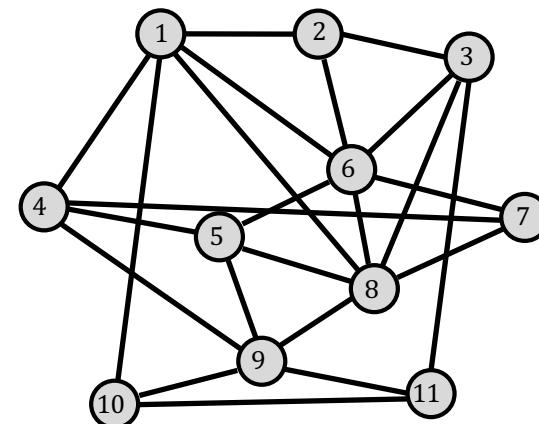
arbitrary graph ordering of nodes

Key Structural Properties of Graphs

Feature
matrix $n \times d$

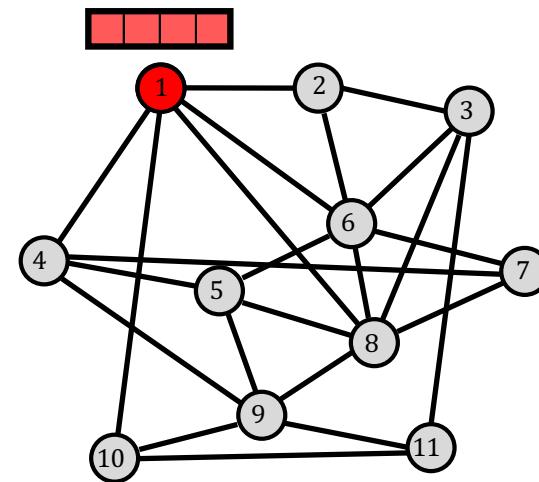
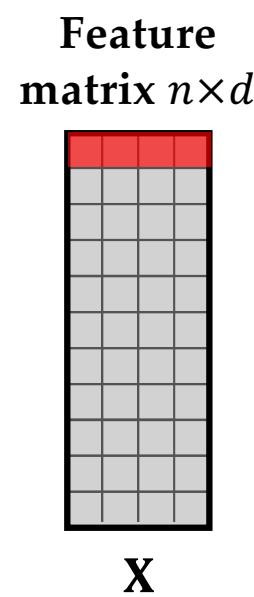
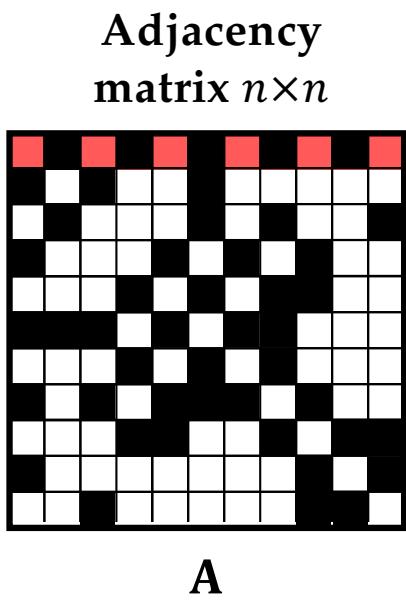


X

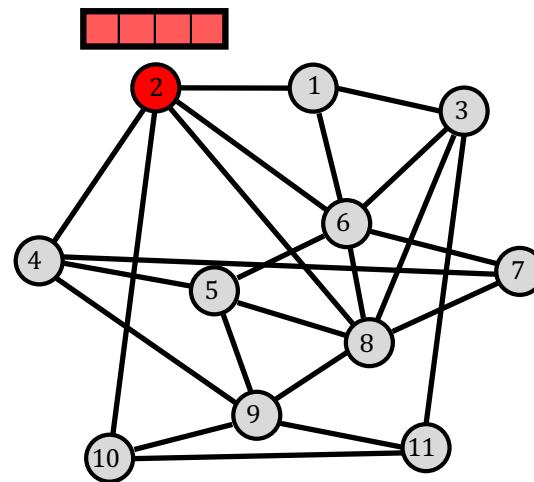
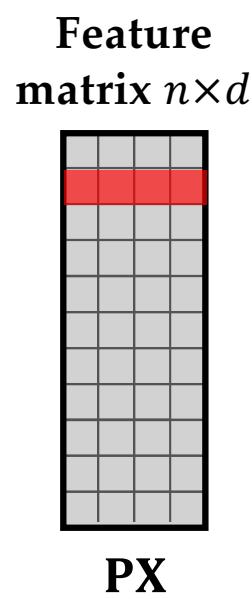
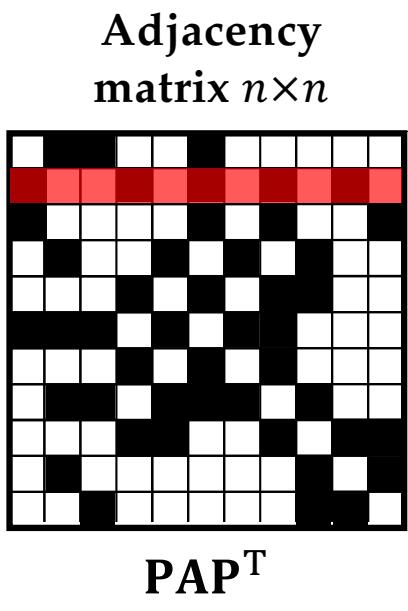


arbitrary ordering of nodes

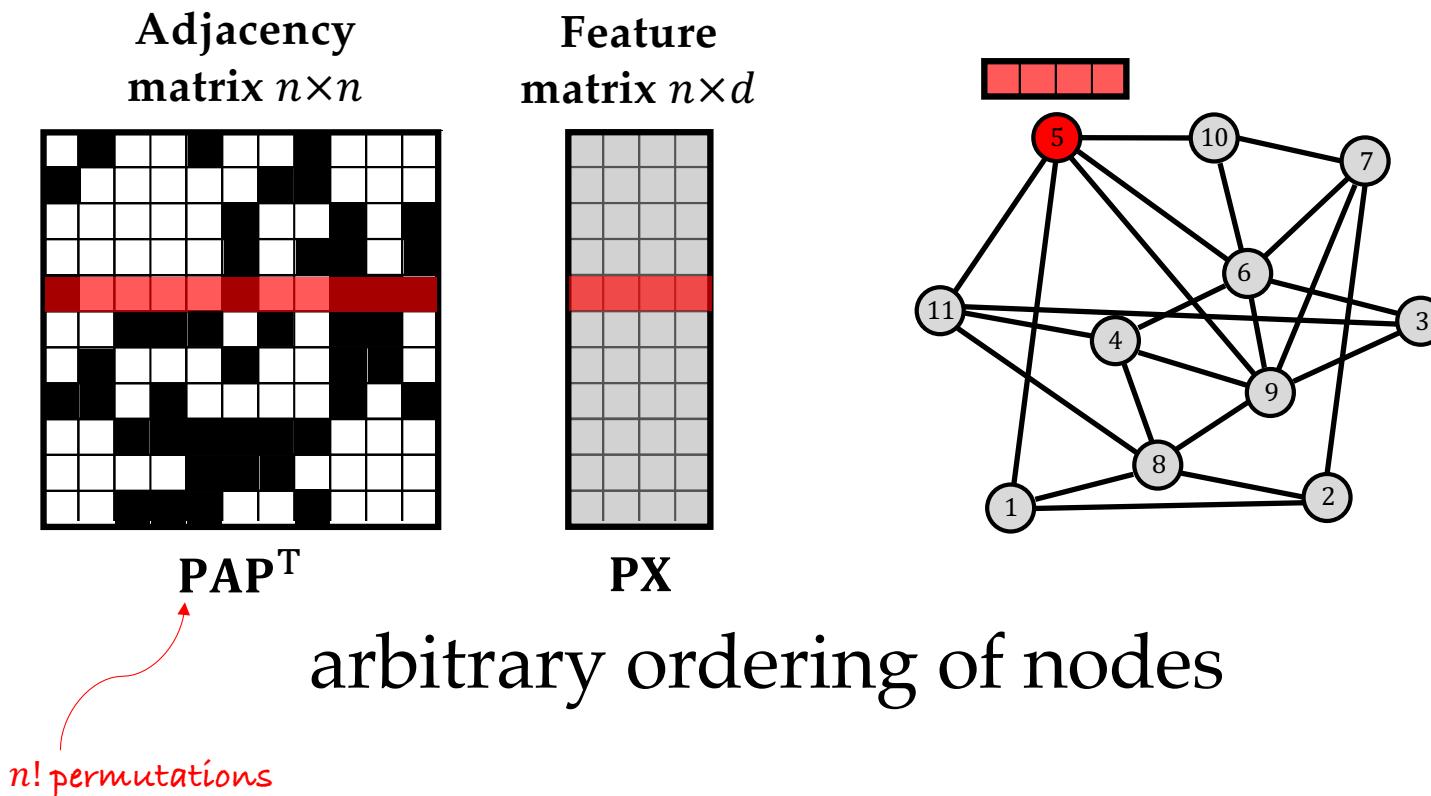
Key Structural Properties of Graphs



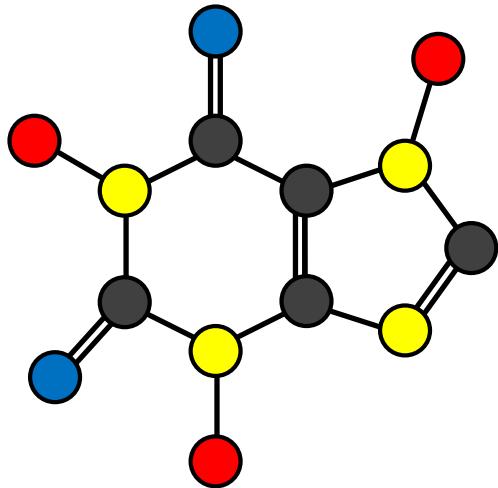
Key Structural Properties of Graphs



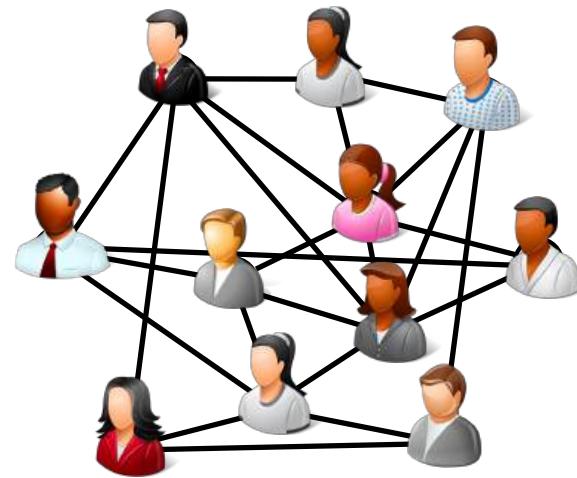
Key Structural Properties of Graphs



Invariant vs Equivariant tasks



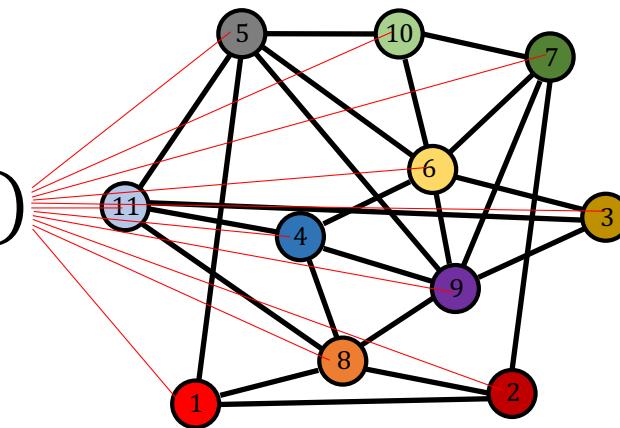
water solubility?



who is a spammer?

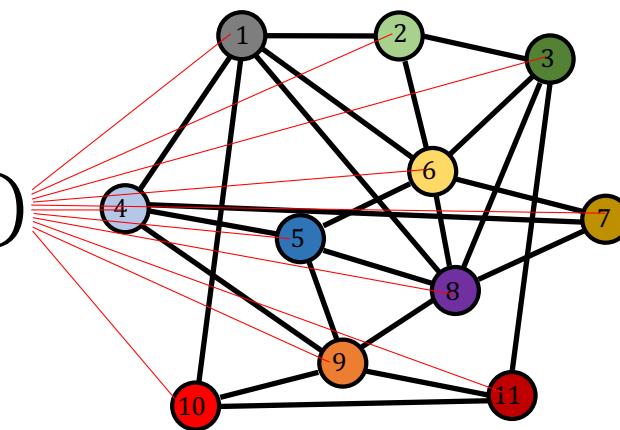
Invariant Graph Functions

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



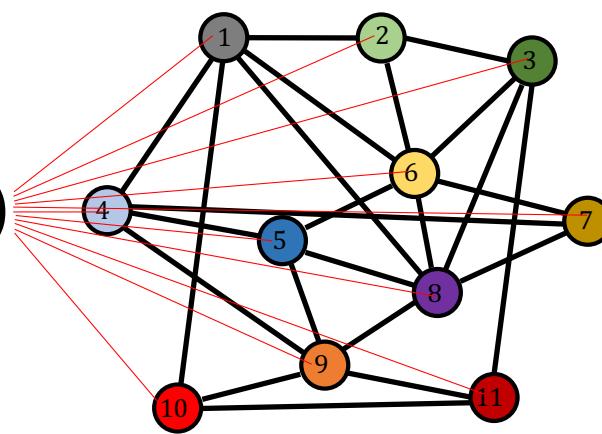
Invariant Graph Functions

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

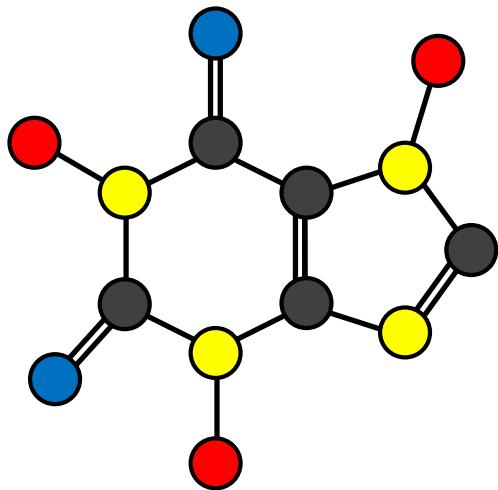


Invariant Graph Functions

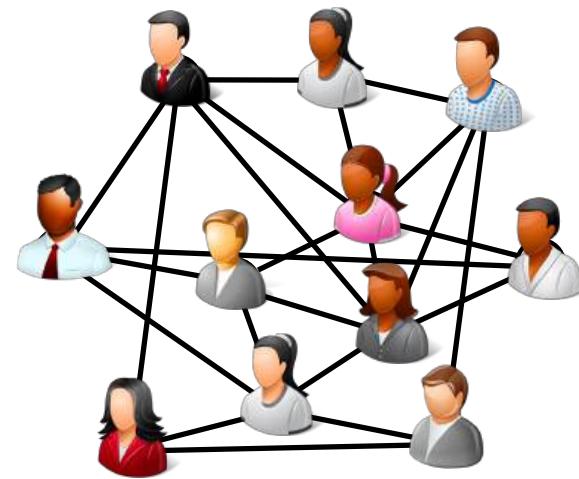
permutation-invariant
 $f(\mathbf{P}\mathbf{X}, \mathbf{P}\mathbf{A}\mathbf{P}^T) = f(\mathbf{X}, \mathbf{A})$



Invariant vs Equivariant tasks



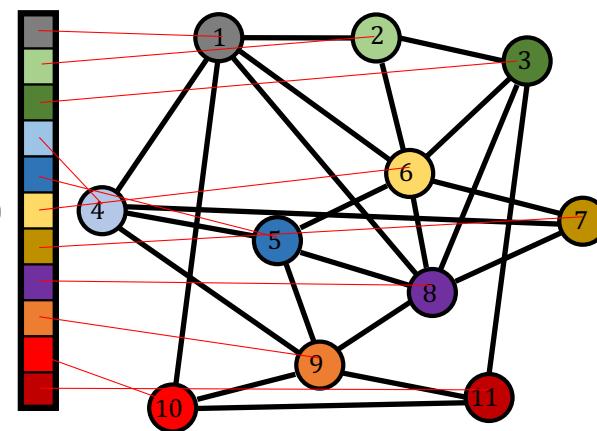
water solubility?



who is a spammer?

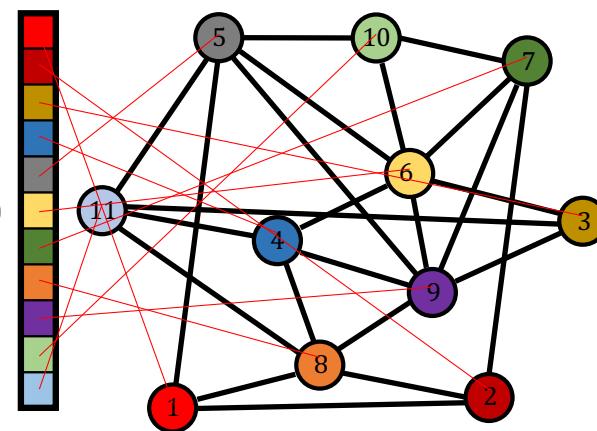
Equivariant Graph Functions

node function $F(X, A)$



Equivariant Graph Functions

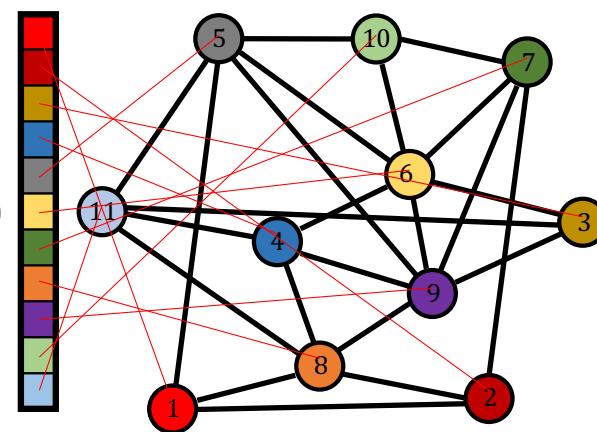
node function $F(X, A)$



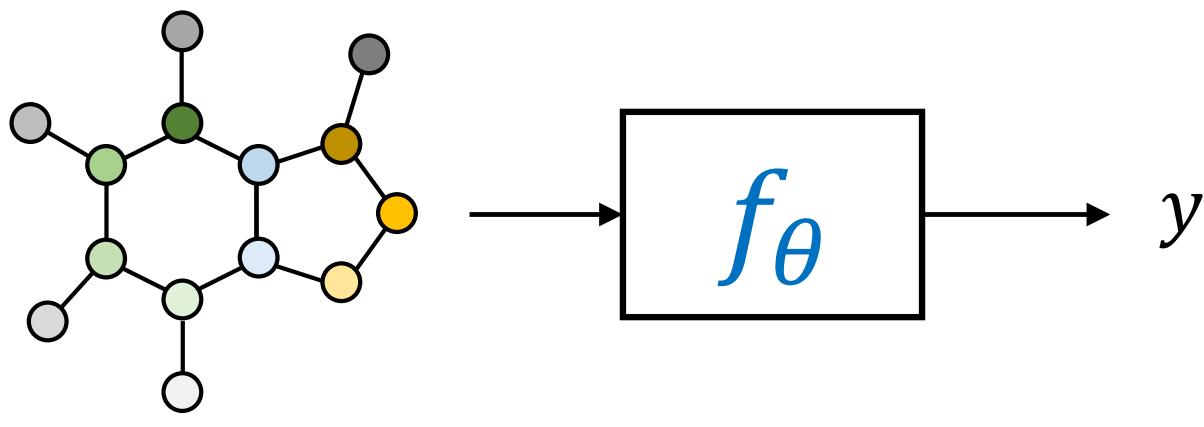
Equivariant Graph Functions

permutation-equivariant

$$F(PX, PAP^\top) = PF(X, A)$$



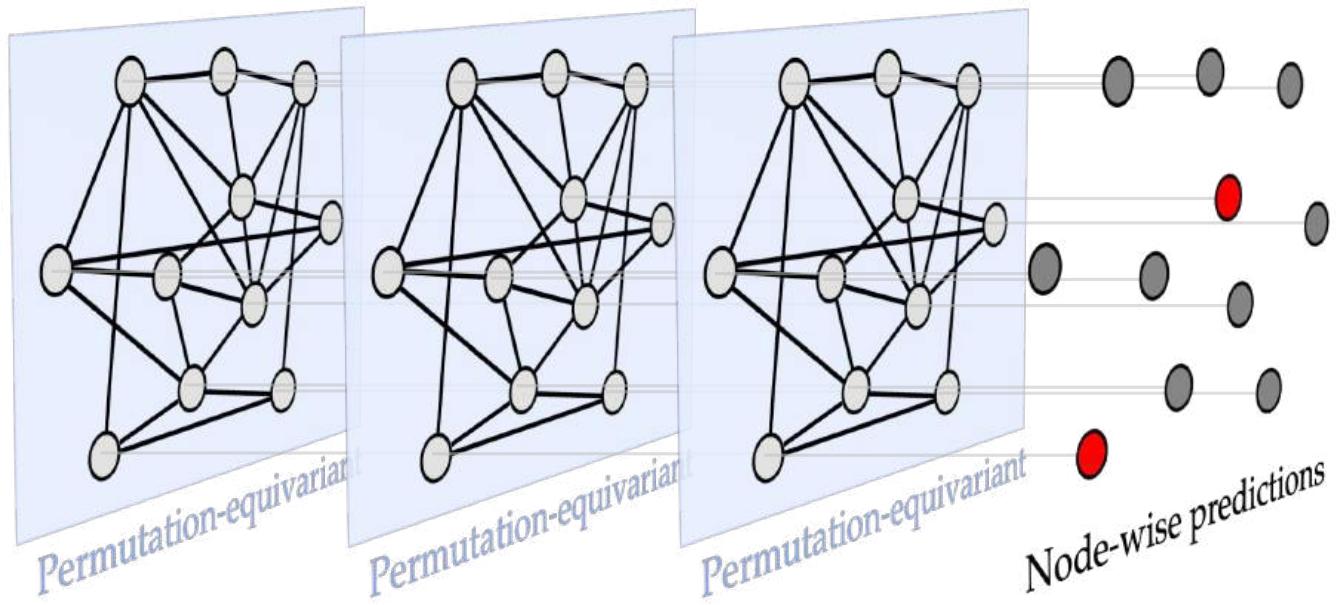
GNNs = Parametric graph functions



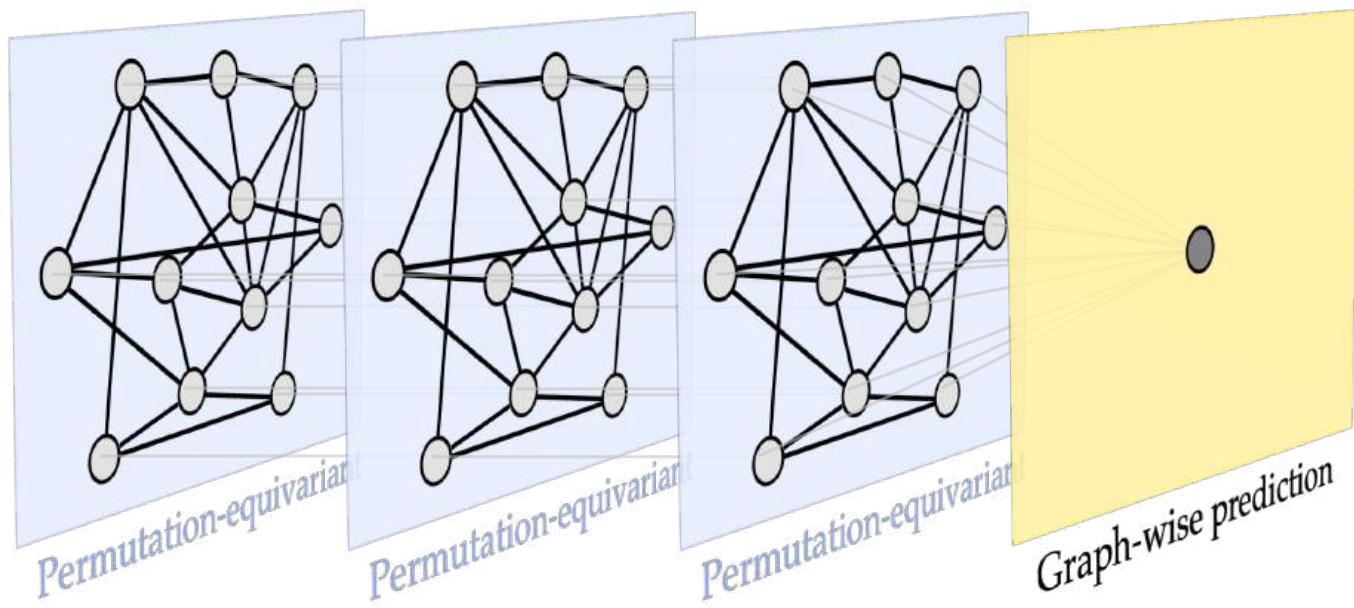
Graph+Features

Note: we use the term “GNN” to refer to general parametric graph functions. The particular class of GNNs we consider are Message Passing Neural Networks (MPNNs). Petar Veličković argues that “it’s all message passing”

Graph Neural Networks: Node tasks



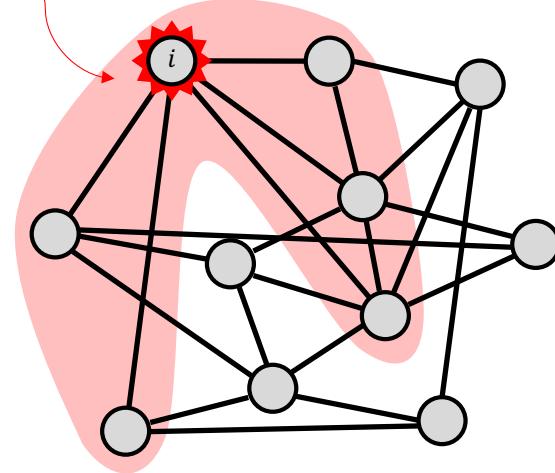
Graph Neural Networks: Graph tasks



Neighbour Aggregation

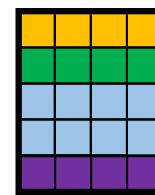
neighbourhood

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



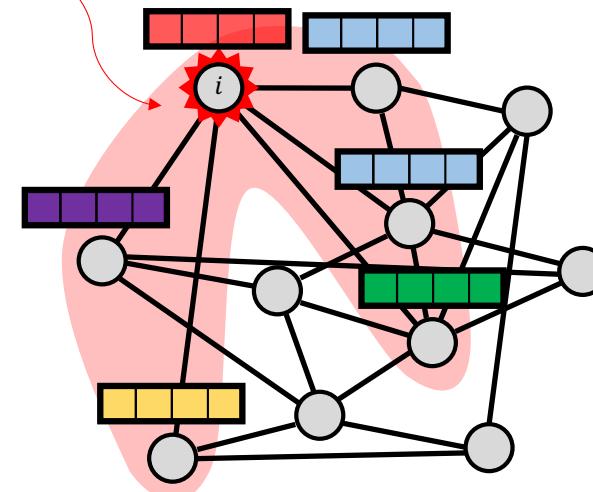
Neighbour Aggregation

multiset of
neighbour features



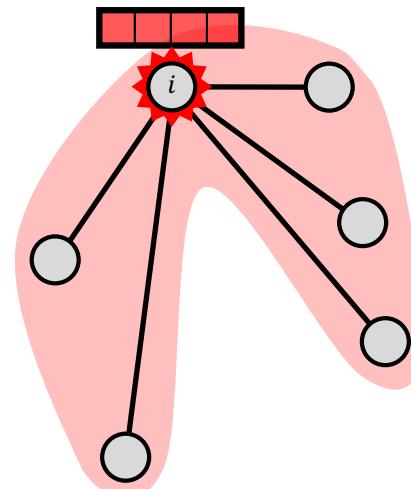
$$\mathbf{X}_{\mathcal{N}_i} = \{\mathbf{x}_{j \in \mathcal{N}_i}\}$$

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



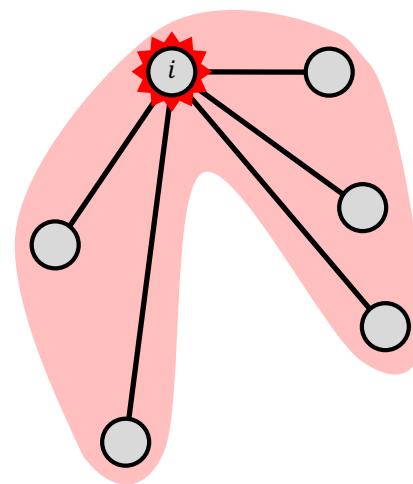
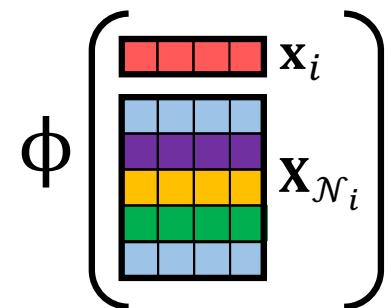
Neighbour Aggregation

multiset of
local function
neighbour features

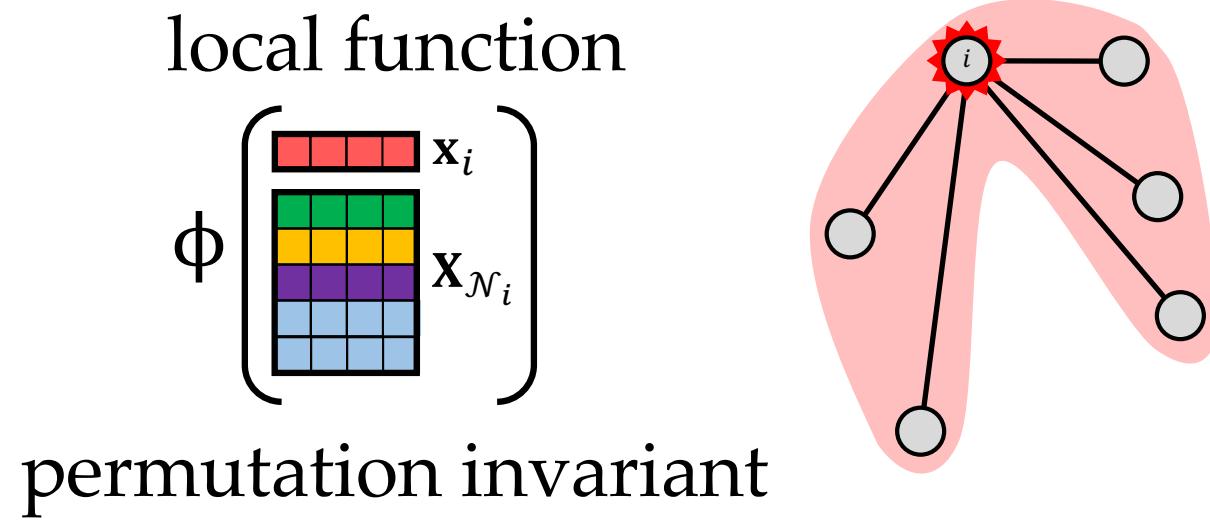
$$\phi \left(\begin{matrix} \text{yellow} \\ \text{green} \\ \text{blue} \\ \text{purple} \end{matrix} \right) \quad \mathbf{x}_i \quad \mathbf{x}_{\mathcal{N}_i}$$
$$\mathbf{X}_{\mathcal{N}_i} = \{\mathbf{x}_{j \in \mathcal{N}_i}\}$$


Neighbour Aggregation

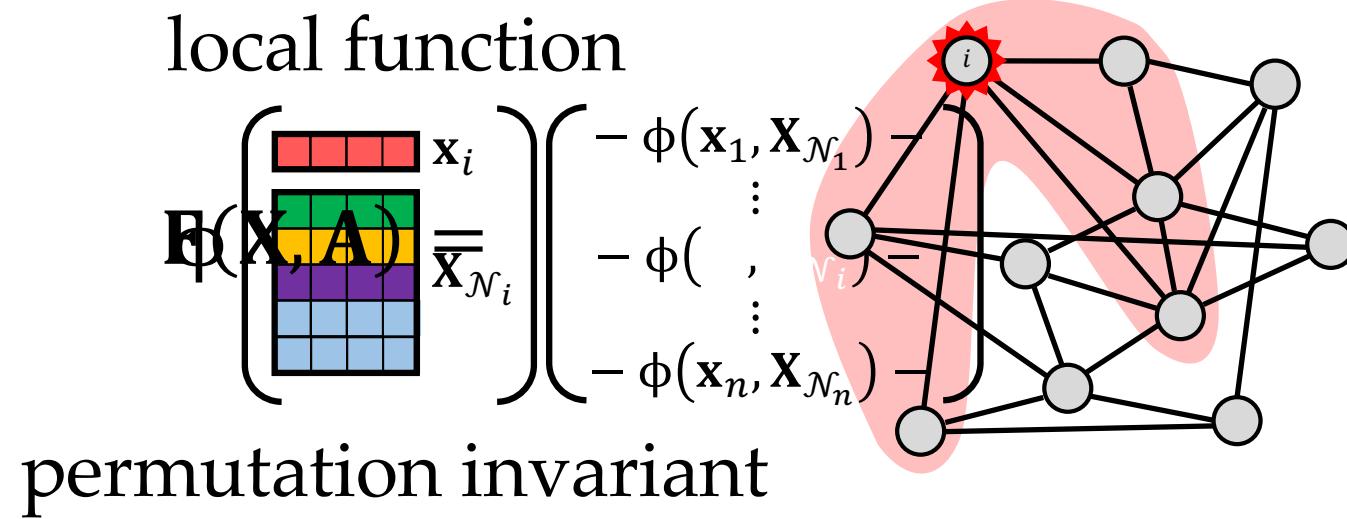
local function



Neighbour Aggregation



Neighbour Aggregation

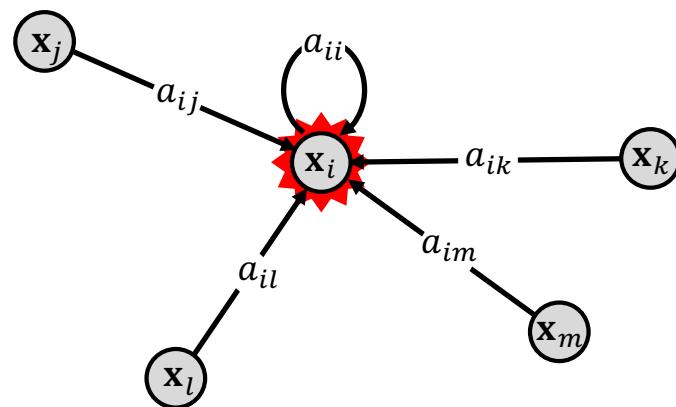


GNN Layer

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

permutation equivariant

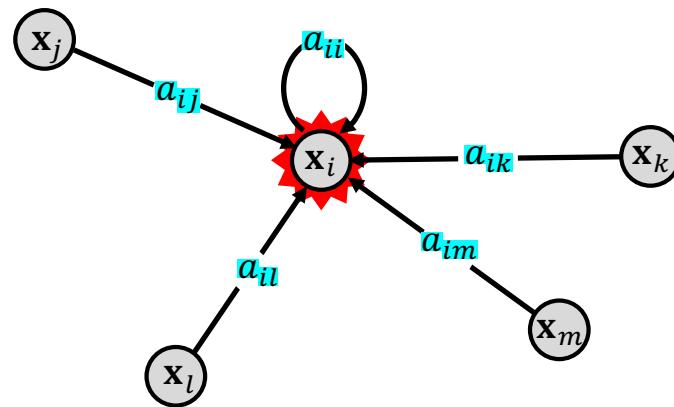
Convolutional GNNs



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

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

Convolutional GNNs

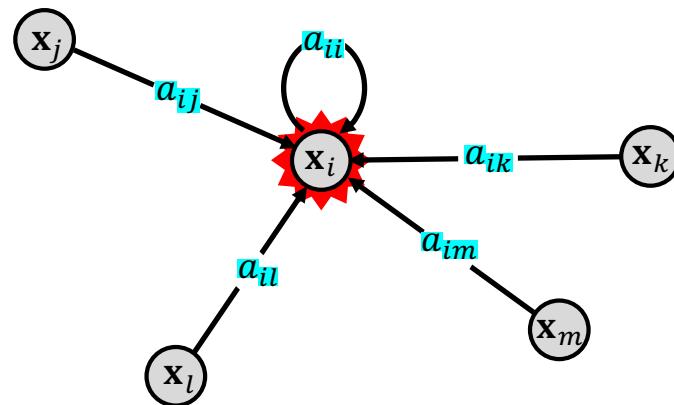


$$\mathbf{x}_i \leftarrow \sigma \left(\sum_{j \in \mathcal{N}_i \cup \{i\}} a_{ij} \psi(\mathbf{x}_j) \right)$$

*graph
adjacency*

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

Convolutional GNNs

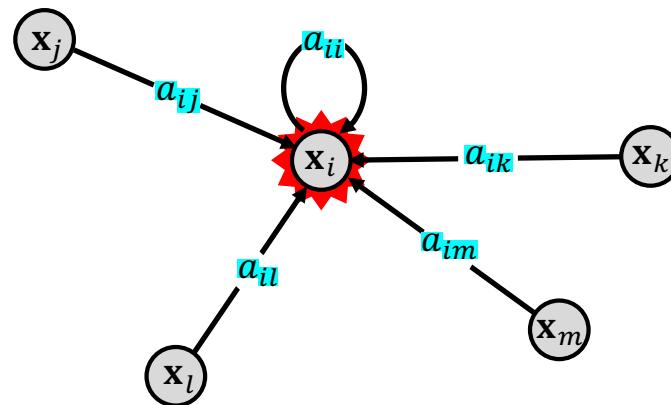


$$\mathbf{x}_i \leftarrow \sigma \left(\sum_{j \in \mathcal{N}_i \cup \{i\}} a_{ij} \psi(\mathbf{x}_j) \right)$$

nonlinear activation graph adjacency node-wise transformation

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

Convolutional GNNs



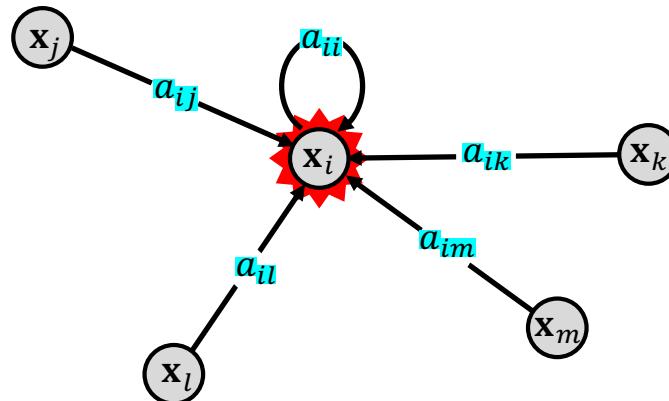
$$\mathbf{x}_i \leftarrow \sigma \left(\sum_{j \in \mathcal{N}_i \cup \{i\}} a_{ij} \mathbf{W} \mathbf{x}_j \right)$$

nonlinear activation graph adjacency node-wise linear transformation

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

Convolutional GNNs

- Simplest GNN
- Highly scalable
- Industrial use cases
- Folklore: works only on homophilic graphs

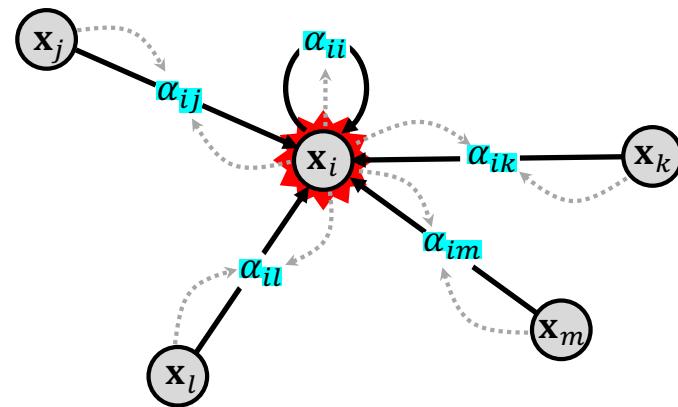


$$\mathbf{X} \leftarrow \sigma(\mathbf{A}\mathbf{X}\mathbf{W})$$

diffusion *channel mixing*
 $n \times n$ $d \times d$

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

Attentional GNNs

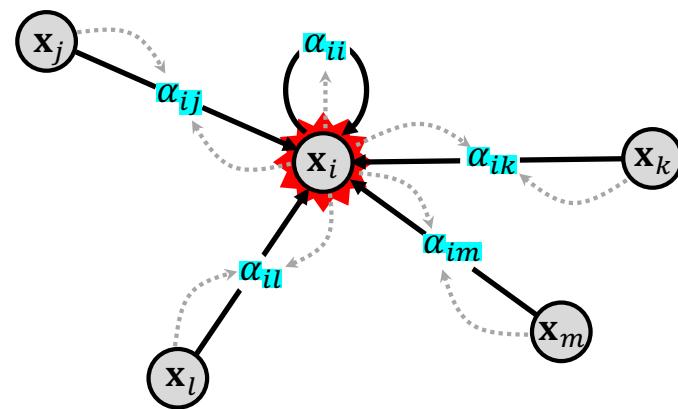


$$\mathbf{x}_i \leftarrow \sigma \left(\sum_{j \in \mathcal{N}_i \cup \{i\}} \alpha_{ij}(\mathbf{x}_i, \mathbf{x}_j) \psi(\mathbf{x}_j) \right)$$

learnable attention weights

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

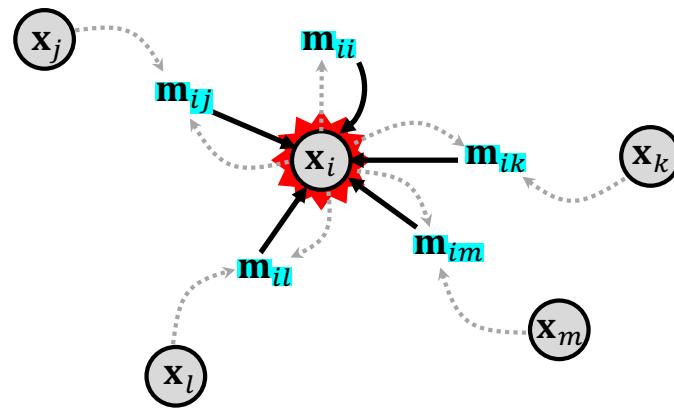
Attentional GNNs



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

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

Message-Passing GNNs



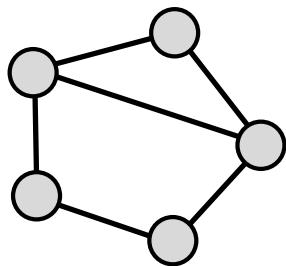
$$\mathbf{x}_i \leftarrow \sigma \left(\sum_{j \in \mathcal{N}_i \cup \{i\}} \psi(\mathbf{x}_i, \mathbf{x}_j) \right)$$

message from
node j to node i

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

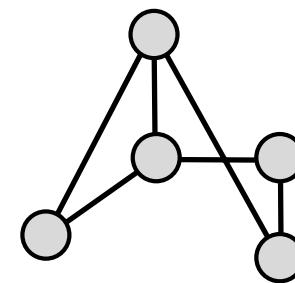
EXPRESSIVE POWER OF GNNS & WEISFEILER-LEHMAN TEST

Graph isomorphism



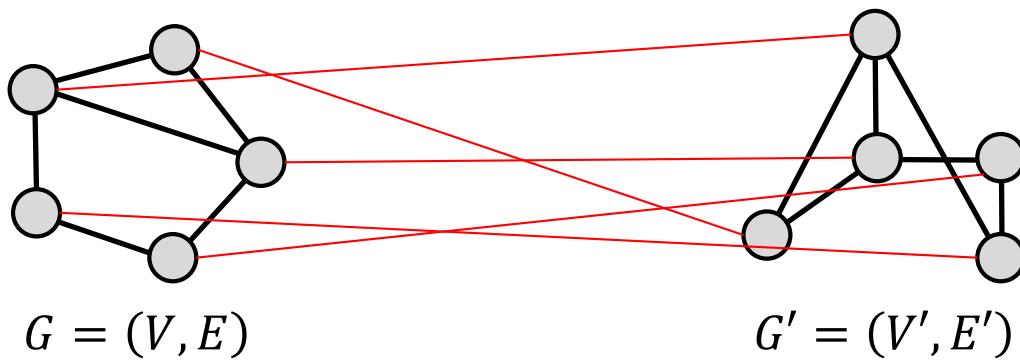
$$G = (V, E)$$

“=”



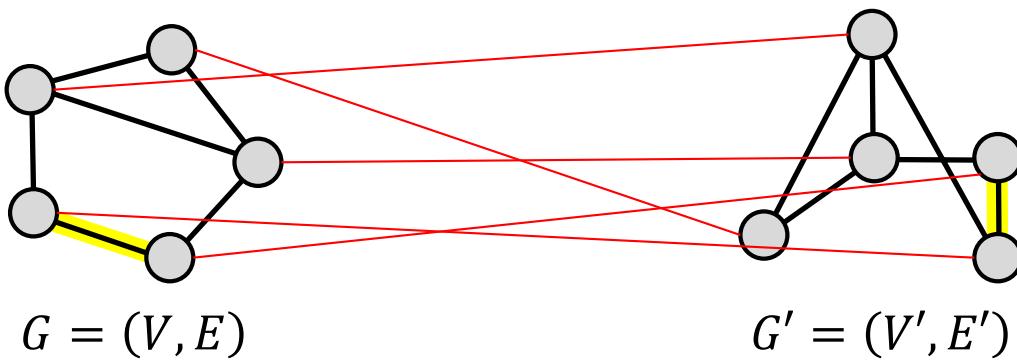
$$G' = (V', E')$$

Graph isomorphism



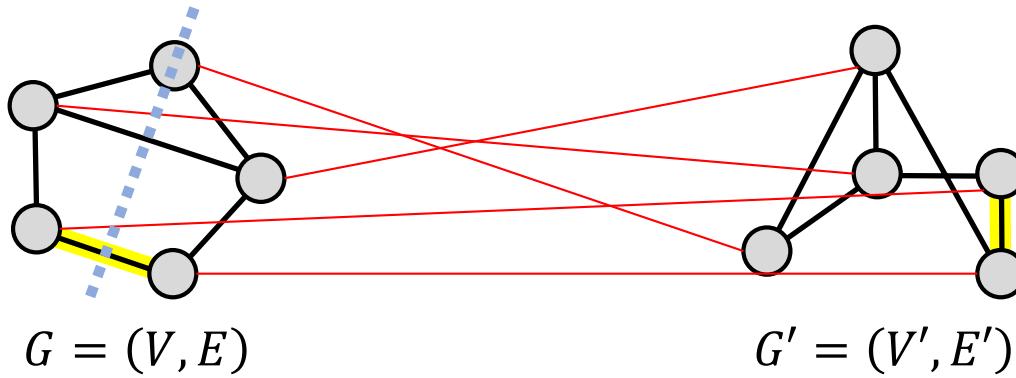
- Two graphs $G = (V, E)$ and $G' = (V', E')$ are **isomorphic** if there exists an *edge-preserving bijection* $\varphi: V \rightarrow V'$ s.t. $u \sim v$ in G iff $\varphi(u) \sim \varphi(v)$ in G'

Graph isomorphism



- Two graphs $G = (V, E)$ and $G' = (V', E')$ are **isomorphic** if there exists an *edge-preserving bijection* $\varphi: V \rightarrow V'$ s.t. $u \sim v$ in G iff $\varphi(u) \sim \varphi(v)$ in G'

Graph isomorphism



- Two graphs $G = (V, E)$ and $G' = (V', E')$ are **isomorphic** if there exists an *edge-preserving bijection* $\varphi: V \rightarrow V'$ s.t. $u \sim v$ in G iff $\varphi(u) \sim \varphi(v)$ in G'

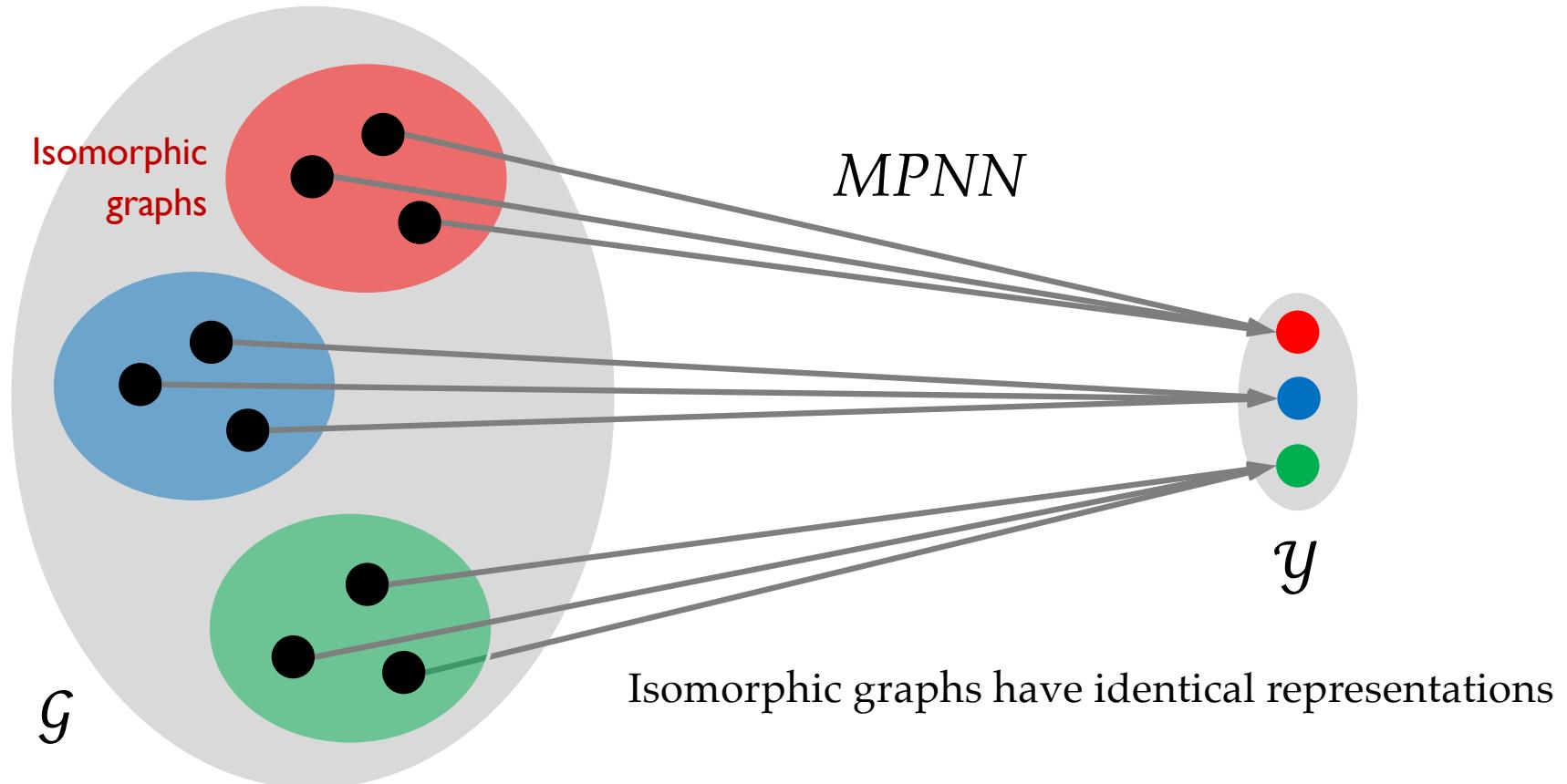
Note: φ is not unique where the graph has symmetries (edge-preserving automorphism)

Universal Approximation on Graphs

Theorem: A class of functions is universally approximating permutation-invariant functions on graphs with finite node features iff it can discriminate graph isomorphisms.

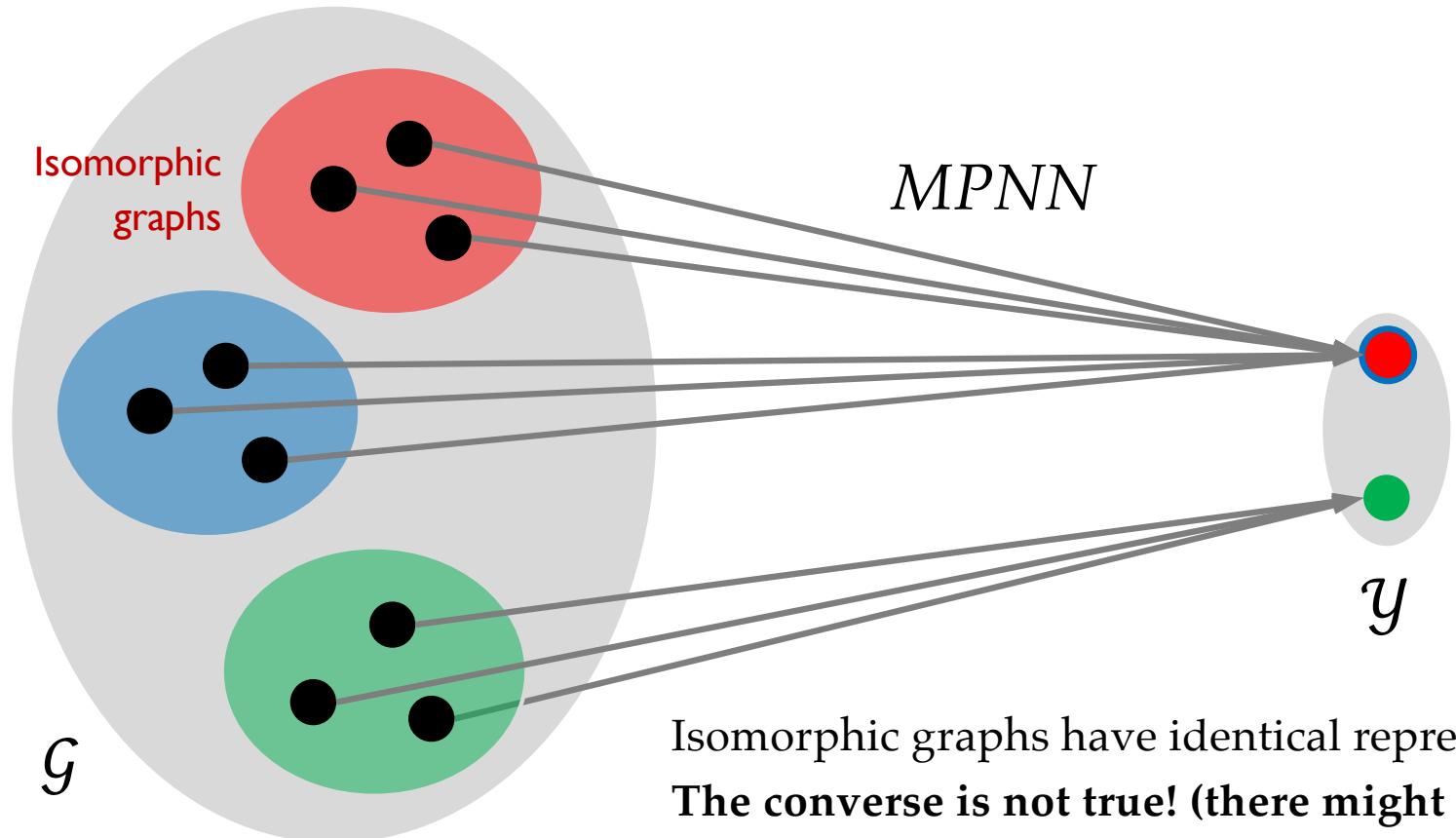
Universal approximation on graphs is equivalent to graph isomorphism testing

What graphs can MPNNs represent?



Adapted from Bevilacqua et al.
(LoG Tutorial 2022)

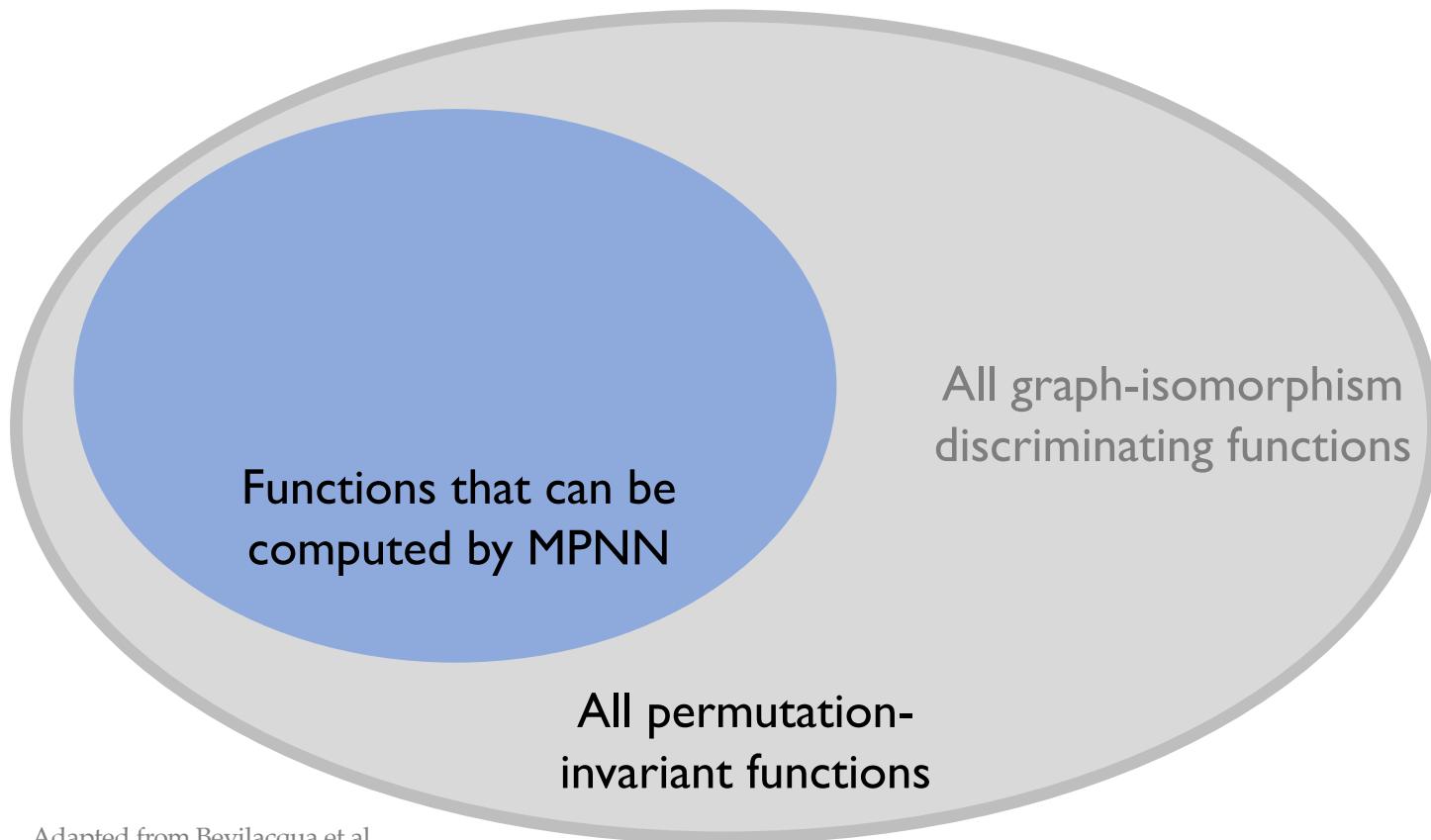
What graphs can MPNNs represent?



Isomorphic graphs have identical representations
The converse is not true! (there might be
indistinguishable non-isomorphic graphs)

Adapted from Bevilacqua et al.
(LoG Tutorial 2022)

Expressive power of MPNNs

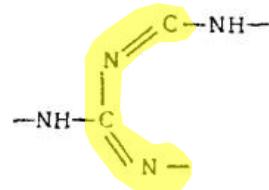
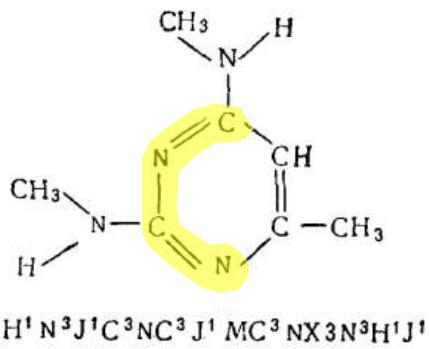


Adapted from Bevilacqua et al.
(LoG Tutorial 2022)

Weisfeiler-Lehman Test & Chemical precursors of GNNs



George Vlăduț



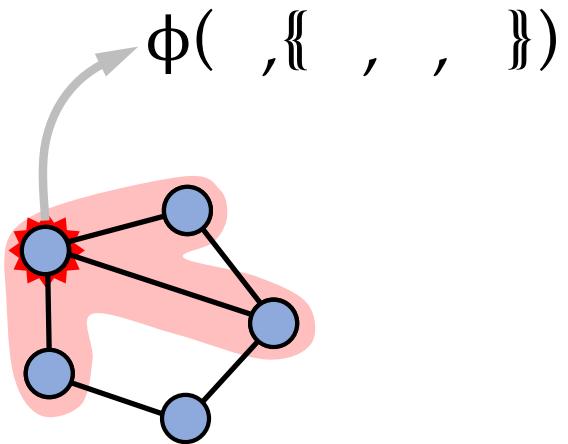
Andrey Lehman



Boris Weisfeiler

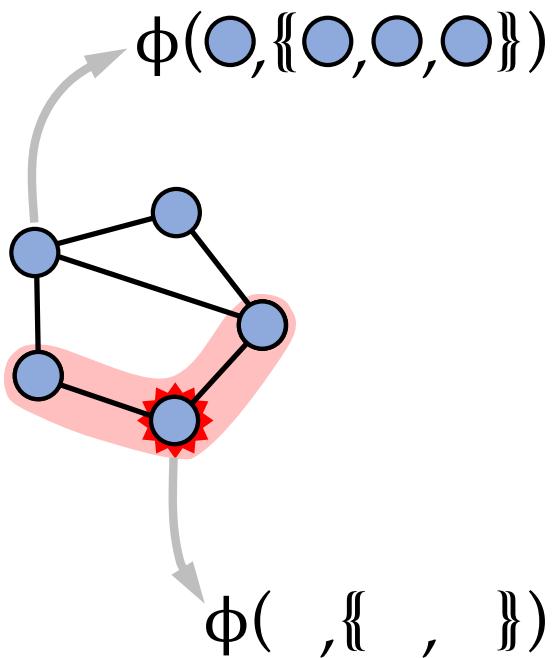
Vlăduț et al. 1959; Weisfeiler, Lehman 1968

Weisfeiler-Lehman Test



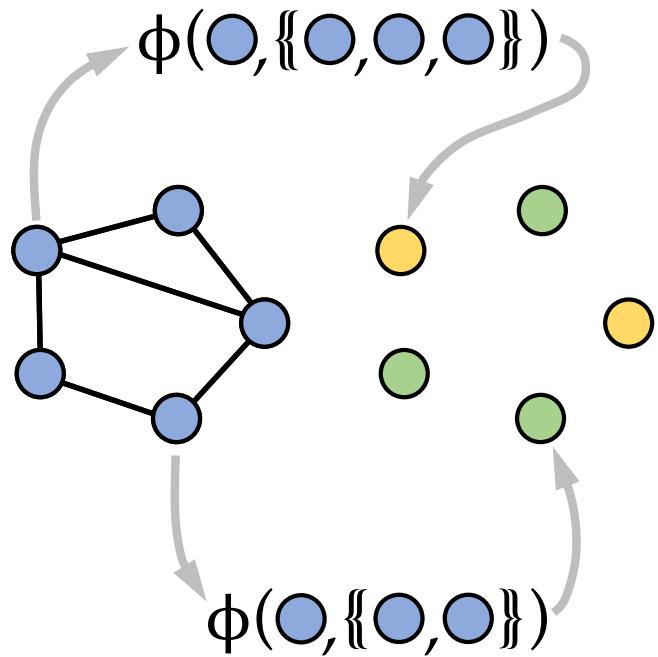
Weisfeiler, Lehman 1968

Weisfeiler-Lehman Test



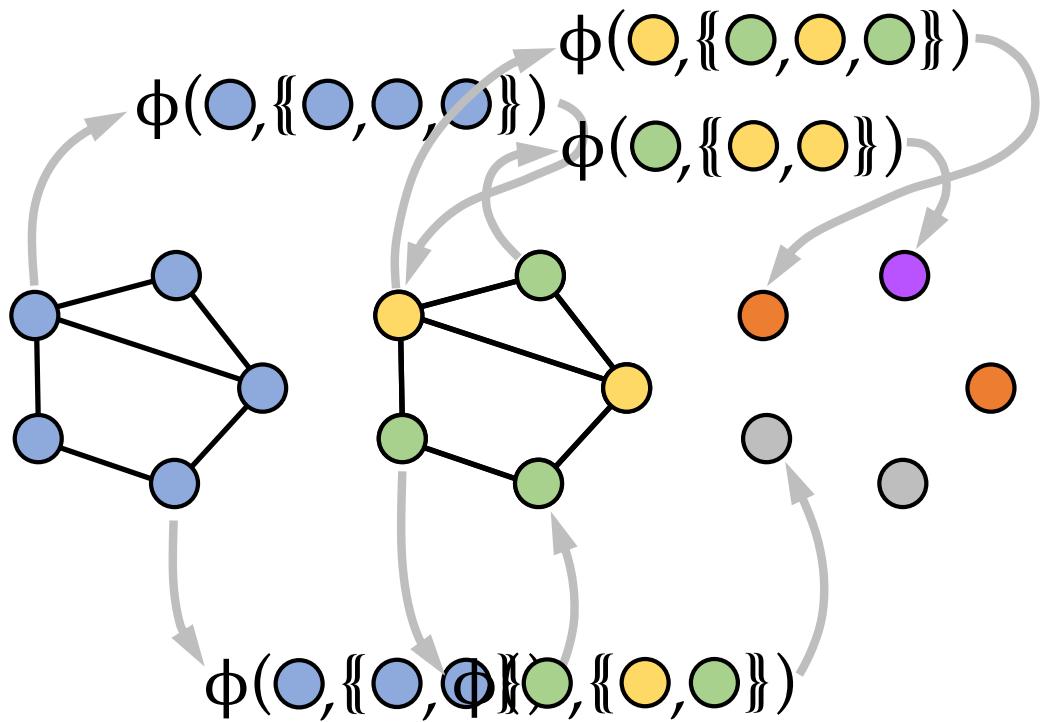
Weisfeiler, Lehman 1968

Weisfeiler-Lehman Test



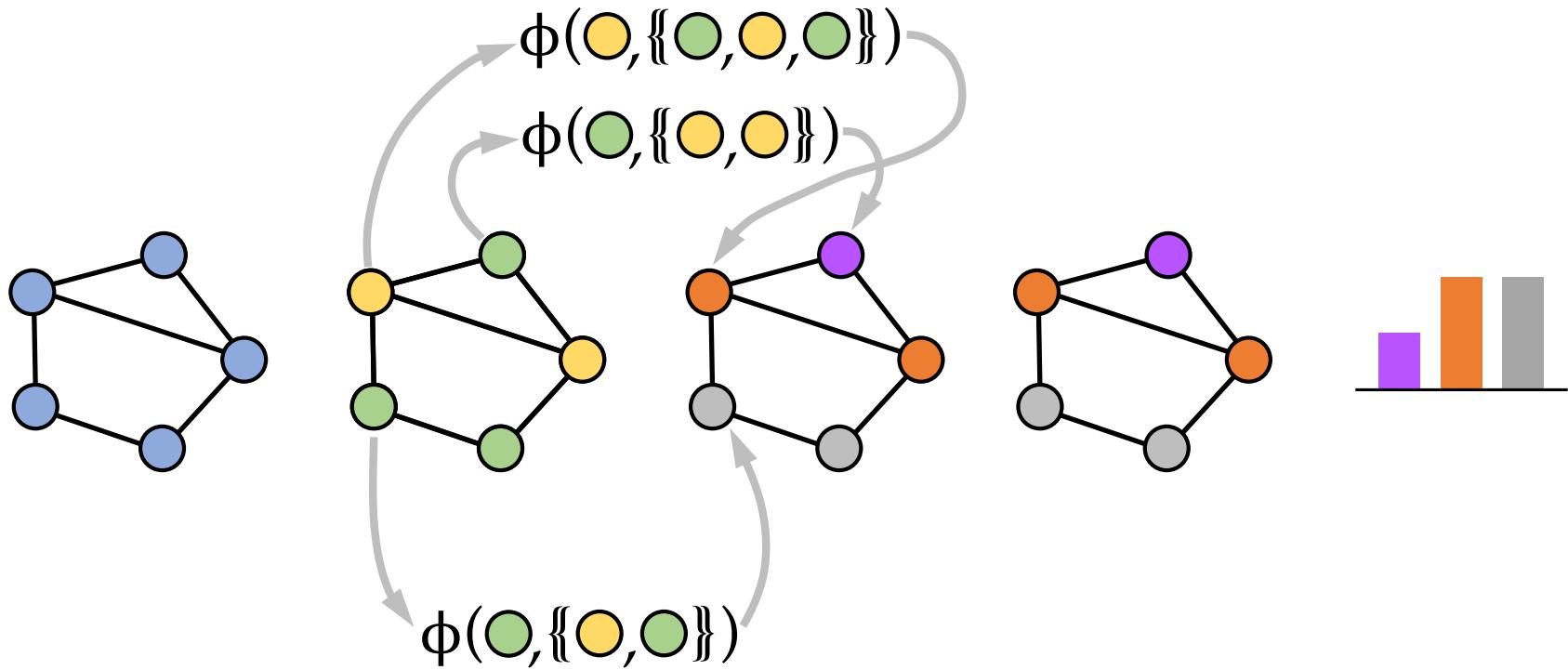
Weisfeiler, Lehman 1968

Weisfeiler-Lehman Test



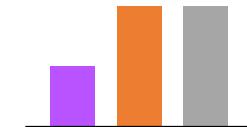
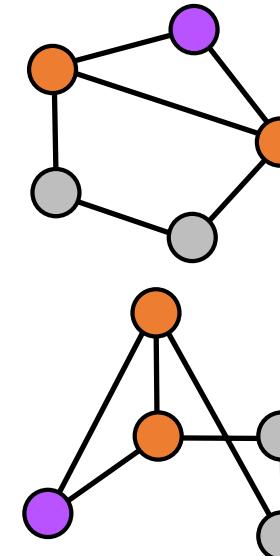
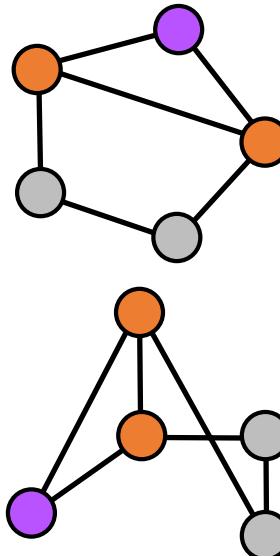
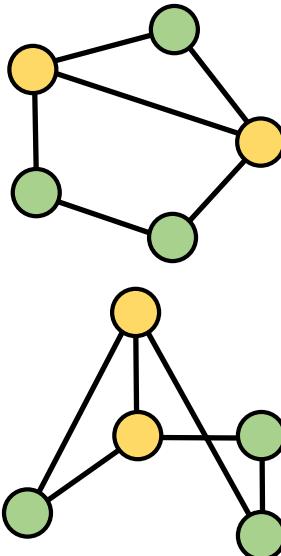
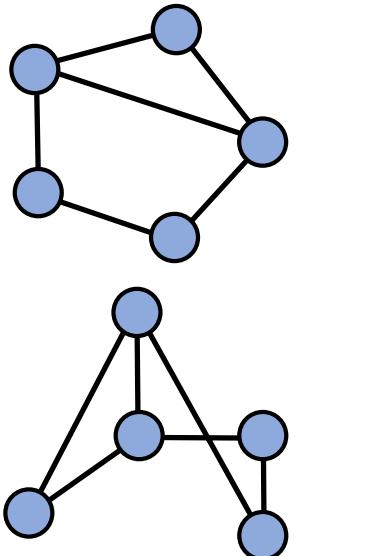
Weisfeiler, Lehman 1968

Weisfeiler-Lehman Test

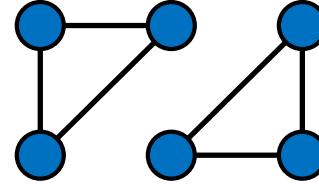
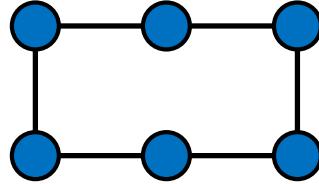


Weisfeiler, Lehman 1968

Weisfeiler-Lehman Test



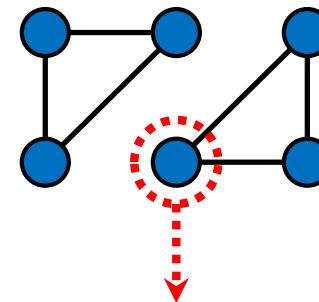
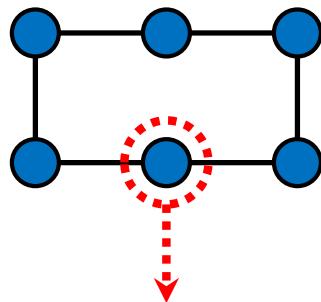
Weisfeiler, Lehman 1968



non-isomorphic graphs that are WL-equivalent

**Necessary but insufficient condition for
graph isomorphism!**

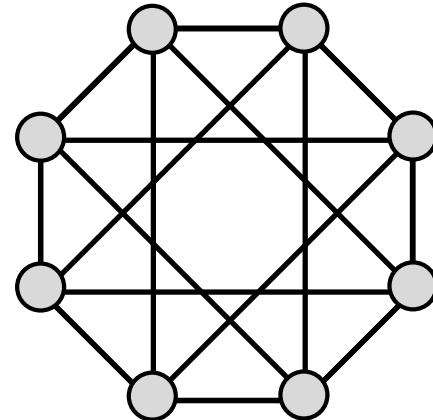
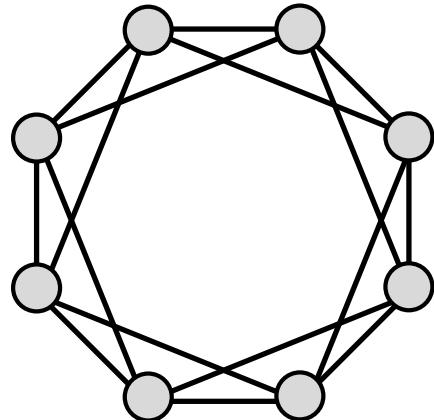
What does WL test see?



