

# Graph Neural Networks

**Mihai Alexe**

ECMWF Bonn

[first.lastname@ecmwf.int](mailto:first.lastname@ecmwf.int)

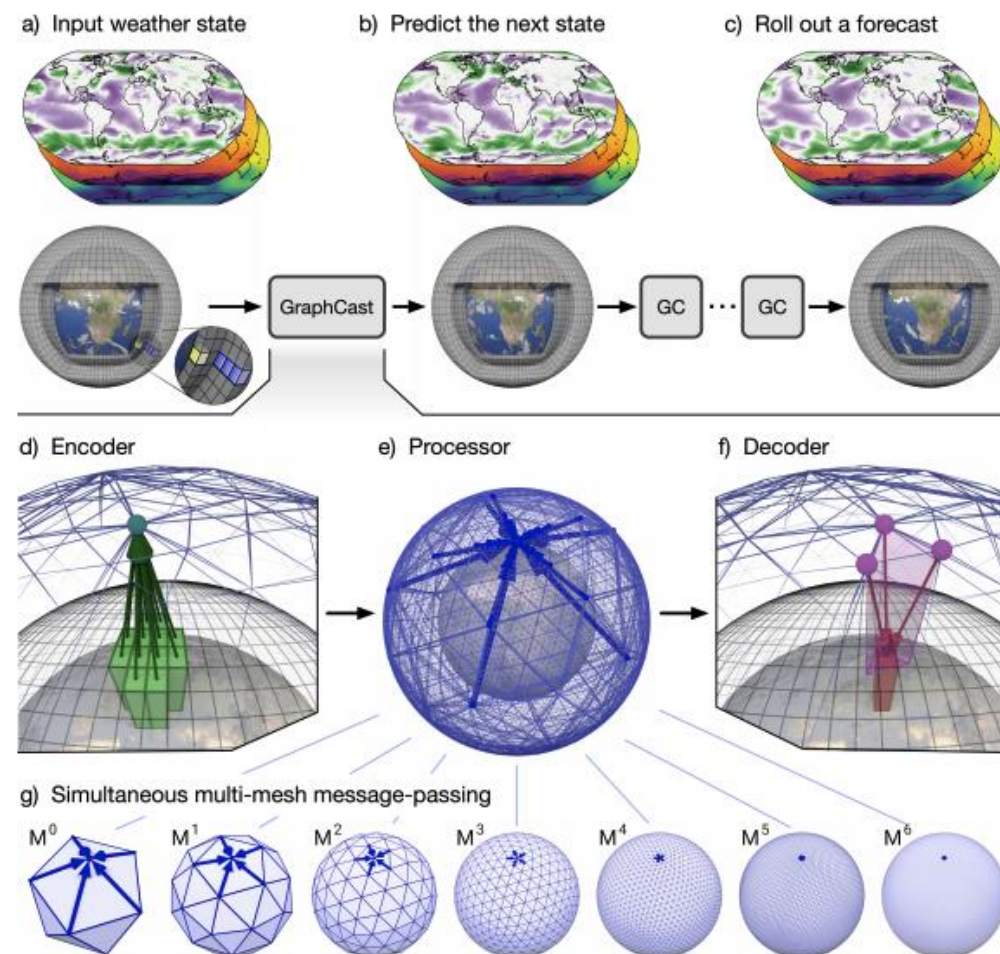


GNNs are **neural networks** built to operate on graph data.

[Submitted on 24 Dec 2022 (v1), last revised 4 Aug 2023 (this version, v2)]

## GraphCast: Learning skillful medium-range global weather forecasting

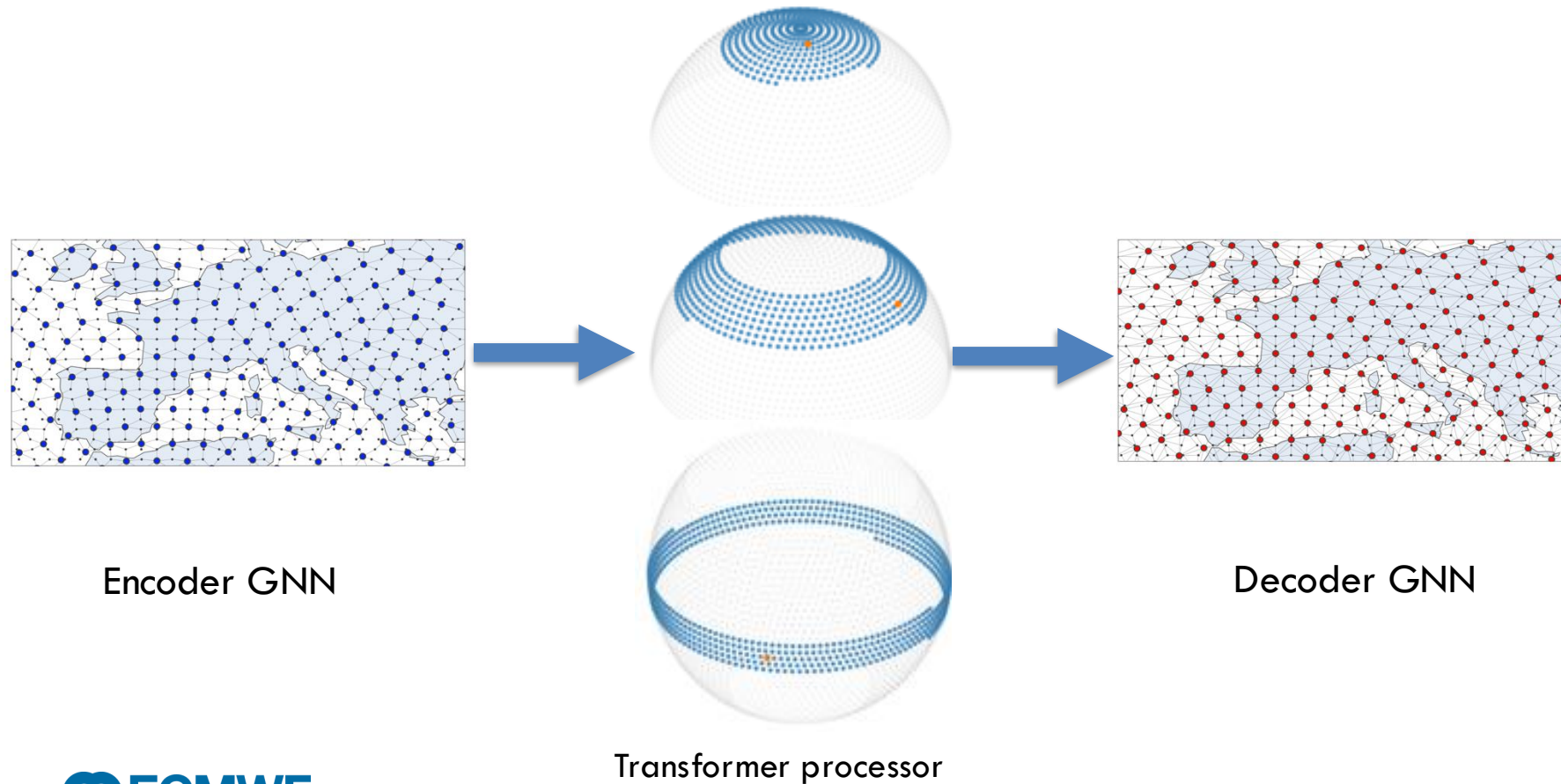
Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirnsberger, Meire Fortunato, Ferran Alet, Suman Ravuri, Timo Ewalds, Zach Eaton-Rosen, Weihua Hu, Alexander Merose, Stephan Hoyer, George Holland, Oriol Vinyals, Jacklynn Stott, Alexander Pritzel, Shakir Mohamed, Peter Battaglia



[Submitted on 3 Jun 2024 (v1), last revised 7 Aug 2024 (this version, v2)]

## AIFS -- ECMWF's data-driven forecasting system

Simon Lang, Mihai Alexe, Matthew Chantry, Jesper Dramsch, Florian Pinault, Baudouin Raoult, Mariana C. A. Clare, Christian Lessig, Michael Maier-Gerber, Linus Magnusson, Zied Ben Bouallègue, Ana Prieto Nemesio, Peter D. Dueben, Andrew Brown, Florian Pappenberger, Florence Rabier



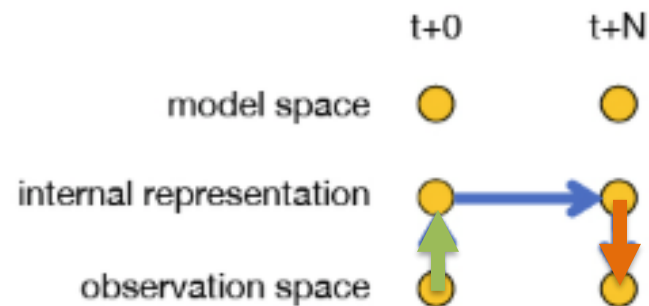
[Submitted on 20 Dec 2024]

## GraphDOP: Towards skilful data-driven medium-range weather forecasts learnt and initialised directly from observations

