

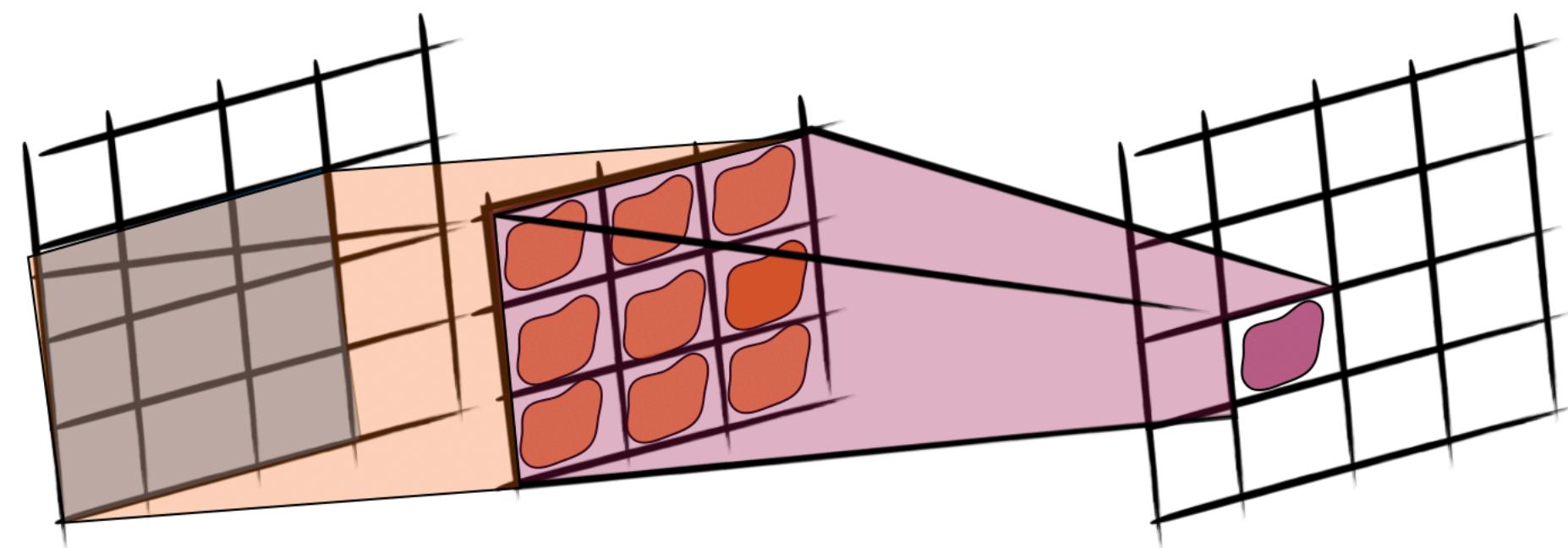
Modern Deep Learning in Physics

Sommersemester 2020/23

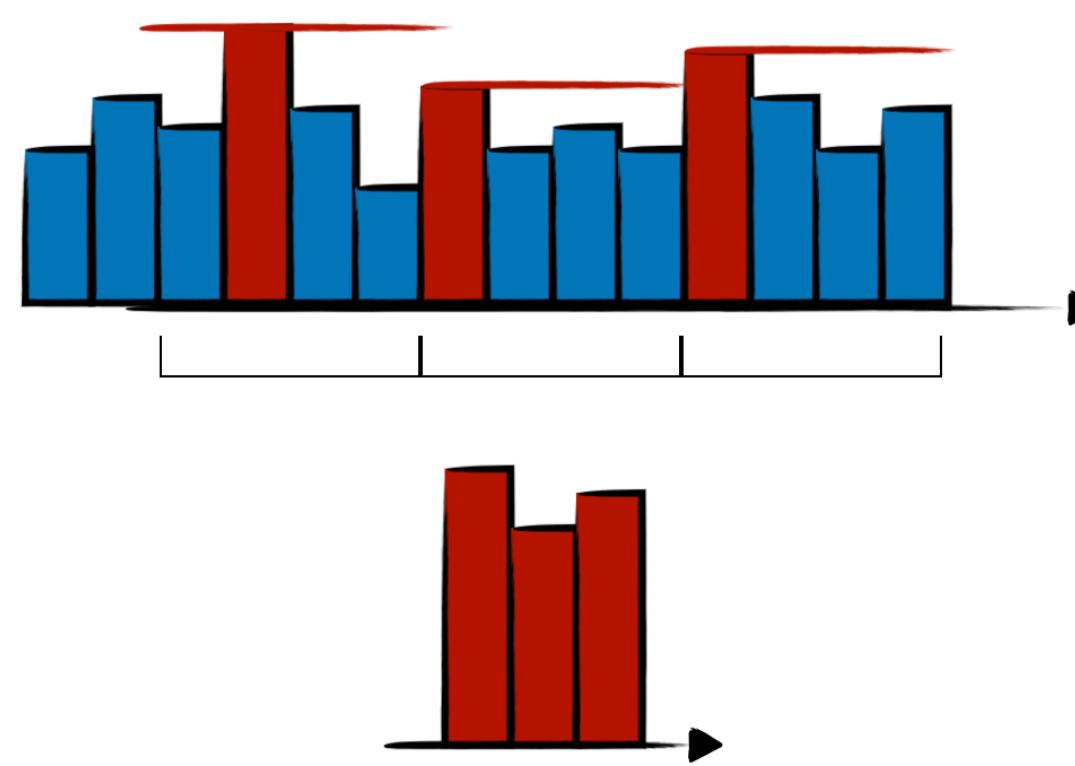
Lecture 8

Recap from Last Week

Convolutions

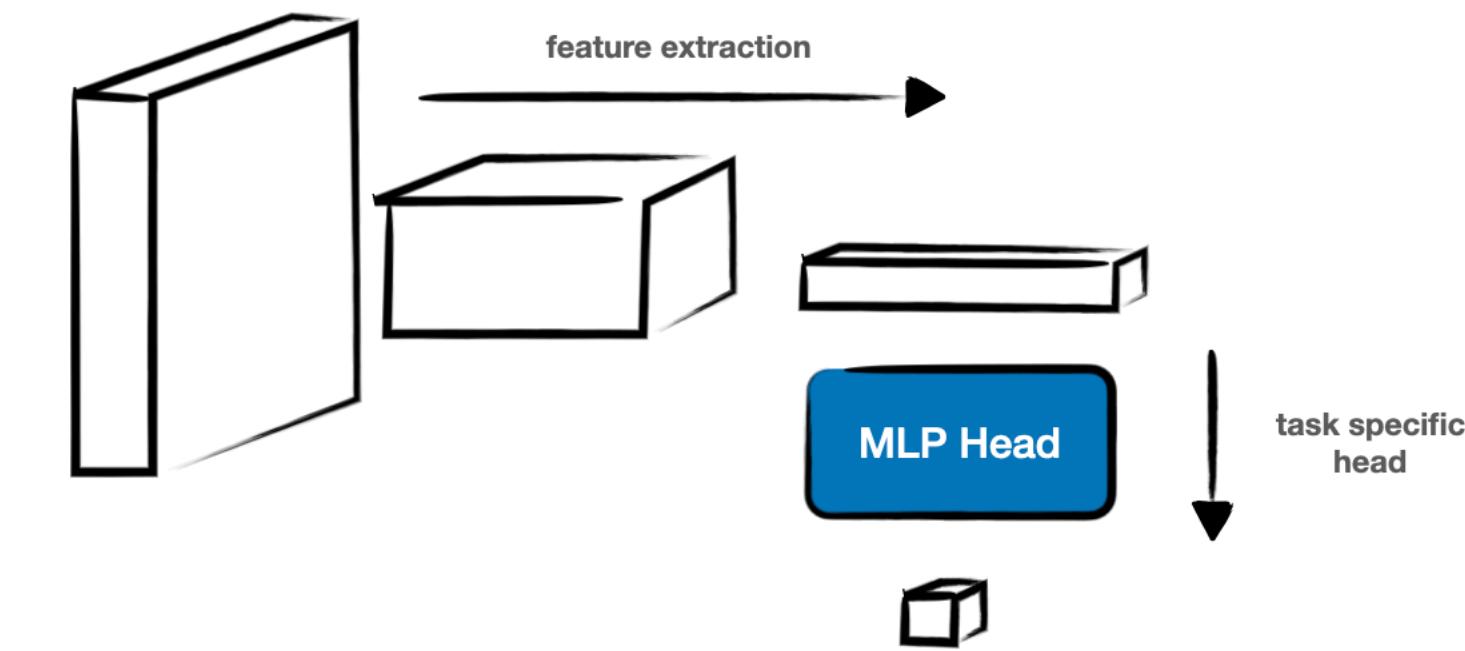


Translational Equi- & Invariance

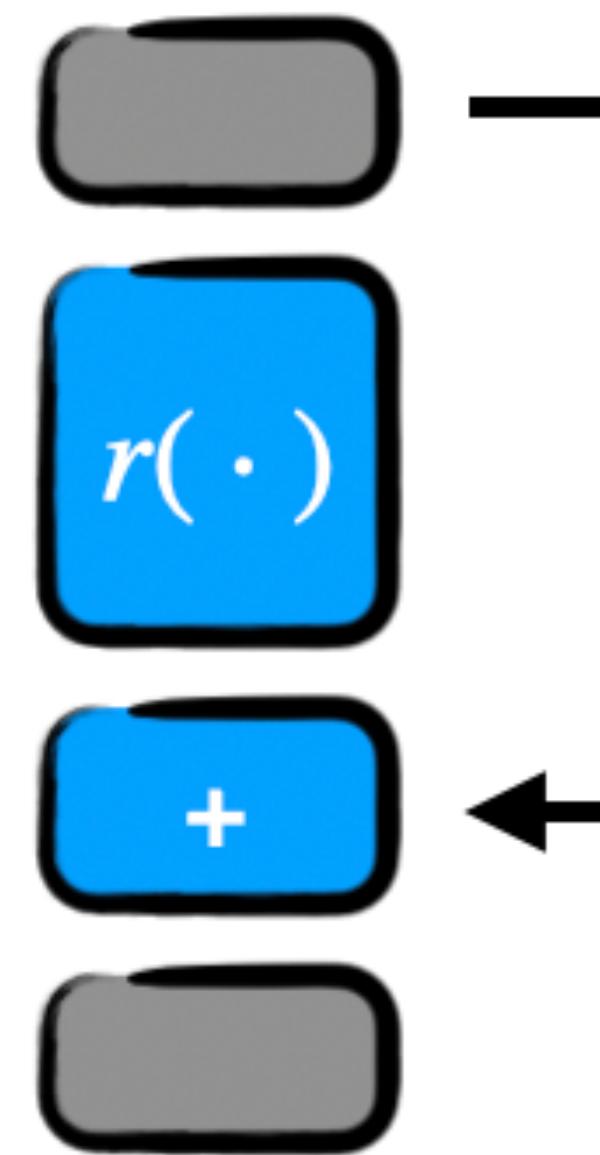


Raw Input

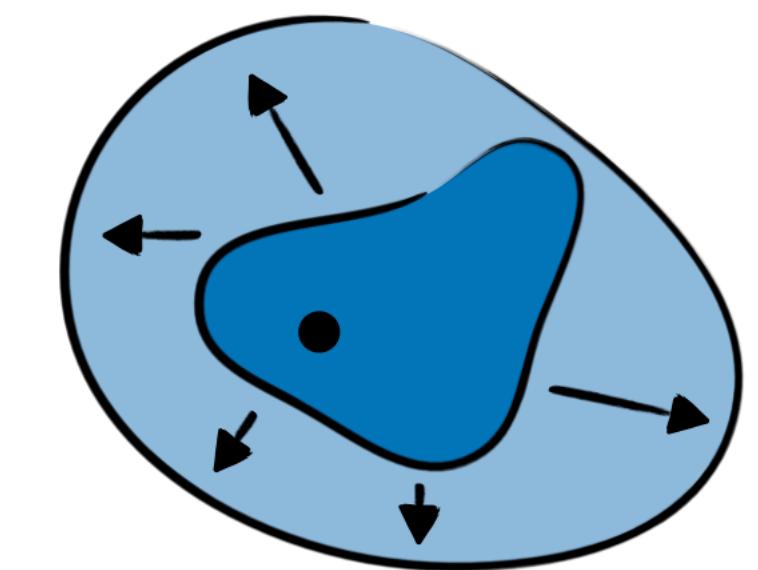
ConvNets



Residual Connections



Nested Function Spaces

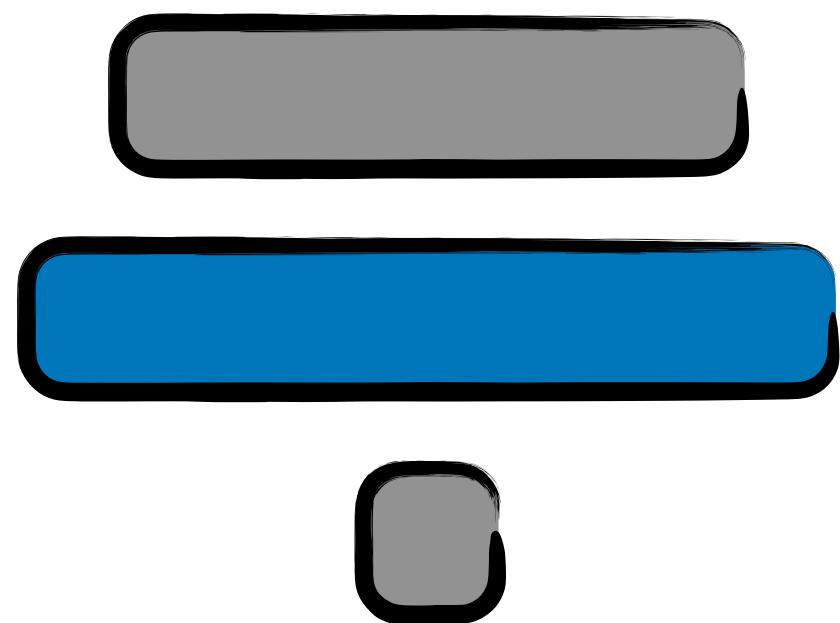


Short Recap: what we have achieved

We started with **universal function approximation**

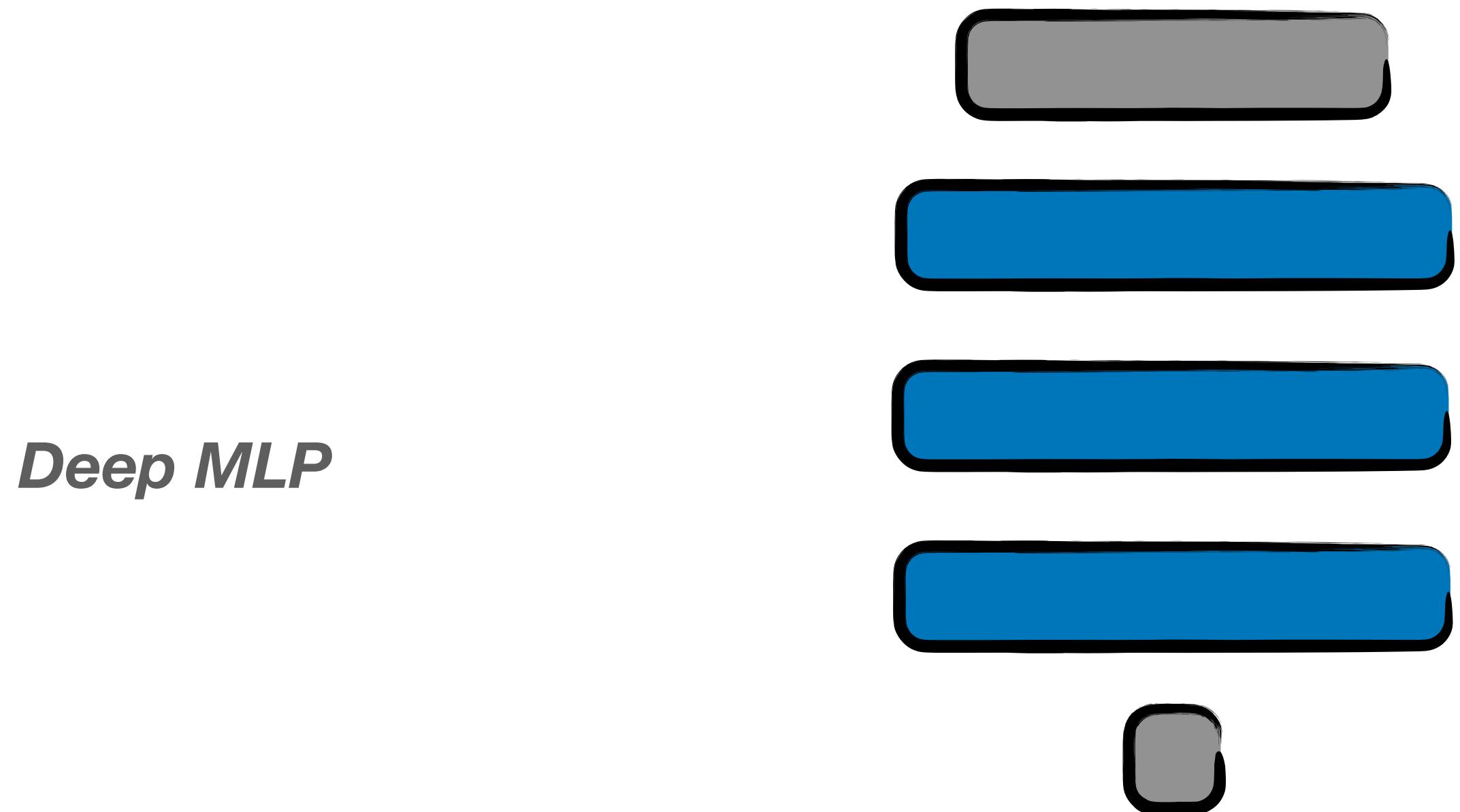
any $f(x) = \mathbb{R}^n \rightarrow \mathbb{R}$ can be represented by an MLP
(given enough data)

Single Layer MLP



Short Recap: what we have achieved

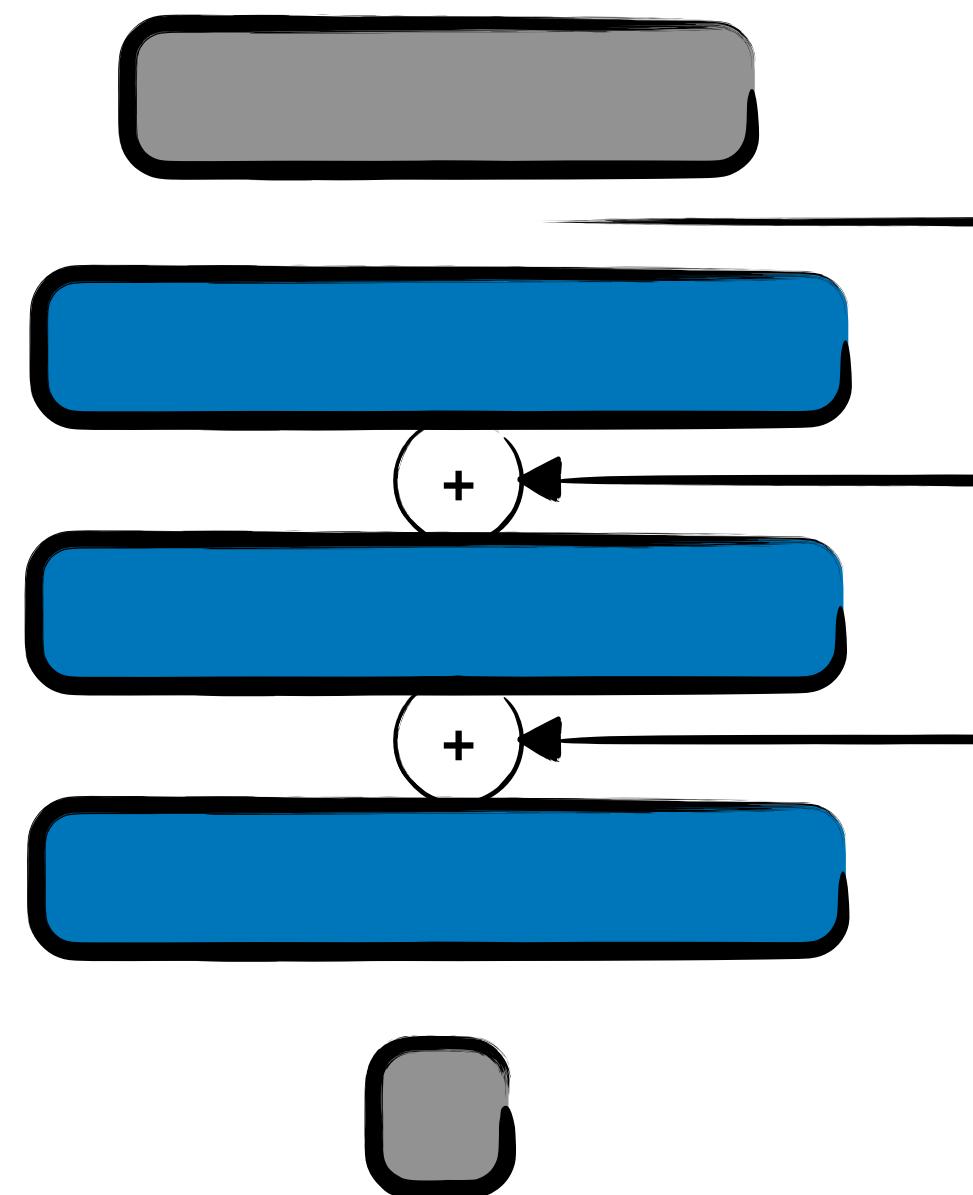
We discussed benefits of depth:
deep networks more expressive than shallow
with same # parameters



Short Recap: what we have achieved

residual connections produce nested function classes and

*Deep MLP
with residual
connections*

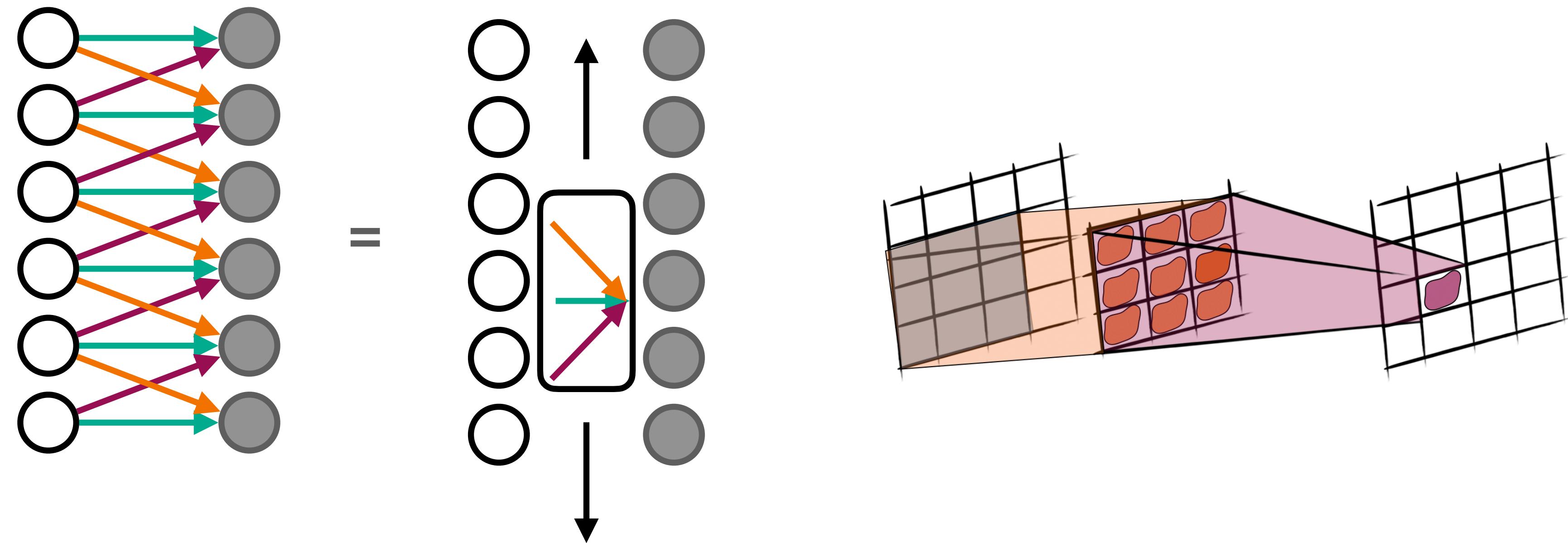


Short Recap: what we have achieved

**With CNN, we adapt the structure of the network
to the data type $\mathbb{R}^n = \mathbb{R}^{w \times h}$ to exploit inductive bias**

