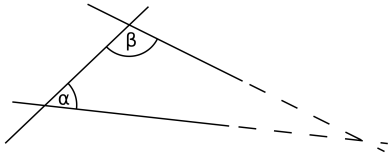


Going back in time...



Going back in time...

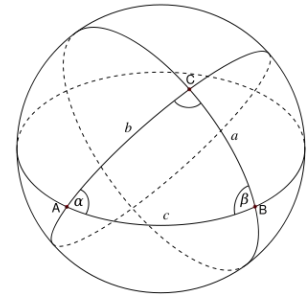
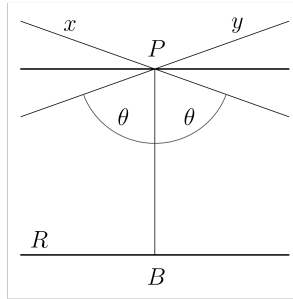
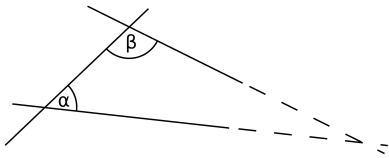


Euclid

~300 B.C.



Going back in time...



Euclid

~300 B.C.



Lobachevsky

~1829



Bolyai

~1830

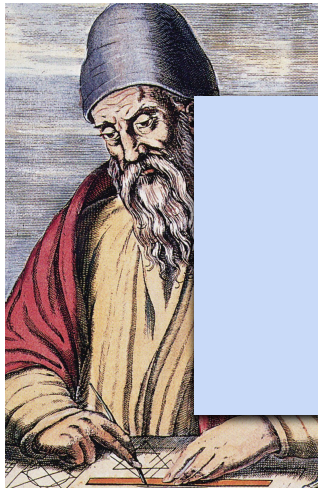
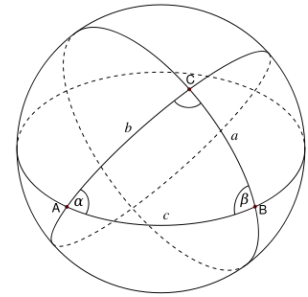
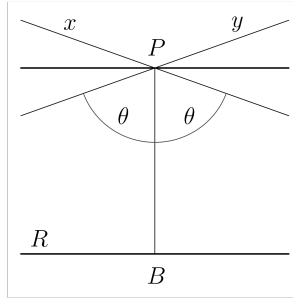
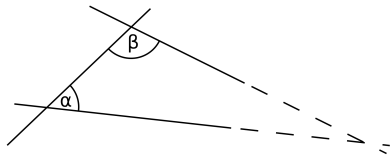


Riemann

~1856

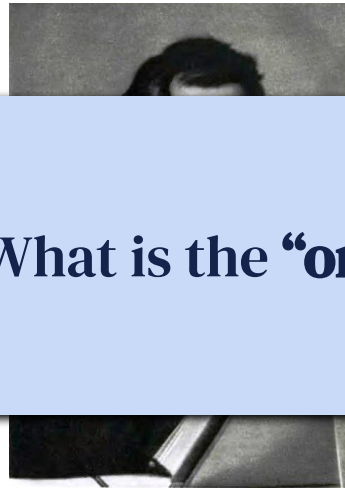


Going back in time...



Euclid

~300 B.C.



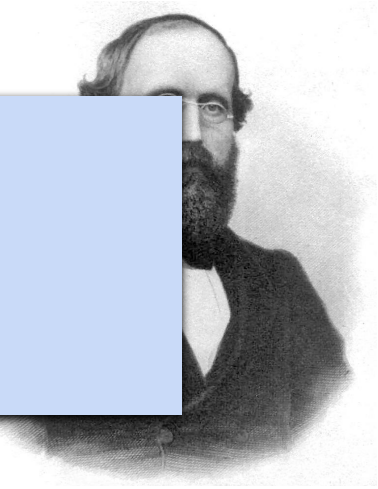
Lobachevsky

~1829



Bolyai

~1830



Riemann

~1856

What is the “one true geometry”?



Felix Klein and the “Erlangen Program” (1872)



Vergleichende Betrachtungen
über
neuere geometrische Forschungen
von
Dr. Felix Klein,
o. ö. Professor der Mathematik an der Universität Erlangen.

Programm
zum Eintritt in die philosophische Facultät und den Senat
der k. Friedrich-Alexanders-Universität
zu Erlangen.

Erlangen.
Verlag von Andreas Deichert
1872.



Felix Klein and the “Erlangen Program” (1872)



Vergleichende Betrachtungen
über
neuere geometrische Forschungen
von
Dr. Felix Klein,
o. ö. Professor der Mathematik an der Universität Erlangen.

Programm
zum Eintritt in die philosophische Facultät und den Senat
der k. Friedrich-Alexanders-Universität
zu Erlangen.

Blueprint for **unifying** geometries

Lens of **invariances** and **symmetries**

Formalised in the language of **group theory!**



Impact of the Erlangen Program

- Strong impact on **geometry**; no longer hunting for the “one true geometry”
 - Formalised by *Élie Cartan* in the 1920s
- *Physics: Noether’s Theorem* (all conservation laws derivable from symmetry!)
 - Even enabled the classification of elementary particles (*irreducible representations*)
- *Category Theory*
 - “can be regarded as a continuation of the Klein Erlangen Program, in the sense that a geometrical space with its group of transformations is generalized to a category with its algebra of mappings”; – Eilenber, Lane (creators of CT)



Deep learning, circa 2020

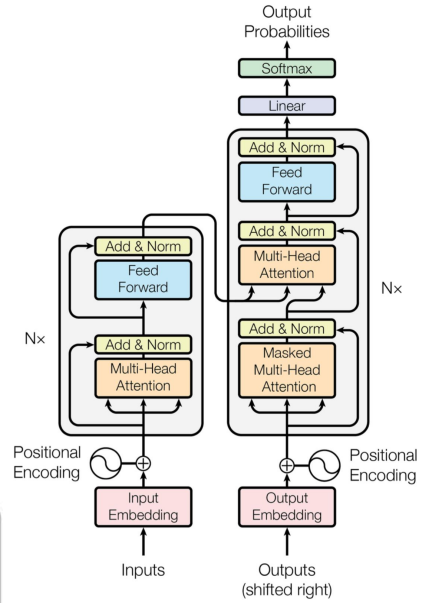
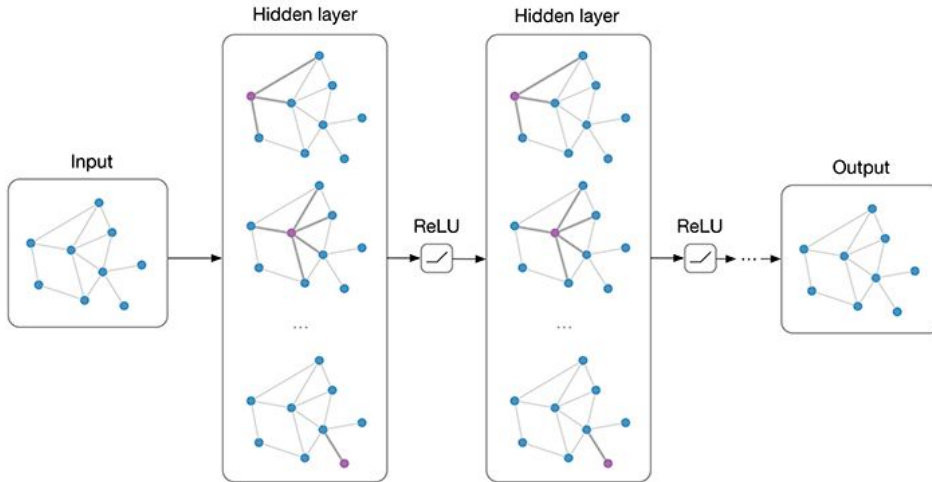
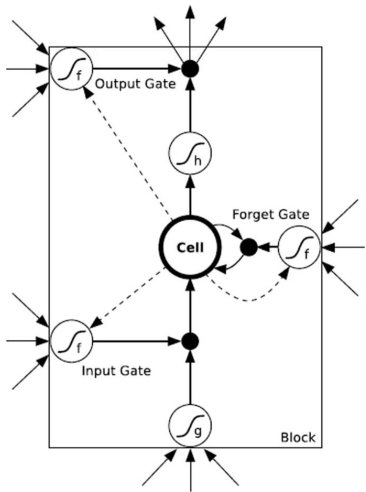
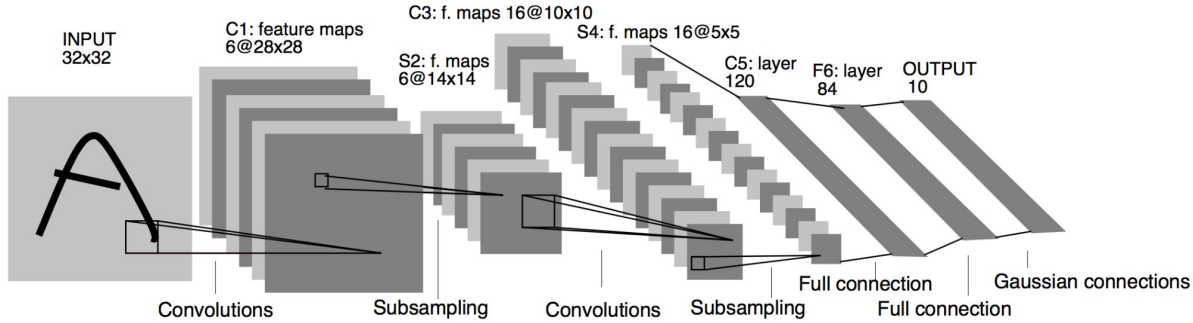


Figure 1: The Transformer - model architecture.



Deep learning, circa 2020

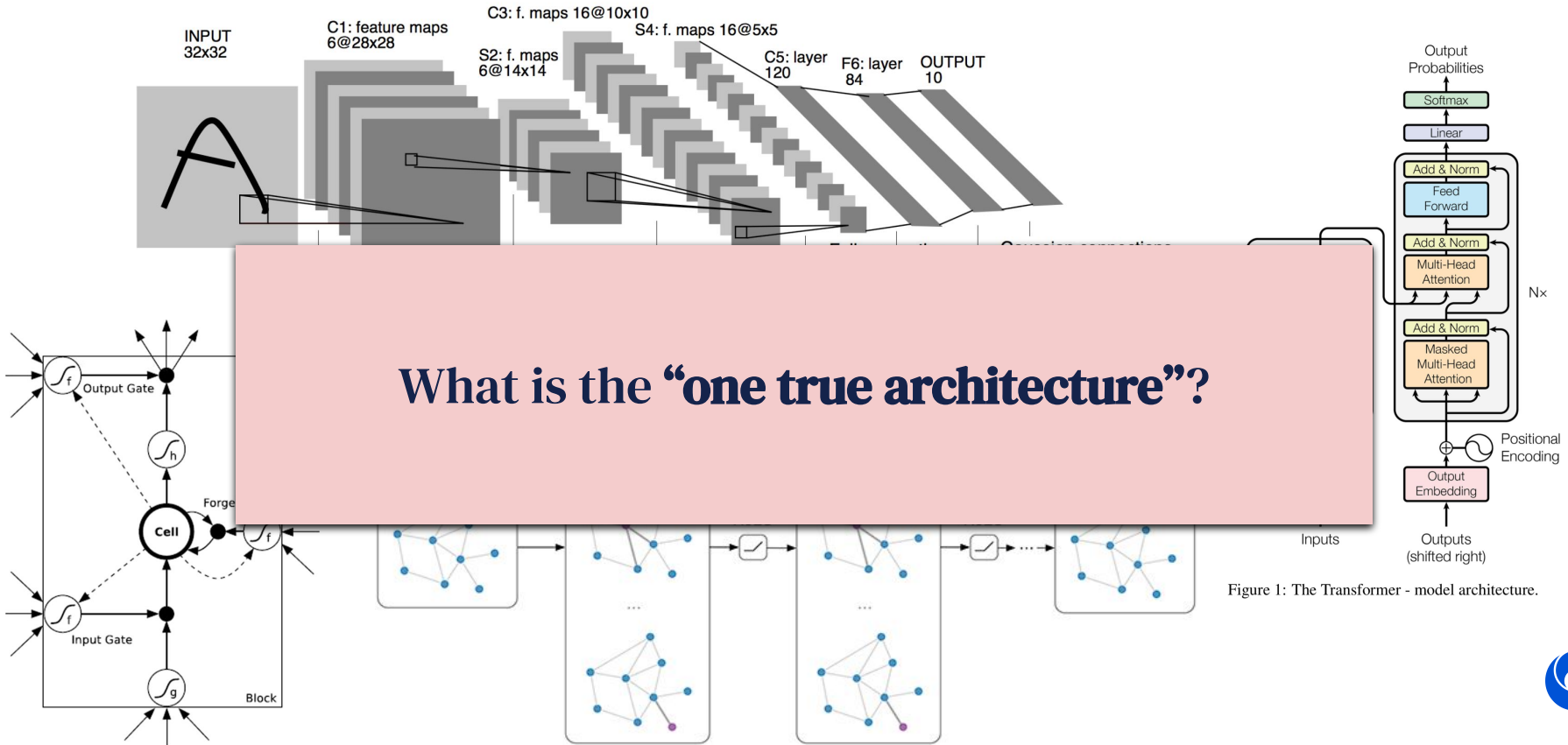
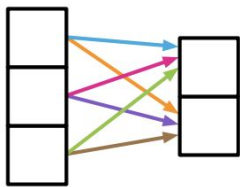


Figure 1: The Transformer - model architecture.

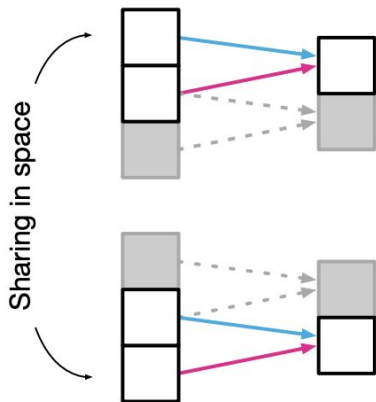


Could GNNs be the answer?

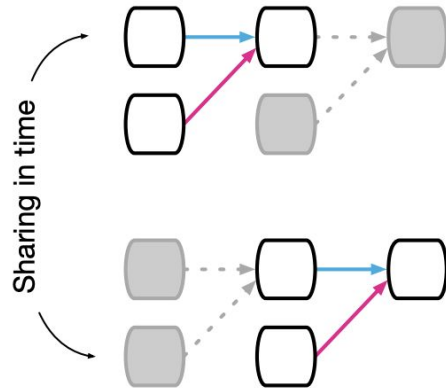
- “If we squint hard enough”, (m)any NNs can be seen as message passing over a graph
 - Further, most data we receive from **nature** is inherently graph-structured
 - So, GNNs likely play a part in the “one true architecture” (motivating this course)



(a) Fully connected



(b) Convolutional



(c) Recurrent

- But to formalise this, we need to understand GNNs **beyond** permutation equivariance!