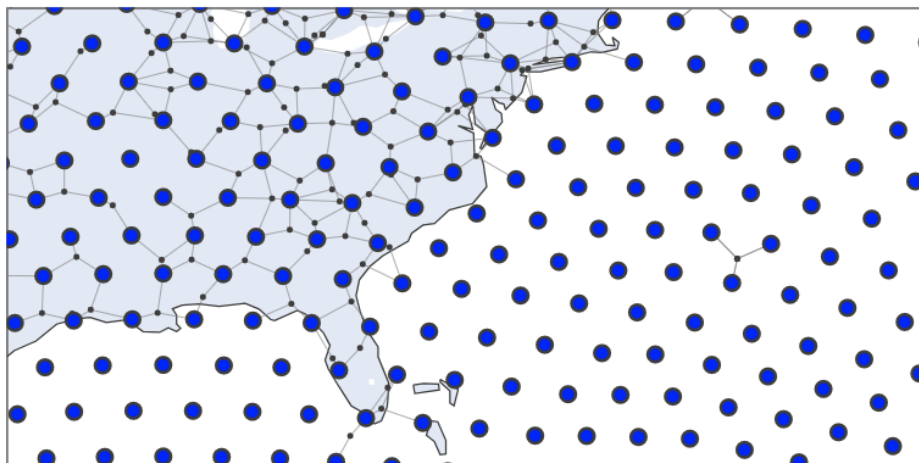
Mihai Alexe, Eulalie Boucher, Peter Lean, Ewan Pinnington, Patrick Laloyaux, Anthony McNally, Simon Lang, Matthew Chantry, Chris Burrows, Marcin Chrust, Florian Pinault, Ethel Villeneuve, Niels Bormann, Sean Healy

### 5 Predict future observations from observations

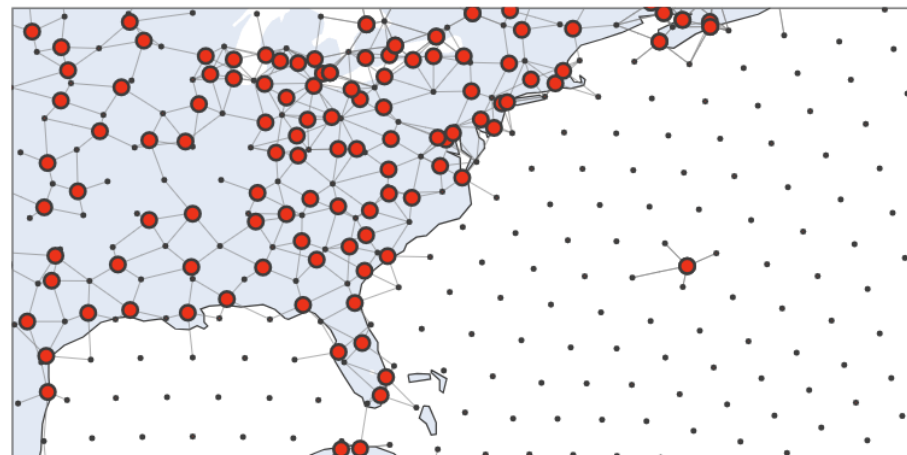
make predictions in observation space,  
use observations as truth



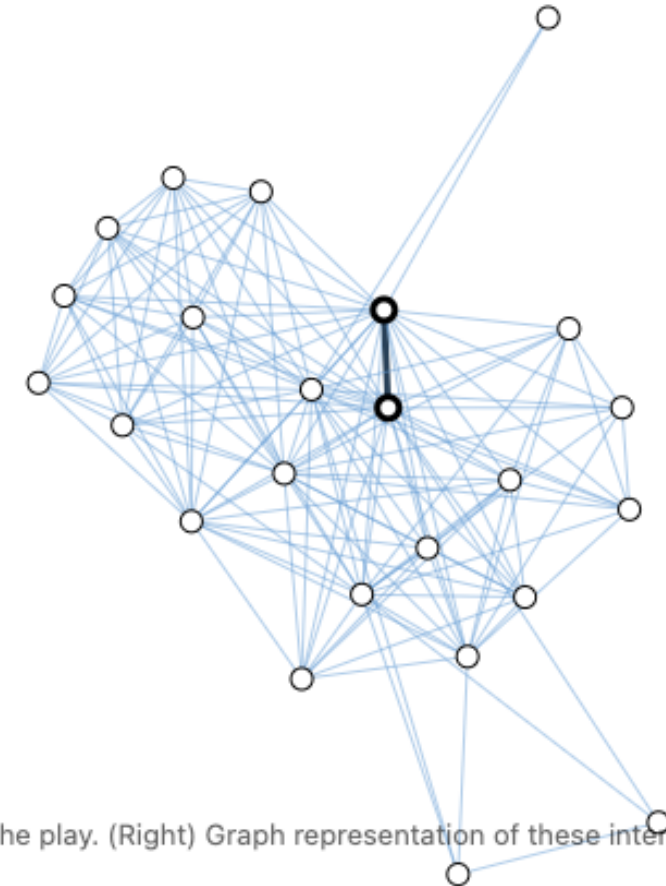
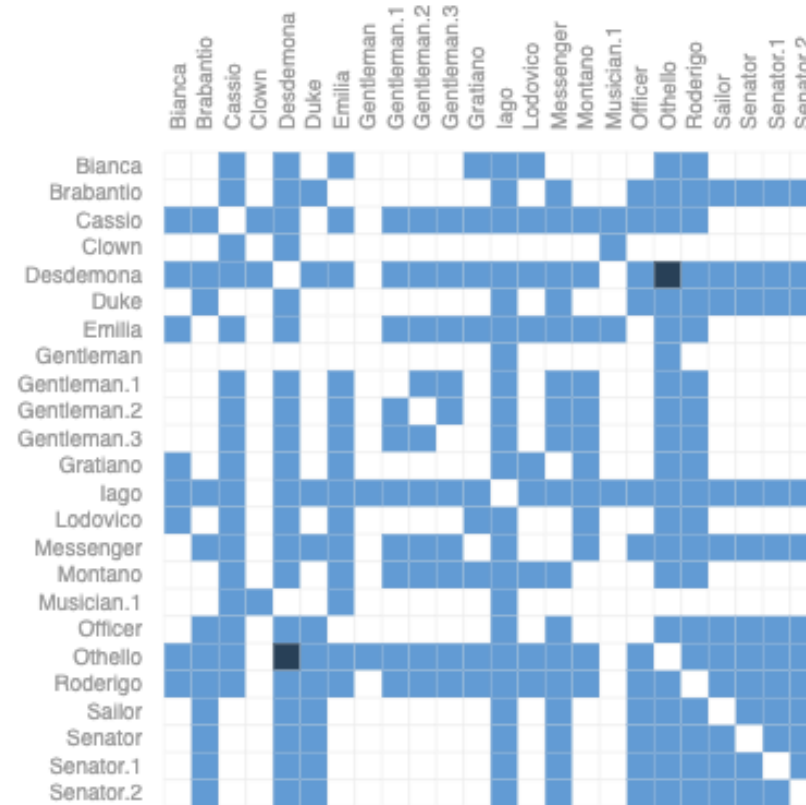
OBS to H graph



H to OBS graph



# Graphs: vertices (nodes), edges (links), connectivity (adjacency) ...



(Left) Image of a scene from the play "Othello". (Center) Adjacency matrix of the interaction between characters in the play. (Right) Graph representation of these interactions.

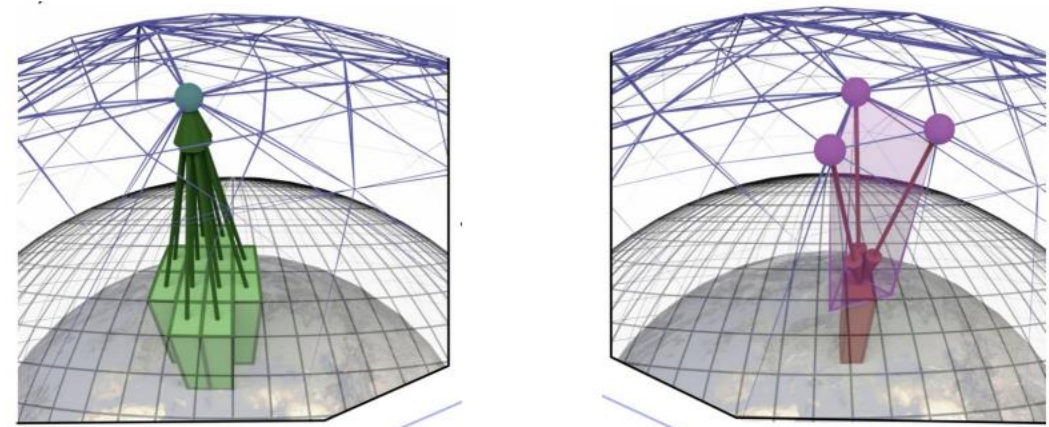
Adjacency matrix **A**

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



# Graph structure representation

**Bipartite graphs:** encode a relation SRC  $\rightarrow$  DST



**ERA5  $\rightarrow$  Latent**

**Latent  $\rightarrow$  ERA5**

Your choice (of representation) matters!