=

What WL cannot test?

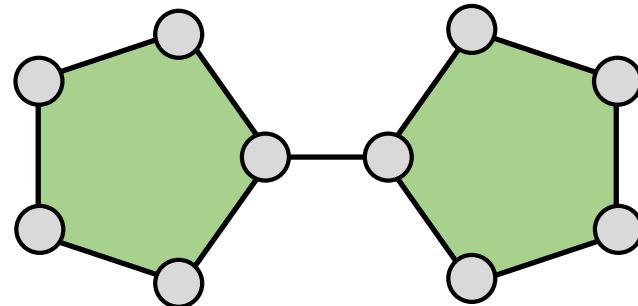
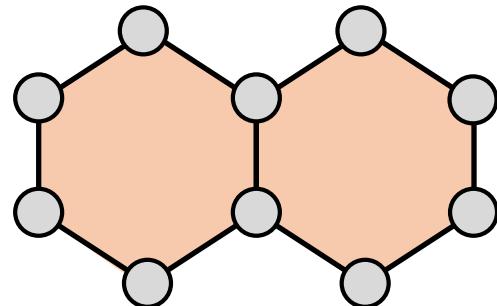
Example of non-isomorphic graphs that cannot be distinguished by Weisfeiler-Lehman test (outputs “possibly isomorphic”)



r-regular graphs (deg=r at every node) with the same number of nodes

What WL cannot test?

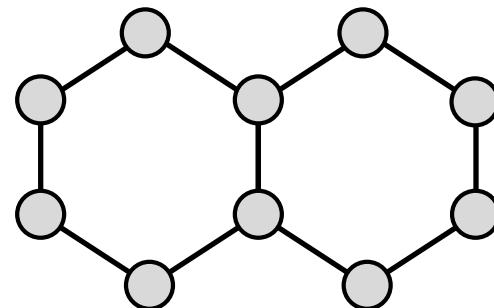
Example of non-isomorphic graphs that cannot be distinguished by Weisfeiler-Lehman test (outputs “possibly isomorphic”)



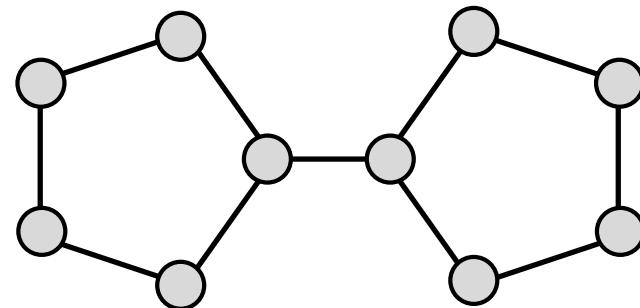
Any induced connected pattern with ≥ 3 nodes (triangles, cycles, etc.)

What WL cannot test?

Example of non-isomorphic graphs that cannot be distinguished by Weisfeiler-Lehman test (outputs “possibly isomorphic”)



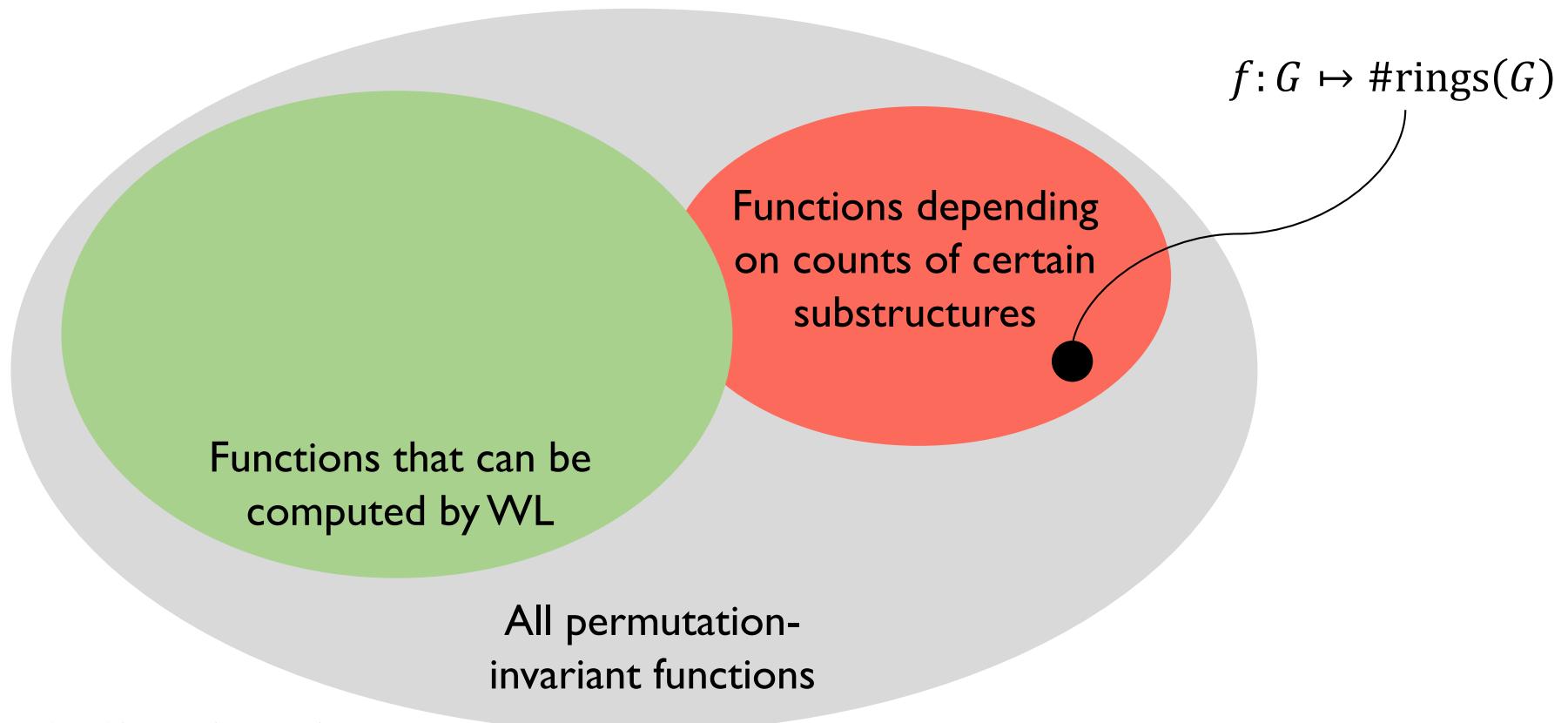
decalin



bicyclopnetyl

Important implications e.g. in chemistry!

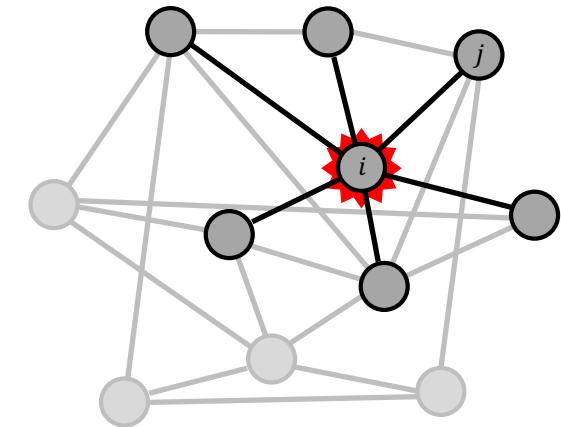
Expressive power of Weisfeiler-Lehman



Adapted from Bevilacqua et al.
(LoG Tutorial 2022)

MPNNs vs WL

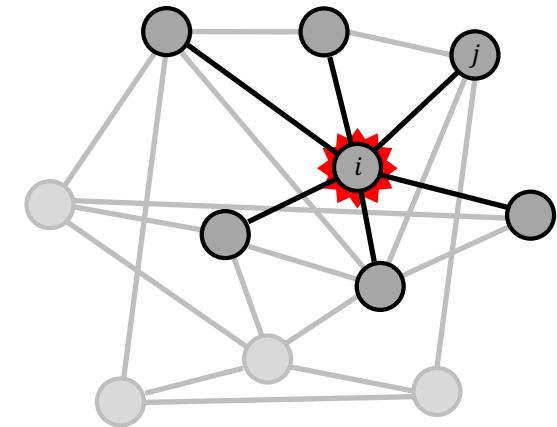
- **WL-test:** $\mathbf{x}_i \leftarrow \phi\left(\mathbf{x}_i, \{\{\mathbf{x}_j : j \in \mathcal{N}_i\}\}\right)$
where ϕ is *injective* (hash function)
- **MPNN:** $\mathbf{x}_i \leftarrow \phi\left(\mathbf{x}_i, \bigcup_{j \in \mathcal{N}_i} \psi(\mathbf{x}_i, \mathbf{x}_j)\right)$



**MPNN expressive power is upper-bounded
by the Weisfeiler-Lehman test**

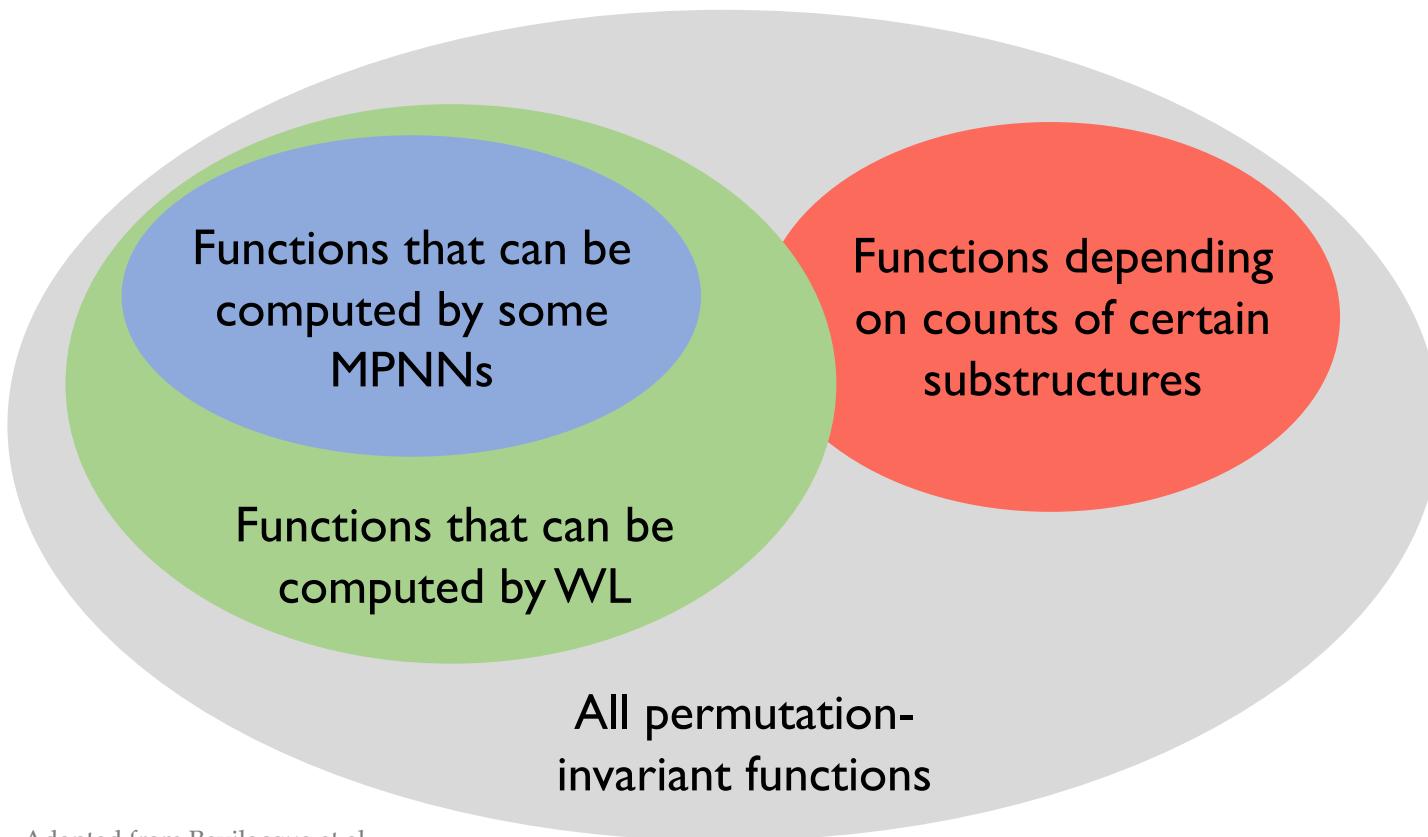
MPNNs vs WL

- **WL-test:** $\mathbf{x}_i \leftarrow \phi(\mathbf{x}_i, \{\{\mathbf{x}_j : j \in \mathcal{N}_i\}\})$
where ϕ is *injective* (hash function)
- **MPNN:** $\mathbf{x}_i \leftarrow \phi\left(\mathbf{x}_i, \bigcup_{j \in \mathcal{N}_i} \psi(\mathbf{x}_i, \mathbf{x}_j)\right)$



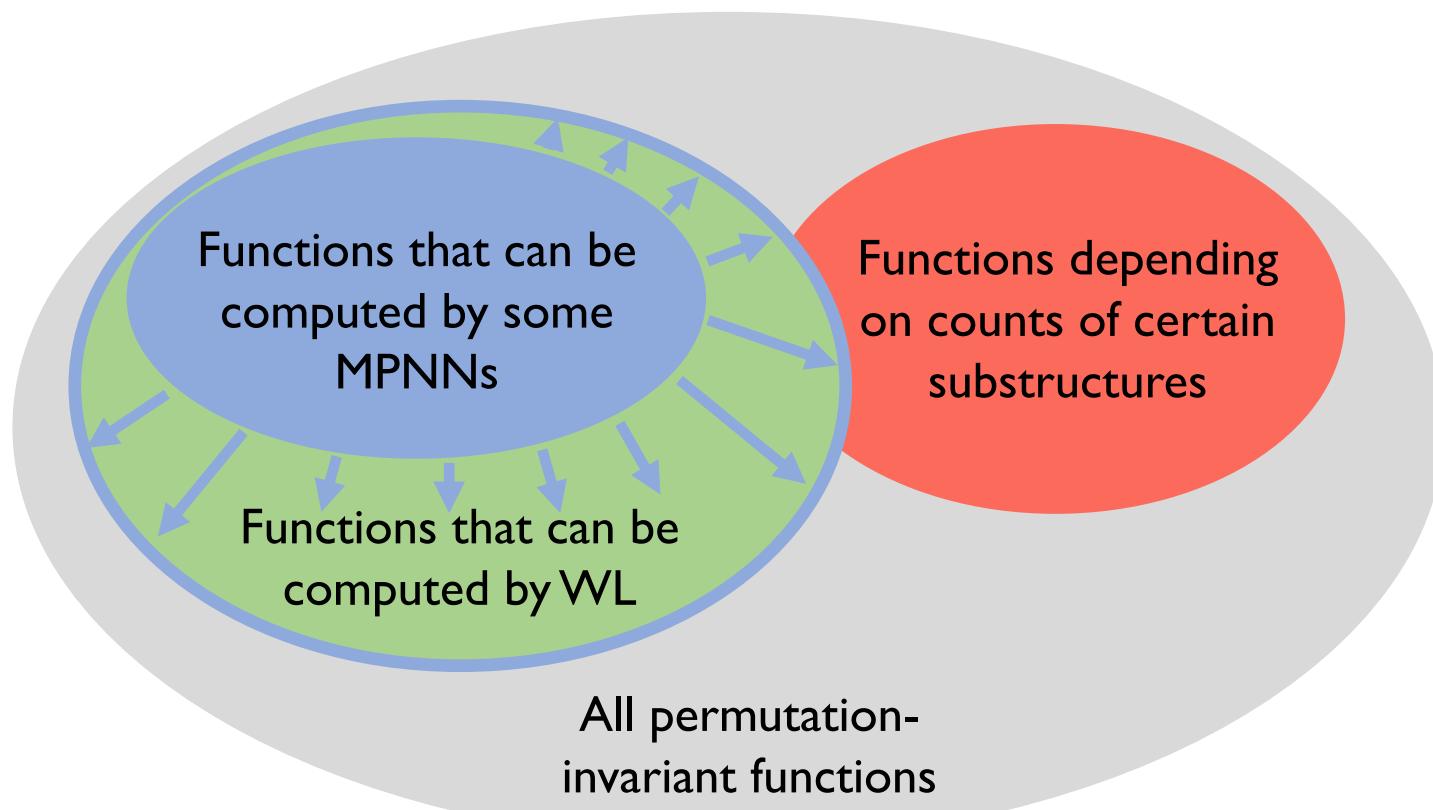
When is MPNN as expressive as the Weisfeiler-Lehman test?

Expressive power of MPNNs



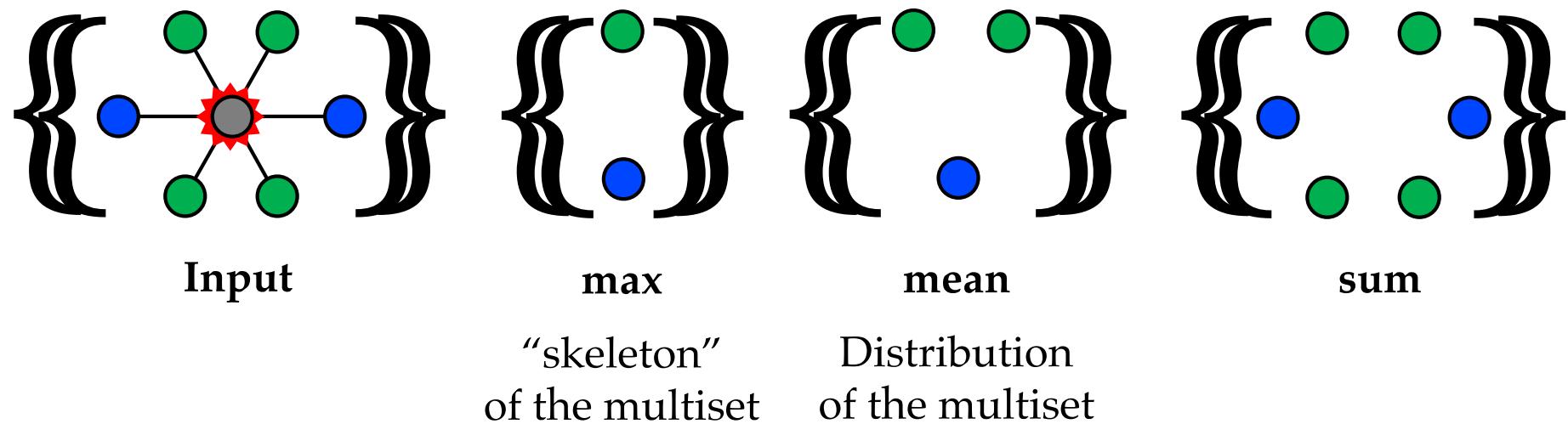
Adapted from Bevilacqua et al.
(LoG Tutorial 2022)

Expressive power of MPNNs

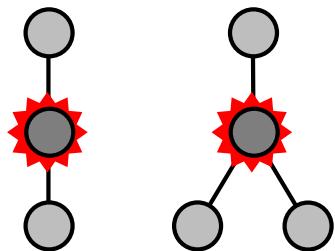


Adapted from Bevilacqua et al.
(LoG Tutorial 2022)

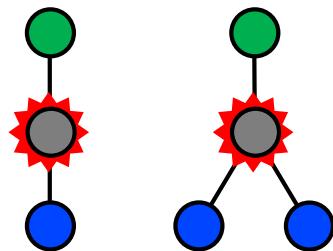
Are all Aggregators the same?



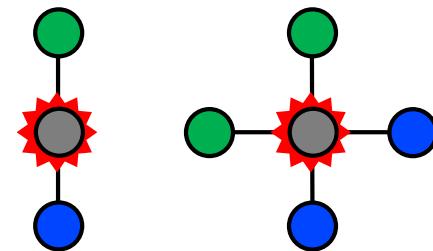
Are all Aggregators the same?



mean and **max**
fail to distinguish



max
fails to distinguish



mean and **max**
fail to distinguish

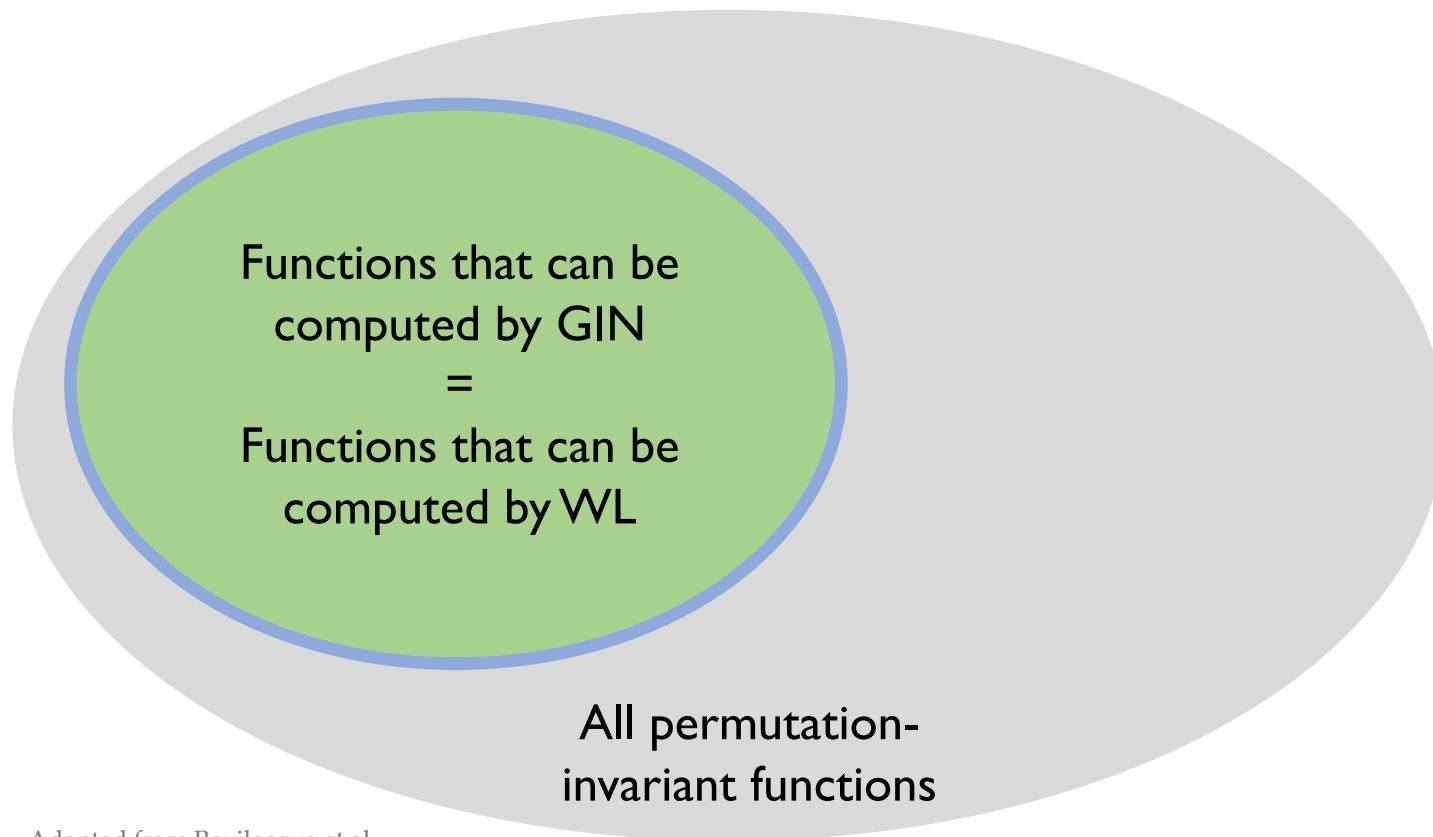
Graph Isomorphism Network (GIN)

Theorem: Assume graph node features are from a countable set. Then, an MPNN with injective aggregator \square , update function ϕ , and graph-wise readout function, is as powerful as the Weisfeiler-Lehman test.

- Assumption of **discrete countable features** (often not the case in practice)
- GIN: uses an injective multiset function of the form

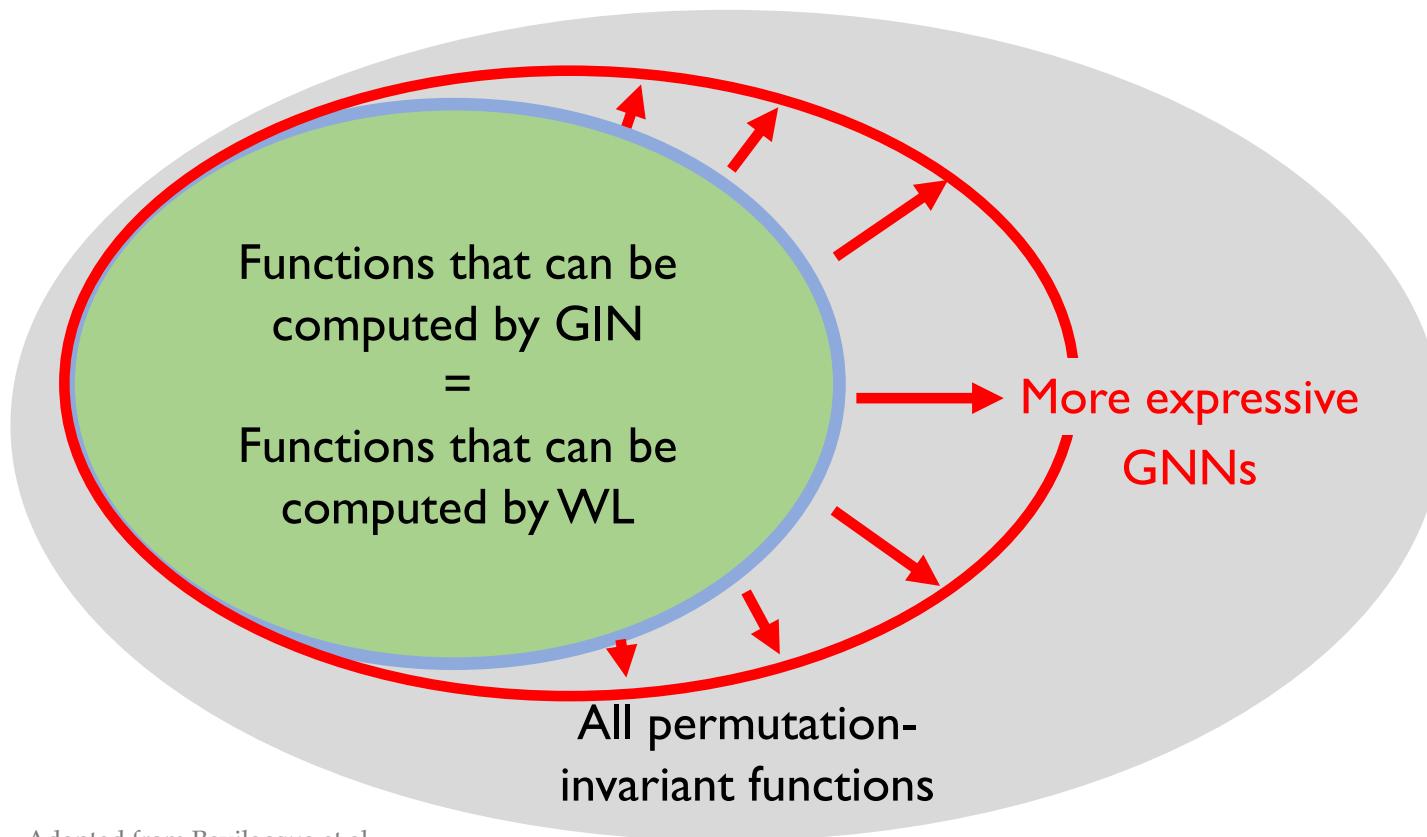
$$\text{MLP} \left((1 + \epsilon) \mathbf{x}_i + \sum_{j \in \mathcal{N}_i} \mathbf{x}_j \right)$$

Expressive power of GIN (“best MPNN”)



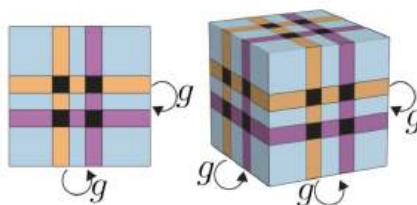
Adapted from Bevilacqua et al.
(LoG Tutorial 2022)

Expressive power of GIN ("best MPNN")



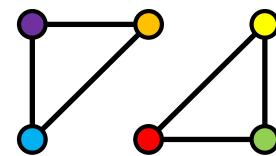
Adapted from Bevilacqua et al.
(LoG Tutorial 2022)

Towards More Expressive GNNs



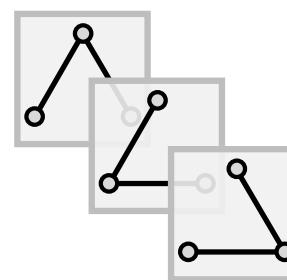
Higher-order
WL tests

Maron et al. 2019
Morris et al. 2019



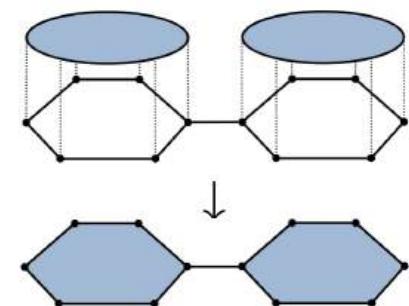
Positional &
Structural encoding

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



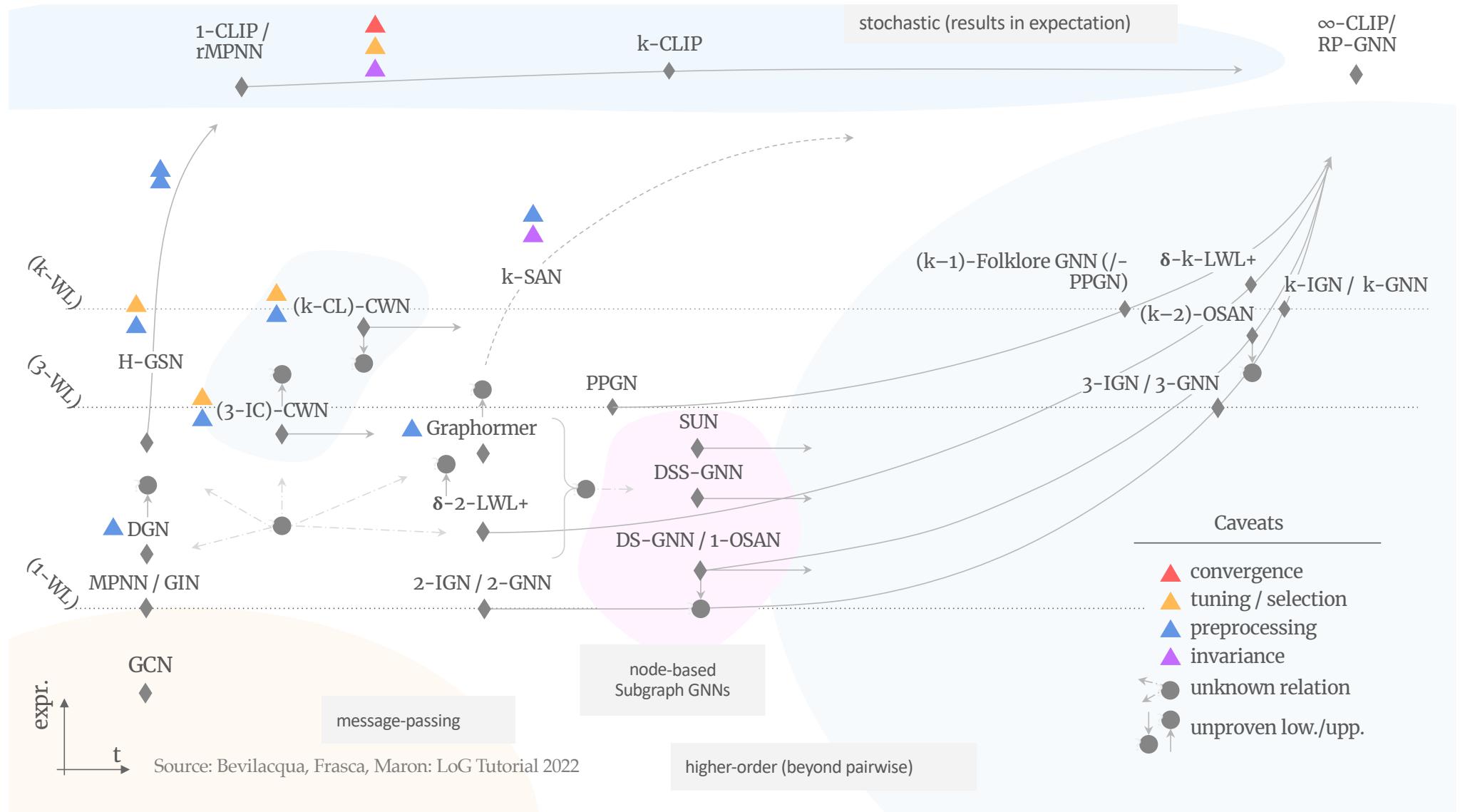
Subgraph
GNNs

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



Topological
message passing

Bodnar, Frasca et B 2021

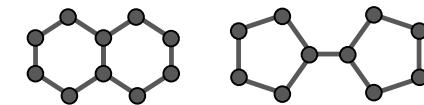


Theory

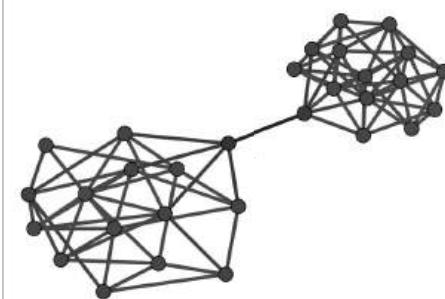


WL test = expressive power

Practice

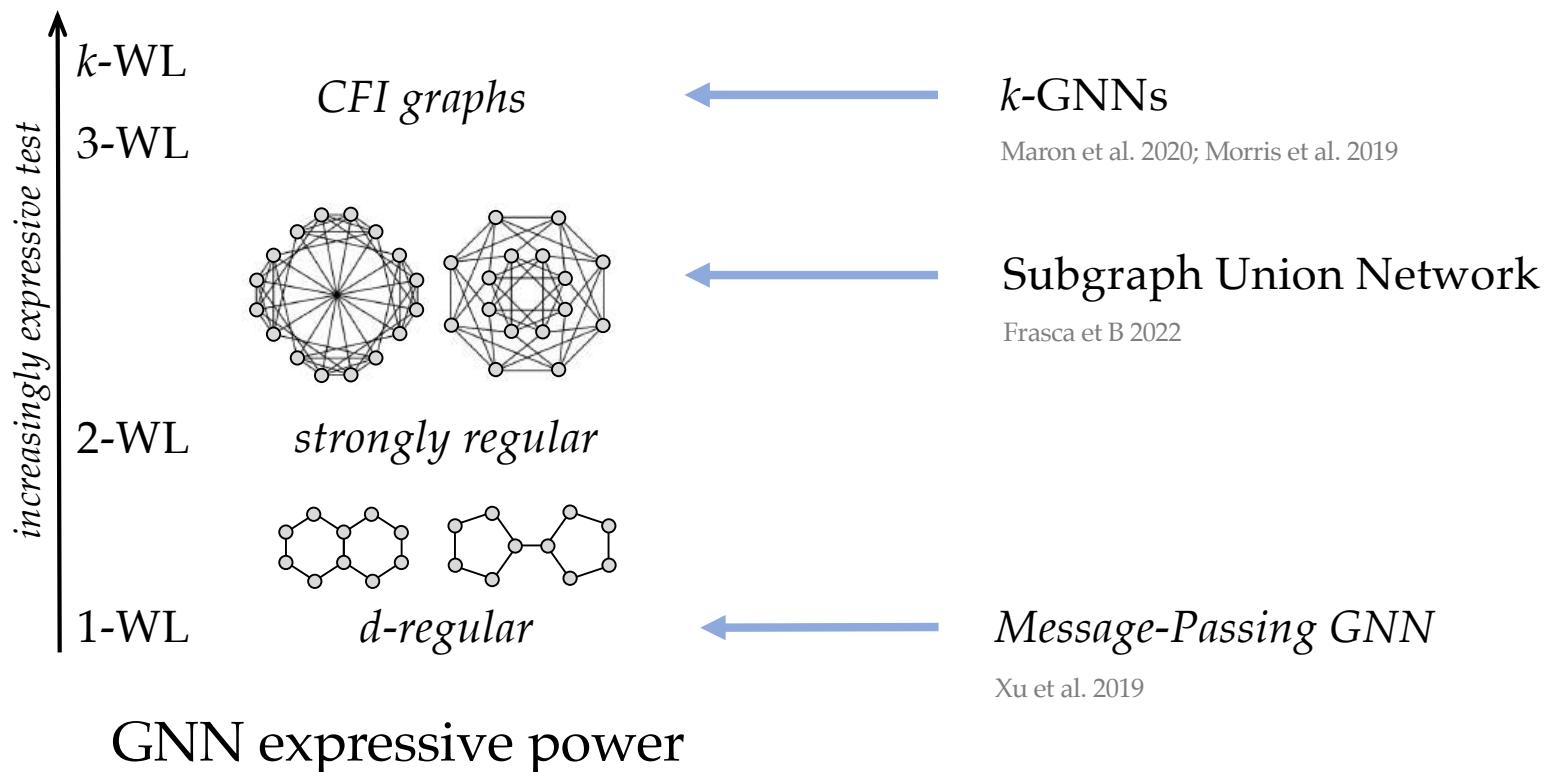


Some non-isomorphic graphs
cannot be tested by WL



Some graphs may be
unfriendly for message passing

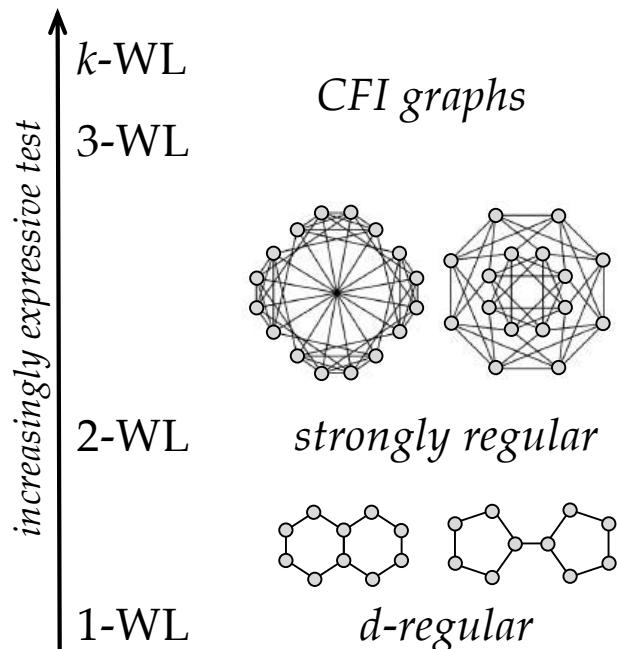
More expressive isomorphism tests (k -WL hierarchy)



GNN expressive power

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

More expressive isomorphism tests (k -WL hierarchy)

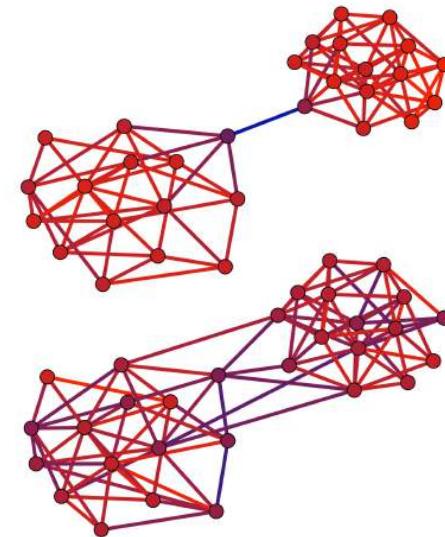


GNN expressive power

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

Decouple input graph from the computational graph

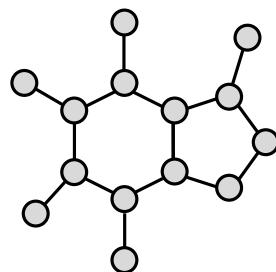
Gap between
Theory & Practice



Graph rewiring

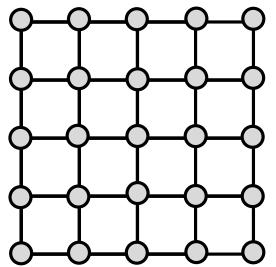
Alon, Yahav 2020 (bottlenecks); Hamilton et al. 2017 (neighbour sampling);
Klicpera et al. 2019 (diffusion); Topping, Di Giovanni et al. 2022 (Ricci flow);
Deac et al. 2022 (expanders); Barbero et al., Di Giovanni 2024 (LASER)

Classical GNNs: propagate information on
the input graph

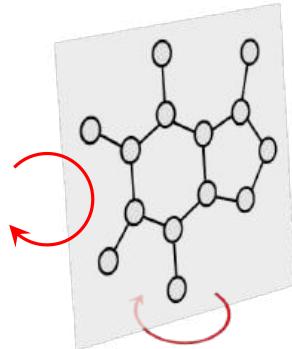


GNN
input graph

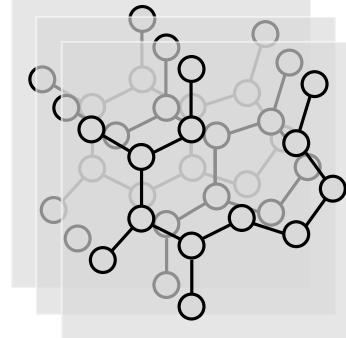
MORE STRUCTURE



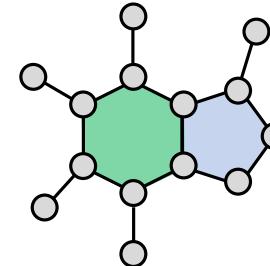
CNN
canonical node
ordering



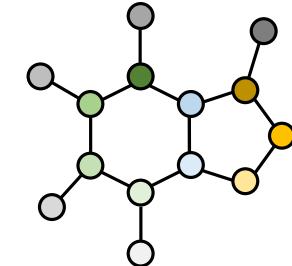
Equivariant GNN
+data symmetry
group



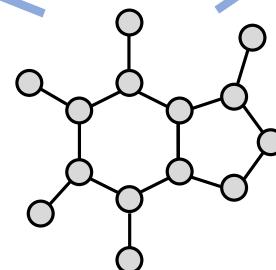
Subgraph GNN
product symmetry
group



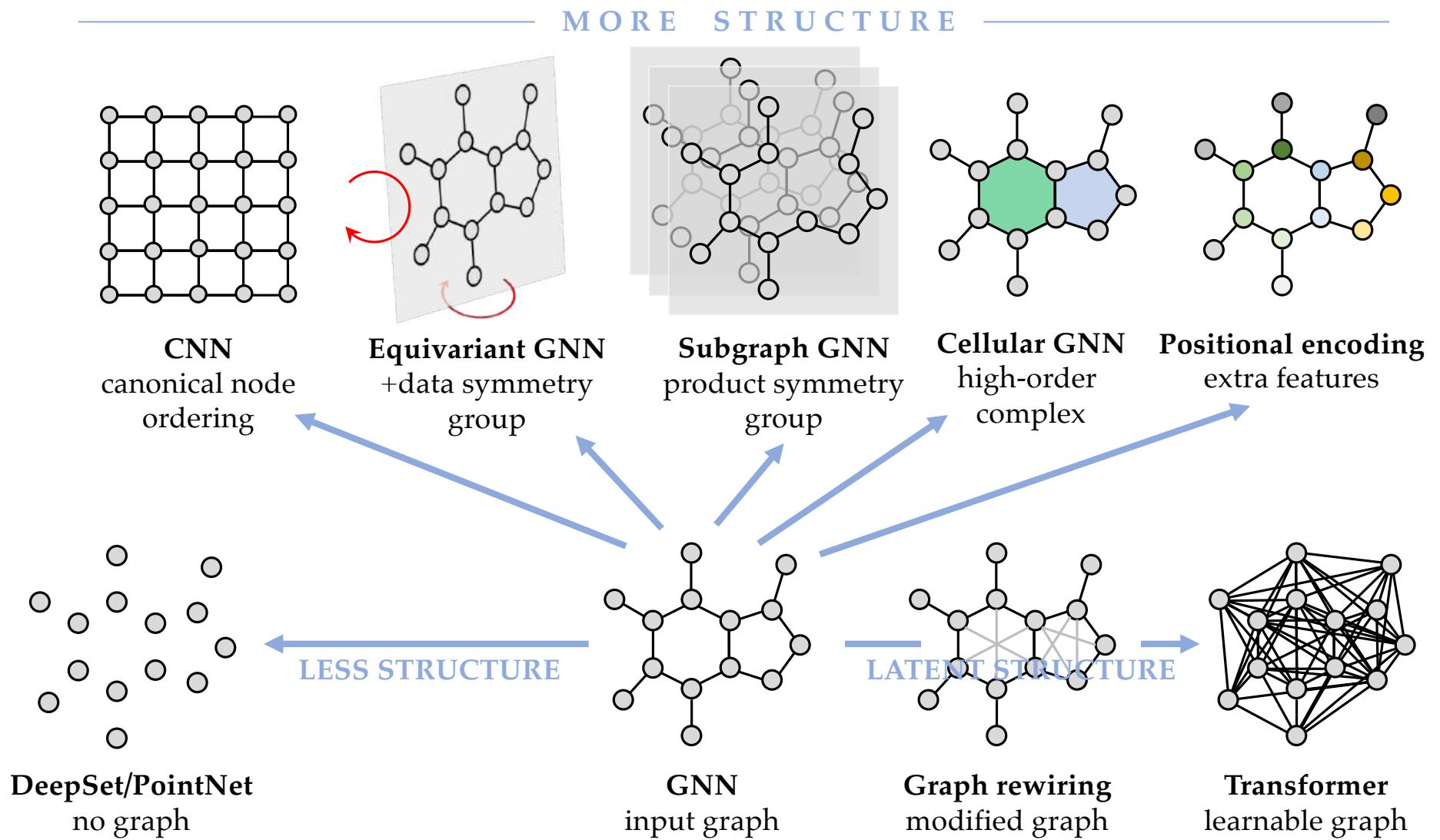
Cellular GNN
high-order
complex



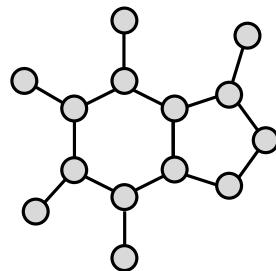
Positional encoding
extra features



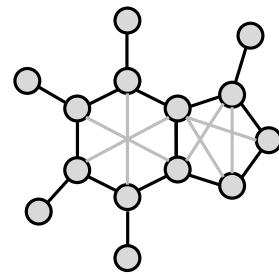
GNN
input graph



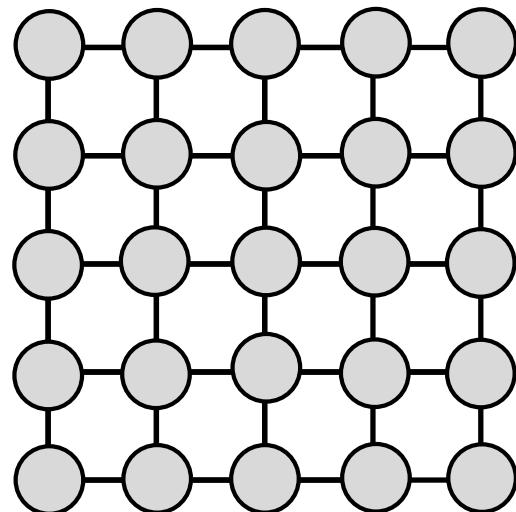
Classical GNNs: propagate information on
the input graph



“Modern” GNNs: propagate information on
some other “better” graph



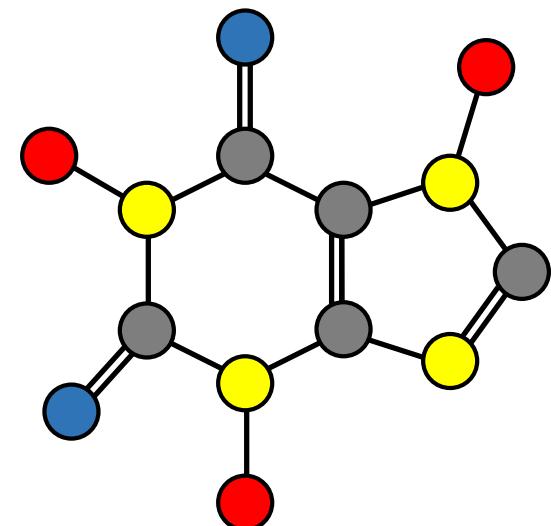
Graphs vs Meshes vs Grids



Grid

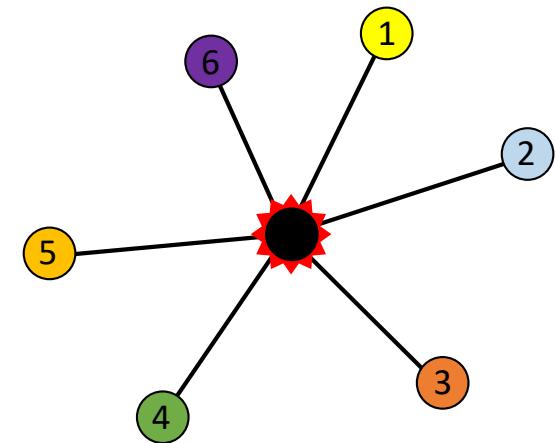
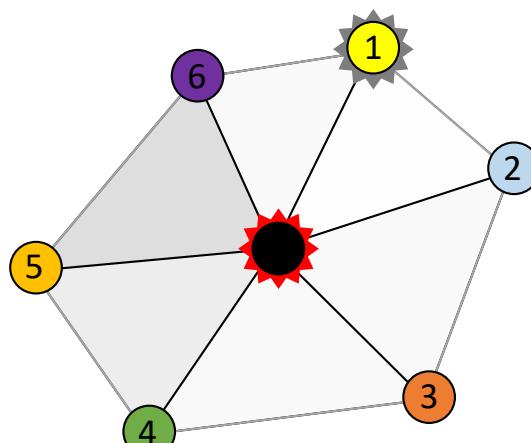
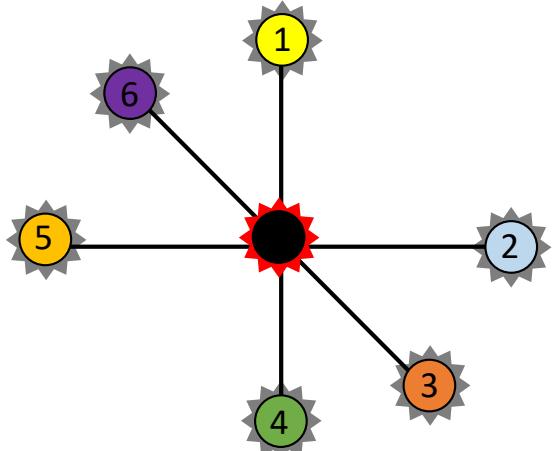


Mesh

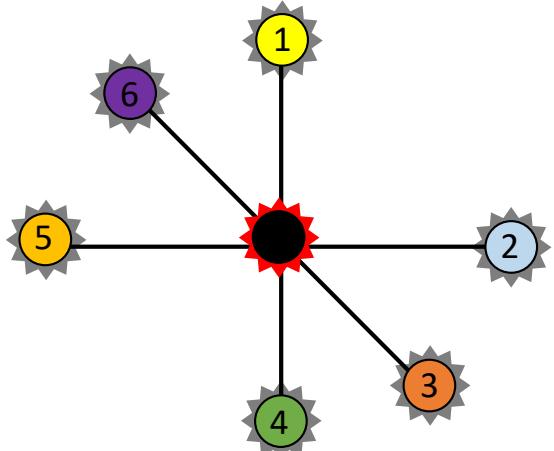


Graph

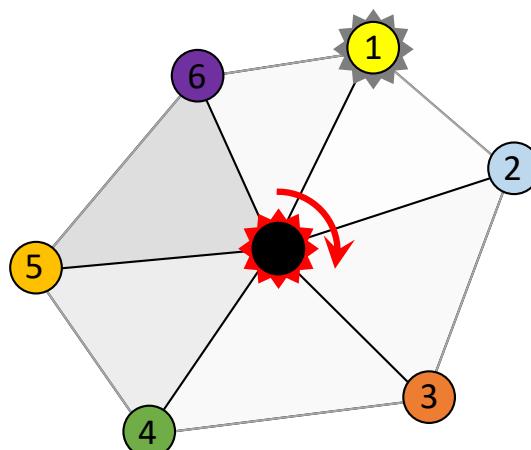
Graphs vs Meshes vs Grids



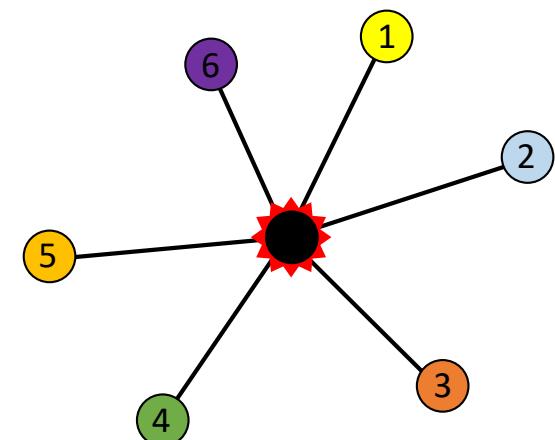
Graphs vs Meshes vs Grids



Grid
Fixed

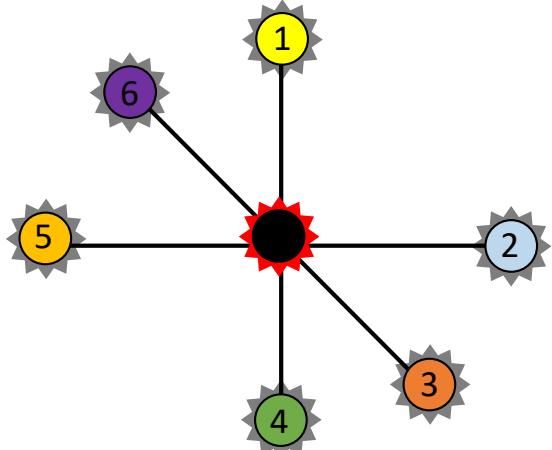


Mesh
Rotation

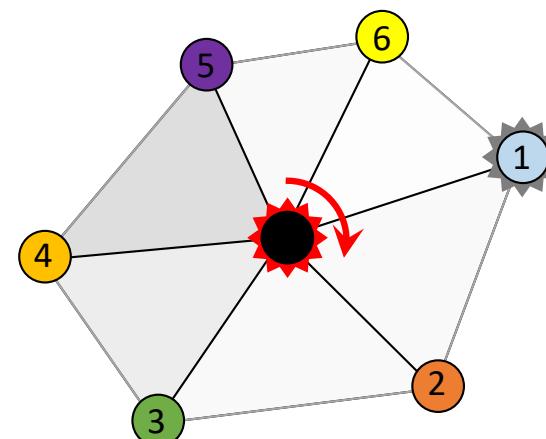


Graph

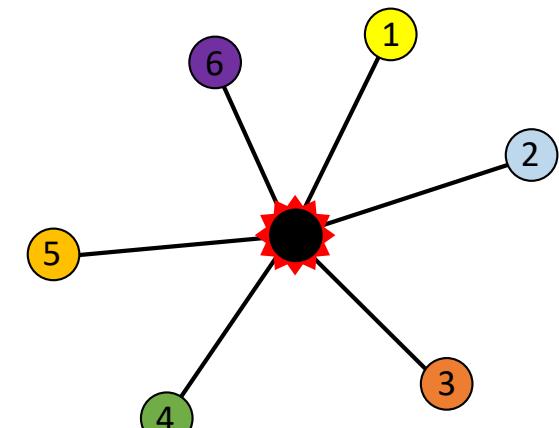
Graphs vs Meshes vs Grids



Grid
Fixed



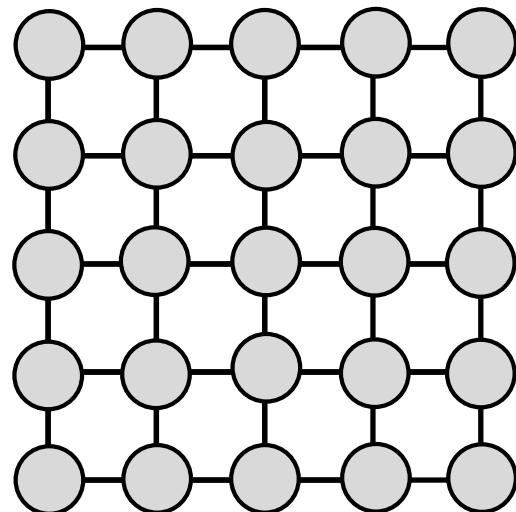
Mesh
Rotation



Graph
Permutation

Graphs have the least structure

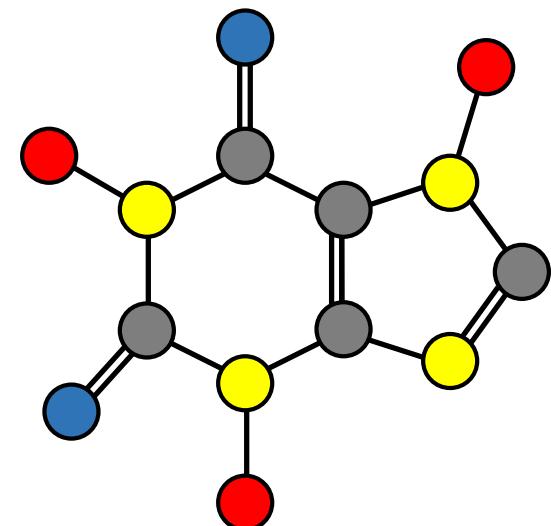
Graphs vs Meshes vs Grids



Grid

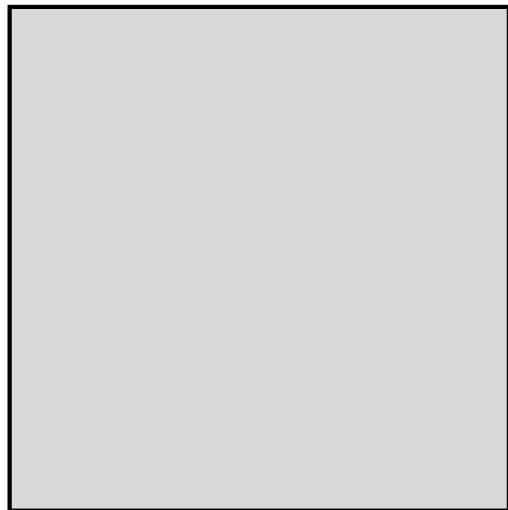


Mesh



Graph

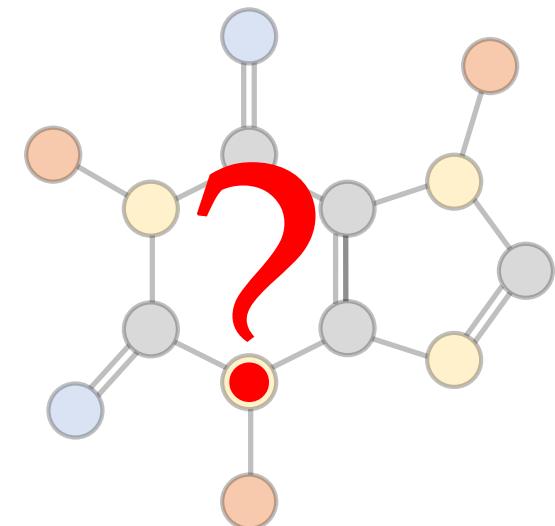
Graphs vs Meshes vs Grids



Grid

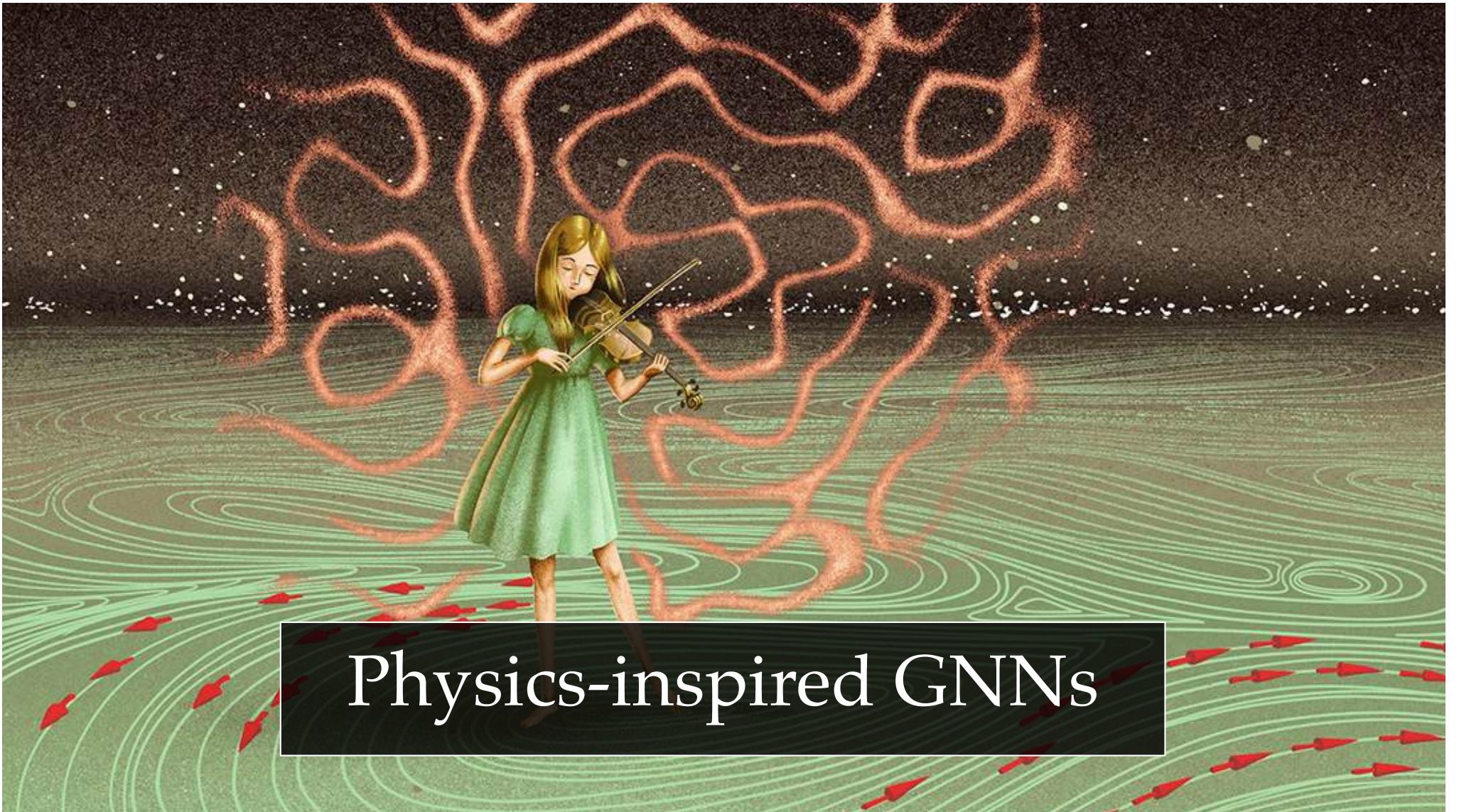


Mesh



Graph

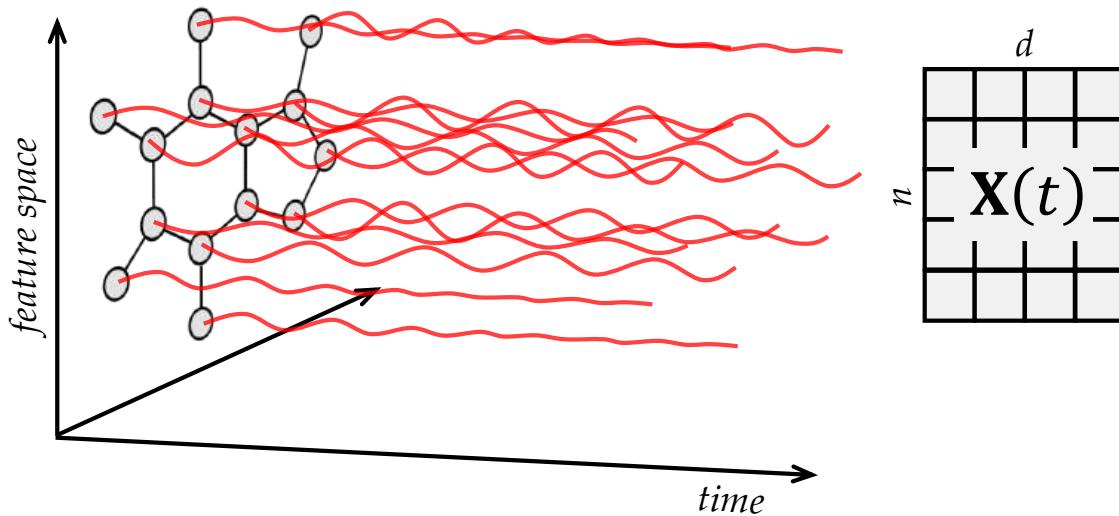
Continuous models for GNNs?