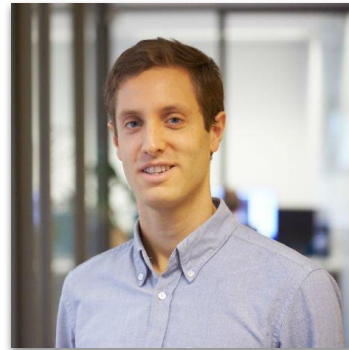
...now it's our turn to study geometry :)



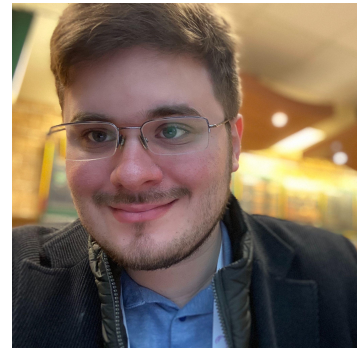
Michael Bronstein
Oxford / Twitter



Joan Bruna
NYU



Taco Cohen
Qualcomm



Petar Veličković
DeepMind / Cambridge



DeepMind

Geometric Deep Learning

GNNs Beyond Permutation Equivariance

Petar Veličković

Stanford University
CS224W
30 November 2021

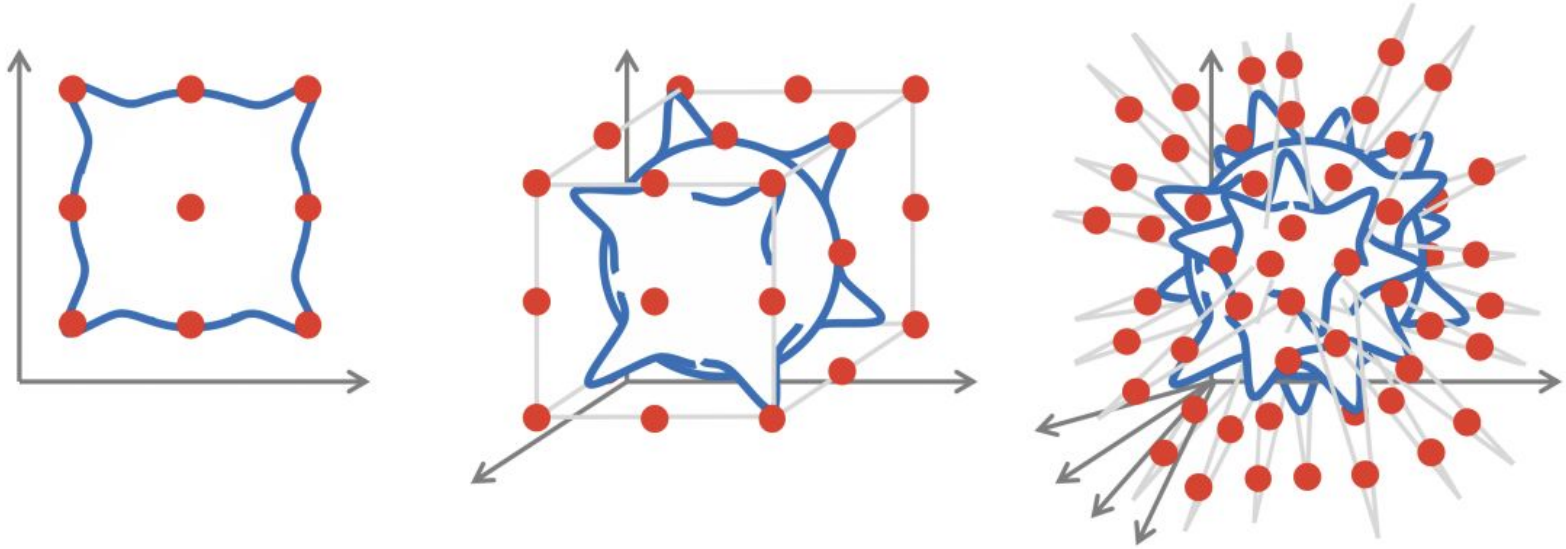


1

Learning in
high dimensions
is hard



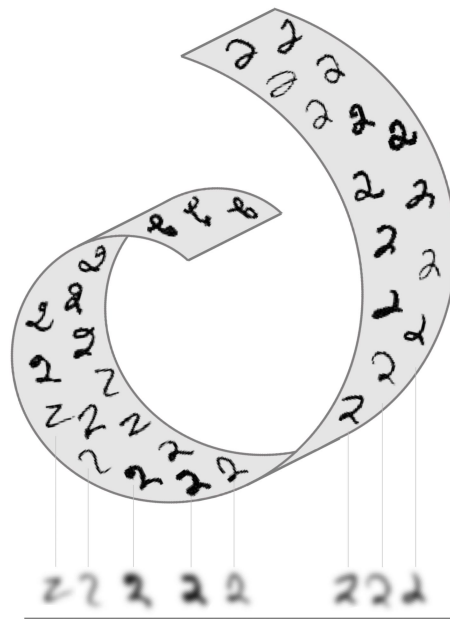
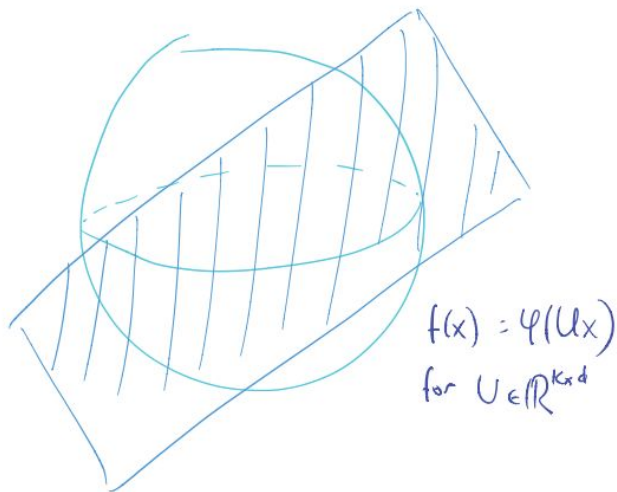
The Curse of Dimensionality



Even “simple” target functions (e.g. 1-Lipschitz) require **exponential** samples in nb. dimensions.



Low-dimensional projections don't necessarily help!



Shallow MLPs can lose a lot of the **fidelity** of their inputs.

What can we do?



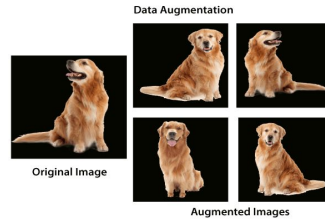
2

Symmetries, Groups and Invariances

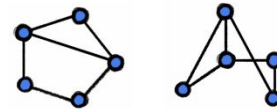
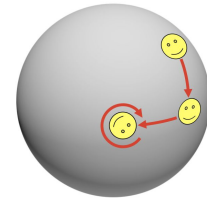


Geometry to the rescue!

- We can inject further assumptions about the **geometry** of through *inductive biases*
 - Restrict the functions in our hypothesis space to ones that *respect* the geometry.
 - This can make the high-dimensional problem more tractable!



- Some popular examples:
 - **Image** data should be processed independently of **shifts**
 - **Spherical** data should be processed independently of **rotations**
 - **Graph** data should be processed independently of **isomorphism**



- We will now attempt to **formalise** this!



Key elements!

- We assume data lives on a **domain**, Ω
 - e.g. for images, $u \in \Omega$ are *pixels*; for graphs they are *nodes*
- We assume a **feature space**, \mathbf{C} , to be stored in elements of a domain
 - For our purposes, $\mathbf{C} = \mathbb{R}^k$
- We can then define *featurised domains* using a space of **feature functions** $\mathbf{X}(\Omega, \mathbf{C})$
 - $x \in \mathbf{X}(\Omega, \mathbf{C})$ is a function s.t. $x(u) \in \mathbf{C}$ gives features of element $u \in \Omega$
 - For discrete environments we can think of \mathbf{X} as a *feature matrix* ($\mathbf{X} \in \mathbb{R}^{|\Omega| \times k}$)



Principle 1: Symmetry groups

- Symmetry is a transformation that leaves an object **invariant** (i.e. *unchanged*)
 - Hence they must be **composable**, **invertible**, contain **identity**...
- In fact, they can be reasoned about using a very elegant mathematical object: the **group**
- Elements of these groups are domain **transformations** (e.g. some functions $g : \Omega \rightarrow \Omega$).



Principle 1: Symmetry groups

- Symmetry is a transformation that leaves an object **invariant** (i.e. *unchanged*)
 - Hence they must be **composable**, **invertible**, contain **identity**...
- In fact, they can be reasoned about using a very elegant mathematical object: the **group**

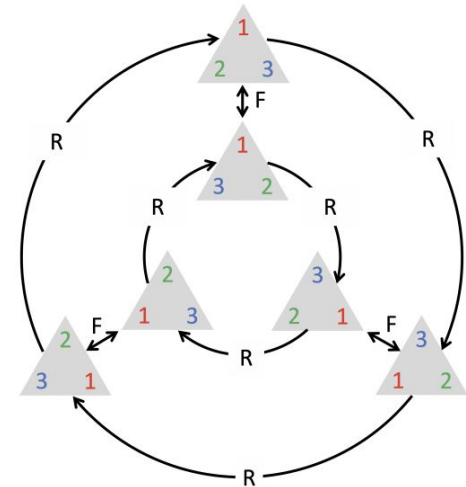
A *group* is a set \mathcal{G} along with a binary operation $\circ : \mathcal{G} \times \mathcal{G} \rightarrow \mathcal{G}$ called *composition* (for brevity, denoted by juxtaposition $g \circ h = gh$) satisfying the following axioms:

Associativity: $(gh)k = g(hk)$ for all $g, h, k \in \mathcal{G}$.

Identity: there exists a unique $e \in \mathcal{G}$ satisfying $eg = ge = g$ for all $g \in \mathcal{G}$.

Inverse: For each $g \in \mathcal{G}$ there is a unique inverse $g^{-1} \in \mathcal{G}$ such that $gg^{-1} = g^{-1}g = e$.

Closure: The group is closed under composition, i.e., for every $g, h \in \mathcal{G}$, we have $gh \in \mathcal{G}$.



Principle 1: Symmetry groups

- Symmetry is a transformation that leaves an object **invariant** (i.e. *unchanged*)
 - Hence they must be **composable, invertible, contain identity...**
- In fact, they can be reasoned about using a very elegant mathematical object: the **group**

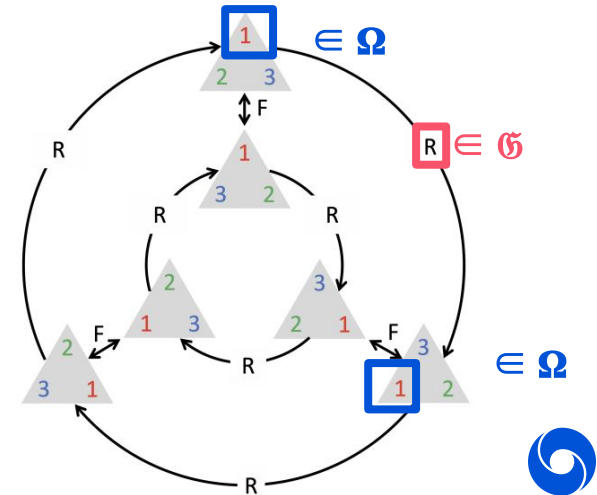
A *group* is a set \mathcal{G} along with a binary operation $\circ : \mathcal{G} \times \mathcal{G} \rightarrow \mathcal{G}$ called *composition* (for brevity, denoted by juxtaposition $g \circ h = gh$) satisfying the following axioms:

Associativity: $(gh)k = g(hk)$ for all $g, h, k \in \mathcal{G}$.

Identity: there exists a unique $e \in \mathcal{G}$ satisfying $eg = ge = g$ for all $g \in \mathcal{G}$.

Inverse: For each $g \in \mathcal{G}$ there is a unique inverse $g^{-1} \in \mathcal{G}$ such that $gg^{-1} = g^{-1}g = e$.

Closure: The group is closed under composition, i.e., for every $g, h \in \mathcal{G}$, we have $gh \in \mathcal{G}$.



Group actions

- We are interested in how these groups affect data
 - Group action $(\mathfrak{g}, u) \mapsto \mathfrak{g}.u$ for a group element \mathfrak{g} , and a domain element u
 - E.g. translating or rotating an image, or permuting a set
- We will be interested in **linear** group actions: $\mathfrak{g}.\alpha x + \beta x' = \alpha(\mathfrak{g}.x) + \beta(\mathfrak{g}.x')$
- This also allows us to represent group actions using **linear algebra**; $\rho : \mathfrak{G} \rightarrow \mathbb{R}^{n \times n}$

A n -dimensional real *representation* of a group \mathfrak{G} is a map $\rho : \mathfrak{G} \rightarrow \mathbb{R}^{n \times n}$, assigning to each $\mathfrak{g} \in \mathfrak{G}$ an *invertible* matrix $\rho(\mathfrak{g})$, and satisfying the condition $\rho(\mathfrak{g}\mathfrak{h}) = \rho(\mathfrak{g})\rho(\mathfrak{h})$ for all $\mathfrak{g}, \mathfrak{h} \in \mathfrak{G}$. A representation is called *unitary* or *orthogonal* if the matrix $\rho(\mathfrak{g})$ is unitary or orthogonal for all $\mathfrak{g} \in \mathfrak{G}$.



Invariance and equivariance

- We can largely simplify **high-dimensional learning** by *exploiting* the symmetries in \mathfrak{G} !

A function $f : \mathcal{X}(\Omega) \rightarrow \mathcal{Y}$ is \mathfrak{G} -*invariant* if $f(\rho(\mathfrak{g})x) = f(x)$ for all $\mathfrak{g} \in \mathfrak{G}$, i.e., its output is unaffected by the group action on the input.

e.g. **image classification**: output class (*likely?*) won't depend on image **shifts**

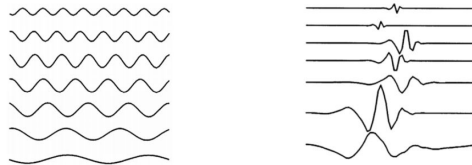
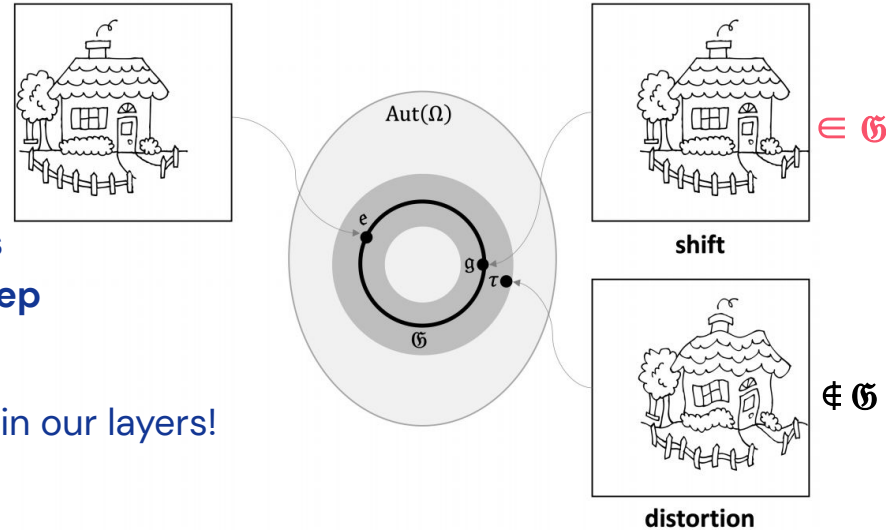
A function $f : \mathcal{X}(\Omega) \rightarrow \mathcal{X}(\Omega)$ is \mathfrak{G} -*equivariant* if $f(\rho(\mathfrak{g})x) = \rho(\mathfrak{g})f(x)$ for all $\mathfrak{g} \in \mathfrak{G}$, i.e., group action on the input affects the output in the same way.

e.g. **image segmentation**: segmentation mask must **follow** any shifts in the input



Principle 2: Scale separation

- Want signal to be **stable** under slight *deformations* of the domain
- We derive: highly beneficial to compose **local** operations to model larger-scale ones
 - local ops won't globally propagate errors
 - e.g. CNNs with 3 x 3 kernels, but **very deep**
- Accordingly, we would like to support **locality** in our layers!
- cf. Fourier Transform vs. Wavelets



DeepMind

3

The Blueprint of Geometric Deep Learning



The key “building blocks” of Geometric Deep Learning

*Linear \mathfrak{G} -equivariant layer $B : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{X}(\Omega', \mathcal{C}')$,
satisfying $B(\mathfrak{g}.x) = \mathfrak{g}.B(x)$ for all $\mathfrak{g} \in \mathfrak{G}$ and $x \in \mathcal{X}(\Omega, \mathcal{C})$.*

Nonlinearity $\sigma : \mathcal{C} \rightarrow \mathcal{C}'$ applied element-wise as $(\sigma(x))(u) = \sigma(x(u))$.

Local pooling (coarsening) $P : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{X}(\Omega', \mathcal{C})$, such that $\Omega' \subseteq \Omega$.