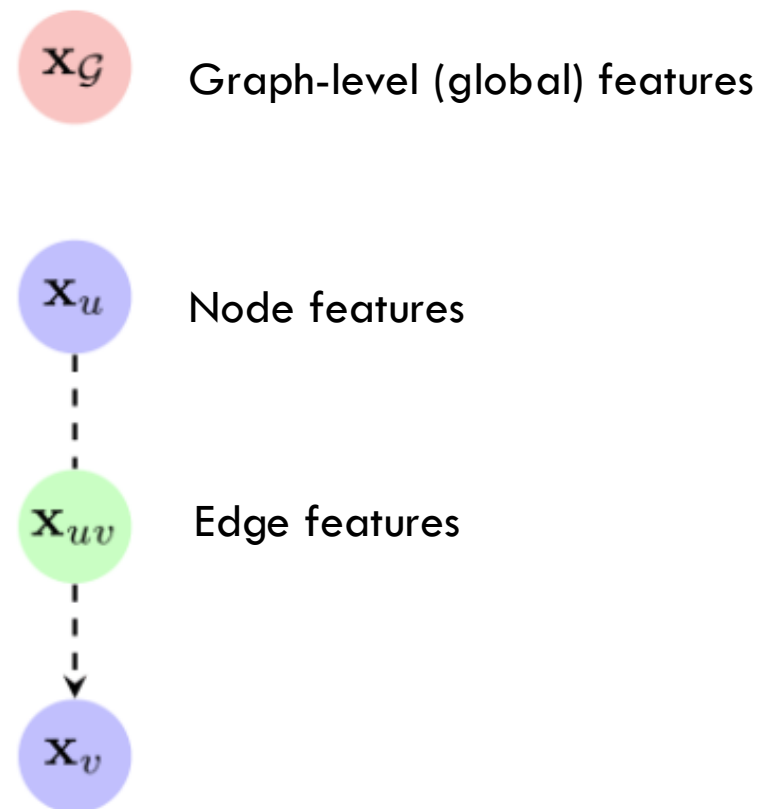
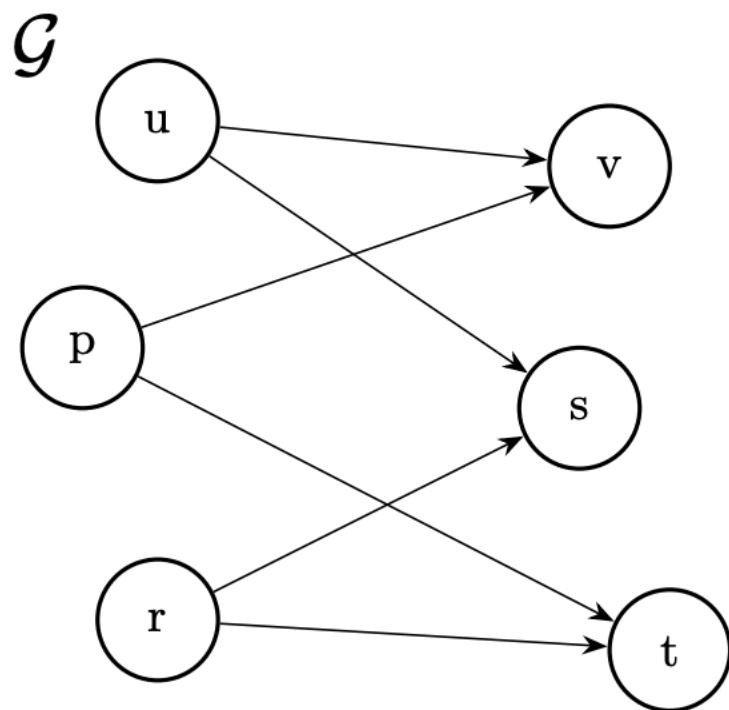
```
src, dst = edge_index
x_src = x[src] # [E, D]
out = torch.empty((N_DST, D), ...)
for e_idx, d in enumerate(dst):
    out[d] += x_src[e_idx]
```

**Edge list**

```
out = torch.empty((N_DST, D), ...)
for d in range(N_DST):
    off_start, off_end = offsets[d], offsets[d+1]
    acc = 0.0
    for e_idx in range(off_start, off_end):
        src = indices[e_idx]
        acc += out[src, :]
    out[d, :] = acc
```

**Compressed format (CSC)**

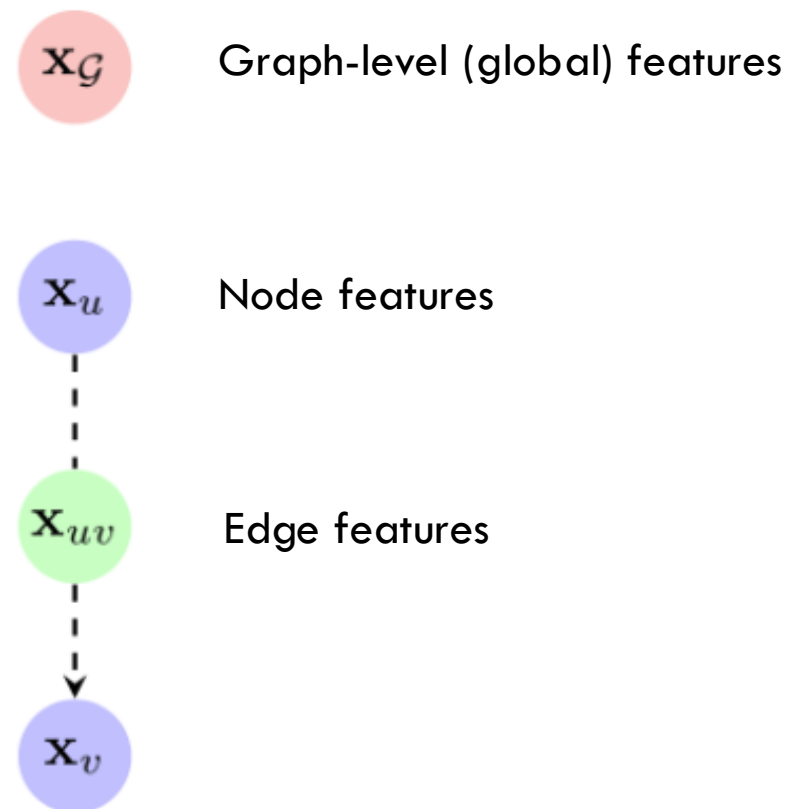
## Graph features (= information associated with elements of our graph)





# Graph neural networks

GNNs are **neural networks**  
built to operate on **graph data**.



## Quick detour: MLPs

### Multi-layer perceptrons

```
from torch import nn

def generate_mlp_module(num_inputs: int = 32, hidden_dim: int = 64, num_outputs: int = 32):
    mlp = nn.Sequential(
        nn.Linear(num_inputs, hidden_dim),
        nn.LeakyReLU(0.1),
        nn.Linear(hidden_dim, hidden_dim),
        nn.LeakyReLU(0.1),
        nn.Linear(hidden_dim, num_outputs),
        nn.LeakyReLU(0.1),
        nn.LayerNorm(num_outputs)
    )
    return mlp
```

( ..., ??? ) = **MLP** ( ..., ??? )

MLPs will be denoted by Greek letters

$\phi$ ,  $\psi$  and  $\rho$

Before we mathematically define a GNN layer...

**Q:** What **inductive biases** should a GNN have?

## Locality

We want the GNN signal to be stable under small domain deformations.

Standard deep NNs (e.g., CNNs) build large-scale ops from small-scale building blocks (e.g., 3x3 convolutions).

GNN layers operate locally, too - over neighborhoods.

We can extract neighborhood features and define local functions (MLPs) operating on them:  
 $\phi$

$$\mathbf{X}_{\mathcal{N}_i} = \{ \{ \mathbf{x}_j : j \in \mathcal{N}_i \} \}$$
$$\phi(\mathbf{x}_i, \mathbf{X}_{\mathcal{N}_i}).$$

# Permutation invariance and equivariance

For certain applications, the specific ordering of nodes and edges should not matter!

Invariance

$$f(PX, PAP^T) = f(X, A)$$

A = adjacency matrix

$$f\left(\begin{array}{c} \text{Graph 1} \end{array}\right) = \mathbf{y} = f\left(\begin{array}{c} \text{Graph 2} \end{array}\right)$$

Examples: **max**, **sum**, **min**, **avg**

$\oplus$  = any permutation-invariant aggregation op acting on one or more graph nodes / edges

## Permutation equivariance

What if we wanted to distinguish between outputs at different nodes? E.g.: node classification

A permutation-invariant aggregator would not allow us to do that ☹️

Instead, we may use functions that **don't change the node ordering**.

That is, if we permute nodes using a permutation matrix  $P$ , it doesn't matter if we do it before or after  $F$ ! 😊

$$\mathbf{F}(\mathbf{X}, \mathbf{A}) = \begin{bmatrix} - & \phi(\mathbf{x}_1, \mathbf{X}_{\mathcal{N}_1}) & - \\ - & \phi(\mathbf{x}_2, \mathbf{X}_{\mathcal{N}_2}) & - \\ & \vdots & \\ - & \phi(\mathbf{x}_n, \mathbf{X}_{\mathcal{N}_n}) & - \end{bmatrix} \quad F(PX, PAP^T) = PF(X, A)$$

If  $\phi$  is permutation invariant over the neighborhood  $\mathbf{x}_u, \mathbf{X}_{\mathcal{N}_u}$   $\mathbf{F}$  is permutation equivariant!