Physics-inspired GNNs

Physical metaphor of Graph ML

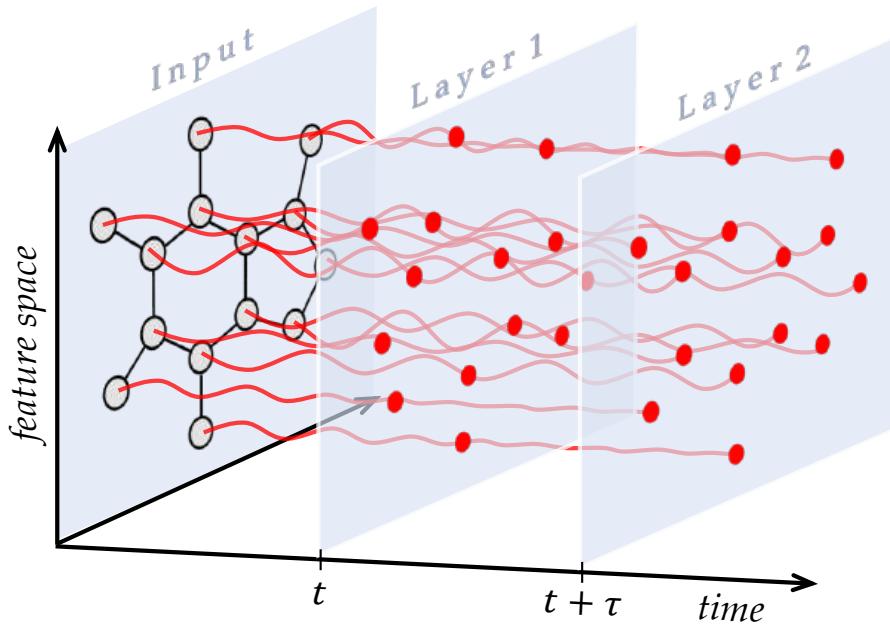
- GNN = dynamic system



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

Haber, Ruthotto 2017; Chen et al. 2019 (Neural ODEs); Xhonneau et al. 2020 (CGNN); Chamberlain et B. 2021 (GRAND, BLEND); Eliasof, Haber 2021 (PDE-GCN); Di Giovanni, Rowbottom et B 2022 (GRAFF), Rusch et B 2022 (GraphCON); Wu et B, Yan 2023; Eliasof et al. 2023; 2024 (TDE-GNN)

Physical metaphor of Graph ML



- GNN = dynamic system
- layers = discretisation of time
- graph = coupling function
(discretisation of space)

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

Haber, Ruthotto 2017; Chen et al. 2019 (Neural ODEs); Xhonneau et al. 2020 (CGNN); Chamberlain et B. 2021 (GRAND, BLEND); Eliasof, Haber 2021 (PDE-GCN); Di Giovanni, Rowbottom et B 2022 (GRAFF), Rusch et B 2022 (GraphCON); Wu et B, Yan 2023; Eliasof et al. 2023; 2024 (TDE-GNN)

Heat Diffusion

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

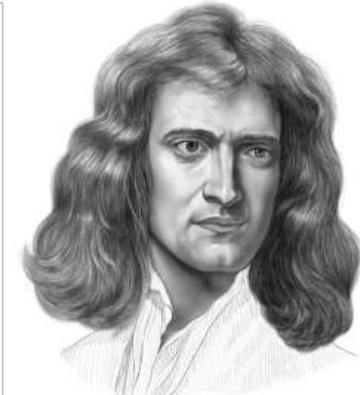
(824)
with a little pressing, I took a drop thereof, and in it discover'd a mighty number of living Creatures. I repeated my observation the same evening with the same success; but the next day I could find none of them alive; and whereas I had laid that drop upon a small Copper Plate, I fancied to my self that the exhalation of the moisture might be the cause of their death, and not the cold weather, which at that time was very moderate.

In the beginning of April I took the Male seed of a Jack or Pike, but could discover nothing more than in that of a Cod-fish, but having added about four times as much Water in quantity as the matter it self was, and then making my remarks, I could perceive that the *Animalcula* did not only wax stronger and swifter, but, to my great amazement, I saw them move with that celerity, that I could compare it to nothing more than what we have seen with our naked Eye, a River Fish chased by its powerful Enemy, which is just ready to devour it: You must observe that this whole Course was not longer than the Diameter of a single Hair of ones Head.

VII. *Scala graduum Caloris.*

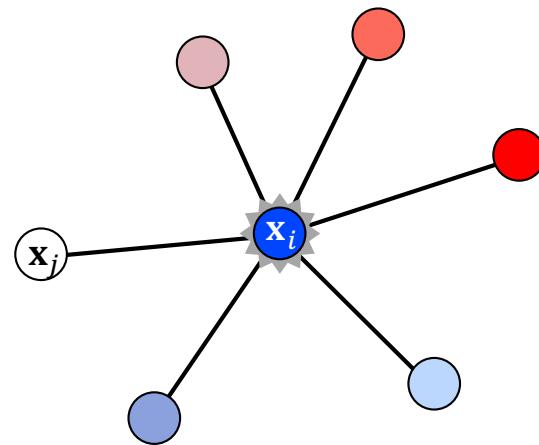
Calorum Descriptiones & signa.

0	Calor aeris hiberni ubi aqua incipit gelu rigescere. Innotescit his calor accurate locando Thermometrum in nive compreßa quo tempore gelu solvitur.
0,1,2,	Calores aeris hiberni.
2,3,4,	Calores aeris verni & autumnalis.
4,5,6,	Calores aeris aestivi.
6	Calor aeris meridiani circa mensem Iulium.
12	Calor maximus quem Thermometer ad contactum



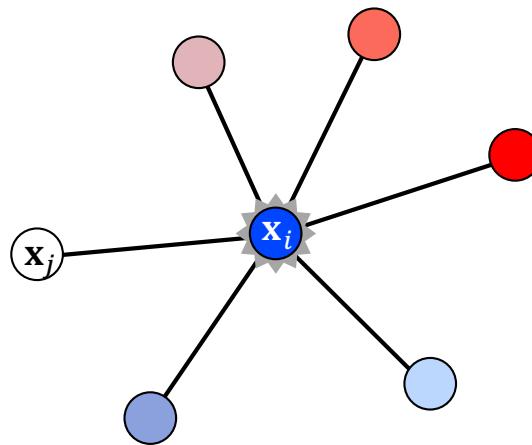
I. Newton

Heat Diffusion Equation on Graphs



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

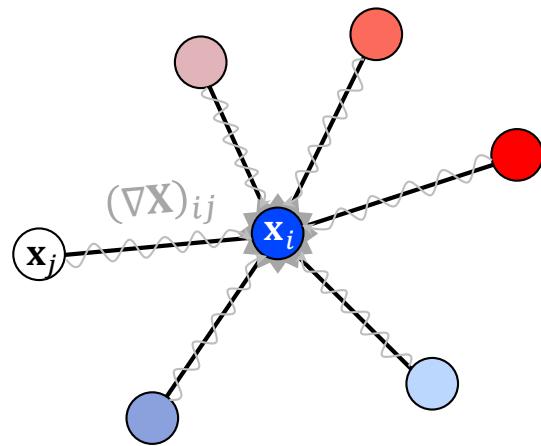
Heat Diffusion Equation on Graphs



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

rate of temperature change self temperature temperature of the environment

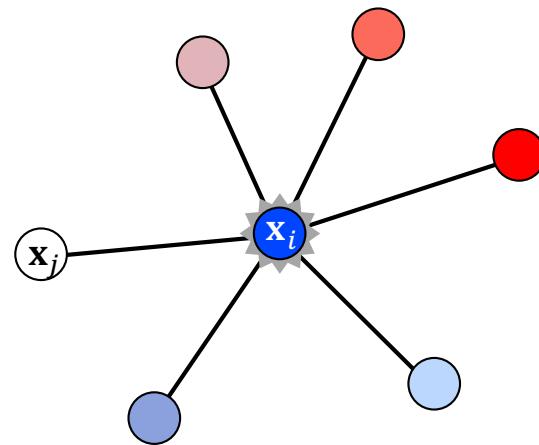
Heat Diffusion Equation on Graphs



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

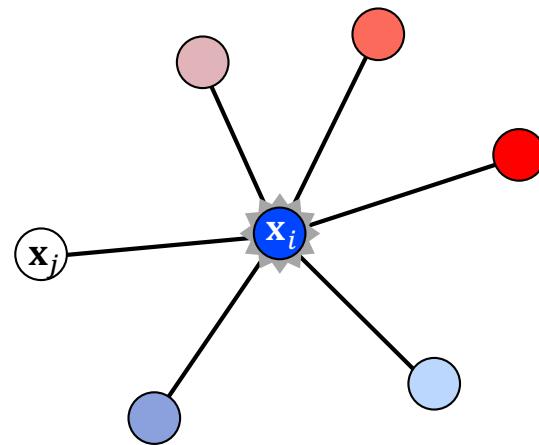
The term $\sum_{j \in \mathcal{N}_i} a_{ij} (\mathbf{x}_i(t) - \mathbf{x}_j(t))$ is shown as a sum of two parts. The first part, $\sum_{j \in \mathcal{N}_i} a_{ij} \mathbf{x}_i(t)$, is labeled "divergence" and "div". The second part, $-\sum_{j \in \mathcal{N}_i} a_{ij} \mathbf{x}_j(t)$, is labeled "gradient" and $-(\nabla \mathbf{X})_{ij}$.

Heat Diffusion Equation on Graphs



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

Heat Diffusion Equation on Graphs



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

Heat Diffusion Equation as a prototypical Gradient Flow

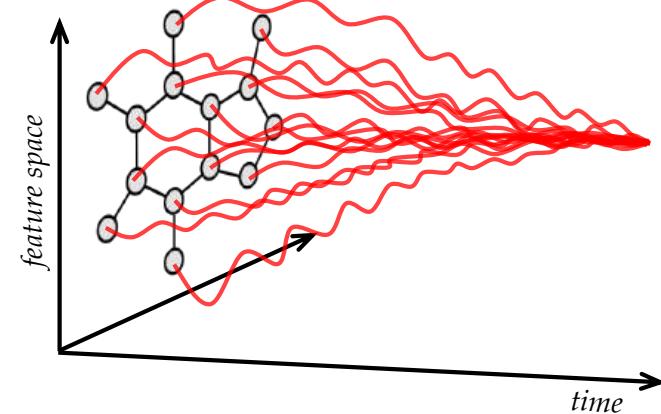
$$\dot{\mathbf{X}}(t) = -\nabla \mathcal{E}(\mathbf{X}(t))$$

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



G. Dirichlet

- Heat equation is the gradient flow of the Dirichlet energy
- “Smoothness” of the node features
- Dirichlet energy **decreases along the flow**
- In the limit $t \rightarrow \infty$ results in “oversmoothing”



Heat Diffusion Equation as a prototypical Gradient Flow

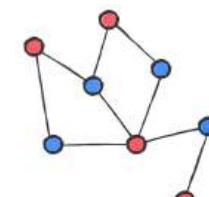
$$\dot{\mathbf{X}}(t) = -\nabla \mathcal{E}(\mathbf{X}(t))$$

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

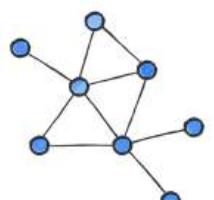


G. Dirichlet

- Heat equation is the gradient flow of the Dirichlet energy
- “Smoothness” of the node features
- Dirichlet energy **decreases along the flow**
- In the limit $t \rightarrow \infty$ results in “oversmoothing”
- Not very expressive: works only in homophilic graphs (“similar neighbours”)



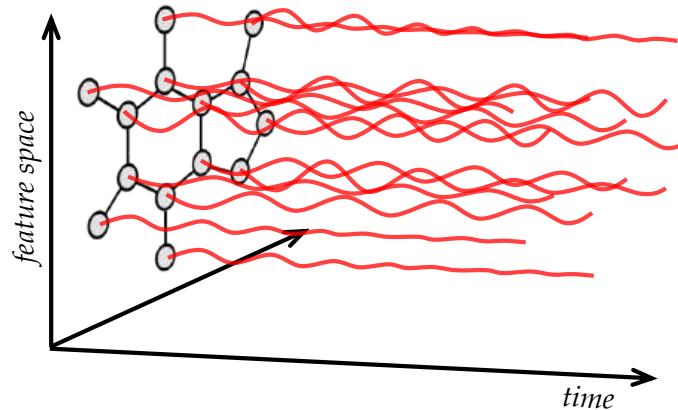
heterophilic



homophilic

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

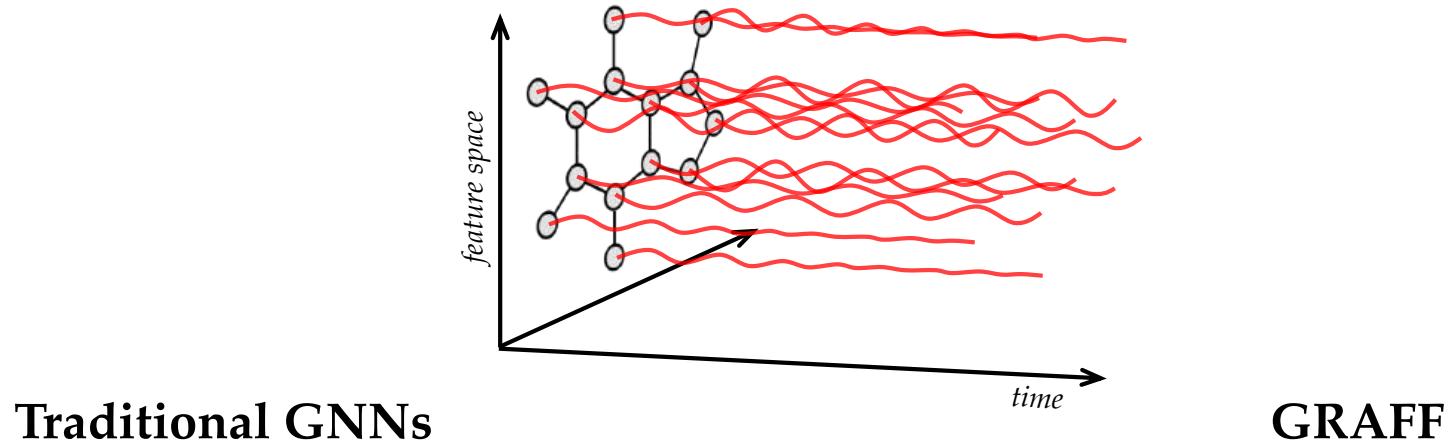
Gradient Flow Framework (GRAFF)



Traditional GNNs

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

Gradient Flow Framework (GRAFF)



Traditional GNNs

$$\mathbf{X}(k+1) = \mathbf{X}(k) + \tau \mathbf{F}_{\theta(k)}(\mathbf{X}(k), \mathcal{G})$$

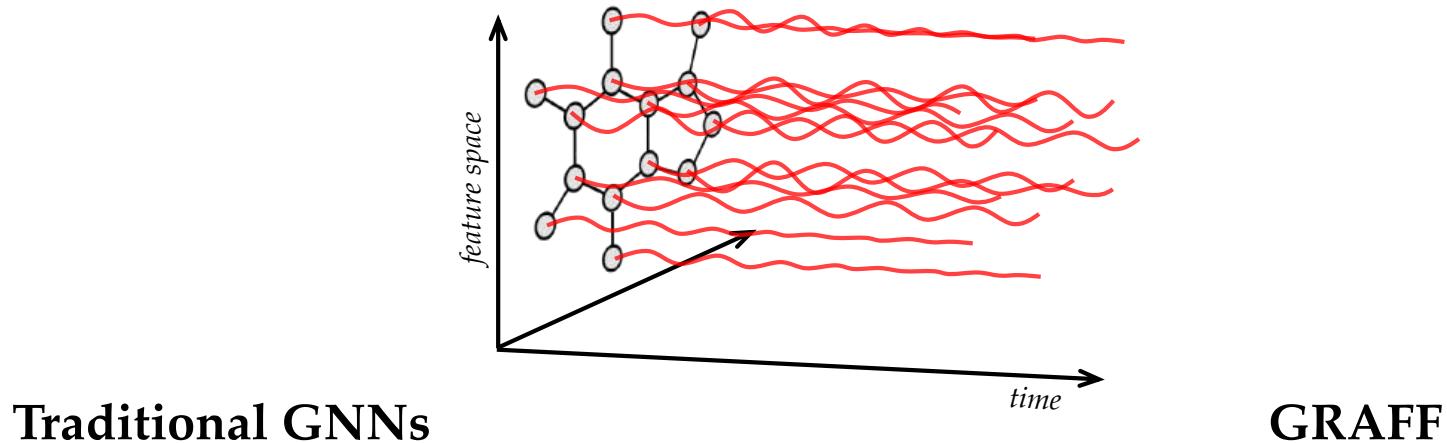
- Parametrize **evolution equations**

GRAFF

$$\mathcal{E}_{\theta(t)}(\mathbf{X}(t), \mathcal{G})$$

- Parametrize **energy**

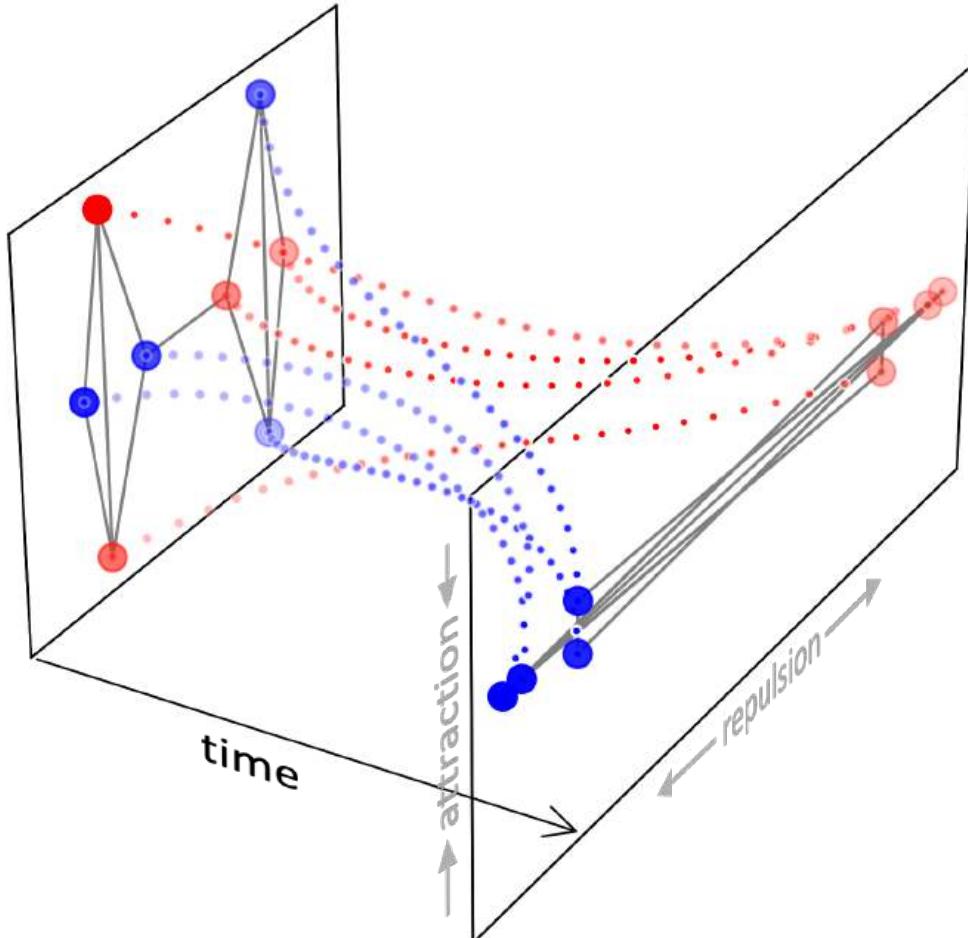
Gradient Flow Framework (GRAFF)



$$\mathbf{X}(k+1) = \mathbf{X}(k) + \tau \mathbf{F}_{\theta(k)}(\mathbf{X}(k), \mathcal{G})$$

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

- Parametrize **evolution equations**
- Parametrize **energy**
- Derive evolution equation as GF
- Better “interpretability”



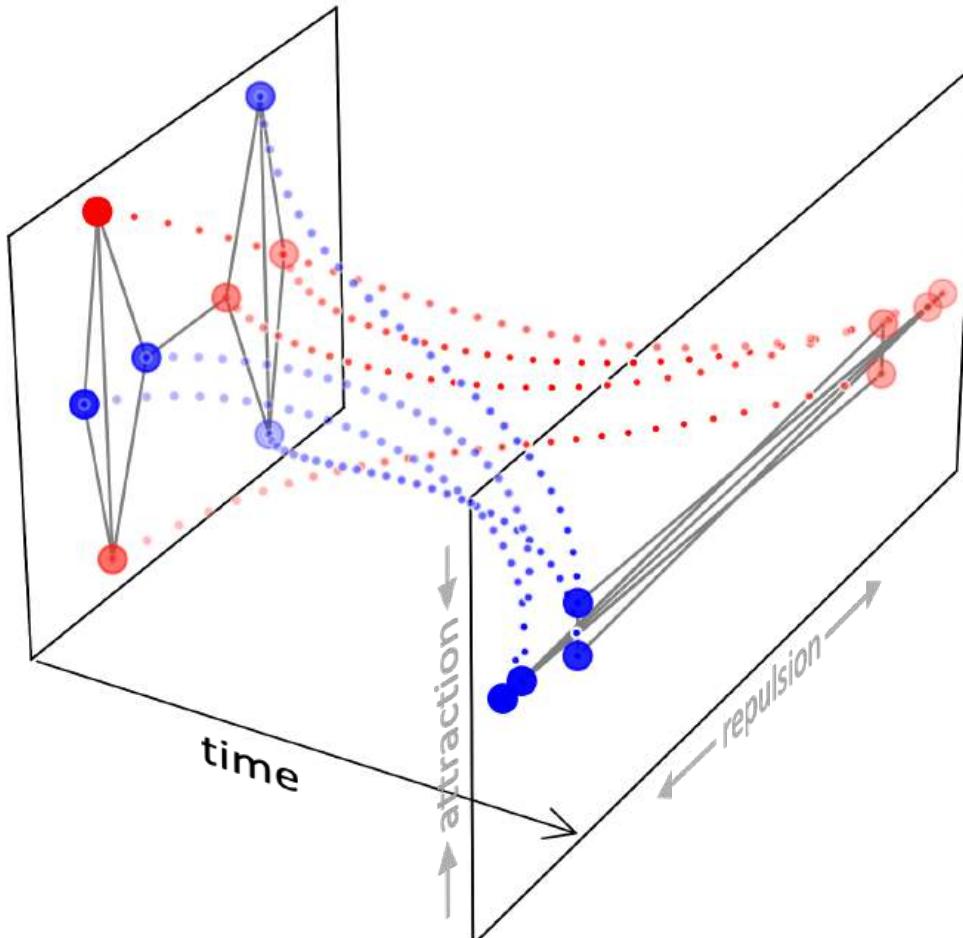
Di Giovanni, Rowbottom et B 2022

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

- **Attraction** along positive eigenvectors of \mathbf{W}
- **Repulsion** along negative eigenvectors of \mathbf{W}

$$\dot{\mathbf{X}}(t) = \bar{\mathbf{A}}\mathbf{X}(t)\mathbf{W}$$

Theorem: Linear graph diffusion (“convolutional GNN”) with appropriately designed channel mixing matrix \mathbf{W} (symmetric & with sufficiently large negative eigenvalues) can provably avoid oversmoothing.



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

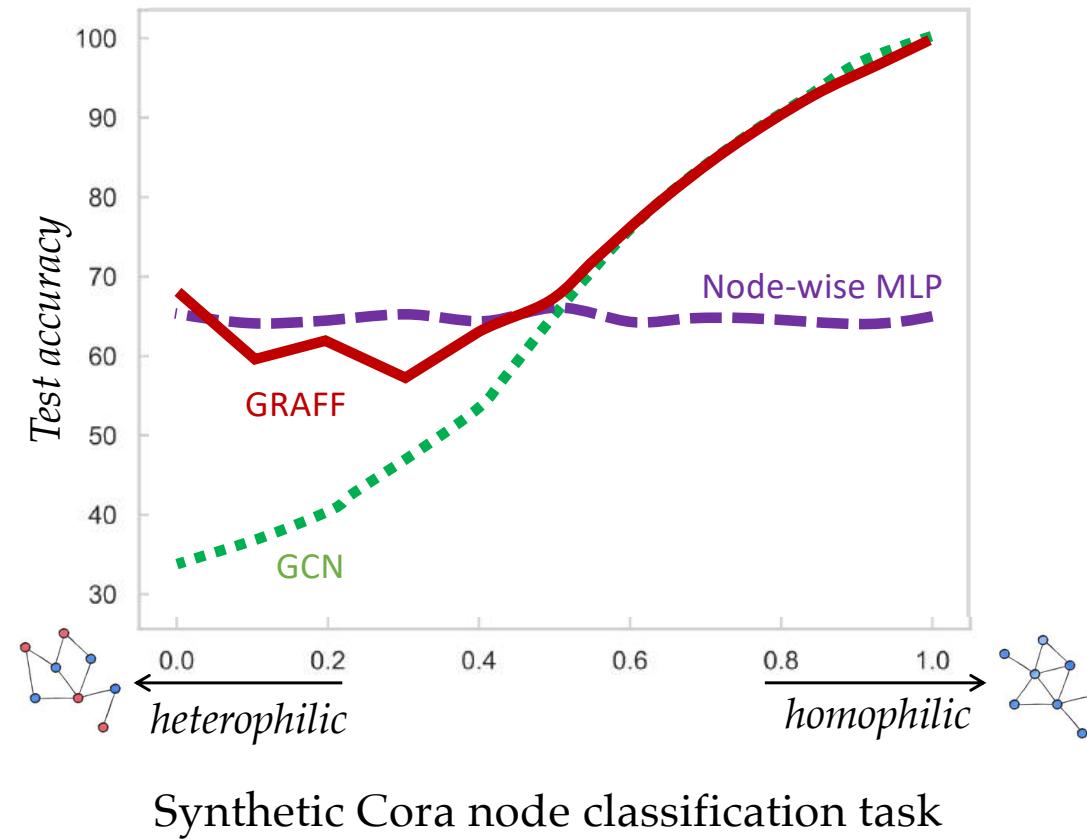
- **Attraction** along positive eigenvectors of \mathbf{W}
- **Repulsion** along negative eigenvectors of \mathbf{W}

$$\dot{\mathbf{X}}(t) = \bar{\mathbf{A}}\mathbf{X}(t)\mathbf{W}$$

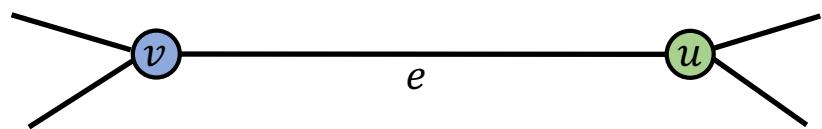
Theorem: Linear graph diffusion (“convolutional GNN”) with appropriately designed channel mixing matrix \mathbf{W} can avoid oversmoothing.

Contradicts GNN “folklore”!

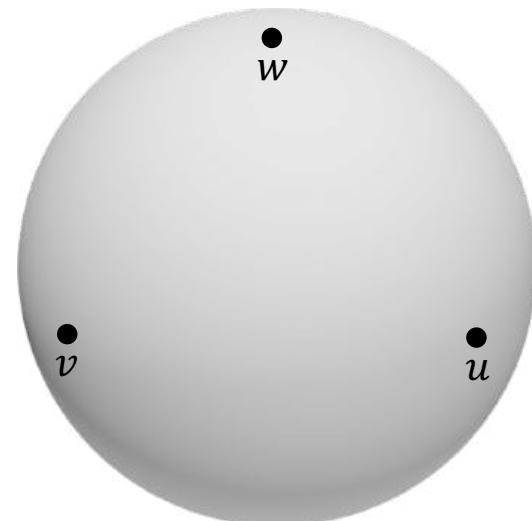
Homophily vs Heterophily



Cellular Sheaves

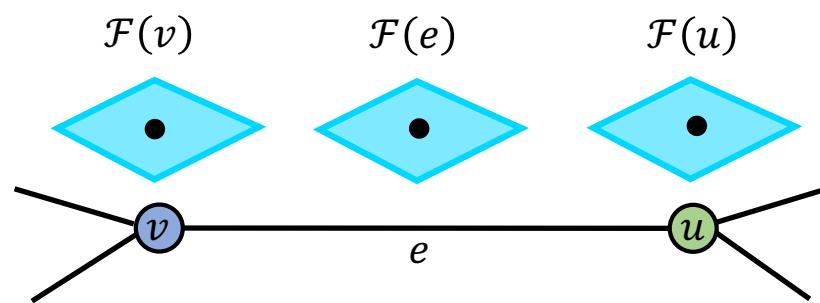


Graph

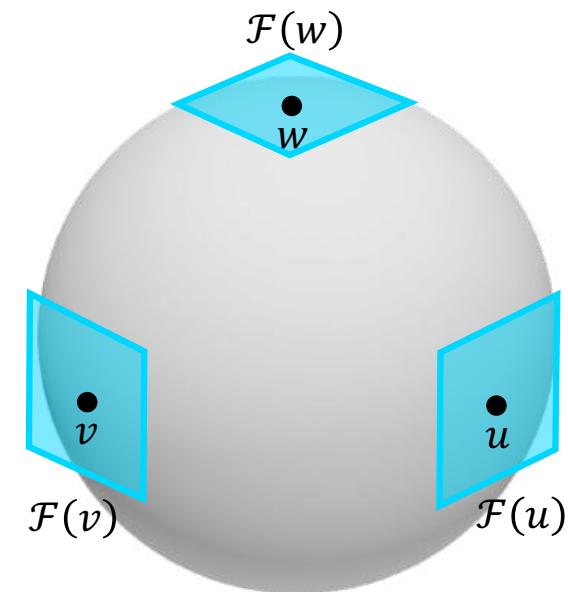


Manifold

Cellular Sheaves

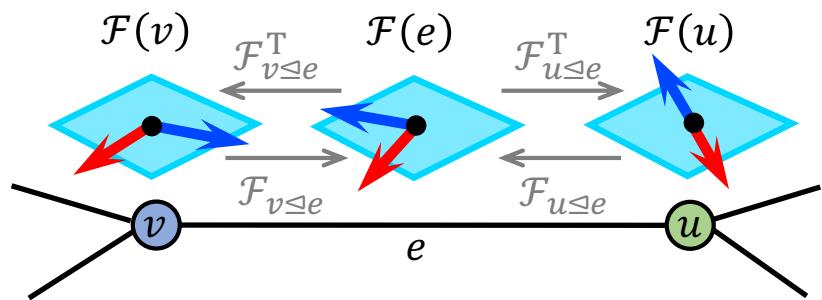


Cellular sheaf \mathcal{F}

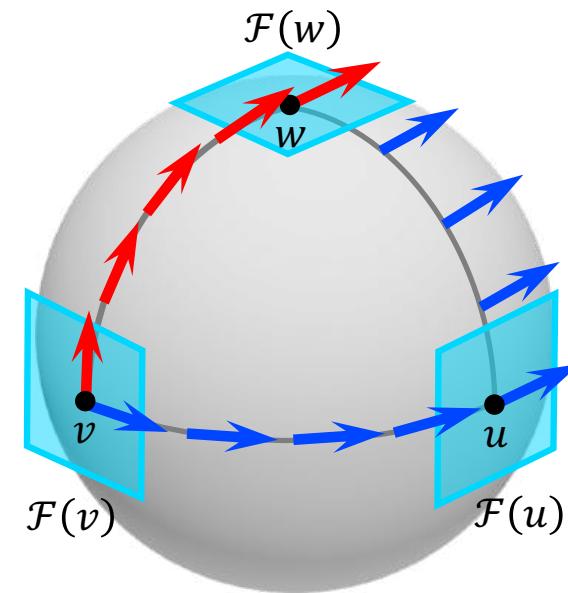


Manifold + Connection

Cellular Sheaves

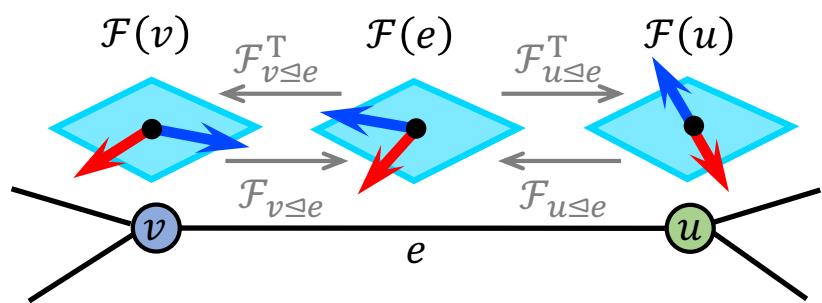


Cellular sheaf \mathcal{F}

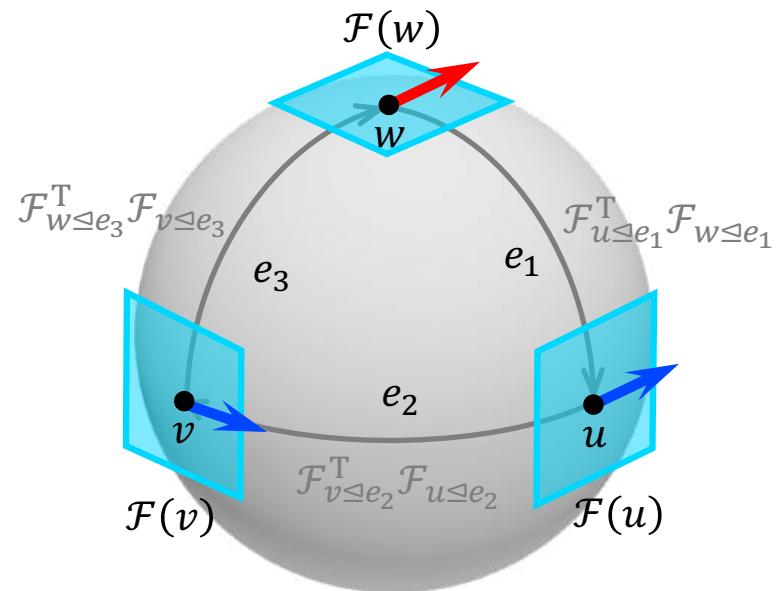


Manifold + Connection

Cellular Sheaves

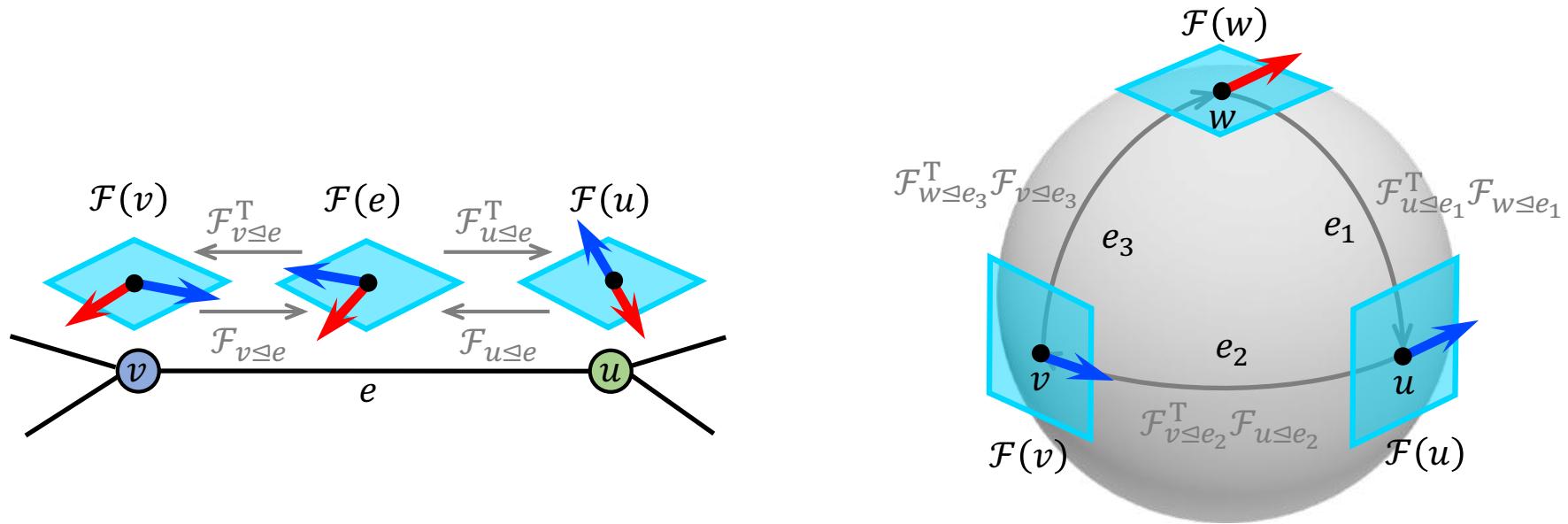


Cellular sheaf \mathcal{F}



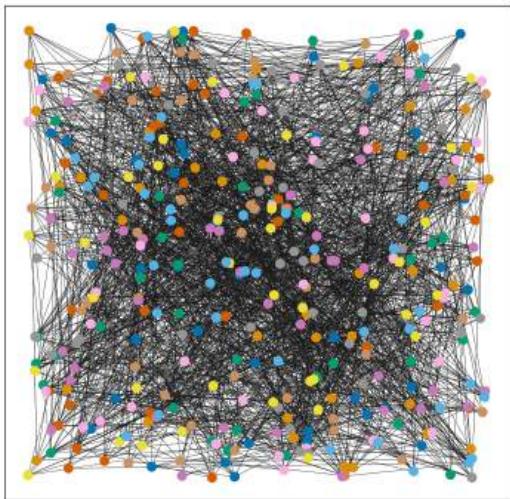
Analogy to parallel transport
on manifolds

Cellular Sheaves



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

Diffusion on Cellular Sheaves



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

$$\dot{\mathbf{X}}(t) = \Delta_{\mathcal{F}} \mathbf{X}(t)$$

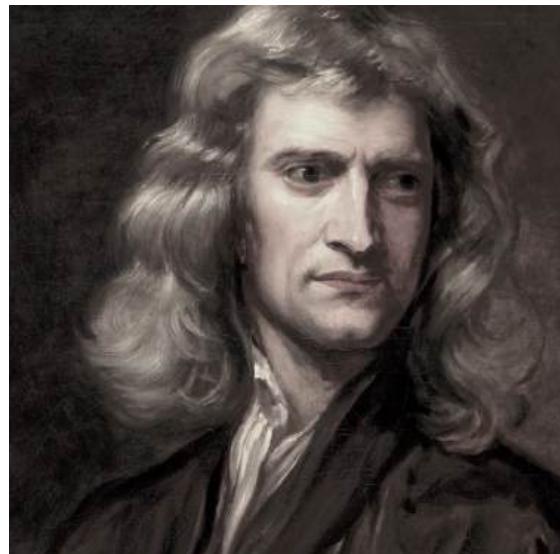
Node classification = limit of sheaf diffusion equation
with an appropriate sheaf class



Michael Galkin

Homogeneous Diffusion in Image Processing

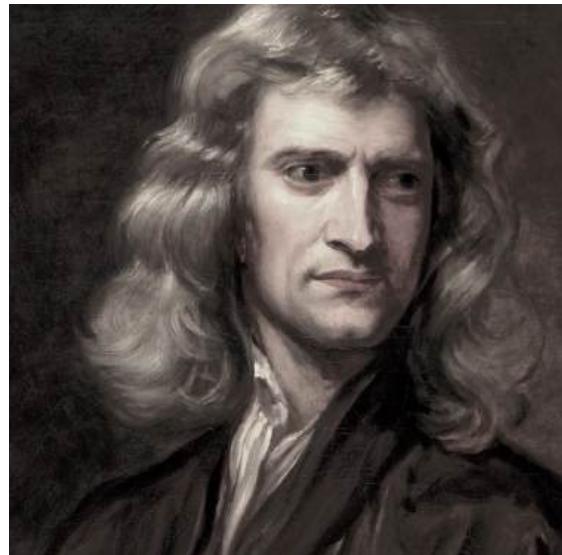
$$\dot{\mathbf{X}}(t) = -\operatorname{div}(c \nabla \mathbf{X}(t))$$



$\mathbf{X}(0)$

Homogeneous Diffusion in Image Processing

$$\dot{\mathbf{X}}(t) = -\operatorname{div}(c \nabla \mathbf{X}(t))$$



$\mathbf{X}(0)$

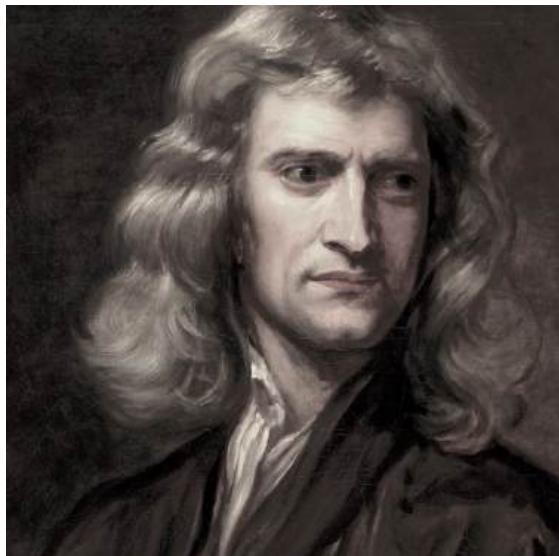


$\mathbf{X}(t) = \mathbf{X}(0) \star \mathbf{G}_{\sigma \propto t}$

Non-homogeneous Diffusion in Image Processing

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

edge indicator



$\mathbf{X}(0)$

Perona, Malik 1990



“Do not diffuse across edges”

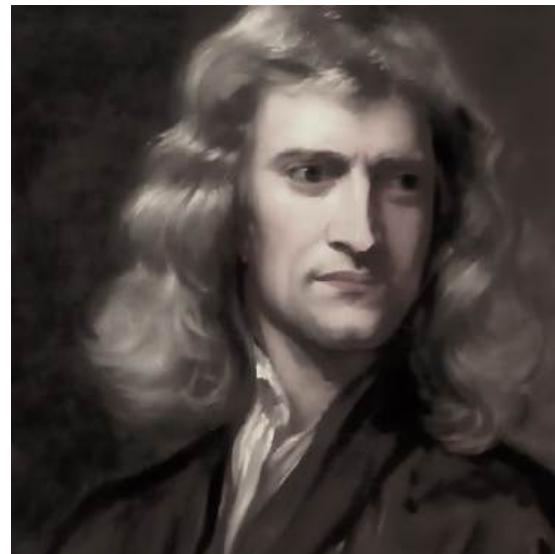
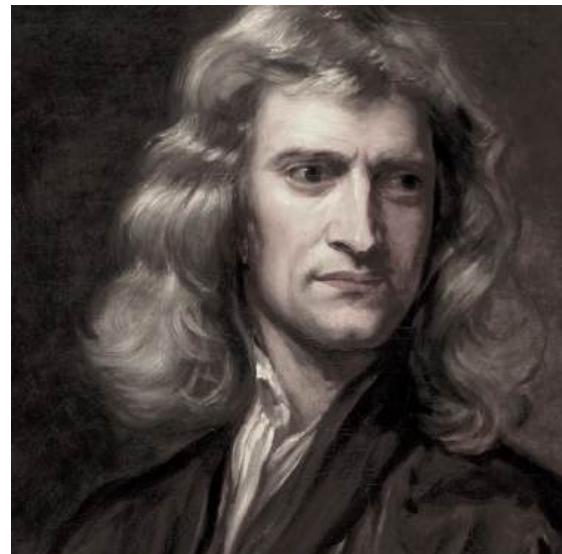


$\mathbf{X}(t) = \mathbf{X}(0) * \mathbf{G}_{\sigma \propto t}$

Non-homogeneous Diffusion in Image Processing

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

edge indicator



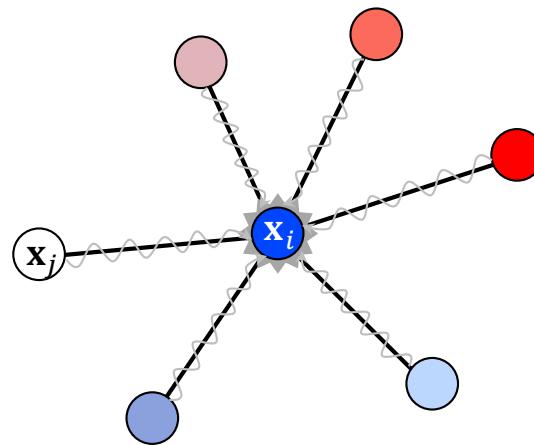
Non-homogeneous

Perona, Malik 1990



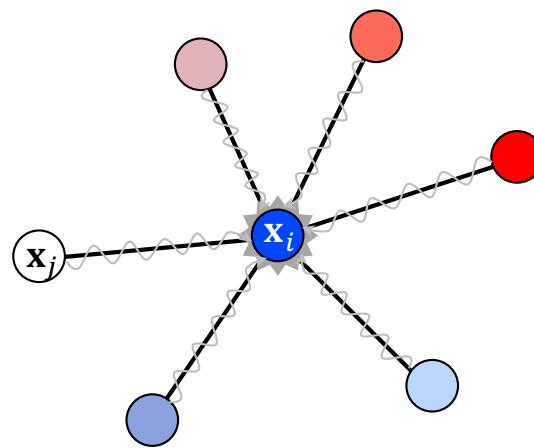
Homogeneous

Non-homogeneous Diffusion on Graphs



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

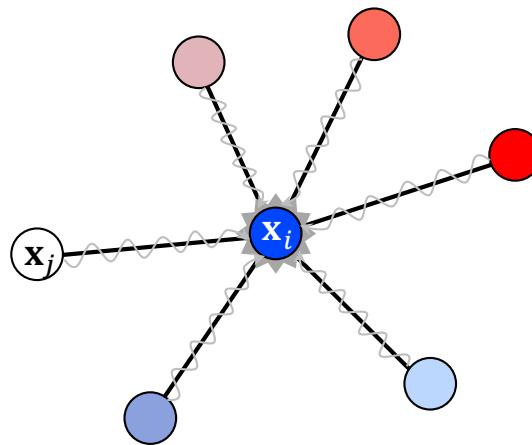
Non-homogeneous Diffusion on Graphs



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

learnable diffusivity

Non-homogeneous Diffusion on Graphs

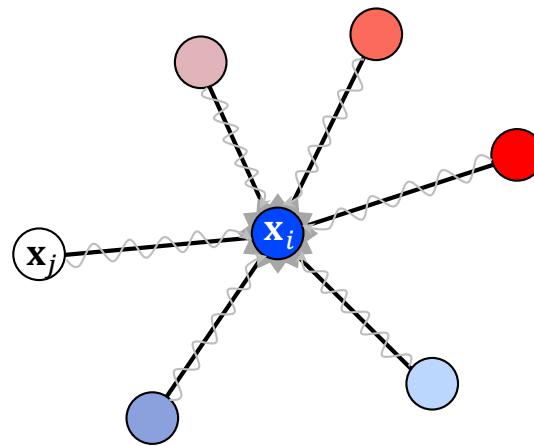


$$\mathbf{x}_i(t + \tau) = \mathbf{x}_i(t) + \tau \sum_{\substack{\text{time} \\ \text{step}}} \sum_{j \in \mathcal{N}_i} a(\mathbf{x}_i(t), \mathbf{x}_j(t)) (\mathbf{x}_j(t) - \mathbf{x}_i(t))$$

Explicit (Forward Euler)
discretization

learnable diffusivity

Non-homogeneous Diffusion on Graphs



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

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

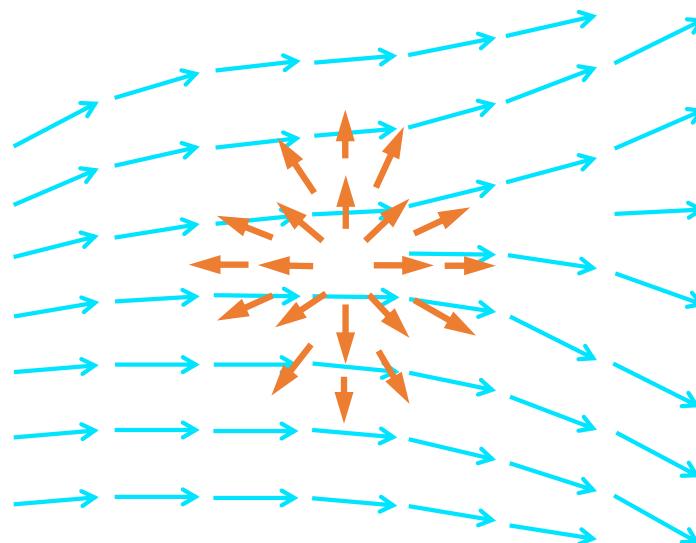
GAT!

Advective Diffusion Transformer

$$\dot{\mathbf{X}}(t) = \operatorname{div}(\mathbf{S}(t)\nabla\mathbf{X}(t)) + \beta\operatorname{div}(\mathbf{V}(t)\mathbf{X}(t))$$

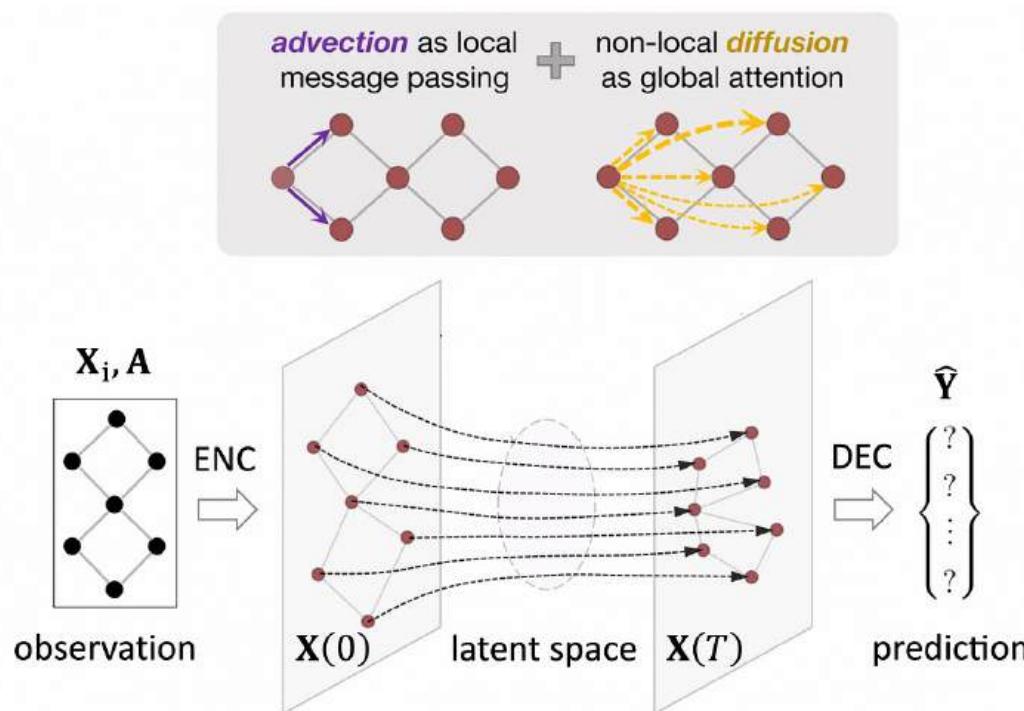
Advective Diffusion Transformer

$$\dot{\mathbf{X}}(t) = \underset{\text{diffusion}}{\operatorname{div}(\mathbf{S}(t)\nabla\mathbf{X}(t))} + \underset{\text{advection}}{\beta\operatorname{div}(\mathbf{V}(t)\mathbf{X}(t))}$$



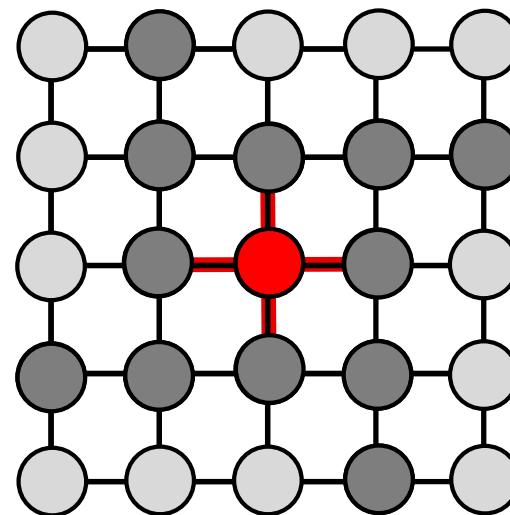
Advective Diffusion Transformer

$$\dot{\mathbf{X}}(t) = \operatorname{div}(\mathbf{S}(t)\nabla\mathbf{X}(t)) + \beta\operatorname{div}(\mathbf{V}(t)\mathbf{X}(t))$$



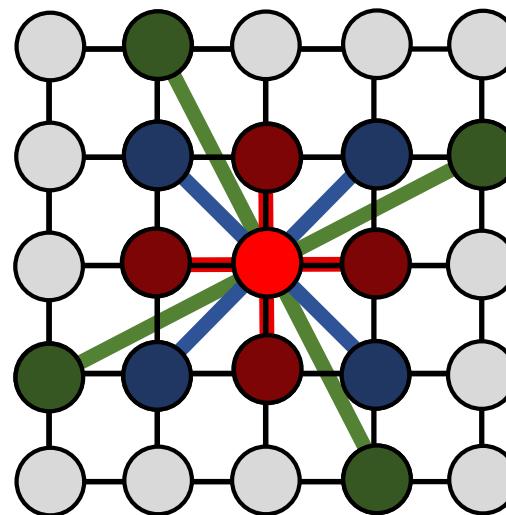
Spatial Derivative: Graph Rewiring?

$$\dot{x}(t) = \Delta x(t)$$



Spatial Derivative: Graph Rewiring?

$$\dot{x}(t) = \Delta x(t)$$



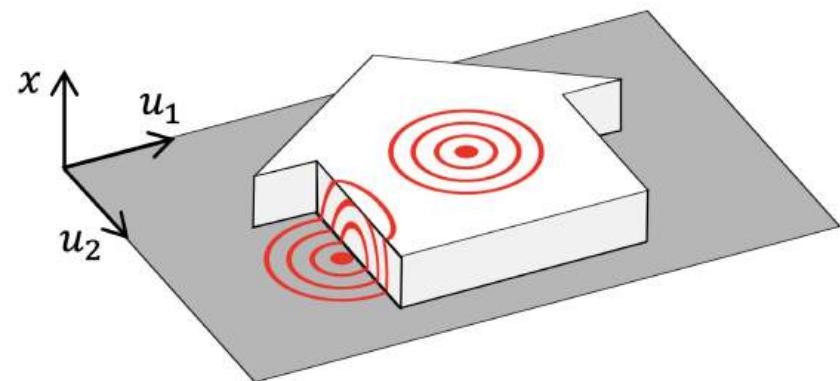
Different discretisations of 2D Laplacian

Images as embedded manifolds



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

Non-linear diffusion



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

Non-Euclidean diffusion

Kimmel et al. 1997; Sochen et al. 1998

Beltrami flow

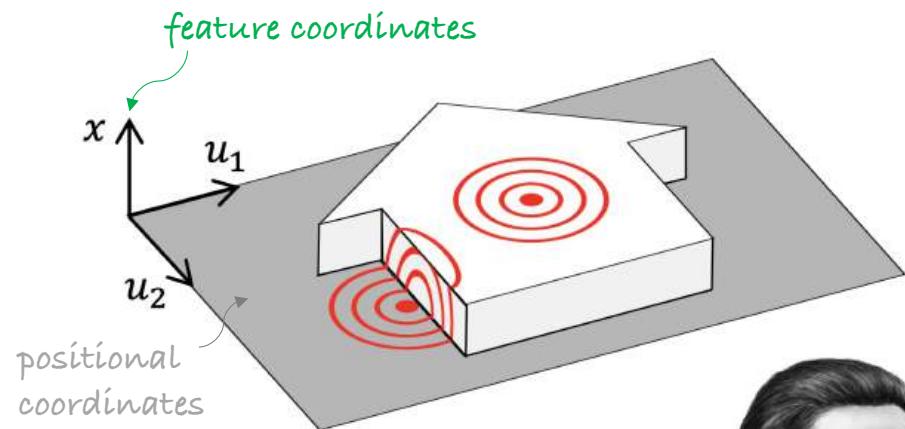
- Consider image as embedded 2-manifold

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

- 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)



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



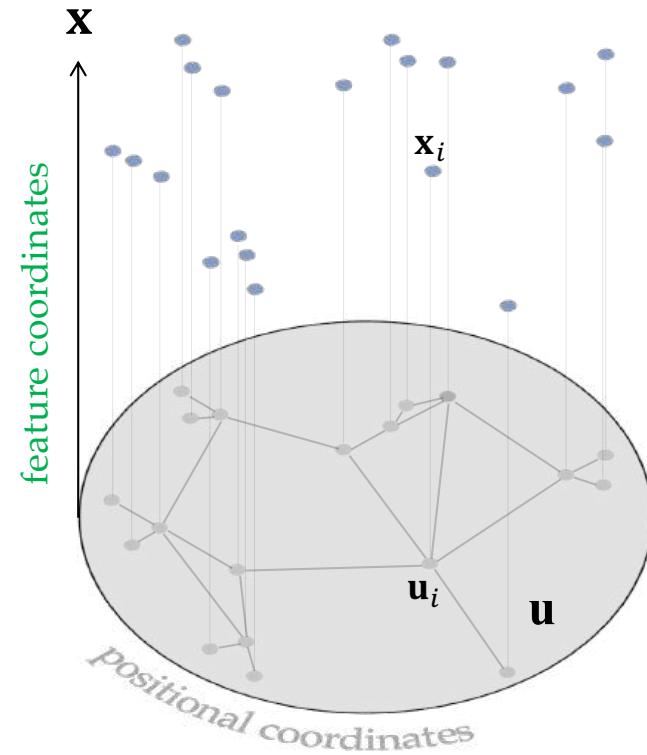
E. Beltrami

Kimmel et al. 1997; Sochen et al. 1998

Graph Beltrami flow

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

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

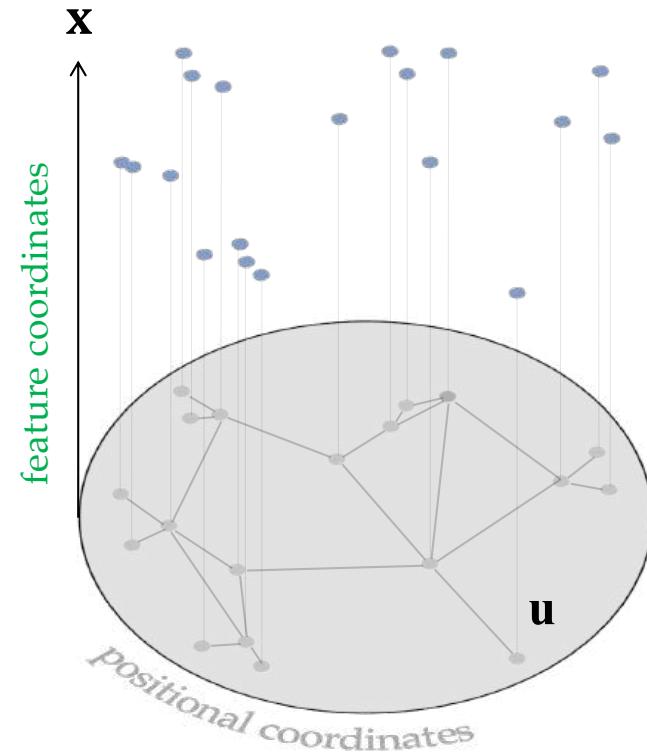


Graph Beltrami flow

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

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

- Evolution of \mathbf{x} = feature diffusion

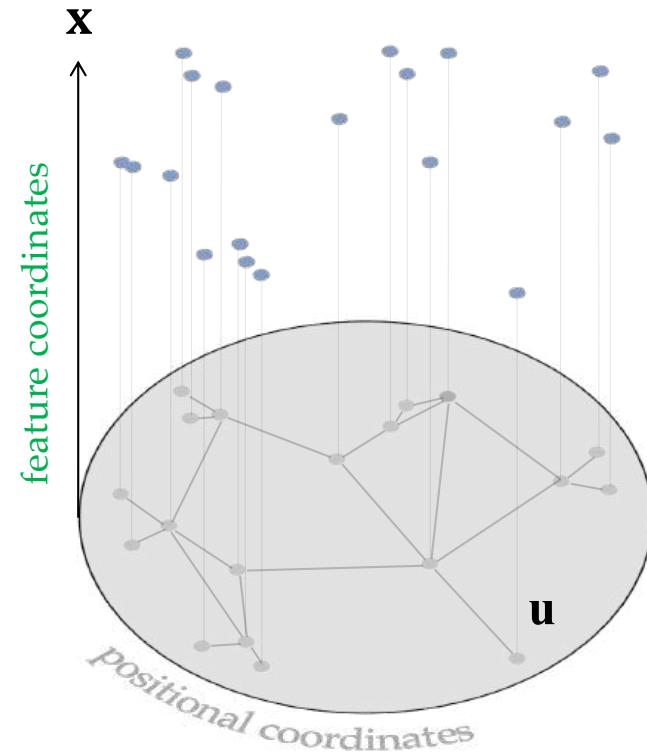


Graph Beltrami flow

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

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

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



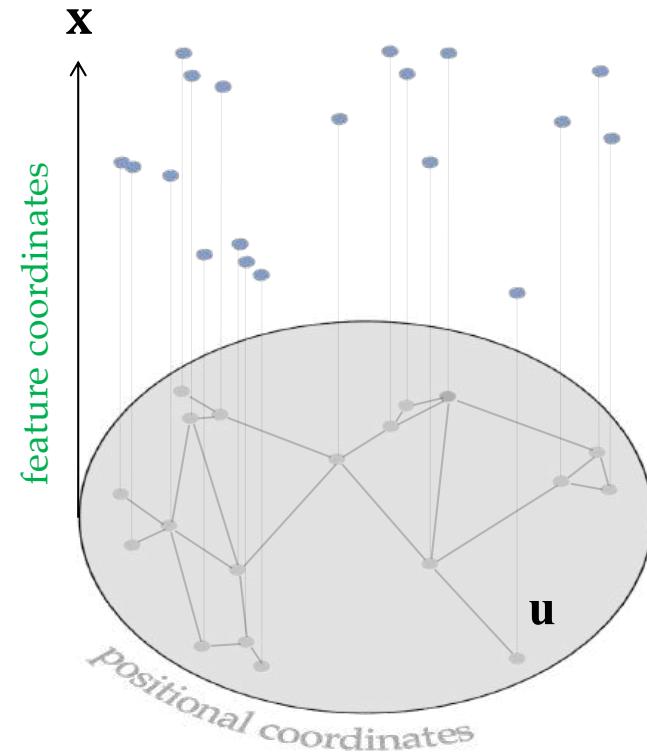
Graph Beltrami flow

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

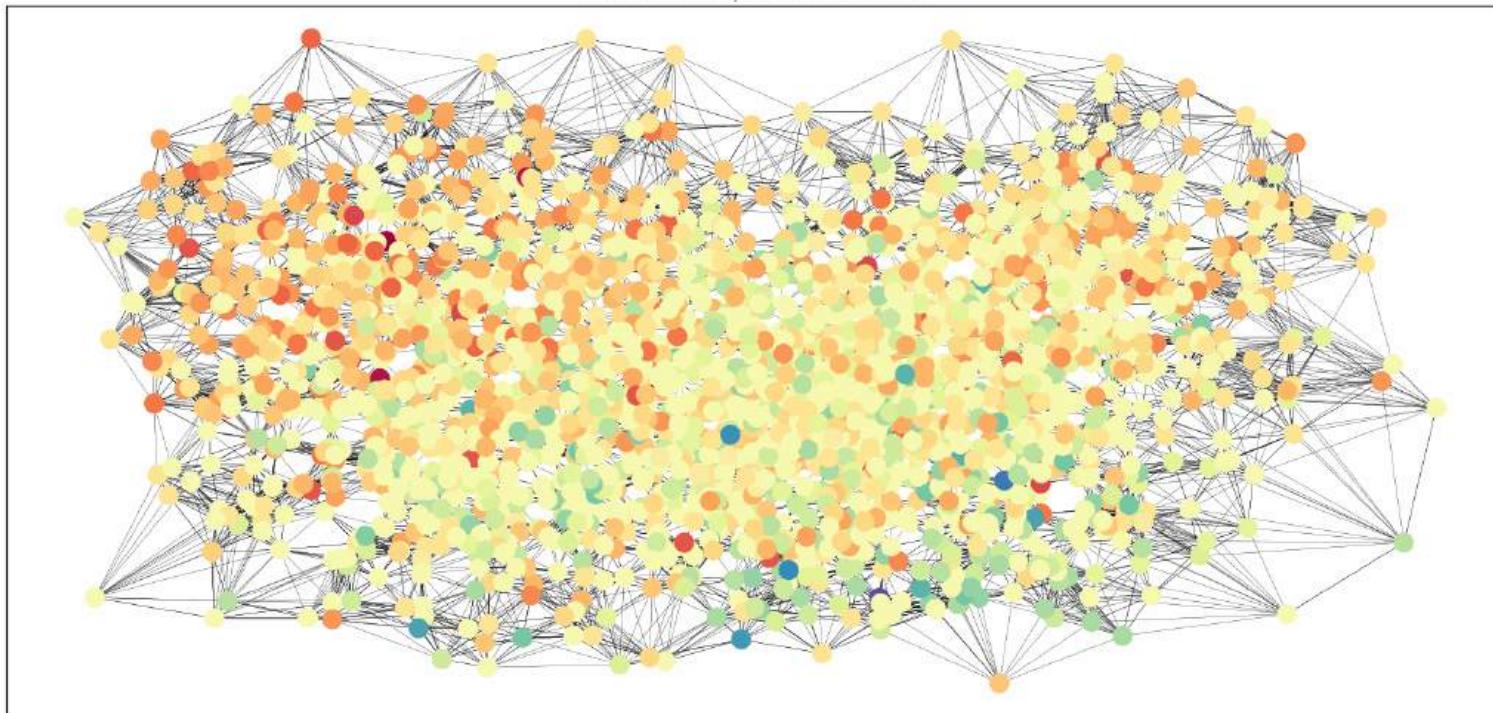
$$\dot{\mathbf{z}}_i(t) = \sum_{j \in \mathcal{N}'_i} a(\mathbf{z}_i(t), \mathbf{z}_j(t)) (\mathbf{z}_j(t) - \mathbf{z}_i(t))$$

rewired graph

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

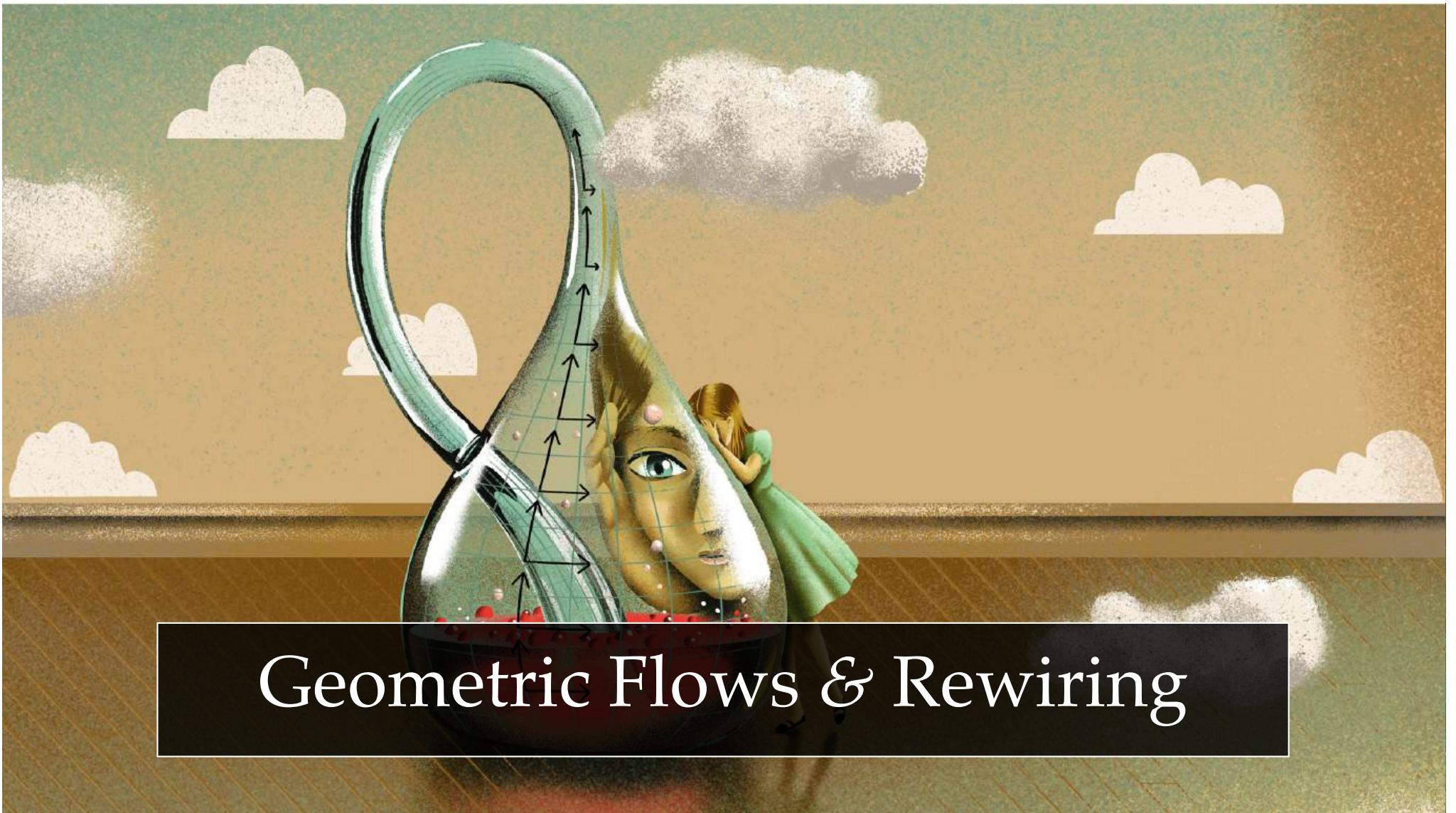


Graph Beltrami flow



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

Chamberlain, Rowbottom, et B. 2021

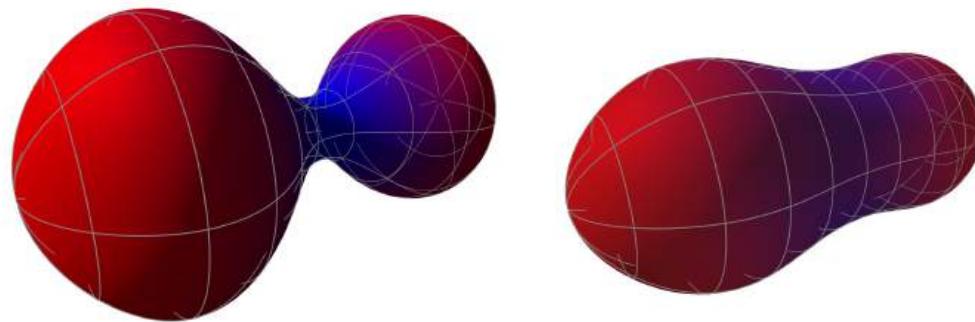


Geometric Flows & Rewiring

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;