*\mathfrak{G} -invariant layer (global pooling) $A : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{Y}$,
satisfying $A(\mathfrak{g}.x) = A(x)$ for all $\mathfrak{g} \in \mathfrak{G}$ and $x \in \mathcal{X}(\Omega, \mathcal{C})$.*



The key “building blocks” of Geometric Deep Learning

Equivariant local layers

*Linear \mathfrak{G} -equivariant layer $B : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{X}(\Omega', \mathcal{C}')$,
satisfying $B(g \cdot x) = g \cdot B(x)$ for all $g \in \mathfrak{G}$ and $x \in \mathcal{X}(\Omega, \mathcal{C})$.*

Nonlinearity $\sigma : \mathcal{C} \rightarrow \mathcal{C}'$ applied element-wise as $(\sigma(x))(u) = \sigma(x(u))$.

Local pooling (coarsening) $P : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{X}(\Omega', \mathcal{C})$, such that $\Omega' \subseteq \Omega$.

*\mathfrak{G} -invariant layer (global pooling) $A : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{Y}$,
satisfying $A(g \cdot x) = A(x)$ for all $g \in \mathfrak{G}$ and $x \in \mathcal{X}(\Omega, \mathcal{C})$.*

Invariant “tail” (if necessary)



The key “building blocks” of Geometric Deep Learning

*Linear \mathfrak{G} -equivariant layer $B : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{X}(\Omega', \mathcal{C}')$,
satisfying $B(g \cdot x) = g \cdot B(x)$ for all $g \in \mathfrak{G}$ and $x \in \mathcal{X}(\Omega, \mathcal{C})$.*

Activation function

Nonlinearity $\sigma : \mathcal{C} \rightarrow \mathcal{C}'$ applied element-wise as $(\sigma(x))(u) = \sigma(x(u))$.

(necessary for deep learning! :))

Local pooling (coarsening) $P : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{X}(\Omega', \mathcal{C})$, such that $\Omega' \subseteq \Omega$.

*\mathfrak{G} -invariant layer (global pooling) $A : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{Y}$,
satisfying $A(g \cdot x) = A(x)$ for all $g \in \mathfrak{G}$ and $x \in \mathcal{X}(\Omega, \mathcal{C})$.*



The key “building blocks” of Geometric Deep Learning

*Linear \mathfrak{G} -equivariant layer $B : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{X}(\Omega', \mathcal{C}')$,
satisfying $B(\mathfrak{g}.x) = \mathfrak{g}.B(x)$ for all $\mathfrak{g} \in \mathfrak{G}$ and $x \in \mathcal{X}(\Omega, \mathcal{C})$.*

Nonlinearity $\sigma : \mathcal{C} \rightarrow \mathcal{C}'$ applied element-wise as $(\sigma(x))(u) = \sigma(x(u))$.

Coarsening layer

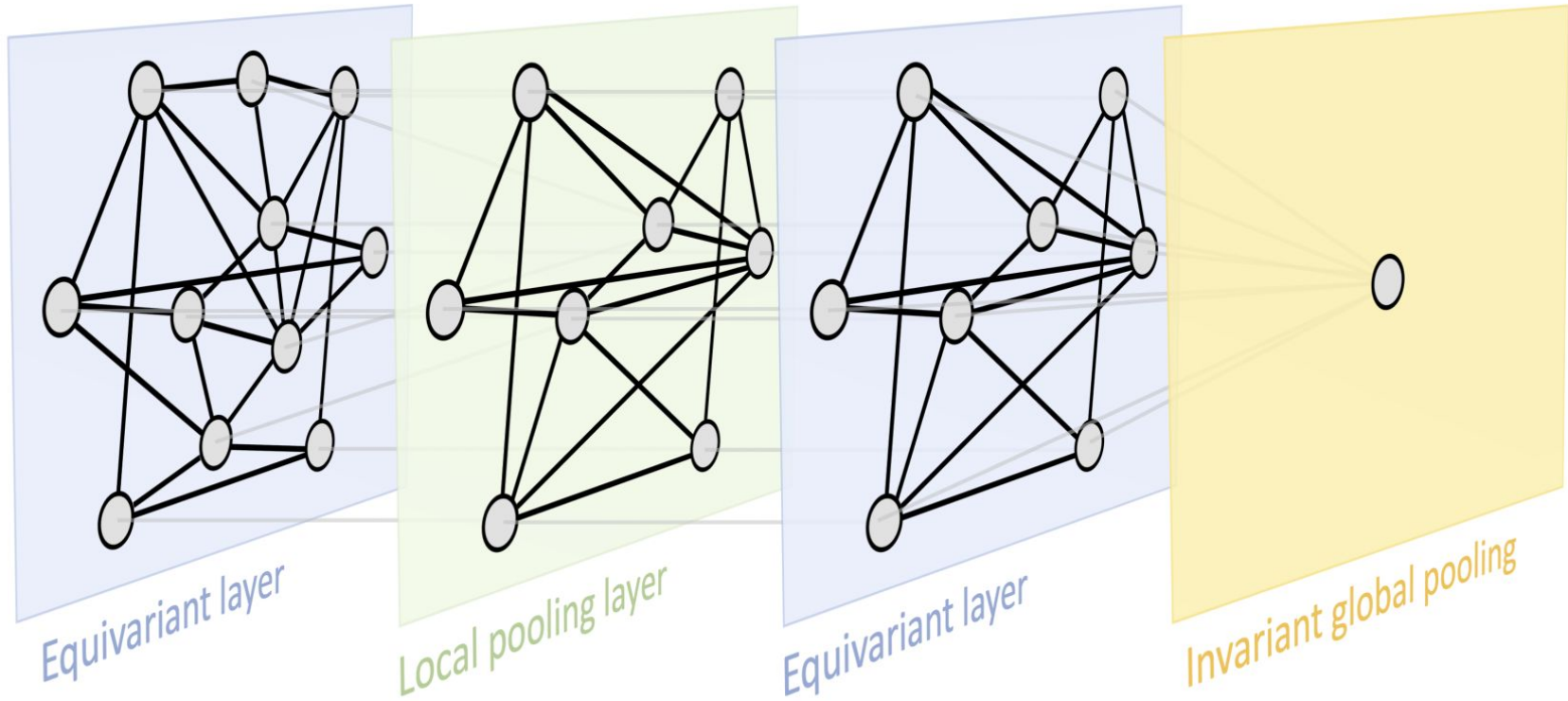
Local pooling (coarsening) $P : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{X}(\Omega', \mathcal{C})$, such that $\Omega' \subseteq \Omega$.

(Not covered here in detail, but

*\mathfrak{G} -invariant layer (global pooling) $A : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{Y}$, follows from scale separation!)
satisfying $A(\mathfrak{g}.x) = A(x)$ for all $\mathfrak{g} \in \mathfrak{G}$ and $x \in \mathcal{X}(\Omega, \mathcal{C})$.*



All you need to build the architectures that are all you need :)



All the fan-favourites are easily derivable

...with many (potentially) unexpected and useful **extras!** :)

Architecture	Domain Ω	Symmetry group \mathcal{G}
<i>CNN</i>	Grid	Translation
<i>Spherical CNN</i>	Sphere / $SO(3)$	Rotation $SO(3)$
<i>Intrinsic / Mesh CNN</i>	Manifold	Isometry $Iso(\Omega)$ / Gauge symmetry $SO(2)$
<i>GNN</i>	Graph	Permutation Σ_n
<i>Deep Sets</i>	Set	Permutation Σ_n
<i>Transformer</i>	Complete Graph	Permutation Σ_n
<i>LSTM</i>	1D Grid	Time warping



DeepMind

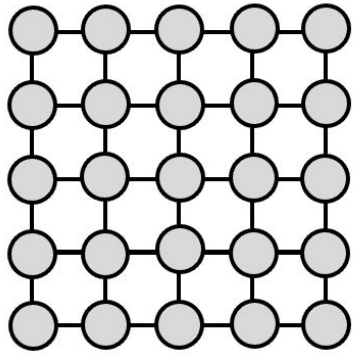
4

The “5G” of Geometric Deep Learning

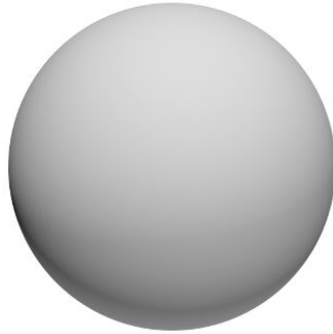


The “5G” of Geometric Deep Learning

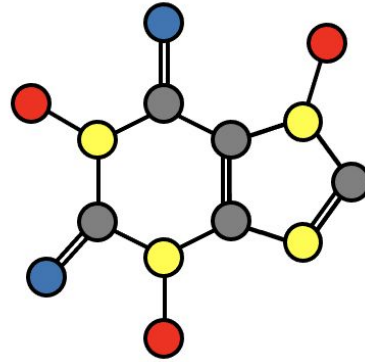
We will use the remainder of this lecture to study a few interesting **instances** of this blueprint



Grids



Groups



Graphs

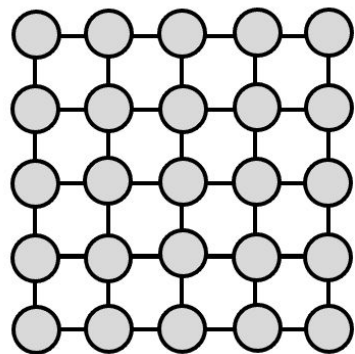


**Geodesics &
Gauges**

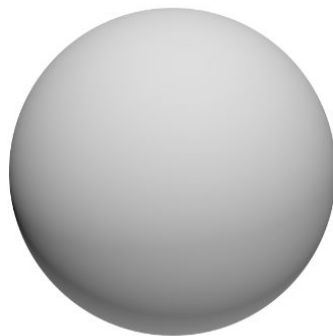


The “5G” of Geometric Deep Learning

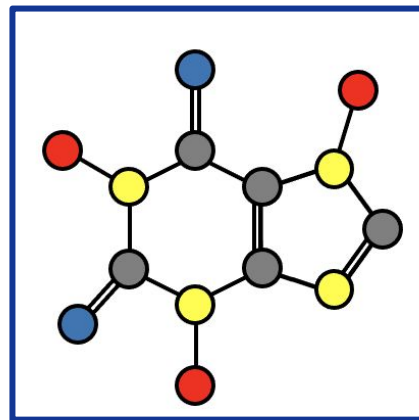
We will use the remainder of this lecture to study a few interesting **instances** of this blueprint



**Images &
Sequences**



**Homogeneous
spaces**



Graphs & Sets



**Manifolds, Meshes &
Geometric graphs**

Throughout CS224W, the focus was primarily on the domain of **graphs**.

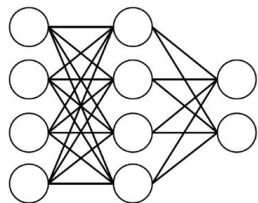


Our strategy for the rest of the lecture

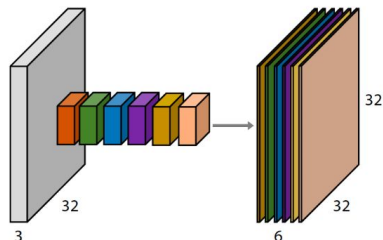
- We start by seeing how **GNNs** fit in this paradigm
 - This will involve re-deriving / re-introducing some concepts you've seen in the course
- Then, we will see how we can use the blueprint to **expand** GNNs into other domains
 - Also, it will give us an insight into “the world beyond”
 - GNNs **beyond** permutation equivariance!
- Our discussion will span many architectures you (may) know of :)
 - Deep Sets
 - Transformers
 - CNNs
 - Spherical CNNs
 - Mesh CNNs



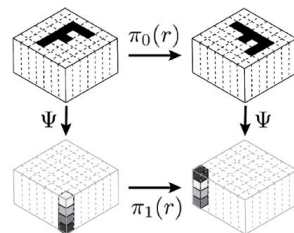
Architectures of interest



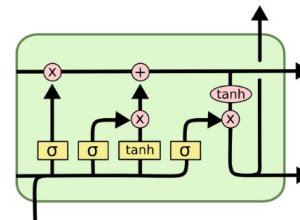
Perceptrons
Function regularity



CNNs
Translation



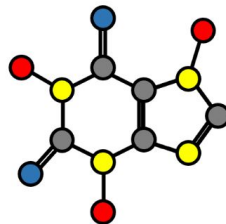
Group-CNNs
Translation+Rotation,
Global groups



LSTMs
Time warping



DeepSets / Transformers
Permutation



GNNs
Permutation



Intrinsic CNNs
Isometry / Gauge choice



DeepMind

5

Geometric DL Perspective on Graph Neural Networks



Learning on sets: Setup

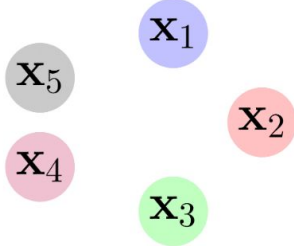
- For now, assume the graph **has no edges** (e.g. our domain is just the set of nodes, V).
- Let $\mathbf{x}_i \in \mathbb{R}^k$ be the features of node i .
- We can stack them into a node feature matrix of shape $n \times k$:

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$$

- That is, the i th row of \mathbf{X} corresponds to \mathbf{x}_i
- Note that, by doing so, we have specified a **node ordering!**
 - We would like the result of any neural networks to not depend on this.



What do we want?



What do we want?

$$f \left(\begin{array}{ccc} & \text{x}_1 & \\ \text{x}_5 & & \\ & & \text{x}_2 \\ \text{x}_4 & & \\ & \text{x}_3 & \end{array} \right) = \mathbf{y}$$



What do we want?

$$f \left(\begin{array}{ccc} \text{x}_5 & \text{x}_1 & \\ \text{x}_4 & & \text{x}_2 \\ & \text{x}_3 & \end{array} \right) = \mathbf{y} = f \left(\begin{array}{ccc} & \text{x}_2 & \\ \text{x}_5 & & \text{x}_4 \\ \text{x}_1 & & \text{x}_3 \end{array} \right)$$



Permutations and permutation matrices

- It will be useful to think about the operations that **change** the node order
 - Such operations are known as **permutations** (there are $n!$ of them)
 - e.g. a permutation (2, 4, 1, 3) means $\mathbf{y}_1 \leftarrow \mathbf{x}_{2'}$, $\mathbf{y}_2 \leftarrow \mathbf{x}_{4'}$, $\mathbf{y}_3 \leftarrow \mathbf{x}_{1'}$, $\mathbf{y}_4 \leftarrow \mathbf{x}_{3'}$.
- To stay within linear algebra, each permutation defines an $n \times n$ **matrix (group action!)**
 - Such matrices are called **permutation matrices**
 - They have exactly one 1 in every row and column, and zeros everywhere else
 - Their effect when left-multiplied is to permute the rows of \mathbf{X} , like so:

$$\mathbf{P}_{(2,4,1,3)}\mathbf{X} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \text{---} & \mathbf{x}_1 & \text{---} \\ \text{---} & \mathbf{x}_2 & \text{---} \\ \text{---} & \mathbf{x}_3 & \text{---} \\ \text{---} & \mathbf{x}_4 & \text{---} \end{bmatrix} = \begin{bmatrix} \text{---} & \mathbf{x}_2 & \text{---} \\ \text{---} & \mathbf{x}_4 & \text{---} \\ \text{---} & \mathbf{x}_1 & \text{---} \\ \text{---} & \mathbf{x}_3 & \text{---} \end{bmatrix}$$



Permutation *invariance*

- We want to design functions $f(\mathbf{X})$ over sets that will not depend on the order
- Equivalently, applying a permutation matrix shouldn't modify the result!
- We arrive at a useful notion of permutation invariance. We say that $f(\mathbf{X})$ is permutation *invariant* if, for *all* permutation matrices \mathbf{P} :

$$f(\mathbf{PX}) = f(\mathbf{X})$$

- One very generic form is the *Deep Sets* model (Zaheer et al., NeurIPS'17): $f(\mathbf{X}) = \phi \left(\sum_{i \in \mathcal{V}} \psi(\mathbf{x}_i) \right)$ where ψ and ϕ are (learnable) functions, e.g. MLPs.
 - The **sum** aggregation is *critical!* (other choices possible, e.g. **max** or **avg**)



Permutation *equivariance*

- Permutation-*invariant* models are a good way to obtain set-level outputs
- What if we would like answers at the **node** level?
 - We want to still be able to **identify** node outputs, which a permutation-invariant aggregator would destroy!
- We may instead seek functions that don't **change** the node order
 - i.e. if we permute the nodes, it doesn't matter if we do it **before** or **after** the function!
- Accordingly, we say that $f(\mathbf{X})$ is permutation equivariant if, for all permutation matrices \mathbf{P} :

$$f(\mathbf{P}\mathbf{X}) = \mathbf{P}f(\mathbf{X})$$



Learning on graphs

- Now we augment the set of nodes with **edges** between them.
 - That is, we consider general $E \subseteq V \times V$.
- We can represent these edges with an **adjacency matrix**, **A**, such that:

$$a_{ij} = \begin{cases} 1 & (i, j) \in \mathcal{E} \\ 0 & \text{otherwise} \end{cases}$$

- Further additions (e.g. *edge features*) are possible but **ignored** for simplicity.
- Our main desiderata (*permutation {in,equi}variance*) still hold!