CNN: Two Key Ideas
• Local Connectivity
• Weight Sharing

*yields translation
equi- and
later in-variance*



A General Trend

This is a general trend in Deep Learning
specific data types \leftrightarrow network architectures

to encode known properties of the in function family

| <i>Data Type</i> | <i>Symmetry</i> | <i>Network</i> |
|-------------------------|---------------------|---------------------|
| <i>Grids/n-D arrays</i> | <i>Translation</i> | <i>CNN</i> |
| <i>Full Graph</i> | <i>Permutation</i> | <i>Transformers</i> |
| <i>Any Graphs</i> | <i>Permutation</i> | <i>GNN</i> |
| <i>Sets of Objects</i> | <i>Permutation</i> | <i>Deep Set</i> |
| <i>Sequences</i> | <i>Time Warping</i> | <i>LSTM</i> |
| ... | ... | ... |



CNN

GNN

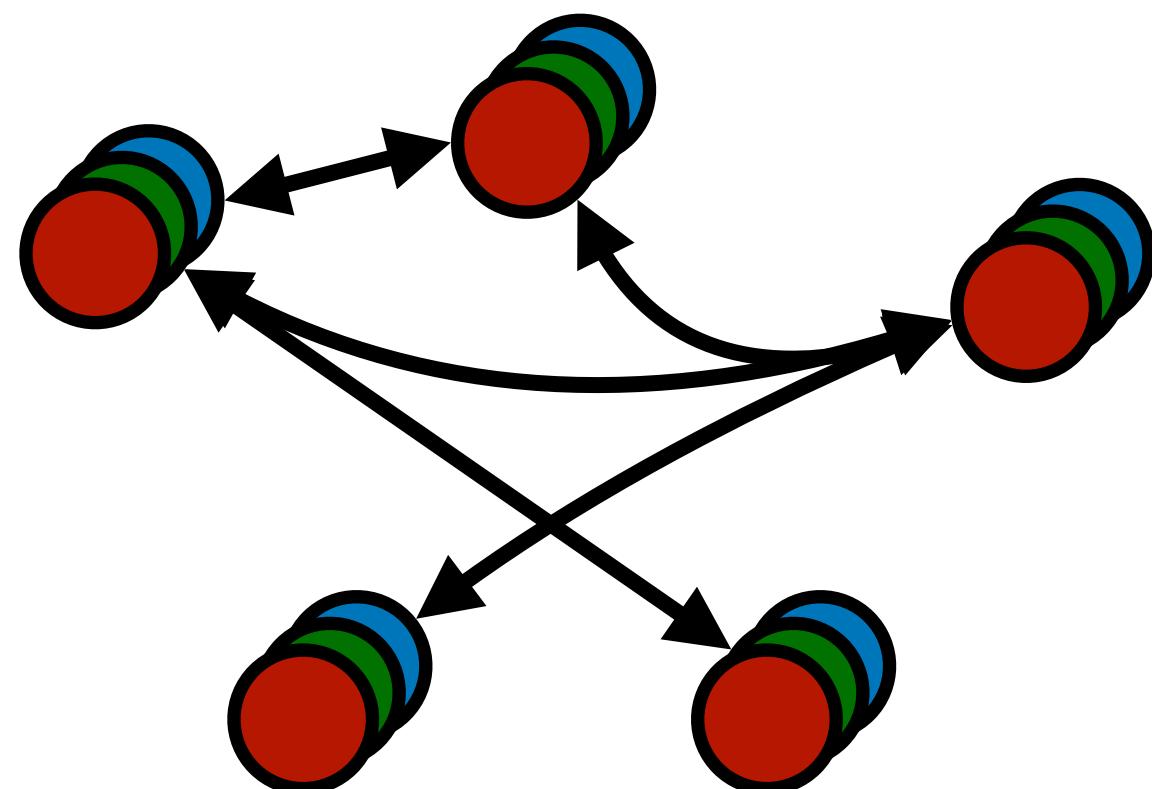
TRANSFORMER

DEEP SETS

RNN

Generalizing CNN

What should we do, if our data is not arranged in the, can we do something similar to CNNs & maintain similar features



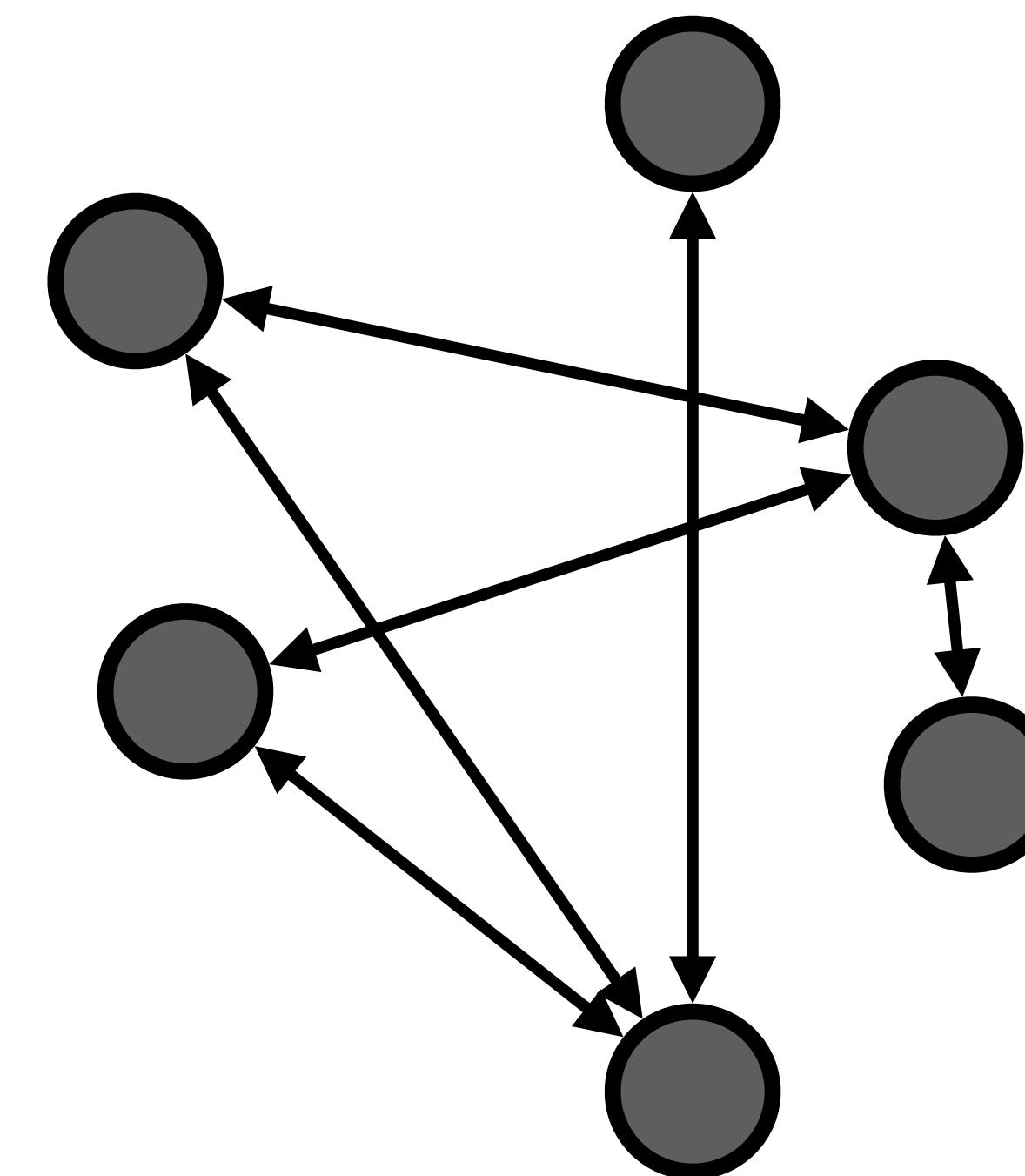
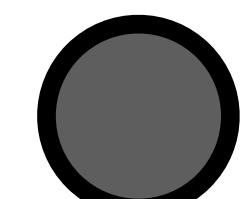
Generalize from *pixels connected spatially* to *sets of generic objects that are connected in some arbitrary way*

Graphs

Sets of non-regularly connected objects are best described as graphs

Defined as tuple (N, E)

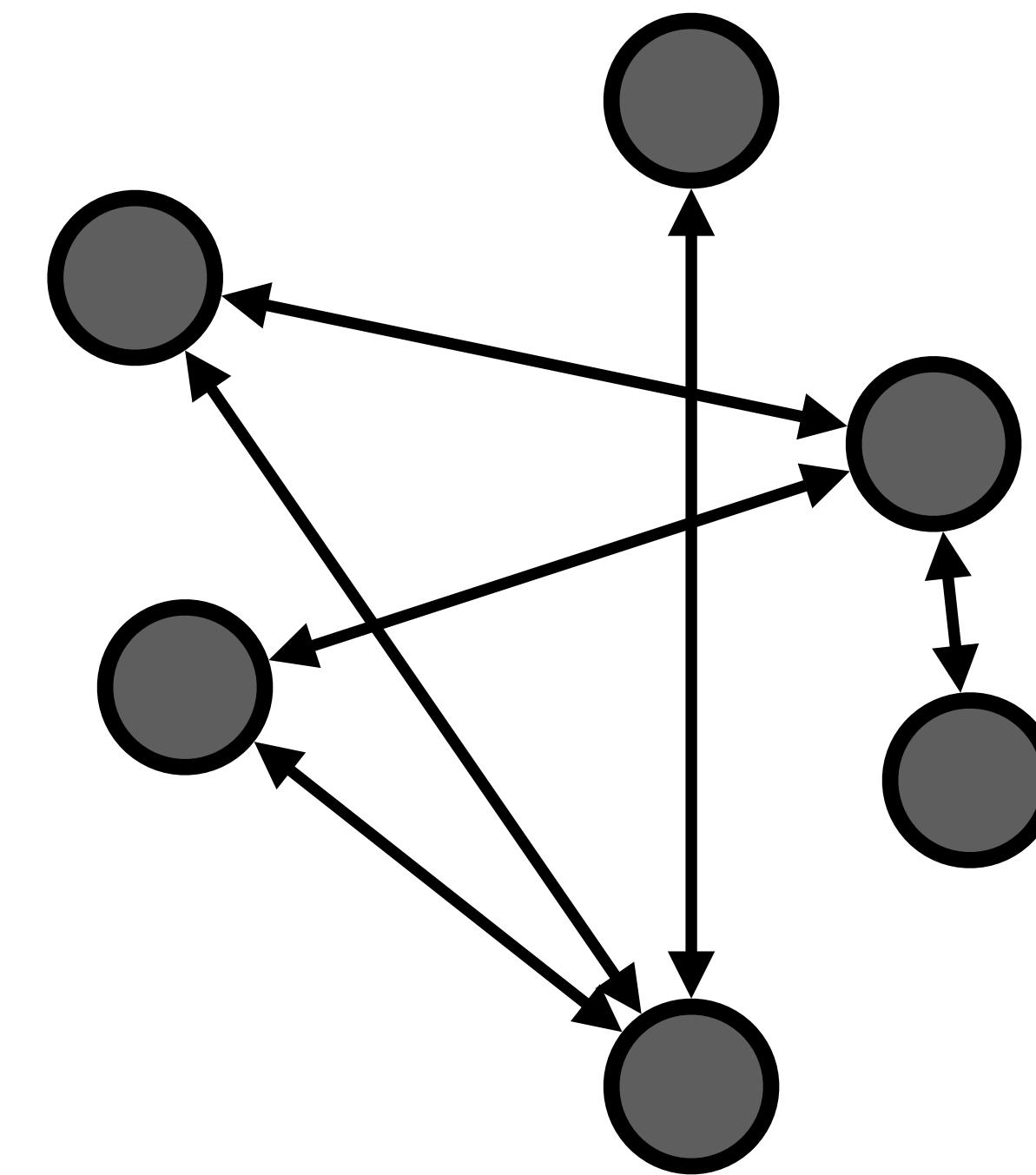
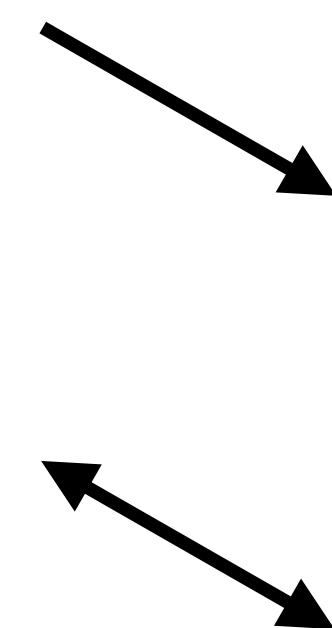
- Nodes (or Vertices)
- Edges



Directed and Undirected Graphs

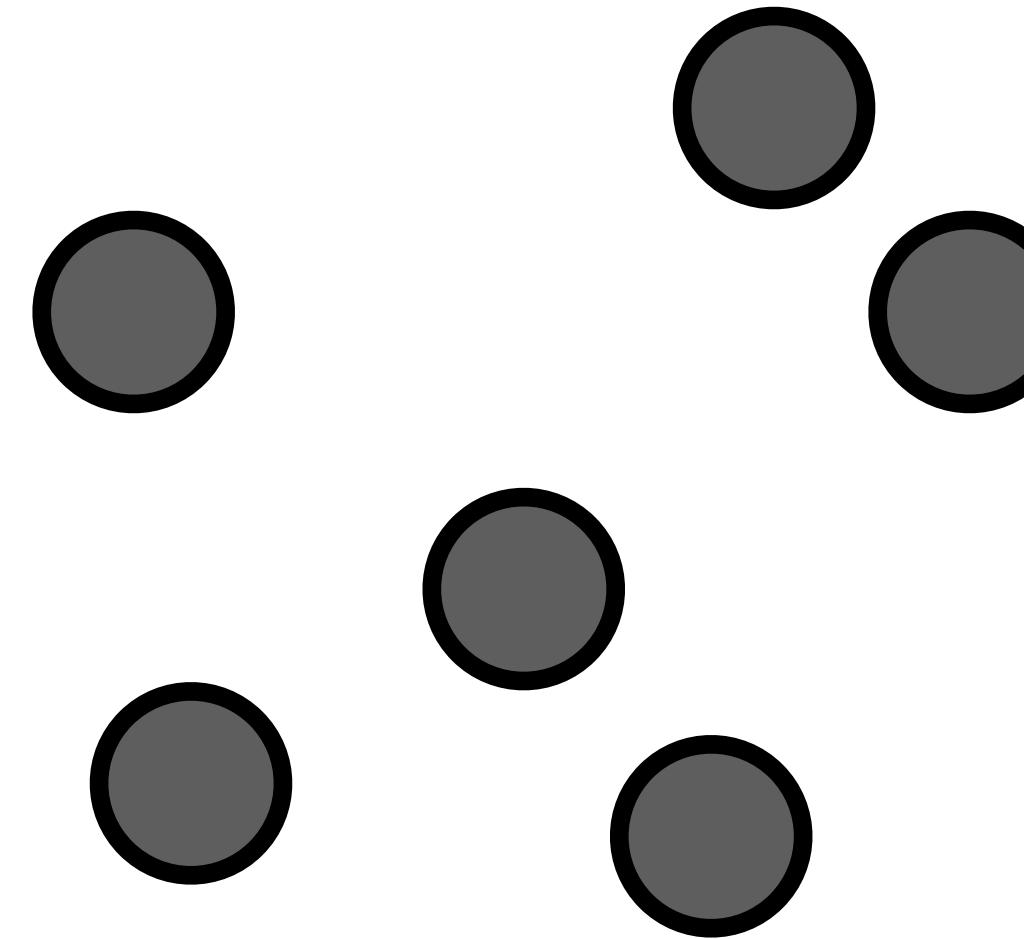
Edges can be either bidirectional or unidirectional

- Directed Edge
- Undirected Edge

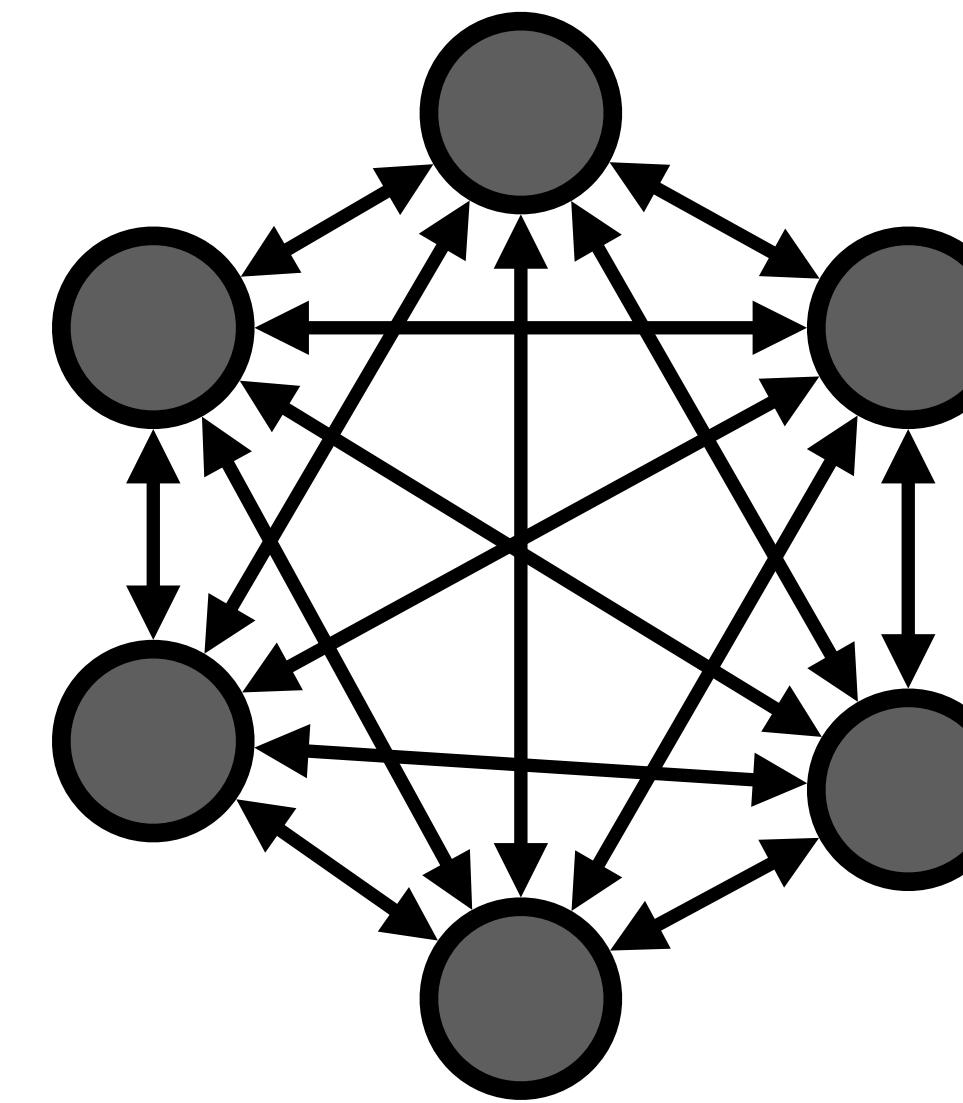


Graphs

Some special cases...