**G. Ricci-
Curbastro**

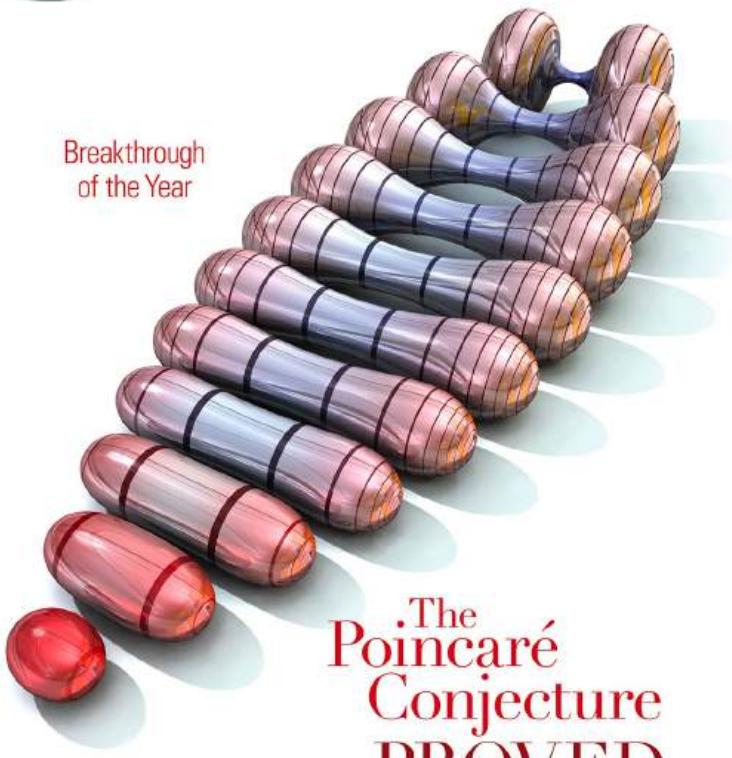


R. Hamilton

Science

22 December 2006 | \$10

Breakthrough
of the Year



The
**Poincaré
Conjecture
PROVED**

Ricci 1903; Hamilton 1988; Perelman 2003



G. Perelman



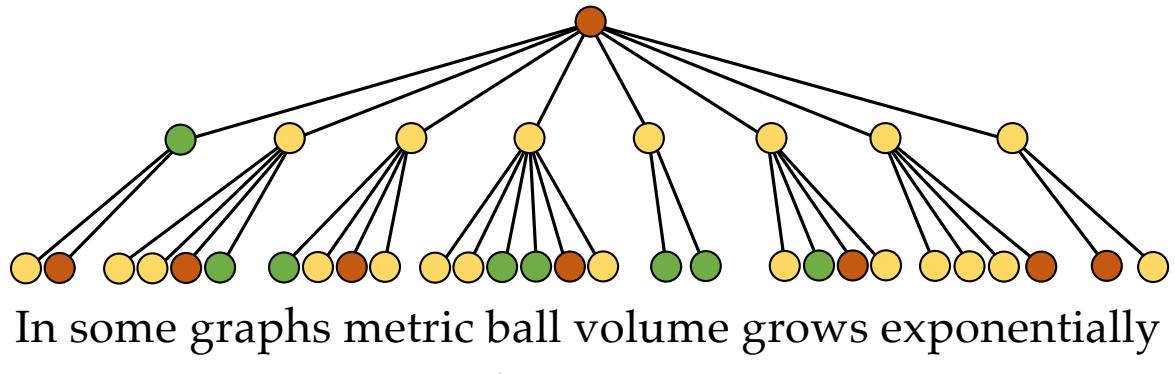
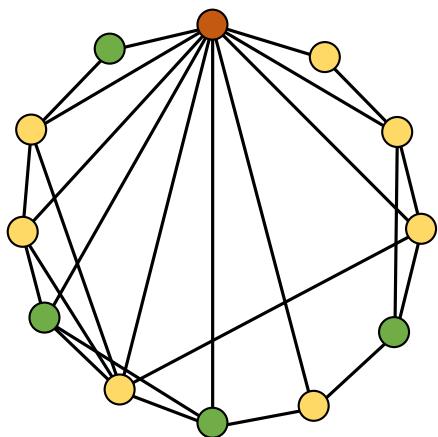
G. Ricci-
Curbastro



R. Hamilton

**“Failure of Message Passing to propagate
information on the graph”**

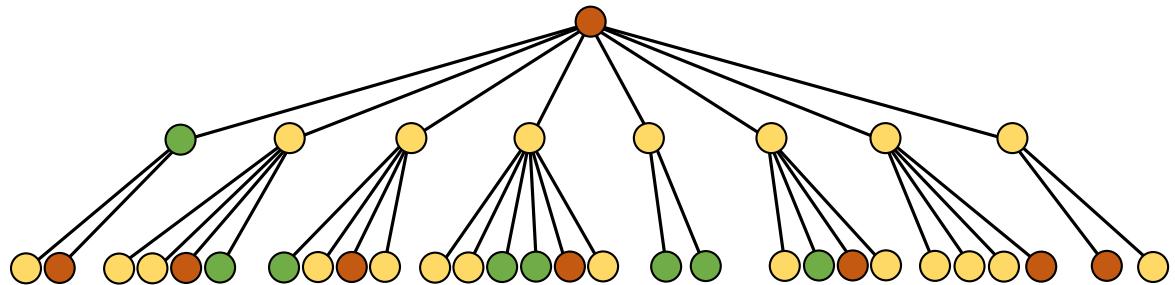
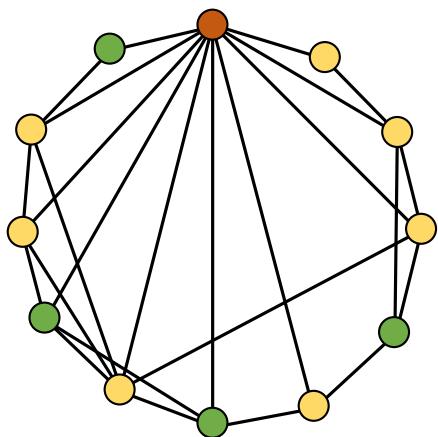
Over-squashing & Bottlenecks



In some graphs metric ball volume grows exponentially with ball radius

Over-squashing = Fast volume growth
+ Long-range interactions

Over-squashing & Bottlenecks



In some graphs metric ball volume grows exponentially with ball radius

graph topology

Over-squashing = **Fast volume growth**
+ **Long-range interactions**

task

Over-squashing

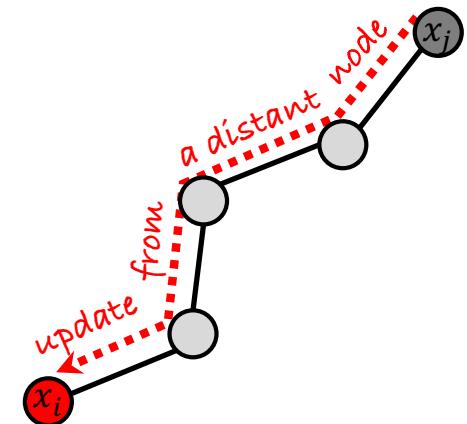
- Consider an MPNN of the form

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

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

Theorem (Sensitivity bound): For any i, j

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



Over-squashing: small Jacobian

$\left\| \frac{\partial \mathbf{x}_i^{(L)}}{\partial \mathbf{x}_j^{(0)}} \right\|$ indicates poor information propagation from input node

Over-squashing

- Consider an MPNN of the form

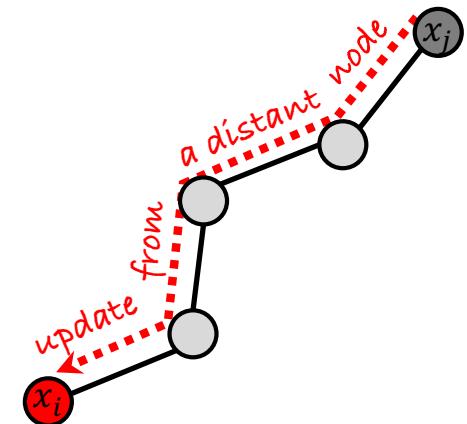
$$\mathbf{x}_i^{(k+1)} = \sigma \left(\mathbf{W}_1 \mathbf{x}_i^{(k)} + \sum_j a_{ij} \mathbf{W}_2 \mathbf{x}_j^{(k)} \right)$$

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

Theorem (Sensitivity bound): For any i, j

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

model topology

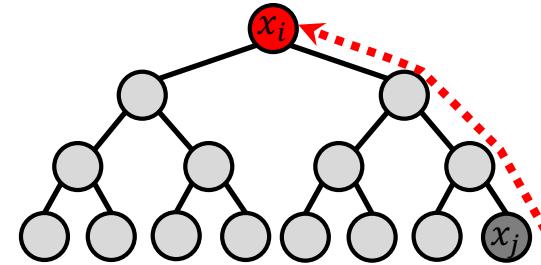
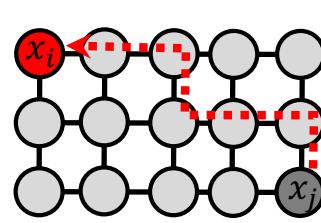


Over-squashing: small Jacobian

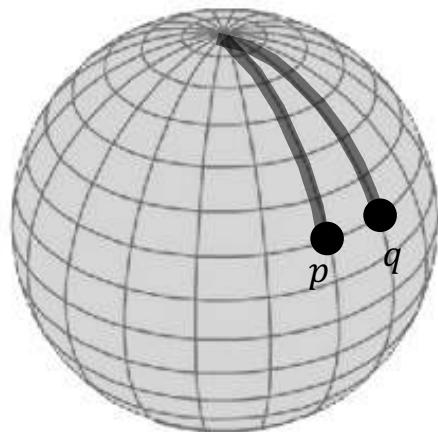
$\left\| \frac{\partial \mathbf{x}_i^{(L)}}{\partial \mathbf{x}_j^{(0)}} \right\|$ indicates poor information propagation from input node

Preventing over-squashing

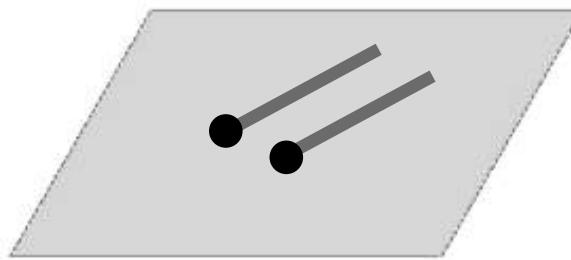
$$\left\| \frac{\partial \mathbf{x}_i^{(L)}}{\partial \mathbf{x}_j^{(0)}} \right\|_1 \leq (\text{model } \mathbf{c}_{\sigma} w p)^{\textcolor{green}{L}} (\mathbf{I} + \mathbf{A})_{ij}^{\textcolor{green}{L}}$$



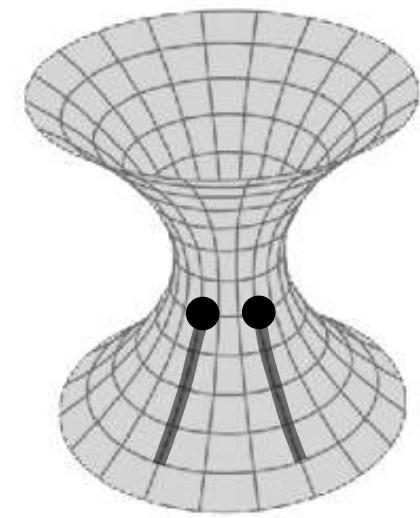
Ricci Curvature on Manifolds



Spherical (>0)



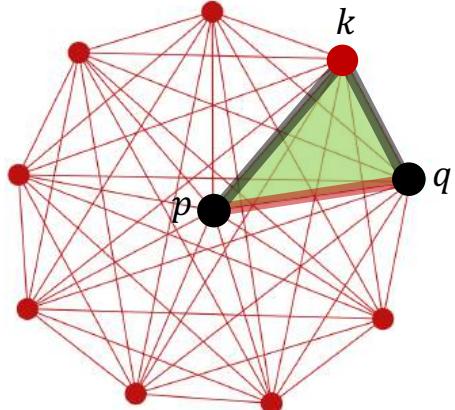
Euclidean ($=0$)



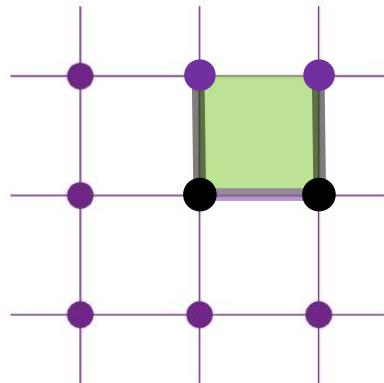
Hyperbolic (<0)

“geodesic dispersion”

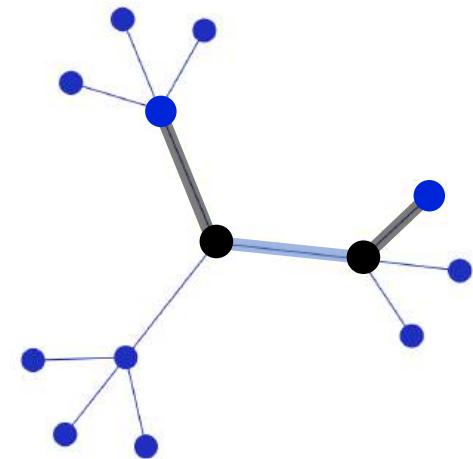
Discrete Ricci Curvature on Graphs



Clique (>0)

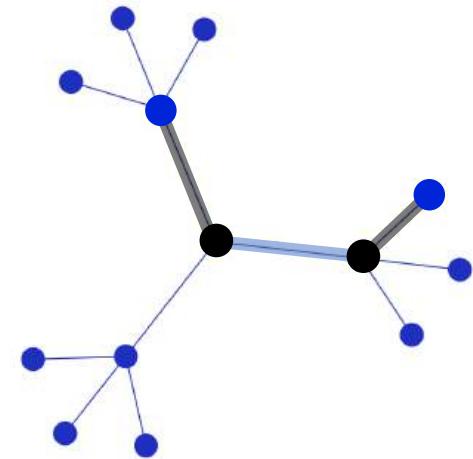
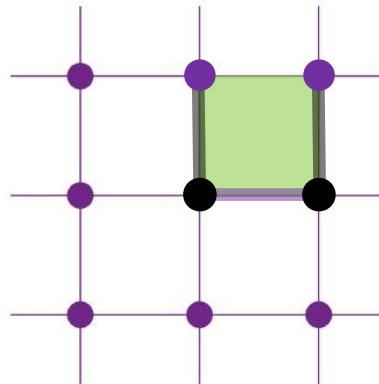
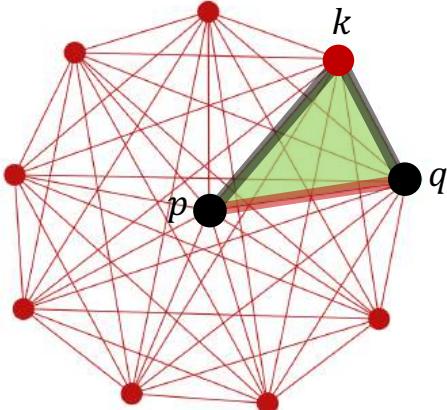


Grid ($=0$)



Tree (<0)

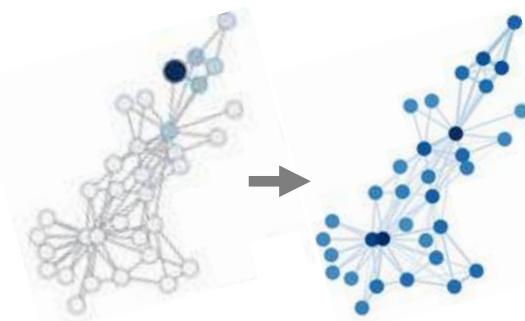
What contributes to over-squashing?



Theorem: Clique (>0) (informal) strong negatively-curved edges contribute to over-squashing.

Grid ($\bar{0}$) Tree ($\lessdot 0$)

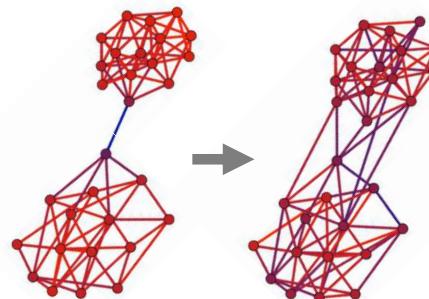
Graph Rewiring Approaches



**Connectivity Diffusion
(DIGL)**

Gasteiger et al. 2019

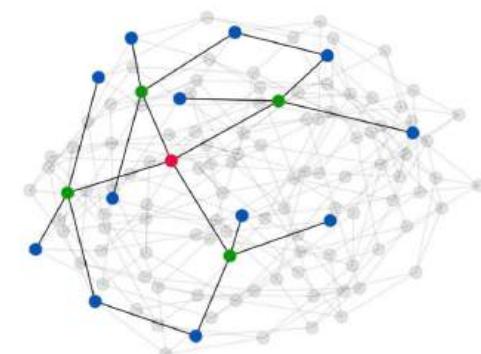
- KNN graph on PPR embedding
- + Good for homophilic graphs
- Bad for heterophilic graphs
- Drastically different graph



**Discrete Ricci flow
(SDRF)**

Topping, Di Giovanni et B. 2022

- Remove negatively-curved edges
- + Good for both homophilic and heterophilic graphs
- + “Surgical” rewiring
- Curvature is computationally expensive



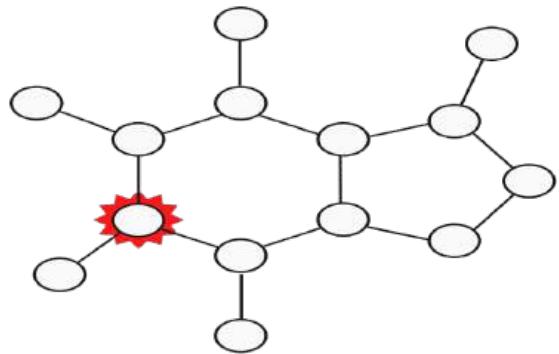
**Expander Propagation
(EGP)**

Deac et al. 2022

- Fixed expander graph
- Alternate MP on original and new graph
- + Optimal graph for MP
- No relation to input graph
- No permutation equivariance

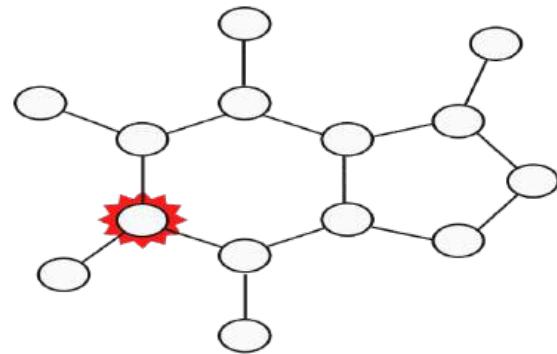
No relation to the task!

Why it is important to consider the task?



Van der Waals interactions

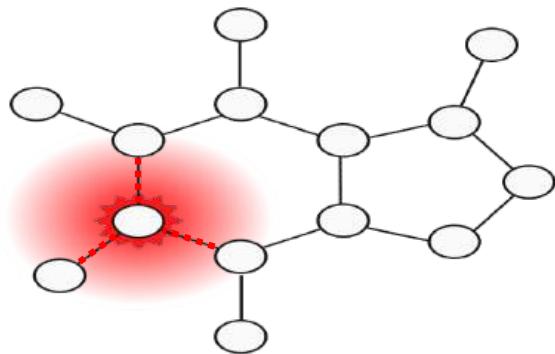
$$\propto r^{-12}$$



Coulomb interactions

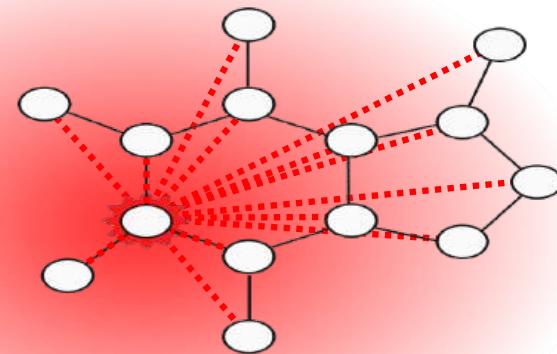
$$\propto r^{-1}$$

Why it is important to consider the task?



Van der Waals interactions

$$\propto r^{-12}$$



Coulomb interactions

$$\propto r^{-1}$$

Same graph+features, different task

Whether the graph is good depends on the task!



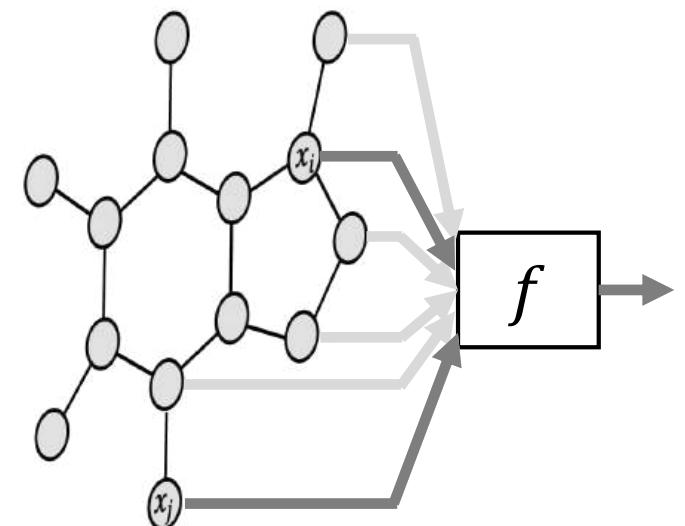
Long-range interactions & Expressivity

Long-range interactions in graph tasks

- **Task** = a function $f(\mathbf{X})$ on the node features of a graph G
- The interaction between features in nodes i and j required for the task is given by

$$\text{Mixing of } f: \text{ mix}_f(i, j) = \max_{\mathbf{X}} \max_{1 \leq \alpha, \beta \leq d} \left| \frac{\partial^2 f(\mathbf{X})}{\partial x_i^\alpha \partial x_j^\beta} \right|$$

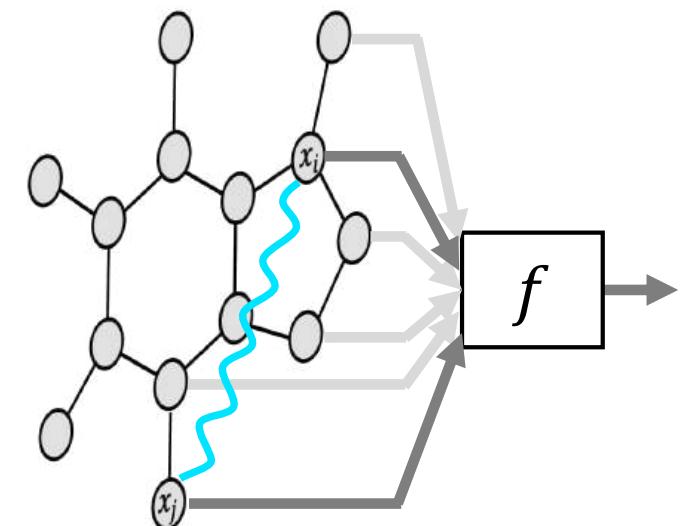
- $f(\mathbf{X}) = \phi(\mathbf{x}_i) + \phi(\mathbf{x}_j)$ is fully separable, thus $\text{mix}_f(i, j) = 0$



Long-range interactions in graph tasks

- **Task** = a function $f(\mathbf{X})$ on the node features of a graph G
- The interaction between features in nodes i and j required for the task is given by

$$\text{Mixing of } f: \text{ mix}_f(i, j) = \max_{\mathbf{X}} \max_{1 \leq \alpha, \beta \leq d} \left| \frac{\partial^2 f(\mathbf{X})}{\partial x_i^\alpha \partial x_j^\beta} \right|$$



- $f(\mathbf{X}) = \phi(\mathbf{x}_i) + \phi(\mathbf{x}_j)$ is fully separable, thus $\text{mix}_f(i, j) = 0$
- $f(\mathbf{X}) = \phi(\langle \mathbf{x}_i, \mathbf{x}_j \rangle)$ mixing depends on how non-linear ϕ is

Capacity bounds

$$\text{mix}_f(i, j) \leq \sum_{k=1}^{L-1} (\mathbf{c}_\sigma \mathbf{w})^{2L-k-1} \left(\mathbf{w}(\mathbf{s}^{L-k})^T \text{diag}(\mathbf{1}^T \mathbf{s}^k) \mathbf{s}^{L-k} + \mathbf{c} \mathbf{Q}_k \right)_{ij}$$

task model topology

What is the capacity of MPNN required for a given task?

Capacity bounds

$$\text{mix}_f(i, j) \leq \sum_{k=1}^{L-1} (\mathbf{c}_\sigma \mathbf{w})^2 \mathbf{s}^{L-k-1} \left(\mathbf{w}(\mathbf{s}^{L-k})^\top \text{diag}(\mathbf{1}^\top \mathbf{s}^k) \mathbf{s}^{L-k} + \mathbf{c} \mathbf{Q}_k \right)_{ij}$$

What is the capacity of MPNN required for a given task?

model + topology

mixing

Capacity bounds

$$\text{mix}_f(i, j) \leq \sum_{k=1}^{L-1} (c_\sigma w)^{2L-k-1} \left(w(\mathbf{S}^{L-k})^T \text{diag}(\mathbf{1}^T \mathbf{S}^k) \mathbf{S}^{L-k} + C \mathbf{Q}_k \right)_{ij}$$

Bound on weights w

$$w \geq \frac{d_{\min}}{c_2} \left(\frac{\text{mix}_f(i, j)}{q} \right)^{1/d(i, j)}$$

- d_{\min} =min node degree
- Fixed depth $L = [d(i, j)/2]$
- q =number of paths of length $d(i, j)$ between i and j

Bound on depth L

$$L \geq \frac{\alpha}{4c_2} + \frac{|E|}{\sqrt{d_i d_j}} (\alpha \text{mix}_f(i, j) - \beta)$$

- d_i =degree of node i
- α, β =model-related constants
- $|E|$ =number of edges
- Bounded weights

Capacity bounds

$$\text{mix}_f(i, j) \leq \sum_{k=1}^{L-1} (c_\sigma w)^{2L-k-1} \left(w(\mathbf{S}^{L-k})^\top \text{diag}(\mathbf{1}^\top \mathbf{S}^k) \mathbf{S}^{L-k} + C \mathbf{Q}_k \right)_{ij}$$

Bound on weights w

$$w \geq \frac{d_{\min}}{c_2} \left(\frac{\text{mix}_f(i, j)}{q} \right)^{1/d(i, j)}$$

“weights need to be large enough to allow mixing”

Bound on depth L

$$L \geq \frac{\tau(i, j)}{4c_2} + \frac{|E|}{\sqrt{d_i d_j}} (\alpha \text{mix}_f(i, j) - \beta)$$

- Depth must be \sim commute time $\tau(i, j)$
- Rewiring tries to improve τ
- τ can be as large as $O(n^3)$, which implies impossibility statements

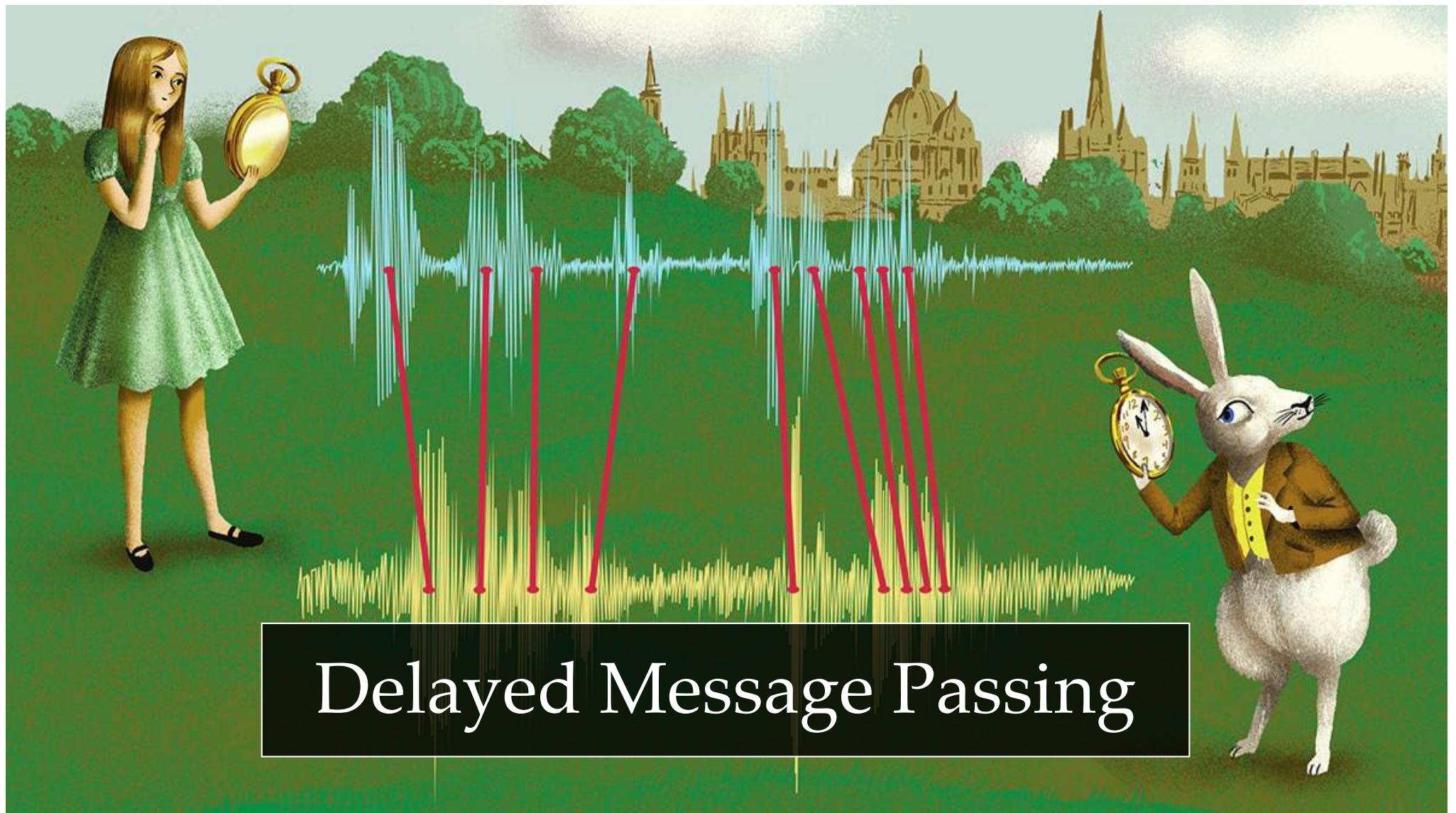
Expressive power without Weisfeiler-Lehman

Expressive power (informal): MPNN with $L \leq n$ layers *cannot learn* tasks that require high mixing among features at nodes with large commute time.



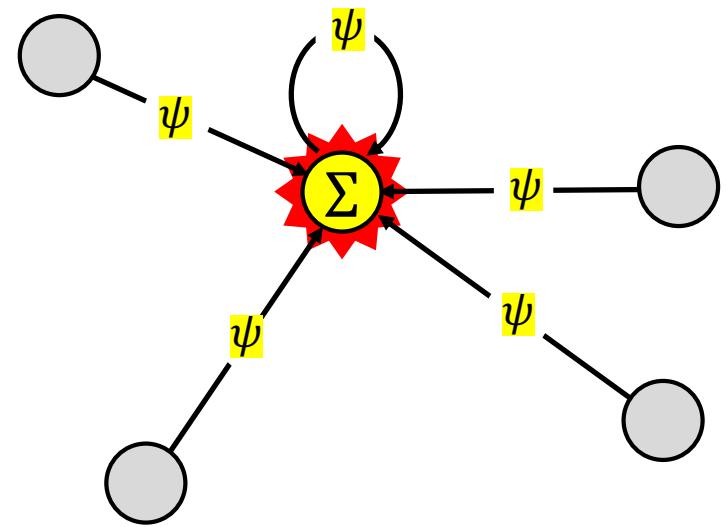
A. Lehman

B. Weisfeiler

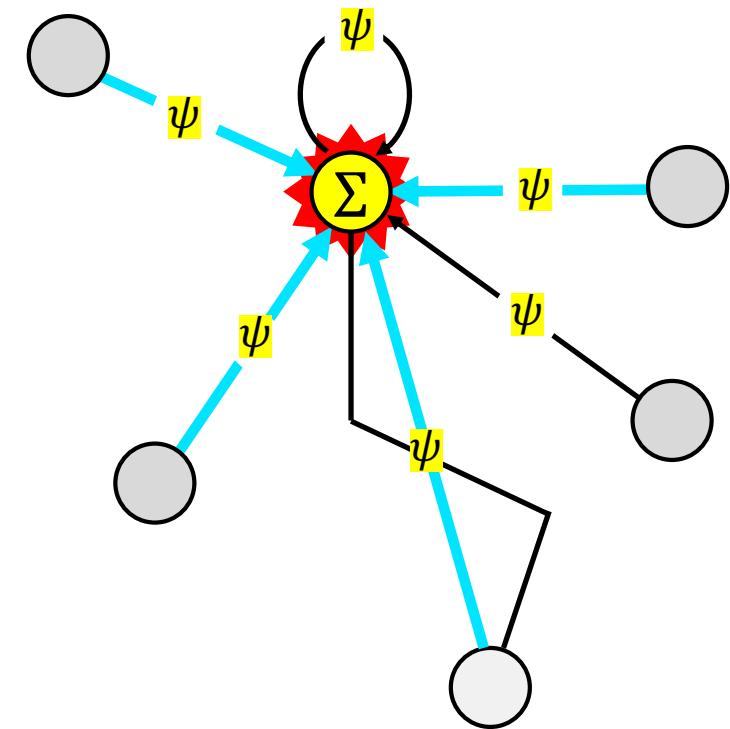


Delayed Message Passing

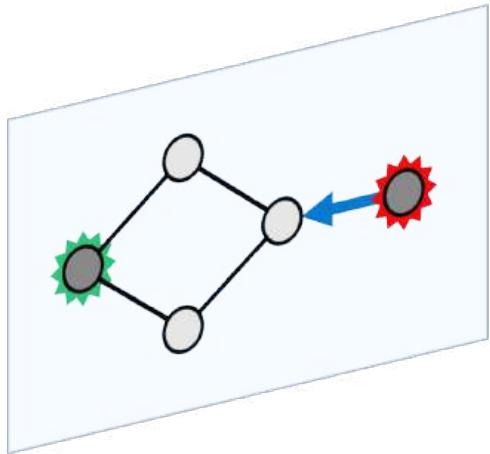
What



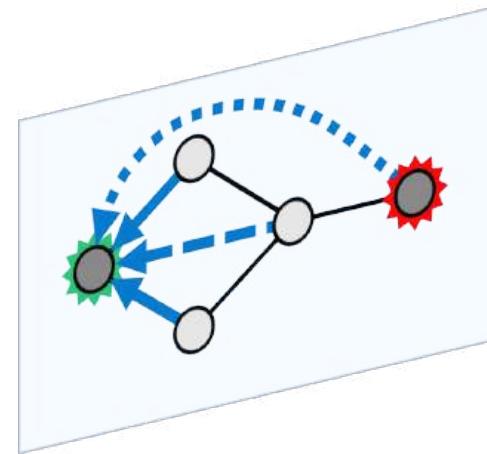
What + Where



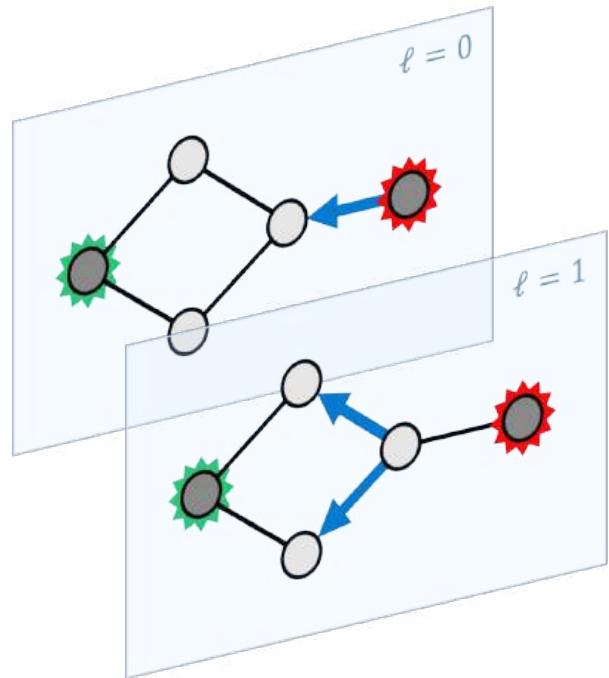
What + Where + When



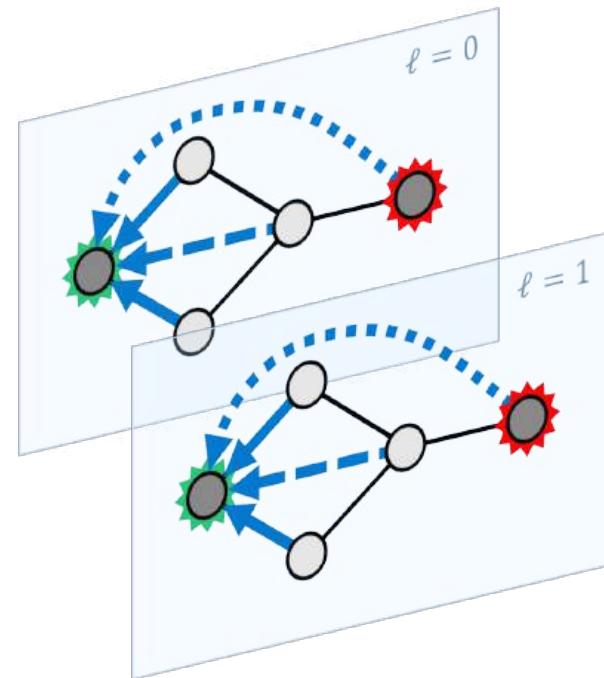
Classical MPNN



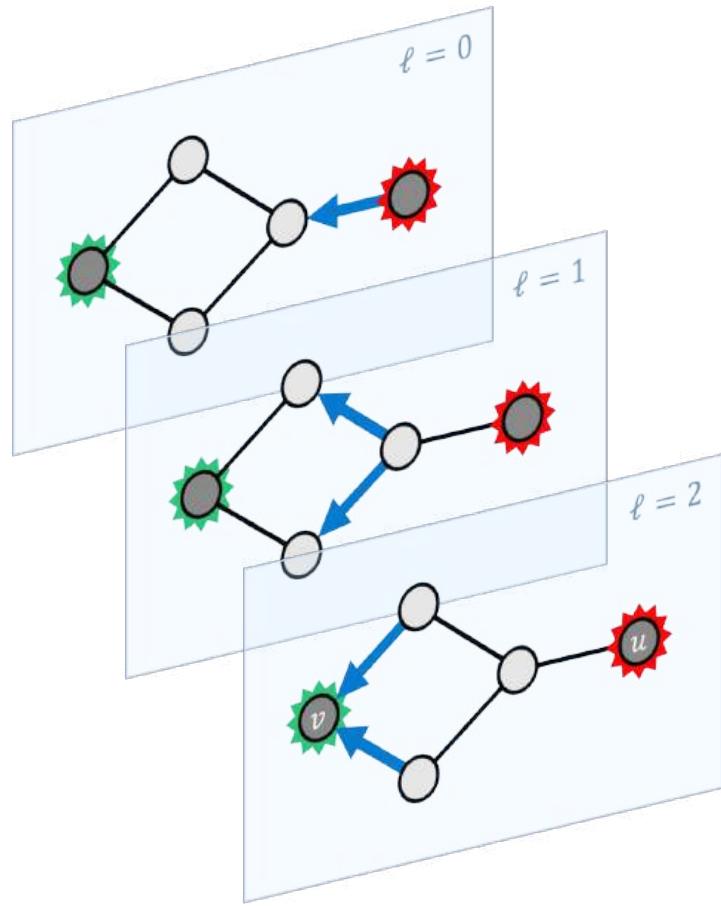
Graph Transformer



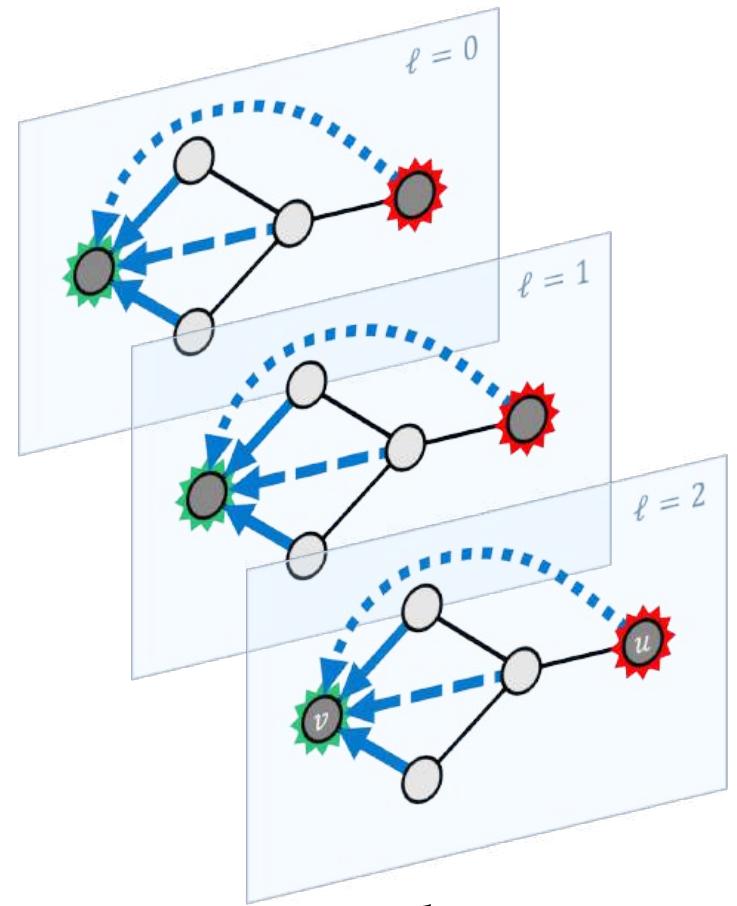
Classical MPNN



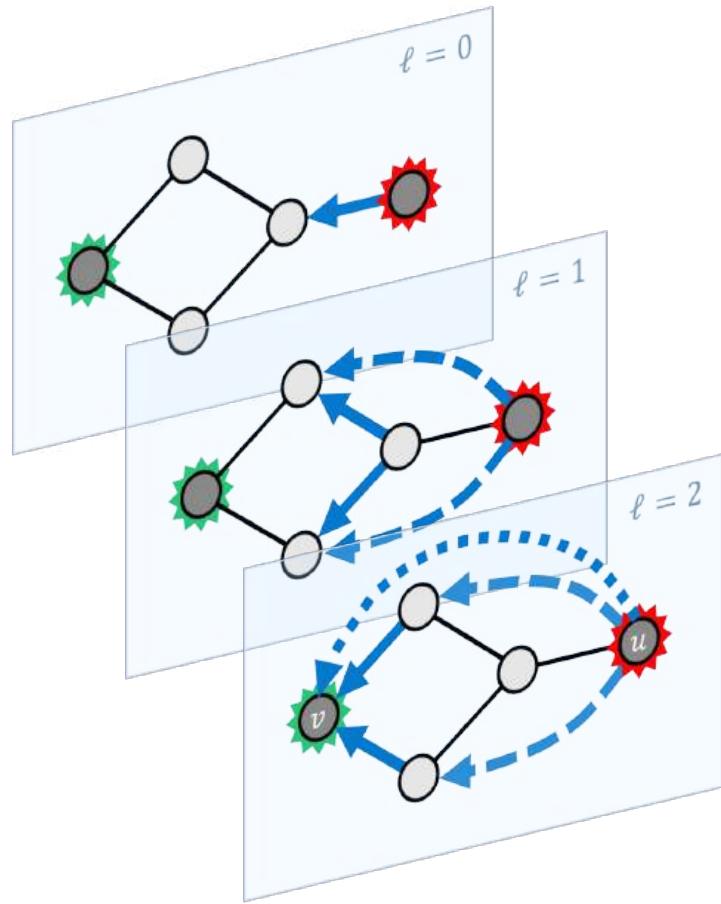
Graph Transformer



Classical MPNN

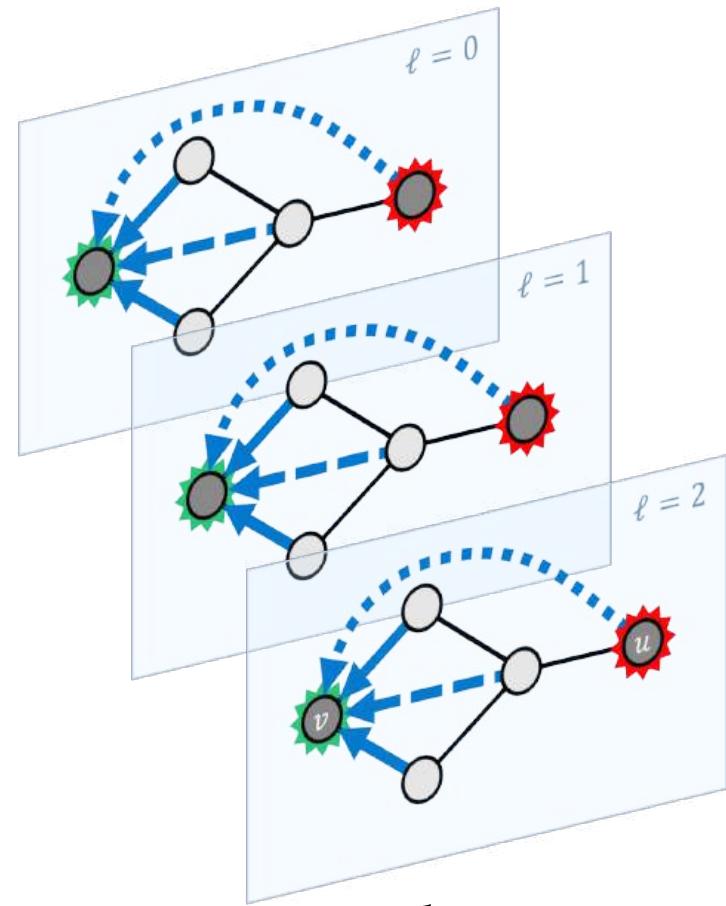


Graph Transformer

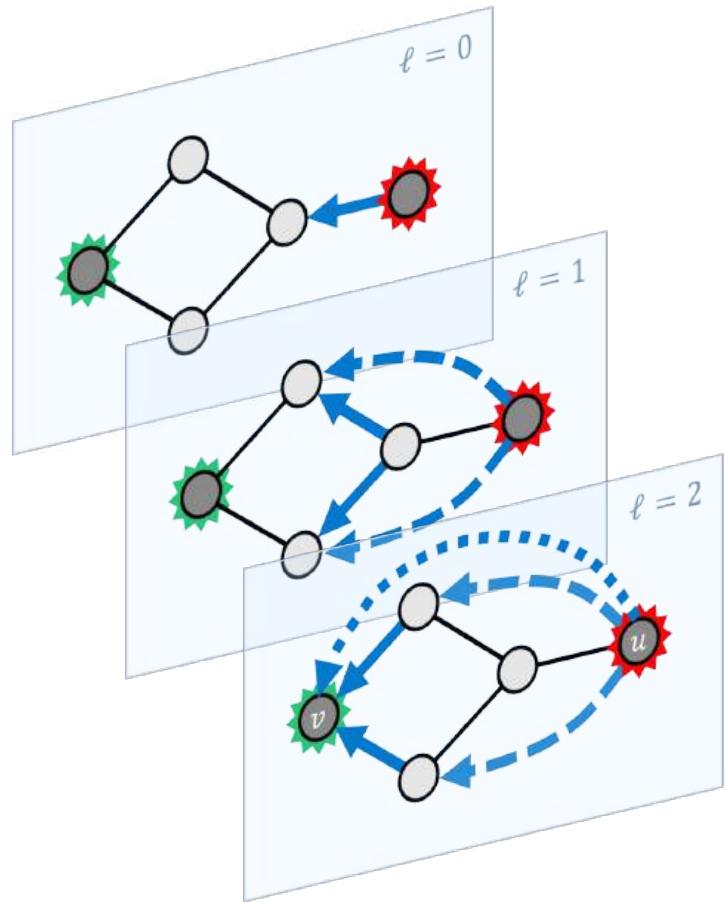


Dynamic Rewiring
(DRew)

Gutteridge, Di Giovanni et B 2023

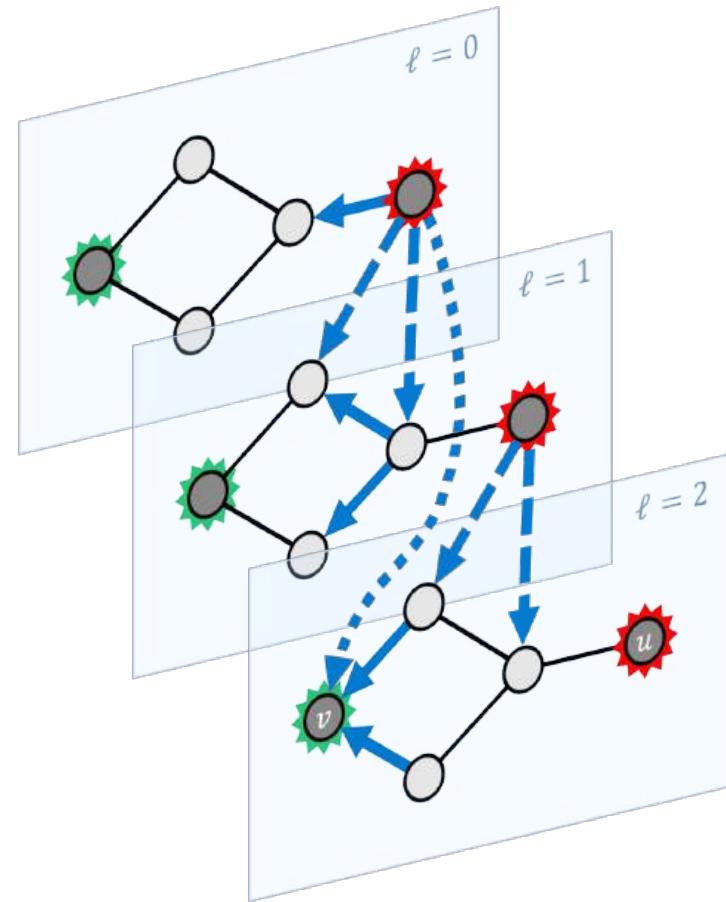


Graph Transformer



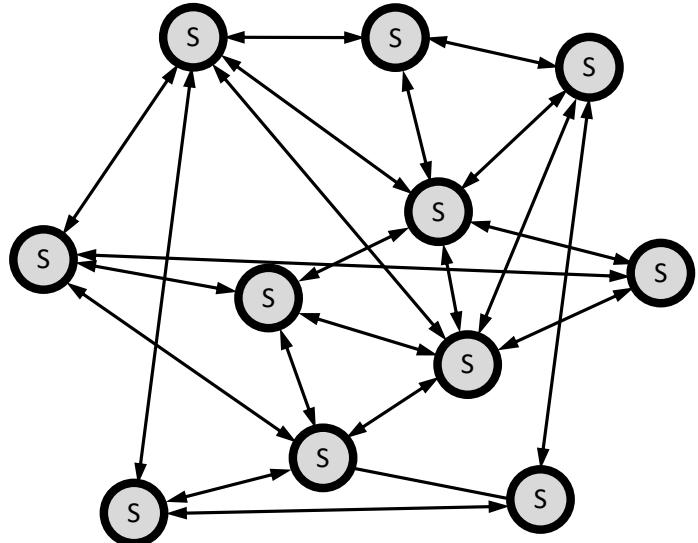
Dynamic Rewiring
(DRew)

Gutteridge, Di Giovanni et B 2023

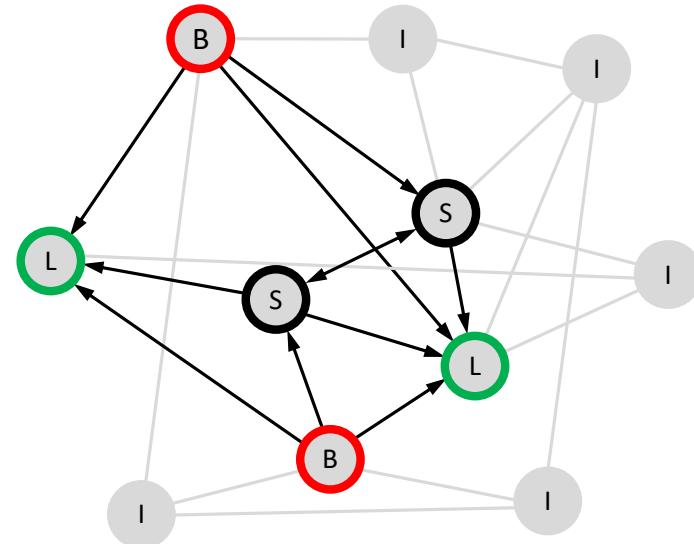


Dynamic Rewiring + delay
(vDRew)

Cooperative Message Passing

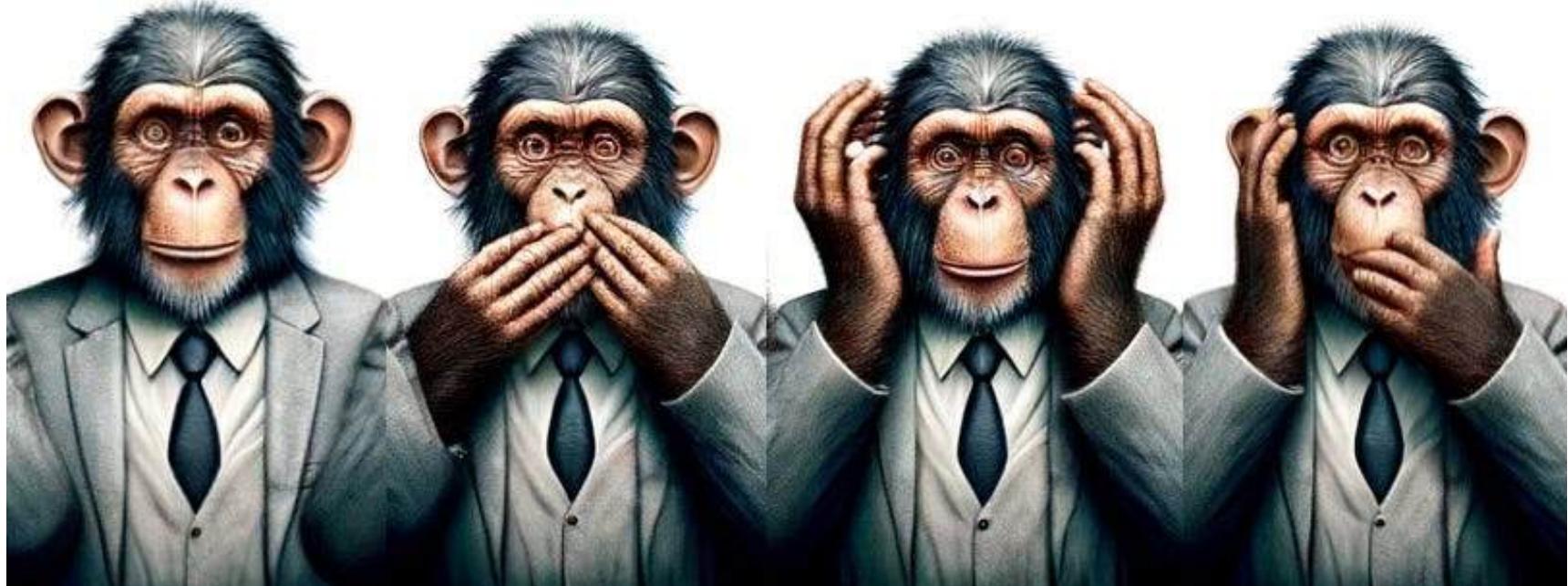


Standard Message Passing
each node Broadcasts & Listens



Cooperative Message Passing
each node individually decides

Cooperative Message Passing



**Broadcast &
Listen**

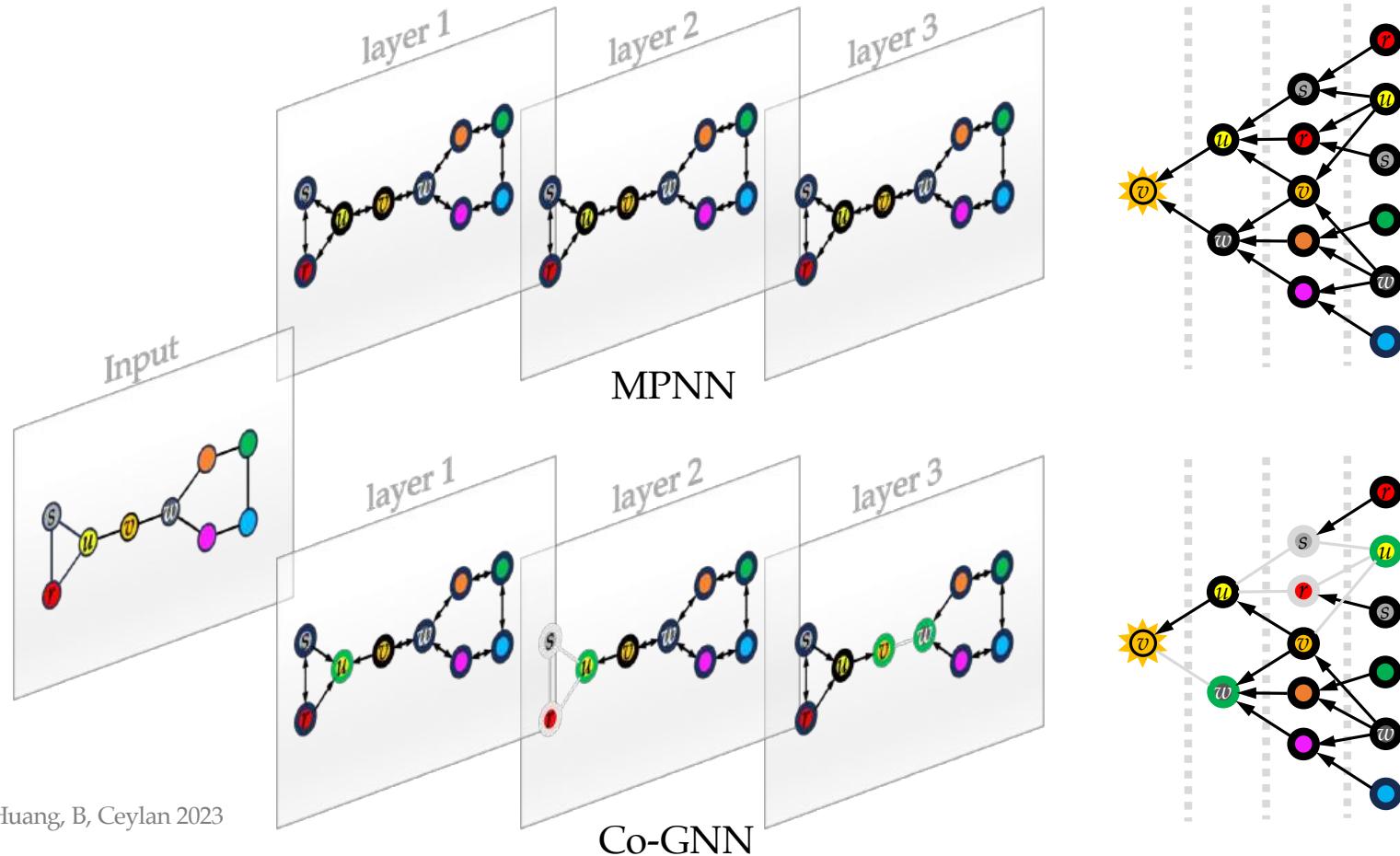
Listen

Broadcast

Isolate

Finkelshtein, Huang, B, Ceylan 2023; Illustration: DALL-E 3 (after a lot of effort)

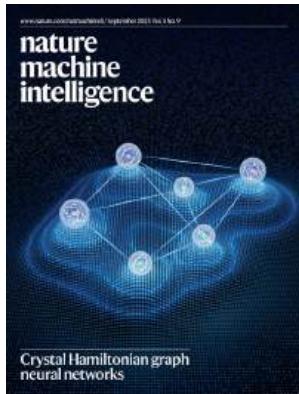
Cooperative Message Passing



Finkelshtein, Huang, B, Ceylan 2023

What do we gain from physics-inspired GNNs?