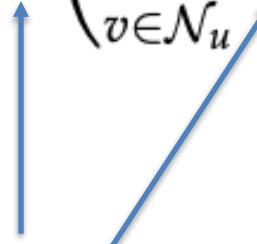


We stack multiple equivariant GNN layers to build large-scale operators:

$$\mathbf{F}(\mathbf{X}, \mathbf{A}) := \phi \left( \bigoplus_{v \in \mathcal{N}_u} \psi(x_u, x_v, x_{uv}) \right)$$

$\bigoplus$  = sum, average ... or any permutation-invariant aggregation op  
acting on one or more graph nodes / edges in a neighborhood

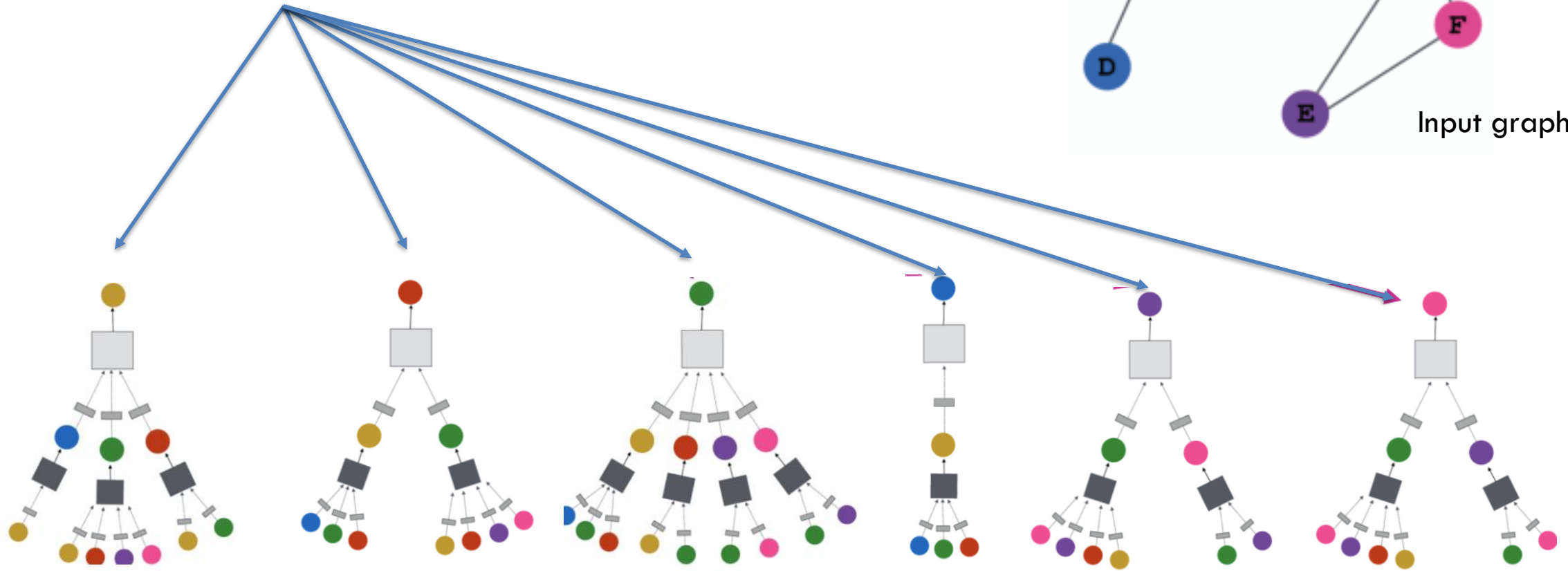
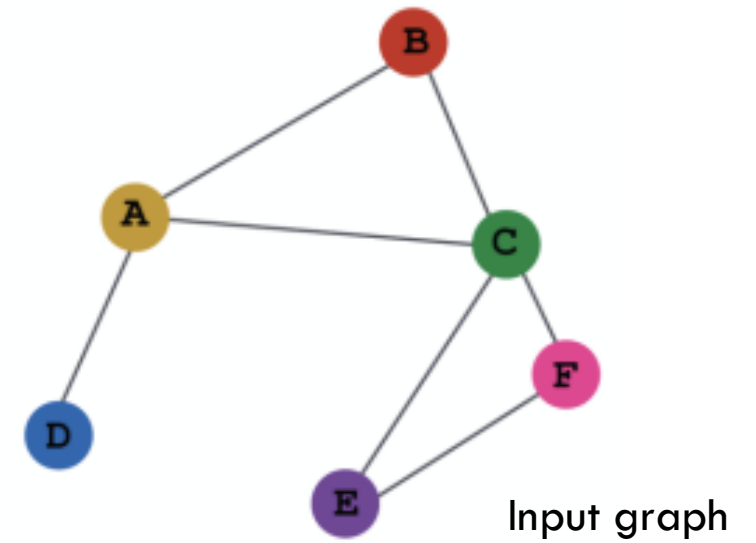
We've just defined a GNN layer!

$$\mathbf{F}(X, A) := \phi \left( \bigoplus_{v \in \mathcal{N}_u} \psi(x_u, x_v, x_{uv}) \right)$$


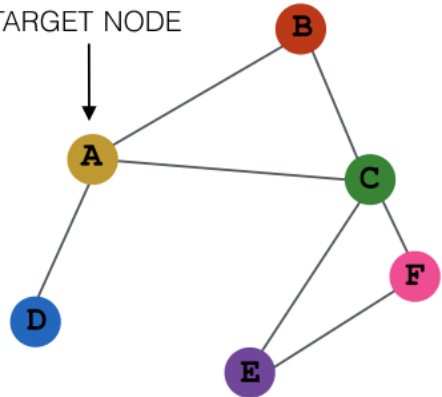
Trainable, shared MLPs

GNN layers are defined by the shared application of local, differentiable and permutation equivariant MLPs

Each node in the graph has its own computational graph, defined by the edge connectivity



TARGET NODE



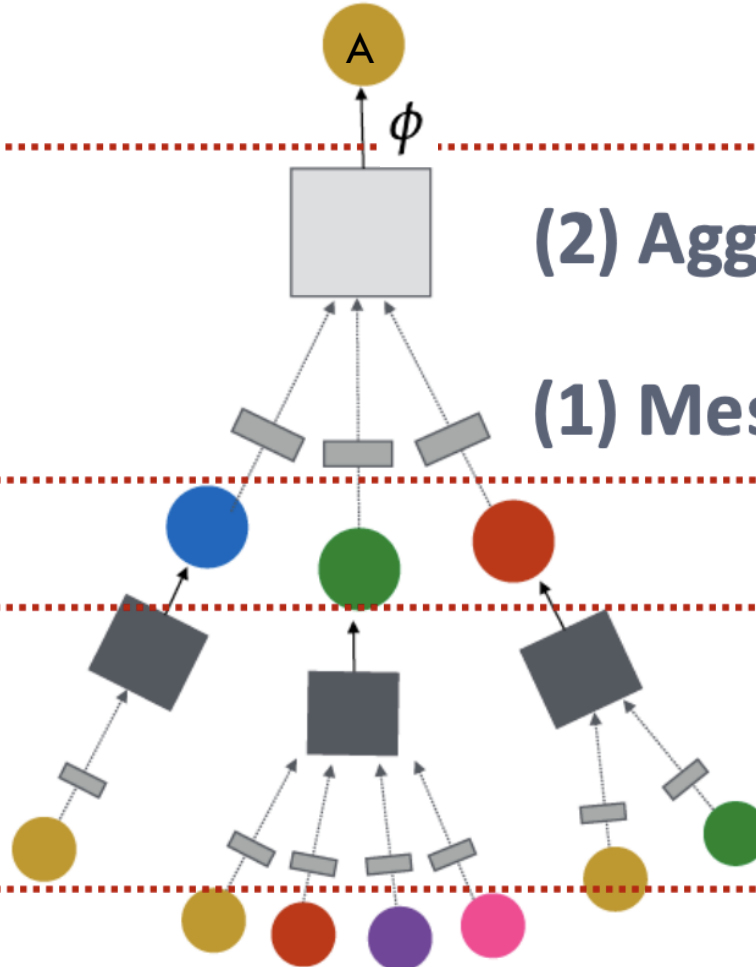
$$\mathbf{F}(\mathbf{X}, \mathbf{A}) := \phi \left( \bigoplus_{v \in \mathcal{N}_u} \psi(x_u, x_v, x_{uv}) \right)$$

**GNN Layer 2**

**(2) Aggregation  $\bigoplus$**

**(1) Message  $\psi$**

**GNN Layer 1**



## Quiz time

$$\mathbf{F}(\mathbf{X}, \mathbf{A}) := \phi \left( \bigoplus_{v \in \mathcal{N}_u} \psi(x_u, x_v, x_{uv}) \right)$$

$\mathbf{x}_{\mathcal{G}}$  Graph-level (global) features

$\mathbf{x}_u$  Node features

$\mathbf{x}_{uv}$  Edge features

$\mathbf{x}_v$

How would the **graph-level feature(s)**

$\mathbf{x}_{\mathcal{G}}$  in this framework?