What's changed?

$$f \left(\begin{array}{ccc} \text{x}_5 & \text{x}_1 & \\ \text{x}_4 & & \text{x}_2 \\ & \text{x}_3 & \end{array} \right) = \mathbf{y} = f \left(\begin{array}{ccc} & \text{x}_2 & \\ \text{x}_5 & & \text{x}_4 \\ \text{x}_1 & & \text{x}_3 \end{array} \right)$$



What's changed?

$$f \left(\begin{array}{cc} \text{x}_5 & \text{x}_1 \\ \text{x}_4 & \text{x}_2 \\ & \text{x}_3 \end{array} \right) = \mathbf{y} = f \left(\begin{array}{cc} & \text{x}_2 \\ \text{x}_5 & \text{x}_4 \\ \text{x}_1 & \text{x}_3 \end{array} \right)$$

$$f \left(\begin{array}{ccc} \text{x}_5 & \text{x}_1 & \text{x}_2 \\ | & / & / \\ \text{x}_4 & & \text{x}_2 \\ | & \backslash & \backslash \\ \text{x}_4 & \text{x}_3 & \text{x}_2 \end{array} \right) = \mathbf{y} = f \left(\begin{array}{ccc} & \text{x}_2 & \\ / & | & / \\ \text{x}_1 & \text{x}_5 & \text{x}_4 \\ \backslash & \backslash & | \\ \text{x}_1 & \text{x}_3 & \text{x}_4 \end{array} \right)$$



Permutation invariance and equivariance on graphs

- The main difference: node permutations now also accordingly act on the **edges**
- We need to appropriately permute both **rows** and **columns** of **A**
 - When applying a permutation matrix **P**, this amounts to **PAP[⊤]**
- We arrive at updated definitions of suitable functions $f(\mathbf{X}, \mathbf{A})$ over graphs:

Invariance: $f(\mathbf{PX}, \mathbf{PAP}^{\top}) = f(\mathbf{X}, \mathbf{A})$

Equivariance: $f(\mathbf{PX}, \mathbf{PAP}^{\top}) = \mathbf{P}f(\mathbf{X}, \mathbf{A})$



Locality on graphs: **neighbourhoods**

- Recall: it is also highly beneficial to design *geometrically stable* (**local**) equivariant layers
- **Graphs** give us a context for locality: a node's **neighbourhood**
 - For a node i , its (1-hop) neighbourhood is commonly defined as follows:

$$\mathcal{N}_i = \{j : (i, j) \in \mathcal{E} \vee (j, i) \in \mathcal{E}\}$$

N.B. we do not explicitly consider *directed* edges, and often we assume $i \in N_i$

- Accordingly, we can extract the *multiset* of **features** in the neighbourhood

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

and define a *local* function, g , as operating over this multiset: $g(\mathbf{x}_i, \mathbf{X}_{\mathcal{N}_i})$.



A recipe for **graph** neural networks

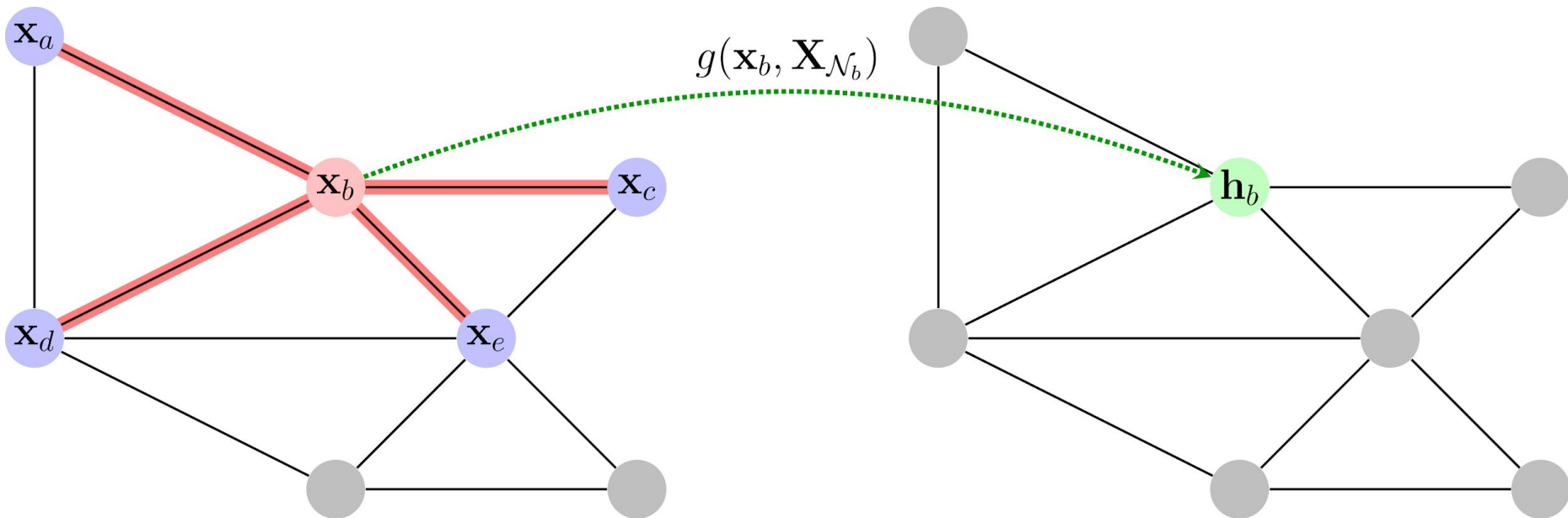
- Now we can construct permutation equivariant functions, $f(\mathbf{X}, \mathbf{A})$, by appropriately applying the local function, g , over *all* neighbourhoods:

$$f(\mathbf{X}, \mathbf{A}) = \begin{bmatrix} \text{---} & g(\mathbf{x}_1, \mathbf{X}_{\mathcal{N}_1}) & \text{---} \\ \text{---} & g(\mathbf{x}_2, \mathbf{X}_{\mathcal{N}_2}) & \text{---} \\ & \vdots & \\ \text{---} & g(\mathbf{x}_n, \mathbf{X}_{\mathcal{N}_n}) & \text{---} \end{bmatrix}$$

- To ensure equivariance, we need g to not depend on the **order** of the vertices in $\mathbf{X}_{\mathcal{N}_i}$
 - Hence, g should be permutation **invariant**!



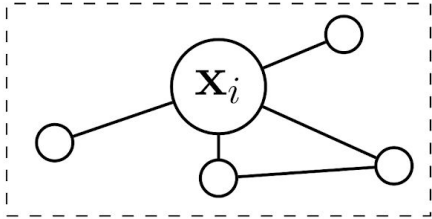
A recipe for graph neural networks, visualised



$$\mathbf{X}_{\mathcal{N}_b} = \{\{\mathbf{x}_a, \mathbf{x}_b, \mathbf{x}_c, \mathbf{x}_d, \mathbf{x}_e\}\}$$



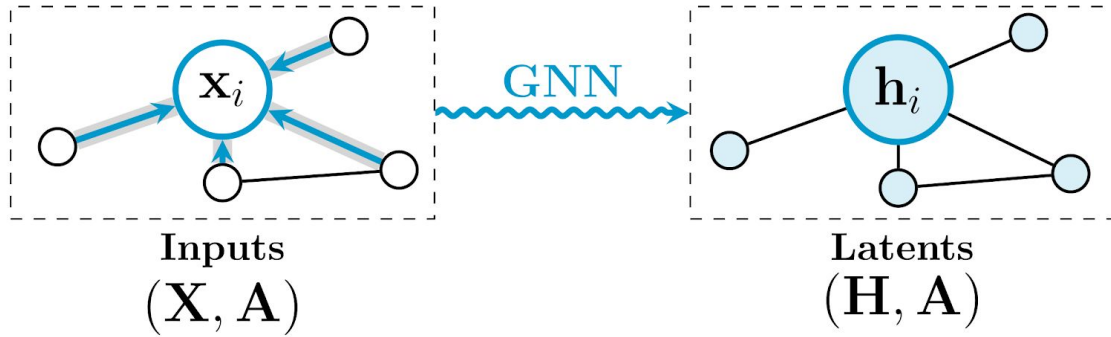
How to use GNNs?



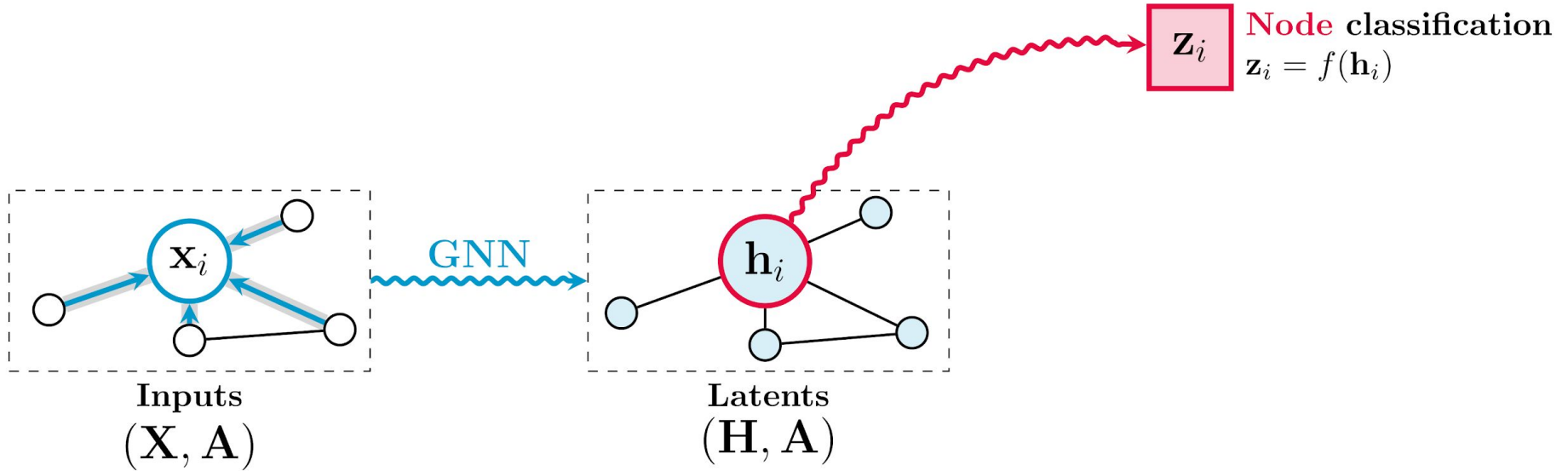
Inputs
(\mathbf{X}, \mathbf{A})



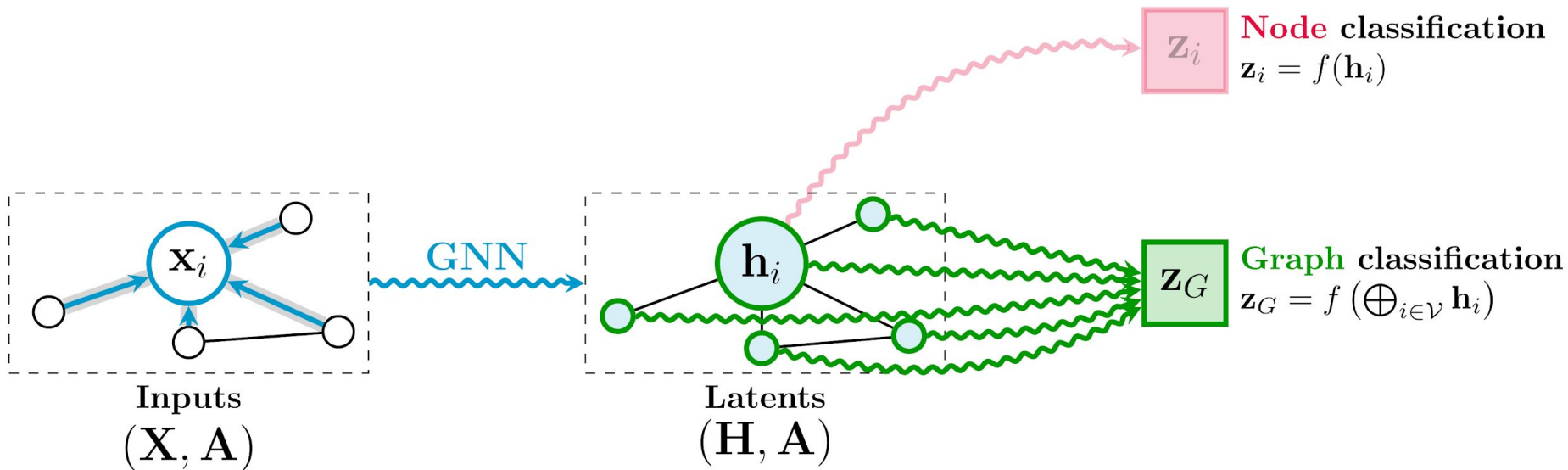
How to use GNNs?



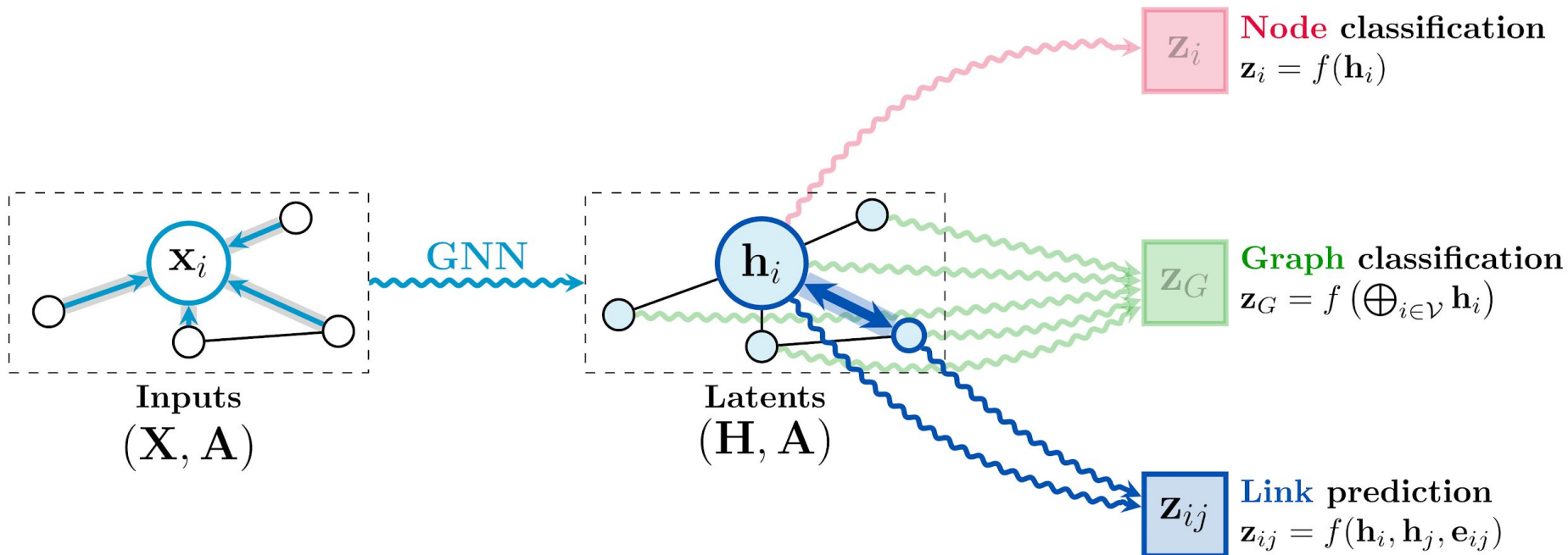
How to use GNNs?