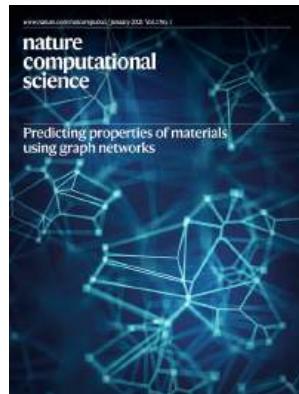
- New perspectives on old problems (e.g. oversmoothing, bottlenecks, etc.)
- Explains old architectures & gives rise to new ones
- Principled architectural choices (residual connection, shared symmetric weights)
- Theoretical guarantees (e.g. stability, convergence, expressive power, etc.)
- Deep links to other fields less known in GNN literature (e.g. differential geometry & algebraic topology)
- In GNNs, the graph is both *input* and *computational device* – not all graphs are good!
- Rewiring tells *what* messages to send *where*
- Dynamic rewiring+delay adds control also *when*



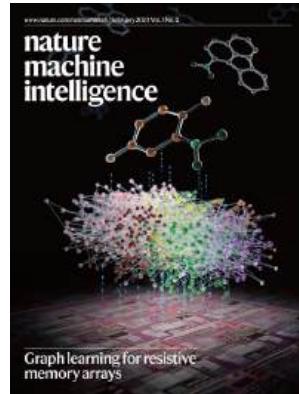
Physics



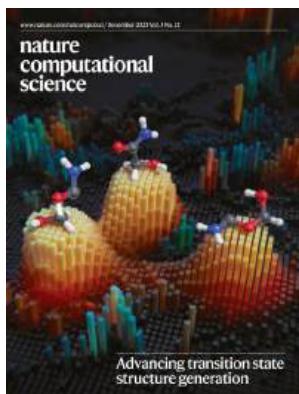
New materials



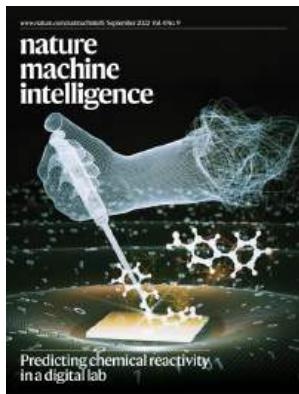
Chip design



Biology



Chemistry



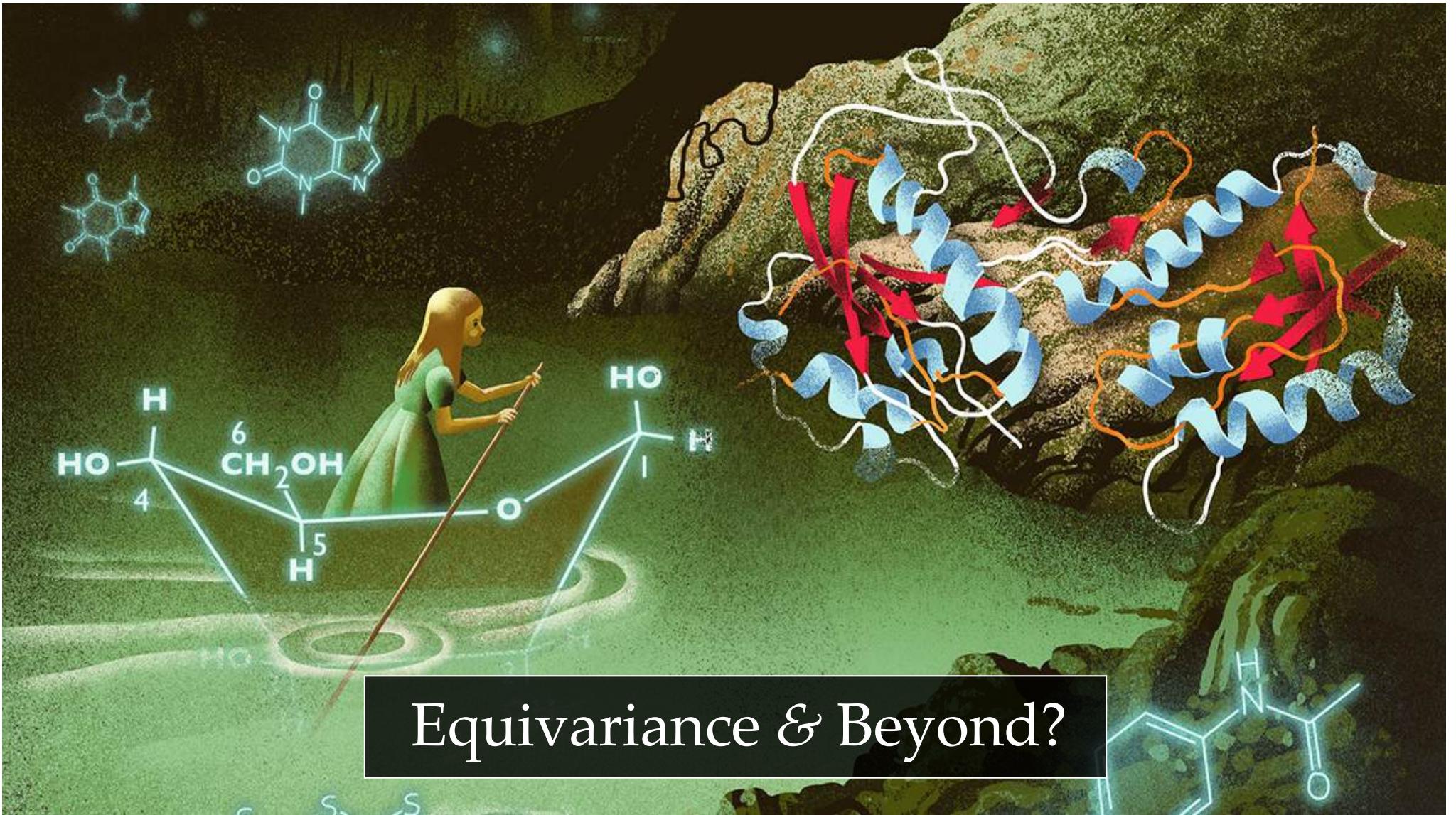
Urban planning



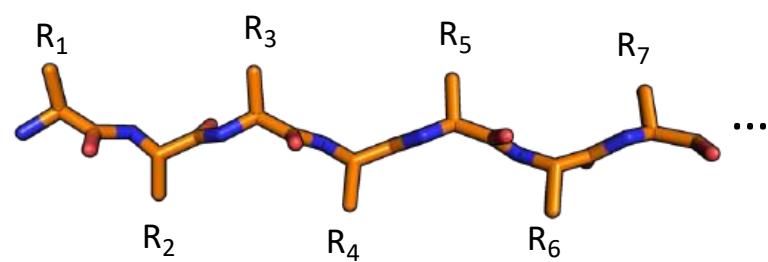
Pure math



Weather

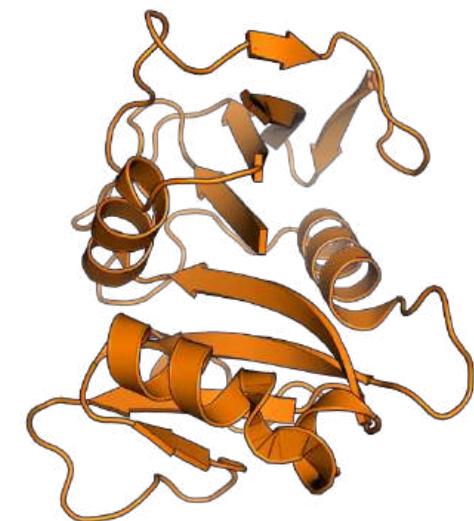


Protein Folding



Sequence of
amino acids

Protein folding



3D structure

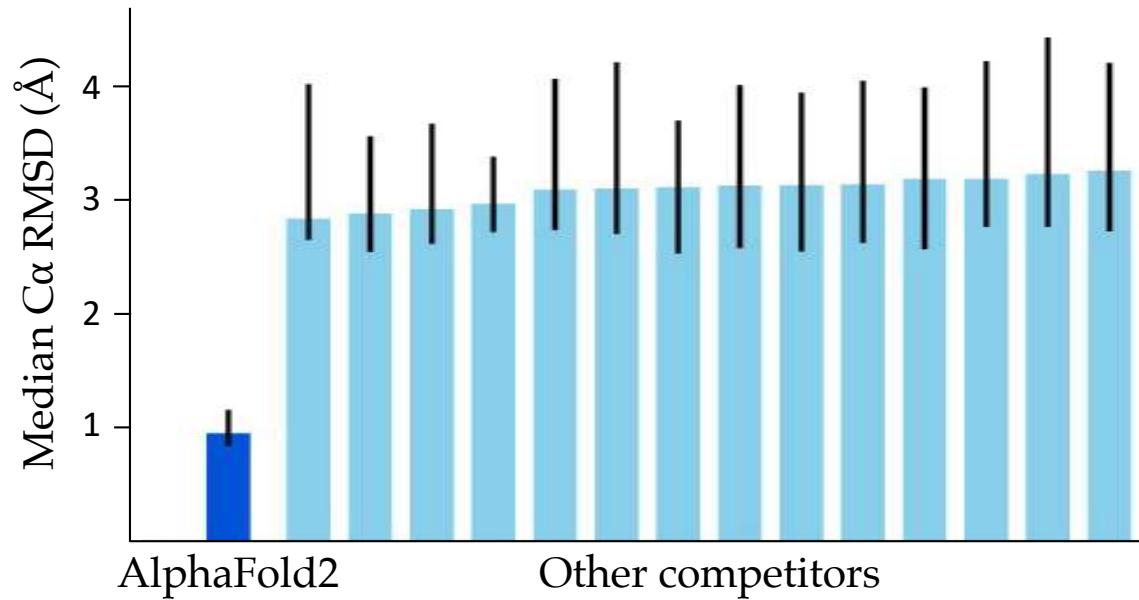
“[protein] conformation is determined by the totality of interatomic interactions and hence by the amino acid sequence”



C. Anfinsen

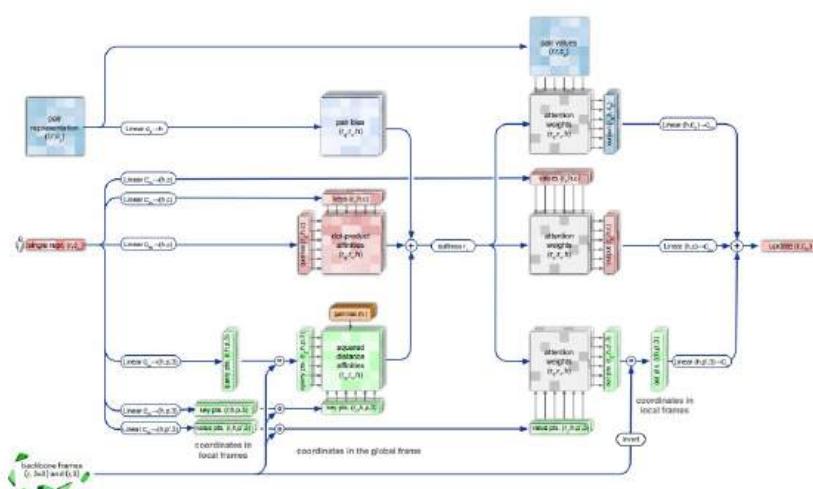
Anfinsen 1972 Nobel Lecture; Portrait: Ihor Gorskyi

AlphaFold2: the “ImageNet moment” of Structural Biology



Jumper et al. 2022; CASP14

Revolution in Structural Biology



Jumper et al. 2021

AlphaFold 2
“Invariant point
attention”



Baek et al. 2021

RosettaFold
SE(3)-equivariant
Transformer



2024 Nobel Prize in Chemistry
for computational protein design

AND THEY LIVED

HAPPILY EVER AFTER



The Bitter Lesson: Equivariance is dead...long live equivariance!

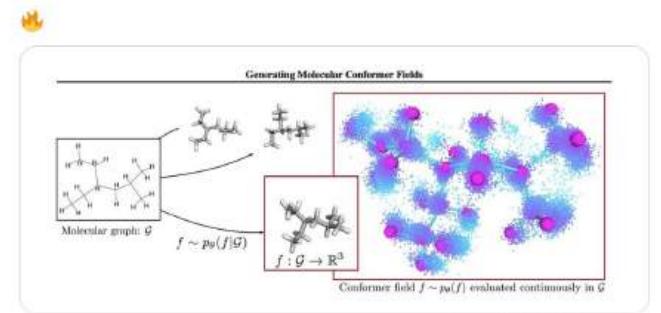
► Equivariance is the idea of giving a model the inductive biases to natively handle rotations, translations and (sometimes) reflections. It has been at the core of Geometric Deep Learning and biomolecular modelling research since AlphaFold 2. However, recent works by top labs have questioned the existing mantra.

- The first shots were fired by Apple, with a paper that obtained SOTA results on predicting the 3D structures of small molecules using a non-equivariant diffusion model with a transformer encoder.
- Remarkably, the authors showed that using the domain-agnostic model did not deleteriously impact generalization and was consistently able to outperform specialist models (assuming sufficient scale was used).
- Next was AlphaFold 3, which infamously dropped all the equivariance and frames constraints from the previous model in favour of another diffusion process coupled with augmentations and, of course, scale.
- Regardless, the greatly improved training efficiency of equivariant models means the practice is likely to stay for a while (at least for academic groups working on large systems such as proteins).



"We [...] empirically show that explicitly enforcing roto-translation equivariance is not a strong requirement for generalization."

"Furthermore, we also show that approaches that do not explicitly enforce roto-translation equivariance (like ours) can match or outperform approaches that do."



stateof.ai 2024



Swallowing the Bitter Pill: Simplified Scalable Conformer Generation

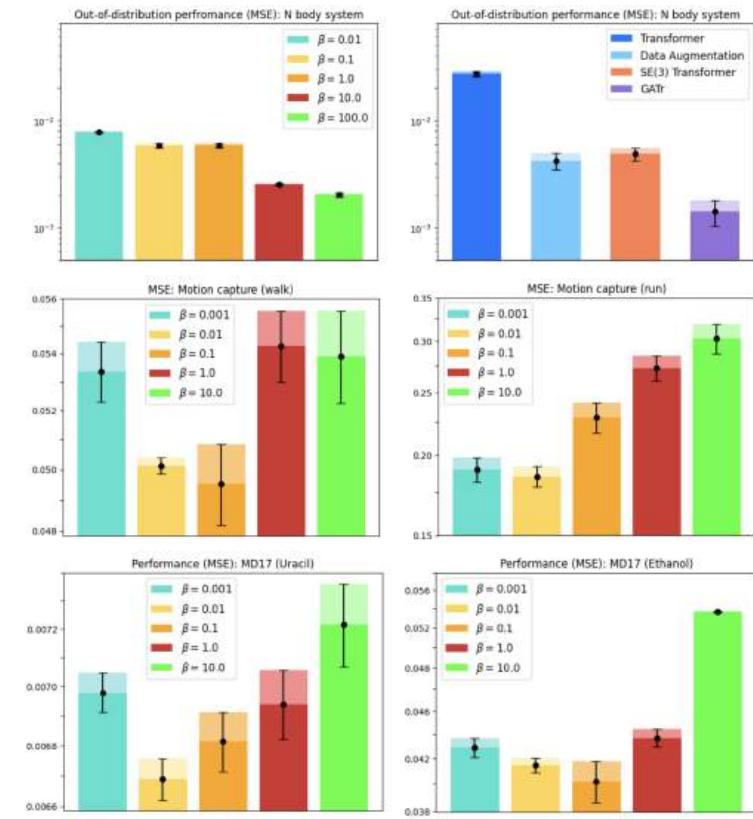
Yuyang Wang¹ Ahmed A. Elhag^{1,2} Navdeep Jaitly¹ Joshua M. Susskind¹ Miguel Ángel Bautista¹

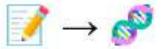
Abstract

We present a novel way to predict molecular conformers through a simple formulation that sidesteps many of the heuristics of prior works and achieves state of the art results by using the advantages of scale. By training a diffusion generative model directly on 3D atomic positions without making assumptions about the explicit structure of molecules (*e.g.* modeling torsional angles) we are able to radically simplify structure learning, and make it trivial to scale up the model sizes. This model, called Molecular Conformer Fields (MCF), works by parameterizing conformer structures as functions that map elements from a molecular graph directly to their 3D location in space. This formulation allows us to boil down the essence of structure prediction to learning a distribution over functions. Experimental results show that scaling up the model capacity leads to large gains in generalization performance *without enforcing inductive biases* like rotational equivariance. MCF represents an advance in extending diffusion models to handle complex scientific problems in a conceptually simple, scalable and effective manner.

is the vast complexity of the 3D structure space, encompassing factors such as bond lengths and torsional angles. Despite the molecular graph dictating potential 3D conformers through specific constraints, such as bond types and spatial arrangements determined by chiral centers, the conformational space experiences exponential growth with the expansion of the graph size and the number of rotatable bonds (Axelrod & Gomez-Bombarelli, 2022). This complicates brute force and exhaustive approaches, making them virtually unfeasible for even moderately small molecules.

Systematic methods, like OMEGA (Hawkins et al., 2010), offer rapid processing through rule-based generators and curated torsion templates. Despite their efficiency, these models typically fail on complex molecules, as they often overlook global interactions and are tricky to extend to inputs like transition states or open-shell molecules. Classic stochastic methods, like molecular dynamics (MD) and Markov chain Monte Carlo (MCMC), rely on extensively exploring the energy landscape to find low-energy conformers. Such techniques suffer from sampling inefficiency for large molecules and struggle to generate diverse representative conformers (Hawkins, 2017; Wilson et al., 1991; Grebner et al., 2011). In the domain of learning-based approaches, several works have looked at conformer generation problems through the lens of probabilistic modeling, using either





The Bitter Lesson: Equivariance is dead...long live equivariance!

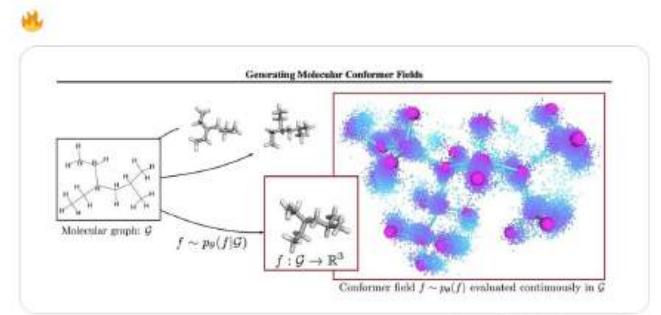
► Equivariance is the idea of giving a model the inductive biases to natively handle rotations, translations and (sometimes) reflections. It has been at the core of Geometric Deep Learning and biomolecular modelling research since AlphaFold 2. However, recent works by top labs have questioned the existing mantra.

- The first shots were fired by Apple, with a paper that obtained SOTA results on predicting the 3D structures of small molecules using a non-equivariant diffusion model with a transformer encoder.
- Remarkably, the authors showed that using the domain-agnostic model did not deleteriously impact generalization and was consistently able to outperform specialist models (assuming sufficient scale was used).
- Next was AlphaFold 3, which infamously dropped all the equivariance and frames constraints from the previous model in favour of another diffusion process coupled with augmentations and, of course, scale.
- Regardless, the greatly improved training efficiency of equivariant models means the practice is likely to stay for a while (at least for academic groups working on large systems such as proteins).



"We [...] empirically show that explicitly enforcing roto-translation equivariance is not a strong requirement for generalization."

"Furthermore, we also show that approaches that do not explicitly enforce roto-translation equivariance (like ours) can match or outperform approaches that do."



stateof.ai 2024





Quo vadimus?

The Hardware Lottery

Sara Hooker

Google Research, Brain Team

Hardware, systems and algorithms research communities have historically had different incentive structures and fluctuating motivation to engage with each other explicitly. This historical treatment is odd given that hardware and software have frequently determined which research ideas succeed (and fail). This essay introduces the term hardware lottery to describe when a research idea wins because it is suited to the available software and hardware and *not* because the idea is superior to alternative research directions.

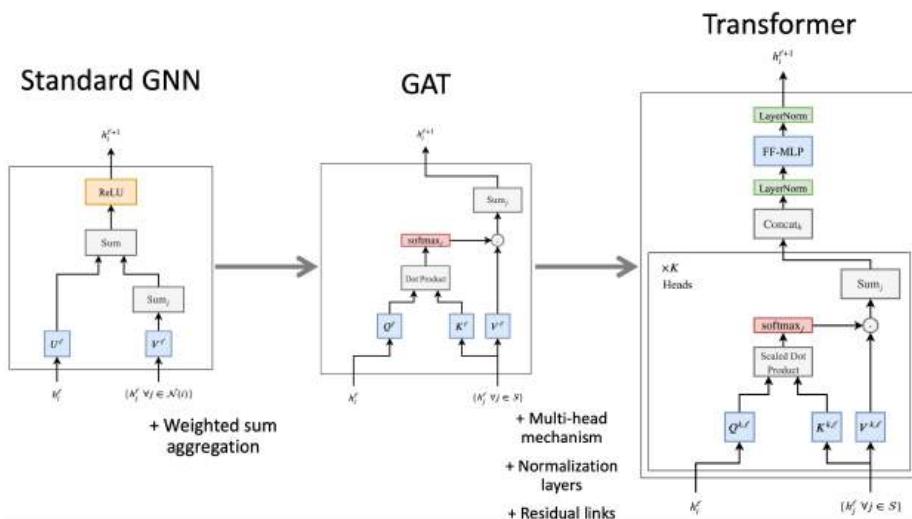
Examples from early computer science history illustrate how hardware lotteries can delay research progress by casting successful ideas as failures. These lessons are particularly salient given the advent of domain specialized hardware which make it increasingly costly to stray off of the beaten path of research ideas. This essay posits that the gains from progress in computing are likely to become even more uneven, with certain research directions moving into the fast-lane while progress on others is further obstructed.

Transformers are Graph Neural Networks

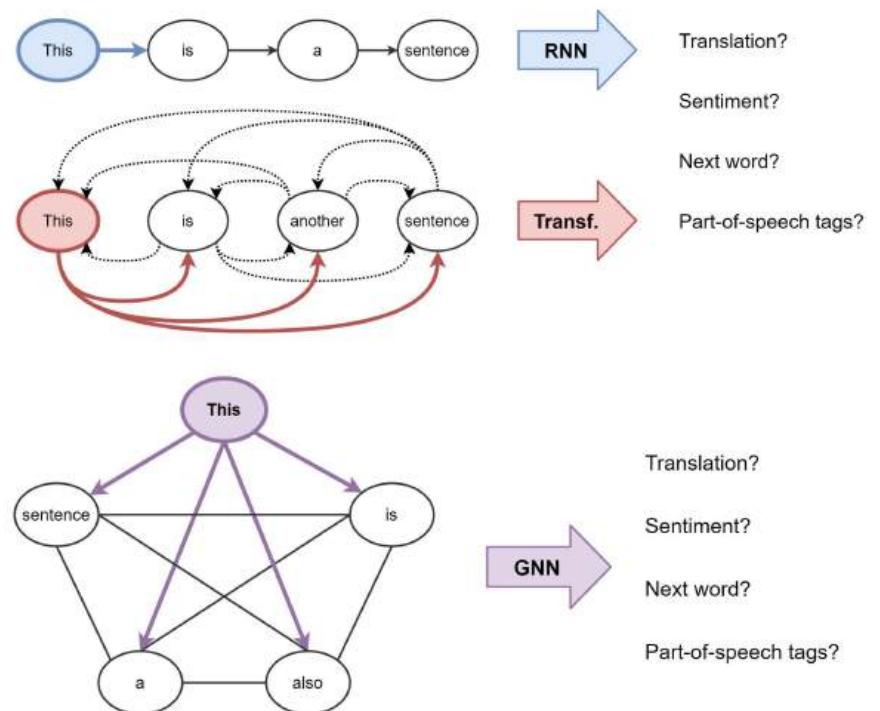
Exploring the connection between Transformer models such as GPT and BERT for Natural Language Processing, and Graph Neural Networks.

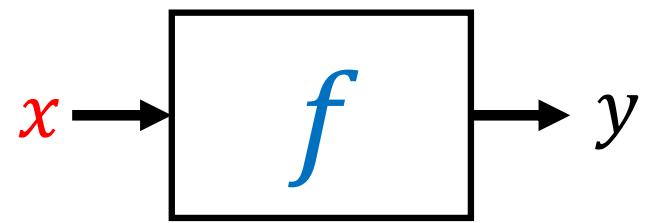
Chaitanya K. Joshi

Last updated on Jun 21, 2021 · 12 min read



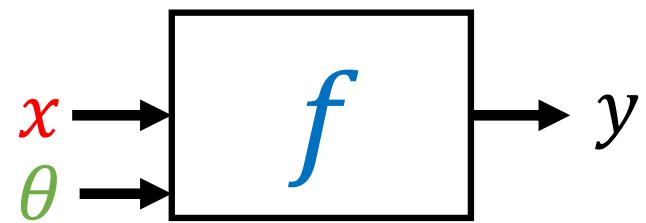
Joshi 2021





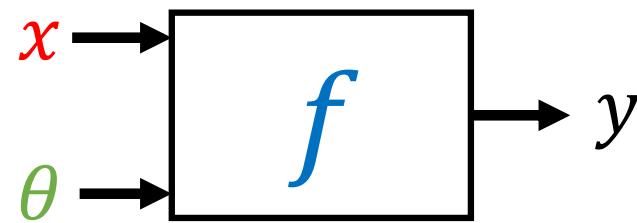
“How f interacts with the group G acting on x ? ”

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



“How f interacts with the group G acting on x ? ”

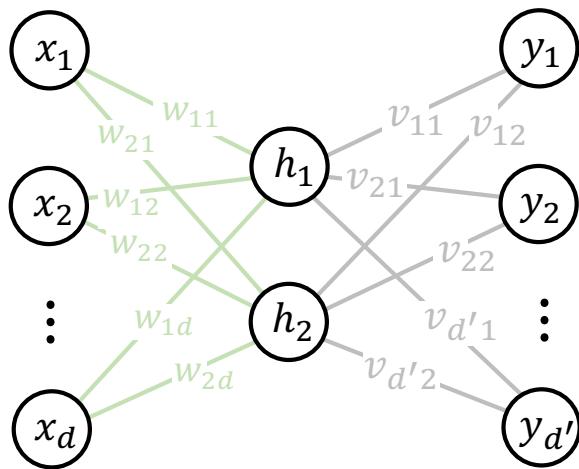
$$f(g \cdot x) = f(x)$$



“How f interacts with the group G acting on x and H acting on θ ? ”

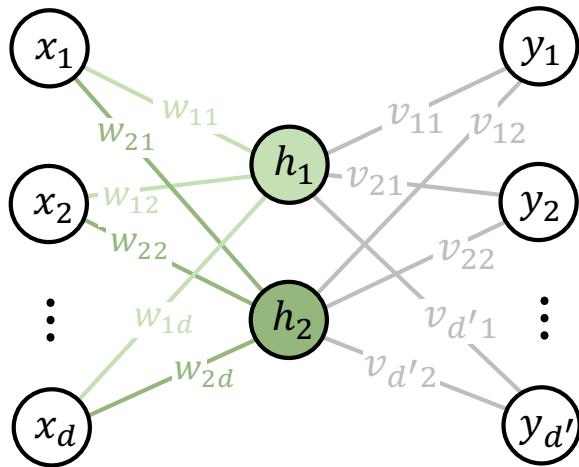
$$f(g.x, h.\theta) = f(x, \theta)$$

Symmetries of the Weights



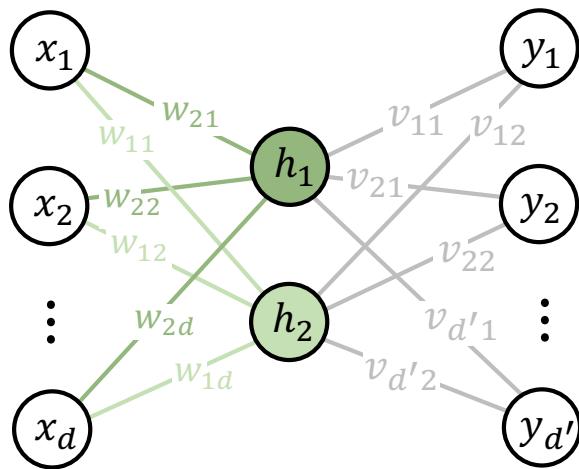
$$\mathbf{y} = \begin{matrix} \mathbf{V} \\ \vdots \end{matrix} \sigma \left(\begin{matrix} \mathbf{W} & \mathbf{x} \\ \vdots & \vdots \end{matrix} \right)$$

Symmetries of the Weights



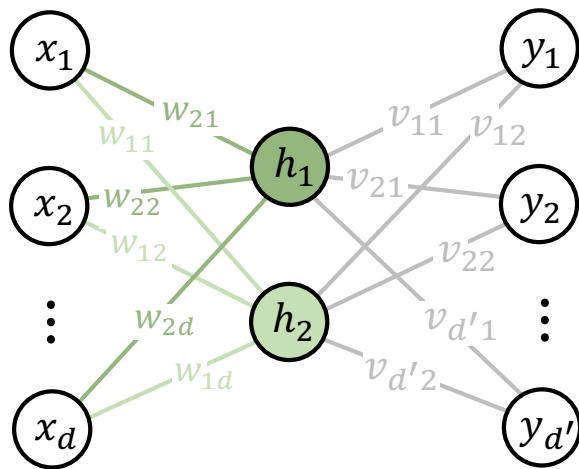
$$\mathbf{y} = \mathbf{V} \sigma \left(\begin{array}{c|c} \mathbf{W} & \mathbf{x} \end{array} \right)$$

Symmetries of the Weights



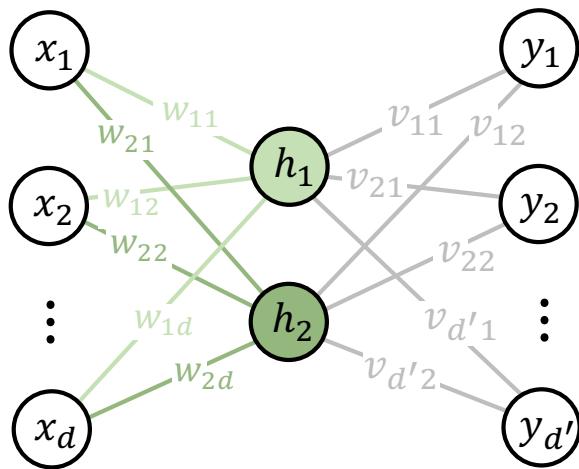
$$\mathbf{y} = \mathbf{V} \sigma \left(\boldsymbol{\Pi} \begin{array}{c|c} \mathbf{W} & \mathbf{x} \end{array} \right)$$

Symmetries of the Weights



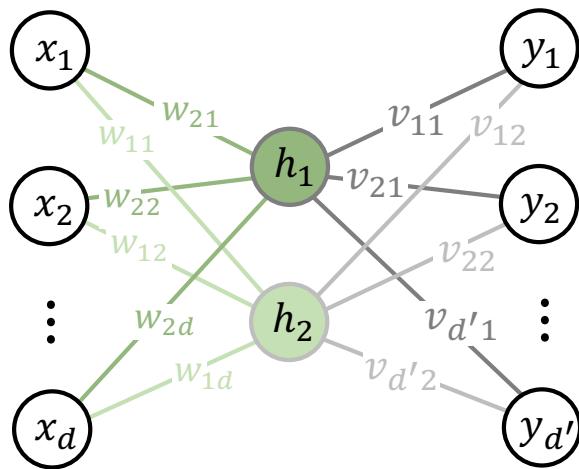
$$\mathbf{y} = \begin{matrix} \mathbf{V} \\ \vdots \end{matrix} \Pi \sigma \left(\begin{matrix} \mathbf{W} & \mathbf{x} \\ \hline \mathbf{w} & \mathbf{x} \end{matrix} \right)$$

Symmetries of the Weights



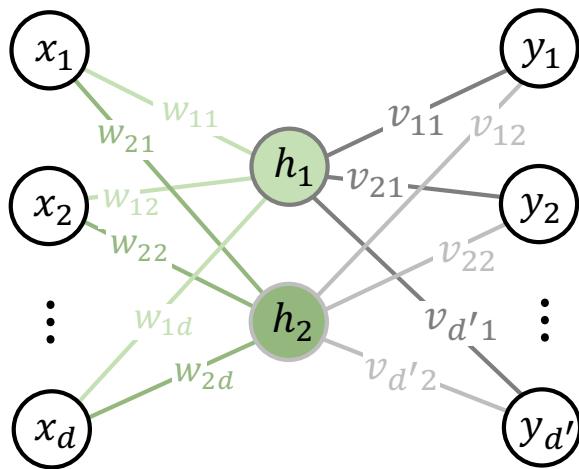
$$\mathbf{y} = \mathbf{V} \Pi \sigma \left(\begin{array}{c|c} \mathbf{W} & \mathbf{x} \end{array} \right)$$

Symmetries of the Weights



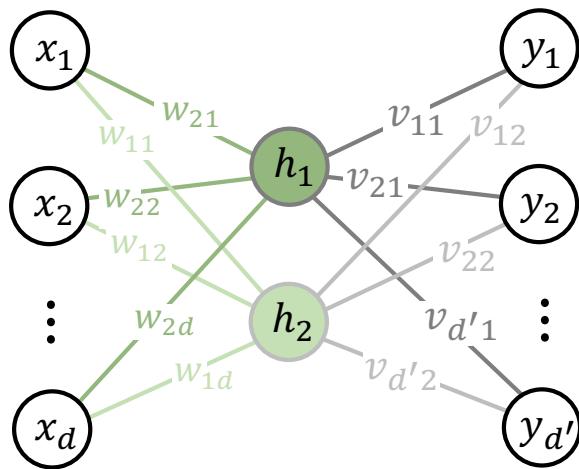
$$\mathbf{y} = \begin{matrix} \mathbf{V} \\ \vdots \end{matrix} \Pi \sigma \left(\begin{matrix} \mathbf{W} & \mathbf{x} \\ \hline \mathbf{w} & \mathbf{x} \end{matrix} \right)$$

Symmetries of the Weights



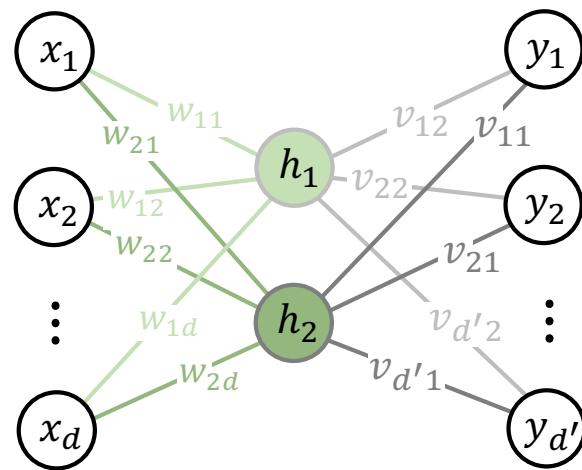
$$\mathbf{y} = \mathbf{V} \Pi \sigma \left(\begin{array}{c|c} \mathbf{W} & \mathbf{x} \end{array} \right)$$

Symmetries of the Weights



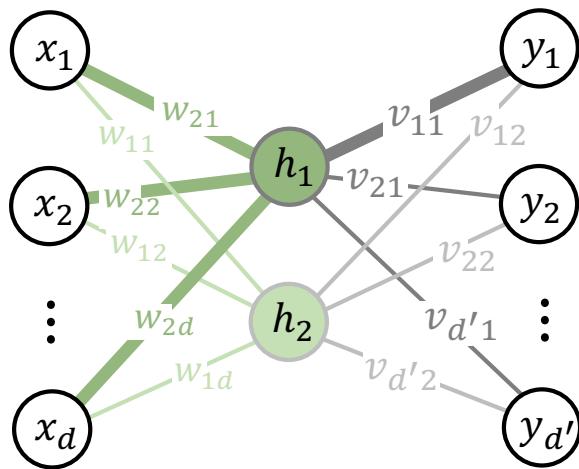
$$\mathbf{y} = \begin{matrix} \mathbf{V} \\ \vdots \end{matrix} \boldsymbol{\Pi}^T \boldsymbol{\Pi} \sigma \left(\begin{matrix} \mathbf{W} & \mathbf{x} \end{matrix} \right)$$

Symmetries of the Weights



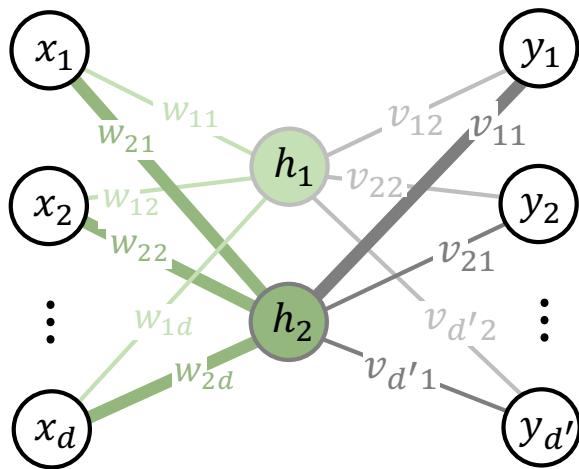
$$\mathbf{y} = \begin{matrix} \mathbf{V} \\ \vdots \end{matrix} \mathbf{\Pi}^T \mathbf{\Pi} \sigma \left(\begin{matrix} \mathbf{W} & \mathbf{x} \\ \vdots & \vdots \end{matrix} \right)$$

Symmetries of the Weights



$$\mathbf{y} = \begin{matrix} \mathbf{V} \\ \vdots \end{matrix} \mathbf{\Pi}^T \mathbf{\Pi} \sigma \left(\begin{matrix} \mathbf{W} & \mathbf{x} \\ \vdots & \vdots \end{matrix} \right)$$

Symmetries of the Weights



$$\mathbf{y} = \begin{matrix} \mathbf{V} \\ \vdots \end{matrix} \boldsymbol{\Pi}^T \boldsymbol{\Pi} \sigma \left(\begin{matrix} \mathbf{W} & \mathbf{x} \\ \vdots & \vdots \end{matrix} \right)$$

Symmetries of the Weights

- L -layer neural network with weights $\boldsymbol{\theta} = (\mathbf{W}_1, \mathbf{b}_1, \dots, \mathbf{W}_L, \mathbf{b}_L)$
- Parameter space symmetry $G = S_{d_1} \times \dots \times S_{d_L}$

$$\mathbf{W}'_1 = \boldsymbol{\Pi}_1^T \mathbf{W}_1$$

$$\mathbf{b}'_1 = \boldsymbol{\Pi}_1^T \mathbf{b}_1$$

$$\mathbf{W}'_l = \boldsymbol{\Pi}_l^T \mathbf{W}_l \boldsymbol{\Pi}_{l-1}$$

$$\mathbf{b}'_l = \boldsymbol{\Pi}_l^T \mathbf{b}_l$$

⋮

⋮

$$\mathbf{W}'_L = \mathbf{W}_L \boldsymbol{\Pi}_{L-1}$$

$$\mathbf{b}'_L = \mathbf{b}_L$$

such that $f(\cdot, g\boldsymbol{\theta}) = f(\cdot, \boldsymbol{\theta})$

Symmetries in the Gradient Space

$$\mathcal{L}_{\mathcal{B}} = \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}} \mathcal{L}_{(\mathbf{x}, \mathbf{y})}$$
$$\ell(f(\mathbf{x}, \boldsymbol{\theta}), \mathbf{y})$$

$$\nabla_{\mathbf{b}_l} \mathcal{L}_{(\mathbf{x}, \mathbf{y})} = \mathbf{g}_l \quad \text{where } \mathbf{g}_l = \frac{\partial \ell(f(\mathbf{x}, \boldsymbol{\theta}), \mathbf{y})}{\partial \mathbf{u}_l}$$

$$\nabla_{\mathbf{W}_l} \mathcal{L}_{(\mathbf{x}, \mathbf{y})} = \mathbf{g}_l \mathbf{a}_{l-1}^T$$

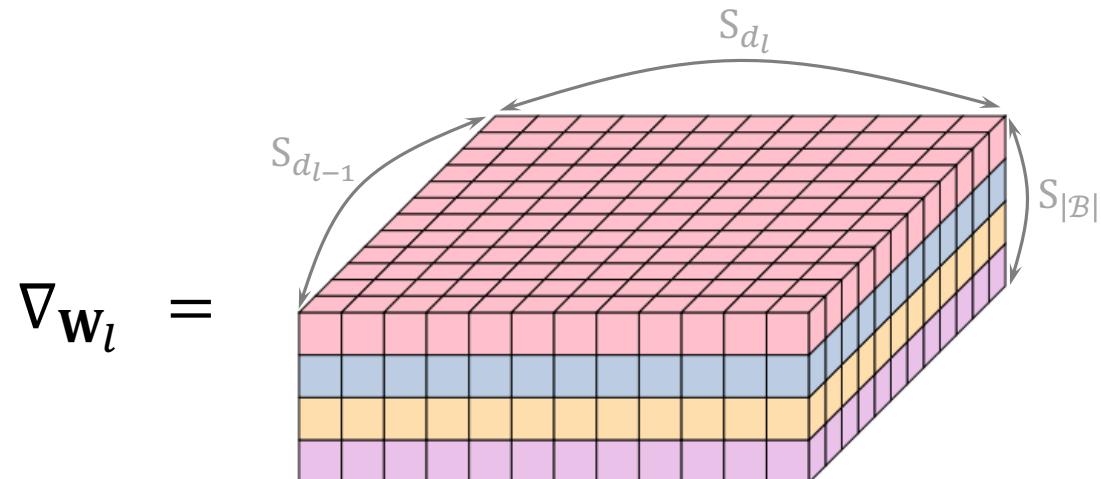
Symmetries in the Gradient Space

$$\mathcal{L}_{\mathcal{B}} = \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}} \ell(f(\mathbf{x}, \boldsymbol{\theta}), \mathbf{y})$$

$$\nabla_{\mathbf{W}_l} = \begin{matrix} d_l \\ d_{l-1} \\ |B| \end{matrix}$$

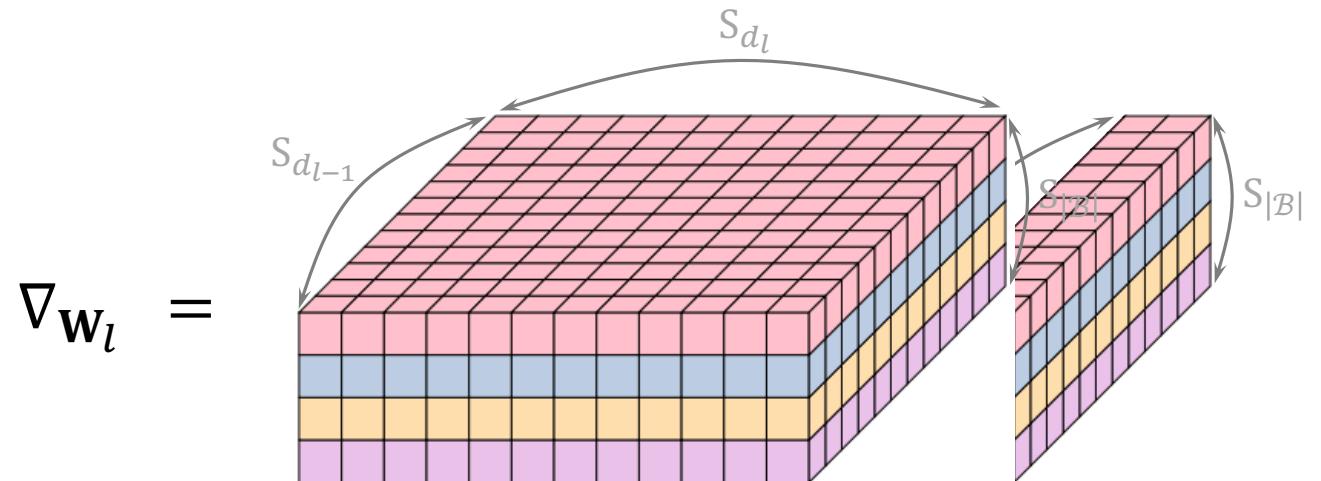
Symmetries in the Gradient Space

$$\mathcal{L}_{\mathcal{B}} = \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}} \ell(f(\mathbf{x}, \boldsymbol{\theta}), \mathbf{y})$$

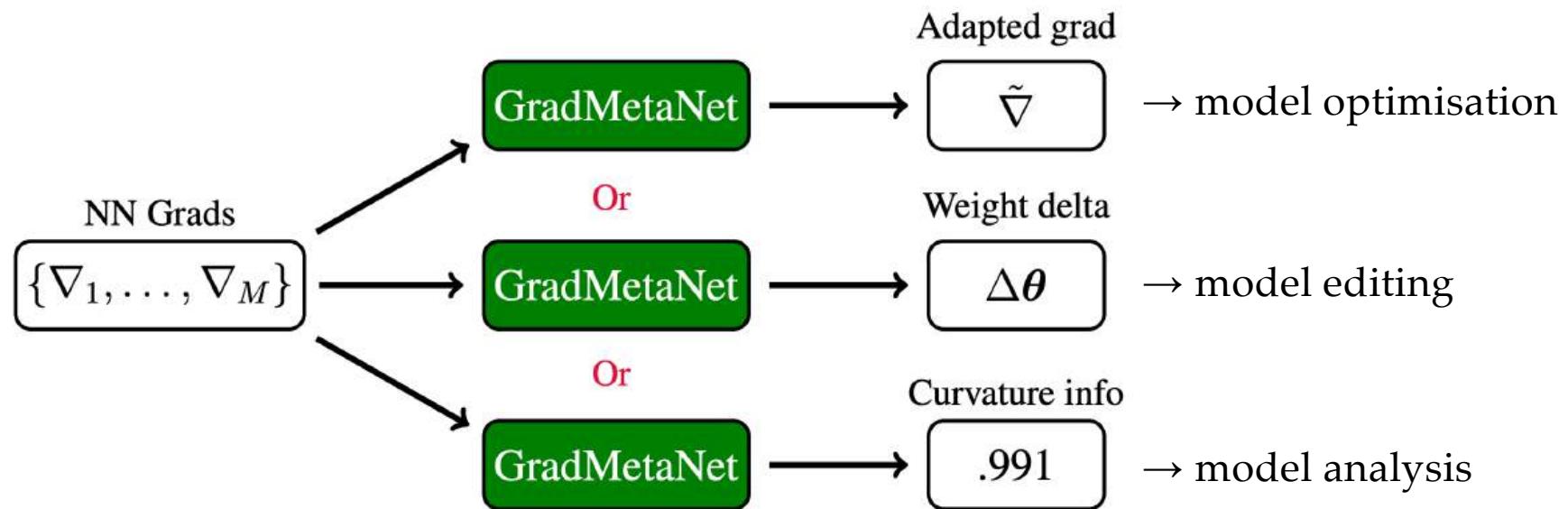


Symmetries in the Gradient Space

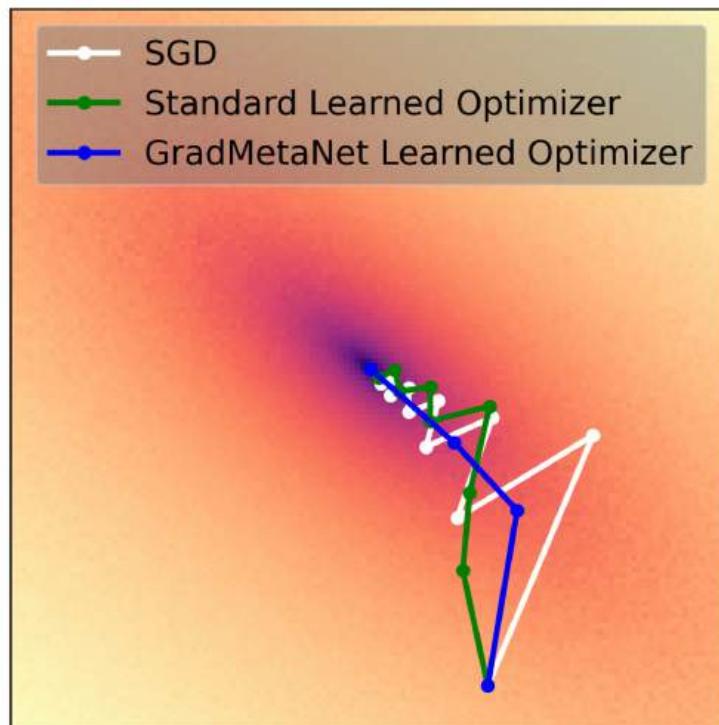
$$\mathcal{L}_{\mathcal{B}} = \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}} \ell(f(\mathbf{x}, \boldsymbol{\theta}), \mathbf{y})$$



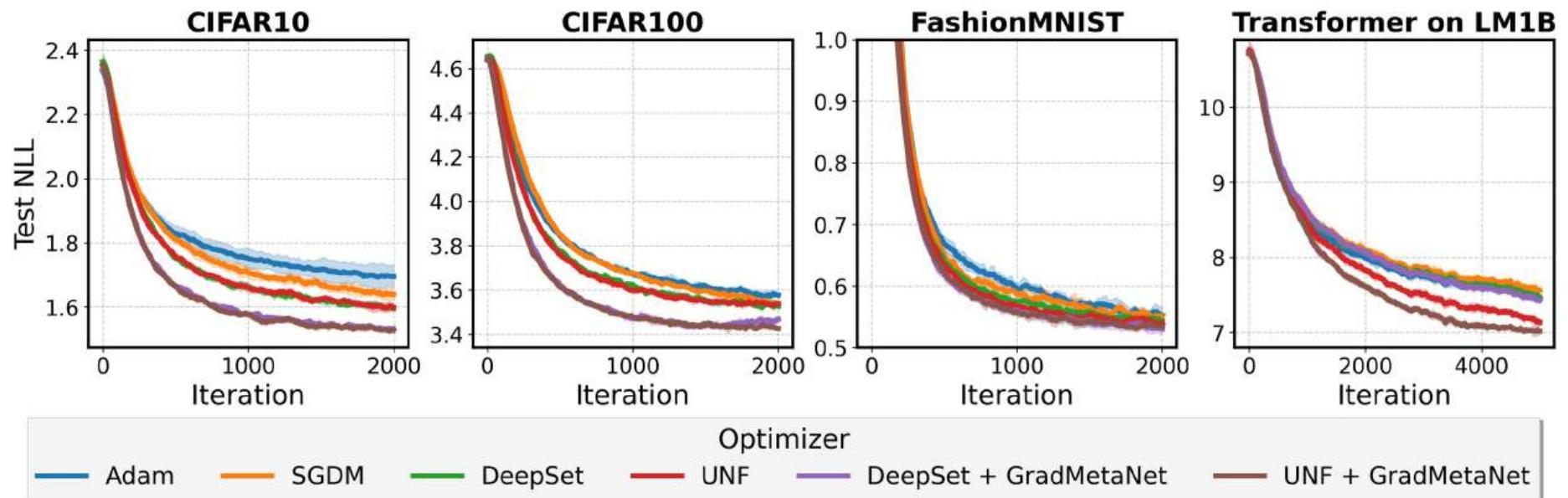
GradMetaNet

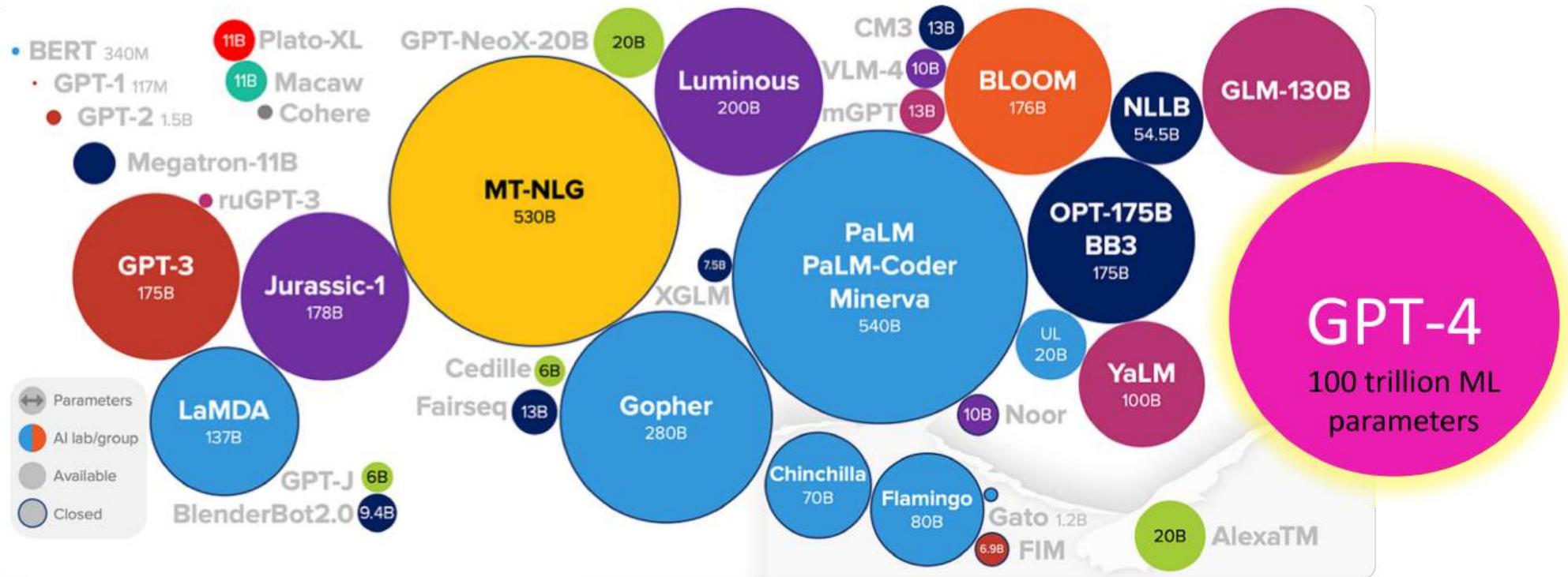


GradMetaNet: Learning Optimisation

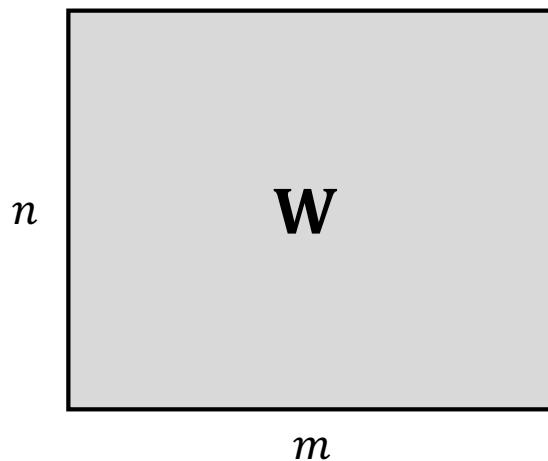


GradMetaNet: Learning Optimisation

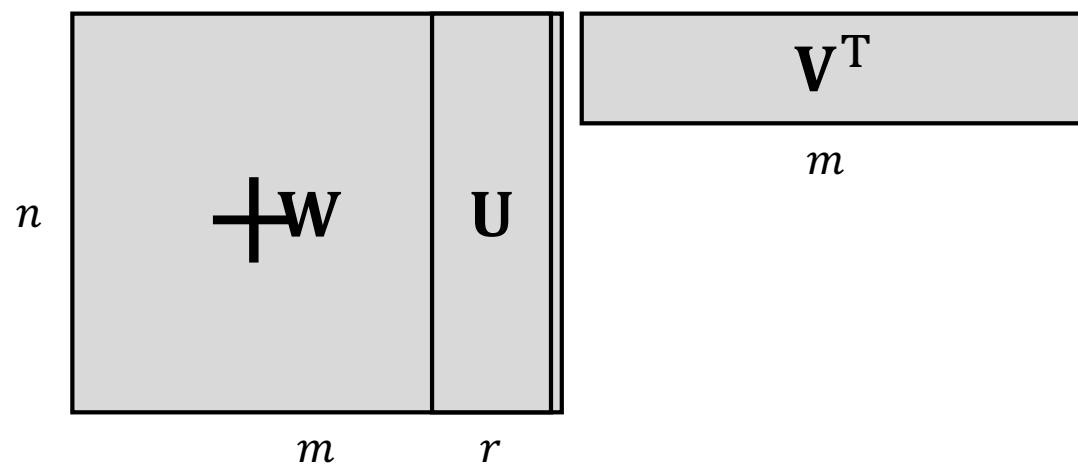




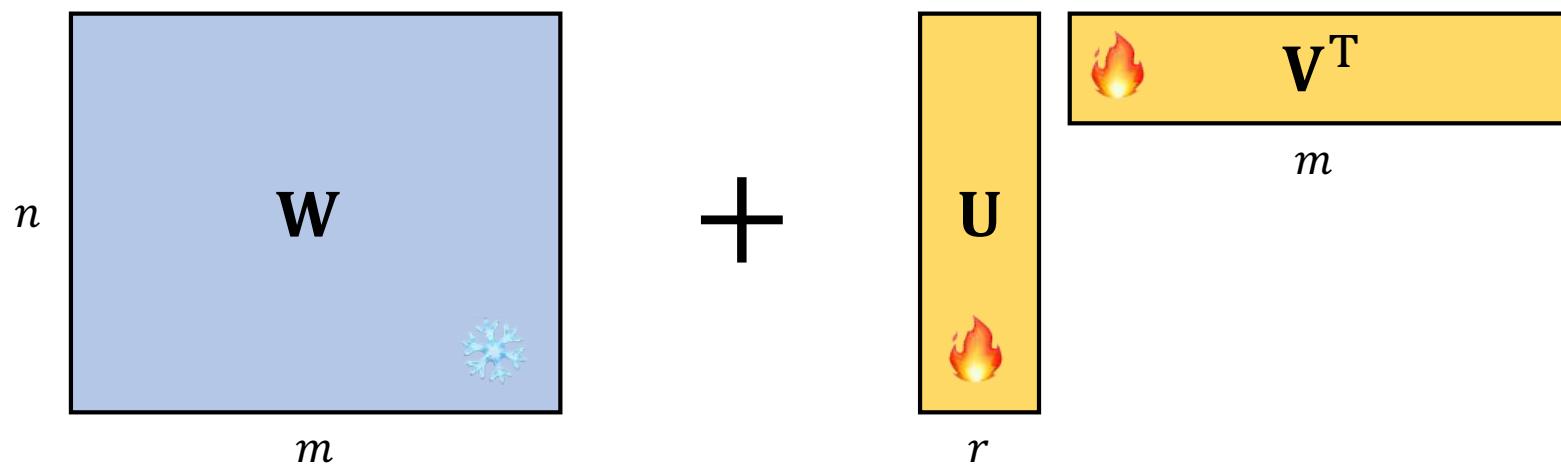
Low-Rank Adaptors (LoRA)



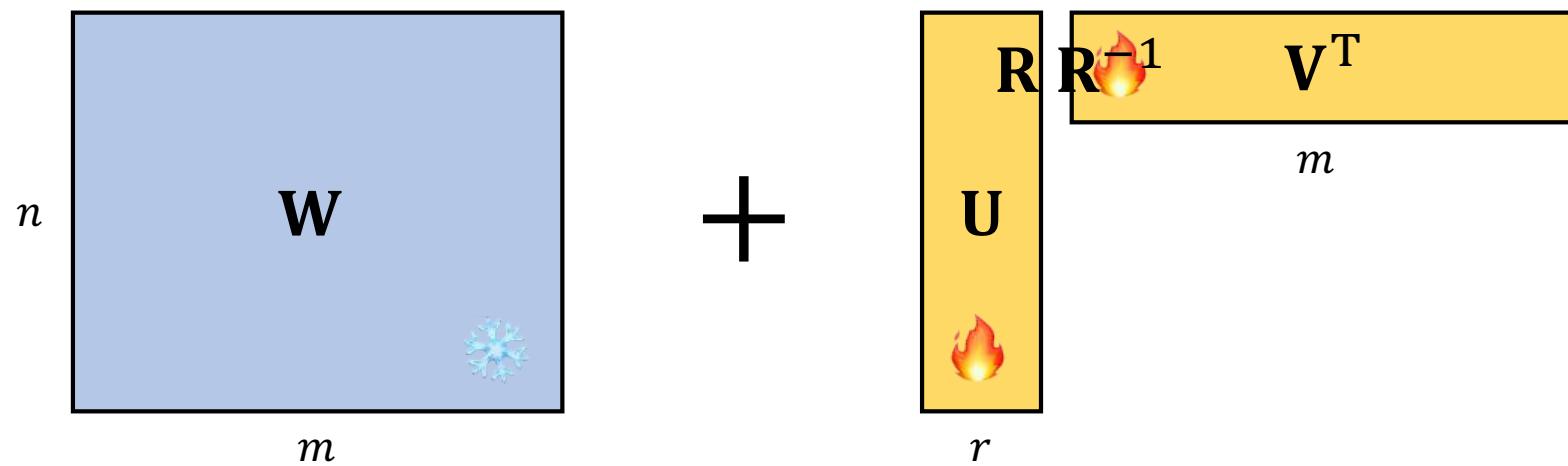
Low-Rank Adaptors (LoRA)



Low-Rank Adaptors (LoRA)

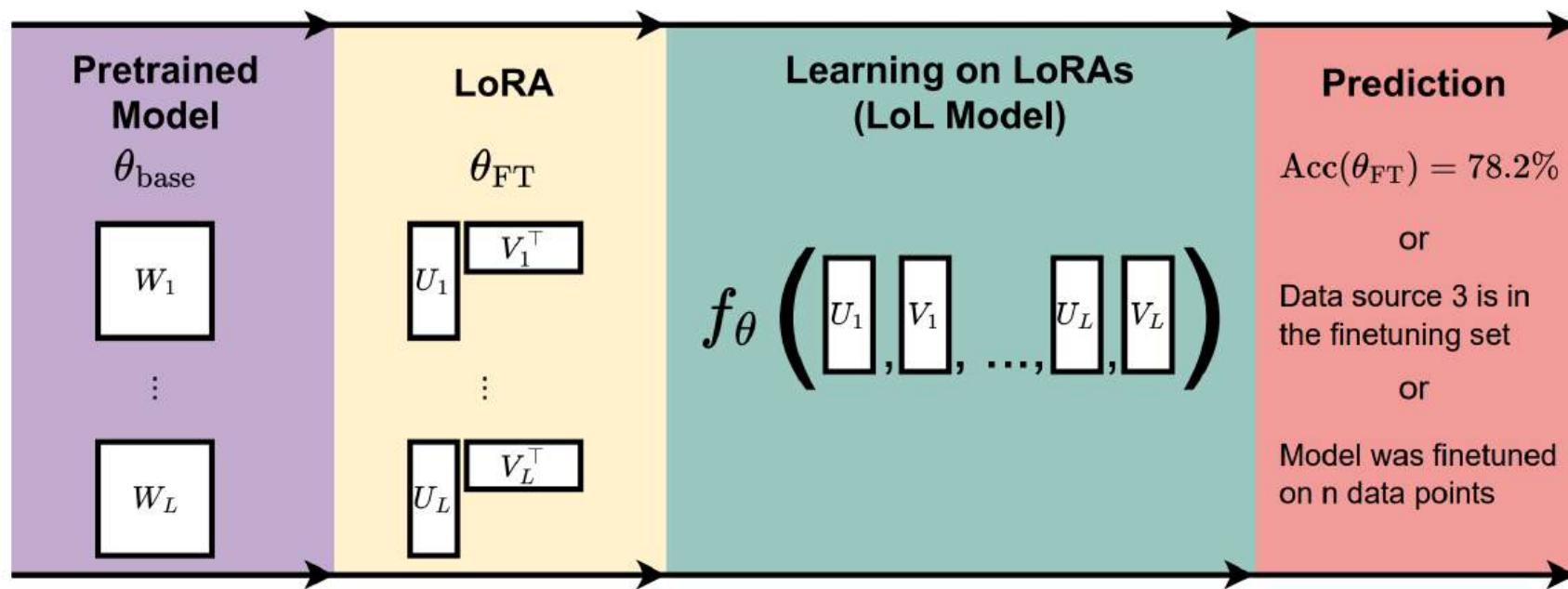


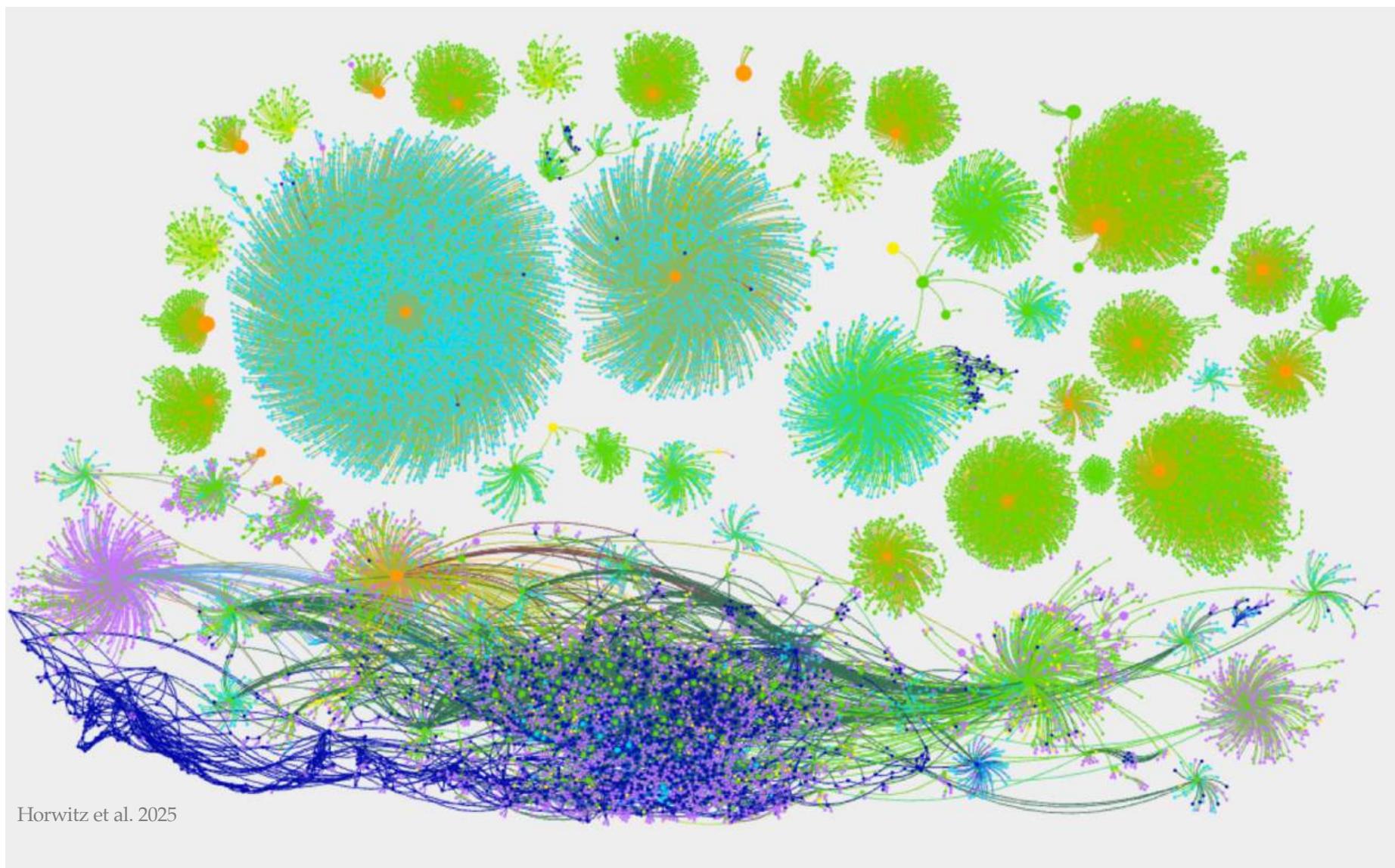
Low-Rank Adaptors (LoRA)



LoRA (\mathbf{U}, \mathbf{V}) defined up to $\text{GL}(r)$

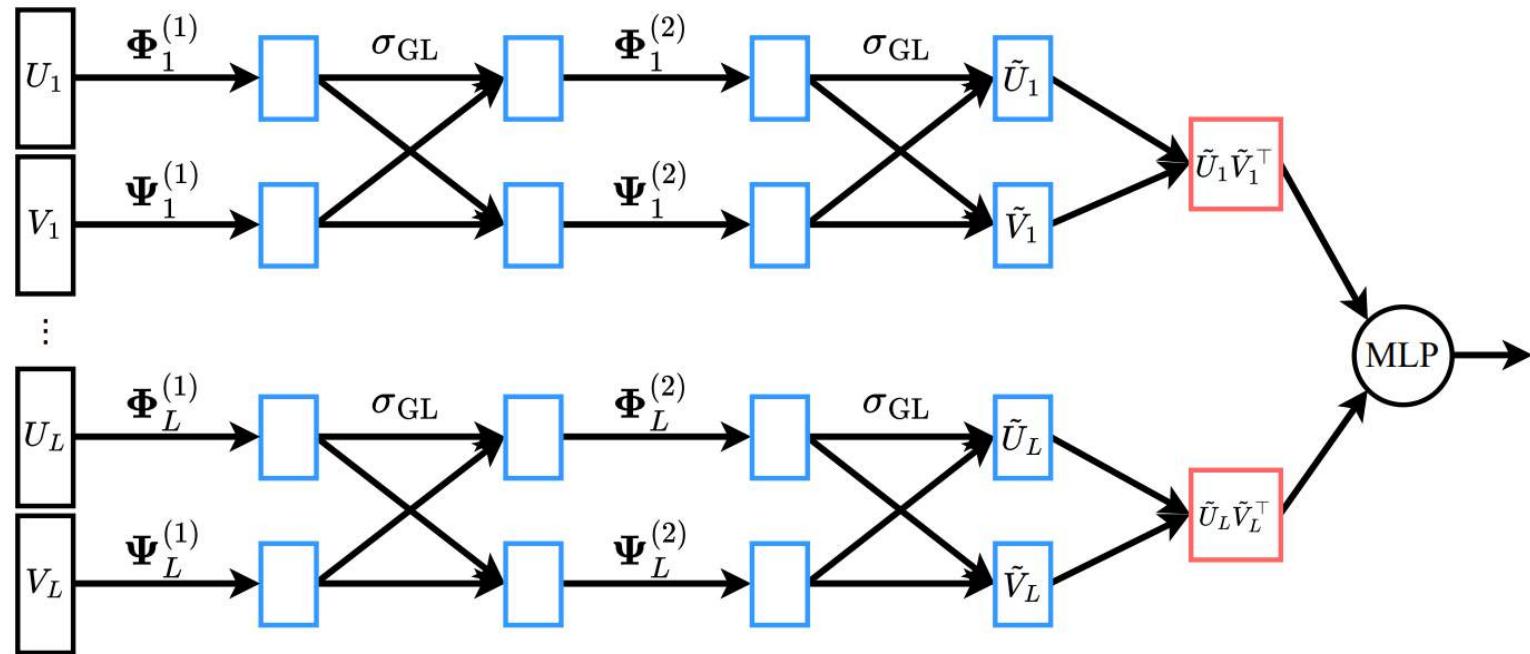
Learning on LoRAs (LoL)





Horwitz et al. 2025

Learning on LoRAs (LOL)



Learning on LoRAs (LOL)

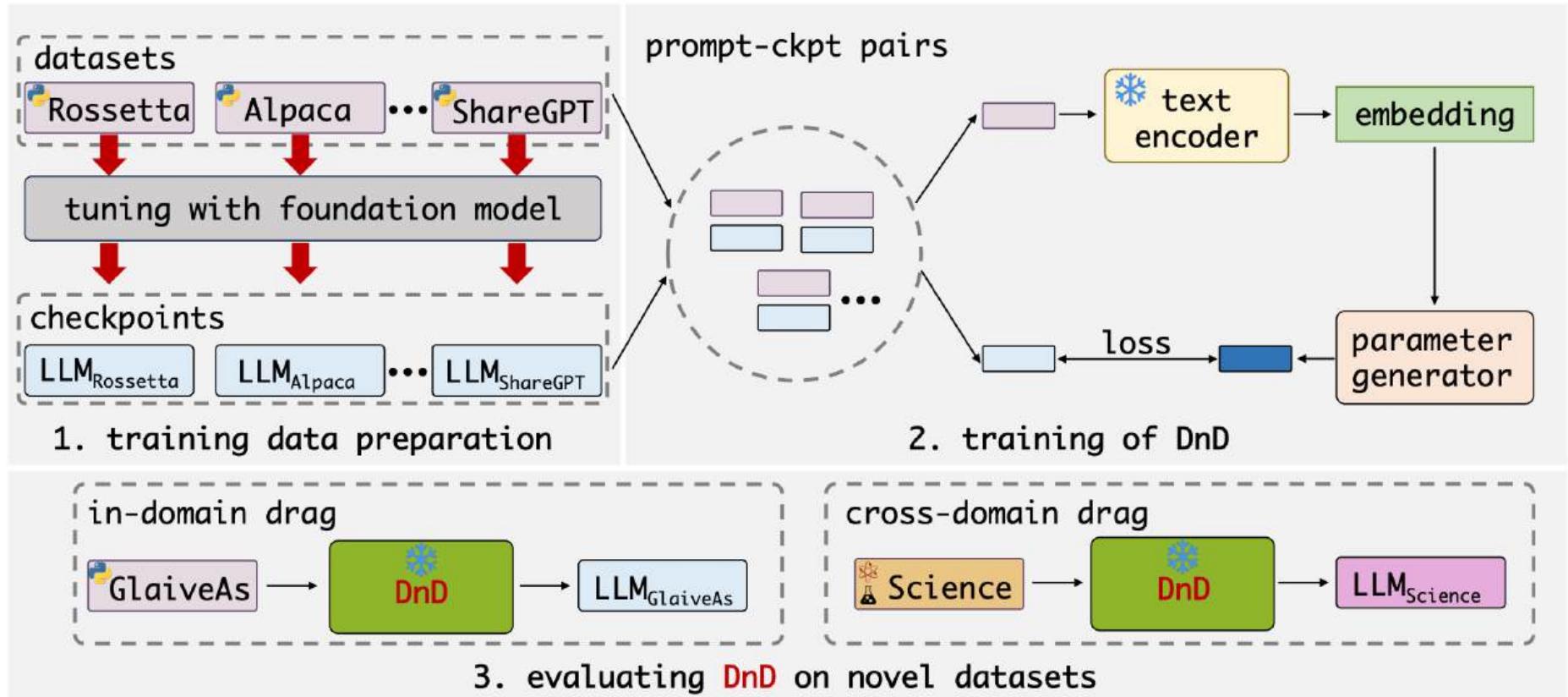
LoL Model	GL-Inv.	O-Inv.	Expressive	Preprocess Time	Forward Time
Naive architectures					
MLP($[U, V]$)	✗	✗	✓	$O((m + n)r)$	$O((m + n)r)$
Transformer($[U, V]$)	✗	✗	✓	$O((m + n)r)$	$O((m + n)r)$
MLP(UV^\top)	✓	✓	✓	$O(mnr)$	$O(mn)^2$
Efficient symmetry-aware architectures					
MLP(O-Align($[U, V]$))	✗	✓	✓	$O((m + n)r^2)$	$O((m + n)r)$
MLP($\sigma(UV^\top)$)	✓	✓	✗	$O((m + n)r^2)$	$O((m + n)r)$
GL-net	✓	✓	✓	$O((m + n)r)$	$O((m + n)r)$

Learning on LoRAs (LoL)

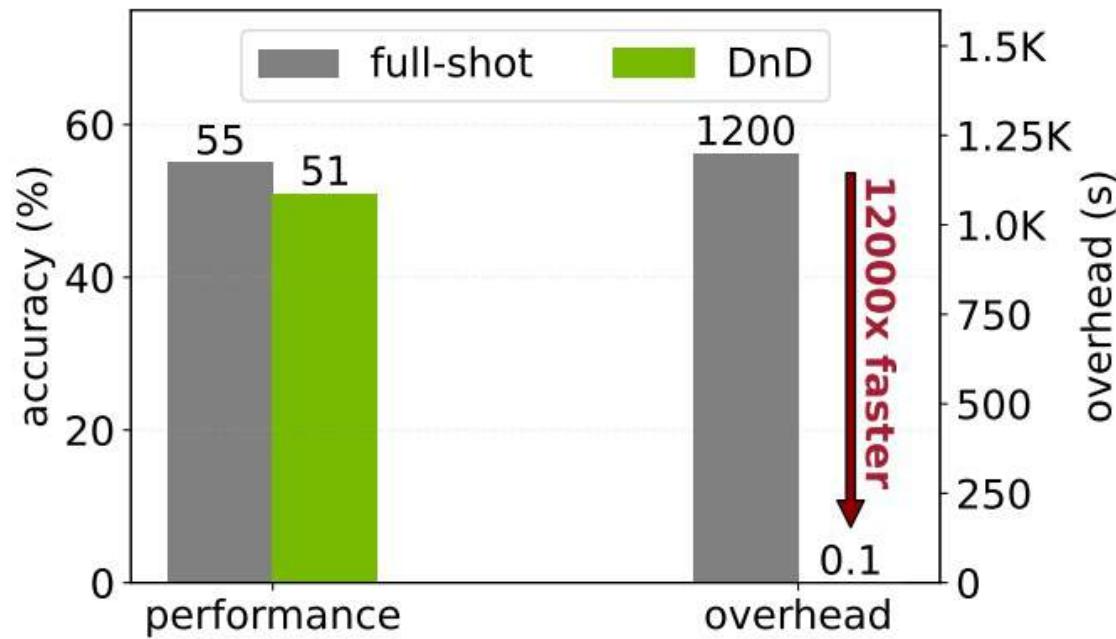
LoL Model	CelebA Attributes		Imagenette Classes		
	Test Loss (↓)	Test Acc (↑)	Test Loss (↓)	Test Acc (↑)	
Naive Models	MLP($[U, V]$)	.554 ± .000	72.4 ± 0.0	.709 ± .004	49.6 ± 1.3
	Transformer($[U, V]$)	.586 ± .014	73.2 ± 0.9	.695 ± .001	50.0 ± 1.3
	MLP(UV^\top)	.267 ± .007	89.1 ± 0.4	.264 ± .011	88.9 ± 0.6
Efficient Invariant	MLP(O-Align($[U, V]$))	.333 ± .008	87.2 ± 0.5	.278 ± .008	87.8 ± 0.3
	MLP($\sigma(UV^\top)$)	.509 ± .013	77.3 ± 1.3	.638 ± .013	65.6 ± 0.6
	GL-net	.232 ± .007	91.3 ± 0.1	.244 ± .005	90.4 ± 0.3

Using LoL models to predict *CelebA* attributes (left) and *Imagenette* classes (right) of the finetuning data of diffusion models, given only the LoRA weights.