## Quiz time

Inductive bias	Task
<b>Locality</b>	<b>All</b> (operators act over neighborhoods)
<b>Invariance</b>	<b>Graph classification</b> (e.g. “bad” vs “good” protein structure)  <b>Graph regression</b> (e.g. molecular energy prediction)  <b>Edge (link) prediction</b>
<b>Equivariance</b>	<b>Node classification / regression</b>  <b>3D molecular dynamics</b> (rotation/translation)

Can we (and should we?) break permutation equivariance?



## Quick detour: MLPs

```
from torch import nn

def generate_mlp_module(num_inputs: int = 32, hidden_dim: int = 64, num_outputs: int = 32):
    mlp = nn.Sequential(
        nn.Linear(num_inputs, hidden_dim),
        nn.LeakyReLU(0.1),
        nn.Linear(hidden_dim, hidden_dim),
        nn.LeakyReLU(0.1),
        nn.Linear(hidden_dim, num_outputs),
        nn.LeakyReLU(0.1),
        nn.LayerNorm(num_outputs)
    )
    return mlp
```

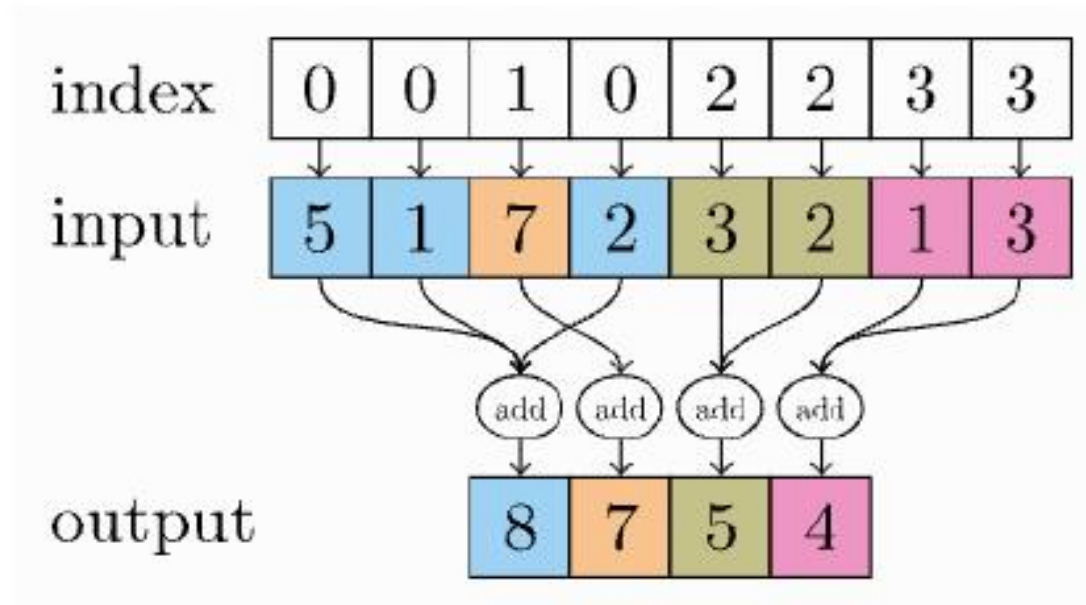
(..., num\_neighbors, **num\_outputs**) = **MLP** (... , num\_neighbors, **num\_inputs**)

Recall: MLPs (1) act on neighborhoods and (2) share the weights.

**Q: what does this imply wrt **num\_neighbors**?**

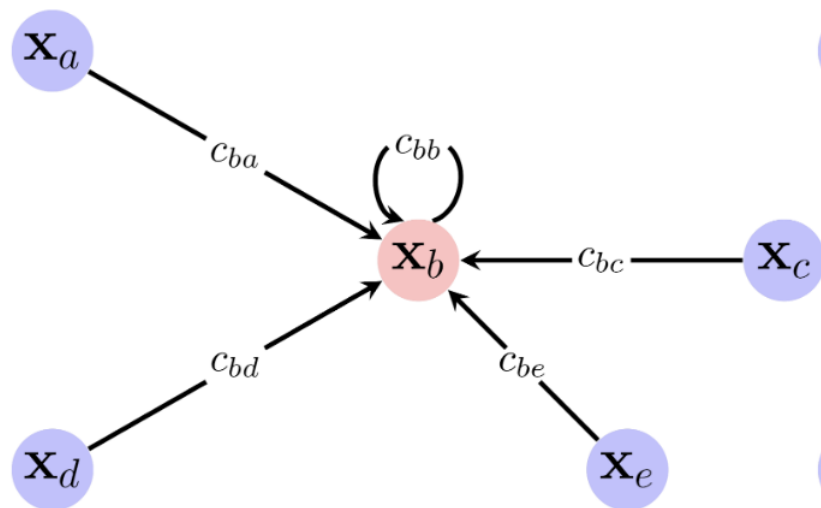
## Two basic ingredients of GNNs

- **Shared MLPs**  $\phi$ ,  $\psi$  and  $\rho$  operate on node and edge features (we'll see how)
- **Sparse index-gather / -scatter operations** over graph neighborhoods (GPU-optimized)

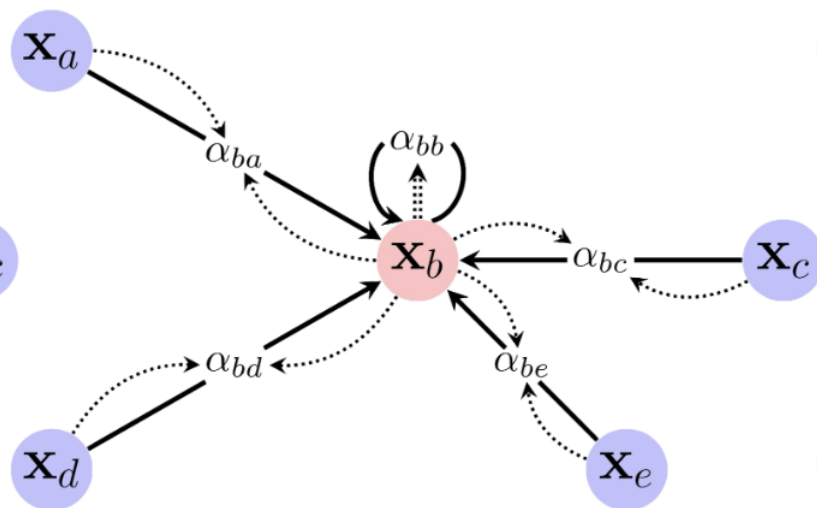


# Flavors of GNNs

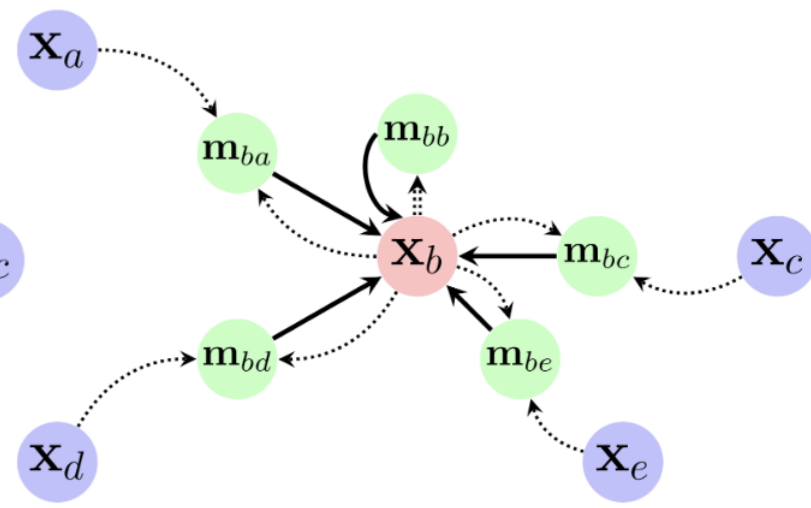
$$h_i = \phi(x_i, \bigoplus_{j \in \mathcal{N}_i} \alpha(x_i, x_j) \psi(x_j))$$