No connectivity (“bag of objects”)

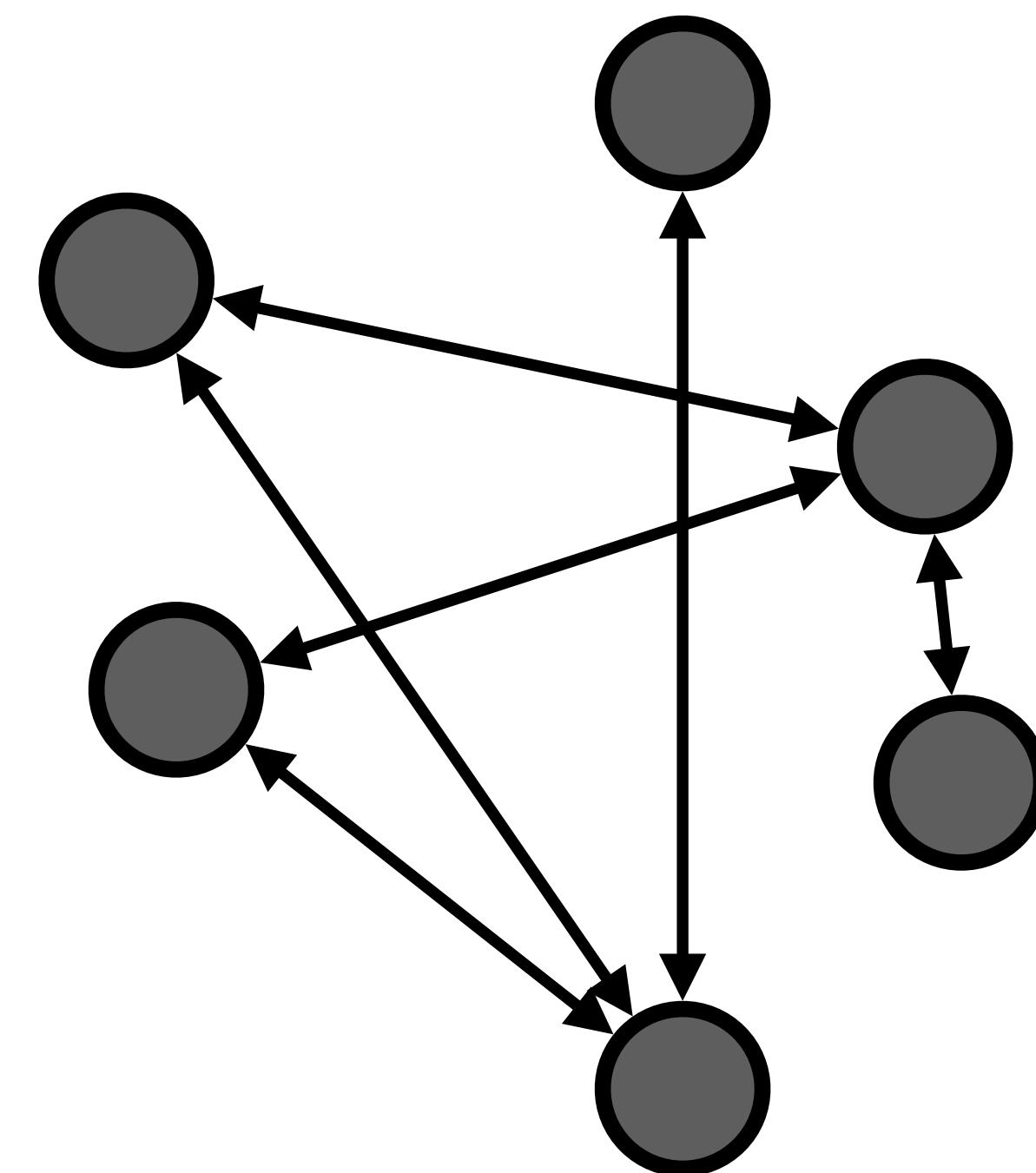


Full Connectivity

Graphs

The **graph structure** is the “substrate” onto which we put out data - it’s not the data itself

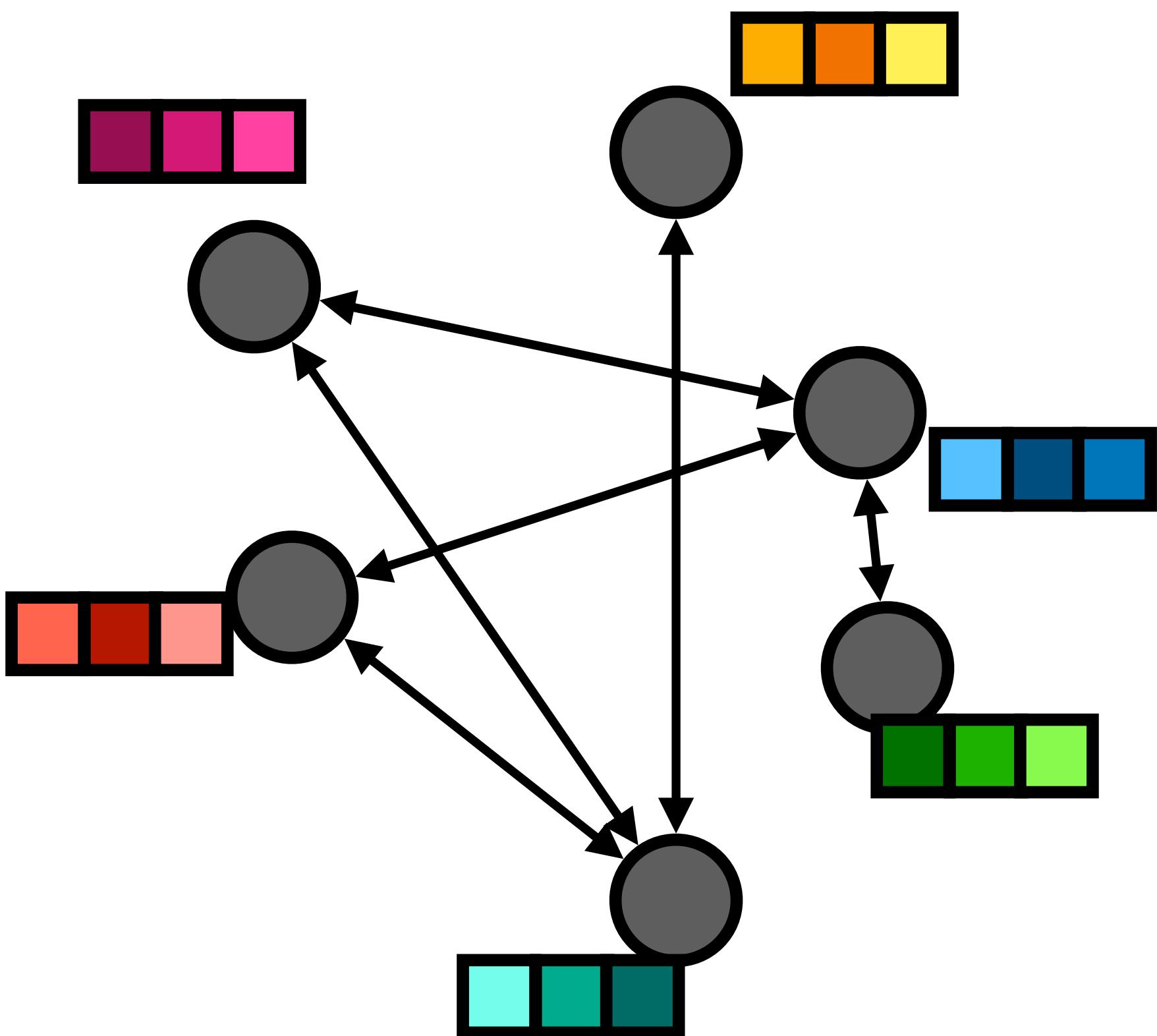
Can attach data to both of the key components of a graph



Graphs

The **graph structure** is the “substrate” onto which we put out data - it’s not the data itself

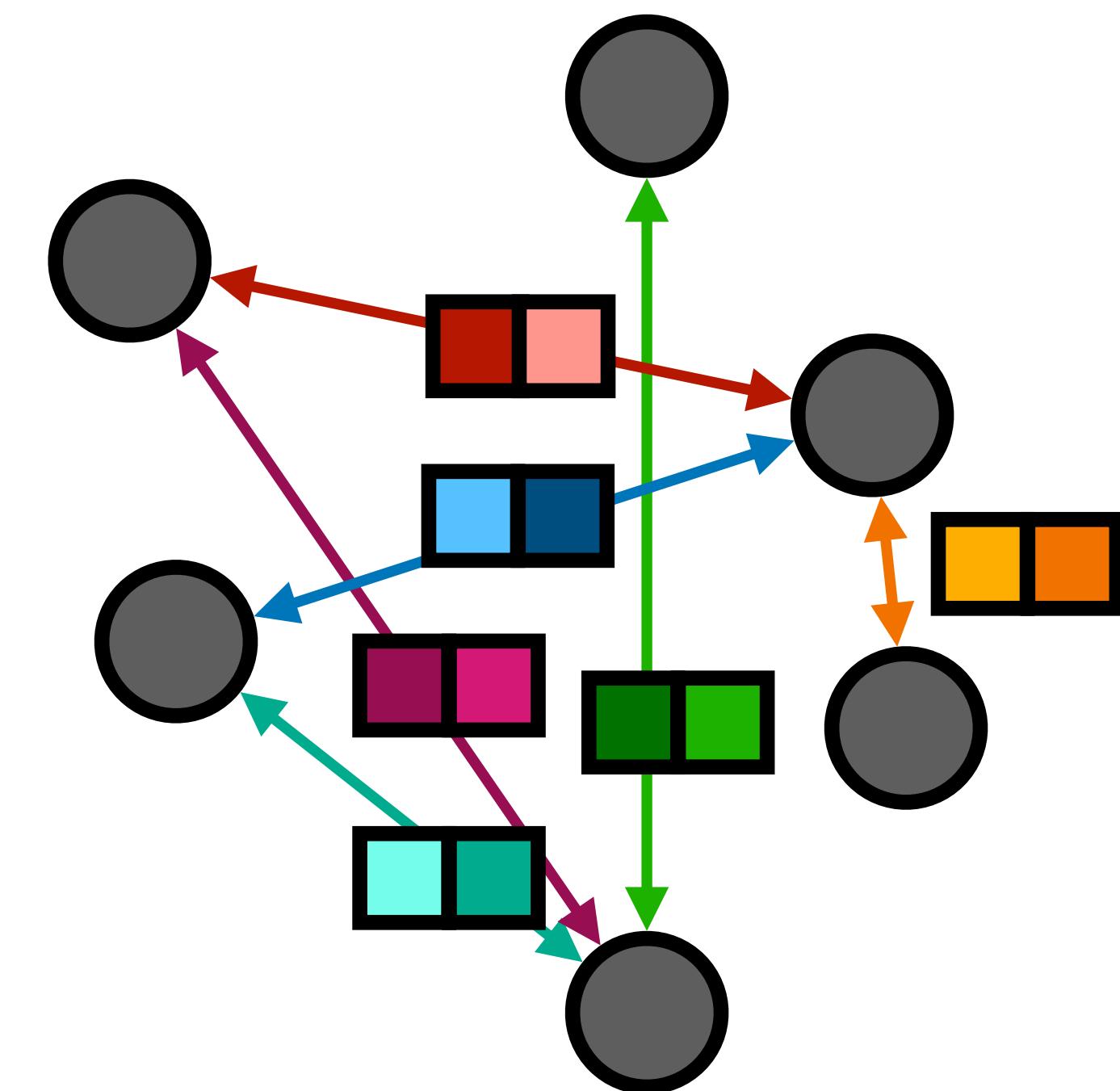
attach features \mathbb{R}^3 to nodes



Graphs

The **graph structure** is the “substrate” onto which we put out data - it’s not the data itself

attach features \mathbb{R}^2 to edges



Examples

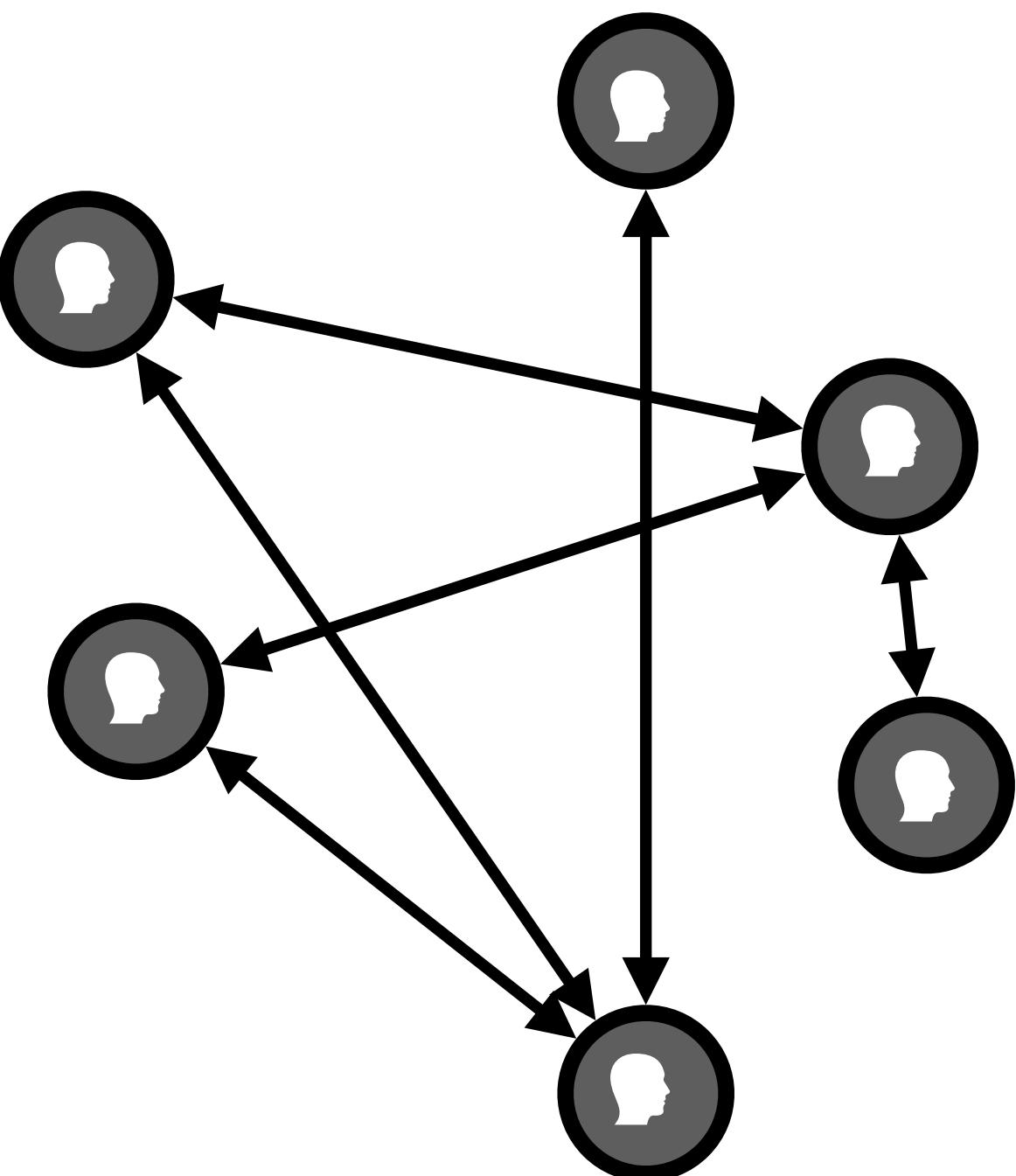
Graph structured data is everywhere!

Node Features

Age Hair Color
Height Location

Edge Features

Type of Relationship
Strength of Relationship
Length of Relationship



Social Networks

Examples

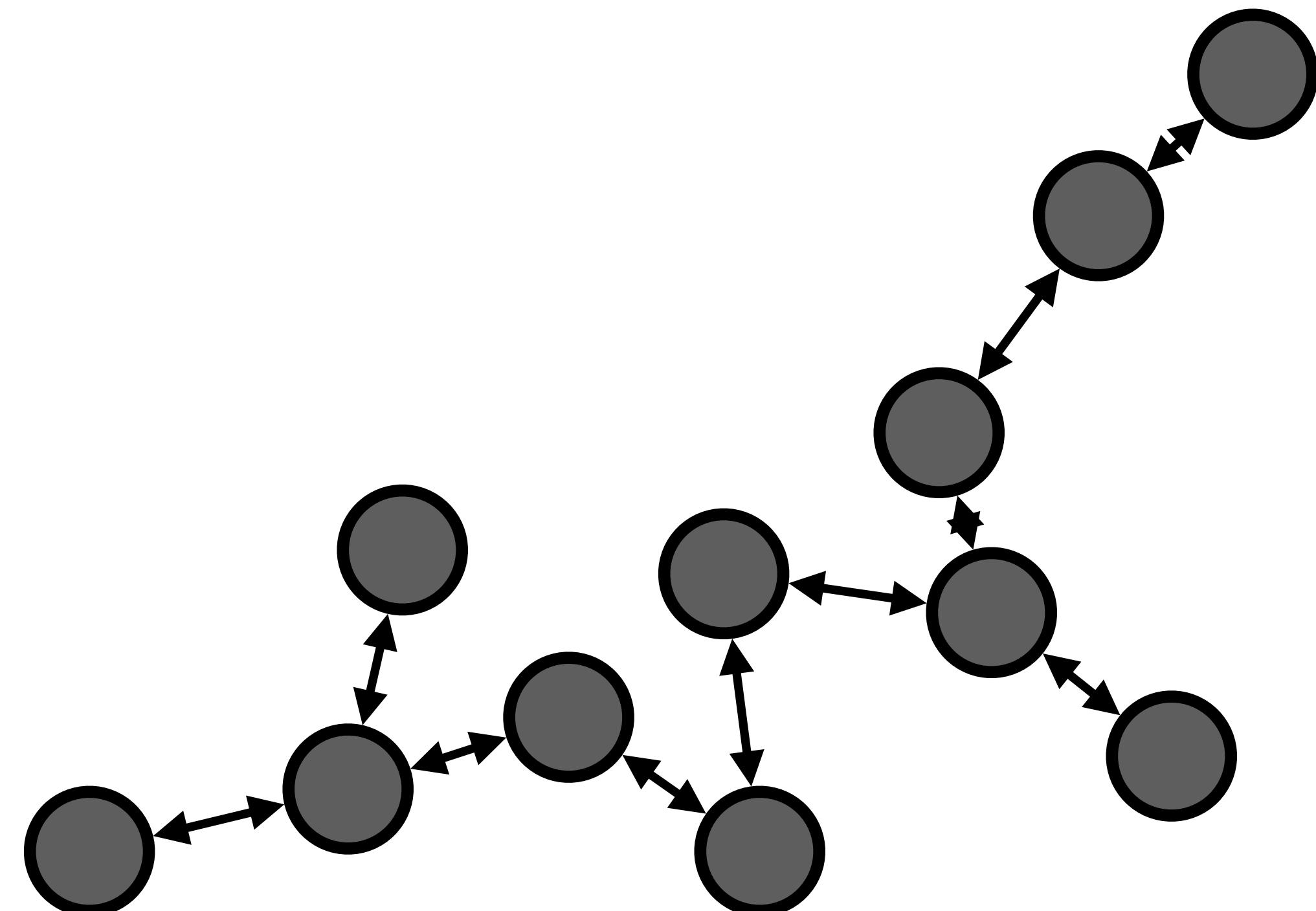
Graph structured data is everywhere!

Node Features

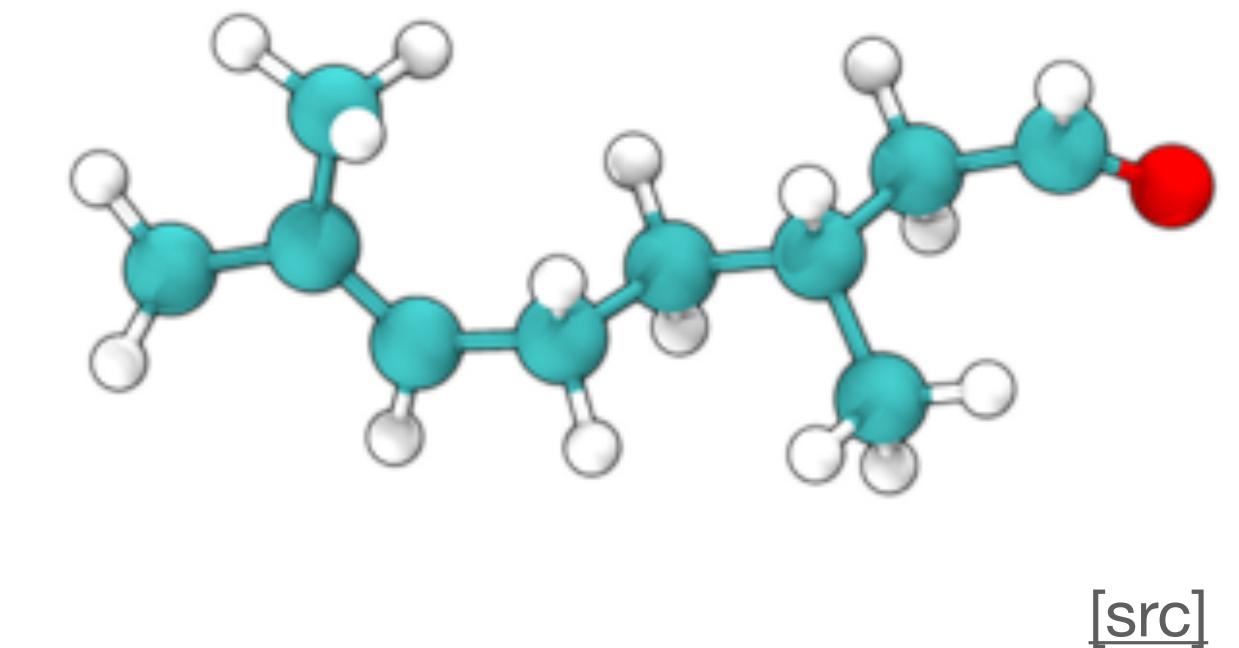
Position

Atom Type

Edge Features



Molecules



Examples

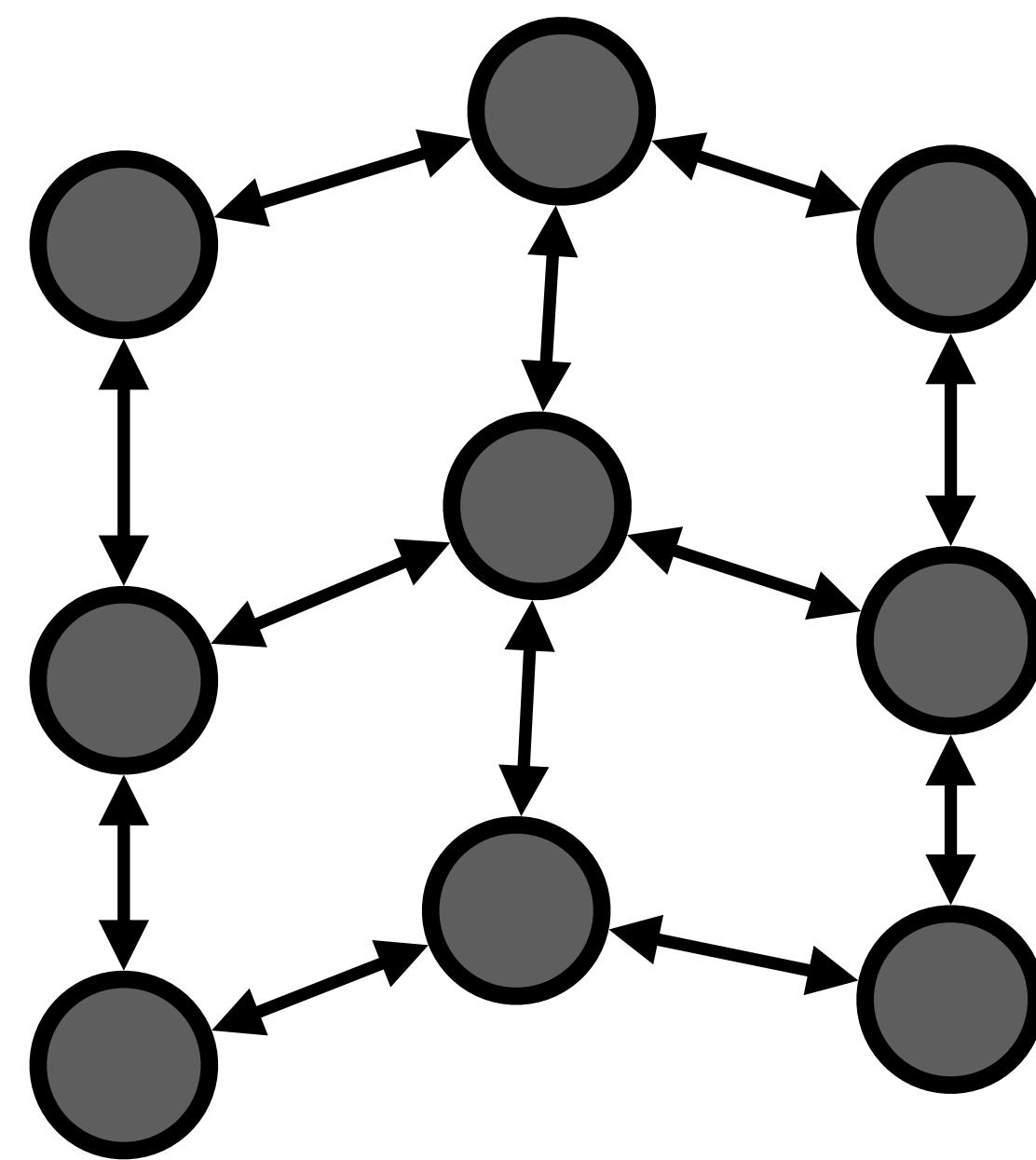
Graph structured data is everywhere!