Drag-and-Drop LLMs



Drag-and-Drop LLMs



THE STORY CONTINUES
IN THE WEIGHT SPACE



ICLR
International Conference On
Learning Representations

First Workshop on Weight Space Learning

<https://weight-space-learning.github.io/>

Invited Speaker



Stella X. Yu

University of Michigan



Michael Mahoney

UC Berkeley, ICSI



Boris Knyazev

Samsung AI Lab (SAIT)



Naomi Saphra

Harvard University



Ludwig Schmit

Stanford University / Anthropic

Organization



**Konstantin
Schürholt**



**Giorgos
Bouritsas**



**Eliahu
Horwitz**



**Derek
Lim**



**Yoav
Gelberg**



**Bo
Zhao**



**Allan
Zhou**



**Damian
Borth**



**Stefanie
Jegelka**

Steering



**Michael
Bronstein**



**Gal
Chechik**



**Stella
X. Yu**



**Haggai
Maron**

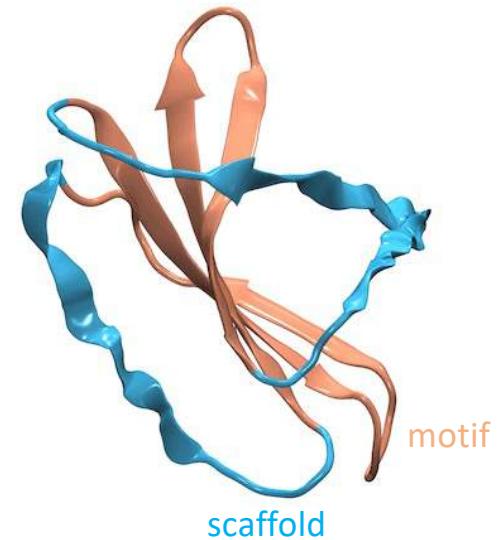
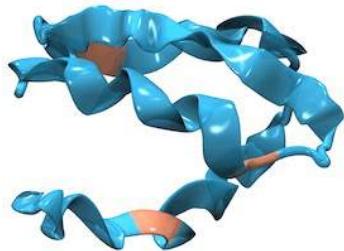
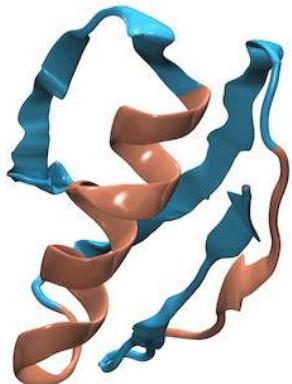


**Yedid
Hoshen**



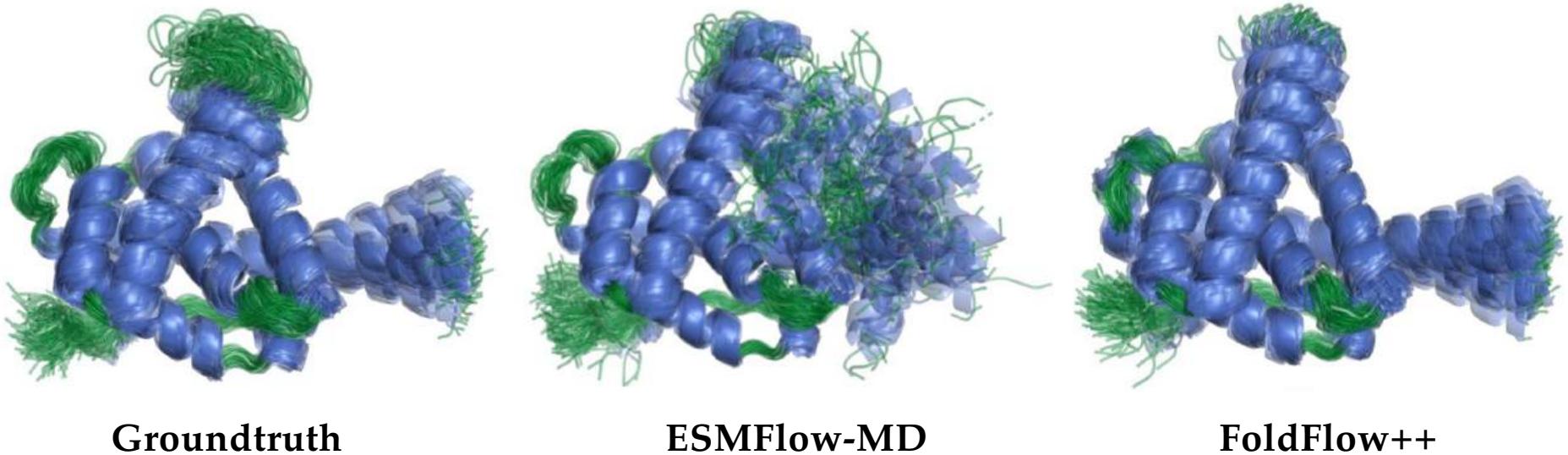
Thank you!

Motif-Scaffolding



Designing novel proteins where functional part (**motif**) is known and
the remaining part (**scaffold**) is generated

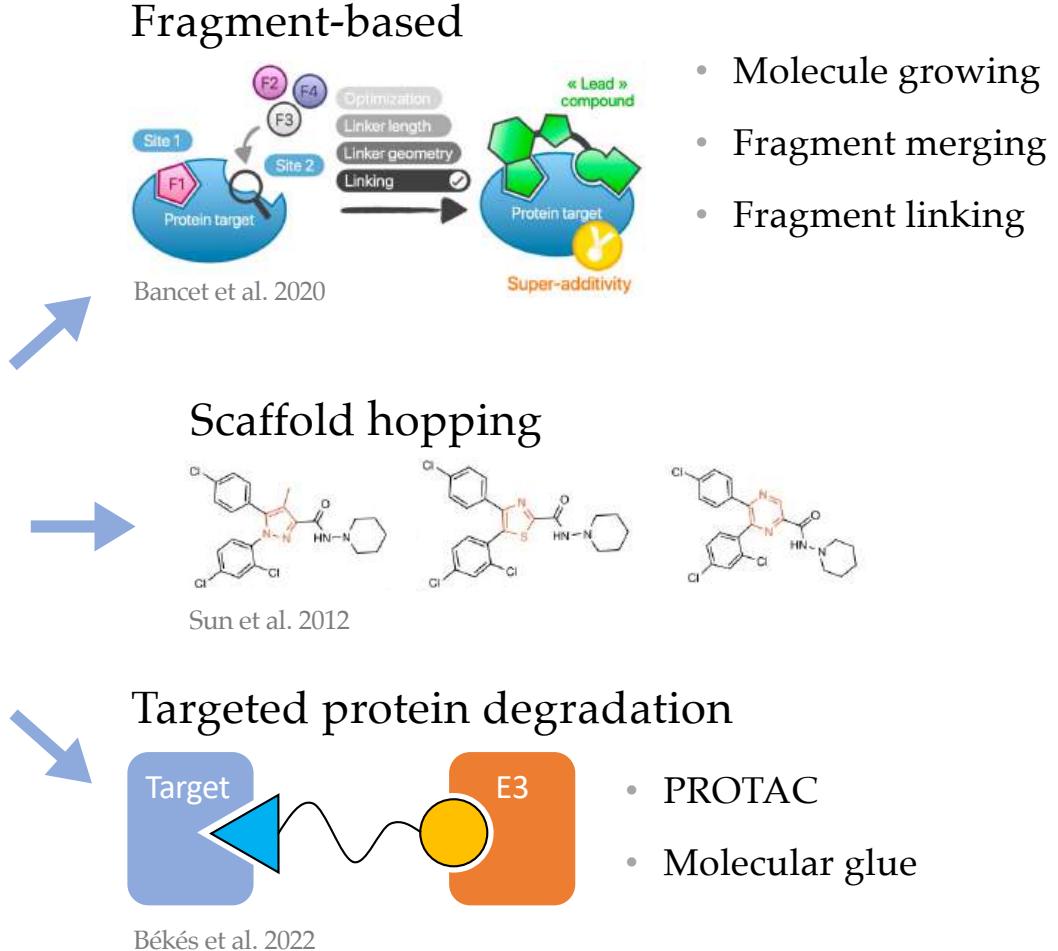
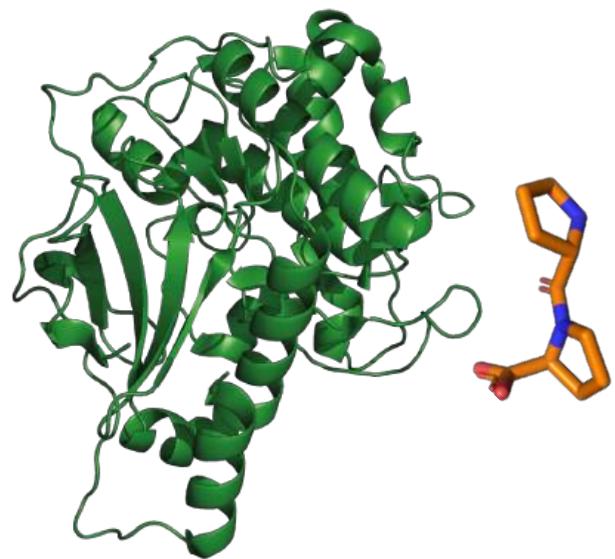
Zero-shot Molecular Dynamics

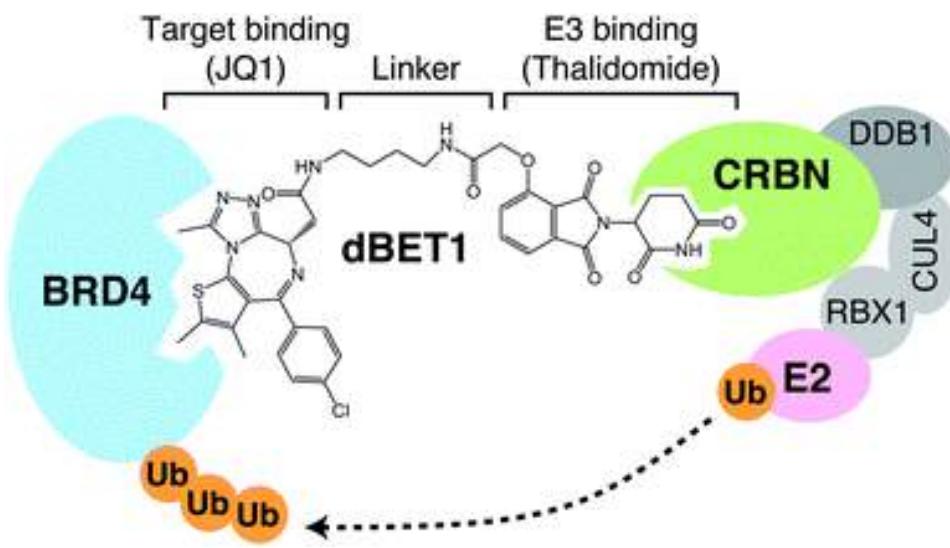


Protein conformations generated by different models

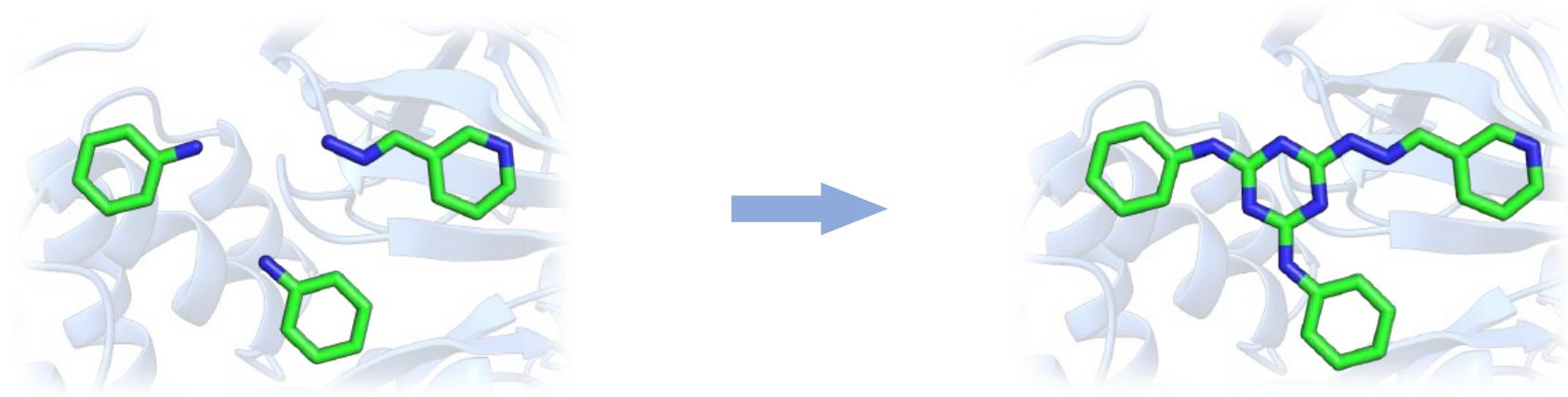
Jing et al. 2024 (ESMFlow); Huguet et Bose, Tong, B 2024 (FoldFlow++)





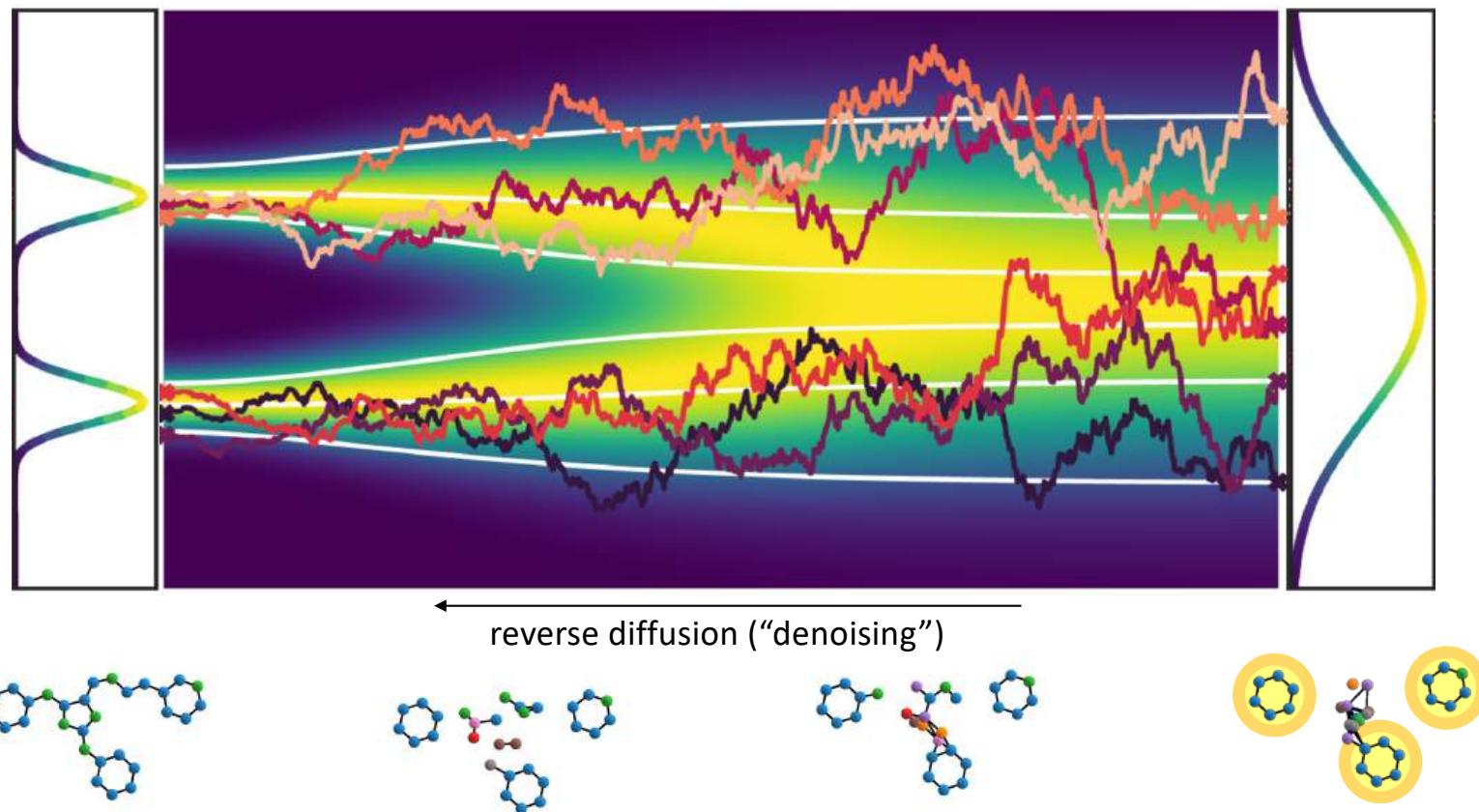


Problem setup: “molecular impainting”



- Fragments placed in 3D space (known fixed orientation)
- Variable linker size
- Unknown attachment atoms (anchors)
- Variable number of fragments
- No clashes with protein pocket

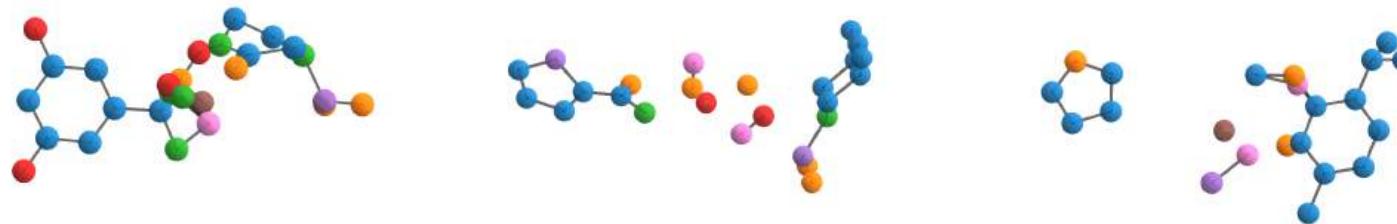
DiffLinker



Igashov, Stärk, Vignac et Welling, B, Correia 2022

Linking pairs of fragments

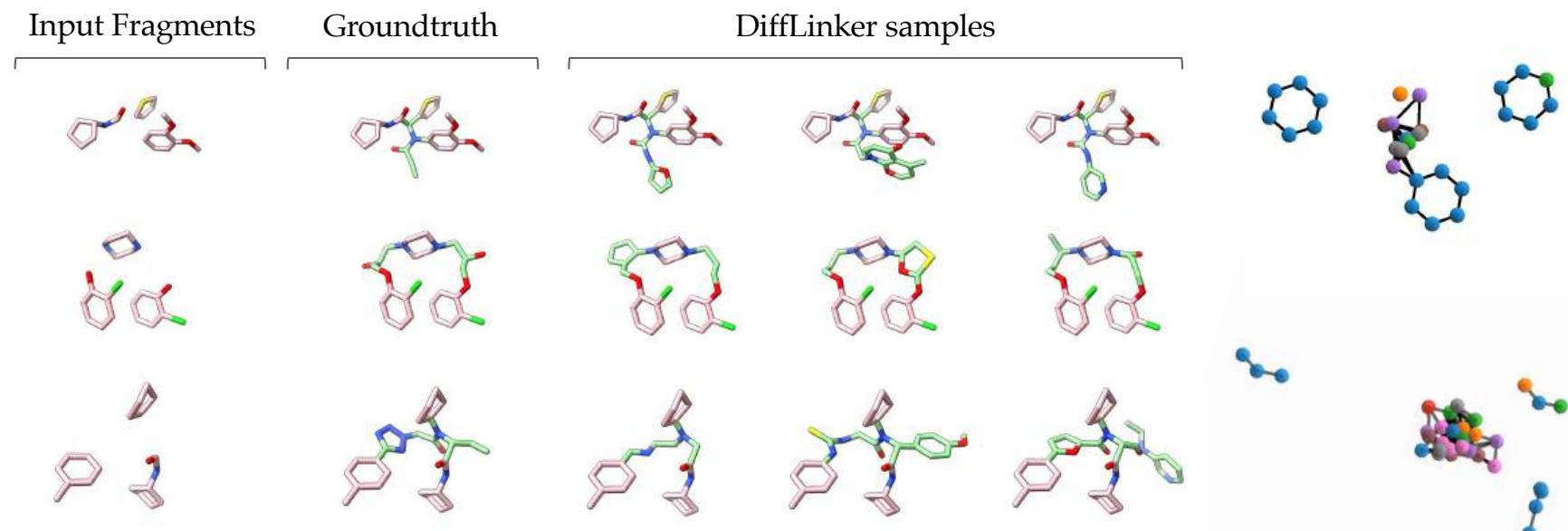
	Method	QED ↑	SA ↓	# Rings ↑	Valid, %	Unique, %	Novel, %
ZINC	DeLinker + ConfVAE + MMFF	0.64	3.11	0.21	98.3	44.2	47.1
	3DLinker (given anchors)	0.65	3.11	0.23	99.3	29.0	41.2
	3DLinker	0.65	3.14	0.24	71.5	29.2	41.9
	DiffLinker	0.68	3.01	0.25	93.8	24.0	30.3
	DiffLinker (given anchors)	0.68	3.03	0.26	97.6	22.7	32.4
	DiffLinker (sampled size)	0.65	3.19	0.32	90.6	51.4	42.9
CASF	DiffLinker (given anchors, sampled size)	0.65	3.24	0.36	94.8	50.9	47.7
	DeLinker + ConfVAE + MMFF	0.35	4.05	0.35	95.7	51.6	55.6
	DiffLinker	0.41	4.00	0.34	85.3	40.5	41.8
	DiffLinker (given anchors)	0.40	4.03	0.38	90.2	37.3	48.4
	DiffLinker (sampled size)	0.40	4.06	0.30	63.7	60.0	49.3
	DiffLinker (given anchors, sampled size)	0.40	4.10	0.38	68.4	57.1	56.9



Igashov, Stärk, Vignac et Welling, B, Correia 2022

Linking multiple fragments

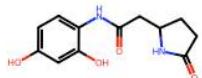
Method	QED ↑	SA ↓	# Rings ↑	Valid, %	Unique, %	Novel, %
GROM	3DLinker	0.36	3.56	0.00	16.3	73.7
	DiffLinker	0.48	2.99	0.75	94.0	34.7
	DiffLinker (given anchors)	0.49	3.01	0.79	94.0	35.0
	DiffLinker (sampled size)	0.45	3.27	0.75	87.2	62.4
	DiffLinker (given anchors, sampled size)	0.46	3.33	0.83	88.8	76.0



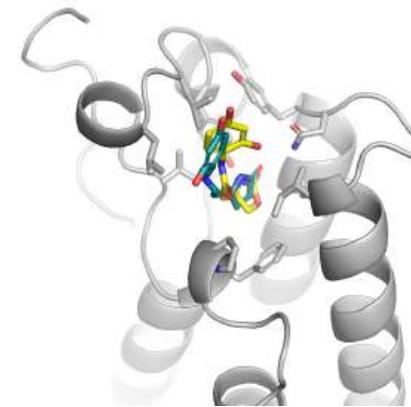
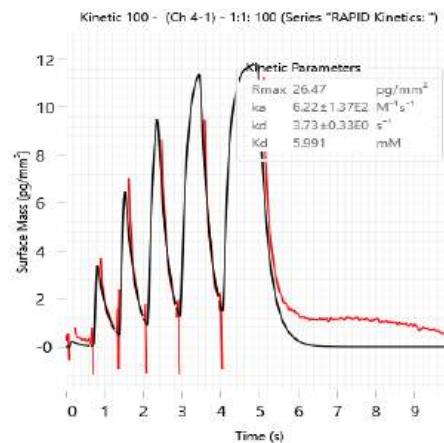
Igashov, Stärk, Vignac et Welling, B, Correia 2022

Experimental results: Binders

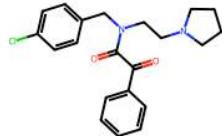
LogP 0.02



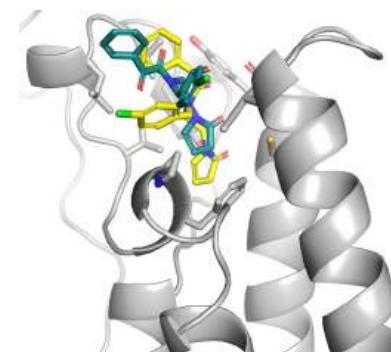
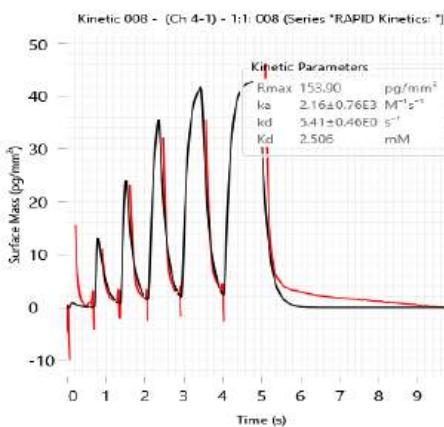
High solubility,
up to 5 mM



LogP 3.88

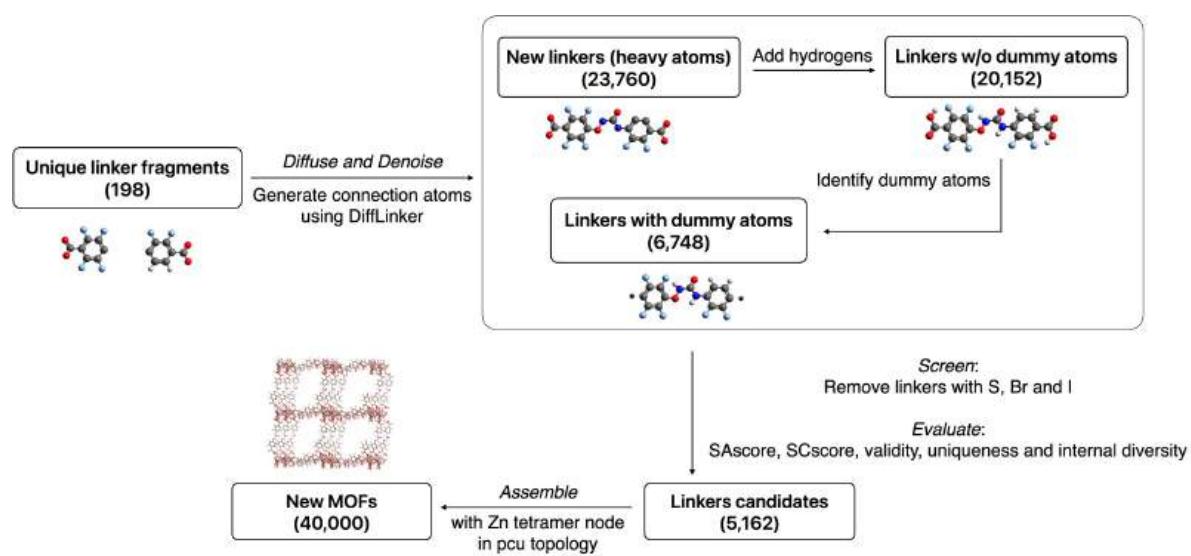


Soluble at 1 mM

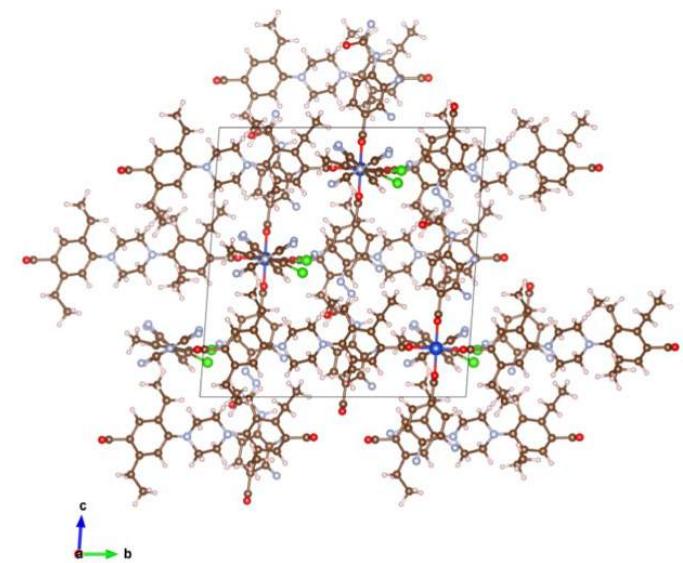


Bromodomain-containing protein 4
(BrD4)

Unexpected Application: Material Design for Carbon capture

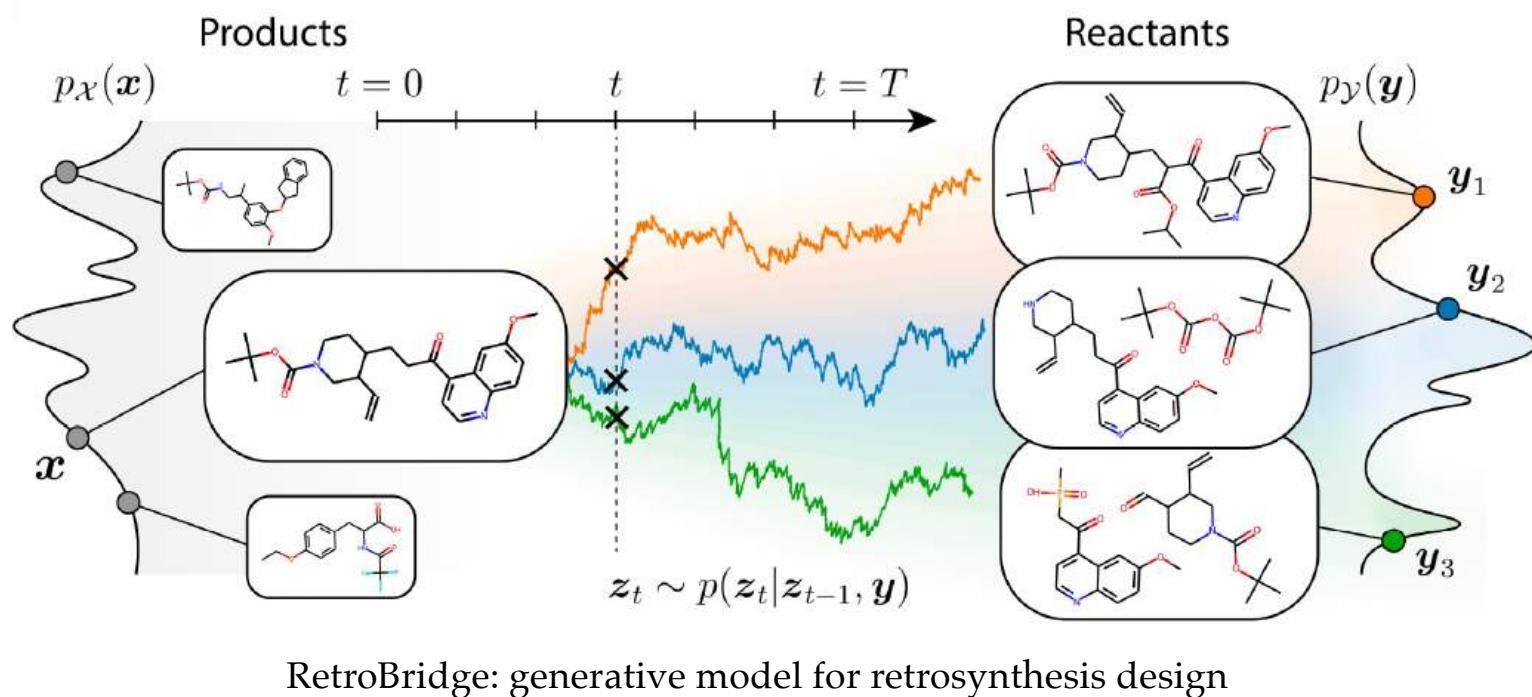


Generation of Metal-Organic Frameworks (MOF) for C capture

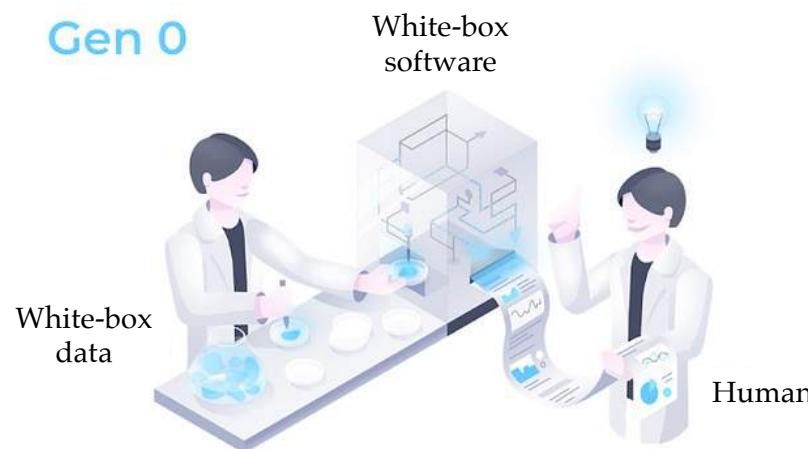


Example of generated MOF

Future directions: Synthesis prediction



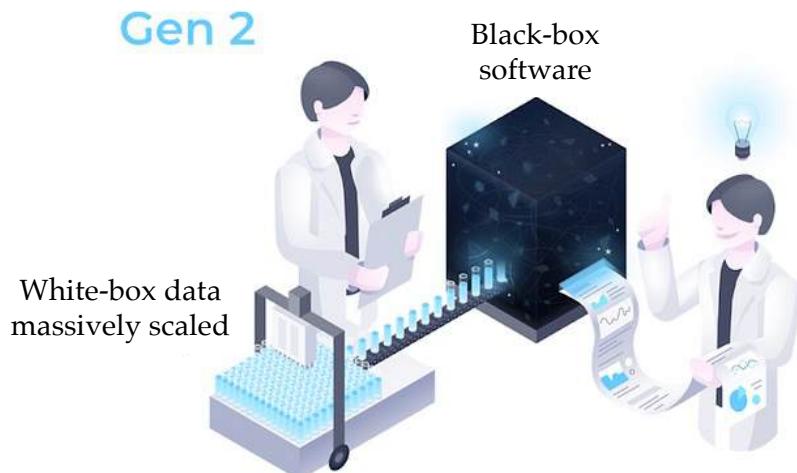
Gen 0



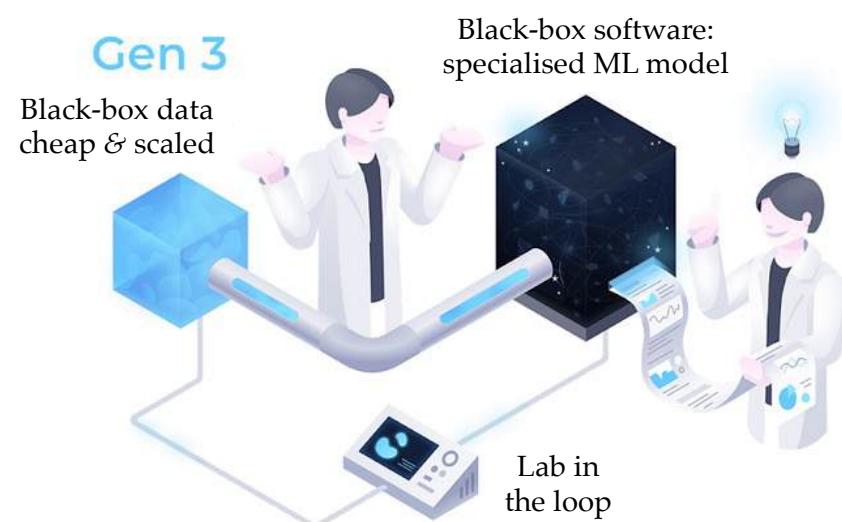
Gen 1



Gen 2



Gen 3



Conclusions

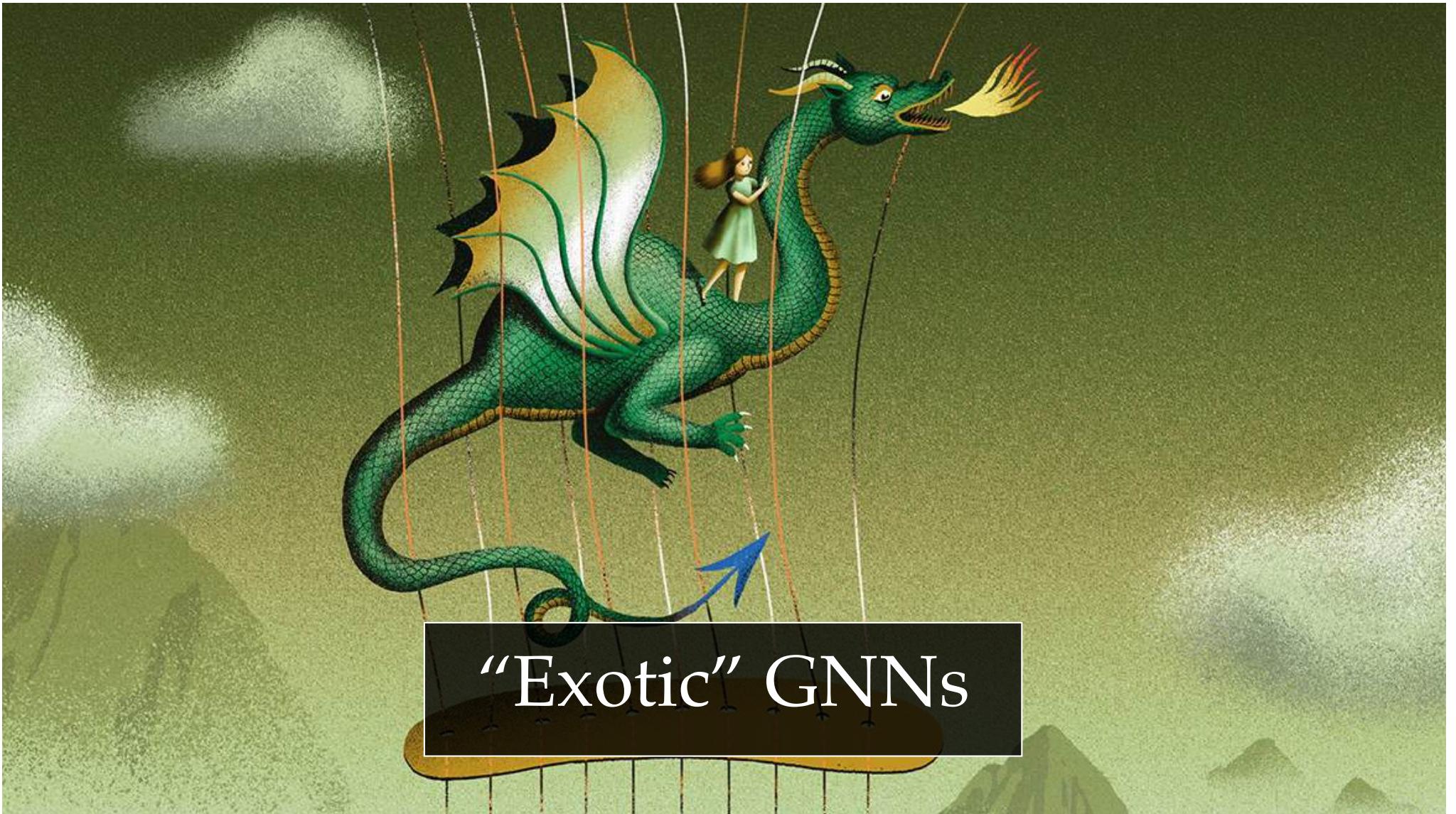
- Geometric deep learning is a powerful tool for molecular modelling
- Roots in fundamental mathematical principles of invariance & symmetry
- Geometric generative models models
- Experimental validation!

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

—Claude Adrien Helvétius

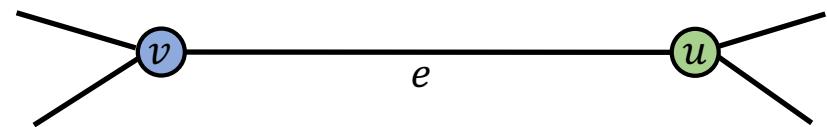


Thank you!

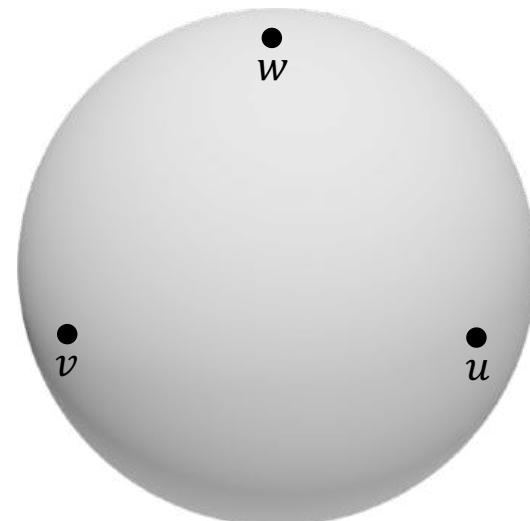


“Exotic” GNNs

Cellular Sheaves

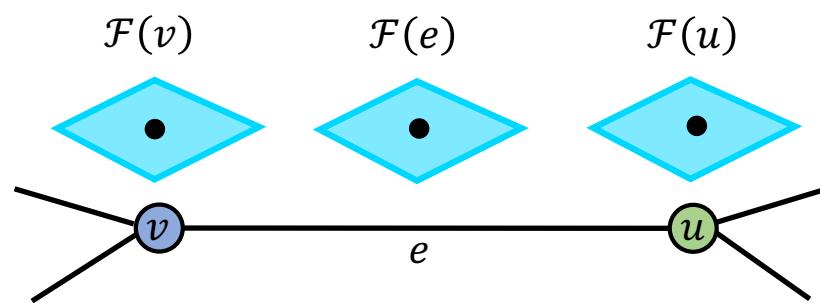


Graph

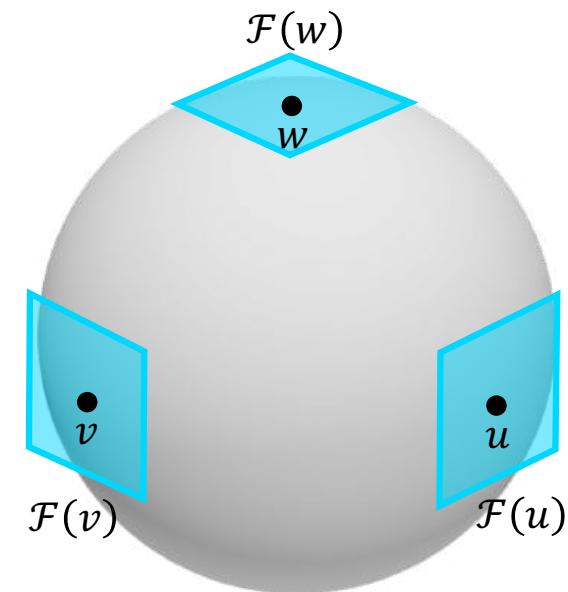


Manifold

Cellular Sheaves

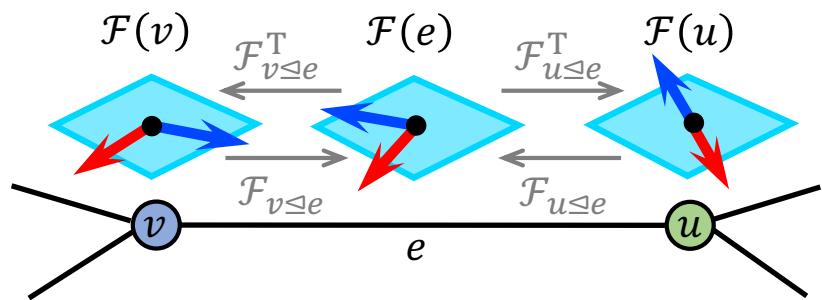


Cellular sheaf \mathcal{F}

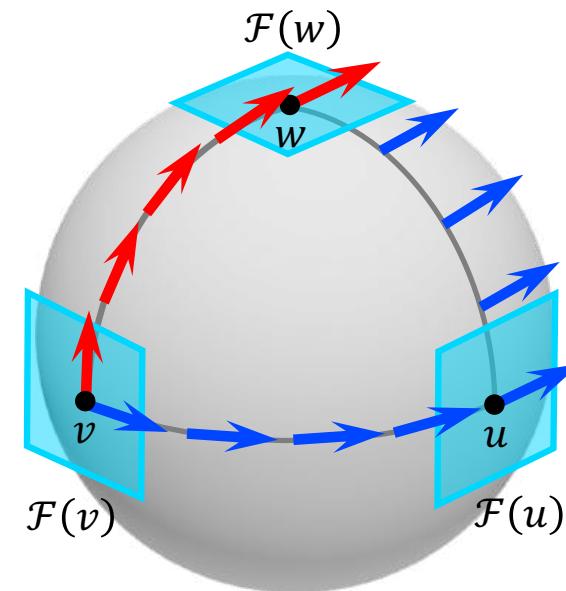


Manifold + Connection

Cellular Sheaves

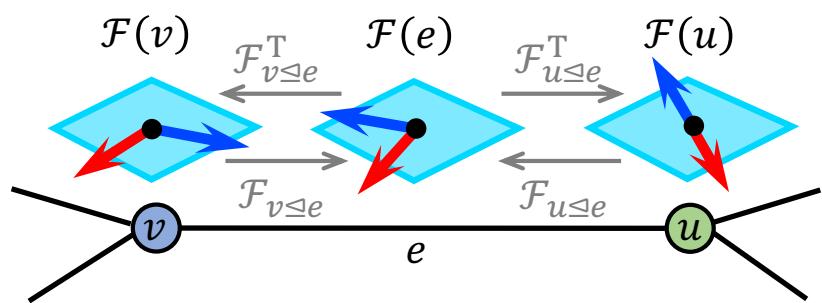


Cellular sheaf \mathcal{F}

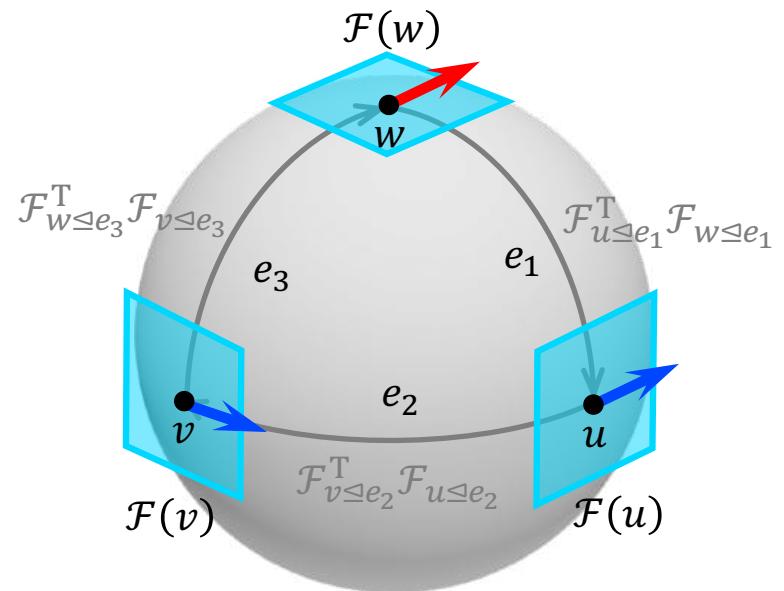


Manifold + Connection

Cellular Sheaves

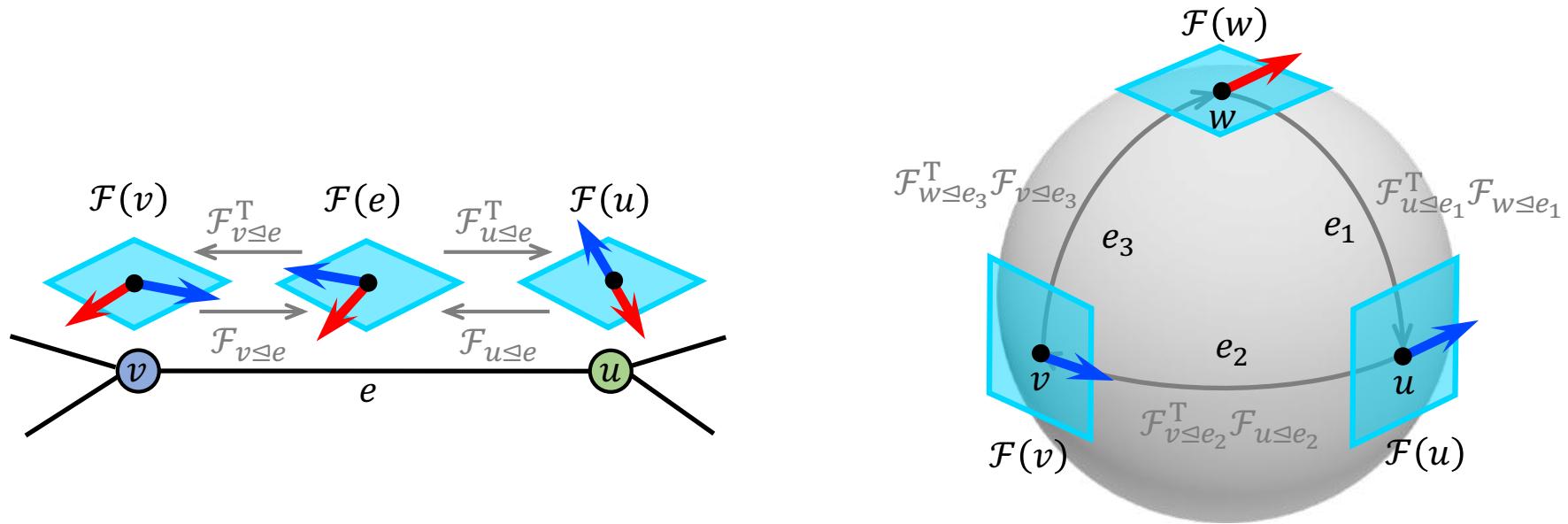


Cellular sheaf \mathcal{F}



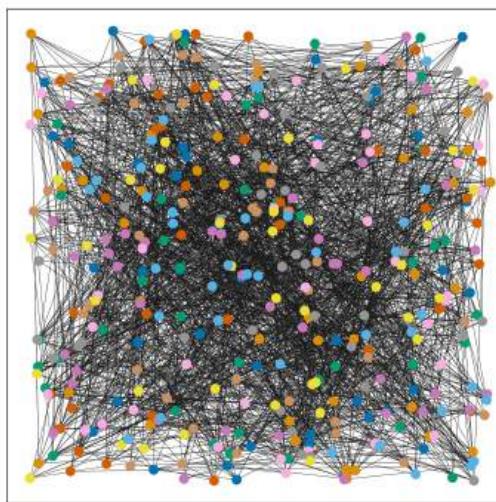
Analogy to parallel transport
on manifolds

Cellular Sheaves



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

Diffusion on Cellular Sheaves



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

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

Alternative to Weisfeiler-Lehman for expressive power?

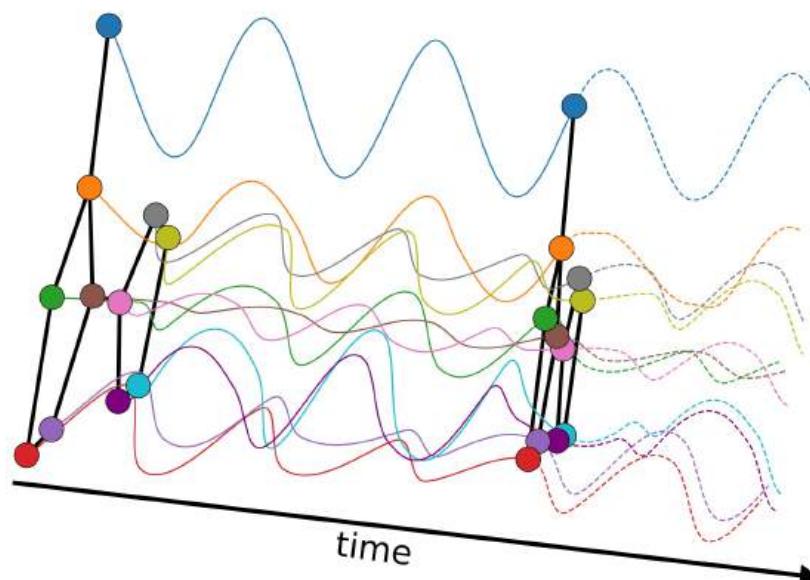
Graph type	# Node classes	Sheaf class \mathcal{F}, dim=d	
Homophilic	2	Symmetric $d=1$	✓
Heterophilic	2	Symmetric $d=1$	✗
	2	Non-symmetric $d=1$	✓
	≥ 3	Non-symmetric $d=1$	✗
	$\leq 2d$	Orthogonal, $d=\{2,4\}$	✓

The capability of sheaf diffusion to solve node classification problem in the limit

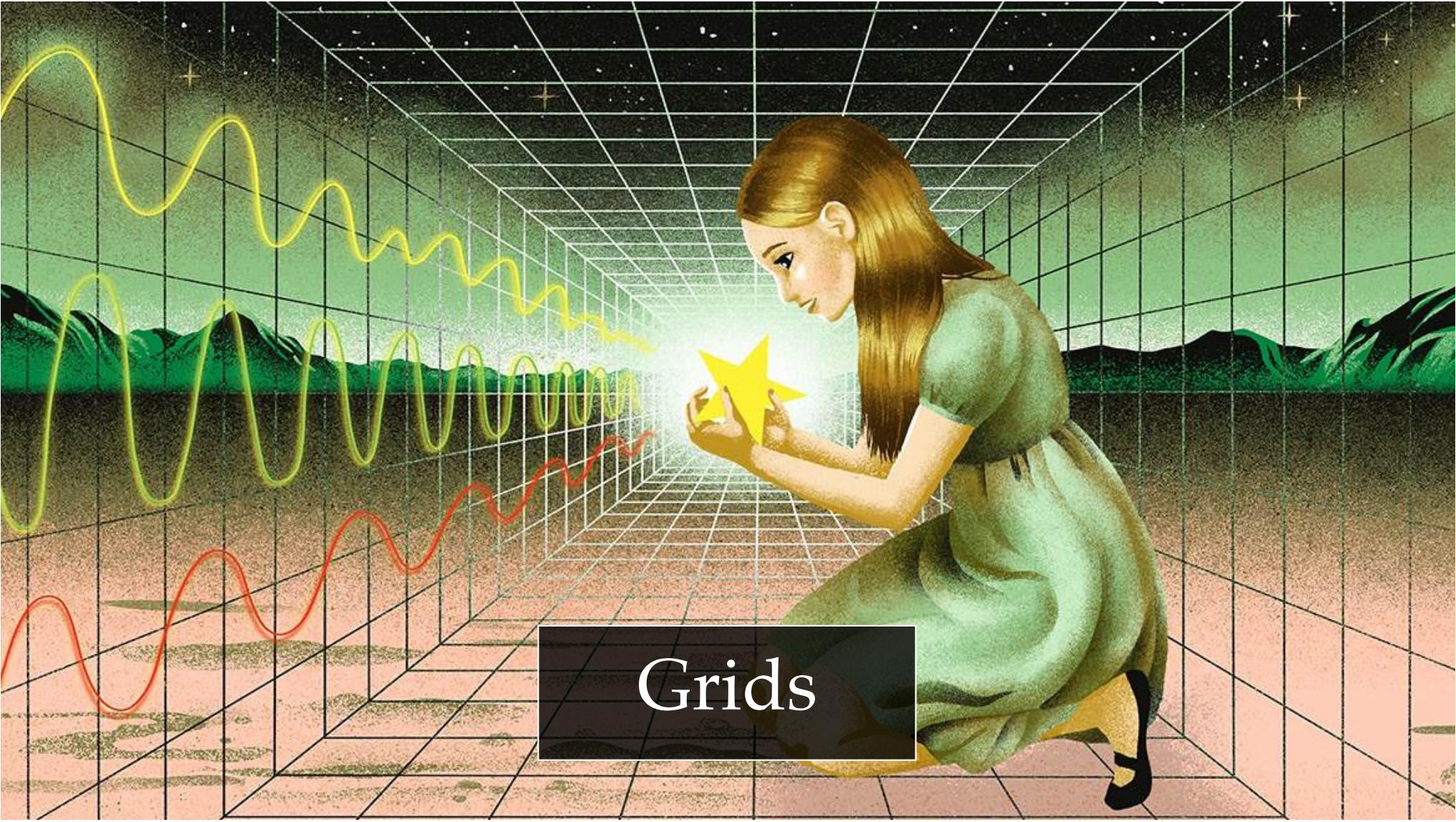
What do we gain from physics-inspired GNNs?

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

Graph-Coupled Oscillators



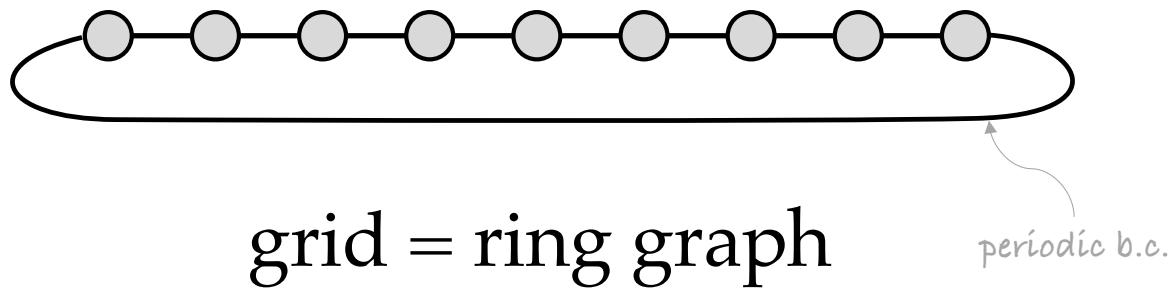
Dynamics of a system of coupled oscillators on a molecular graph



A girl with long blonde hair, wearing a green dress, is kneeling in a 3D grid landscape. She is holding a bright yellow star in her hands. The landscape features rolling hills and mountains in the background under a night sky with stars. A large yellow wavy line and a red wavy line are visible in the foreground.

Grids

Grids vs Graphs

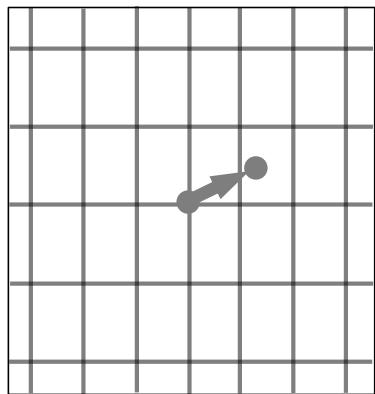
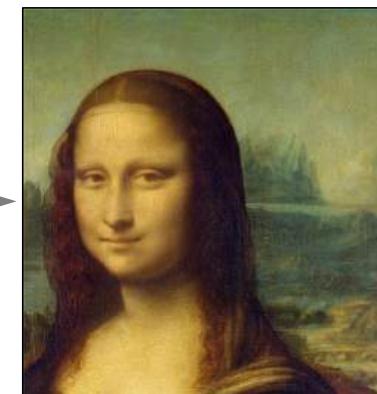


“Lifting”



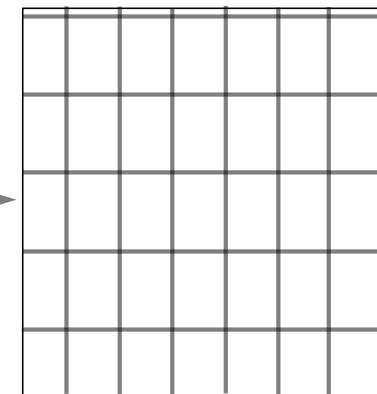
Group representation

$$-\rho(g):\mathcal{X}(\Omega) \rightarrow \mathcal{X}(\Omega)-\rightarrow$$



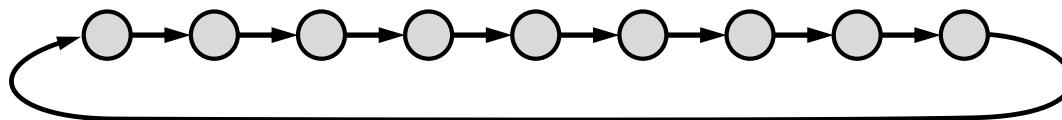
group

$$g:\Omega \rightarrow \Omega$$



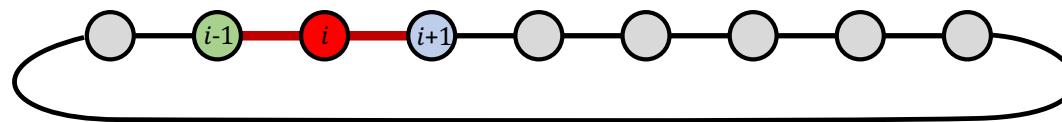
Shift operator

$$\begin{array}{c|c} \begin{array}{c} \text{dark grey} \\ \text{light green} \\ \text{dark green} \\ \text{orange} \\ \text{light green} \end{array} & = \\ \hline & \left[\begin{array}{ccccc} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ 1 & & & & 1 \end{array} \right] \begin{array}{c} \text{light green} \\ \text{dark green} \\ \text{orange} \\ \text{light green} \\ \text{dark grey} \end{array} \end{array}$$



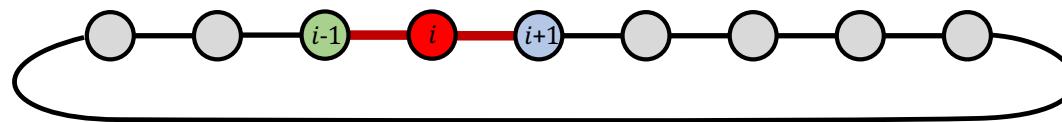
Note: Observe that the adjacency matrix of the directed ring graph is exactly the shift operator. We will use it when defining graph convolutions.