Convolutional



Attentional



Message-passing

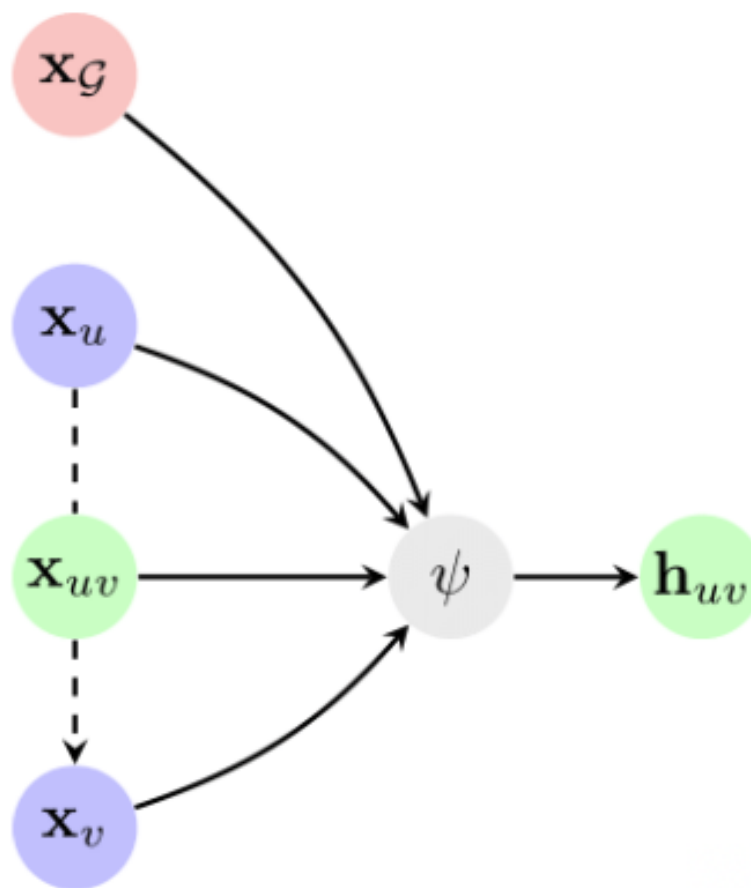
$$h_i = \phi(x_i, \bigoplus_{j \in \mathcal{N}_i} c_{ij} \psi(x_j))$$

$$h_i = \phi(x_i, \bigoplus_{j \in \mathcal{N}_i} \psi(x_i, x_j, e_{ij}))$$

$$\mathbf{m}_{ij} := \psi(x_i, x_j, e_{ij})$$

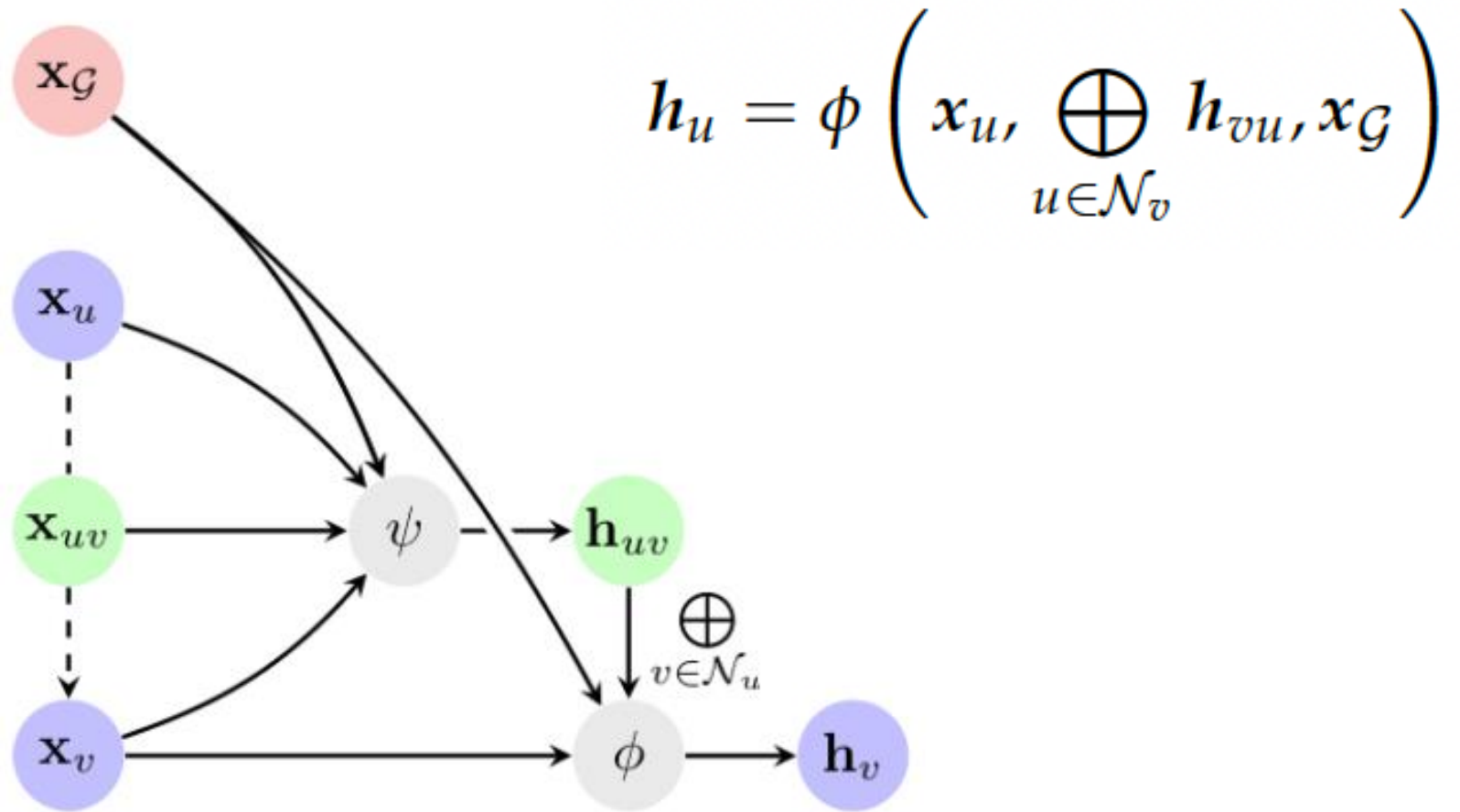
## Message-passing GNNs

## Step 1: Edge updates



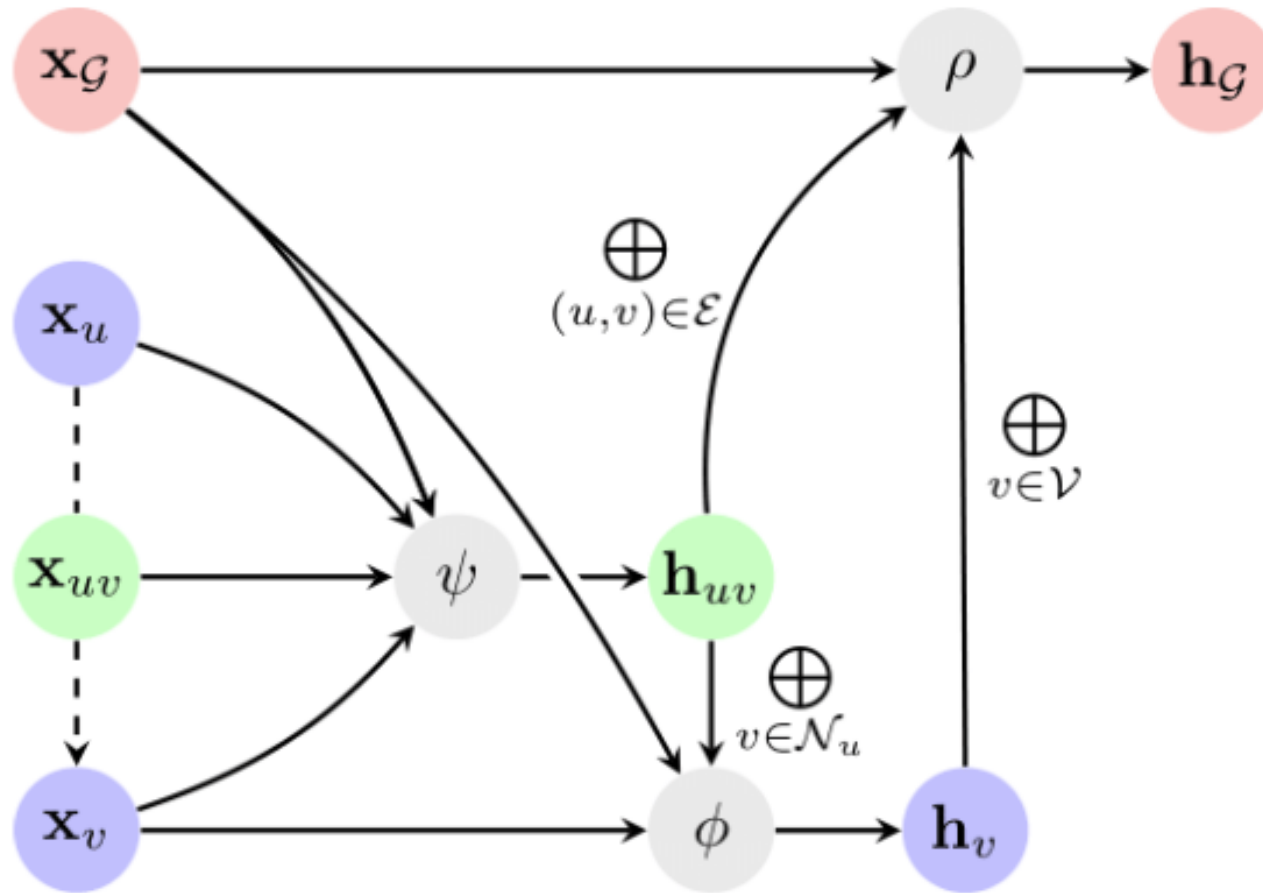
$$h_{uv} = \psi(x_u, x_v, x_{uv}, x_G)$$

## Step 2: Node updates





### Step 3: Graph feature updates



$$h_G = \rho \left( \bigoplus_{u \in \mathcal{V}} h_u, \bigoplus_{(u,v) \in \mathcal{E}} h_{uv}, x_G \right)$$