Node Features

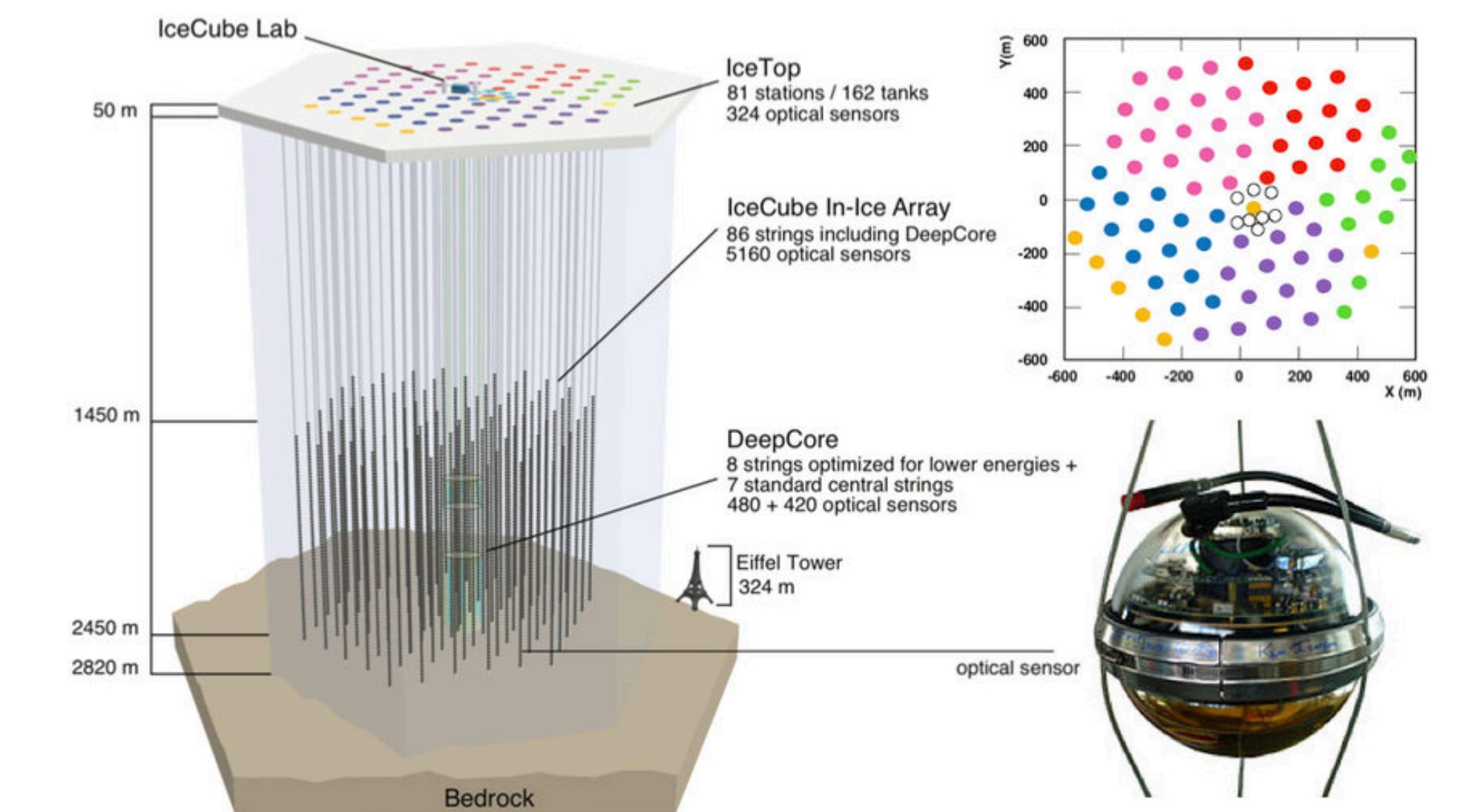
Position

Atom Type

Edge Features



Particle Detectors



Examples

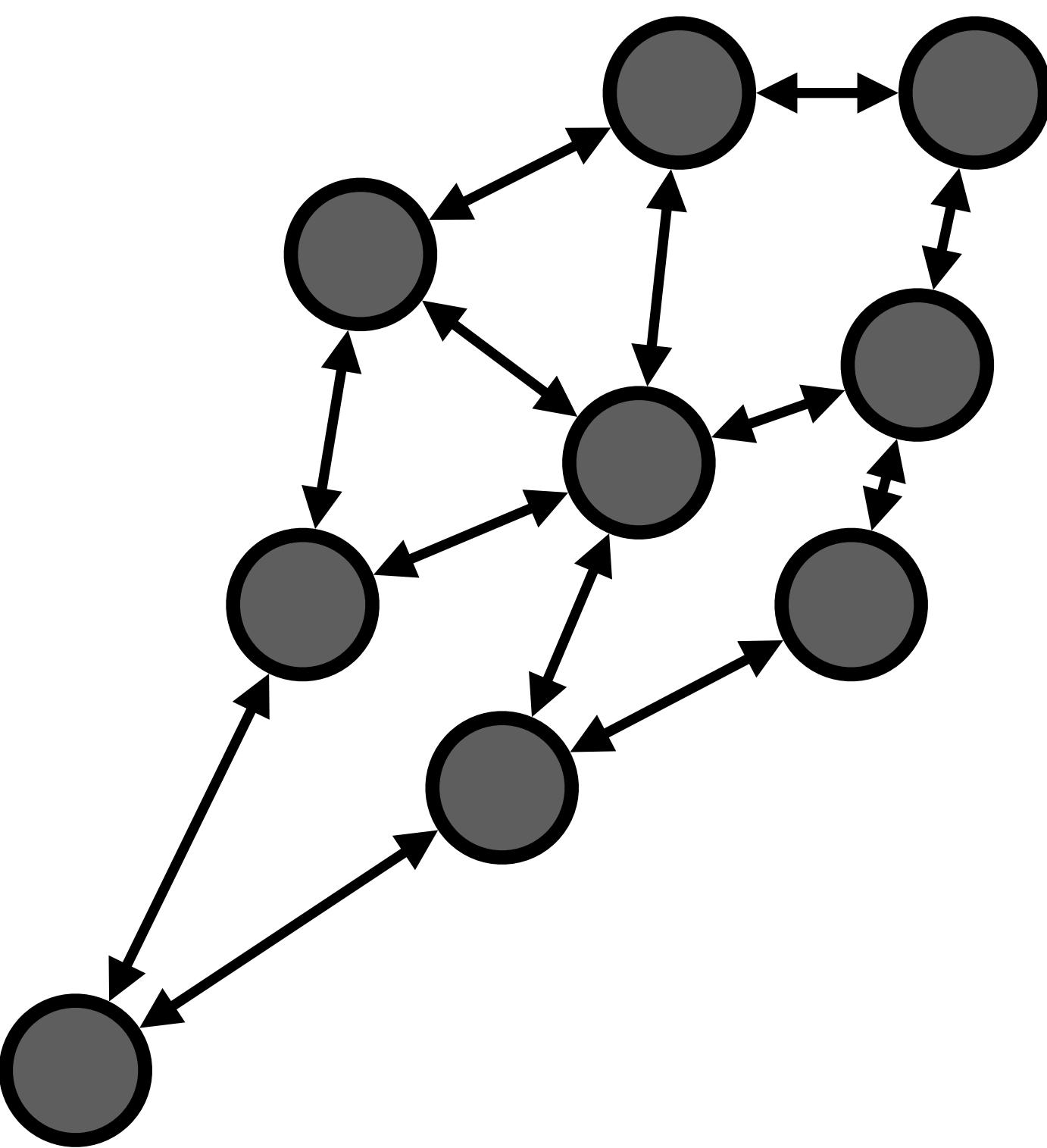
Graph structured data is everywhere!

Node Features

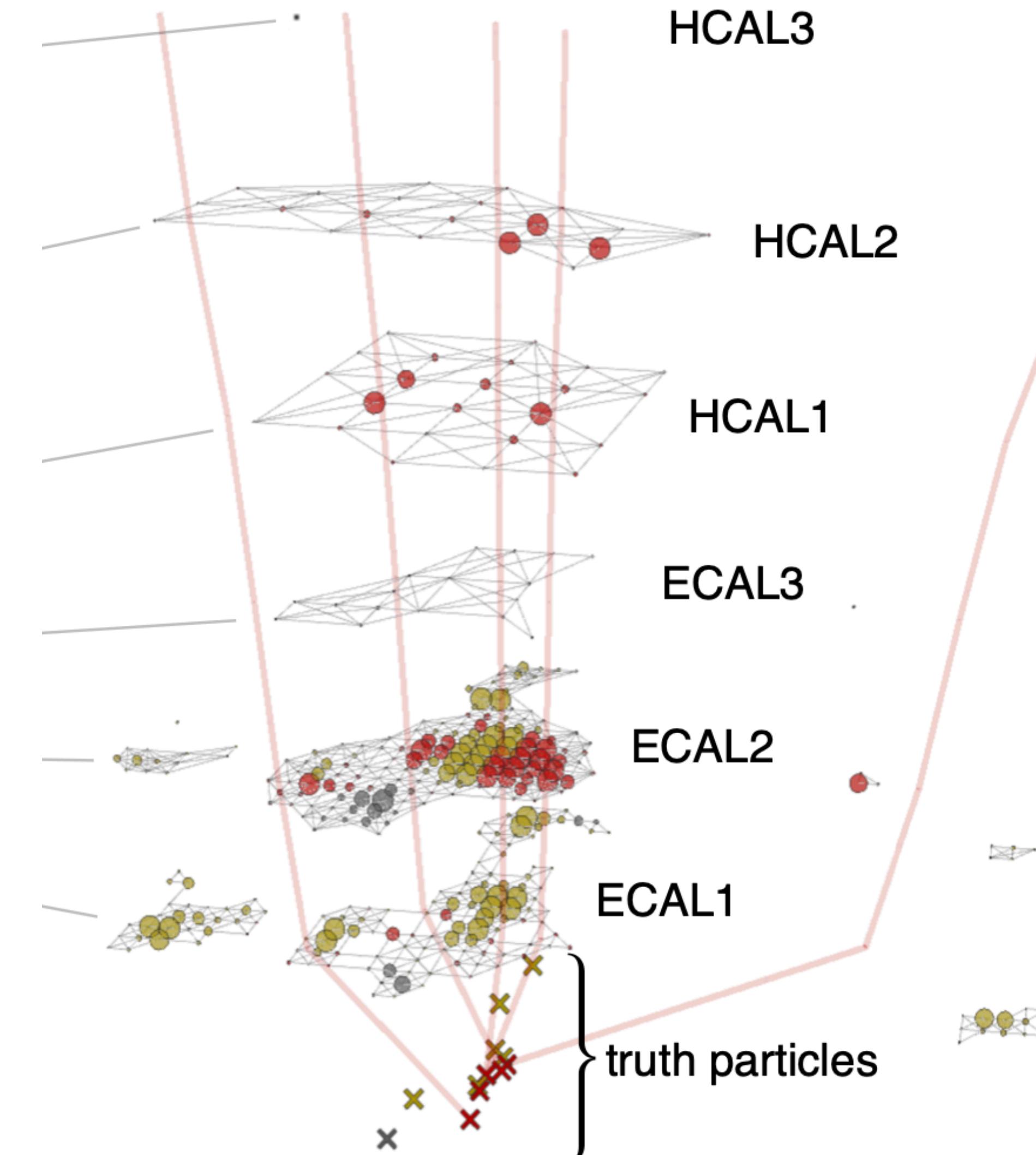
Position

Observed Signal

Edge Features



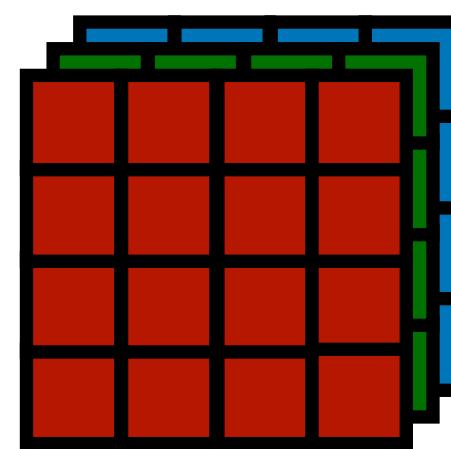
Particle Detectors



- ✖● photon
- ✖● neutral hadron
- ✖● charged hadron
- track, extrapolation

Examples

In fact, even images can be thought of as graphs



Examples

In fact, even images can be thought of as graphs

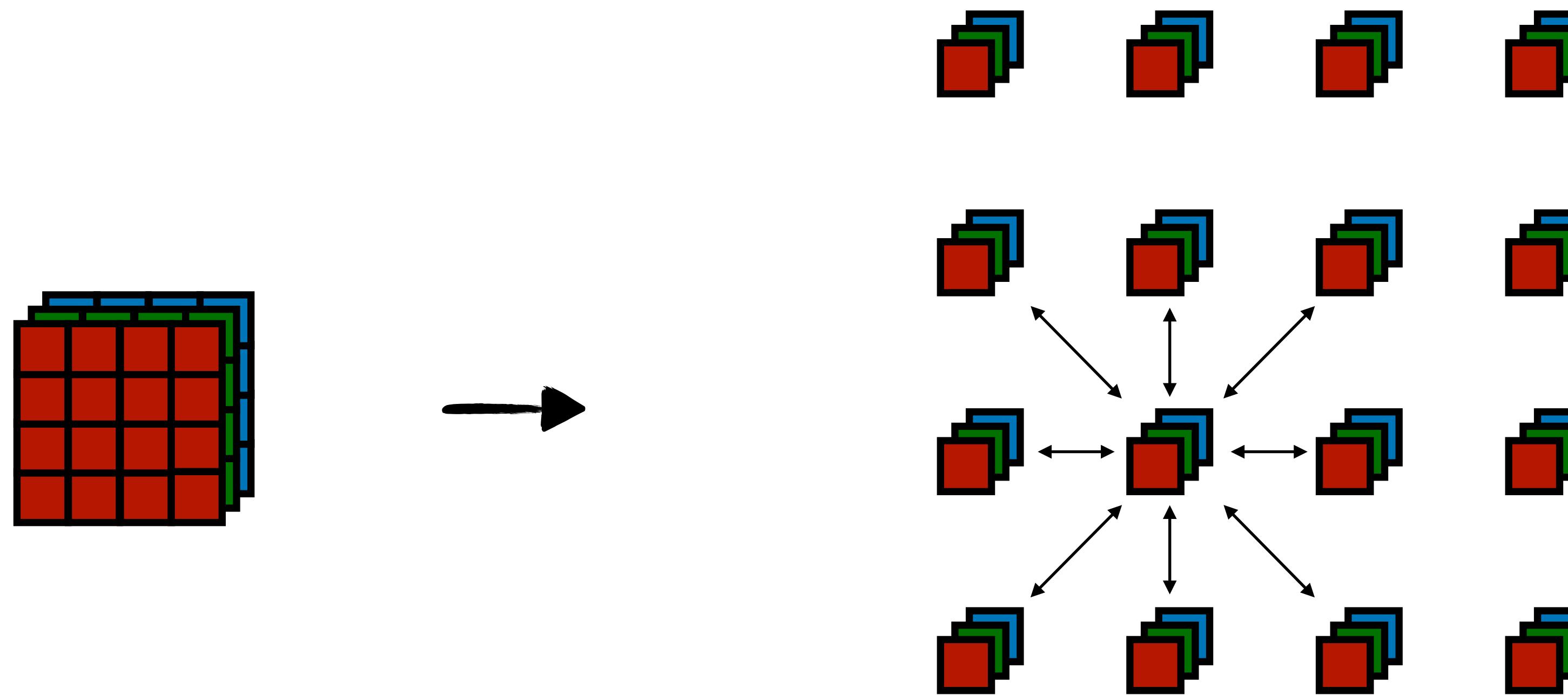


Each pixel
is its own object

Node Features: RGB intensities

Examples

In fact, even images can be thought of as graphs

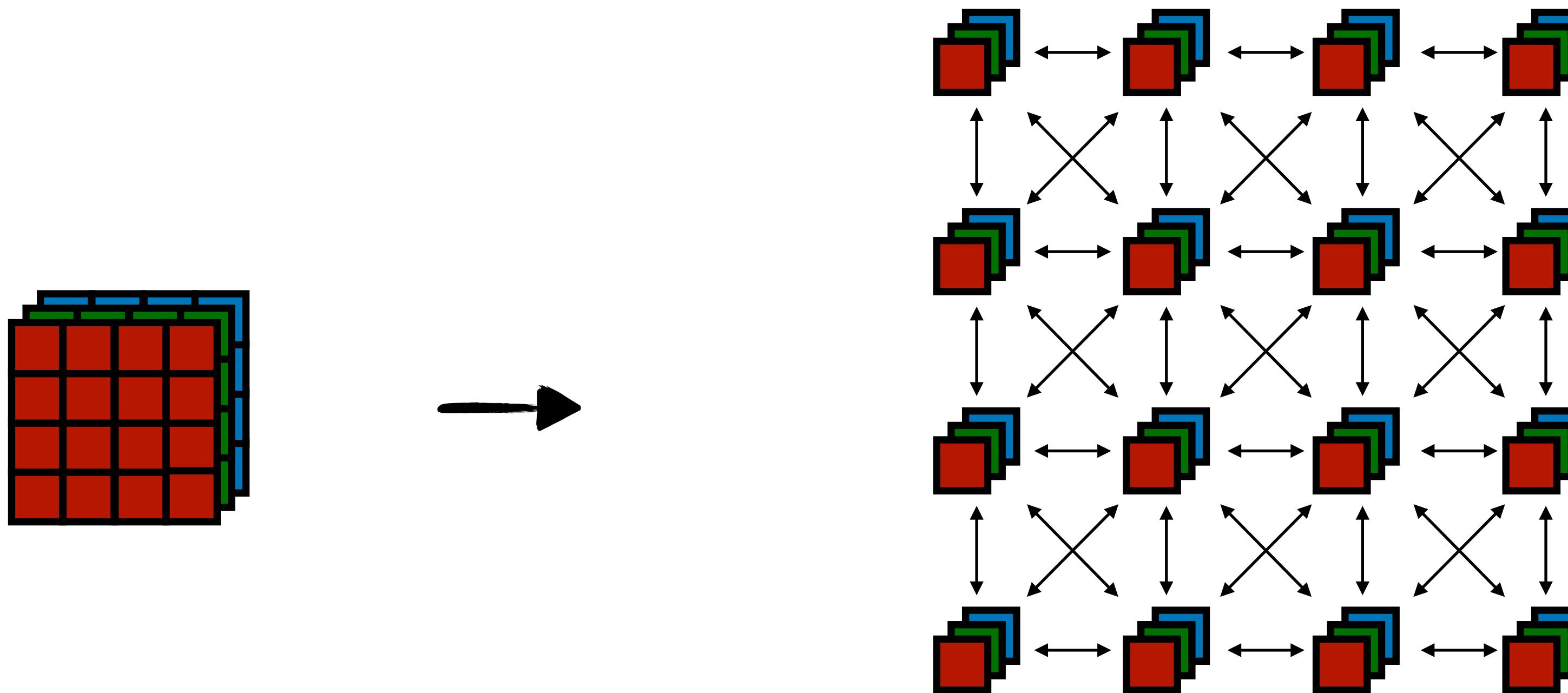


Images: Grids with very regular local connectivity

Edge Features: orientation (to the left / to the right, ...)

Examples

In fact, even images can be thought of as graphs



Learning Problems on Graphs

There are a variety of possible learning problems on graphs

Node-Level

Edge-Level

Graph-Level

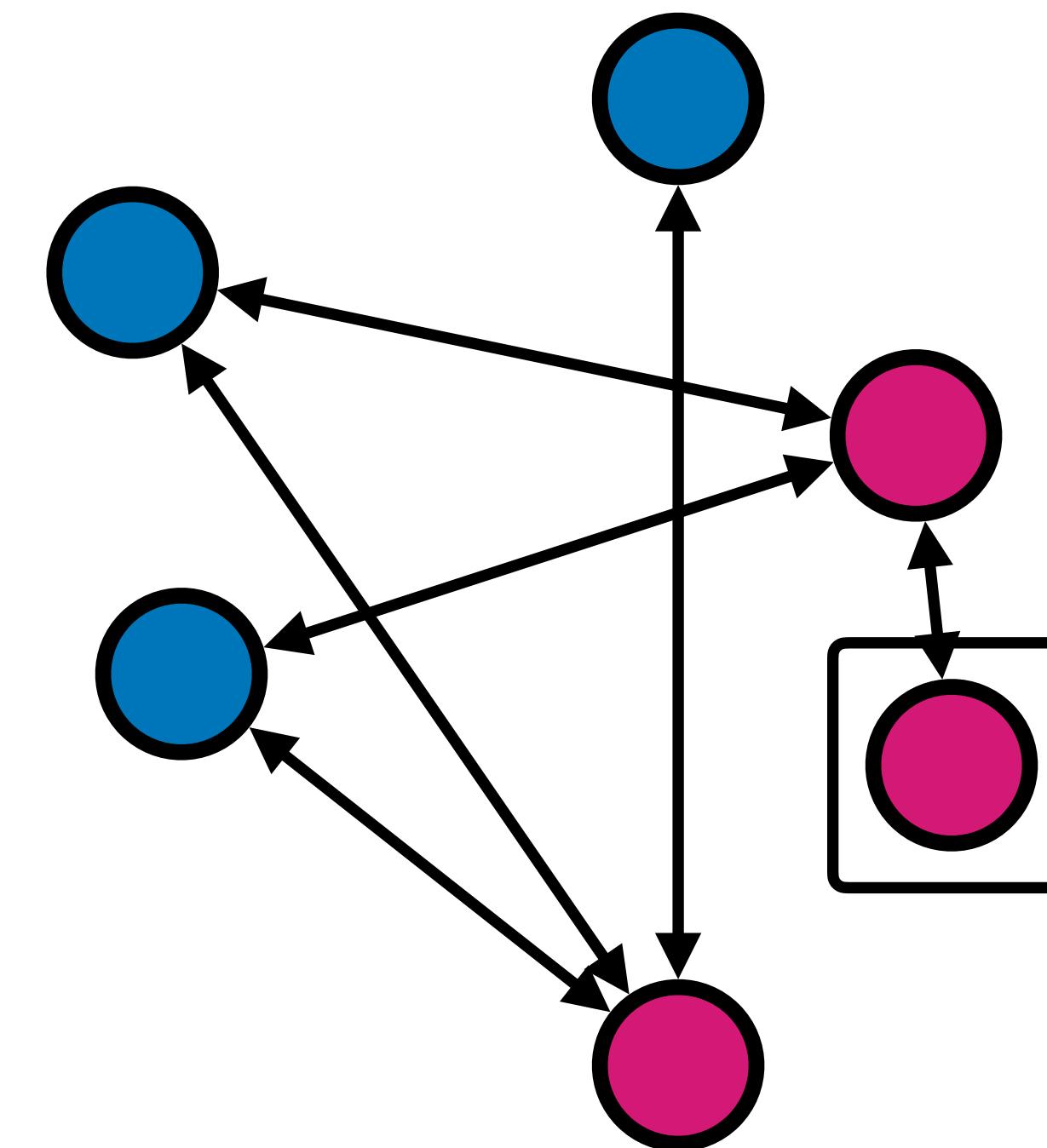
Learning Problems on Graphs

Node-level prediction:

Predict some latent property of a given node, given the graph.