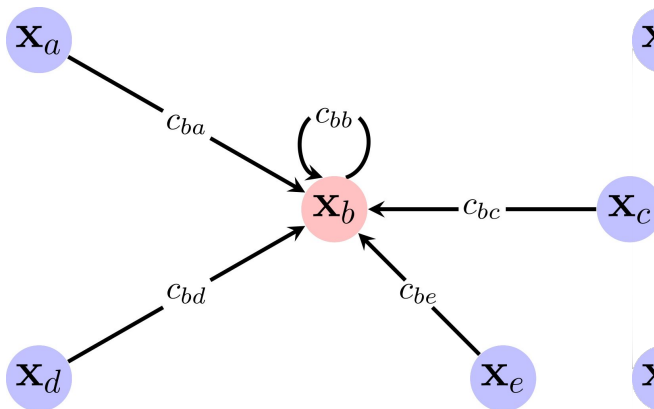
How to use GNNs?



How to use GNNs?



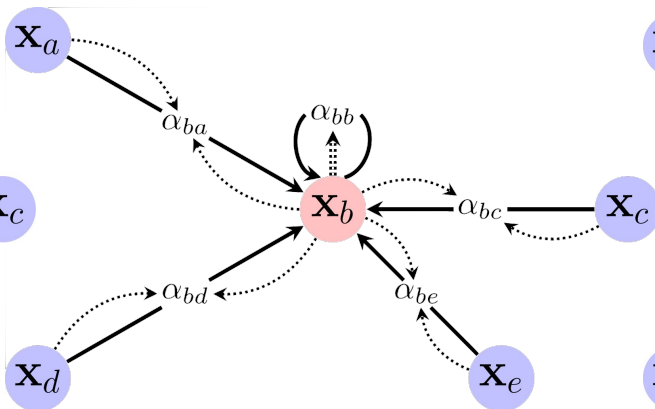
The three “flavours” of GNN layers



Convolutional

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

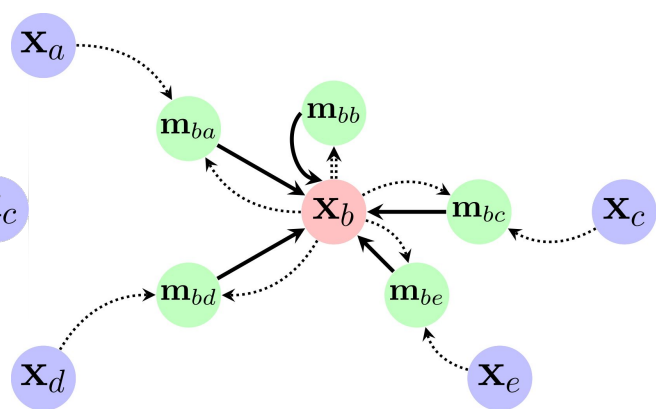
e.g. GraphSAGE, GCN, SGC



Attentional

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

e.g. MoNet, GAT, GATv2



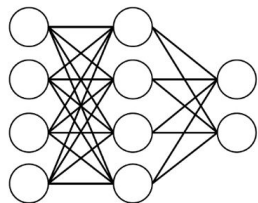
Message-passing

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

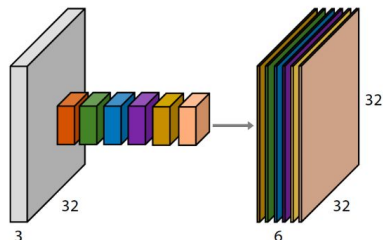
e.g. IN, MPNN, GraphNet



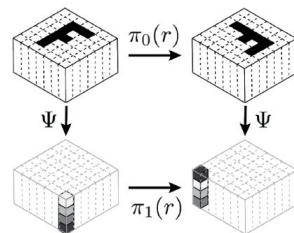
Architectures of interest



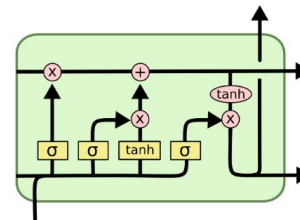
Perceptrons
Function regularity



CNNs
Translation



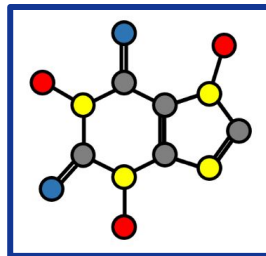
Group-CNNs
Translation+Rotation,
Global groups



LSTMs
Time warping



DeepSets / Transformers
Permutation



GNNs
Permutation



Intrinsic CNNs
Isometry / Gauge choice



6

Permutation equivariant NNs are GNNs



What to do when there is **no** graph?

- Before we go beyond permutation equivariance, a brief note on how we can use the GDL blueprint to immediately relate permutation equivariant NNs as **special cases** of GNNs
- So far, we've assumed something (seemingly) very innocent: that *the **graph** is given to us!*



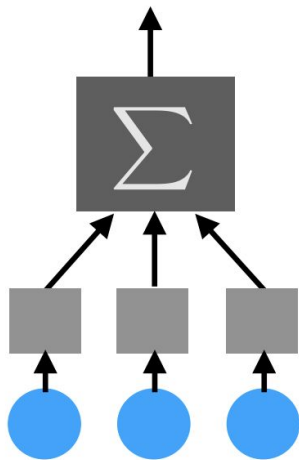
What to do when there is **no** graph?

- Before we go beyond permutation equivariance, a brief note on how we can use the GDL blueprint to immediately relate permutation equivariant NNs as **special cases** of GNNs
- So far, we've assumed something (seemingly) very innocent: that *the **graph** is given to us!*
- In practice, the given graph may often be *suboptimal* for the task we're solving
 - For various connectivity querying on graphs, maintaining the right set of edges can make a difference between **linear-time** and (amortised) **constant-time** complexity!
- Taken to the extreme: *what to do when there is **no** graph?*
 - Assume we're given a node feature matrix, but no adjacency
- Let's briefly cover two "special case" solutions...



Option 1: Assume *no edges*

Deep Sets (Zaheer *et al.*, NeurIPS'17)



No edges (set)

$$\mathcal{N}_i = \{i\}$$

$$\mathbf{h}_i = \psi(\mathbf{x}_i)$$



Option 2: Assume *all* edges

$$\mathcal{N}_i = \mathcal{V}$$

Let the GNN decide which edges matter!

Using conv-GNNs no longer makes sense.

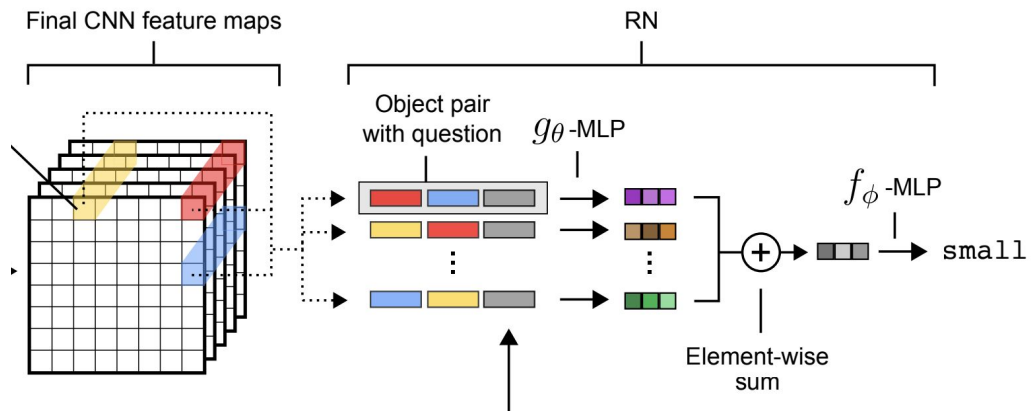
If we use **attentional** GNNs we recover:

$$\mathbf{h}_i = \phi \left(\mathbf{x}_i, \bigoplus_{j \in \mathcal{V}} a(\mathbf{x}_i, \mathbf{x}_j) \psi(\mathbf{x}_j) \right)$$

Does this look *familiar*?

Interaction Nets (Battaglia *et al.*, NeurIPS'16)

Relational Nets (Santoro *et al.*, NeurIPS'17)



All edges (*fully-connected graph*)



A note on Transformers

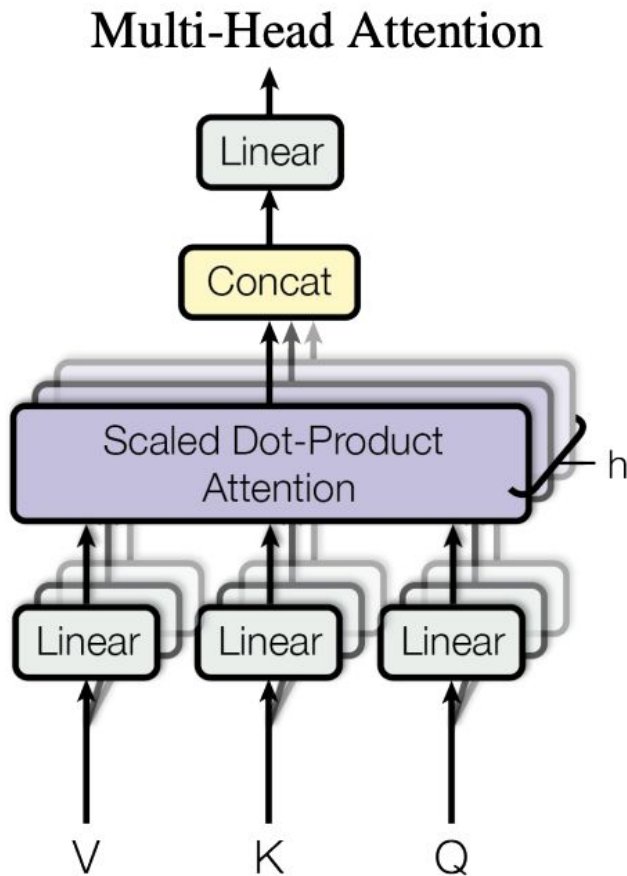
Transformers **are** Graph Neural Networks!

- Fully-connected graph
- Attentional flavour

The sequential structural information is injected through the **positional embeddings**. Dropping them yields a fully-connected GAT model.

Attention can be seen as inferring **soft adjacency**.

See Joshi (The Gradient; 2020).

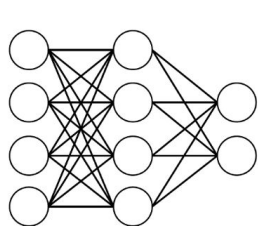


Remark: the “truth” likely lies *in between*

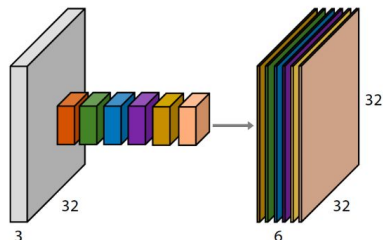
- Empty graph ignores a potential **wealth** of information
- Full graph can be hard to scale (*quadratic* complexity, large neighbourhoods)
- Ideally, we want to **infer** the adjacency matrix **A**, then use it as **edges** for a GNN!
 - Commonly termed “*latent graph inference*”.
 - Choosing edges is a **discrete decision** -- inherently hard to backpropagate!
- Out of scope for this lecture. Some interesting pointers include:
 - Neural Relational Inference (Kipf, Fetaya *et al.*, ICML'18)
 - Dynamic Graph CNN (Wang *et al.*, ACM TOG'18)
 - Differentiable Graph Module (Kazi *et al.*, MICCAI'20)
 - Pointer Graph Networks (Veličković *et al.*, NeurIPS'20)



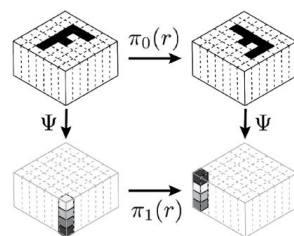
Architectures of interest



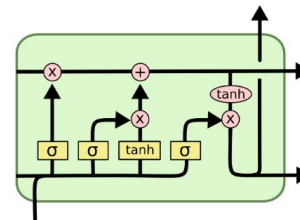
Perceptrons
Function regularity



CNNs
Translation



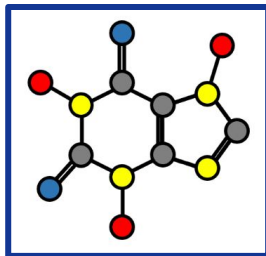
Group-CNNs
Translation+Rotation,
Global groups



LSTMs
Time warping



DeepSets / Transformers
Permutation



GNNs
Permutation



Intrinsic CNNs
Isometry / Gauge choice



7

The Fourier connection: CNNs and Spectral GNNs



Look to the Fourier transform

- The **convolution theorem** defines a very attractive identity:

$$(x \star y)(\xi) = \hat{x}(\xi) \cdot \hat{y}(\xi) \qquad \hat{x}(\xi) = \int_{-\infty}^{+\infty} x(u)e^{-i\xi u} du$$

“convolution in the time domain is multiplication in the frequency domain”

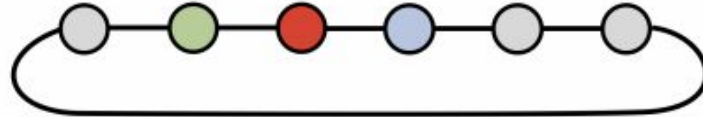
- This could give us a ‘detour’ to defining convolutions on graphs
 - Pointwise multiplication is **easy!**
 - But what are the ‘domains’ in this case?
- We will first see how graphs arise in **discrete** sequences (grids).



What's changed?

- Grids (e.g. images, text, speech) can still be seen as a **graph**
 - Pixels connected to immediate neighbours in the grid
- What would a graph neural network look like in this case?

- Simplifying assumption: **cyclical** grids



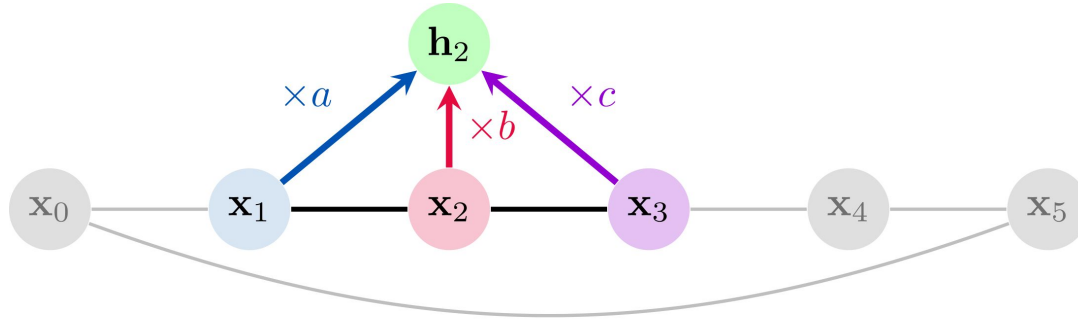
- Now every node has *identical* structure and degree
 - Allows us to more directly define a *convolution* over a neighbourhood
- Also, **coordinates matter**: we get a much stronger **translation** equivariance requirement!