## The message-passing algorithm

---

**Input:** Graph  $\mathcal{G}(\mathcal{V}, \mathcal{E})$  with  $\{x_{\mathcal{G}}, h_u, h_{uv}\}$ .

**for** each edge  $e_{uv}$  **do**

    Gather sender and receiver nodes  $x_u, x_v$

    Update edge  $h_{uv} \leftarrow \psi(x_u, x_v, x_{uv}, x_{\mathcal{G}})$

**end for**

**for** each node  $u$  **do**

    Aggregate all incoming edges to  $u$ :  $h_u^* := \bigoplus_{v, (v,u) \in \mathcal{E}} h_{vu}$

    Compute node-wise features  $h_u \leftarrow \phi(x_u, h_u^*, x_{\mathcal{G}})$

**end for**

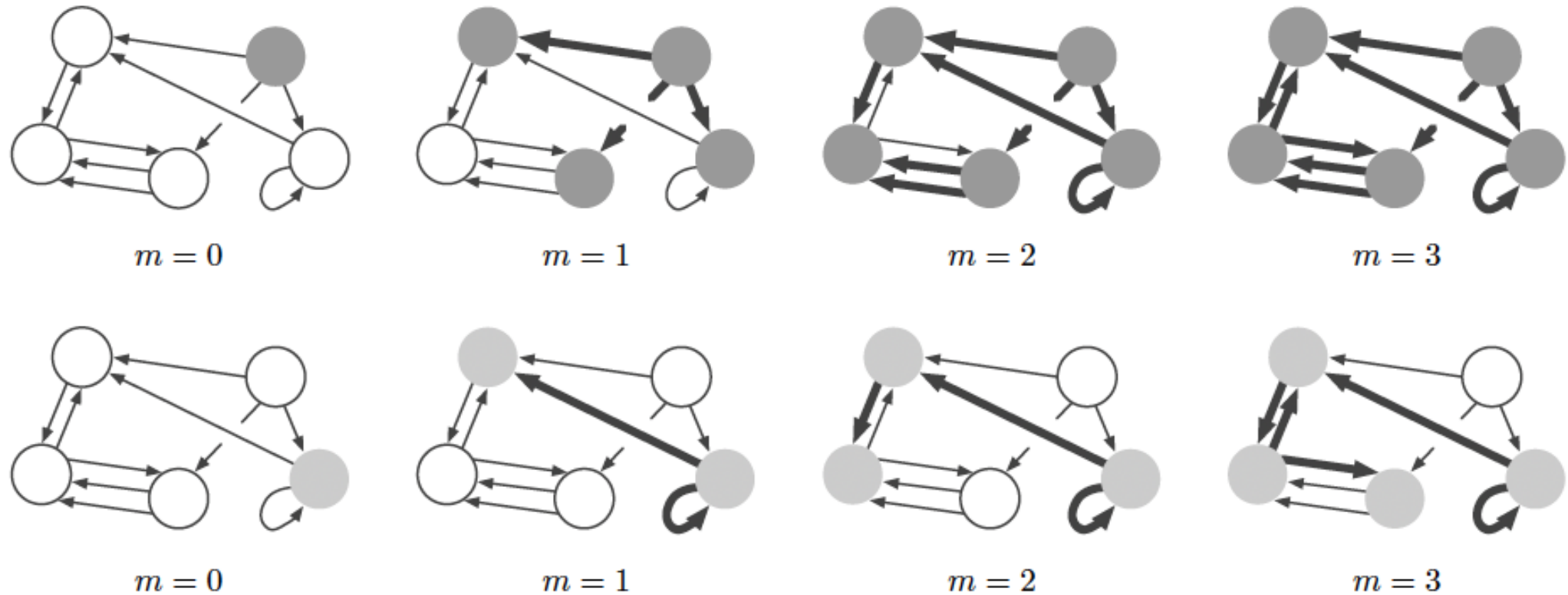
Aggregate all edges and nodes  $u^* := \bigoplus_{u \in \mathcal{V}} h_u, \mathbf{e}^* := \bigoplus_{(u,v) \in \mathcal{E}} h_{uv}$

Compute global features  $h_{\mathcal{G}} \leftarrow \rho(x_{\mathcal{G}}, u^*, \mathbf{e}^*)$

**Output:** Graph  $\mathcal{G}$  with new  $\{x_{\mathcal{G}}, h_u, h_{uv}\}$ .

---

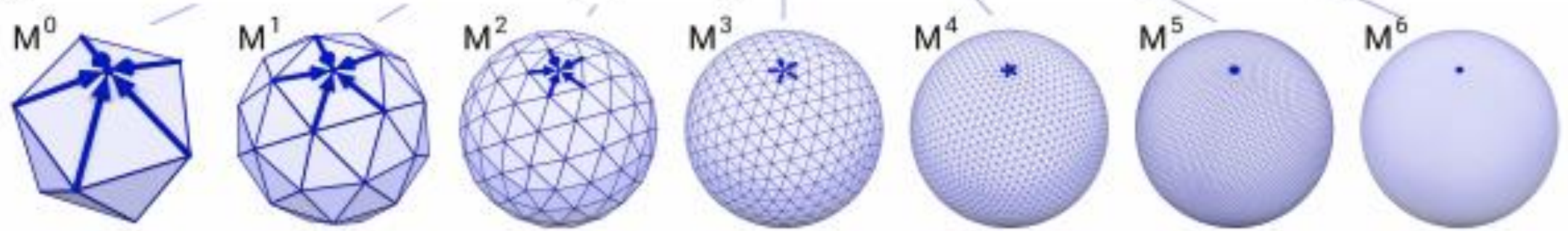
## Message passing: information propagation



NB: This happens **simultaneously** for all nodes in the graph!

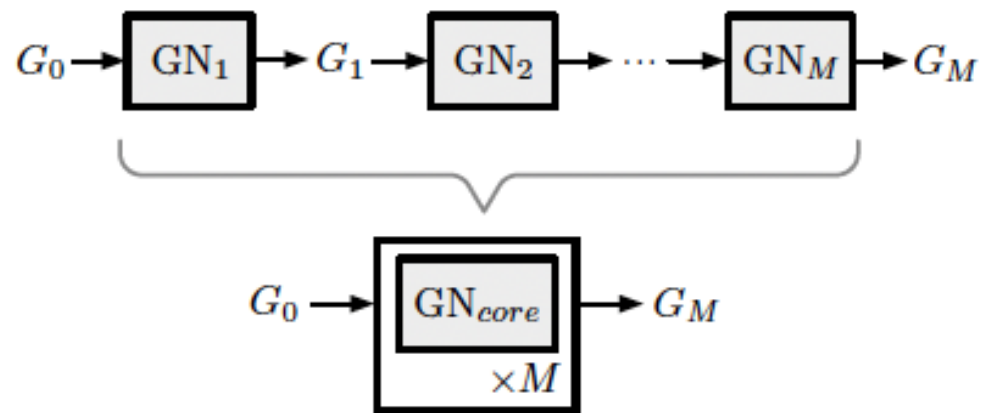
Figure from ([Battaglia et al., 2018](#))

g) Simultaneous multi-mesh message-passing

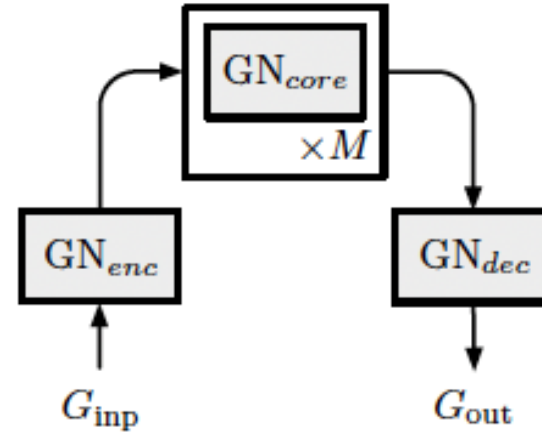


The multi-mesh allows information to propagate faster, across longer distances

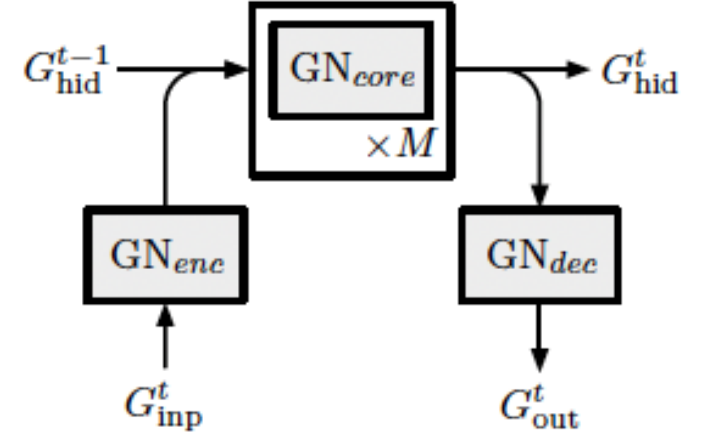
## GNN block structures



(a) Composition of GN blocks



(b) Encode-process-decode



(c) Recurrent GN architecture

GraphCast, AIFS and GraphDOP use both (a) and (b)



[https://github.com/pyg-team/pytorch\\_geometric](https://github.com/pyg-team/pytorch_geometric)



<https://github.com/ecmwf/anemoi-core/tree/main/graphs>





# anemoi graphs

- Create a new graph:

```
>>> anemoi-graphs create recipe.yaml graph.pt
```

- Describe an existing graph:

```
>>> anemoi-graphs describe graph.pt
```

- Inspect visually an existing graph:


```
>>> anemoi-graphs inspect graph.pt graph_viz/
```

```
graph_viz
├── data_to_hidden.html
├── distribution_edge_attributes.png
├── distribution_node_adjacency.png
├── distribution_node_attributes.png
├── hidden_to_data.html
├── hidden_to_hidden.html
└── isolated_nodes.html
```


Local files generated to inspect graphs.