Example:

What's a person's favorite sports club given a social network and observed node features



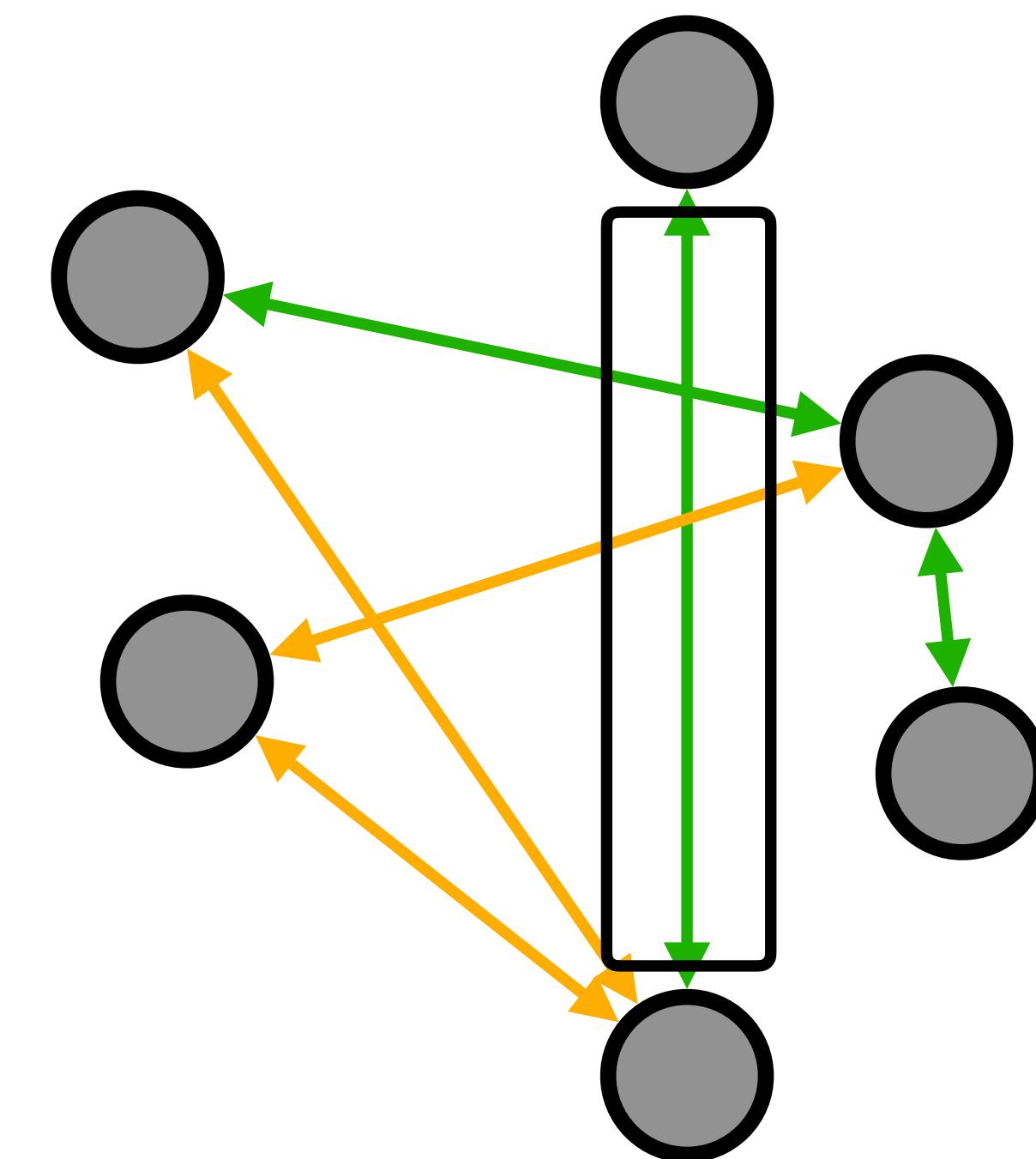
Learning Problems on Graphs

Edge-level prediction:

Predict some latent property
of a given edge, given the graph.

Example:

Given a social network,
which edge corresponds to a
first-degree family relationship?



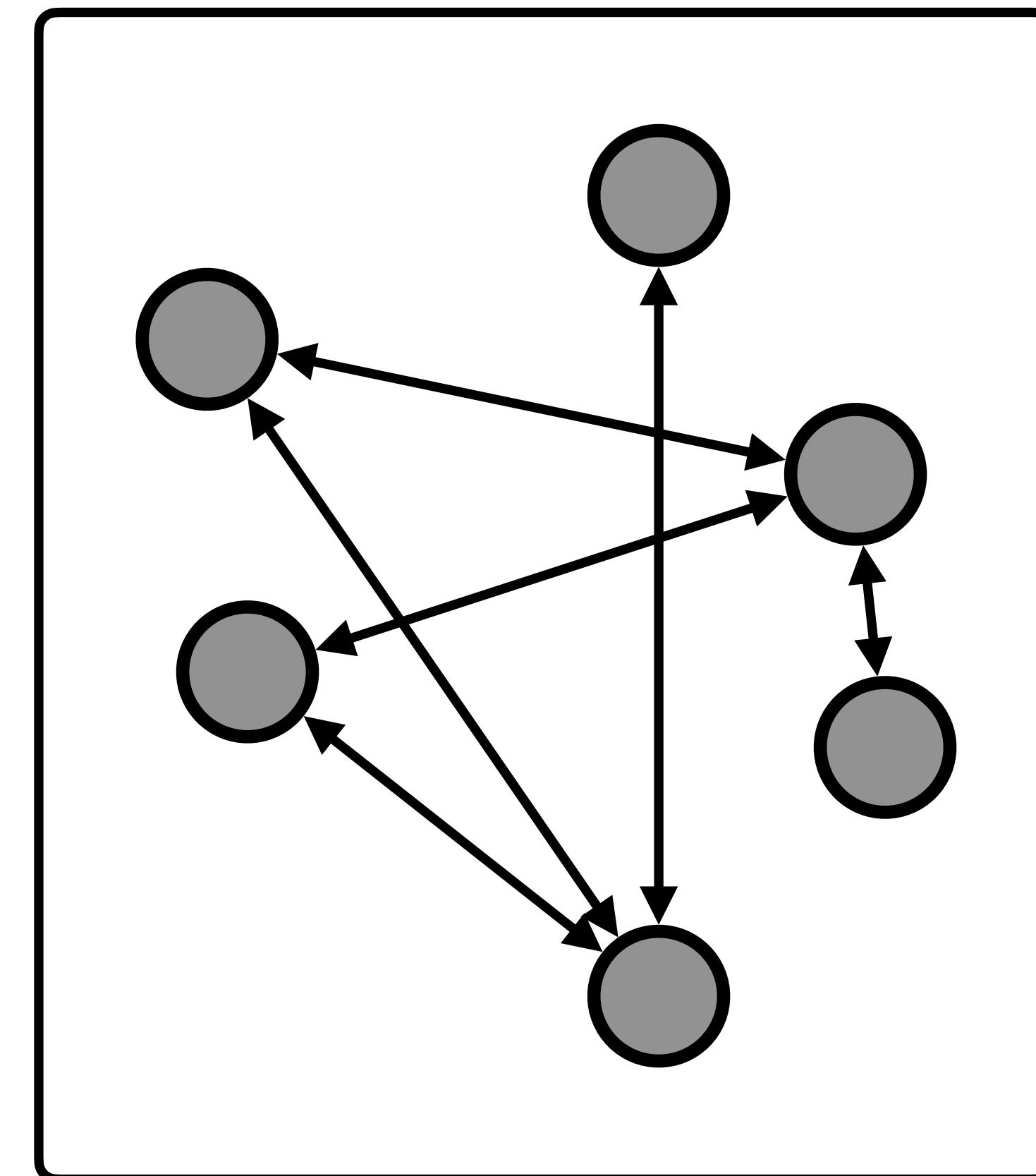
Learning Problems on Graphs

Graph-level prediction:

Predict some latent property
of a given edge, given the graph.

Example:

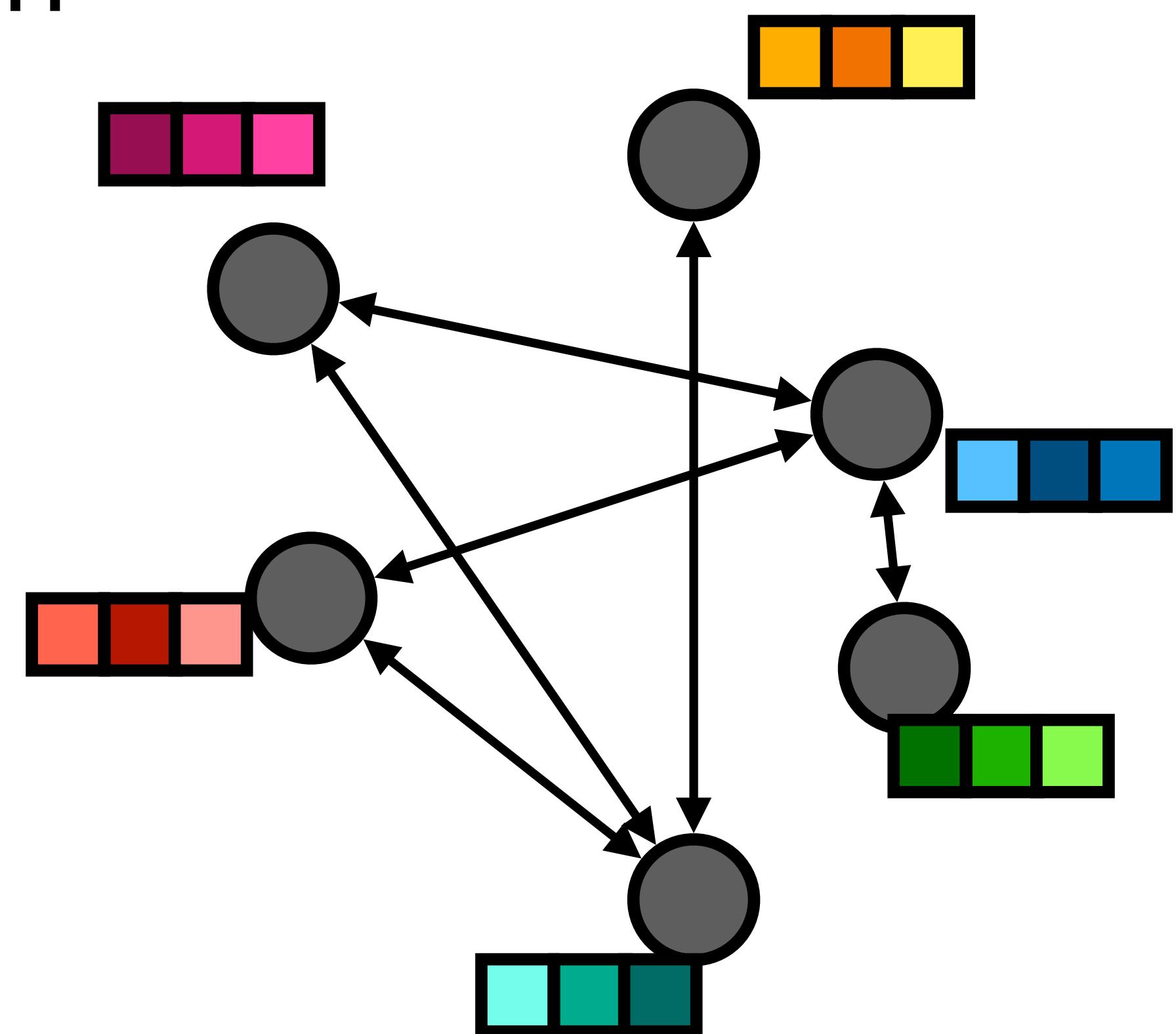
Given a social network,
which edge corresponds to a
first-degree family relationship?



Data Representation of Graphs

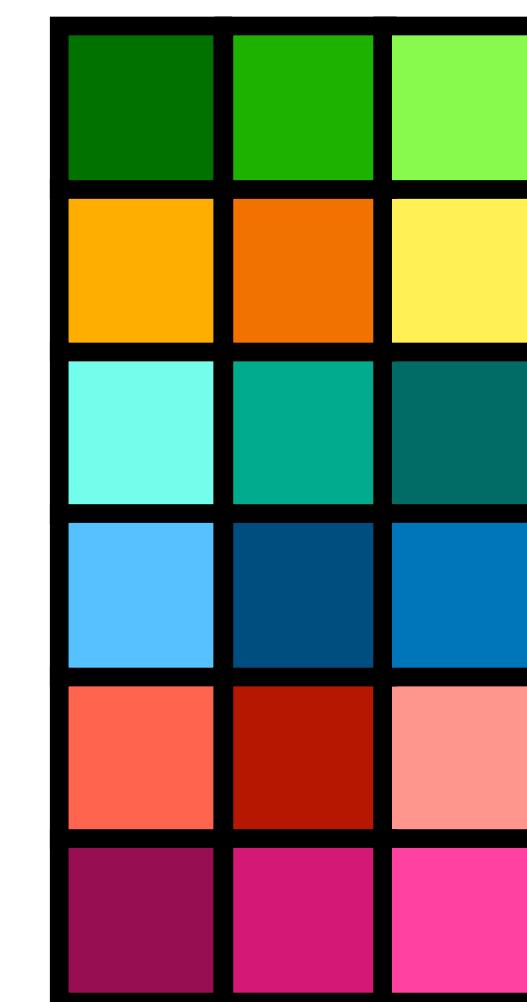
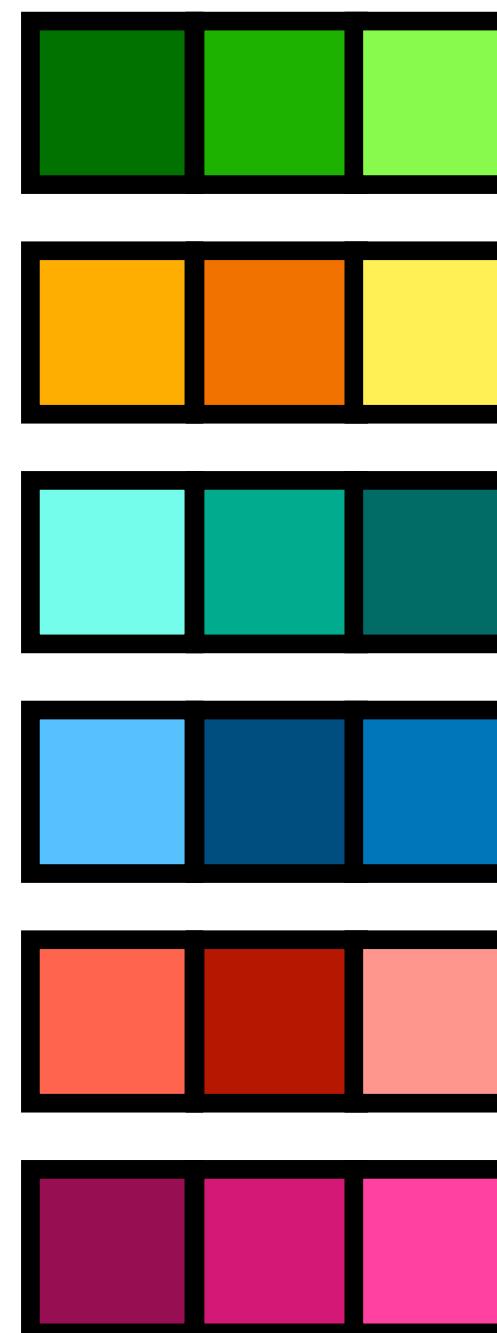
We can compactly represent the graph with matrices, but for this we need to impose an node order

Important this order is arbitrary

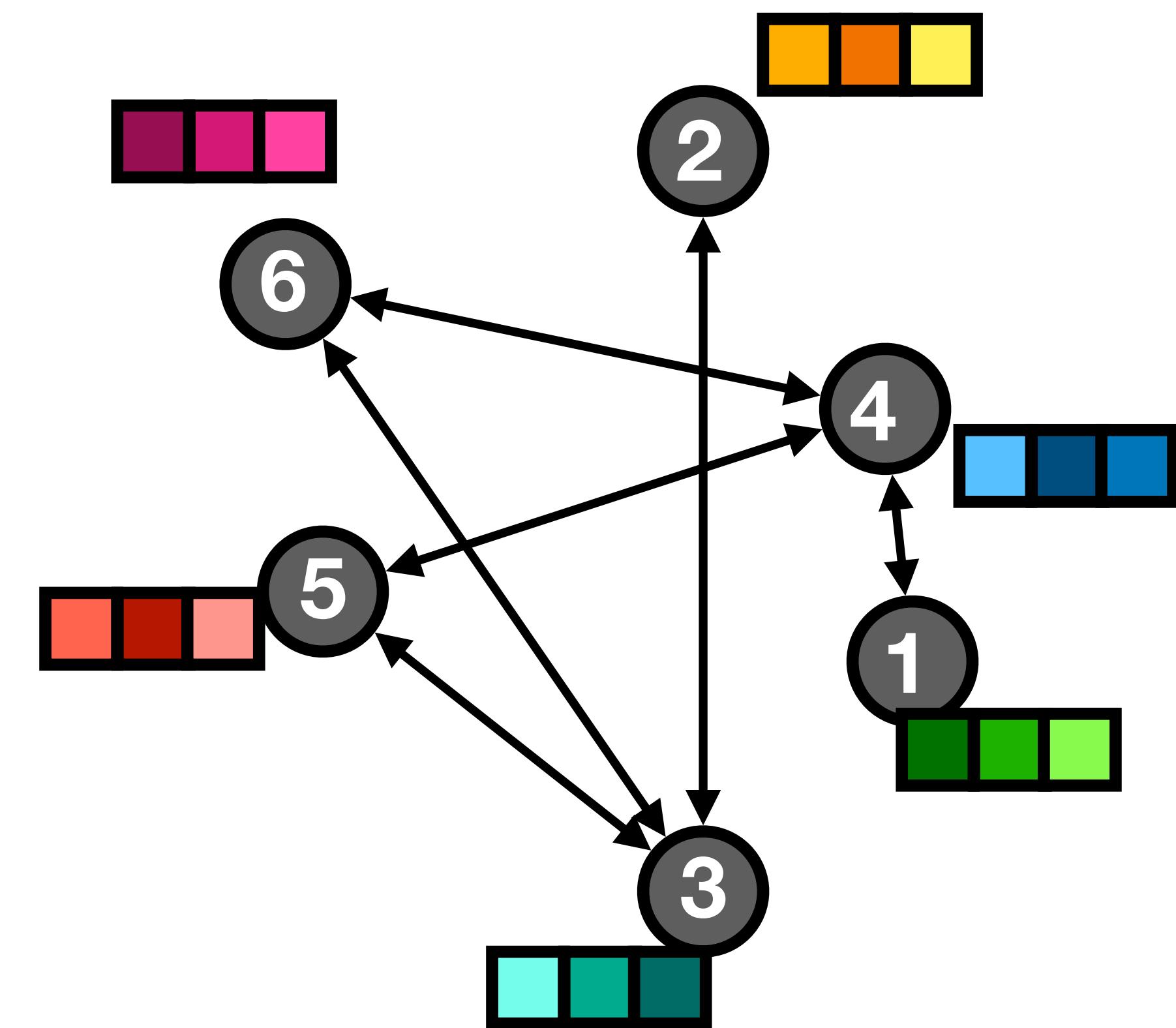


Data Representation of Graphs

With an ordering we can represent the data as a matrix

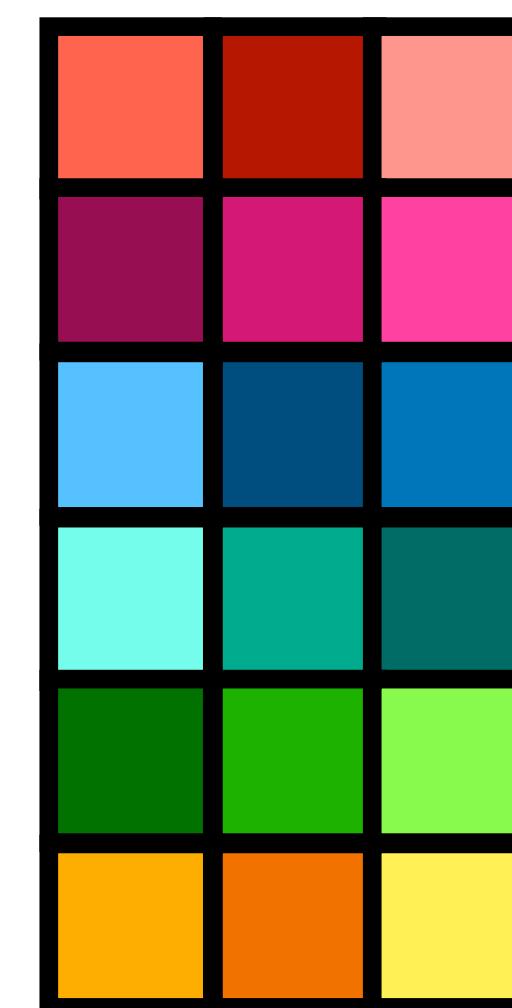
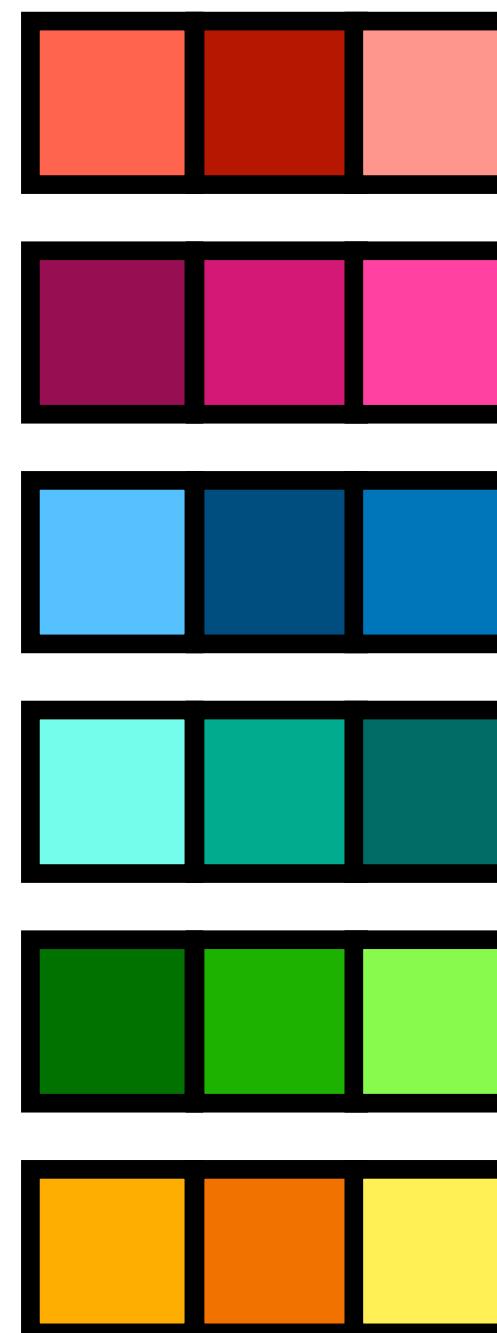