Rethinking the convolution on *sequences*

*for easier handling of boundary conditions

- We can imagine a sequence as a *cyclical grid graph*, and a **convolution** over it:



- NB this defines a **circulant** matrix $\mathbf{C}([b, c, 0, 0, \dots, 0, a])$ s.t. $\mathbf{H} = f(\mathbf{X}) = \mathbf{C}\mathbf{X}$

$$f(\mathbf{X}) = \begin{bmatrix} b & c & & & a \\ a & b & c & & \\ & \ddots & \ddots & \ddots & \\ & & a & b & c \\ c & & & a & b \end{bmatrix} \begin{bmatrix} \text{---} & \mathbf{x}_0 & \text{---} \\ \text{---} & \mathbf{x}_1 & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{x}_{n-2} & \text{---} \\ \text{---} & \mathbf{x}_{n-1} & \text{---} \end{bmatrix}$$



Properties of circulants, and their eigenvectors

- Circulant matrices **commute!**
 - That is, $\mathbf{C}(\mathbf{v})\mathbf{C}(\mathbf{w}) = \mathbf{C}(\mathbf{w})\mathbf{C}(\mathbf{v})$, for *any* parameter vectors \mathbf{v}, \mathbf{w} .
- Matrices that commute are **jointly diagonalisable**.
 - That is, the eigenvectors of one are eigenvectors of *all* of them!
- Conveniently, the eigenvectors of circulants are the *discrete Fourier basis*

$$\phi_\ell = \frac{1}{\sqrt{n}} \left(1, e^{\frac{2\pi i \ell}{n}}, e^{\frac{4\pi i \ell}{n}}, \dots, e^{\frac{2\pi i (n-1)\ell}{n}} \right)^\top, \quad \ell = 0, 1, \dots, n-1$$

- This can be easily computed by studying $\mathbf{C}([0, 1, 0, 0, 0, \dots])$, which is the **shift** matrix.



The DFT and the convolution theorem

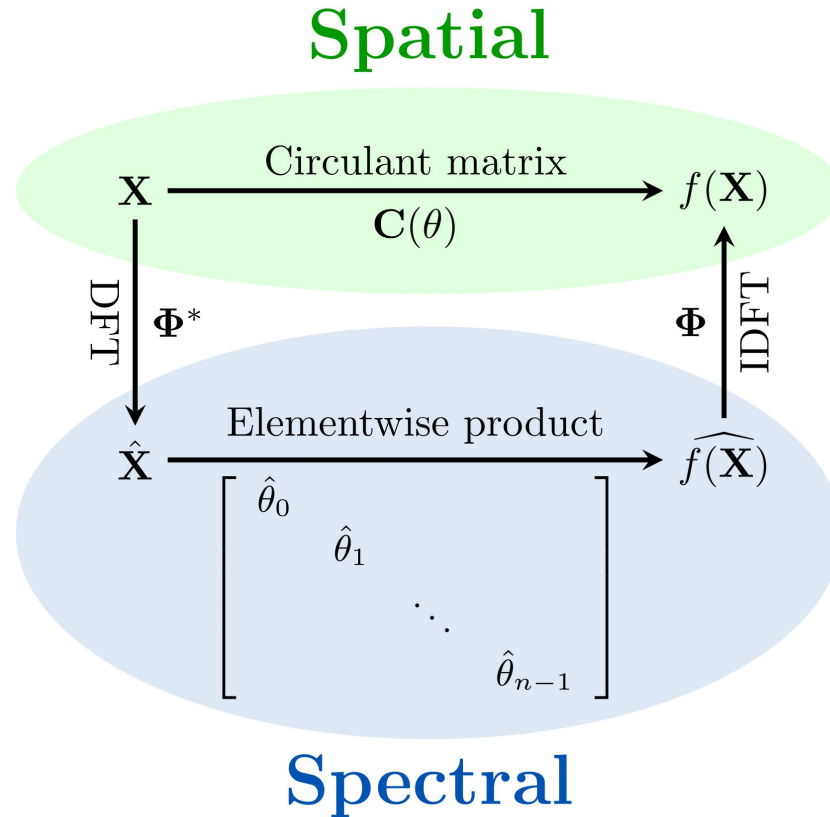
- If we stack these Fourier basis vectors into a matrix: $\Phi = (\phi_0, \dots, \phi_{n-1})$
 - We recover the *discrete Fourier transform* (DFT), as multiplication by Φ^* . (conjugate transpose)
- We can now eigendecompose any circulant as $\mathbf{C}(\theta) = \Phi \Lambda \Phi^*$
 - Where Λ is a diagonal matrix of its eigenvalues, $\hat{\theta}$
- The convolution theorem naturally follows:

$$f(\mathbf{X}) = \mathbf{C}(\theta)\mathbf{X} = \Phi \Lambda \Phi^* \mathbf{X} = \Phi \begin{bmatrix} \hat{\theta}_0 & & \\ & \ddots & \\ & & \hat{\theta}_{n-1} \end{bmatrix} \Phi^* \mathbf{X} = \Phi(\hat{\theta} \circ \hat{\mathbf{X}})$$

- Now, as long as we know Φ , we can express our convolution using $\hat{\theta}$ rather than θ



What we have covered so far



Key idea: we don't **need** to know the circulant if we know its eigenvalues!

Credits to **Michael Bronstein**

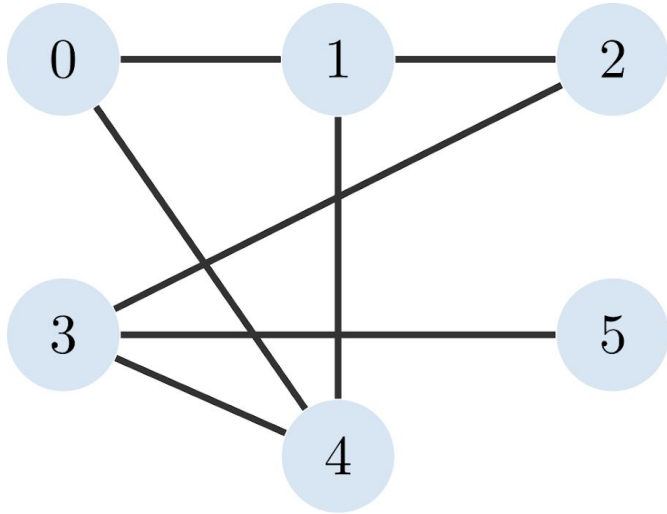


What about *graphs*?

- On graphs, convolutions of interest need to be more generic than circulants.
 - But we can still use the concept of **joint eigenbases**!
 - If we know a “graph Fourier basis”, Φ , we can only focus on learning the eigenvalues.
- For grids, we wanted our convolutions to commute with *shifts*.
 - We can think of the shift matrix as an **adjacency matrix** of the grid
 - This generalises to non-grids!
 - For the grid convolution on n nodes, Φ was always the same (n -way DFT).
 - Now **every graph** will have its own Φ !
- Want our convolution to commute with \mathbf{A} , but we cannot always eigendecompose \mathbf{A} !
- Instead, use the **graph Laplacian matrix**, $\mathbf{L} = \mathbf{D} - \mathbf{A}$, where \mathbf{D} is the degree matrix.
 - Captures all adjacency properties in mathematically convenient way!



Example Laplacian



$$\mathbf{L} = \begin{bmatrix} 2 & -1 & 0 & 0 & -1 & 0 \\ -1 & 3 & -1 & 0 & -1 & 0 \\ 0 & -1 & 2 & -1 & 0 & 0 \\ 0 & 0 & -1 & 3 & -1 & -1 \\ -1 & -1 & 0 & -1 & 3 & 0 \\ 0 & 0 & 0 & -1 & 0 & 1 \end{bmatrix}$$



Graph Fourier Transform

- Assuming undirected graphs, \mathbf{L} is:
 - **Symmetric** ($\mathbf{L}^T = \mathbf{L}$)
 - **Positive semi-definite** ($\mathbf{x}^T \mathbf{L} \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^{|\mathcal{V}|}$)
 - This means we will be able to *eigendecompose* it!
- This allows us to re-express $\mathbf{L} = \mathbf{\Phi} \mathbf{\Lambda} \mathbf{\Phi}^*$, as before.
 - Changing the eigenvalues in $\mathbf{\Lambda}$ expresses **any** operation that commutes with \mathbf{L} .
 - Commonly referred to as the **graph Fourier transform** (Bruna *et al.*, ICLR'14)
- Now, to convolve with some feature matrix \mathbf{X} , do as follows (the diagonal can be **learnable**):

$$f(\mathbf{X}) = \mathbf{\Phi} \begin{bmatrix} \hat{\theta}_0 & & \\ & \ddots & \\ & & \hat{\theta}_{n-1} \end{bmatrix} \mathbf{\Phi}^* \mathbf{X}$$



Spectral GNNs in practice

- However, directly learning the eigenvalues is typically inappropriate:
 - Not **localised**, doesn't **generalise** to other graphs, computationally **expensive**, etc.
- Instead, a common solution is to make the eigenvalues related to Λ , the eigenvalues of L
 - Commonly by a degree- k polynomial function, p_k
 - Yielding $f(\mathbf{x}) = \Phi p_k(\Lambda) \Phi^* \mathbf{x} = p_k(L) \mathbf{x}$
 - Popular choices include:
 - *Cubic splines* (Bruna *et al.*, ICLR'14)
 - *Chebyshev polynomials* (Defferrard *et al.*, NeurIPS'16)
 - *Cayley polynomials* (Levie *et al.*, Trans. Sig. Proc.'18)
- **NB** by using a polynomial in L , we have defined a **conv-GNN!**
 - With coefficients defined by $c_{ij} = (p_k(L))_{ij}$
 - Most efficient spectral approaches “*spatialise*” themselves in similar ways
 - The “*spatial-spectral*” divide is often ***not really a divide!***

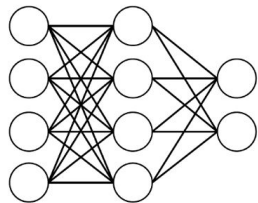


The Transformer positional encodings and beyond

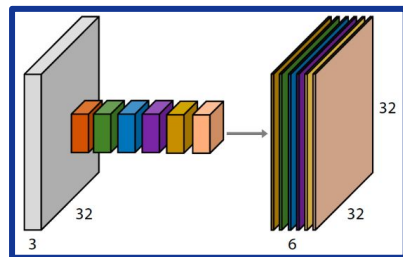
- Lastly, another look at Transformers.
- Transformers signal that the input is a **sequence** of words by using *positional embeddings*
 - Sines/cosines sampled depending on position
$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$
- **Very** similar to the DFT eigenvectors!
- Positional embeddings could hence be interpreted as eigenvectors of the grid graph
 - Which is the only assumed 'underlying' connectivity between the words
- We can use this idea to run Transformers over *general* graph structures!
 - Just feed some eigenvectors of the graph Laplacian (columns of Φ)
 - See the **Graph Transformer** from Dwivedi & Bresson (2021)
 - Rapidly emerging area of research!



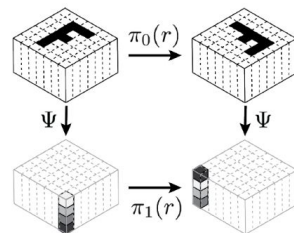
Architectures of interest



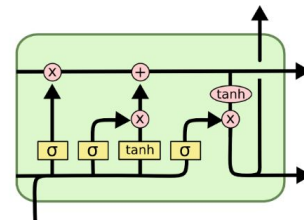
Perceptrons
Function regularity



CNNs
Translation



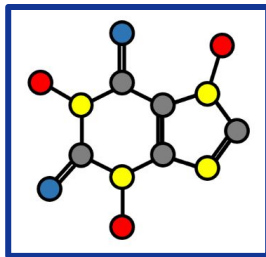
Group-CNNs
Translation+Rotation,
Global groups



LSTMs
Time warping



DeepSets / Transformers
Permutation



GNNs
Permutation



Intrinsic CNNs
Isometry / Gauge choice



8

The Group connection: Spherical CNNs



Convolutions in the Euclidean domain

*NB all our findings will be applicable to discretisations

- Now consider the *continuous* Euclidean domain $\Omega = \mathbb{R}$, where convolution takes the form:

$$(x \star \psi)(u) = \langle x, T_u \psi \rangle = \int_{\mathbb{R}} x(v) \psi(u + v) dv.$$

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

- The output of the **grid** convolution is another function on the grid
 - But this is **not** the case for **every** domain!
- Now we can see how to generalise the convolution to more general \mathbb{G}



Group convolutions

- We can define **inner products** on general domains Ω : $\langle x, \psi \rangle = \int_{\Omega} x(u)\psi(u)du$

$$(x \star \psi)(u) = \langle x, T_u \psi \rangle = \int_{\mathbb{R}} x(v)\psi(u+v)dv. \quad \text{Grids}$$

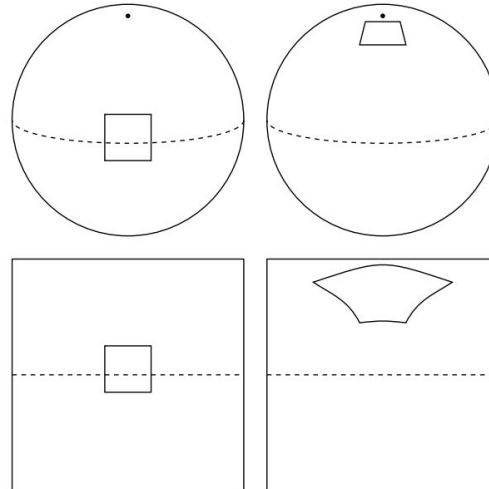
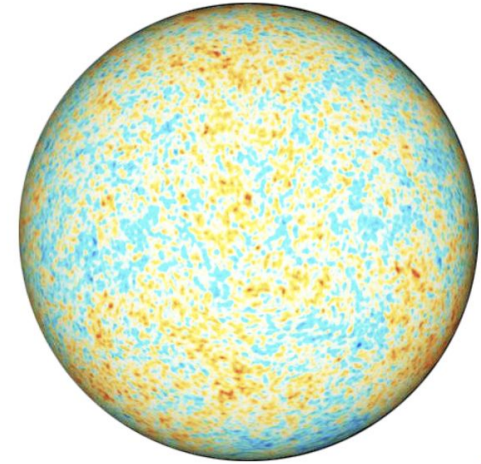
$$(x \star \psi)(\mathbf{g}) = \langle x, \rho(\mathbf{g})\psi \rangle = \int_{\Omega} x(u)\psi(\mathbf{g}^{-1}u)du \quad \text{Groups}$$

- **NB:** the group convolution is, generally, a function over elements of \mathbb{G} !!!
 - It was also a real-input function under the translation group...
 - ...because the group itself is also $\mathbb{G} = \mathbb{R}$ (scalar **shifts**!!!)



Example: Spherical convolution

- Consider signal defined on a **sphere**, $\Omega = S^2$
 - Very relevant, e.g. for earth maps, astrophysics, ...
- We want to be **rotation equivariant**
 - This means using the rotation group $\mathfrak{G} = \text{SO}(3)$
 - Image CNNs cannot support this!

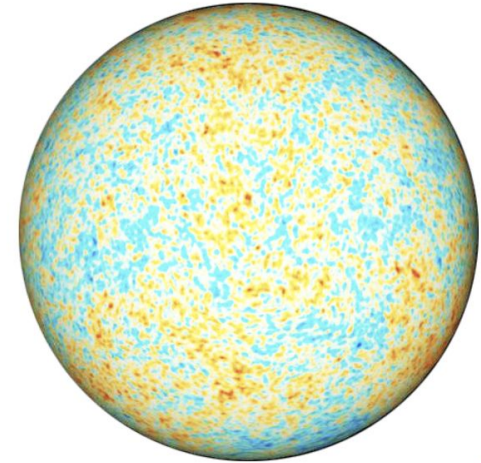


Example: Spherical convolution

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

- Consider signal defined on a **sphere**, $\Omega = S^2$
 - Very relevant, e.g. for earth maps, astrophysics, ...
- We want to be **rotation equivariant**
 - This means using the rotation group $\mathfrak{G} = \text{SO}(3)$
- We can represent the points on a sphere with 3-dim unit vectors, \mathbf{u}
 - The **group action** (rotation) transforms it with a 3 x 3 orthogonal matrix, \mathbf{R}
- We recover our spherical convolution:

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



Example: SO(3) convolution

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

- The output of a spherical convolution is a function over **all 3D rotations** in $\text{SO}(3)$
 - So, to “stack more layers”, we need to define a **convolution over $\Omega = \mathbb{G} = \text{SO}(3)$**
- Luckily, our blueprint still works: we can define a \mathbb{G} -action over elements of \mathbb{G} by function composition: $(\mathbf{g}, \mathbf{h}) \mapsto \mathbf{gh}$
- The corresponding group action can be defined as follows: $(\rho(\mathbf{g})x)(\mathbf{h}) = x(\mathbf{g}^{-1}\mathbf{h})$
- Consequently, we can build \mathbb{G} -equivariant layers over \mathbb{G} !
- Start with spherical convolution at the input layer, then **stack** $\text{SO}(3)$ -convolutions!