<https://anemoi.readthedocs.io/projects/graphs/en/latest/>


<https://github.com/ecmwf/anemoi-core/tree/main/graphs>




Path




Format version: 0.0.1



Size



Graph ready, last update 17 seconds ago.



Statistics ready.

</

Console log when describing/inspecting a graph with anemoi-graphs.

Note: The inspection tools provided are designed for testing different graph configuration but it is not recommended for high-resolution graphs with a high number of nodes/edges.



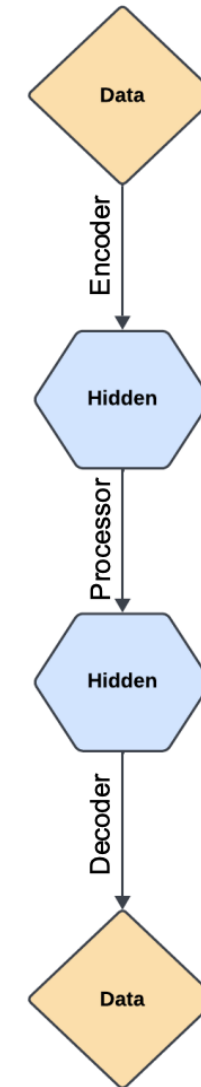
# anemoi graphs

## Graph recipe

recipe.yaml

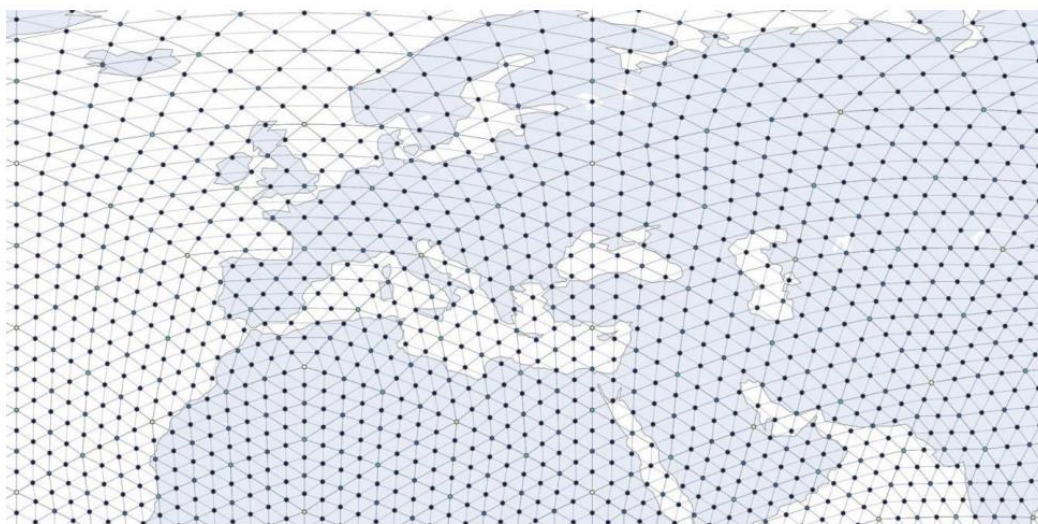
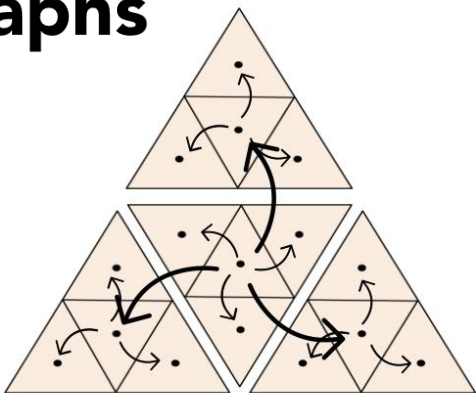
```
nodes:
  data:
    node_builder:
      _target_: anemoi.graphs.nodes.ZarrDatasetNodes
      dataset: my_zarr_dataset
  hidden:
    node_builder:
      _target_: anemoi.graphs.nodes.TriNodes
      resolution: 5 # num of refinements

edges:
  # Encoder configuration
  - source_name: data
    target_name: hidden
    edge_builders:
      - _target_: anemoi.graphs.edges.CutOffEdges
        cutoff_factor: 0.6
  # Processor configuration
  - source_name: hidden
    target_name: hidden
    edge_builders:
      - _target_: anemoi.graphs.edges.MultiScaleEdges
        x_hops: 1
  # Decoder configuration
  - source_name: hidden
    target_name: data
    edge_builders:
      - _target_: anemoi.graphs.edges.KNNEdges
        num_nearest_neighbours: 3
```





**anemoi**  
graphs



Multi-scale graph edges

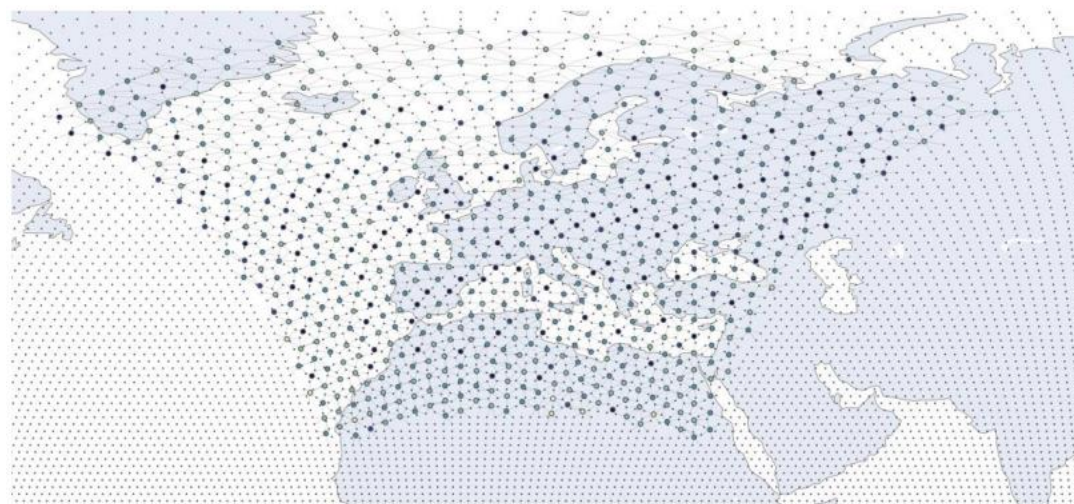
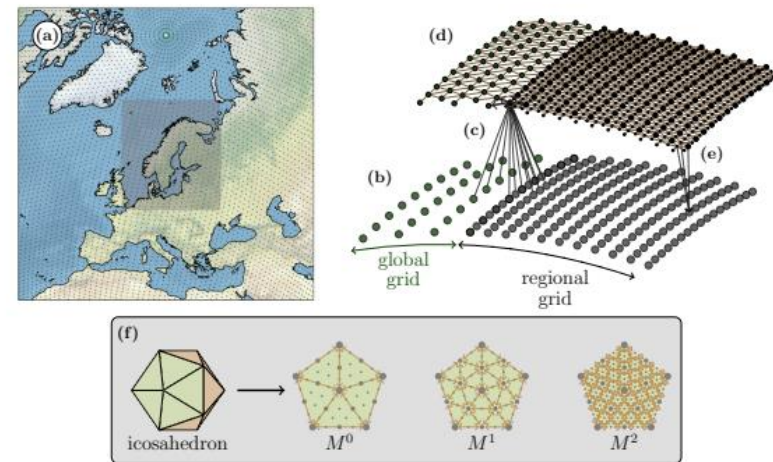


Diagram of encoder connections.

Regional or limited-area modeling

<https://arxiv.org/abs/2409.02891>

<https://arxiv.org/abs/2507.18378>

## Further references

(Veličković, 2023) <https://arxiv.org/pdf/2301.08210.pdf>

(Keisler, 2022) <https://arxiv.org/abs/2202.07575>

(Lam et al., 2023) <https://arxiv.org/abs/2212.12794>

(Sanchez-Lengeling et al., 2021) <https://distill.pub/2021/gnn-intro/>

(Veličković, 2023) <https://geometricdeeplearning.com/lectures/>

(Battaglia et al., 2018) <https://arxiv.org/abs/1806.01261>

(Sanchez-Gonzalez et al., 2020) <https://arxiv.org/abs/2002.09405>

Transformers are fully connected attentional GNNs  
(+ a positional embedding)

$$A = \mathbb{1}\mathbb{1}^T$$

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

$$h_u = \phi \left( x_u, \bigoplus_{v \in \mathcal{V}} \alpha(x_u, x_v) \psi(x_v) \right)$$

Attention transformers learn a “soft adjacency”