$$X \in \mathbb{R}^{n \times f}$$

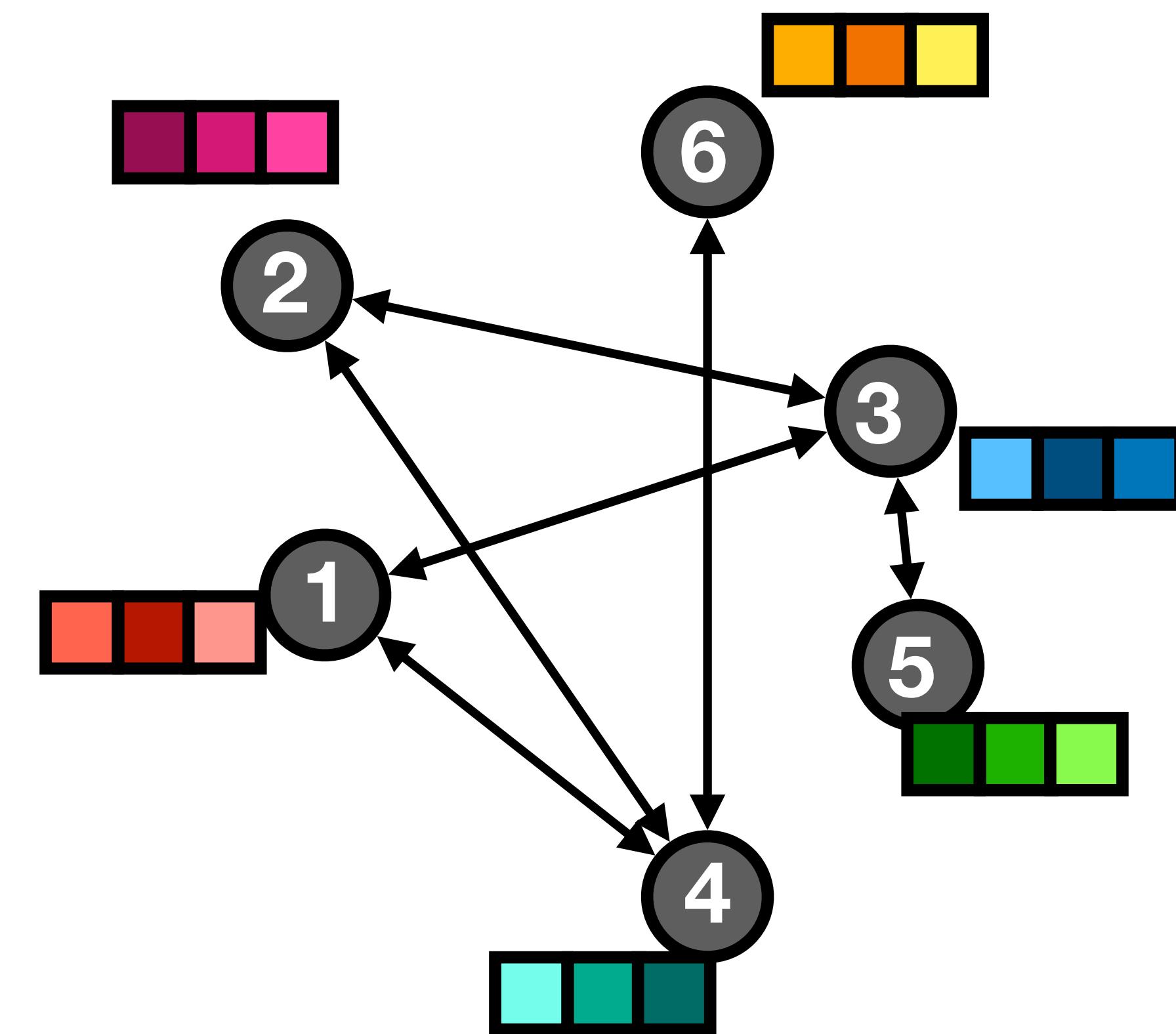


Data Representation of Graphs

For a different ordering...



$$X \in \mathbb{R}^{n \times f}$$

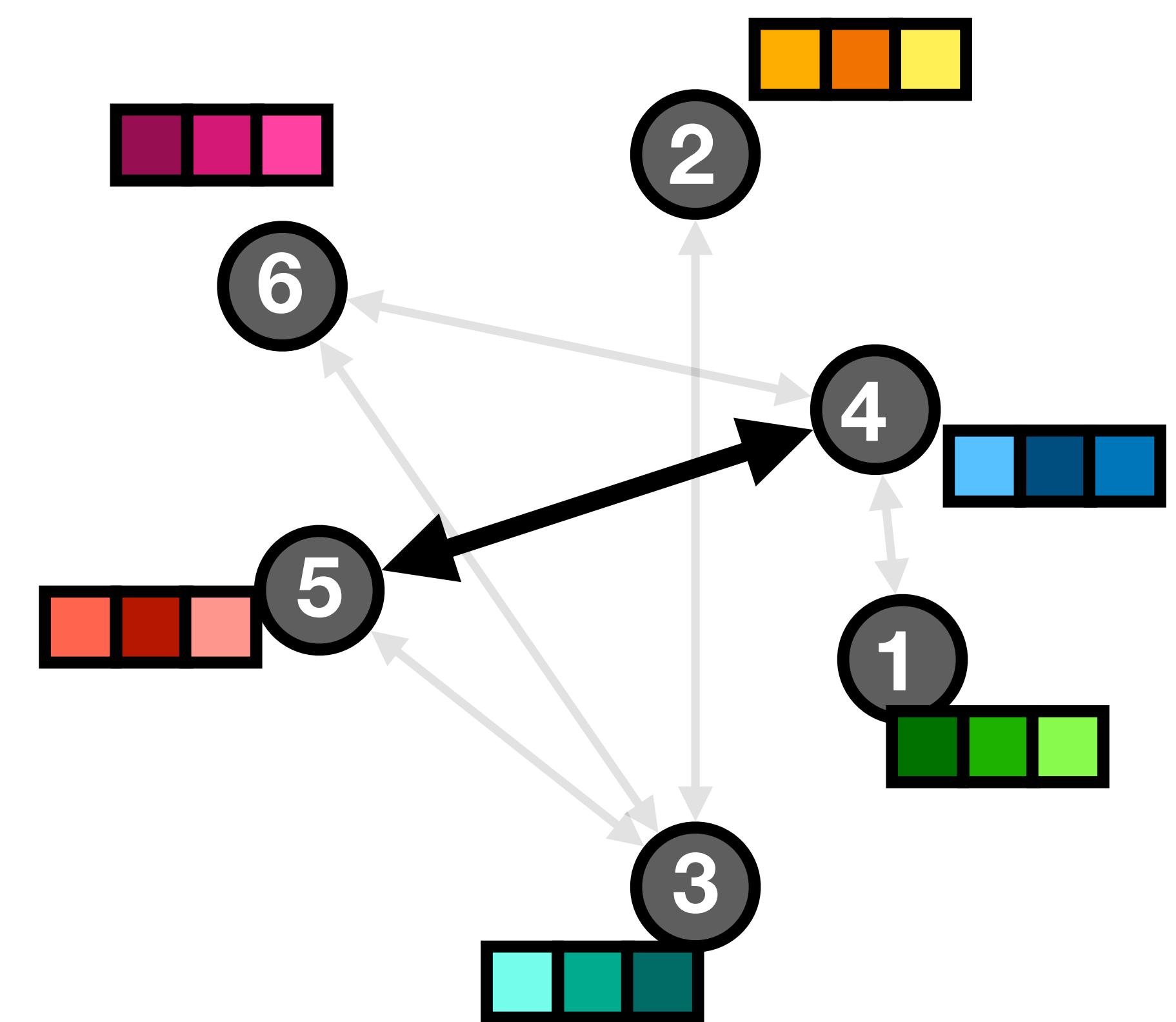


Data Representation of Graphs

The edges can be represented as an

Adjacency Matrix $A \in \mathbb{R}^{n \times n}$

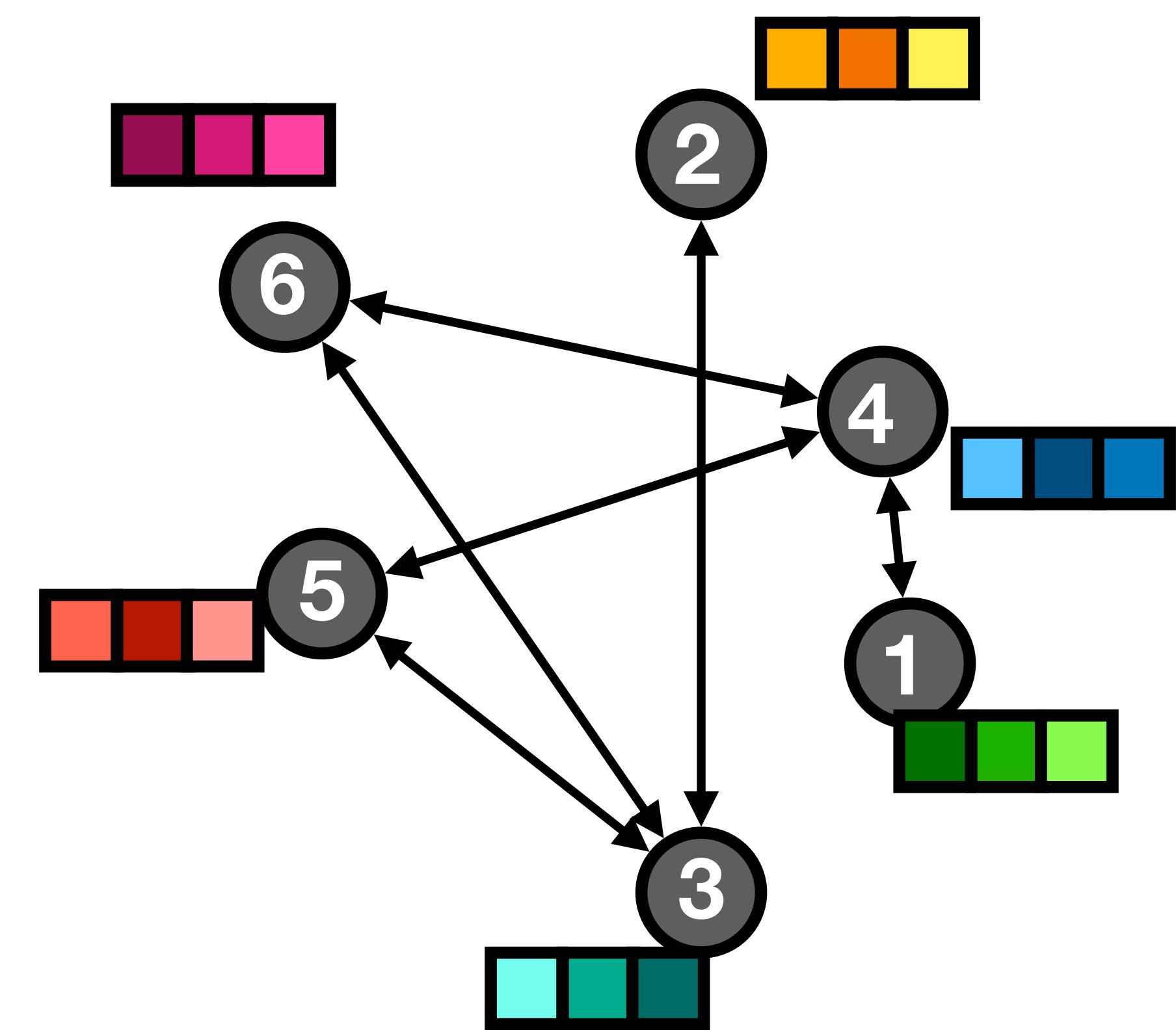
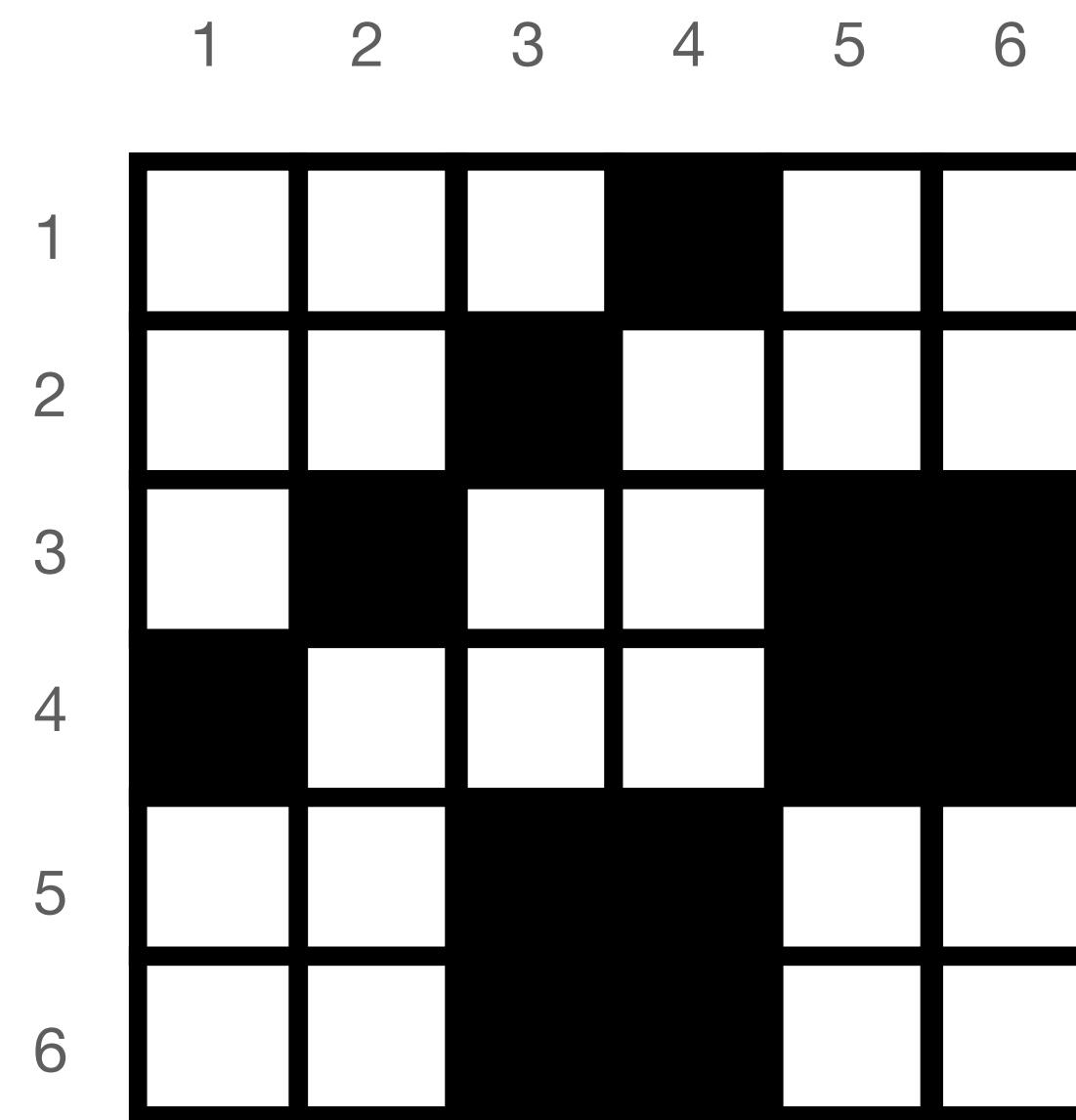
| | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 1 | | | | | | |
| 2 | | | | | | |
| 3 | | | | | | |
| 4 | | | | | | |
| 5 | | | | | | |
| 6 | | | | | | |



Data Representation of Graphs

The edges can be represented as an

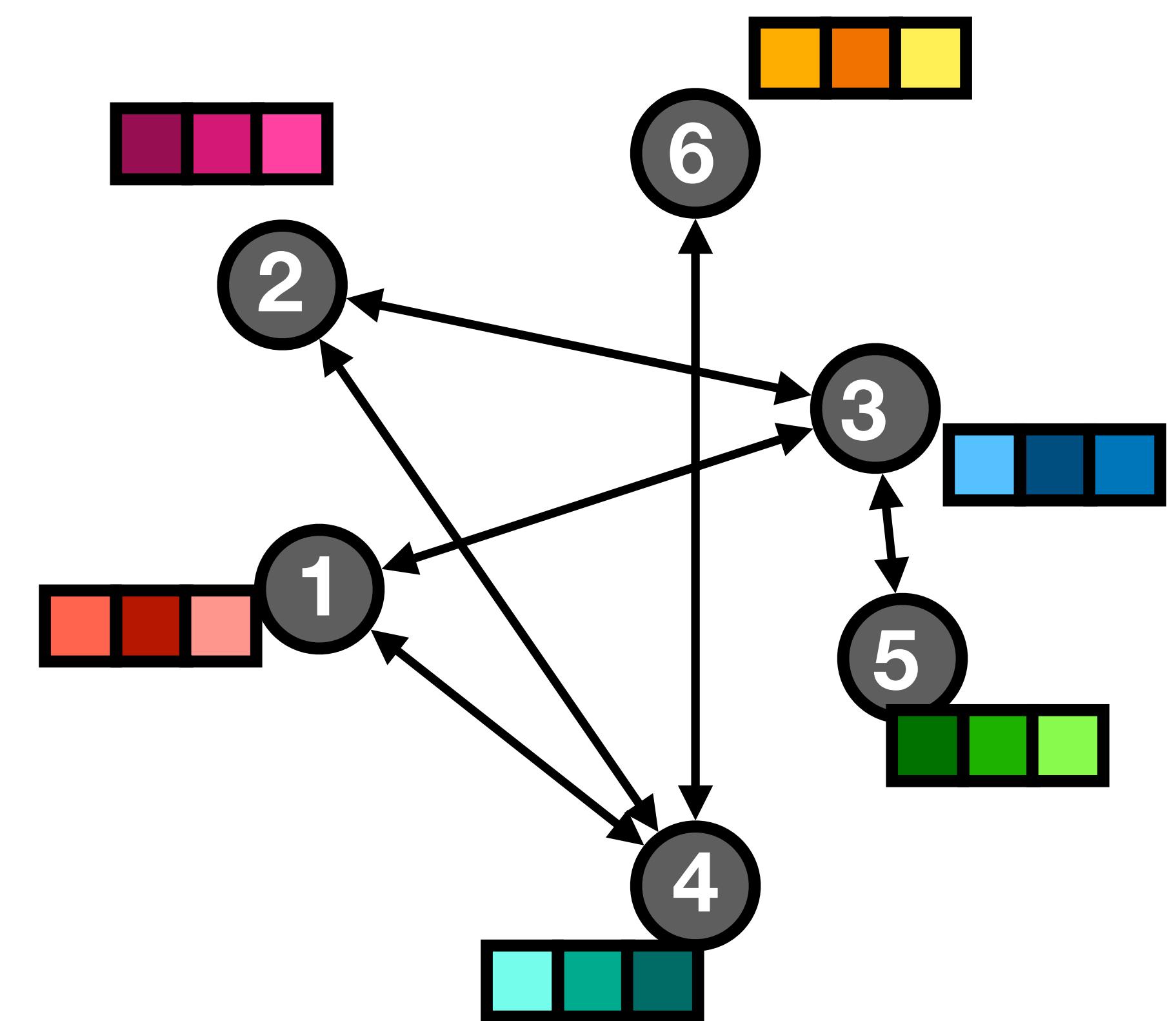
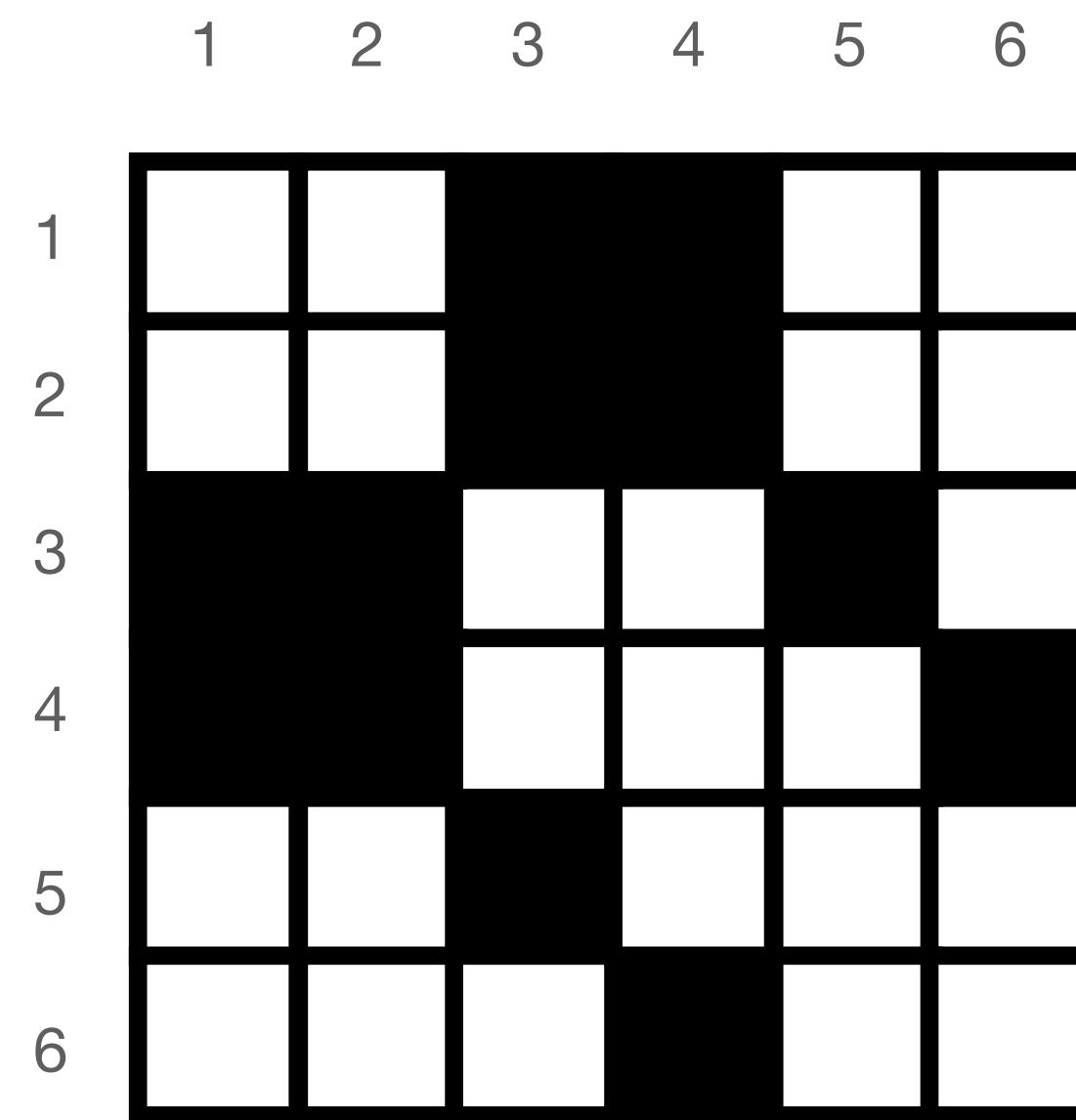
Adjacency Matrix $A \in \mathbb{R}^{n \times n}$