Grids vs Graphs



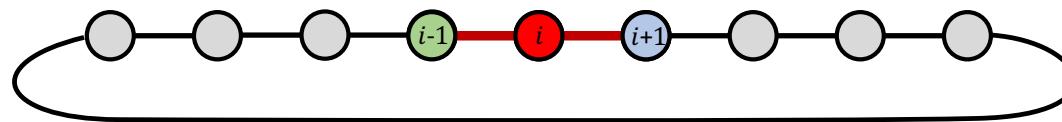
fixed neighbourhood structure

Grids vs Graphs



fixed neighbourhood structure

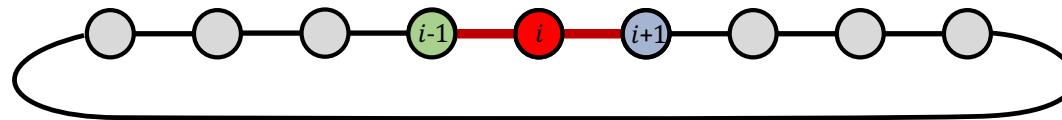
Grids vs Graphs



local aggregation function

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

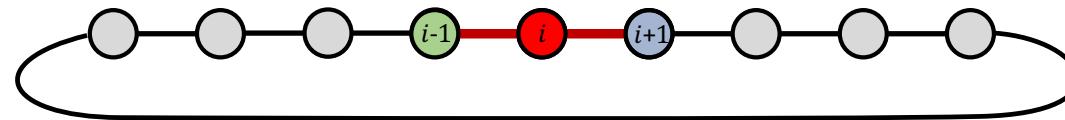
Grids vs Graphs



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

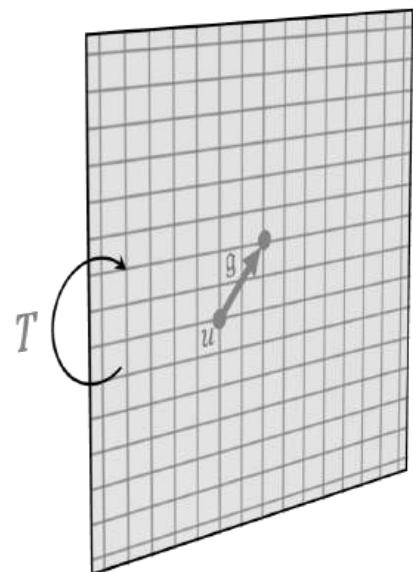
$$\begin{bmatrix} b & c & a \\ a & b & c \\ c & a & b \end{bmatrix} \left[\begin{array}{c} 1 \\ 1 \\ 1 \end{array} \right] = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} b & c & a \\ a & b & c \\ c & a & b \end{bmatrix}$$

convolution \Leftrightarrow shift-equivariant

convolution emerges from translation symmetry

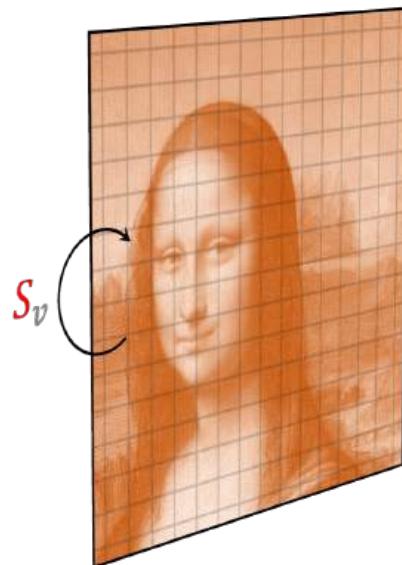
CNNs as an Instance of Geometric Deep Learning Blueprint

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

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

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

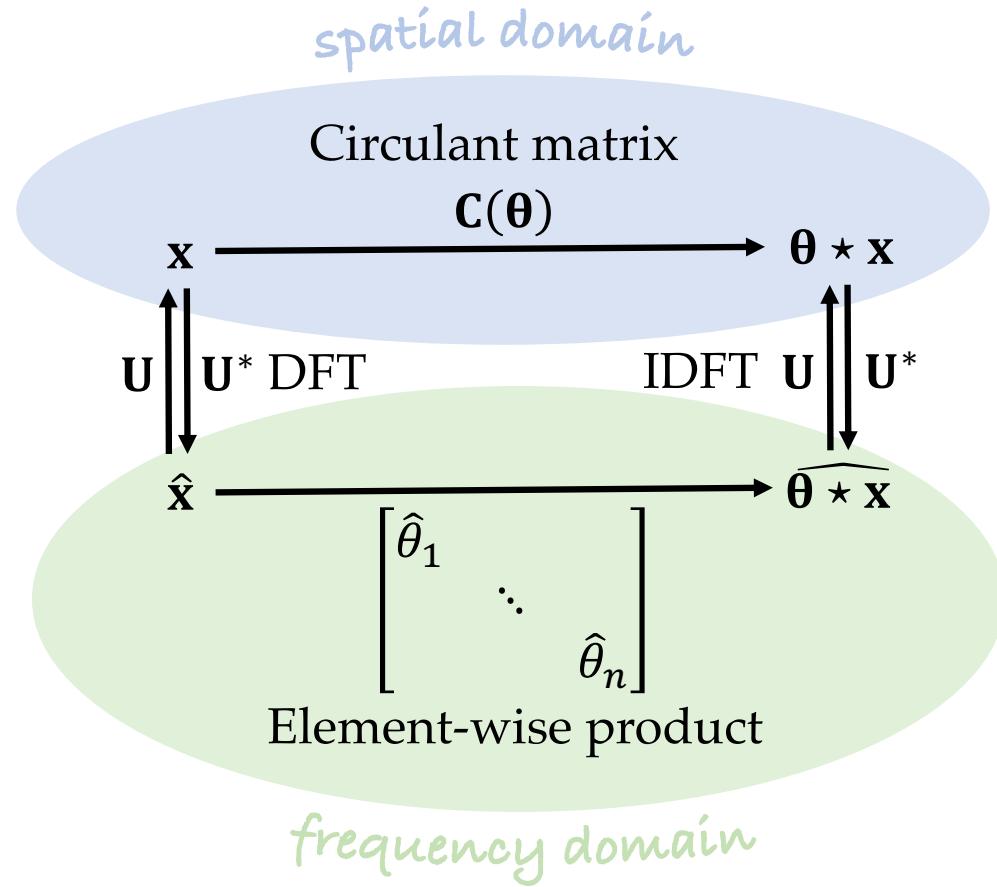
Fourier basis

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

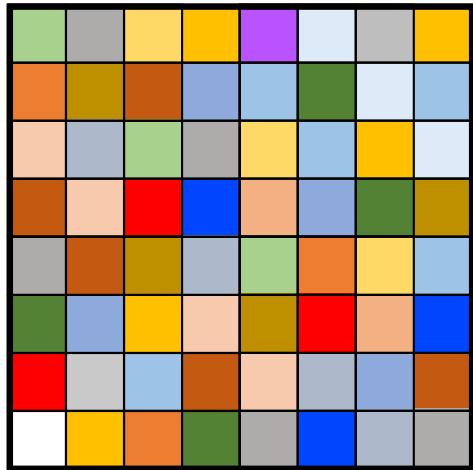
$$\begin{bmatrix} b & c & a & b & c & a \\ a & b & c & a & b & c \\ c & a & b & c & a & b \\ a & b & c & a & b & c \\ b & c & a & b & c & a \\ c & a & b & c & a & b \\ a & b & c & a & b & c \\ b & c & a & b & c & a \end{bmatrix} = \begin{bmatrix} | & & & & | \\ \mathbf{u}_1 & & \dots & & \mathbf{u}_n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

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

commuting matrices are jointly
diagonalisable by Fourier Transform

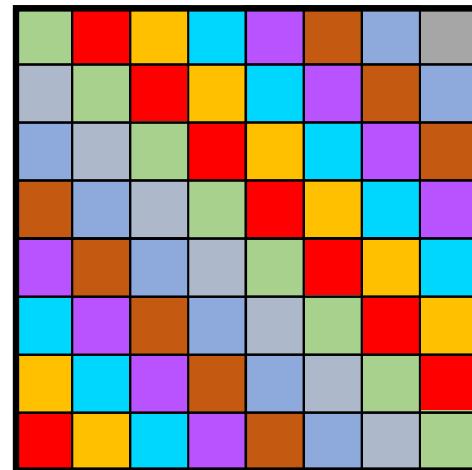


Fully connected



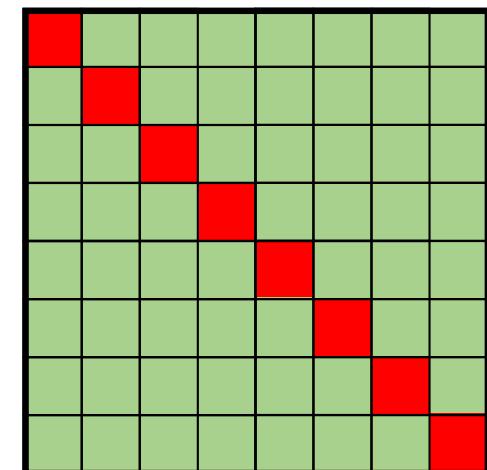
Trivial $G = \{e\}$
 $\mathcal{O}(n^2)$ DOF

Grids



Translation
 $\mathcal{O}(n)$ DOF

Graphs / Sets



Permutation
 $\mathcal{O}(1)$ DOF



Groups

Convolution, revisited



Convolution, revisited

convolution = matching shifted filter

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


A diagram illustrating the components of the convolution formula. On the left, a curved arrow points from the variable u in the expression $(x \star \psi)(u)$ to the term $S_u \psi$ in the integral. On the right, another curved arrow points from the same u to the term $(u - v)$ inside the integral, indicating its role as a shift operator.

domain $\Omega \cong$ symmetry group G
i.e., $\star \psi: \mathcal{X}(\Omega) \rightarrow \mathcal{X}(\Omega)$

Group Convolution

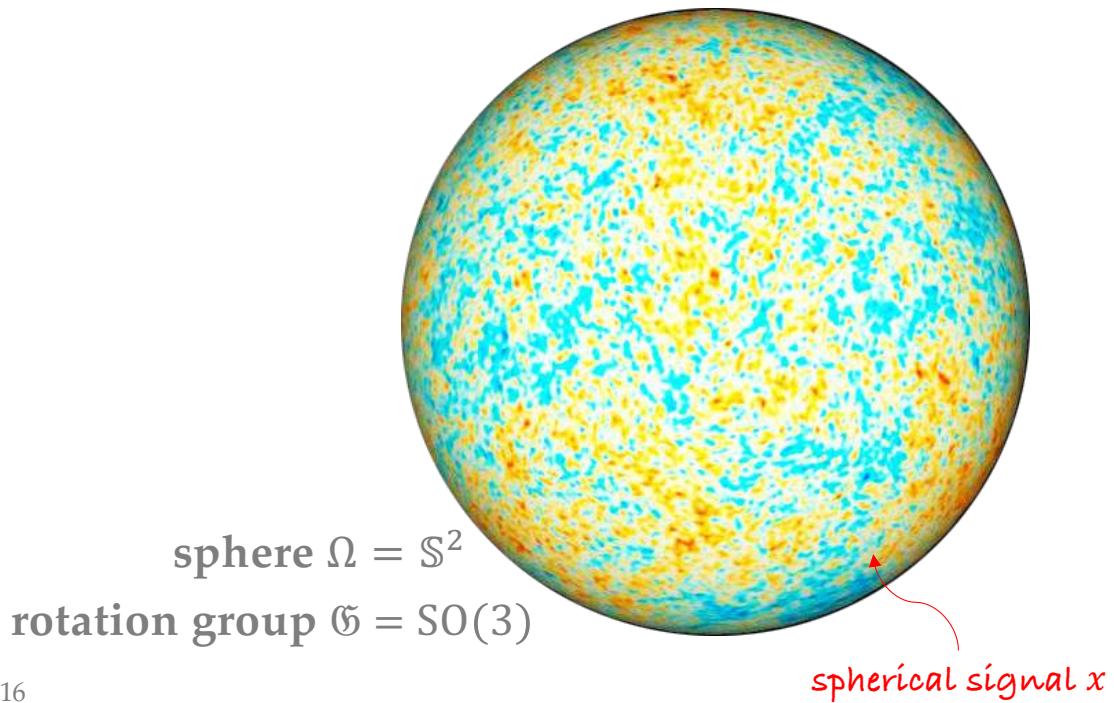
convolution = matching transformed filter

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


group element group representation

The convolution outputs a signal on the group G
i.e., $\star \psi: \mathcal{X}(\Omega) \rightarrow \mathcal{X}(G)$

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 $SO(3)$

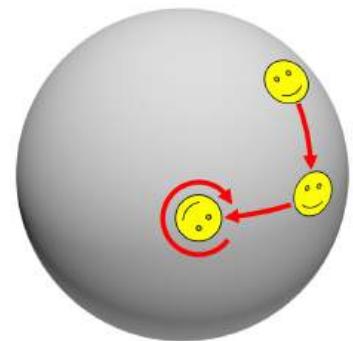
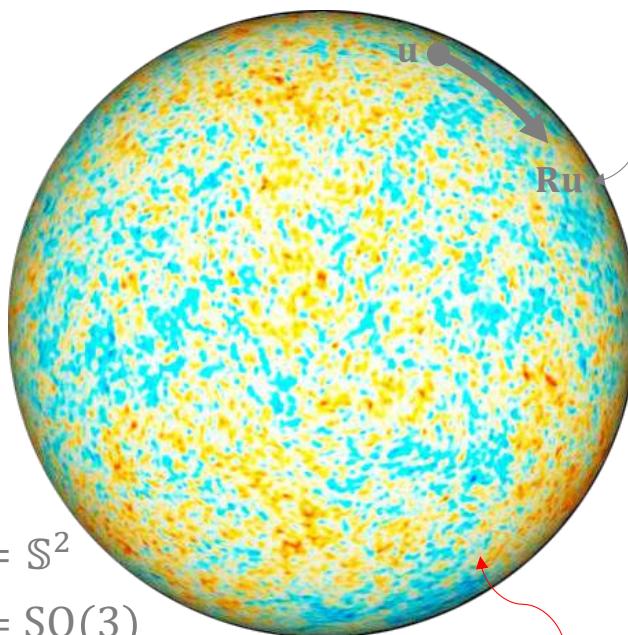
sphere $\Omega = \mathbb{S}^2$

rotation group $G = SO(3)$

spherical signal x

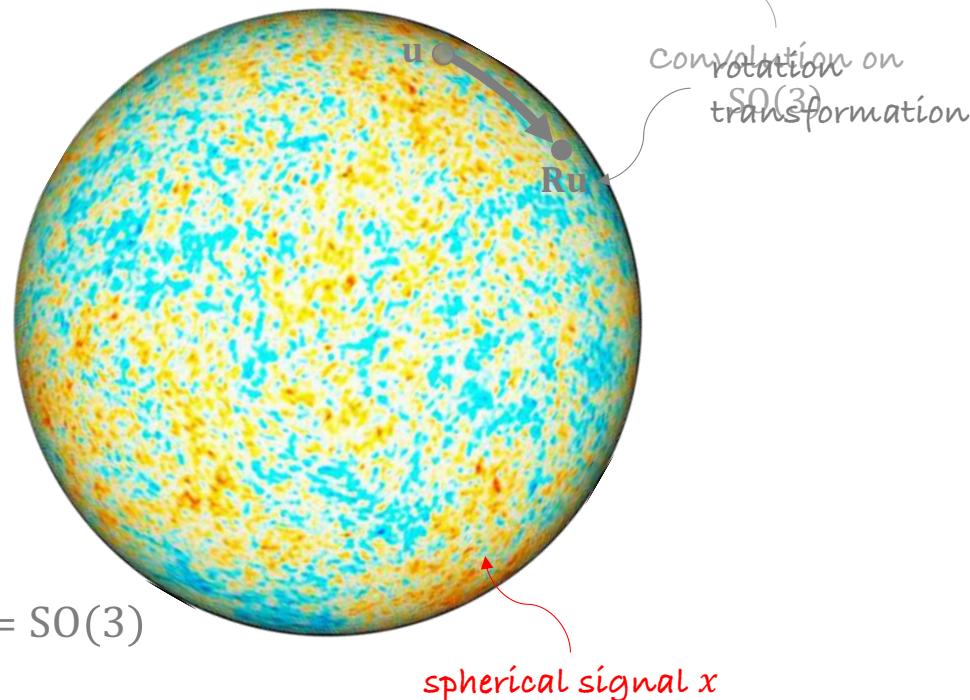
rotation transformation

Cohen, Welling 2016



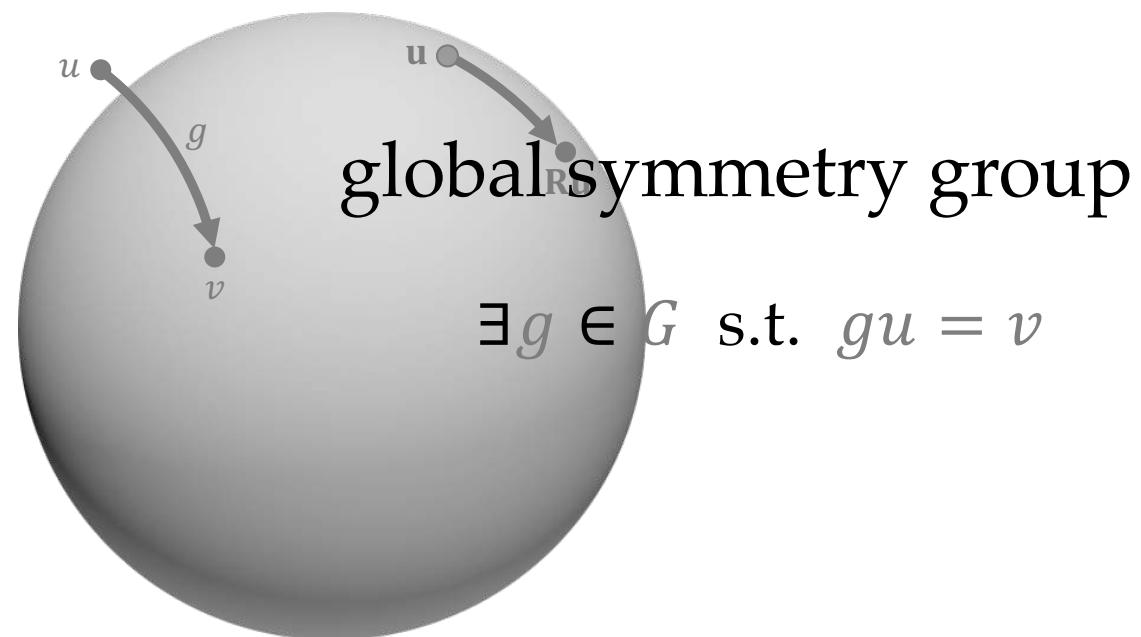
Convolution on the Sphere

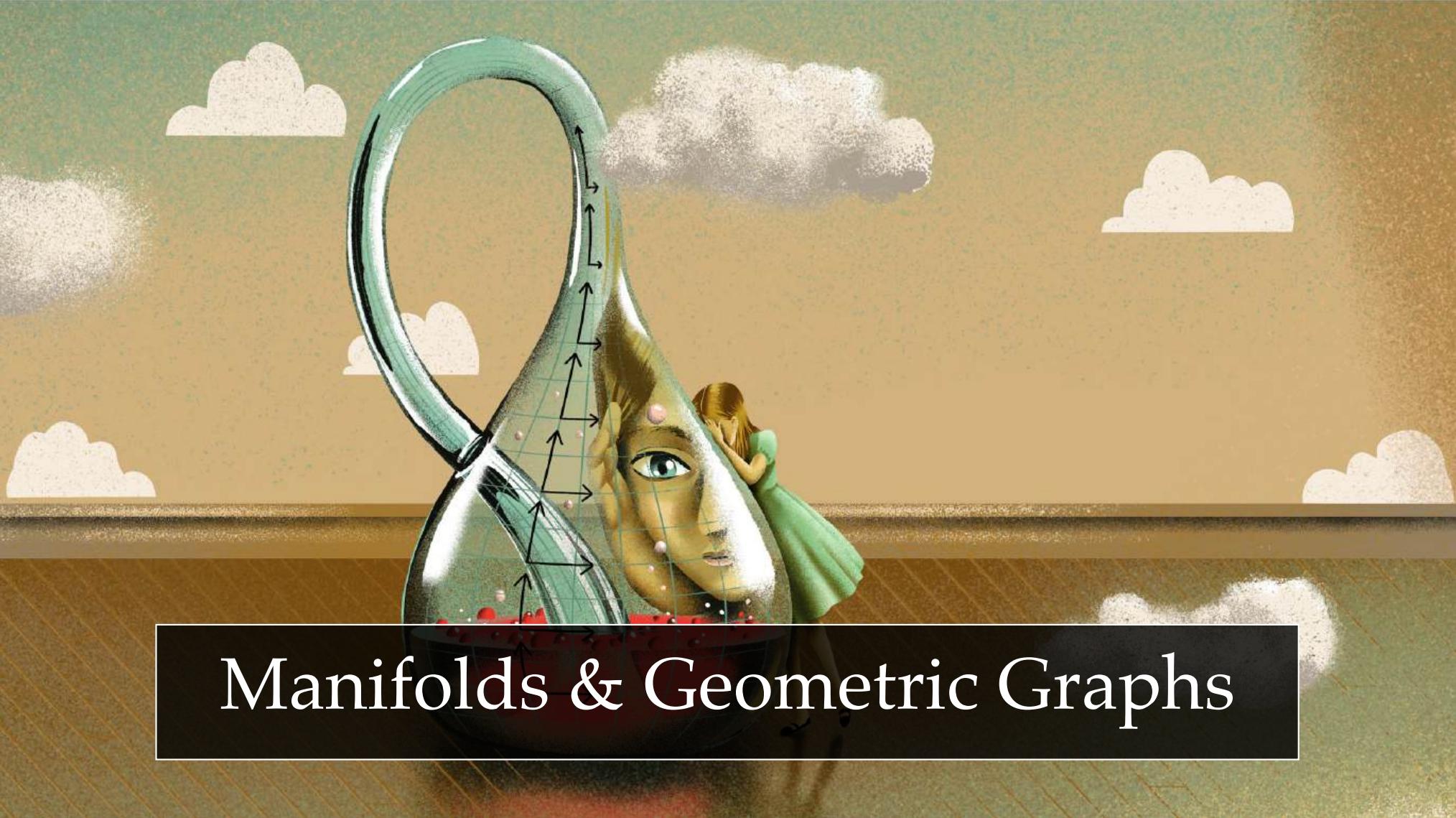
$$((x \star \psi) \star \phi)(R) = \int_{SO(3)} (x \star \psi)(Q) \phi(R^{-1}Q)dQ$$



Cohen, Welling 2016

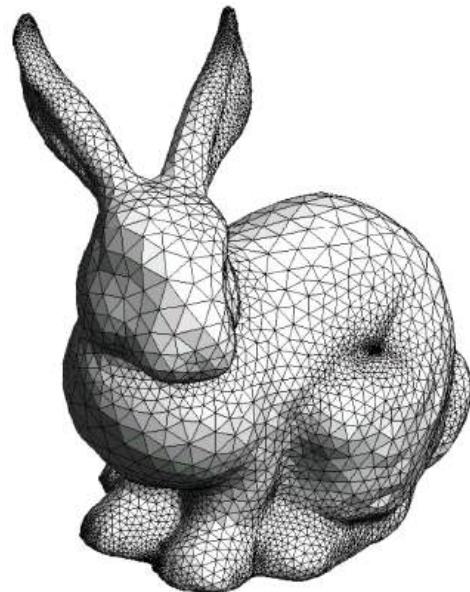
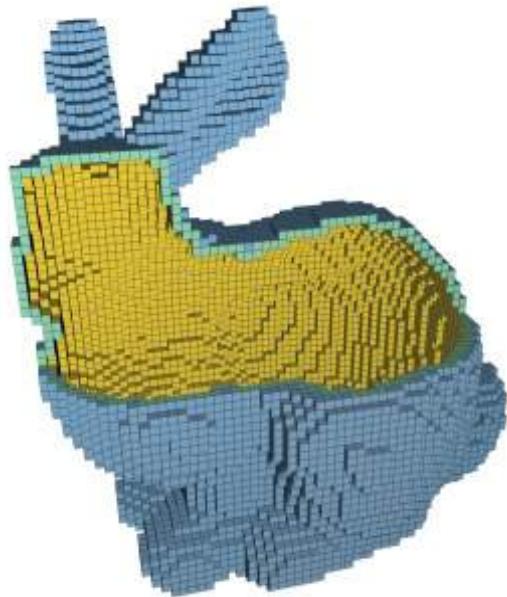
Homogeneous Spaces





Manifolds & Geometric Graphs

Why Manifolds?



More efficient representation: no “waste”
for internal structures

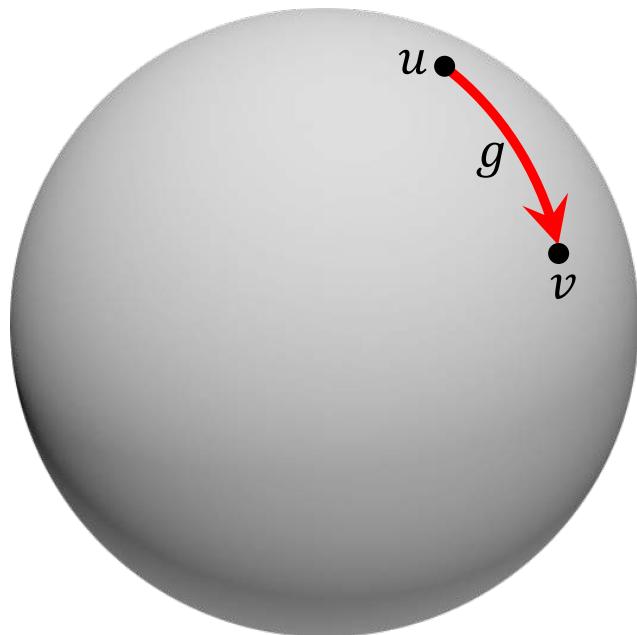
Natural model for
deformable shapes

Why Manifolds?

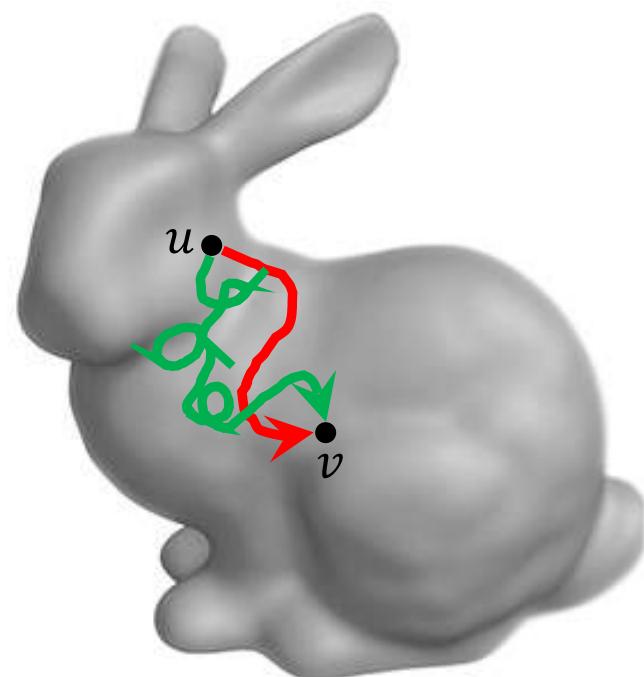
In protein modeling,
abstract out internal
structure that is irrelevant
for interactions + allow
some conformation
changes



Homogeneous Spaces

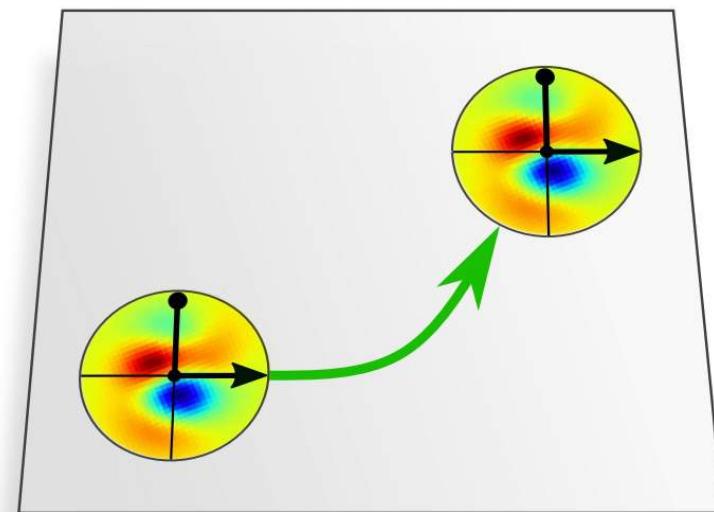
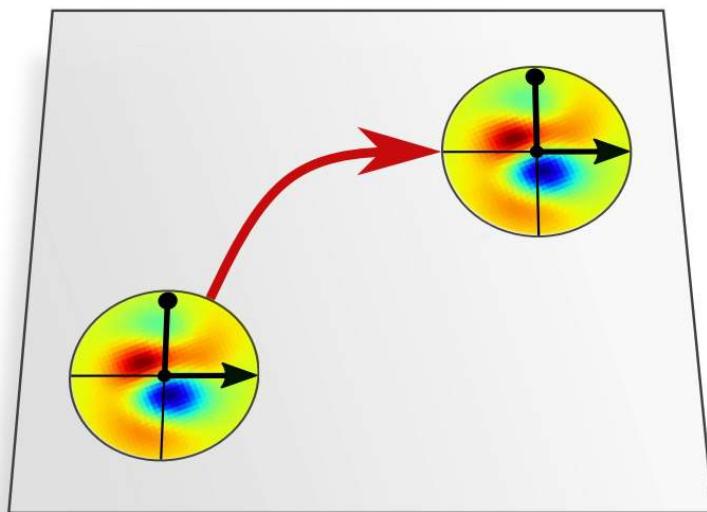


global symmetry group
 $\exists g \in G$ s.t. $gu = v$



no global symmetry group

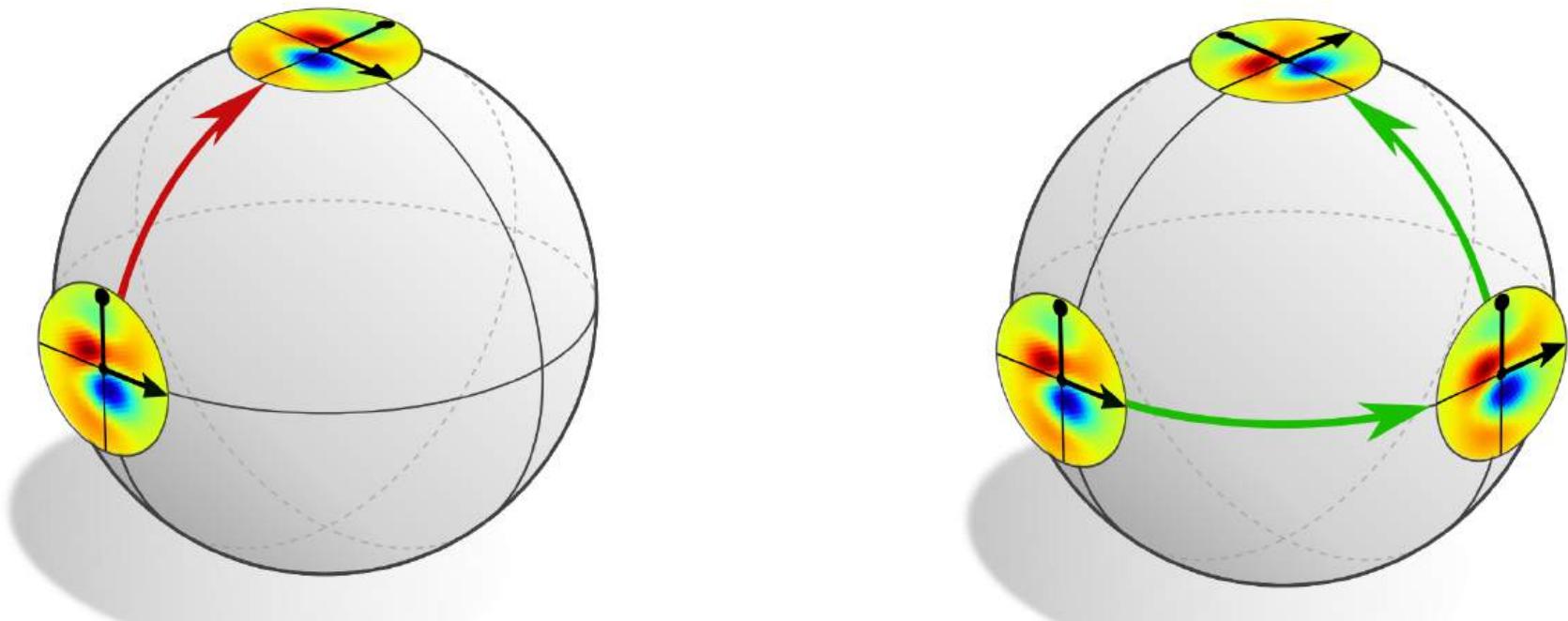
Euclidean Convolution



Euclidean space: Transport the filter around the domain

Figure: M. Weiler et al. 2021

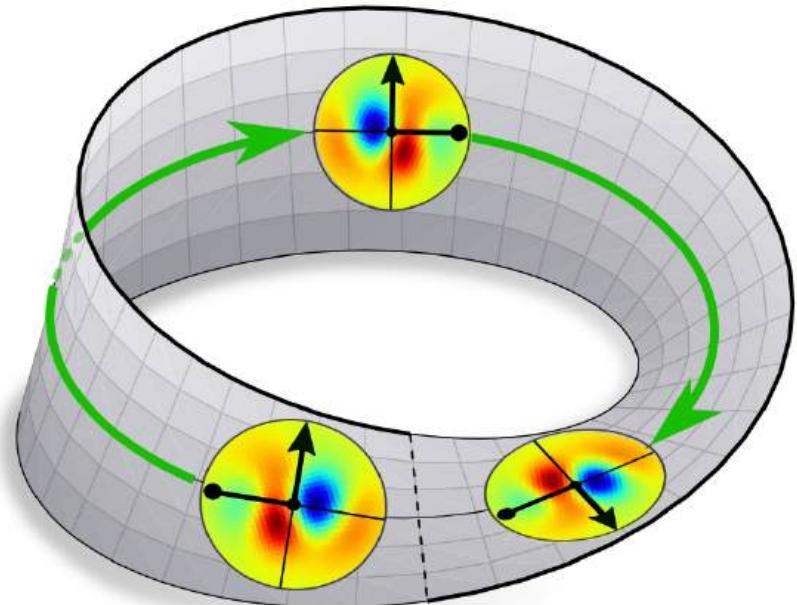
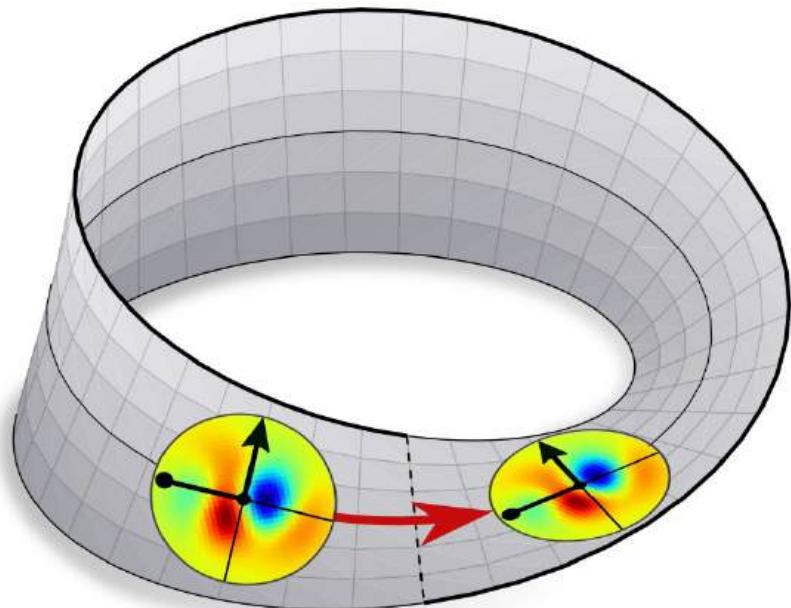
Non-Euclidean Convolution



Manifold: Result of transport is *path dependent*

Figure: M. Weiler et al. 2021

Non-Euclidean Convolution

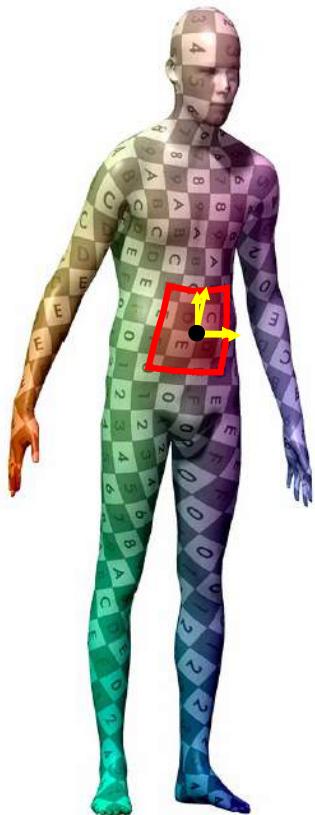


Manifold: Result of transport is *path dependent*

Figure: M. Weiler et al. 2021

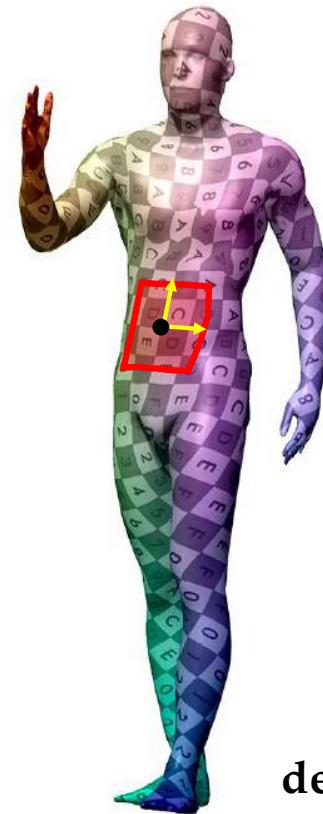
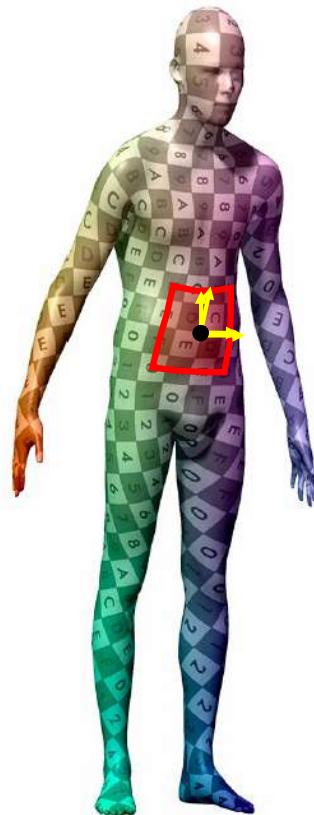
Two Types of Invariance

**Local gauge
transformation**



Two Types of Invariance

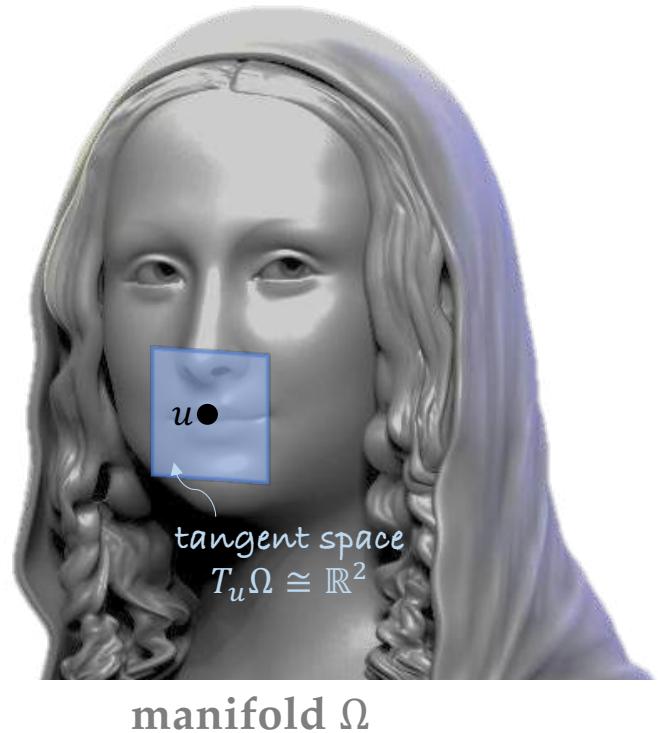
**Local gauge
transformation**



**Global
deformation**



Manifolds



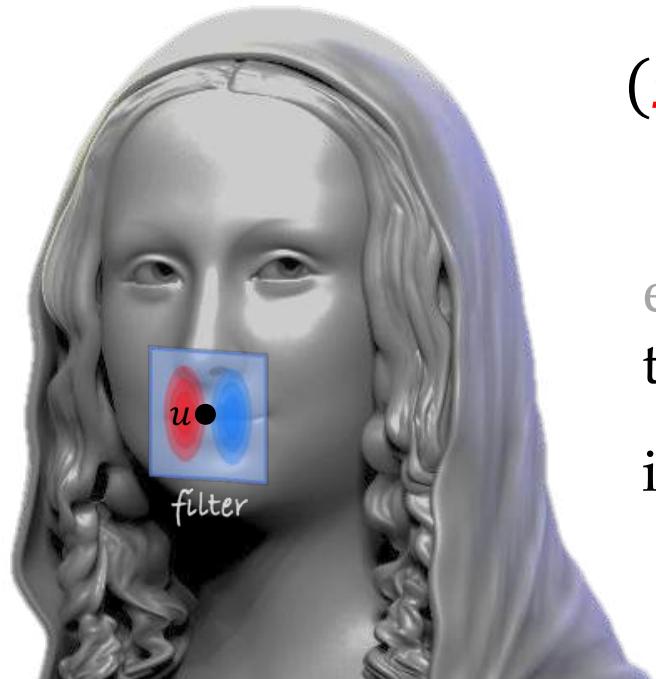
manifold = locally Euclidean space

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

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

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

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

intrinsic filter = invariant to isometries

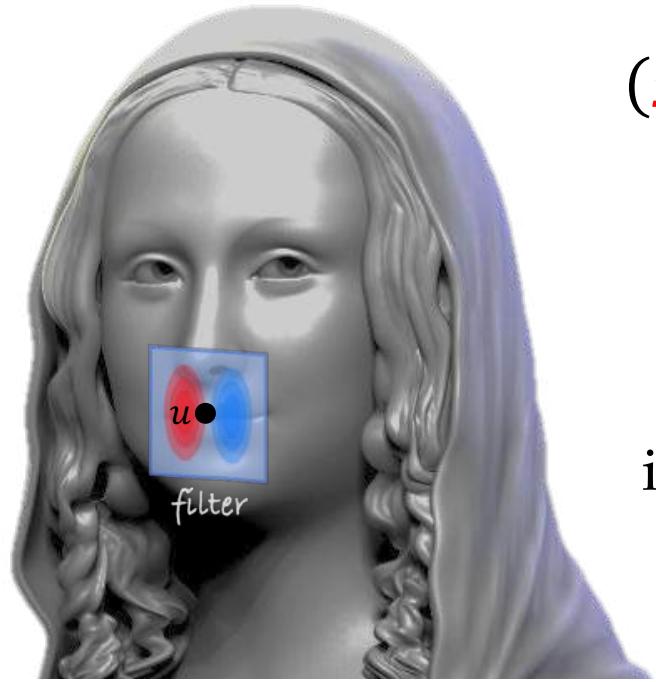
Geodesic CNNs



Anisotropic intrinsic filters on a manifold

Masci et B 2015; Boscaini et B 2016; Monti et B 2017

Geodesic CNNs



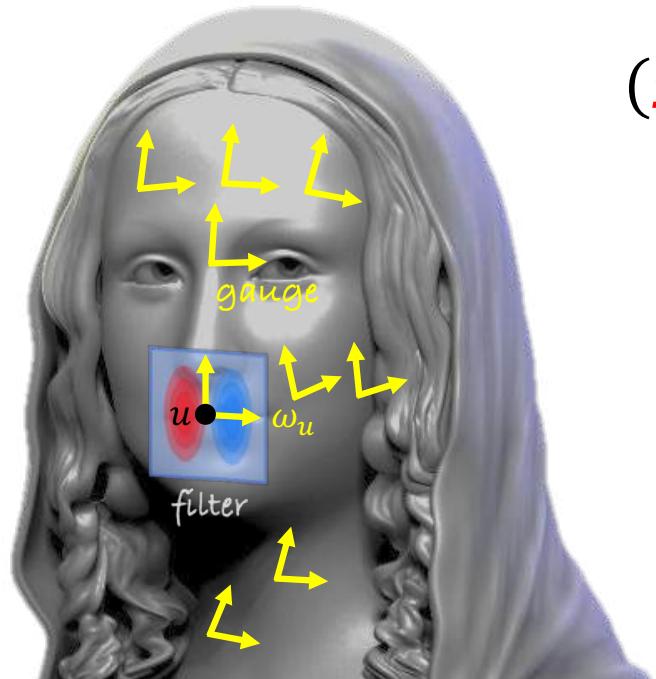
manifold Ω
isometry group $\text{Iso}(\Omega)$

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

Problem: these are abstract vectors!

intrinsic filter = invariant to isometries

Geodesic CNNs

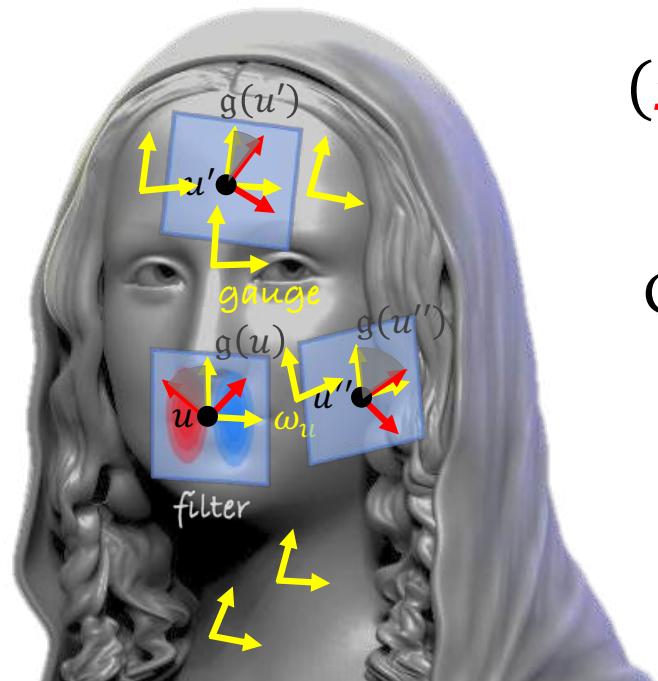


manifold Ω
isometry group $\text{Iso}(\Omega)$

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

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

Gauge Transformations



manifold Ω
structure group G

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

Gauge defined up to gauge transformation

$$g: \Omega \rightarrow G$$

Cohen et al. 2019; Weiler et al. 2021

Structure Group

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

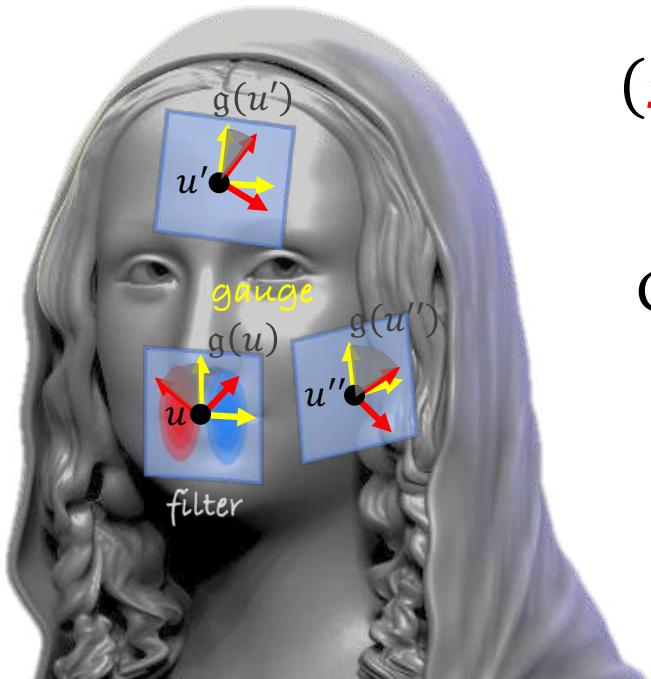
- | | | |
|-------------------------------|-----------------|-------------------------------------|
| • “Naked” manifold | $GL(s)$ | invertible matrices |
| • Manifold+orientation | $GL^+(s)$ | invertible matrices with $\det > 0$ |
| • Manifold+volume | $SL(s)$ | matrices with $\det = 1$ |
| • Manifold+metric | $O(s)$ | 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* $g: \Omega \rightarrow G$

• “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 $G = \text{SO}(2)$

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

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

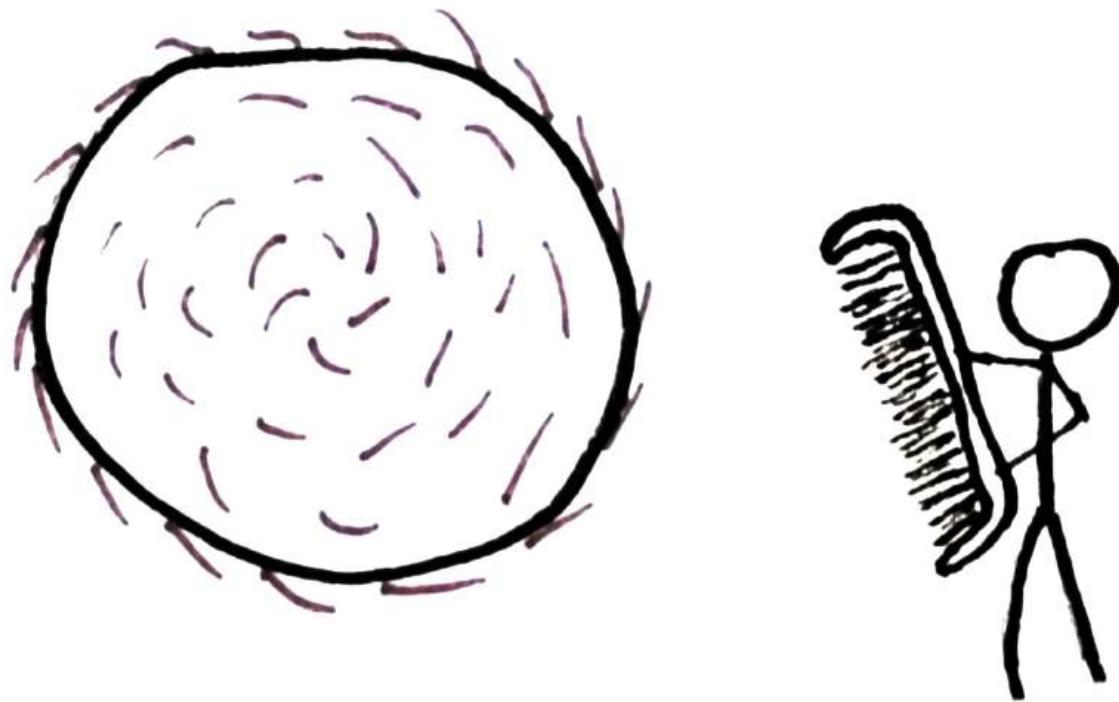
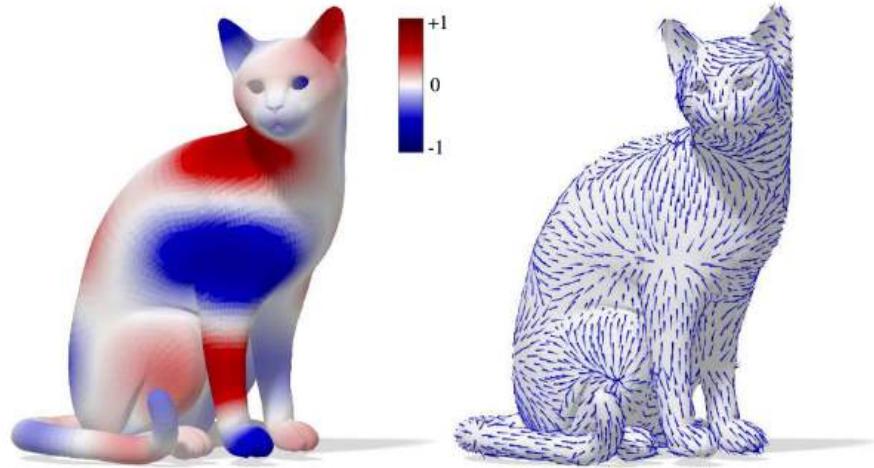
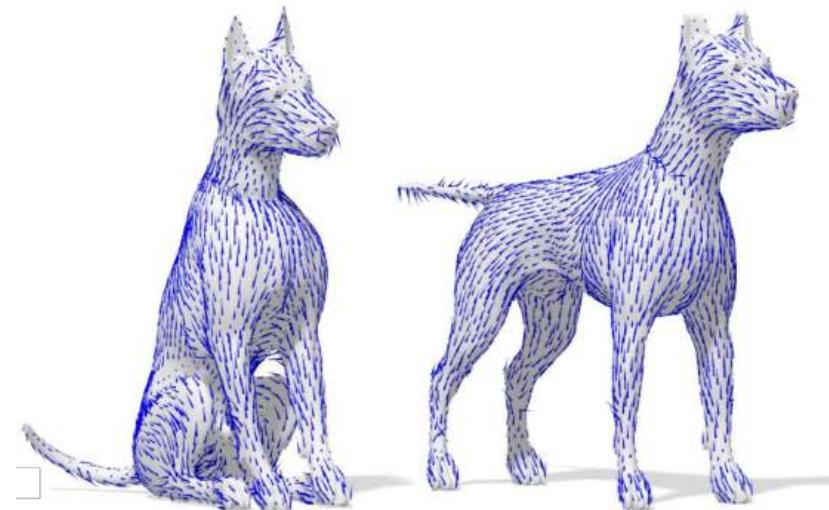


Image: Minutephysics

Theory vs Practice: Stable Gauges



Gradient of intrinsic function

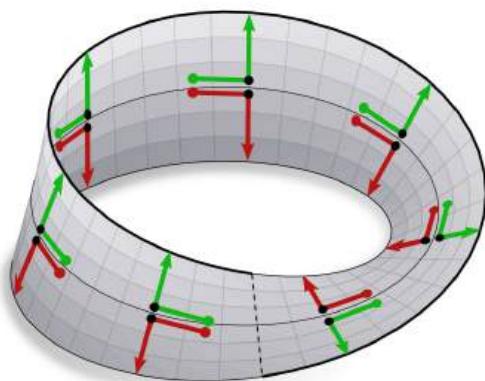


Deformation-invariant
stable gauge

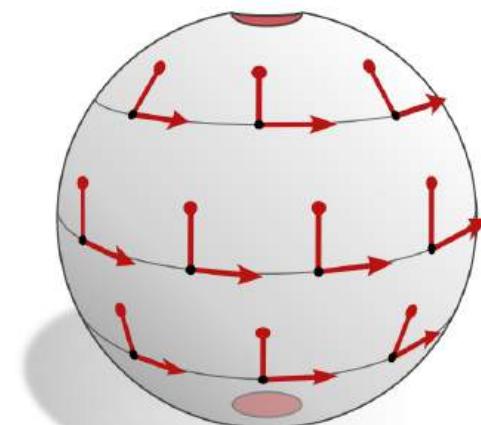
Structure Group



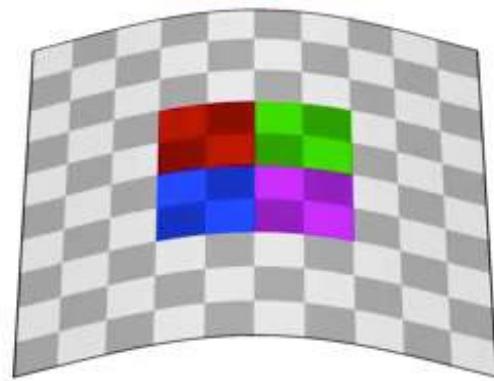
rotation
 $\text{SO}(2)$



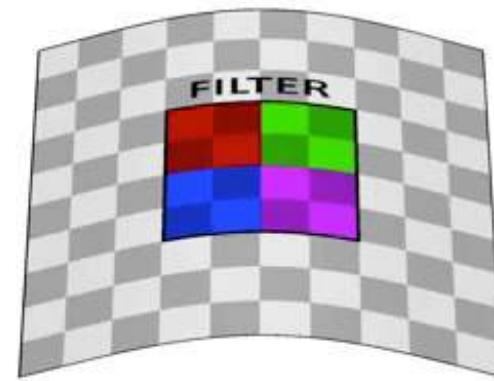
reflection
 R



fixed gauge
 $\{\text{id}\}$



**Euclidean (extrinsic)
convolution**



**Geometric (intrinsic)
convolution**



Applications

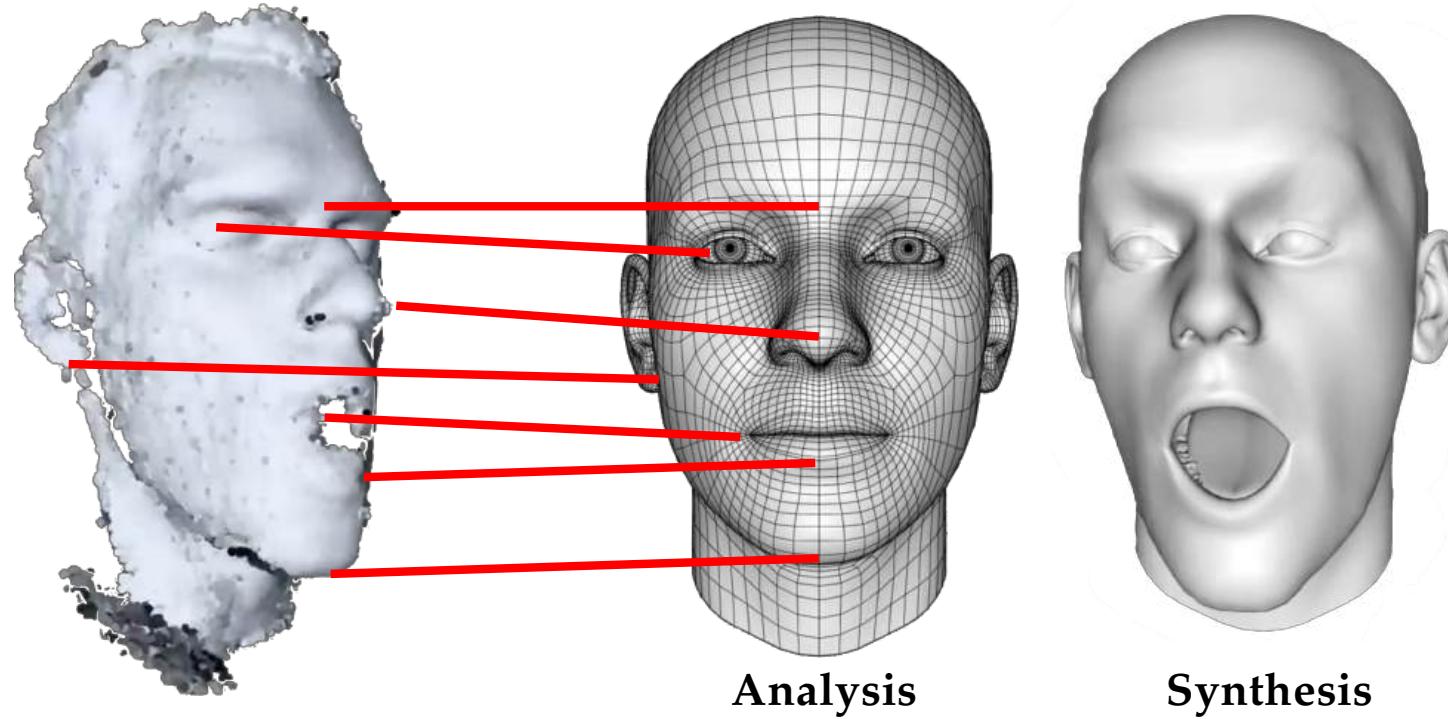
FaceShift 2015



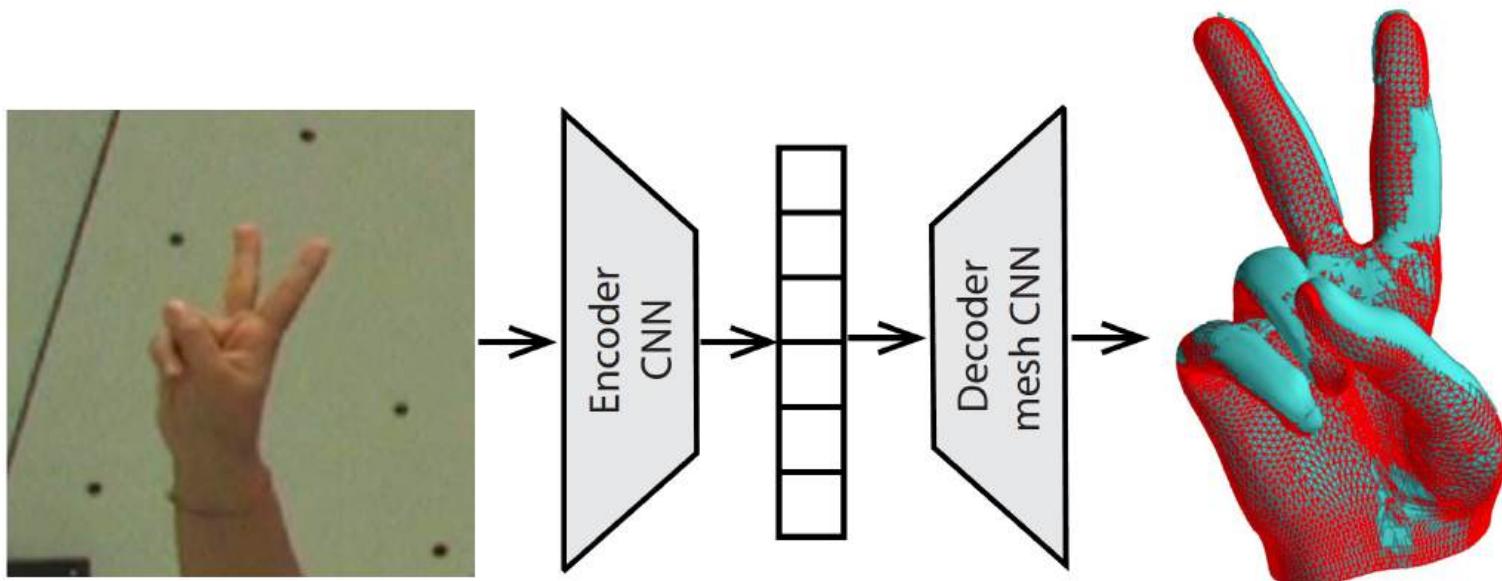
GDC

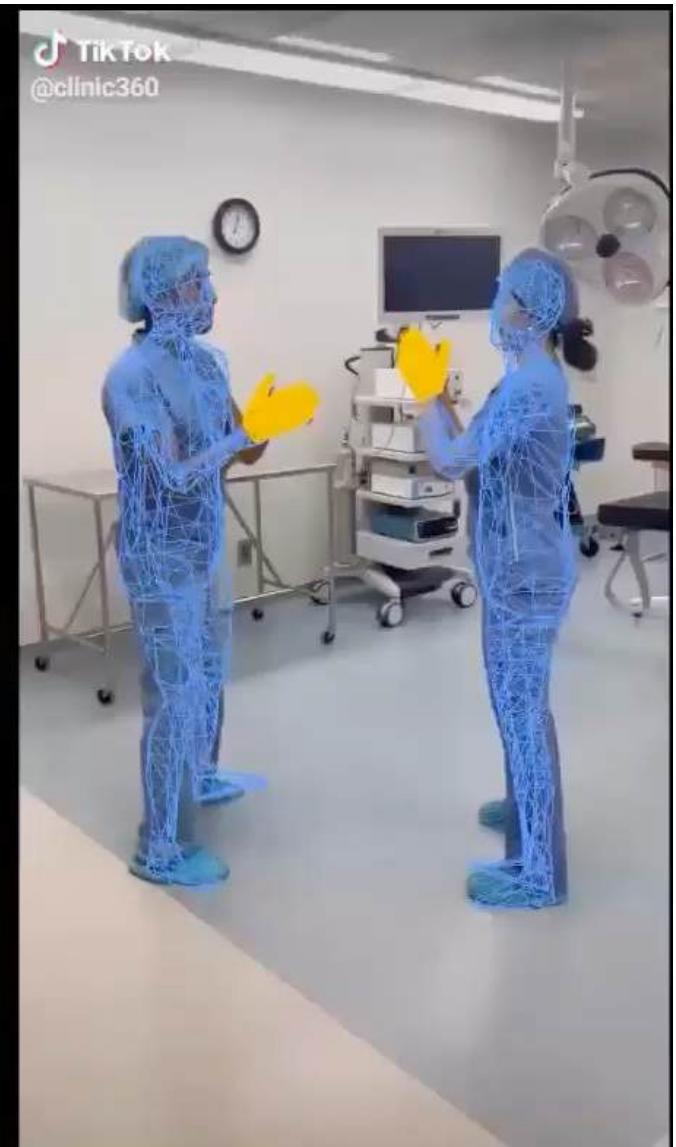


Shape Analysis & Synthesis



3D Hand Reconstruction



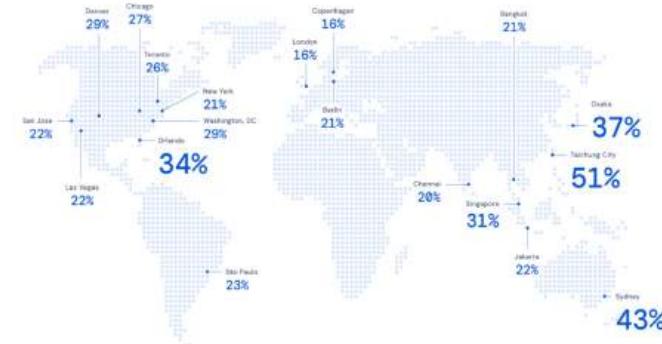
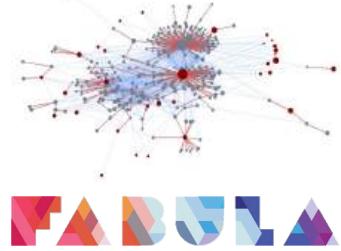


Snap Acquires Ariel AI To Enhance AR Features

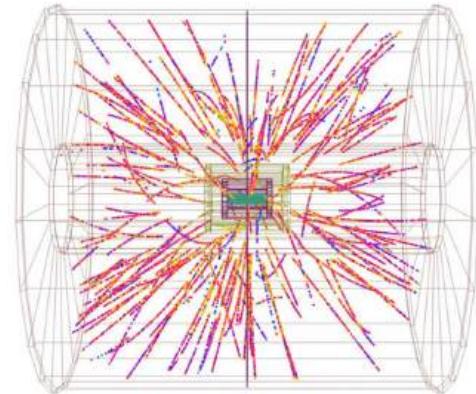




Fake news detection



Navigation



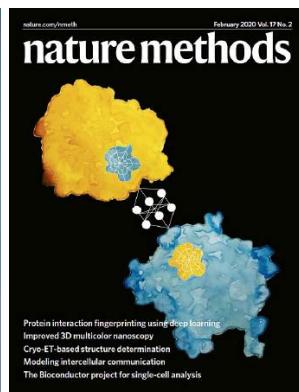
Particle physics



Pure math



Structural biology



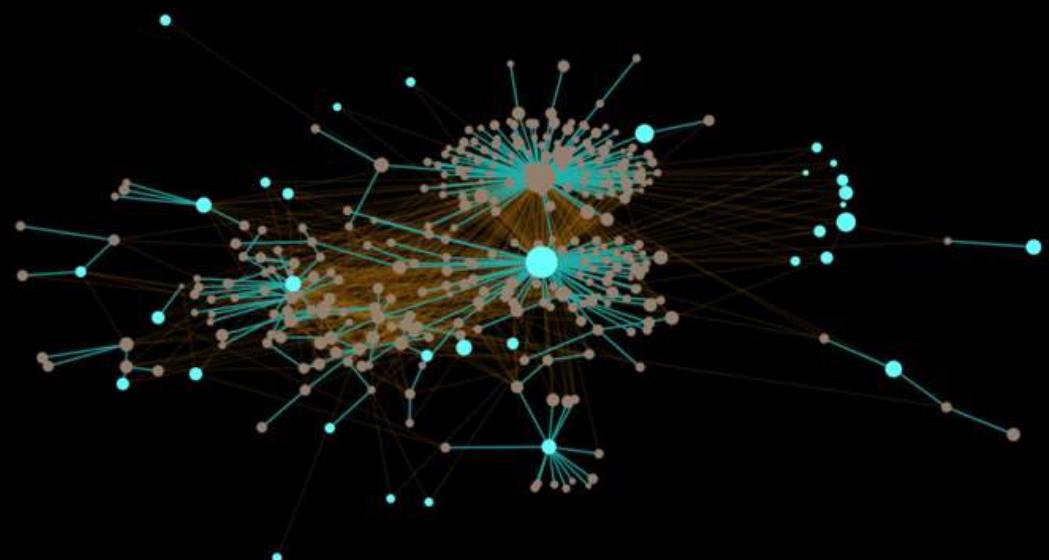
Drug discovery



Safe spaces for South
Asia's vultures p. 1086

Synthetic Notch receptors
enhance T cell therapies p. 1112

Dietary fiber fights
diabetes p. 1151

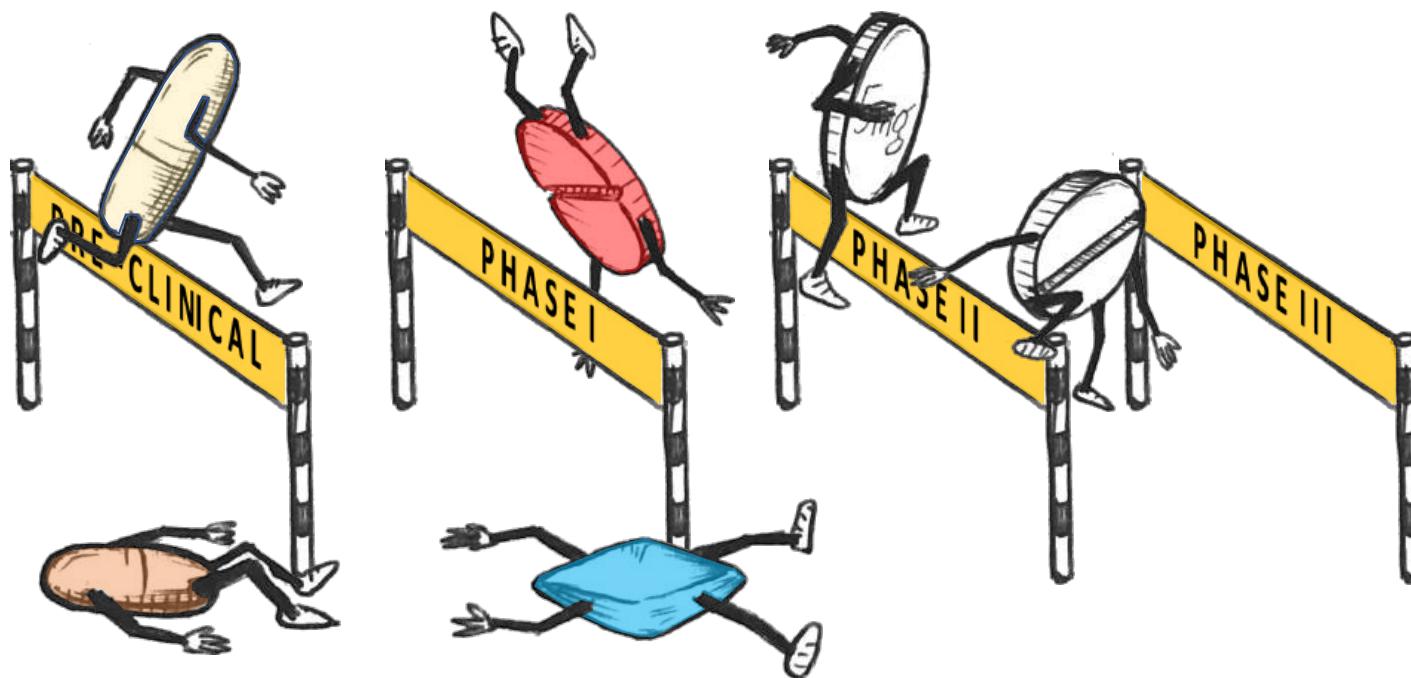


Monti et B 2019

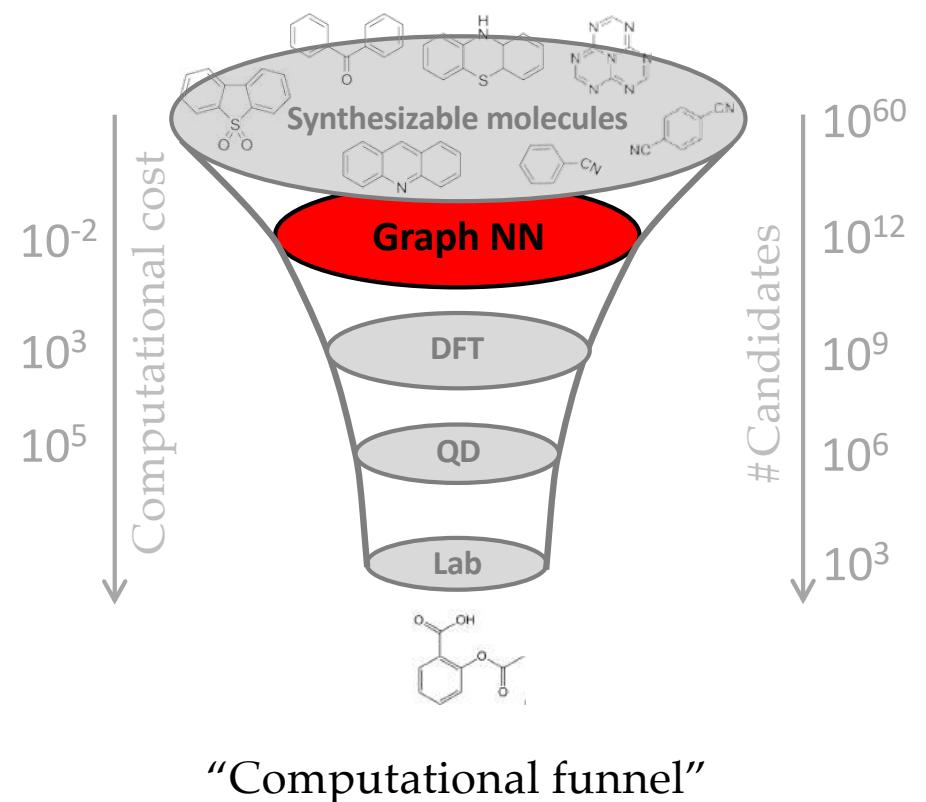
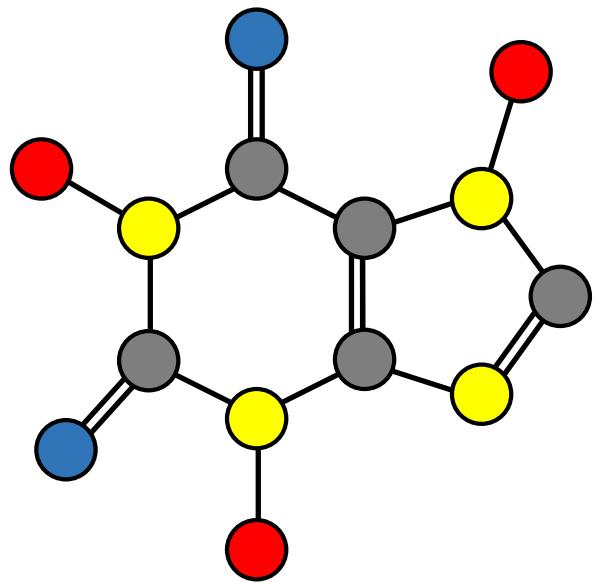