Example: SO(3) convolution, cont'd

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

- SO(3) group actions defined by 3x3 rotation matrices \mathbf{R}
- Hence, group action $(\rho(\mathbf{g})x)(\mathbf{h}) = x(\mathbf{g}^{-1}\mathbf{h})$ is expressible as $x(\mathbf{R}^{-1}\mathbf{Q})$
 - for some two rotation matrices \mathbf{R} and \mathbf{Q} .
- Recall the expression for the group convolution: $(x \star \psi)(\mathbf{g}) = \langle x, \rho(\mathbf{g})\psi \rangle = \int_{\Omega} x(u)\psi(\mathbf{g}^{-1}u)du$
- Putting it all together, we obtain the following two-layer convolution over spheres:

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

- This forms the essence of **Spherical CNNs** (Cohen et al., ICLR'18)



Relationship to GNNs

- The spherical CNN can still be related to GNNs
 - In practice, the sphere is *discretised* and the filter ψ is *local*
 - The non-zero elements of ψ are the neighbourhood of point \mathbf{u}
 - The integration is akin to message passing! (as many points will have zero product)

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

- Further, the \mathbb{G} -convolution described here works over **any** global symmetry group \mathbb{G}



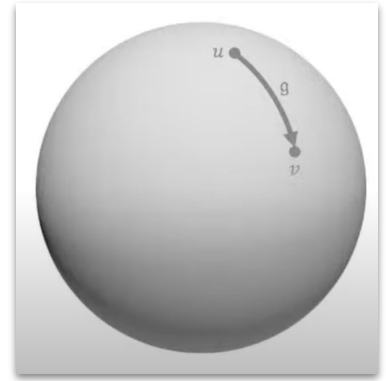
Caveat 1: Tractability

- The \mathcal{G} -convolution described here works over **any** symmetry group \mathcal{G}
- But it is only **tractable** for very small groups
 - e.g. $SO(3)$ was describable with a 3x3 orthogonal matrices
- It may be tempting to apply this idea to graphs
 - But the permutation group Σ_n has $n!$ entries to maintain
- This hints at existence of “graph convolutions” not captured by our ‘flavours’
 - (though it likely captures the most tractable ones :))
- Some of these convolutions may tradeoff expressivity for stability / complexity
 - Suggests a “way out” of the Weisfeiler–Lehman expressivity limit...

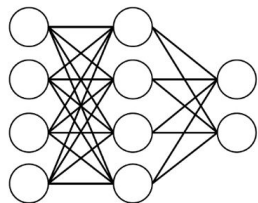


Caveat 2: Homogeneity

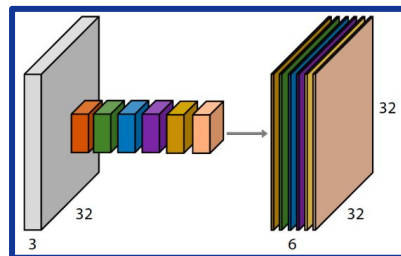
- Both the spherical and the circular grid domains are **homogeneous**
- For any two points $u, v \in \Omega$, there exists some $g \in \mathcal{G}$ s.t. $g.u = v$
- In a sense, “all points look the same”
 - This does **not** hold for graphs!
- This allowed ‘sliding’ filters across Ω to build convolutions
- What can we do in the more general case?



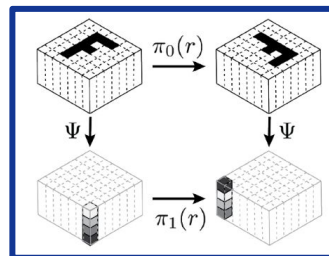
Architectures of interest



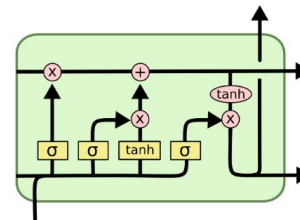
Perceptrons
Function regularity



CNNs
Translation



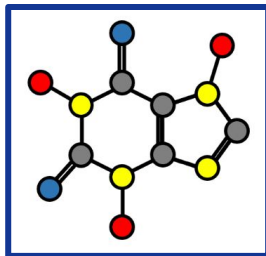
Group-CNNs
Translation+Rotation,
Global groups



LSTMs
Time warping



DeepSets / Transformers
Permutation



GNNs
Permutation



Intrinsic CNNs
Isometry / Gauge choice



DeepMind

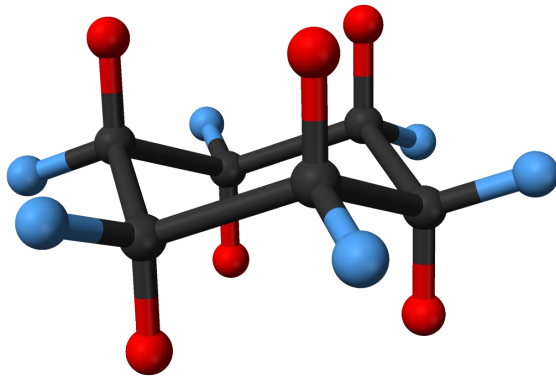
9

Geometric Graphs, Geodesics and Gauges



Geometric graphs

- Thus far, we have assumed our graphs to be a discrete, unordered, set of nodes and edges
- In many cases, this is not the entire story!
 - The graph, in fact, may often be endowed with some spatial geometry
 - It will be useful for us to **exploit** this geometry!
- Molecules are a classical case (with their three-dimensional *conformers*)



Learning over geometric graphs

- Simplified setup: nodes endowed with **features**, \mathbf{f}_u and **coordinates** $\mathbf{x}_u \in \mathbb{R}^3$
- An **equivariant message passing** layer transforms them separately
 - Yielding updated features \mathbf{f}'_u and coordinates \mathbf{x}'_u
- We can now express a group of symmetries \mathcal{G} we would like to be resistant to
 - In the case of molecules, a standard group is the Euclidean group, $E(3)$
 - Roto-translations and reflections
- For any 3D orthogonal matrix \mathbf{R} and translation vector \mathbf{b} , we define a group action $g \in E(3)$:

$$\rho(g)\mathbf{x} = \mathbf{R}\mathbf{x} + \mathbf{b}$$

and applying them to coordinates typically should not affect how features are processed!



E(n)-equivariant GNNs

- As for permutation equivariance, there exist **many** GNNs that obey E(n)-equivariance
- One elegant solution was exposed by Satorras *et al.* (ICML'21):

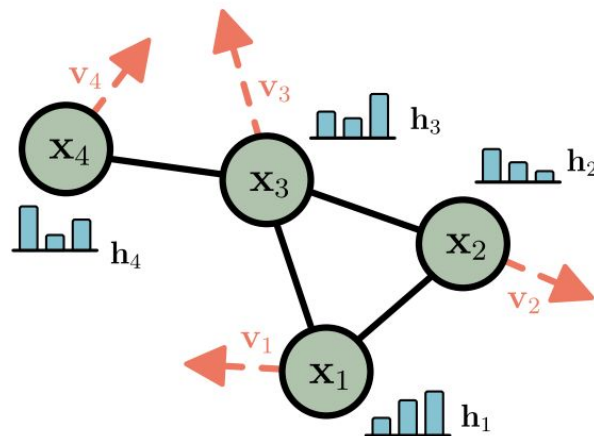
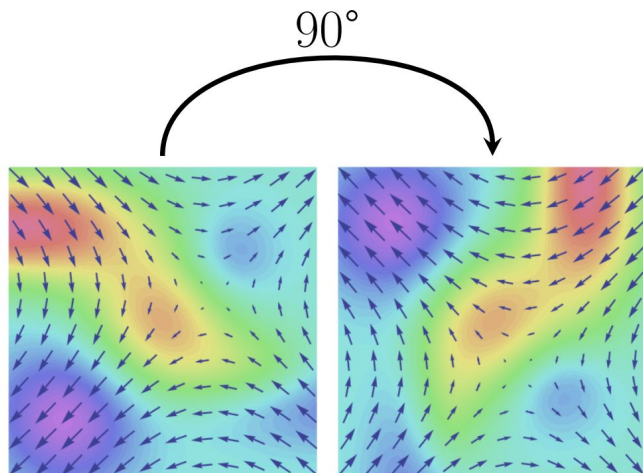
$$\mathbf{f}'_u = \phi \left(\mathbf{f}_u, \bigoplus_{v \in \mathcal{N}_u} \psi_f(\mathbf{f}_u, \mathbf{f}_v, \|\mathbf{x}_u - \mathbf{x}_v\|^2) \right),$$
$$\mathbf{x}'_u = \mathbf{x}_u + \sum_{v \neq u} (\mathbf{x}_u - \mathbf{x}_v) \psi_c(\mathbf{f}_u, \mathbf{f}_v, \|\mathbf{x}_u - \mathbf{x}_v\|^2)$$

- The actions of E(n) do not change **distance** between nodes (they are *isometric*)
 - Hence if $\mathbf{x}_u \leftarrow \mathbf{R}\mathbf{x}_u + \mathbf{b}$...
 - \mathbf{f}'_u does not change, and $\mathbf{x}'_u \leftarrow \mathbf{R}\mathbf{x}'_u + \mathbf{b}$, as expected!
- Hence, this layer is E(n)-equivariant over scalar features \mathbf{f}_u



Vector-structured features

- However, what if some of the node features (f_u) depend on the geometry?
 - e.g. they could be vector **forces**
 - Rotating the molecule should rotate these vectors too!
- The model on the previous slide does **not** take this into account!



Vector-structured features

- However, what if some of the node features (\mathbf{f}_u) depend on the geometry?
 - e.g. they could be vector **forces**
 - Rotating the molecule should rotate these vectors too!
- The model on the previous slide does **not** take this into account!
- Satorras *et al.* do propose a variant of their model that works with vectors:

$$\mathbf{v}_i^{l+1} = \phi_v(\mathbf{h}_i^l) \mathbf{v}_i^l + C \sum_{j \neq i} (\mathbf{x}_i^l - \mathbf{x}_j^l) \phi_x(\mathbf{m}_{ij})$$

$$\mathbf{x}_i^{l+1} = \mathbf{x}_i^l + \mathbf{v}_i^{l+1}$$

- But the issue will keep re-appearing as we “tensor up” our features!
 - Is there a way to talk about the general “set of solutions” for E(n)-equivariance?



Towards a general solution: Tensor Field Networks

- We actually can factorise the group action of SE(3) on inputs of *any* tensor order !!

$$\rho(g) = \mathbf{Q}^\top \left[\bigoplus_{\ell} \mathbf{D}_{\ell}(g) \right] \mathbf{Q}$$

where \mathbf{D}_{ℓ} are the Wigner D-matrices, and \mathbf{Q} is a \mathfrak{g} -specific change-of-basis matrix.

- This can be exploited to write a generic SE(3)-equivariant model (the TFN)

$$\mathbf{f}_{\text{out},i}^{\ell} = \sum_{k \geq 0} \underbrace{\int \mathbf{W}^{\ell k}(\mathbf{x}' - \mathbf{x}_i) \mathbf{f}_{\text{in}}^k(\mathbf{x}') d\mathbf{x}'}_{k \rightarrow \ell \text{ convolution}} = \sum_{k \geq 0} \sum_{j=1}^n \underbrace{\mathbf{W}^{\ell k}(\mathbf{x}_j - \mathbf{x}_i) \mathbf{f}_{\text{in},j}^k}_{\text{node } j \rightarrow \text{node } i \text{ message}}$$

$$\mathbf{W}^{\ell k}(\mathbf{x}) = \sum_{J=|k-\ell|}^{k+\ell} \varphi_J^{\ell k}(\|\mathbf{x}\|) \mathbf{W}_J^{\ell k}(\mathbf{x}), \quad \text{where } \mathbf{W}_J^{\ell k}(\mathbf{x}) = \sum_{m=-J}^J Y_{Jm}(\mathbf{x}/\|\mathbf{x}\|) \mathbf{Q}_{Jm}^{\ell k}$$

where the matrices \mathbf{W}_J are precomputed based on the above parametrisation



SE(3)-Transformers (Fuchs, Worrall, *et al.*, NeurIPS'20)

Restrict TFNs to only act over neighbourhood N_i . Also adds Transformer attention.

$$\mathbf{f}_{\text{out},i}^{\ell} = \underbrace{\mathbf{W}_V^{\ell\ell} \mathbf{f}_{\text{in},i}^{\ell}}_{\textcircled{3} \text{ self-interaction}} + \sum_{k \geq 0} \sum_{j \in \mathcal{N}_i \setminus i} \underbrace{\alpha_{ij}}_{\textcircled{1} \text{ attention}} \underbrace{\mathbf{W}_V^{\ell k} (\mathbf{x}_j - \mathbf{x}_i) \mathbf{f}_{\text{in},j}^k}_{\textcircled{2} \text{ value message}}.$$

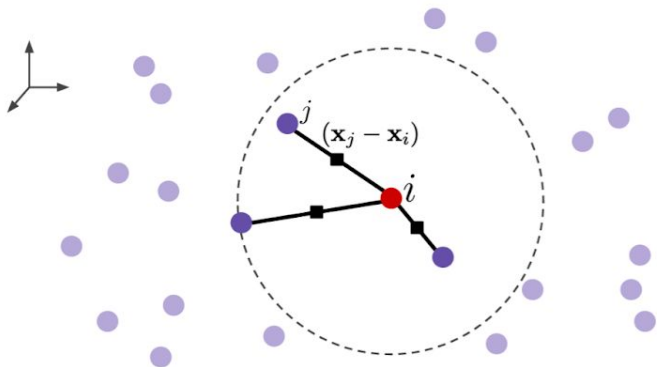
$$\alpha_{ij} = \frac{\exp(\mathbf{q}_i^{\top} \mathbf{k}_{ij})}{\sum_{j' \in \mathcal{N}_i \setminus i} \exp(\mathbf{q}_i^{\top} \mathbf{k}_{ij'})}, \quad \mathbf{q}_i = \bigoplus_{\ell \geq 0} \sum_{k \geq 0} \mathbf{W}_Q^{\ell k} \mathbf{f}_{\text{in},i}^k, \quad \mathbf{k}_{ij} = \bigoplus_{\ell \geq 0} \sum_{k \geq 0} \mathbf{W}_K^{\ell k} (\mathbf{x}_j - \mathbf{x}_i) \mathbf{f}_{\text{in},j}^k.$$

Gives all possible equivariant (attentional) GNNs, but we still need to pre-compute \mathbf{W}_j

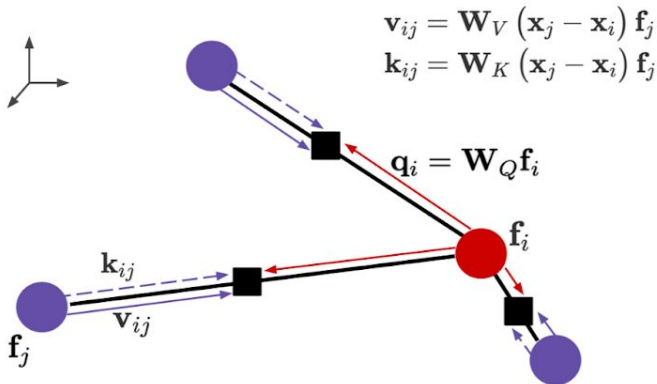


SE(3)-Transformers (Fuchs, Worrall, *et al.*, NeurIPS'20)

Step 1: Get nearest neighbours and relative positions



Step 3: Propagate queries, keys, and values to edges



Step 2: Get SO(3)-equivariant weight matrices



Clebsch-Gordan Coeff.