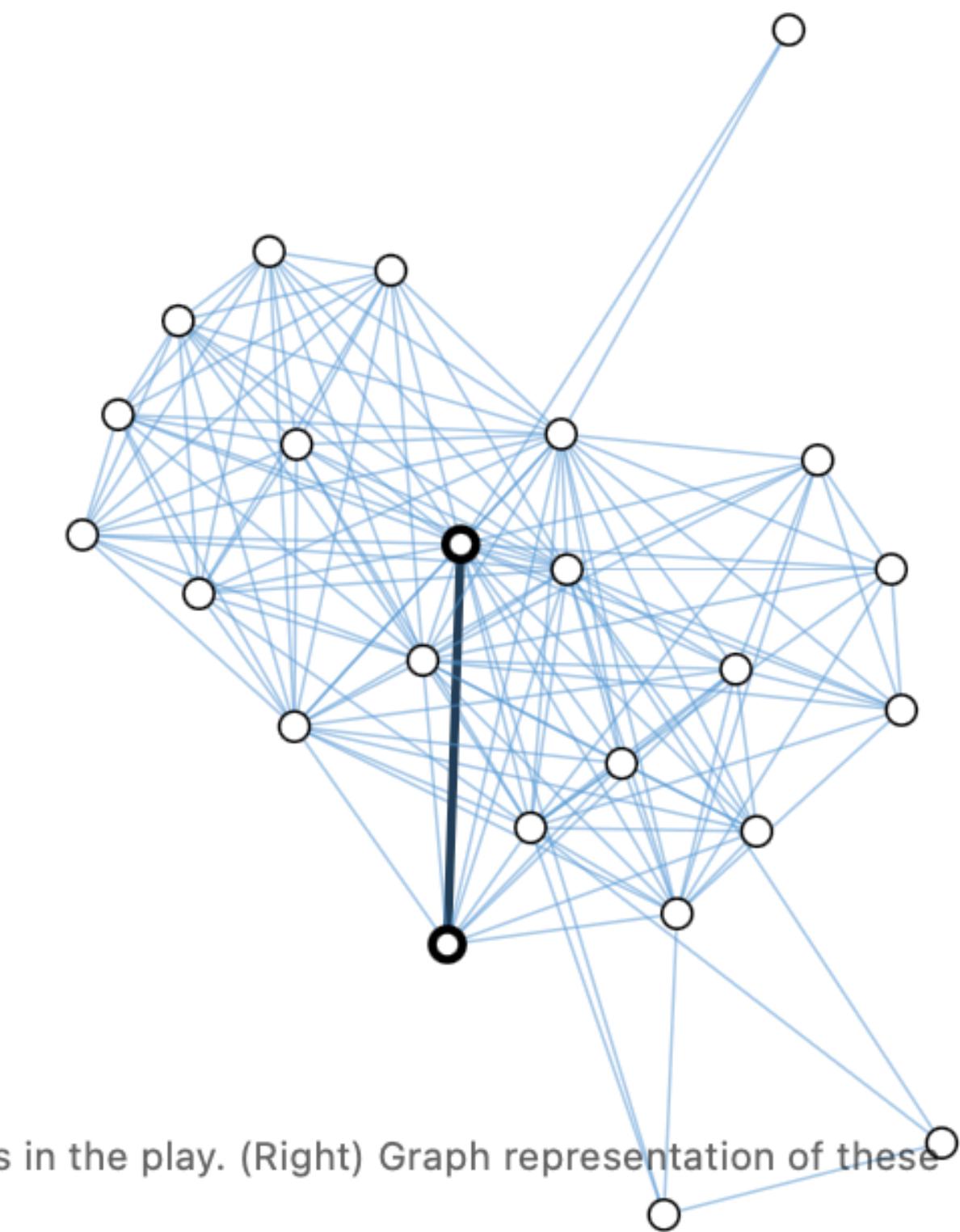
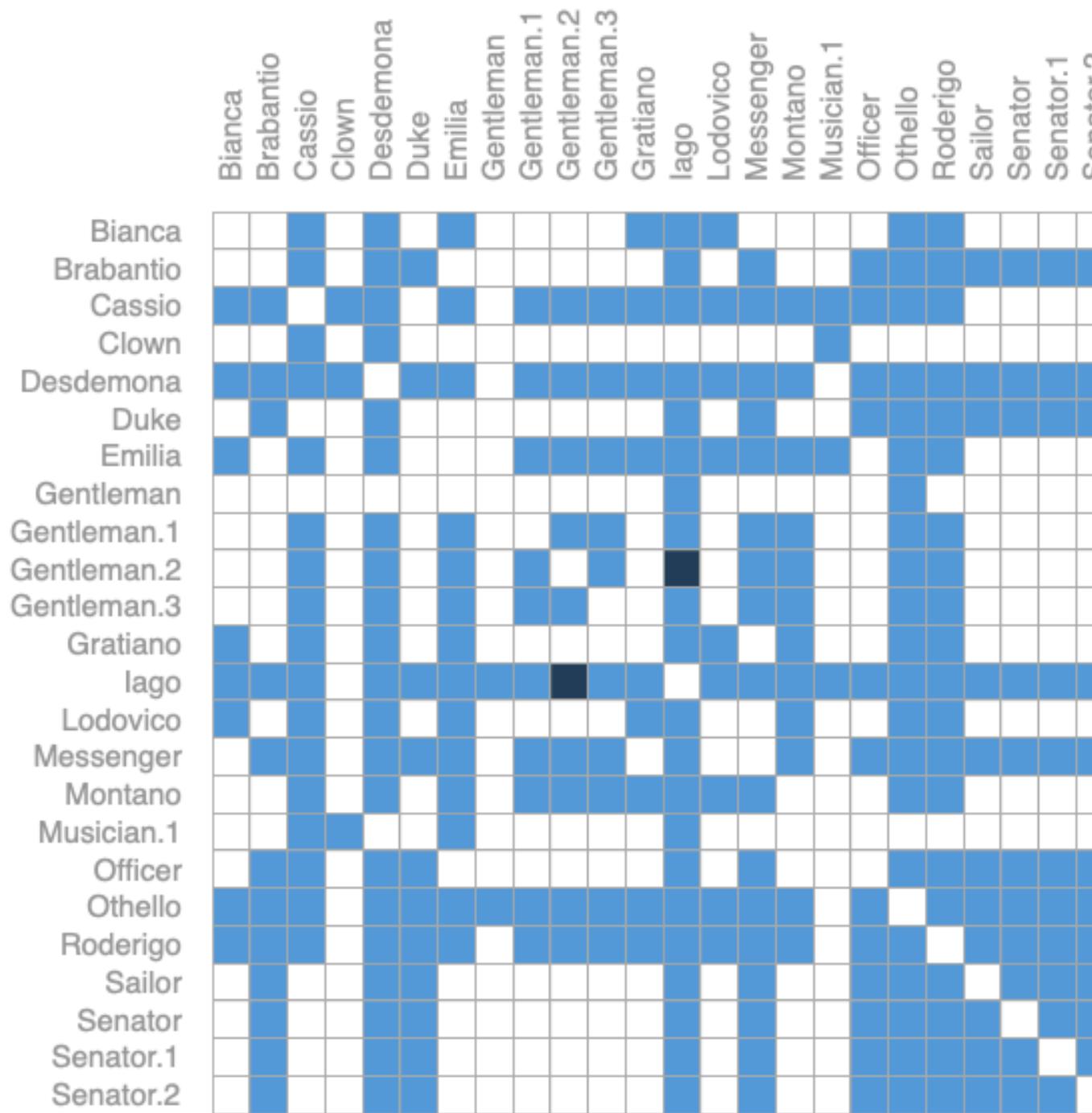
Data Representation of Graphs

Ordered differently...

Adjacency Matrix $A \in \mathbb{R}^{n \times n}$

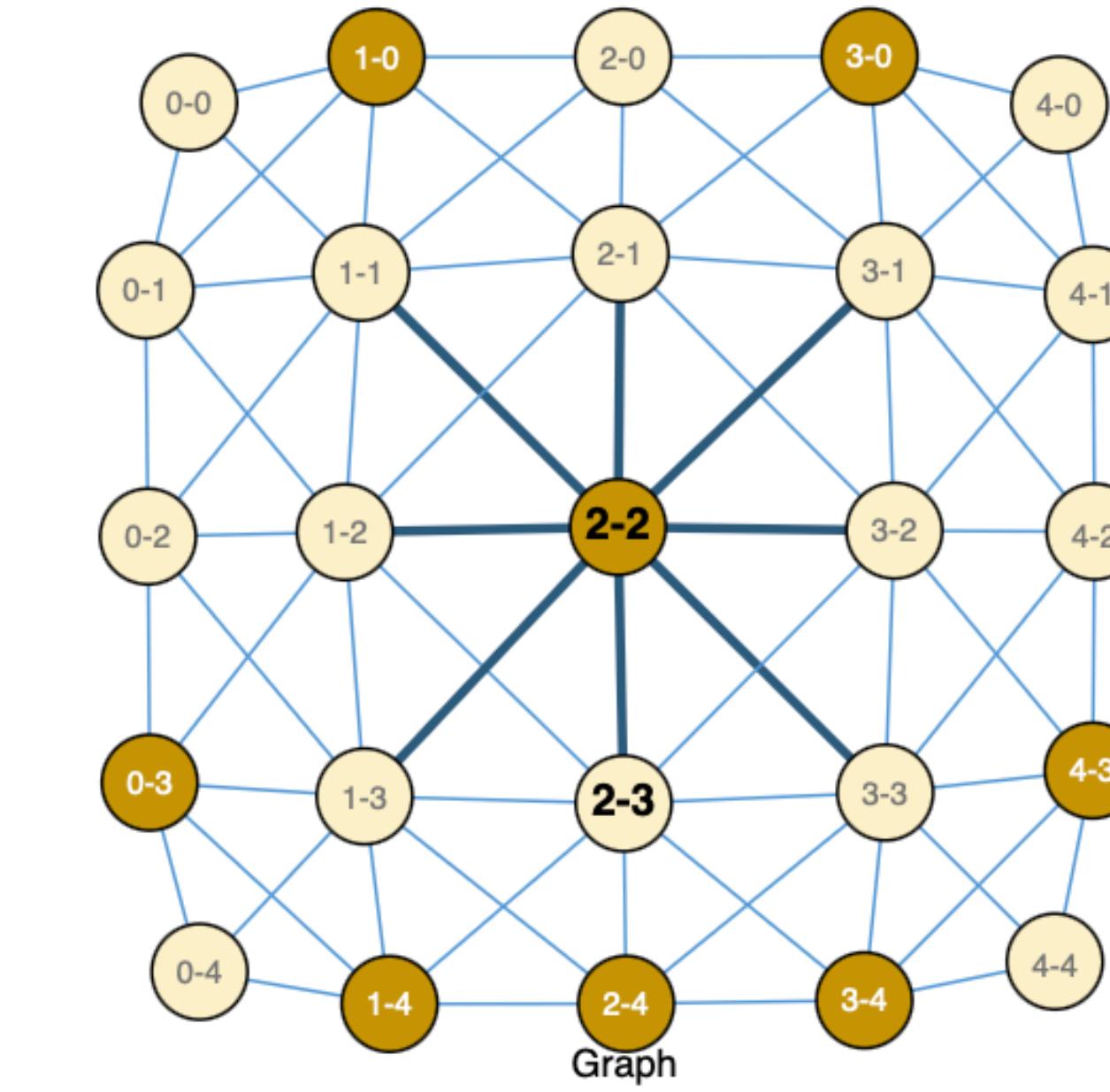
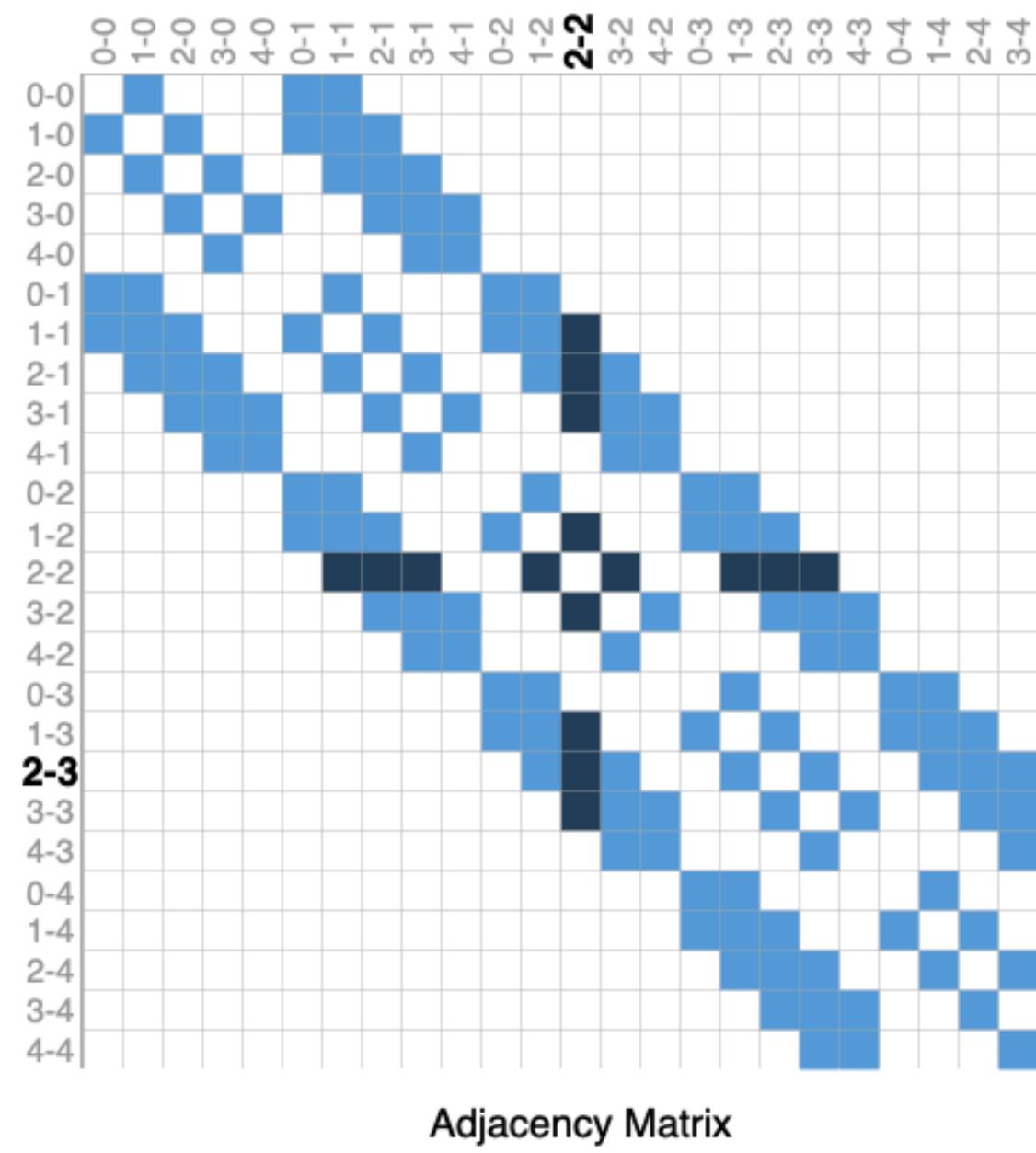
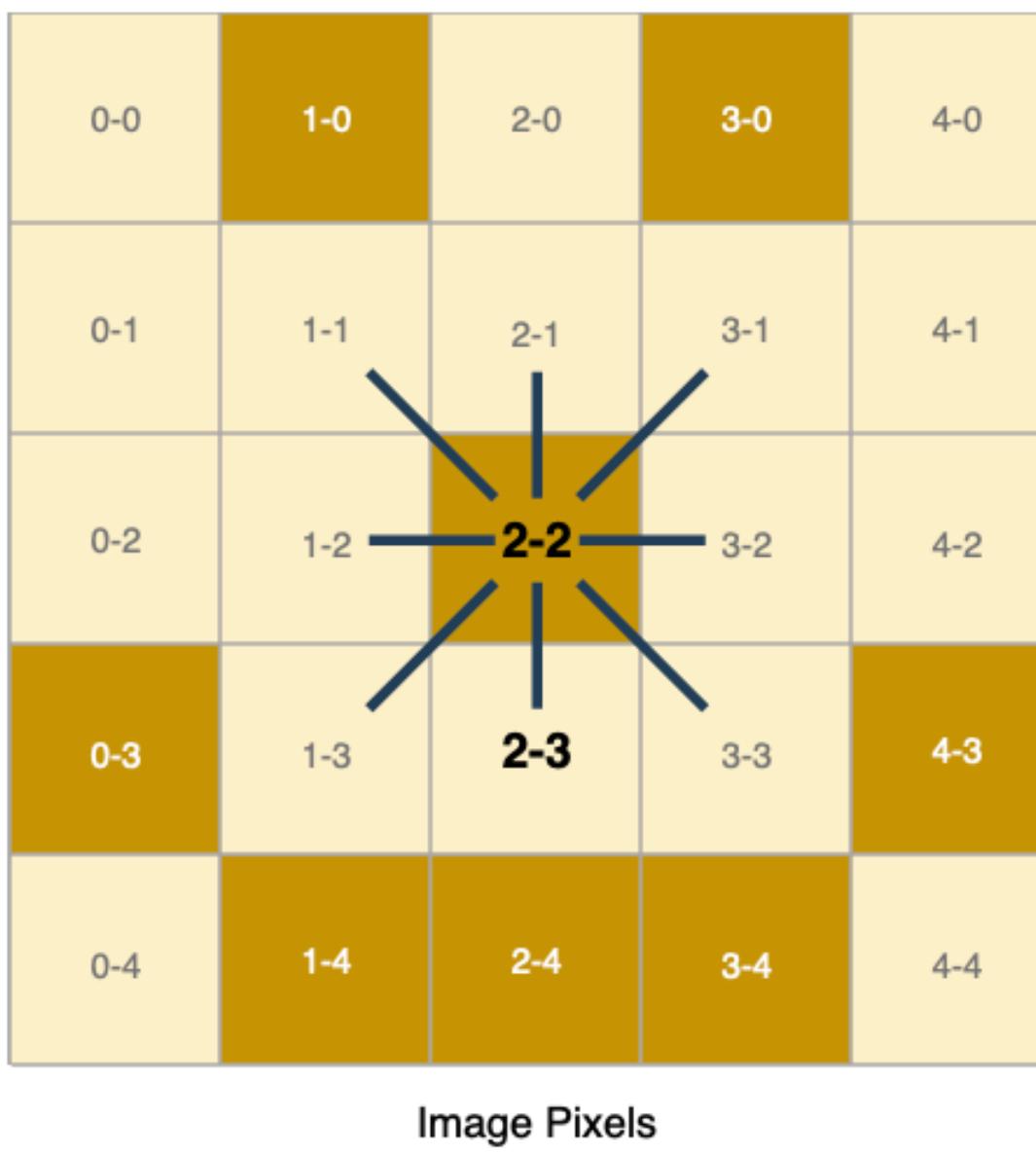


Examples



(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.

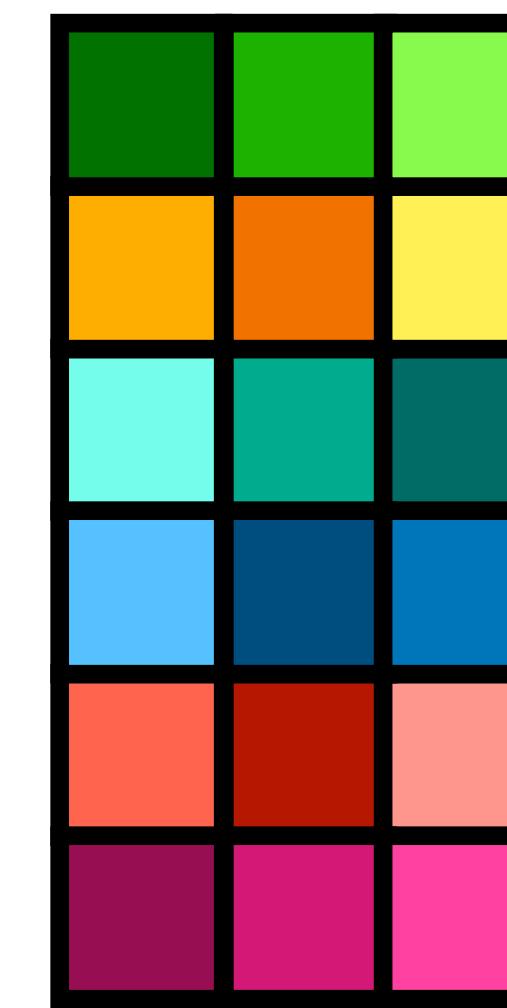
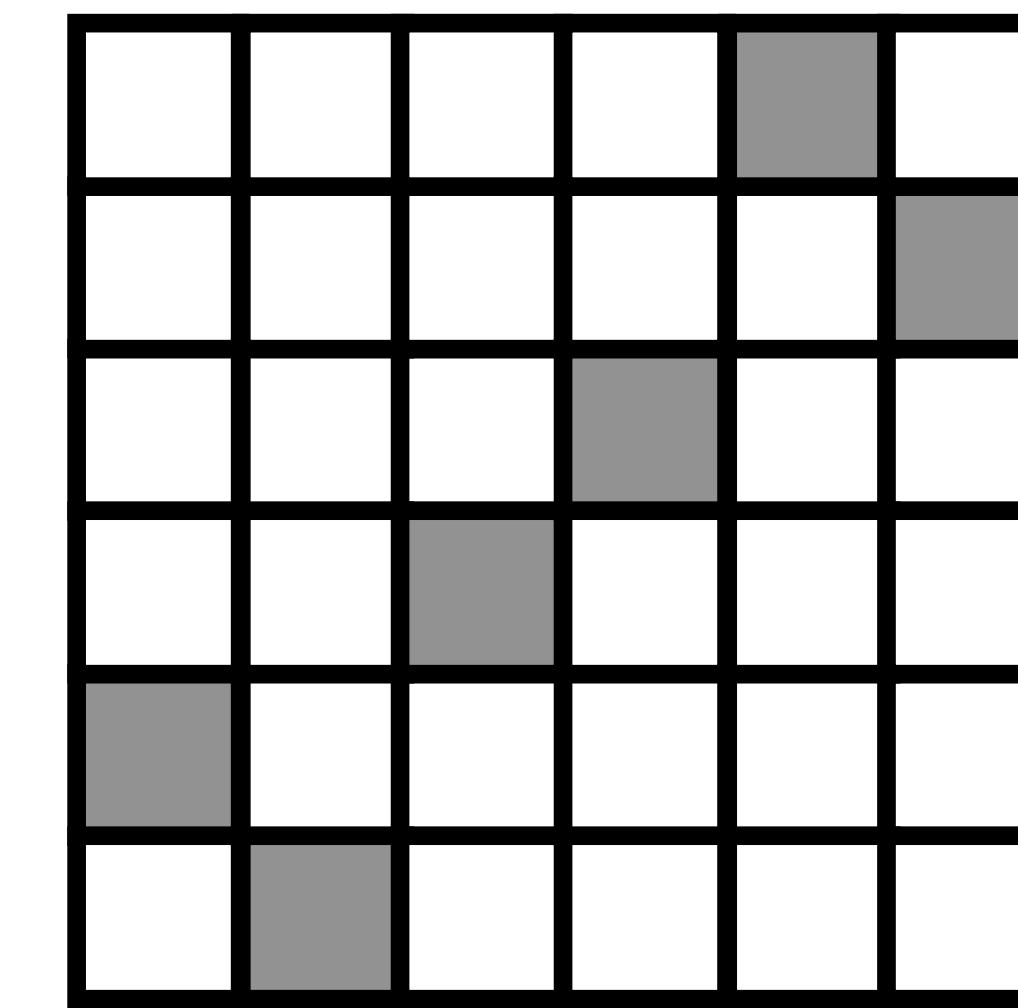
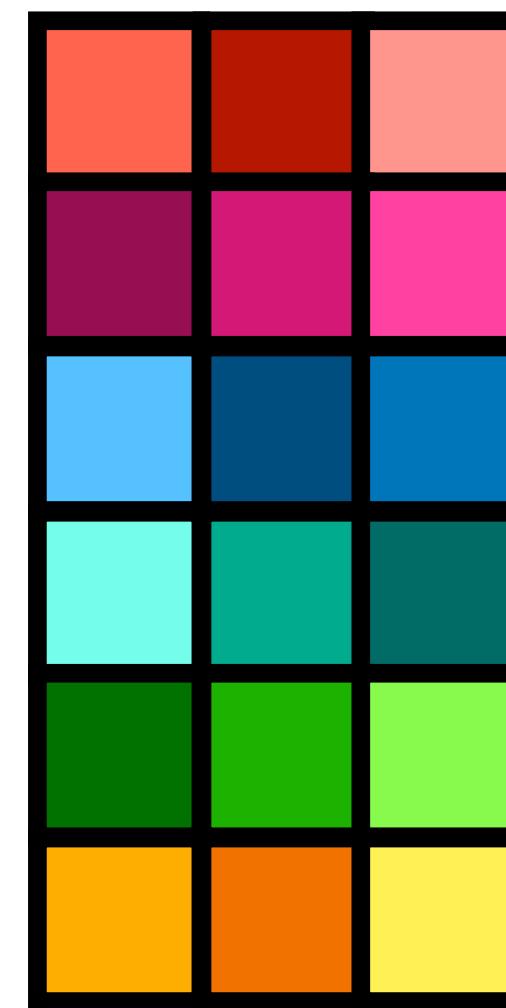
Examples



Click on an image pixel to toggle its value, and see how the graph representation changes.