Acquired by Twitter in 2019

Drug Discovery & Design

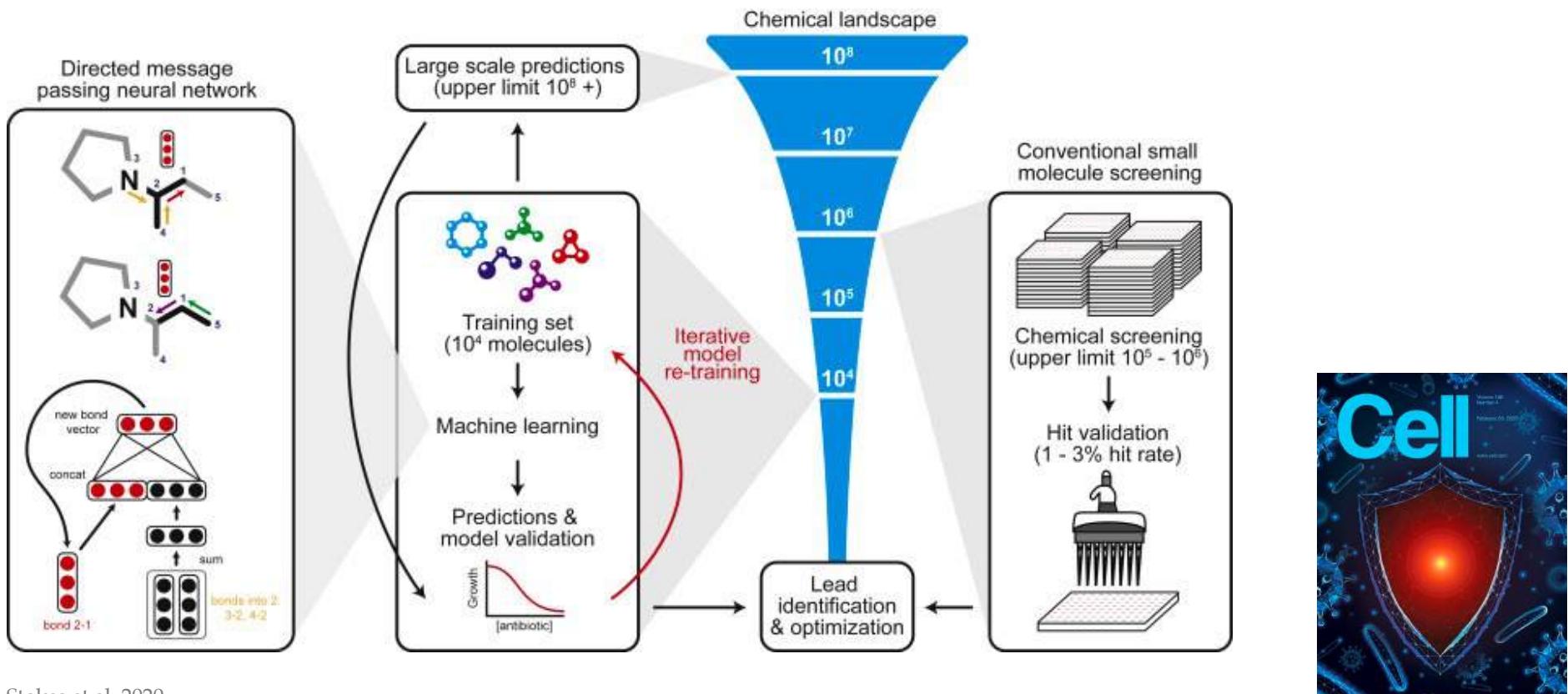


Virtual Drug Screening

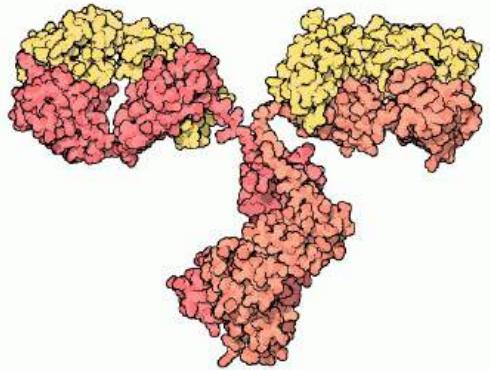


Duvenaud et al. 2015; Gilmer et al. 2017; Jin et al. 2020

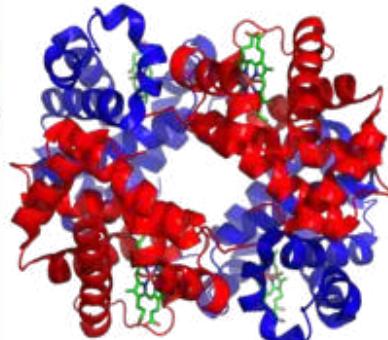
New Antibiotic Discovery



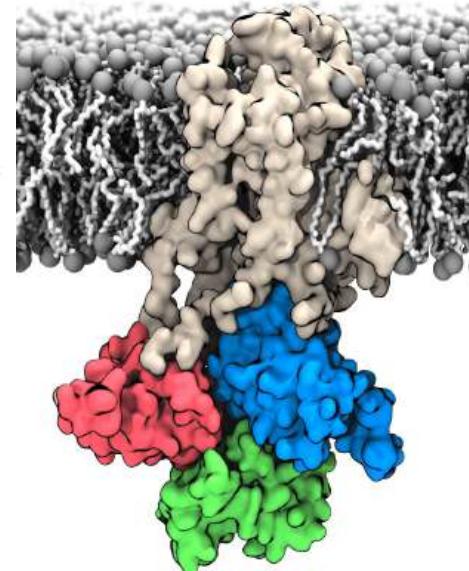
Stokes et al. 2020



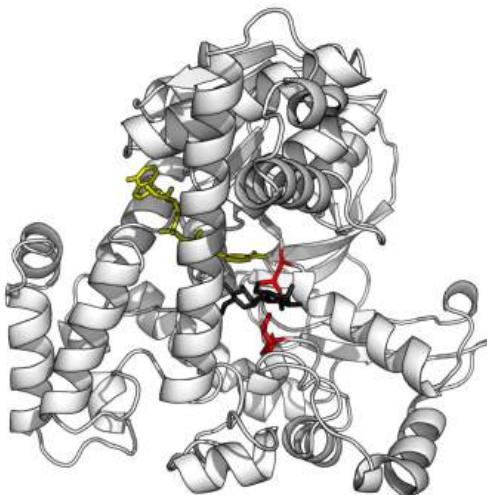
Defense (antibody)



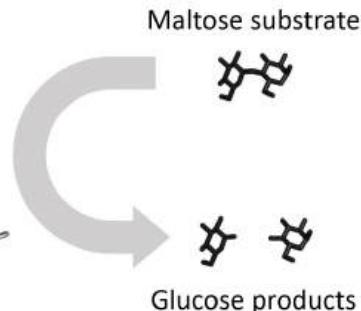
Storage (haemoglobin)



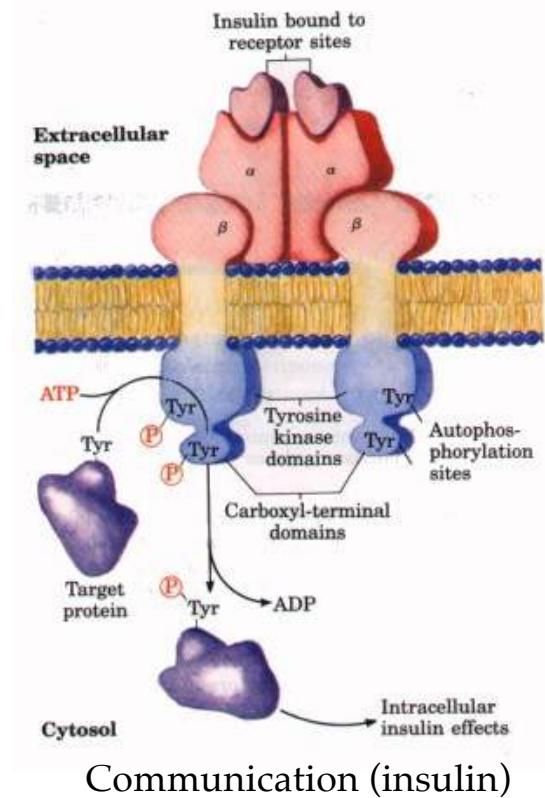
Transport (calcium pump)



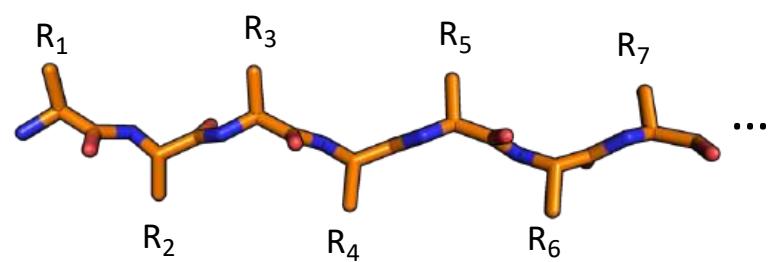
Catalysis (enzyme)



Structure (collagen)

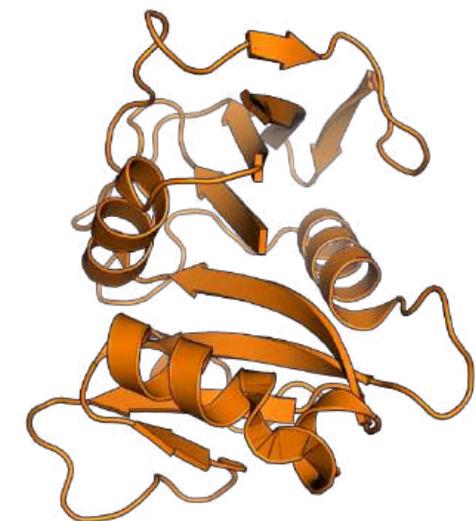


Protein Folding



Primary protein
structure

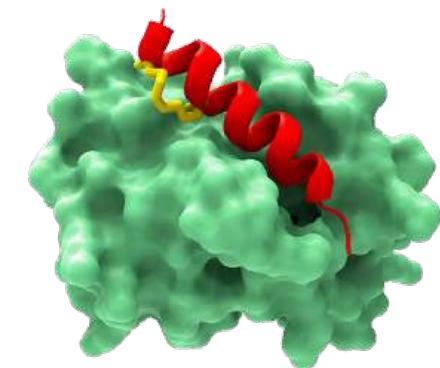
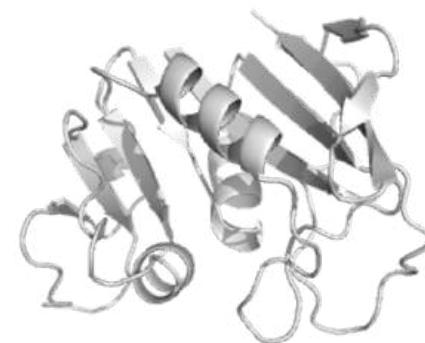
Protein folding



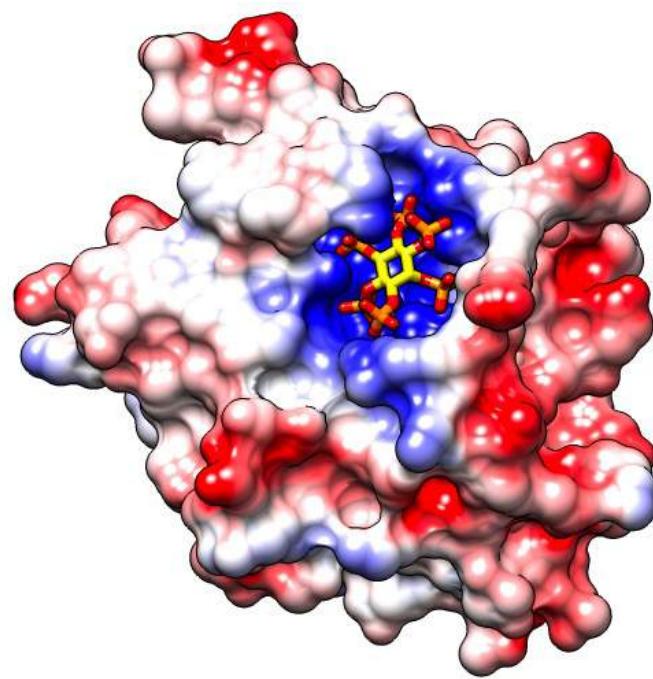
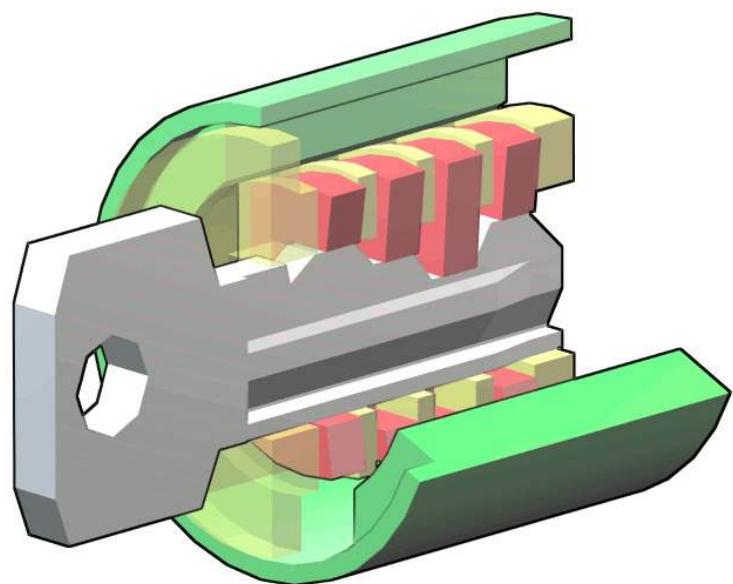
Tertiary protein
structure

Sequence → Structure → Function

L G C V L A W N T N S K

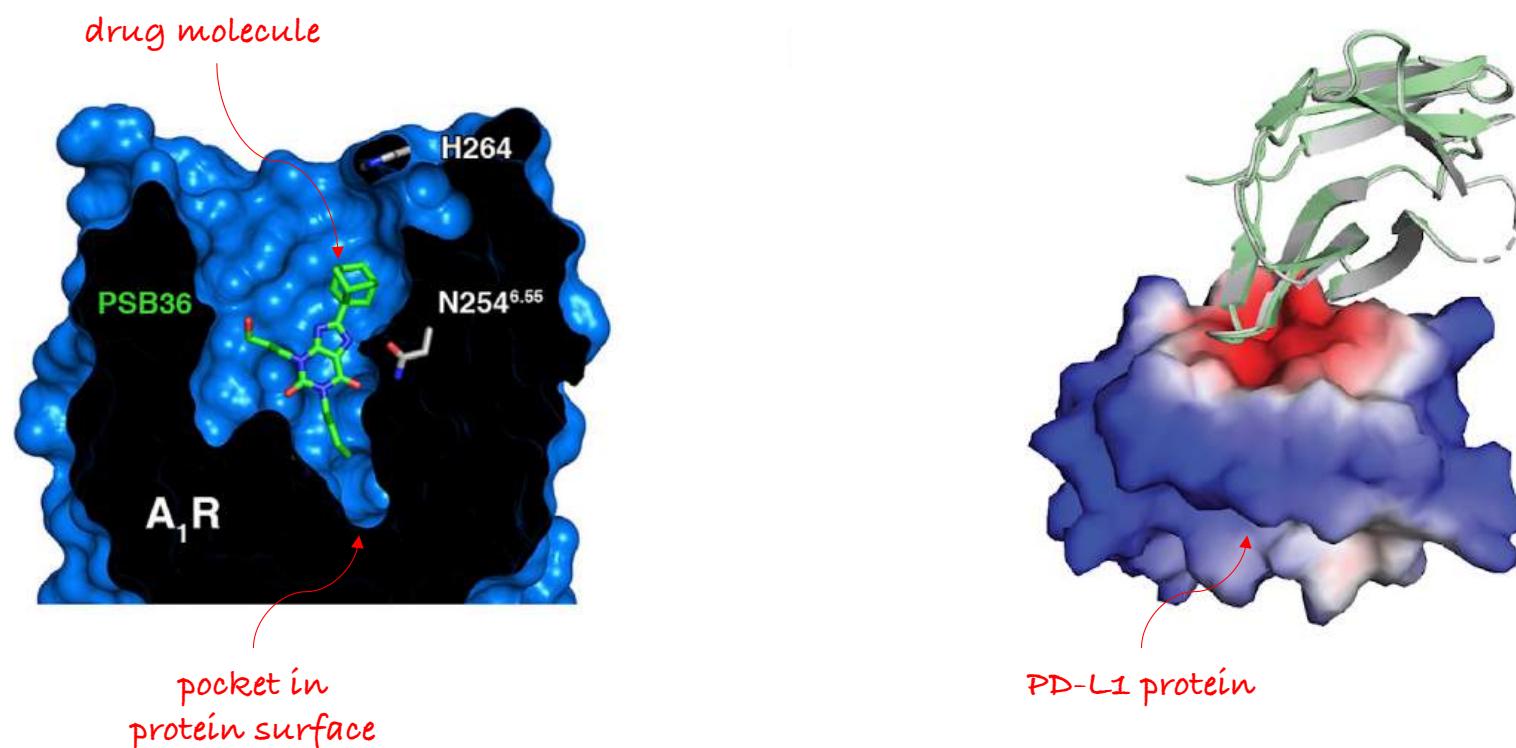


Lock-Key Metaphor

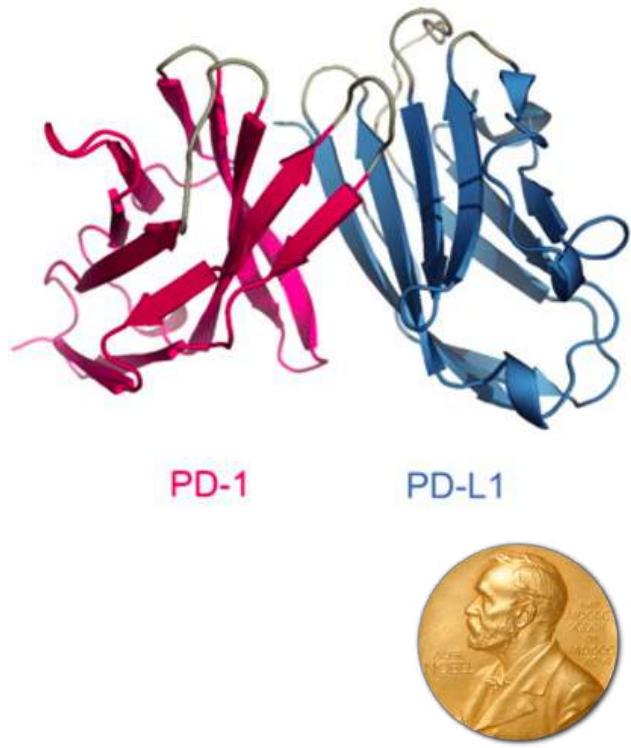


Emil Fischer "Schlüssel-Schloss-Prinzip" 1894

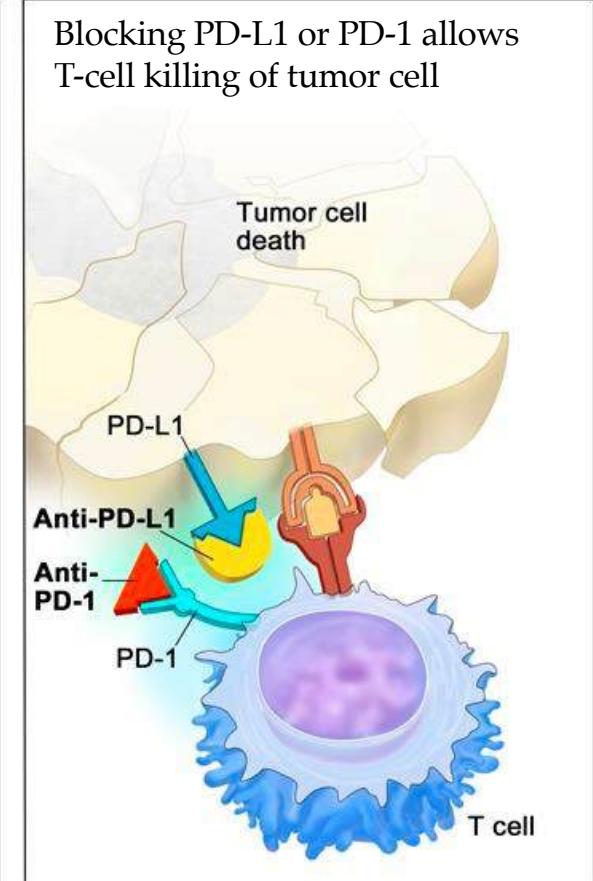
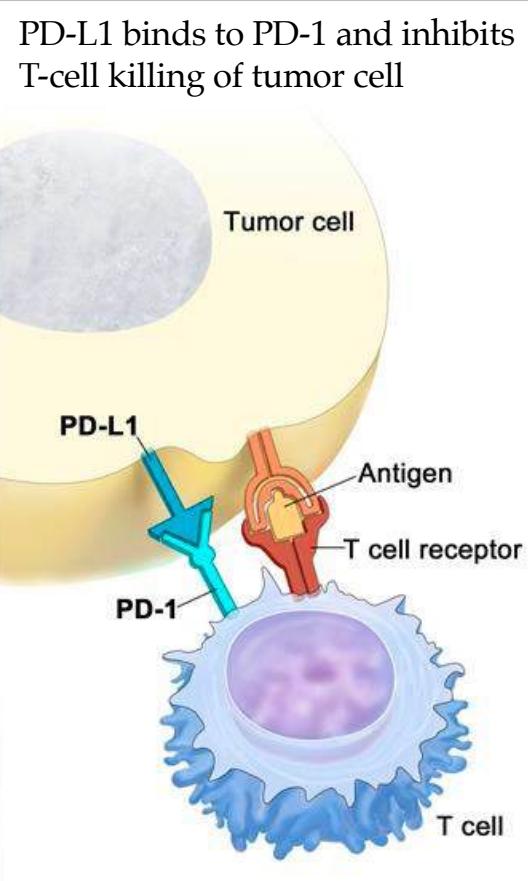
Biologic Drug Design



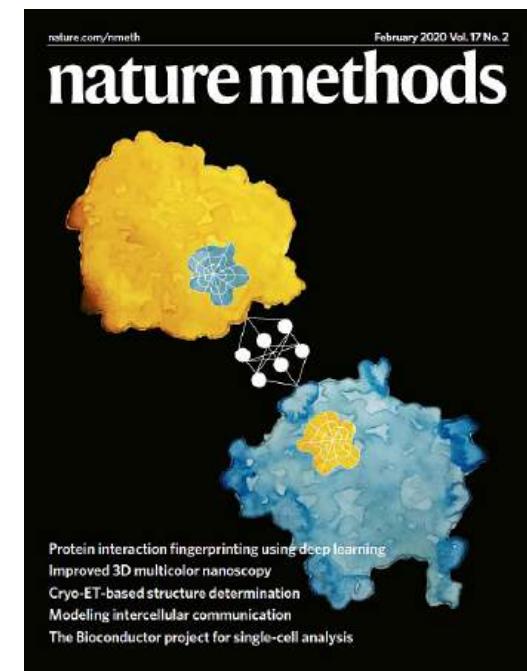
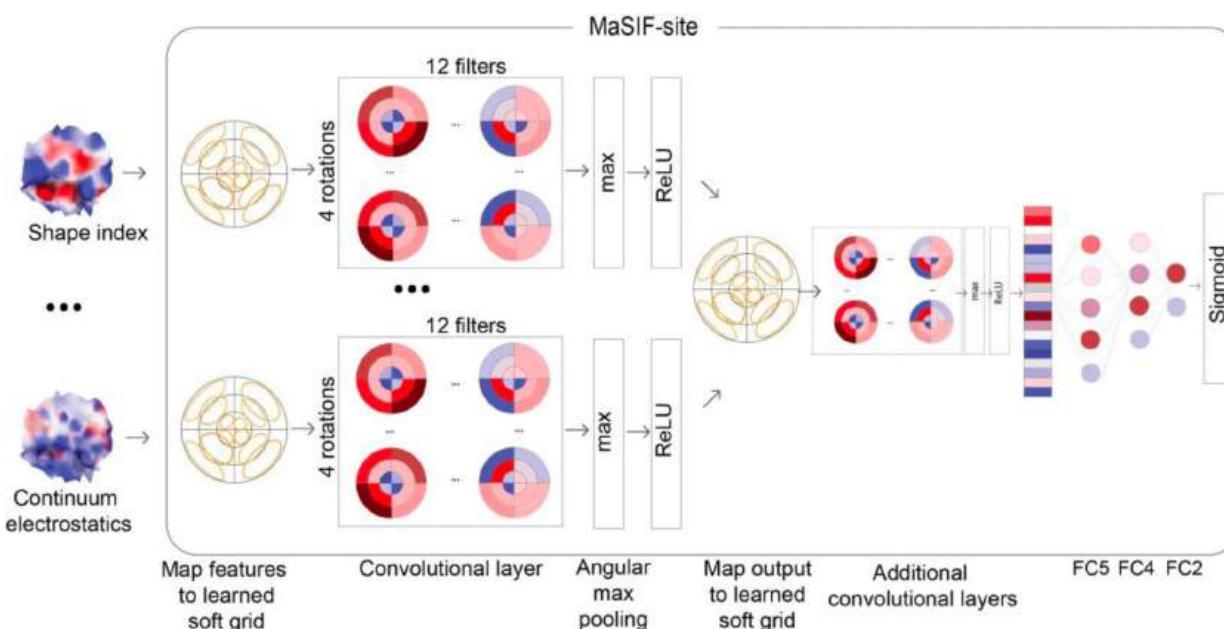
Immunotherapy



2018 Nobel Prize
PD-proteins role in
immunotherapy

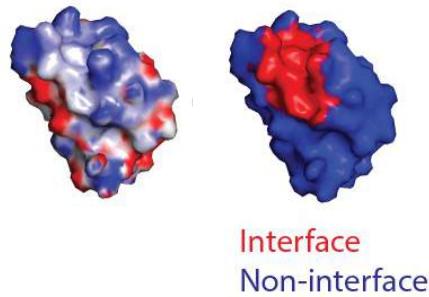


MaSIF: Geometric ML for Protein Function Prediction

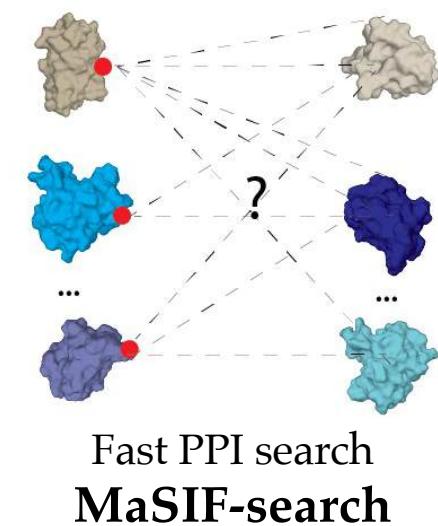
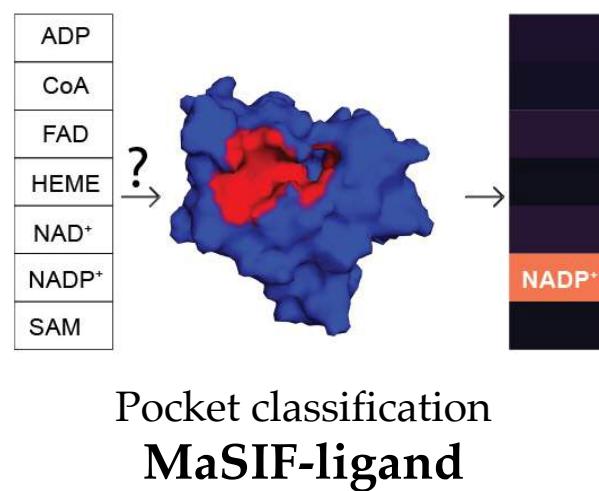


Gainza, Sverrisson et B, Correia 2020

MaSIF Applications

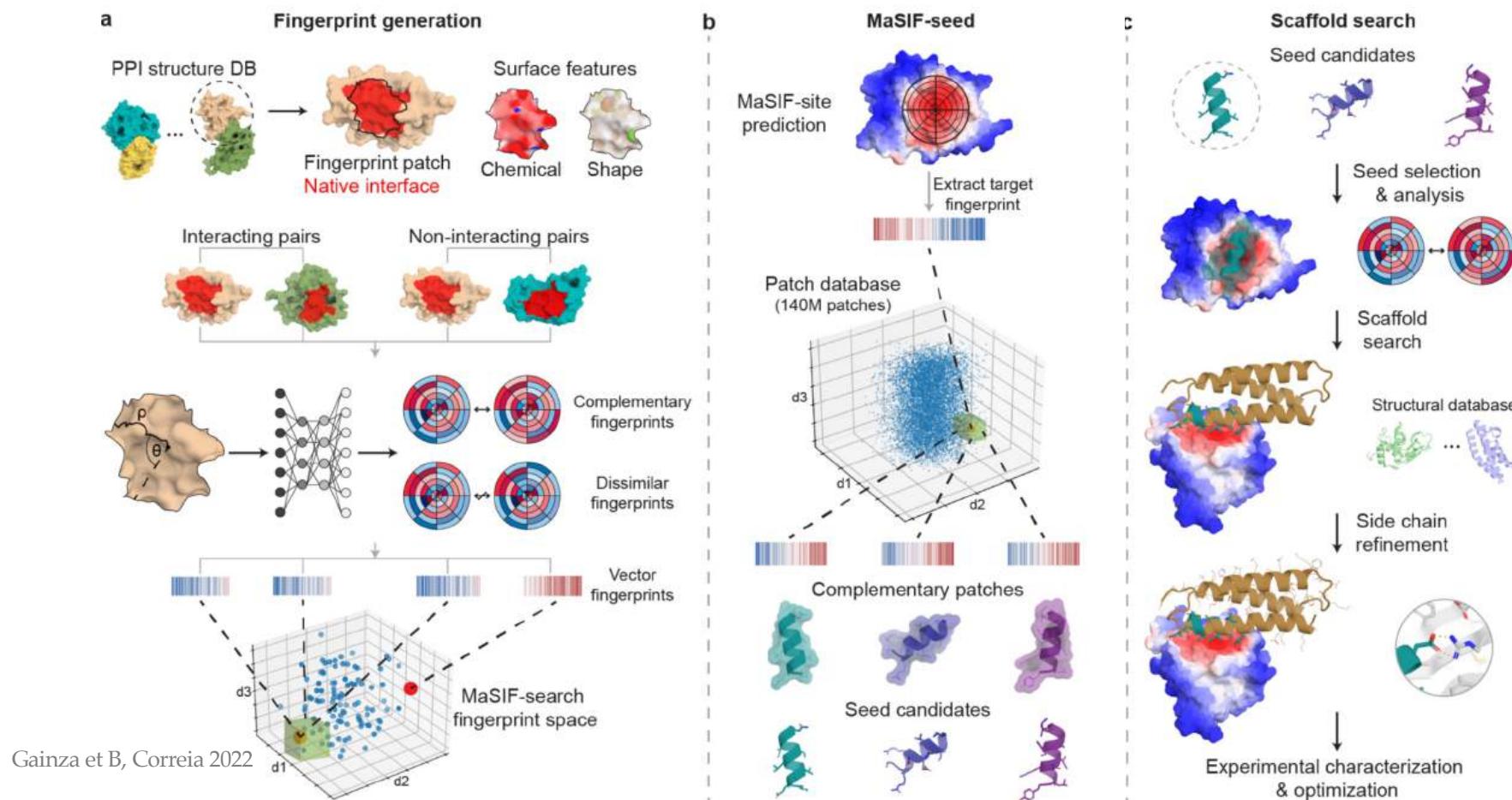


Interface site prediction
MaSIF-site

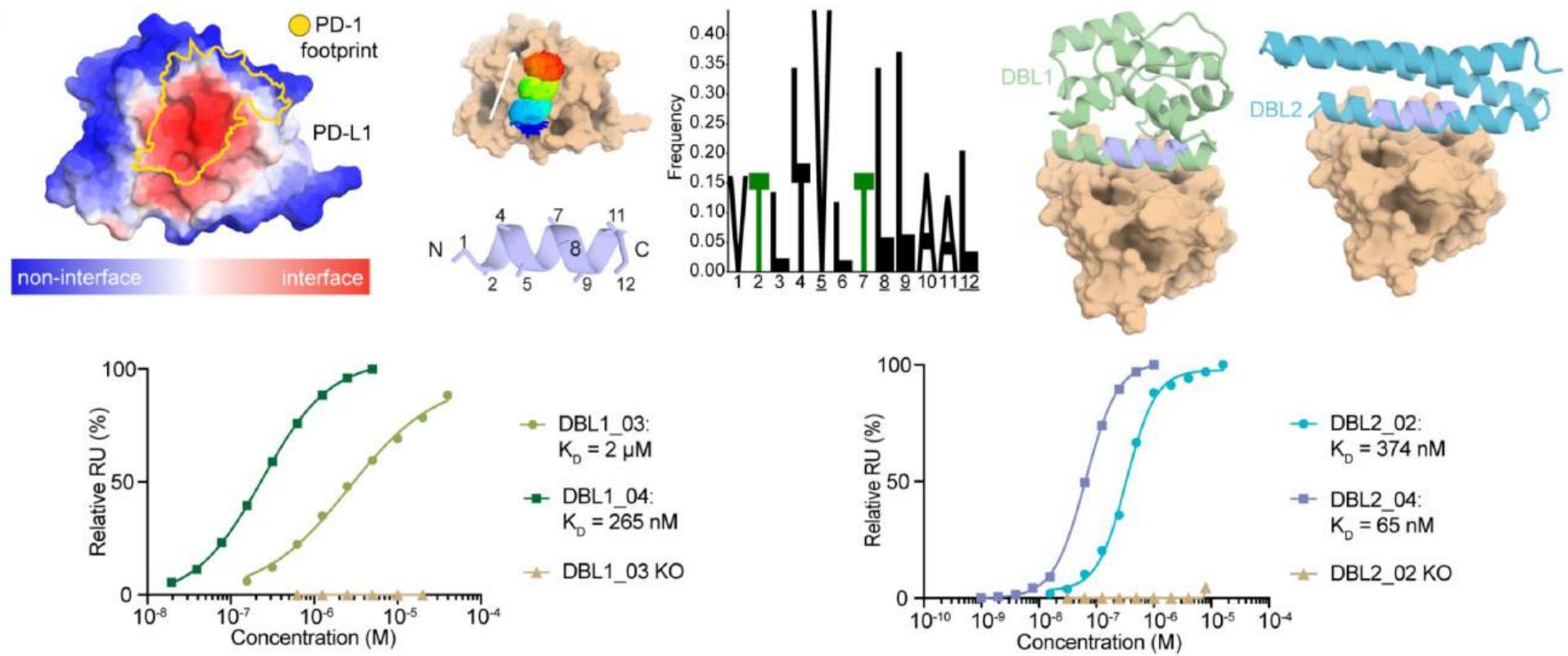




De novo protein design with MaSIF

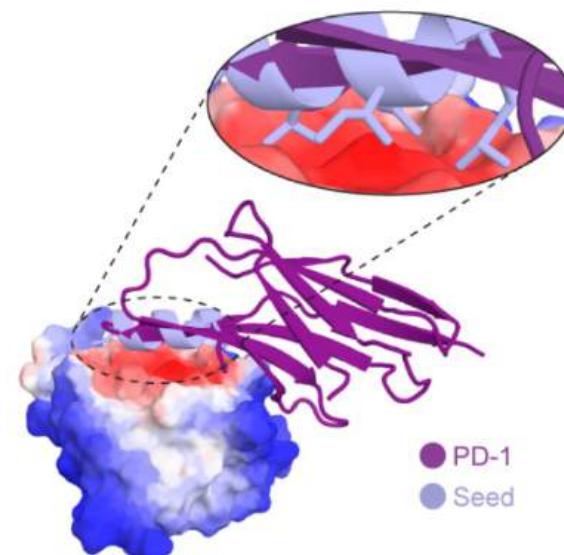


PD-L1 binder design

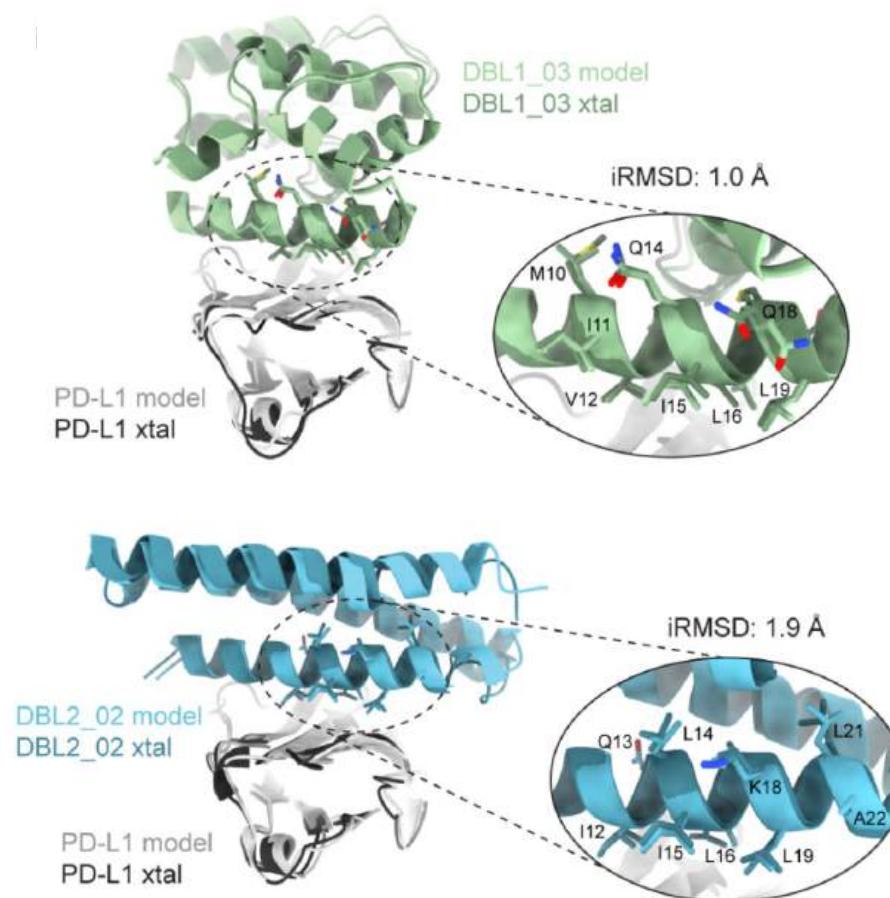


Gainza et al., Correia 2022

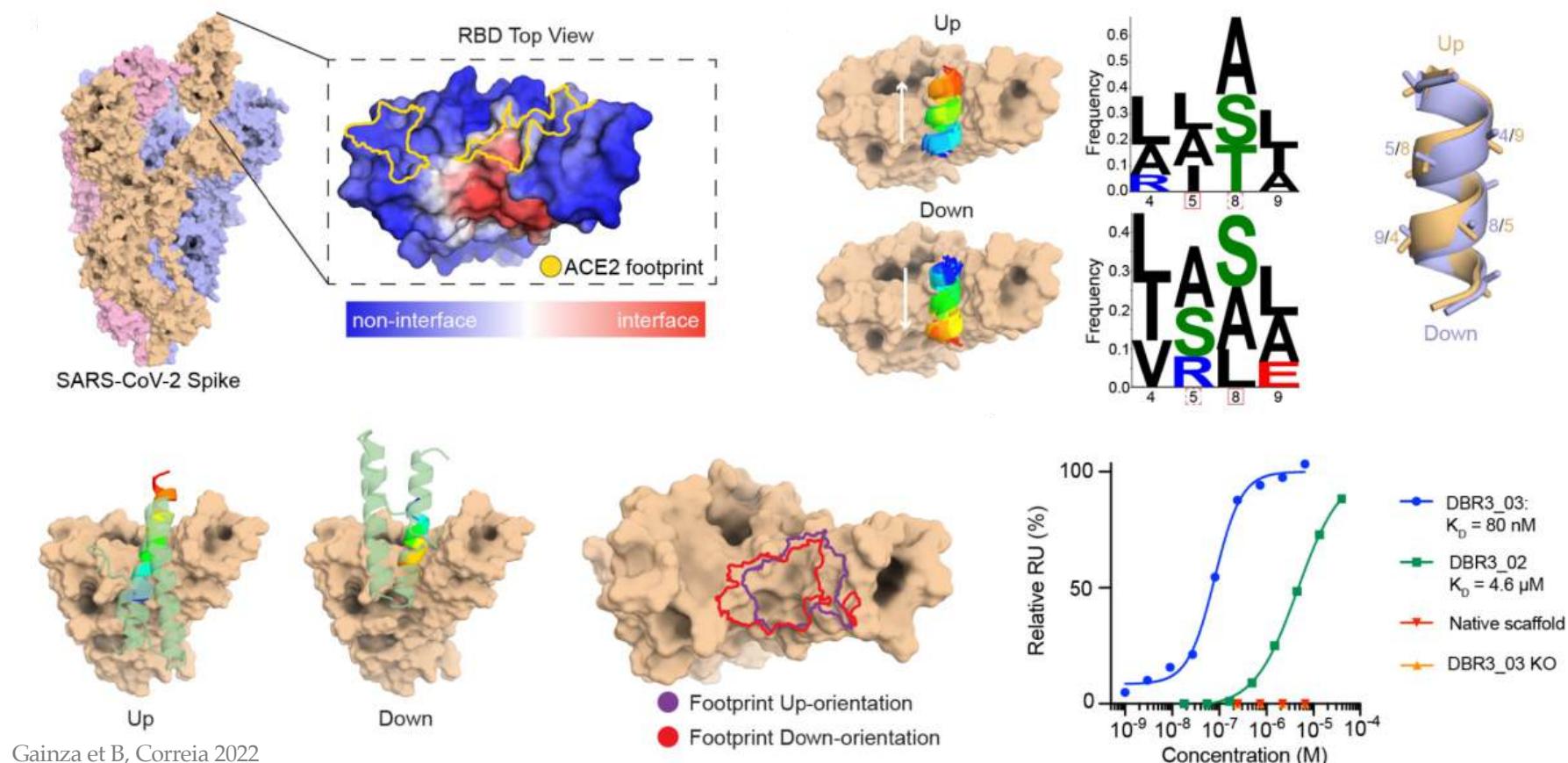
PD-L1 binder structure



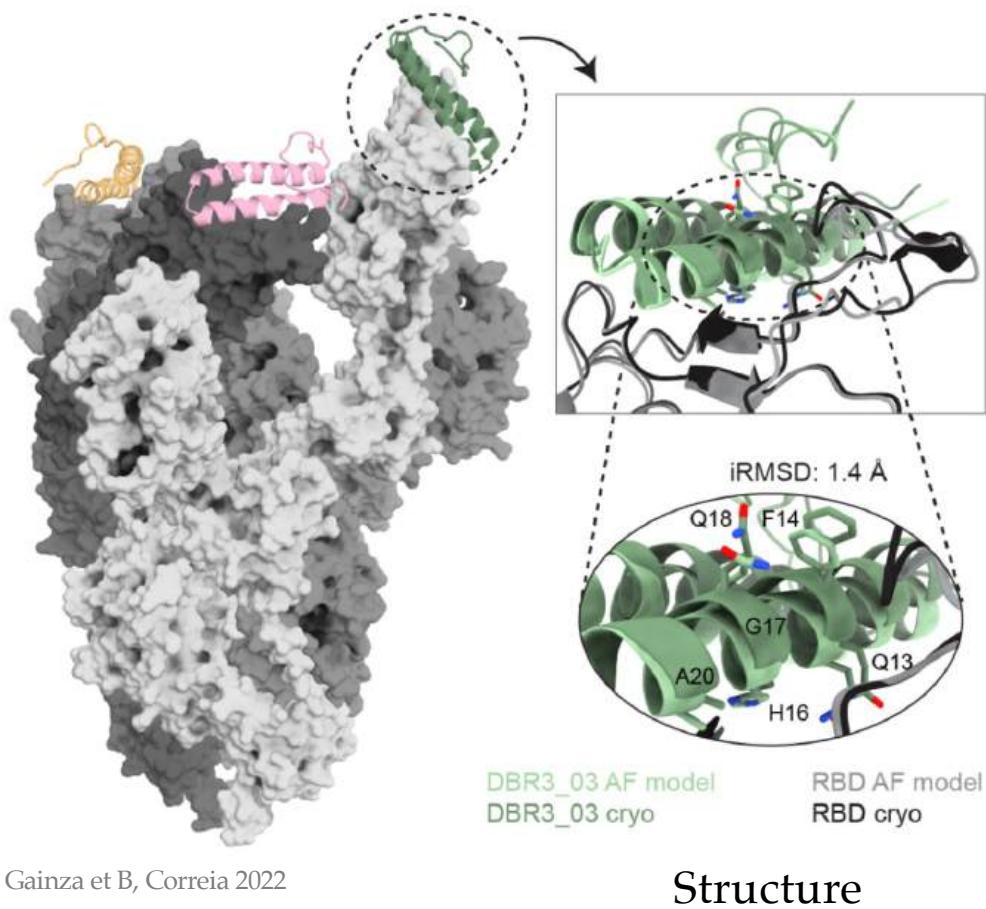
Gainza et B, Correia 2022



SARS-CoV-2 spike binder design

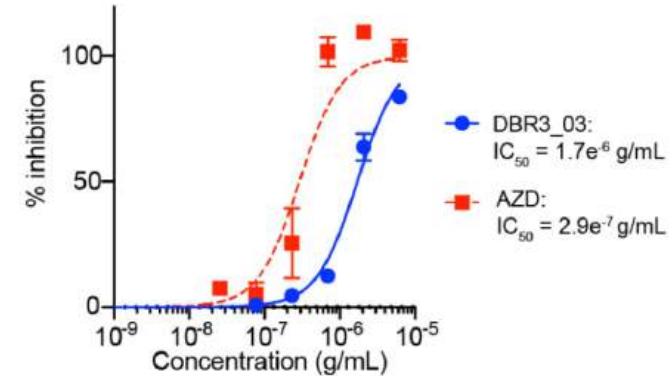
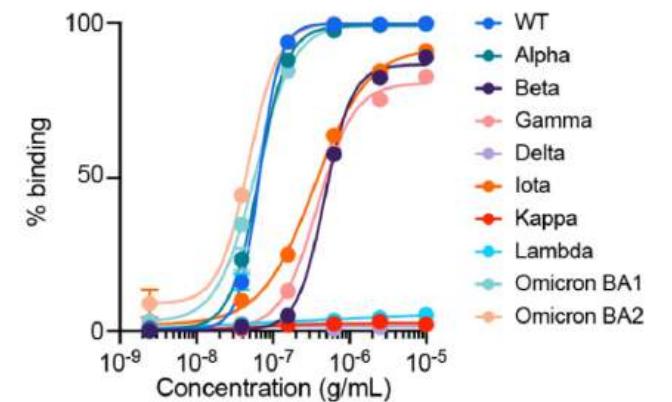


SARS-CoV-2 spike binder function



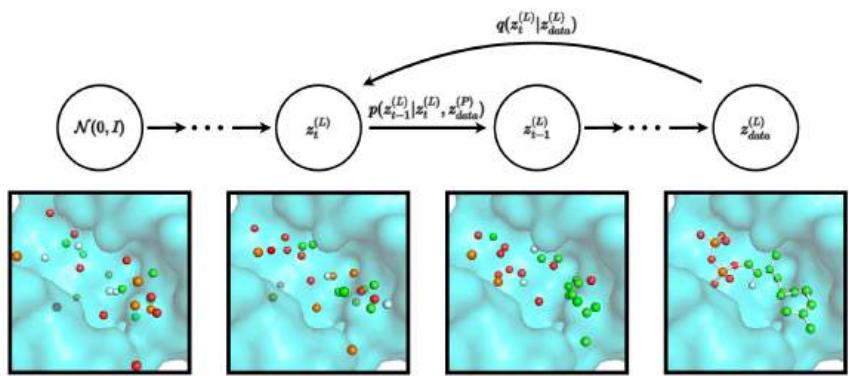
Gainza et al., Correia 2022

Binding spike protein in multiple virus variants



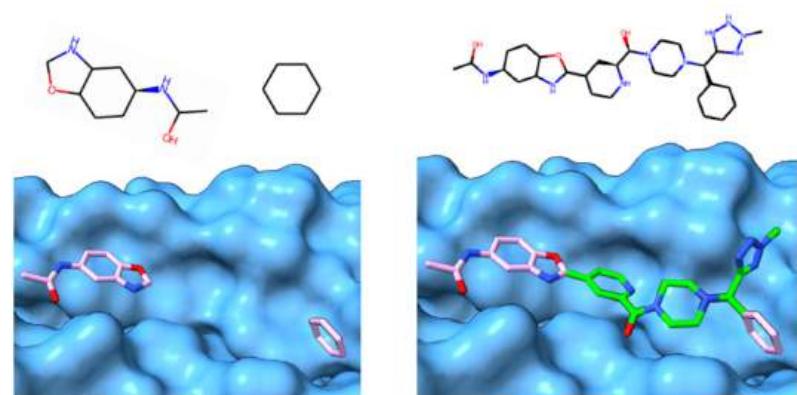
Pseudovirus neutralisation

Structure-Based Drug Design



Equivariant diffusion model conditioned
on target pocket 3D structure

Schneuing et B, Welling, Correia 2022



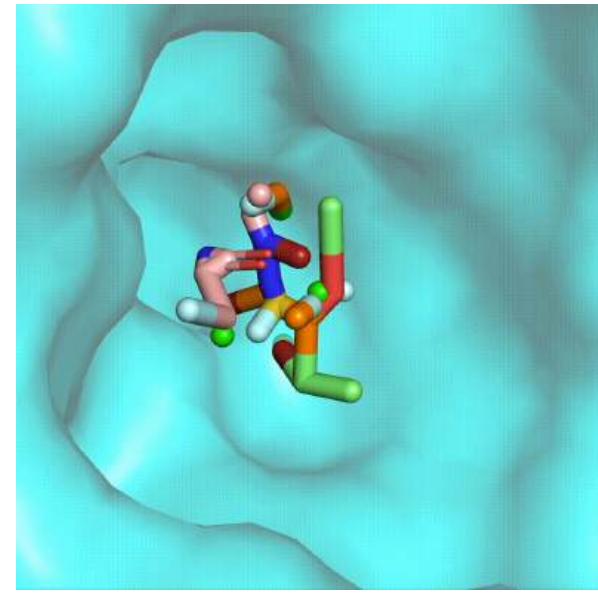
DiffLinker: Impainting of molecular
fragments conditioned on target pocket

Igashov et B, Welling, Correia 2022



“Painting of an astronaut
riding a dog on the Moon”

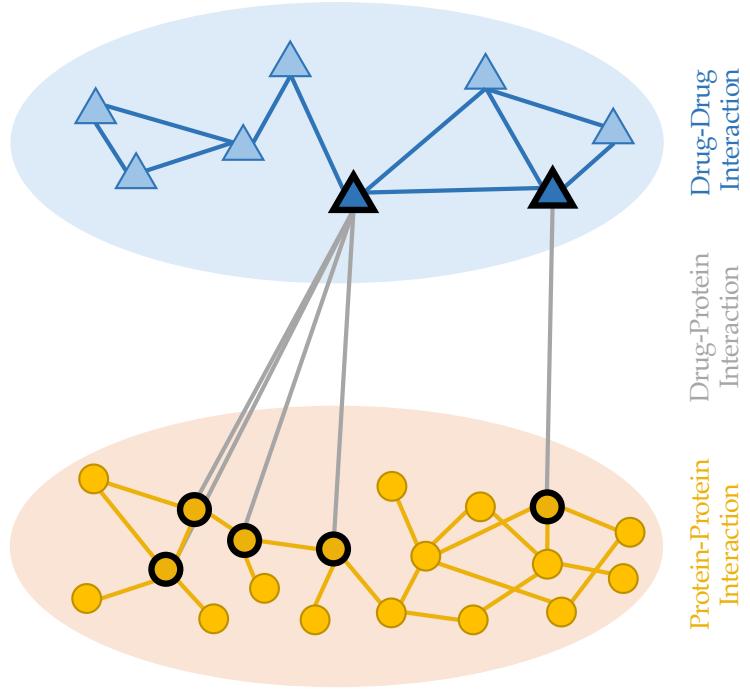
DALL-E 2023 (prompt by B)



“Drug-like molecule
binding a protein pocket”

Schneuing et Welling, B, Correia 2022

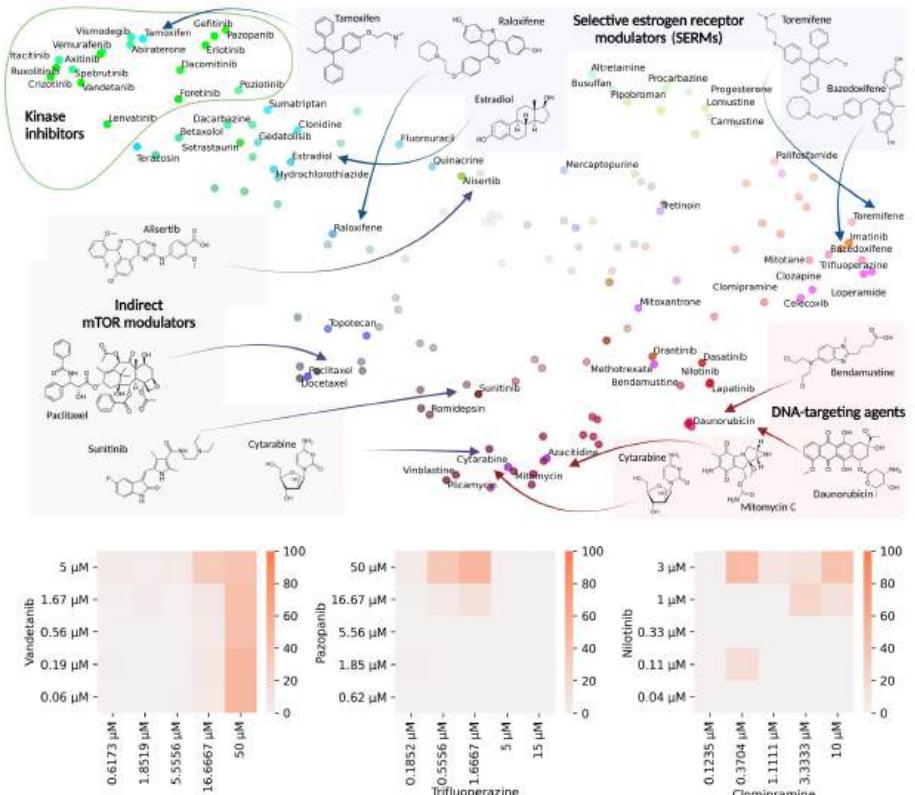
Drug Repositioning



Drug side effect prediction

Zitnik et al. 2018

Drug-Drug Interaction
Drug-Protein Interaction
Protein-Protein Interaction



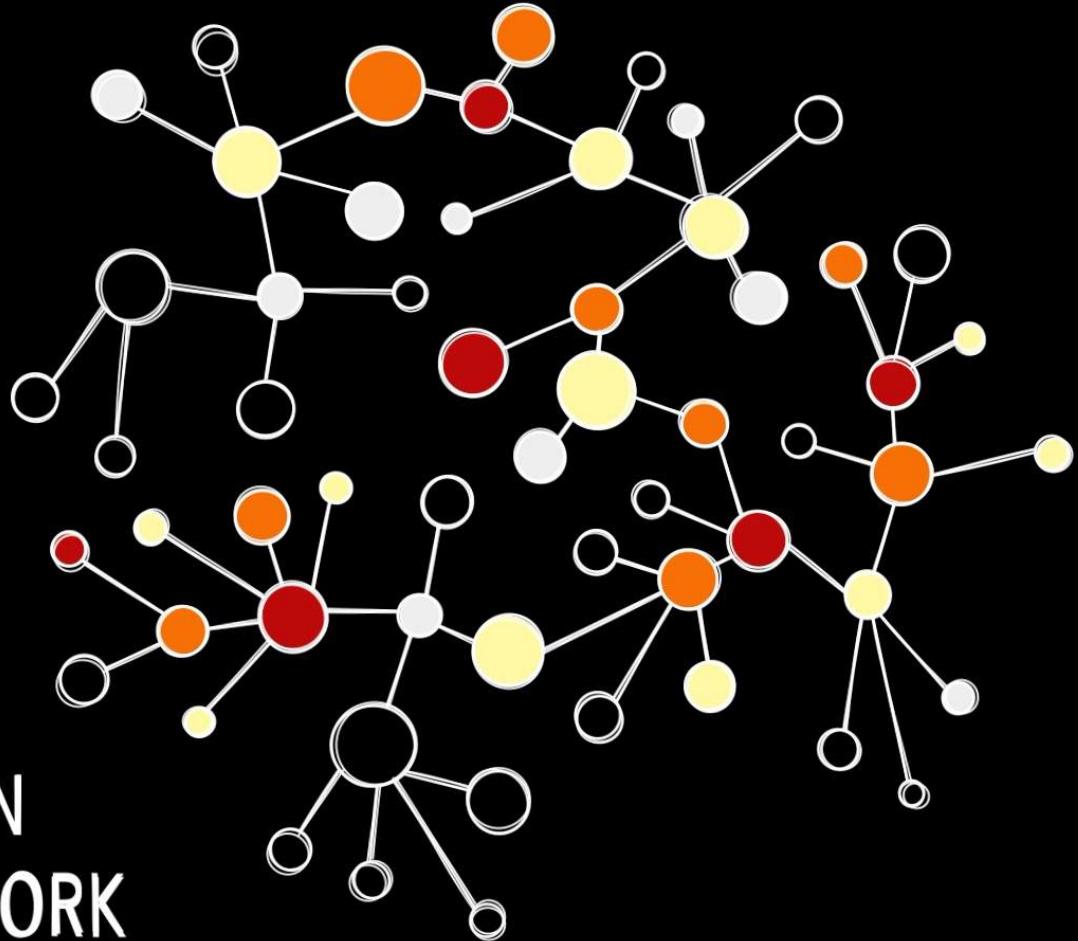
Drug synergy prediction

Bertin et B, Bengio 2021 (RECOVER & Gates Foundation)



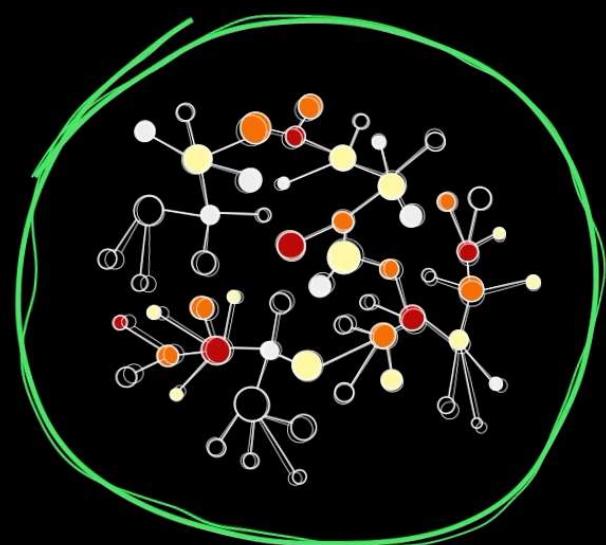
Veselkov et B 2019

Hyperfoods



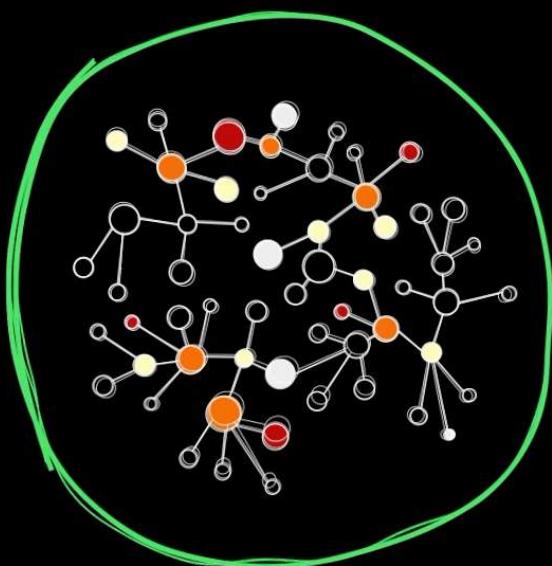
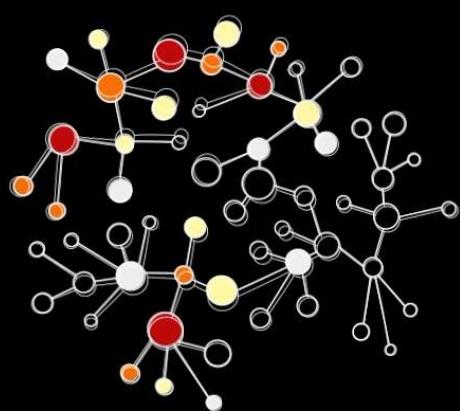
PROTEIN-PROTEIN INTERACTION NETWORK

Veselkov et B 2019



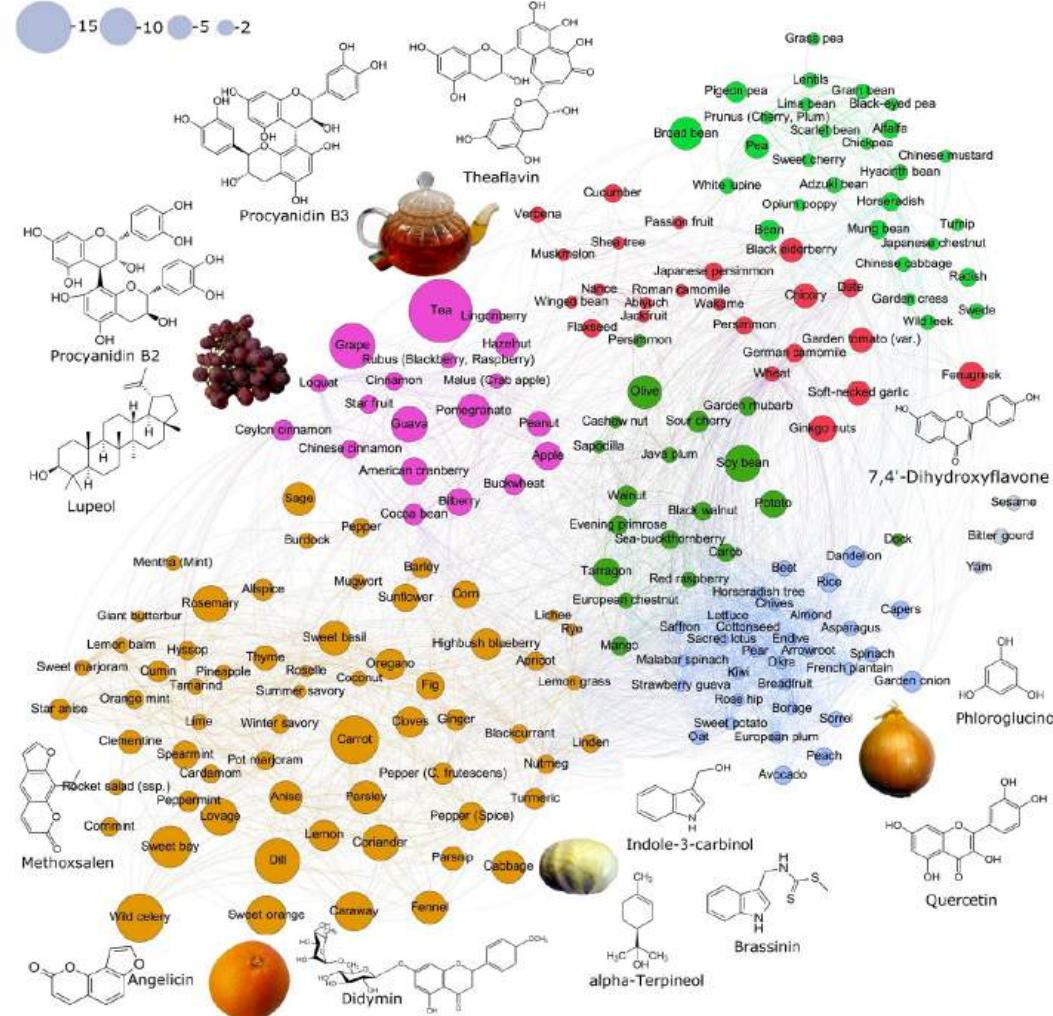
ANTICANCER

non ANTICANCER



ANTICANCER

Hyperfoods



Veselkov et B 2019

Video: Vodafone Foundation / Recipes: Bruno Barbieri
Based on Laponogov et B 2020

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

—Claude Adrien Helvétius



Thank you!