$$\mathbf{Q}_{Jm}^{\ell k}$$



Radial Neural Network

$$\varphi_J^{\ell k}(\|\mathbf{x}\|)$$



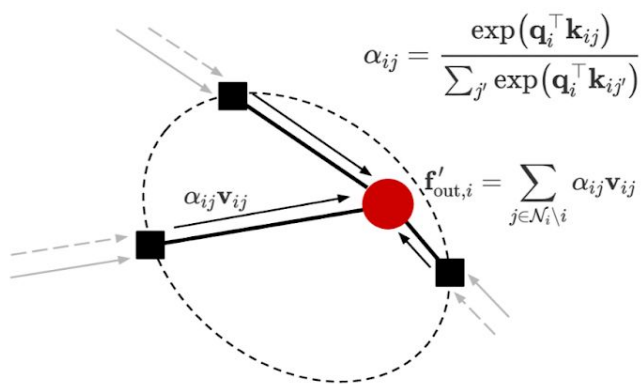
Spherical Harmonics

$$Y_{Jm}\left(\frac{\mathbf{x}}{\|\mathbf{x}\|}\right)$$

Matrix \mathbf{W} consists of blocks mapping between degrees

$$\mathbf{W}(\mathbf{x}) = \mathbf{W}\left(\left\{\mathbf{Q}_{Jm}^{\ell k}, \varphi_J^{\ell k}(\|\mathbf{x}\|), Y_{Jm}\left(\frac{\mathbf{x}}{\|\mathbf{x}\|}\right)\right\}_{J,m,\ell,k}\right)$$

Step 4: Compute attention and aggregate



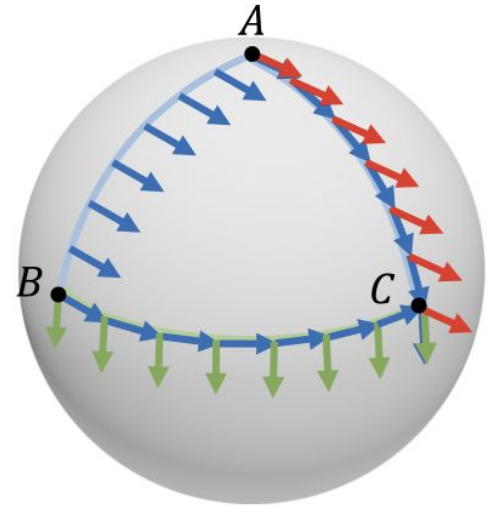
Manifolds

- More generally, we can define **manifold** domains
 - Highly relevant e.g. for computer graphics, protein design, fMRI processing
- In fact, manifold-oriented models end up not at all dissimilar to equivariant message passing!
- Detailed exposition and extended theory are deferred to the draft book



Manifolds

- More generally, we can define **manifold** domains
 - Highly relevant e.g. for computer graphics, protein design, fMRI processing
- In fact, manifold-oriented models end up not at all dissimilar to equivariant message passing!
- Key concept: **parallel transport**, which allows us to transport tangent vectors along a curve to preserve **direction**
 - NB. It is **path-dependent** (e.g. $BC \neq ABC$)



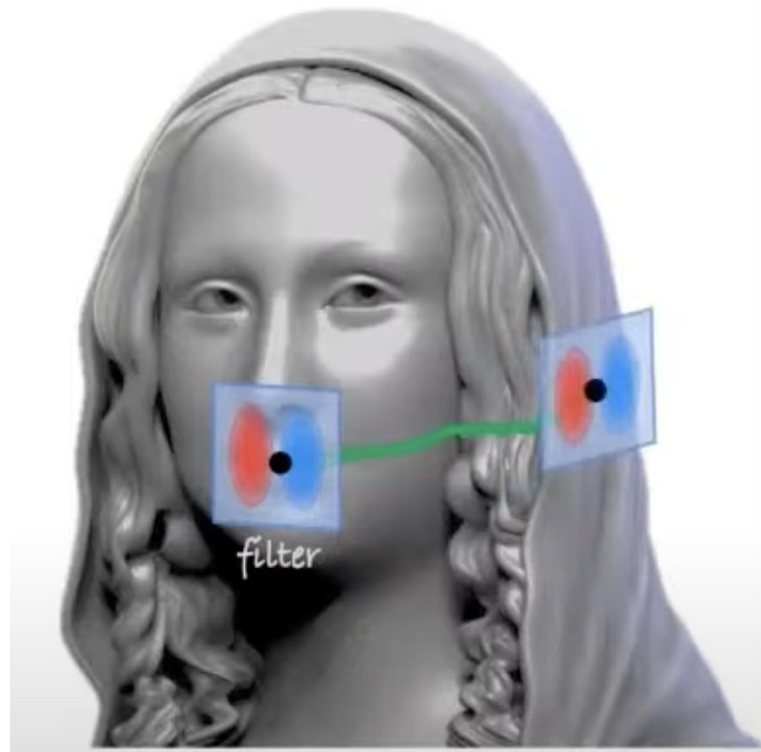
Manifolds

- More generally, we can define **manifold** domains
 - Highly relevant e.g. for computer graphics, protein design, fMRI processing
- In fact, manifold-oriented models end up not at all dissimilar to equivariant message passing!
- Key concept: ***parallel transport***, which allows us to transport tangent vectors along a curve to preserve **direction**
 - NB. It is **path-dependent**
- Now we can define (with some caveats) “sliding” a filter along more general surfaces



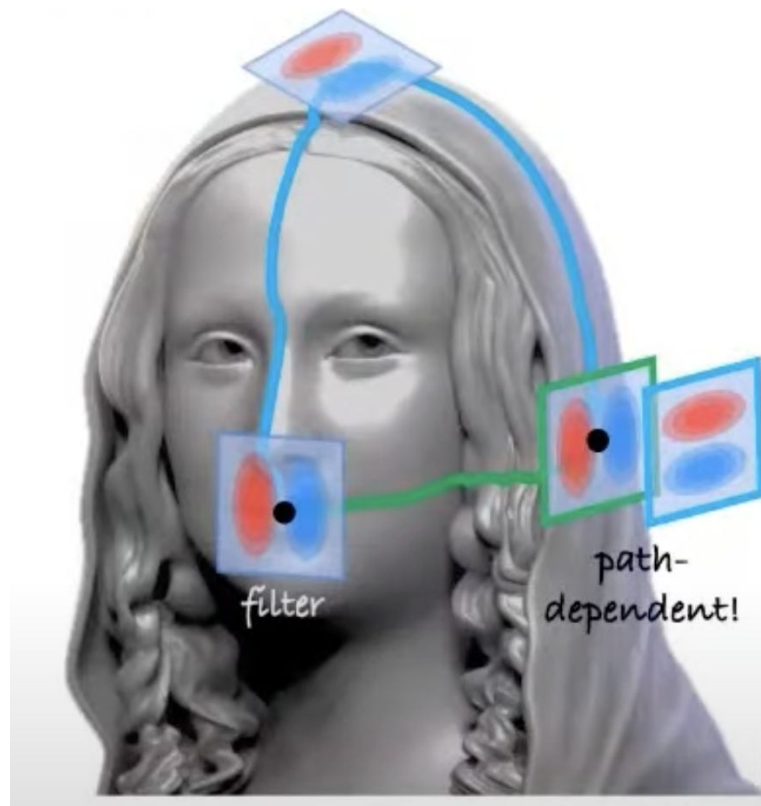
Manifolds

- More generally, we can define **manifold** domains
 - Highly relevant e.g. for computer graphics, protein design, fMRI processing
- In fact, manifold-oriented models end up not at all dissimilar to equivariant message passing!
- Key concept: **parallel transport**, which allows us to transport tangent vectors along a curve to preserve **direction**
 - NB. It is **path-dependent**
- Now we can define (with some caveats) “sliding” a filter along more general surfaces



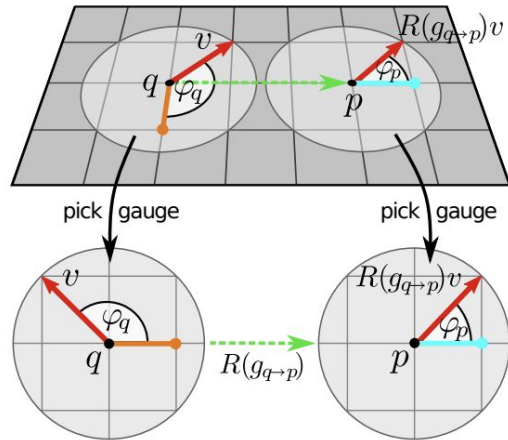
Manifolds

- More generally, we can define **manifold** domains
 - Highly relevant e.g. for computer graphics, protein design, fMRI processing
- In fact, manifold-oriented models end up not at all dissimilar to equivariant message passing!
- Key concept: **parallel transport**, which allows us to transport tangent vectors along a curve to preserve **direction**
 - NB. It is **path-dependent**
- Now we can define (with some caveats) “sliding” a filter along more general surfaces

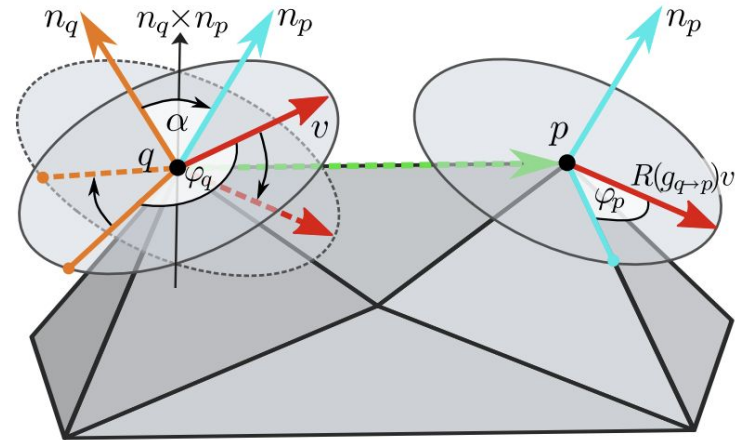


Some architectures recovered from blueprint

- Geodesic CNN (Masci *et al.*, CVPR'15)
- Gauge-equivariant Mesh CNN (de Haan *et al.*, ICLR'21)



(a) Parallel transport on a flat mesh.



(b) Parallel transport along an edge of a general mesh.



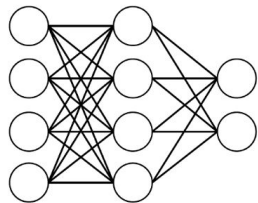
Some architectures recovered from blueprint

- Geodesic CNN (Masci *et al.*, CVPR'15)
- Gauge-equivariant Mesh CNN (de Haan *et al.*, ICLR'21)
 - GEM-CNNs feel a lot like **message passing!** (with *precomputed* transport matrices)

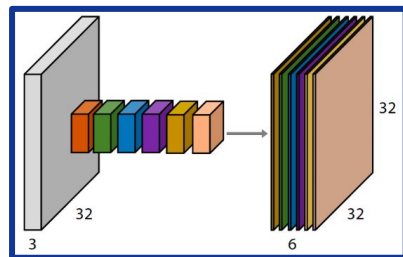
$$\mathbf{h}_u = \Theta_{\text{self}} \mathbf{x}_u + \sum_{v \in \mathcal{N}_u} \Theta_{\text{neigh}}(\vartheta_{uv}) \rho(\mathfrak{g}_{v \rightarrow u}) \mathbf{x}_v$$



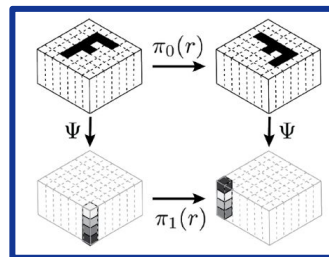
Architectures of interest



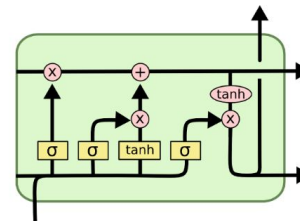
Perceptrons
Function regularity



CNNs
Translation



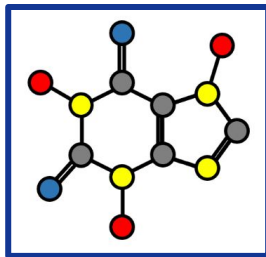
Group-CNNs
Translation+Rotation,
Global groups



LSTMs
Time warping



DeepSets / Transformers
Permutation



GNNs
Permutation



Intrinsic CNNs
Isometry / Gauge choice



DeepMind

A

**Further
resources**



GEOMETRIC DEEP LEARNING

Grids, Groups, Graphs, Geodesics, and Gauges

Michael M. Bronstein, Joan Bruna, Taco Cohen, Petar Veličković

[Read the paper](#)

[Read the blog post](#)

[Watch the ICLR'21 keynote](#)

[Watch the Erlangen keynote](#)

[Follow the Lectures \(AMMI 2021\)](#)

[Contact the authors](#)



156-page proto book on arXiv

Geometric Deep Learning Grids, Groups, Graphs, Geodesics, and Gauges

Michael M. Bronstein¹, Joan Bruna², Taco Cohen³, Petar Veličković⁴

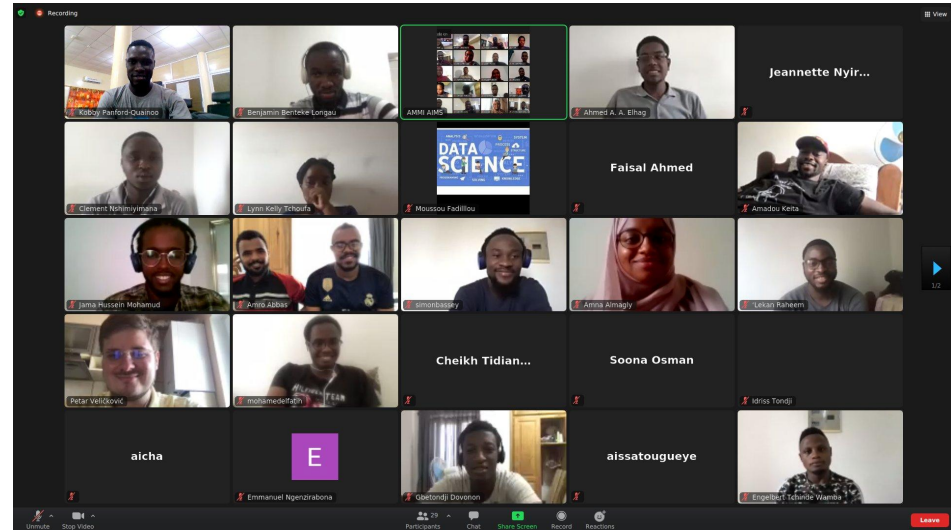
May 4, 2021



African Master's in Machine Intelligence course (12h)

All lecture recordings

Lecture 1: <i>Introduction</i>	Michael M. Bronstein	Recording	Slides
Lecture 2: <i>High-Dimensional Learning</i>	Joan Bruna	Recording	Slides
Lecture 3: <i>Geometric Priors I</i>	Taco Cohen	Recording	Slides
Lecture 4: <i>Geometric Priors II</i>	Joan Bruna	Recording	Slides
Lecture 5: <i>Graphs & Sets I</i>	Petar Veličković	Recording	Slides
Lecture 6: <i>Graphs & Sets II</i>	Petar Veličković	Recording	Slides
Lecture 7: <i>Grids</i>	Joan Bruna	Recording	Slides
Lecture 8: <i>Groups</i>	Taco Cohen	Recording	Slides
Lecture 9: <i>Geodesics & Manifolds</i>	Michael M. Bronstein	Recording	Slides
Lecture 10: <i>Gauges</i>	Taco Cohen	Recording	Slides
Lecture 11: <i>Sequences & Time Warping</i>	Petar Veličković	Recording	Slides
Lecture 12: <i>Conclusions</i>	Michael M. Bronstein	Recording	Slides



DeepMind

Thank you!

Questions / Feedback?

`petarv@deepmind.com | https://petar-v.com`

With many thanks to Joan Bruna, Michael Bronstein and Taco Cohen