Permutation Matrices

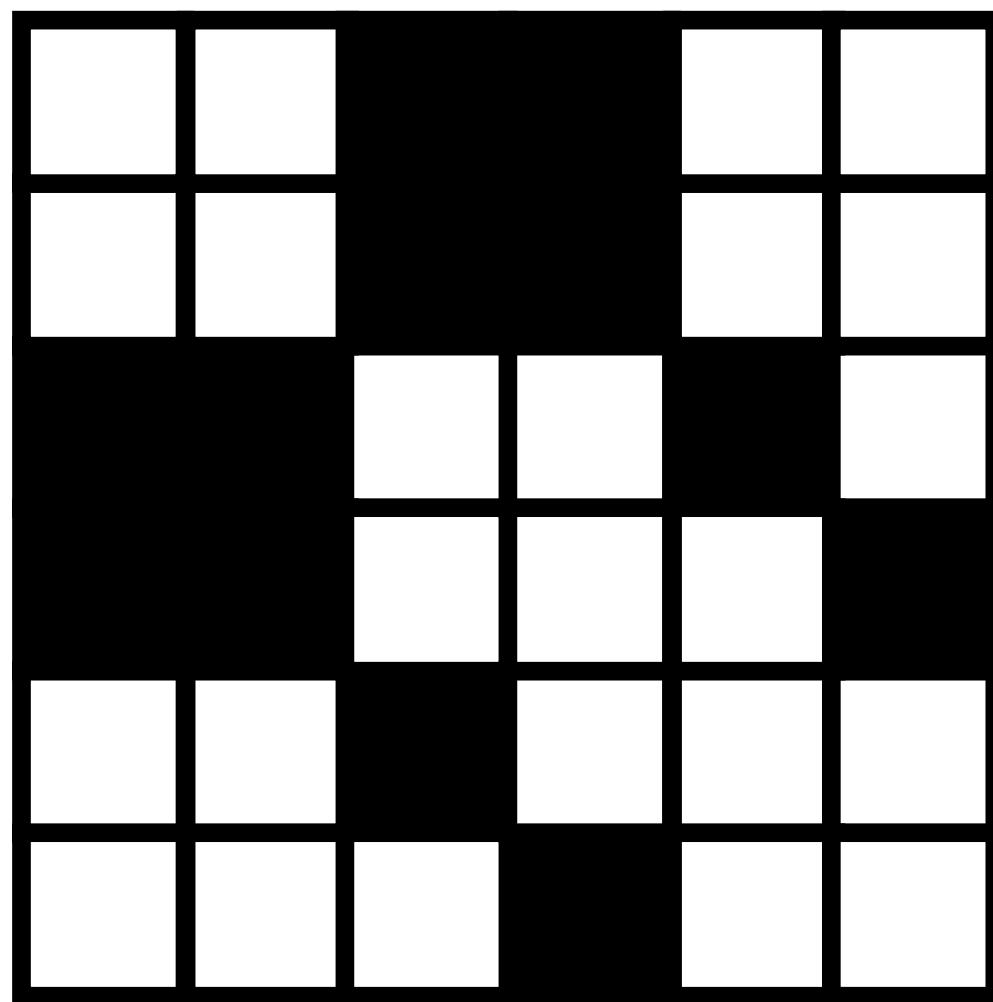
We can represent permutations using a **permutation matrix**



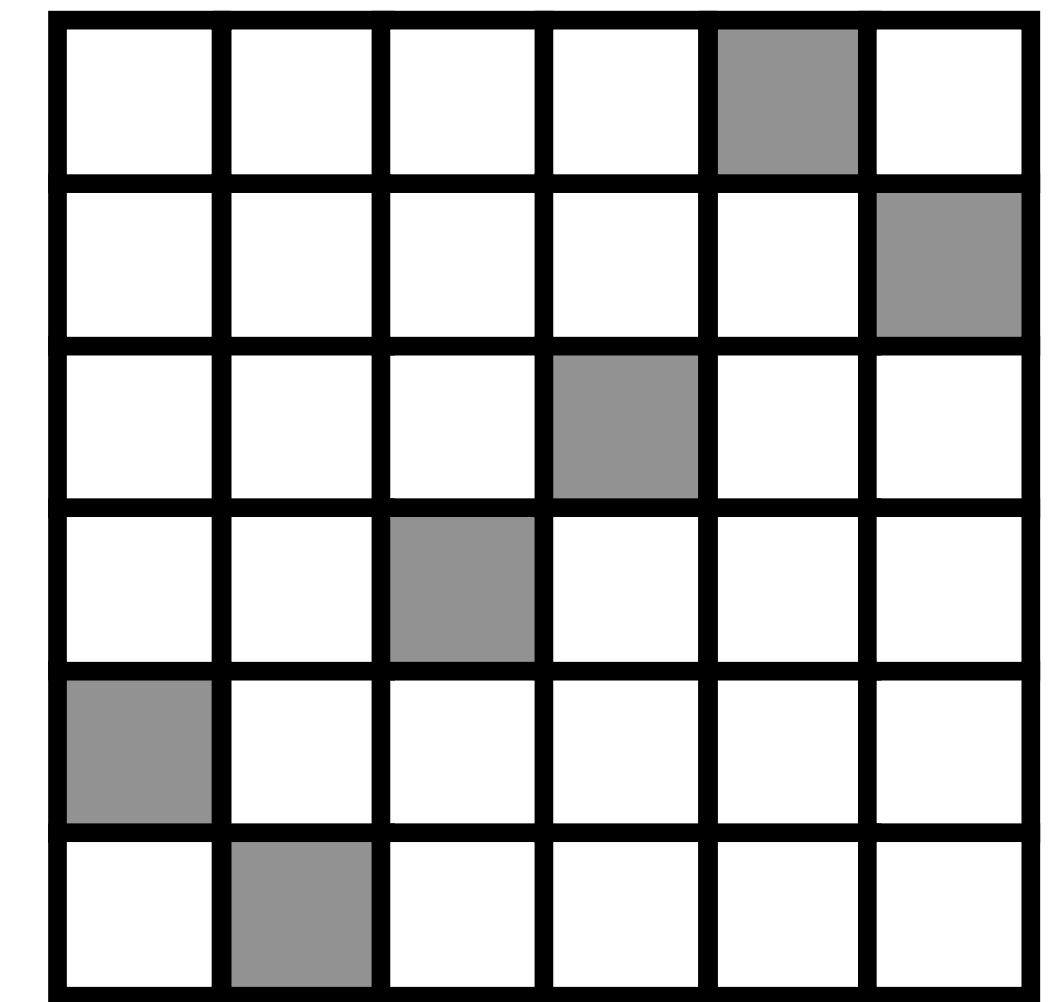
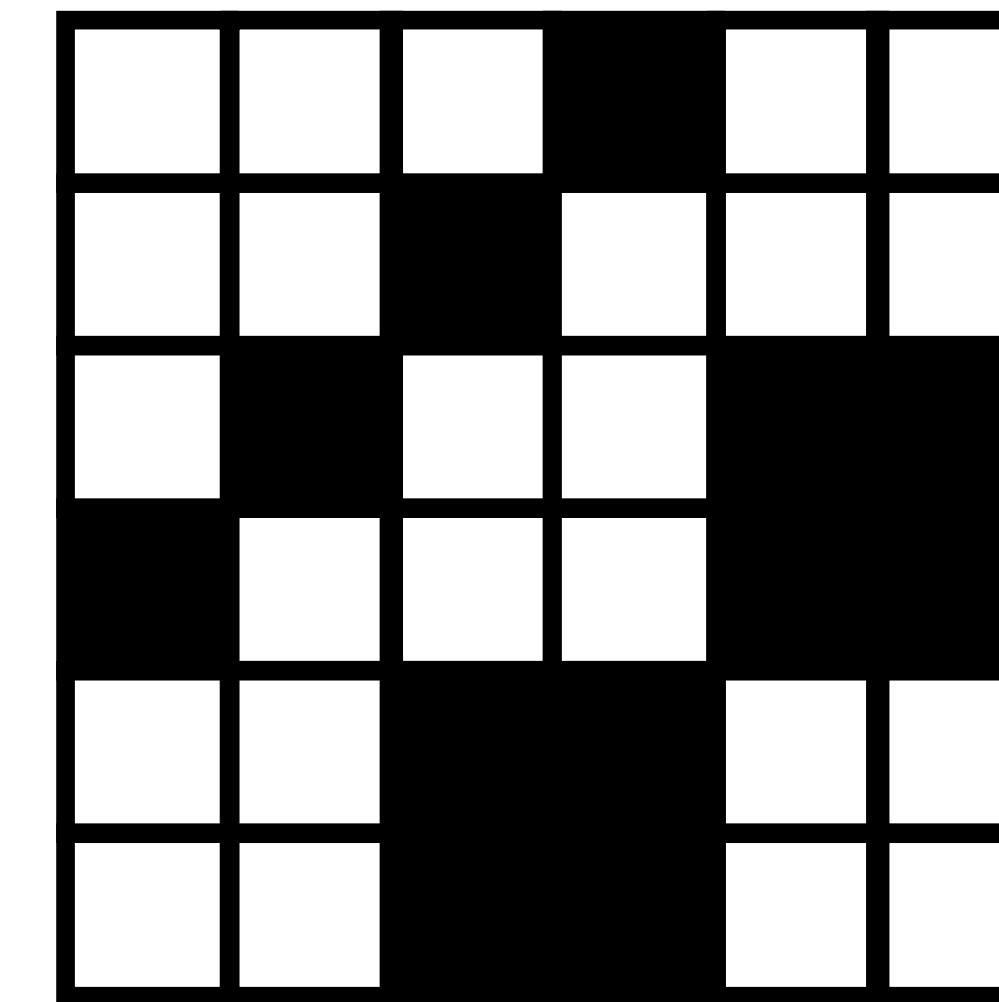
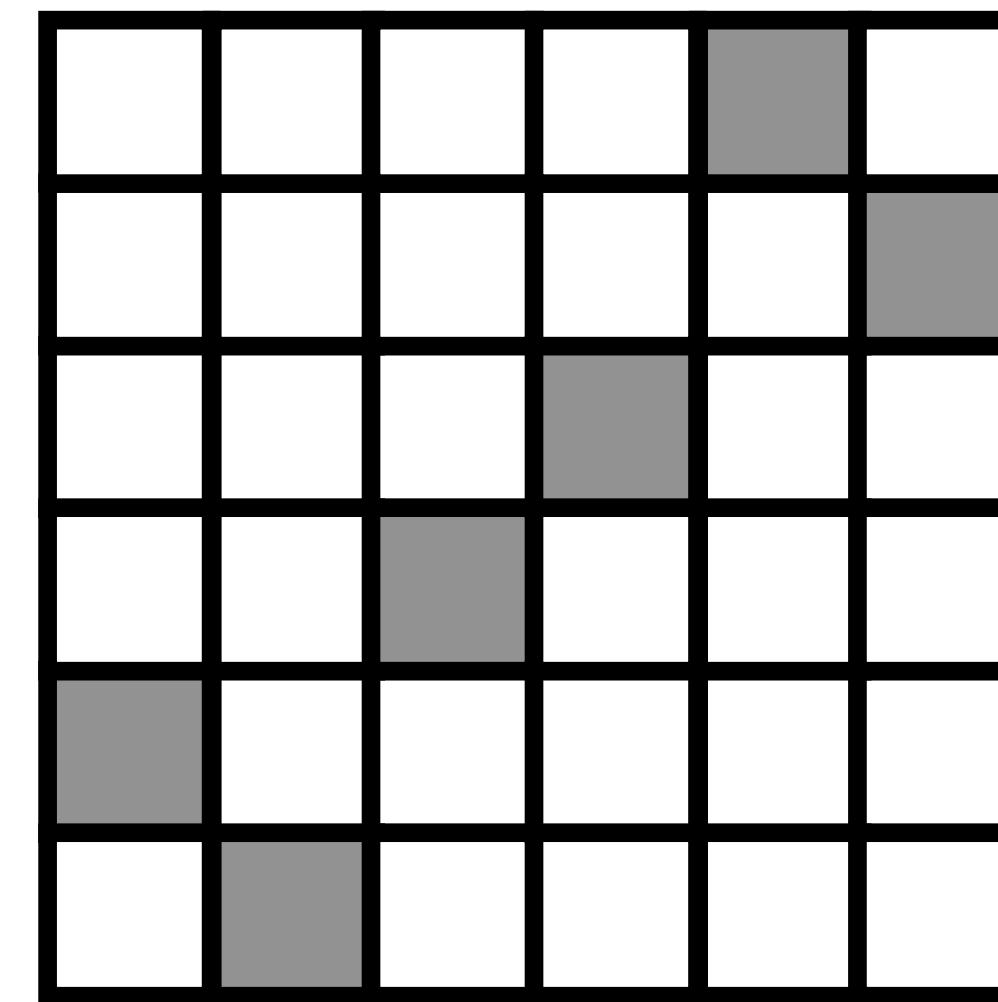
$$X' = PX$$

Permutation Matrices

We can represent permutations using a **permutation matrix**



=

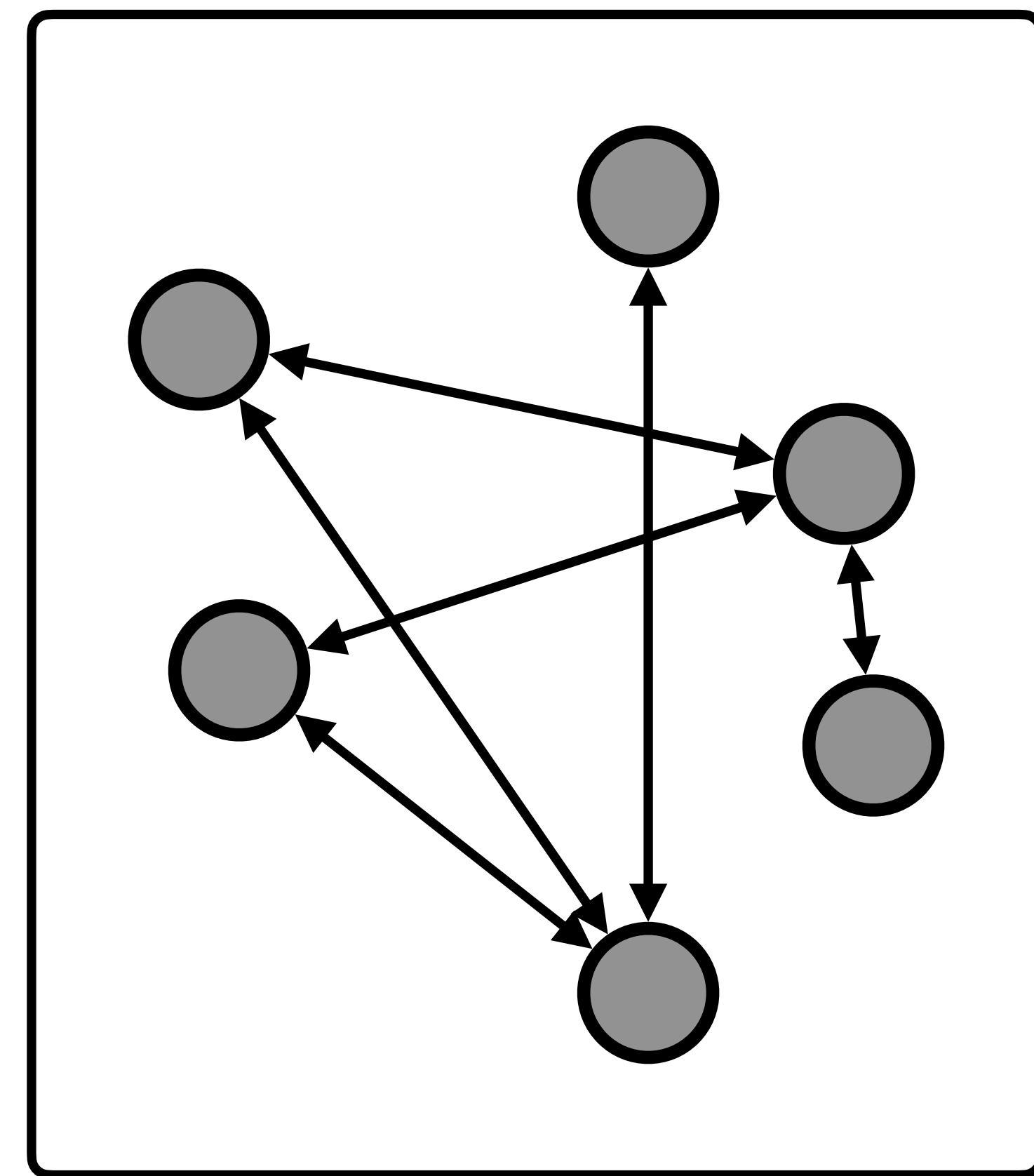


$$A' = PAP^T$$

Function Symmetries

Node-level functions create new node data from old node data

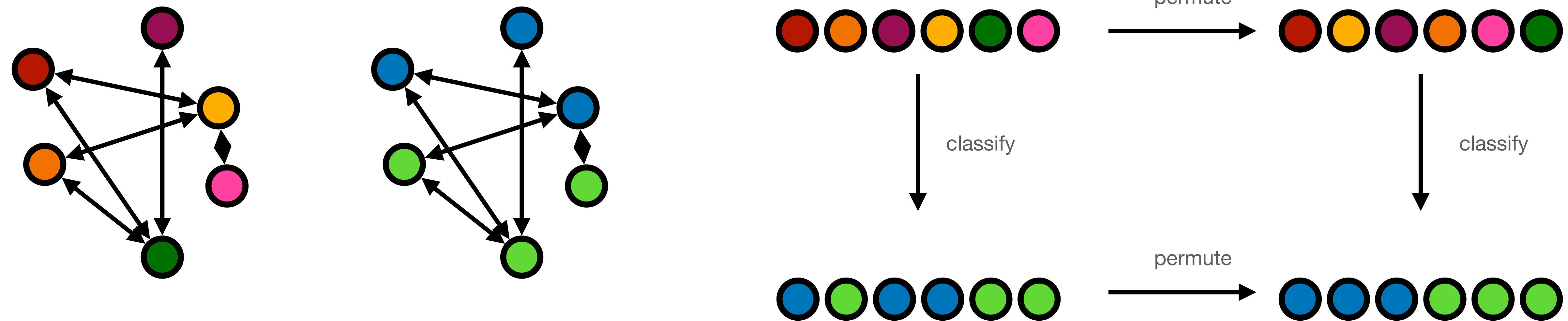
$$Y = f(X)$$



Function Symmetries

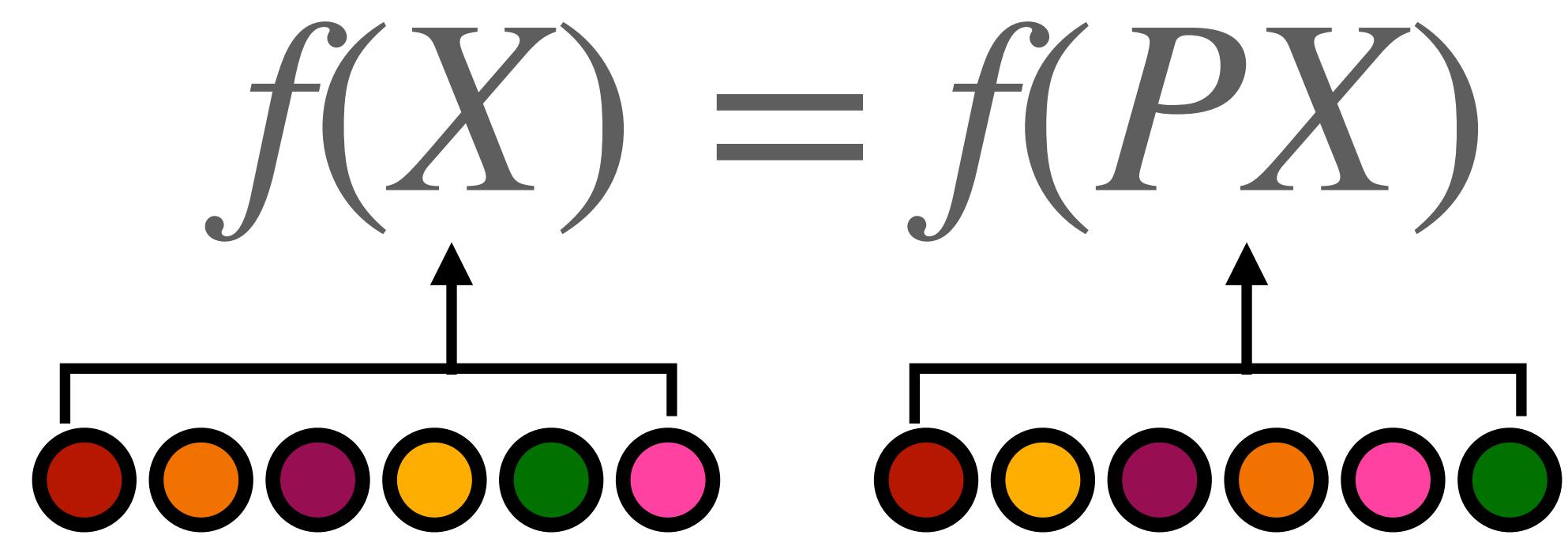
Node-level functions should
be permutation equivariant

$$Pf(X) = f(PX)$$

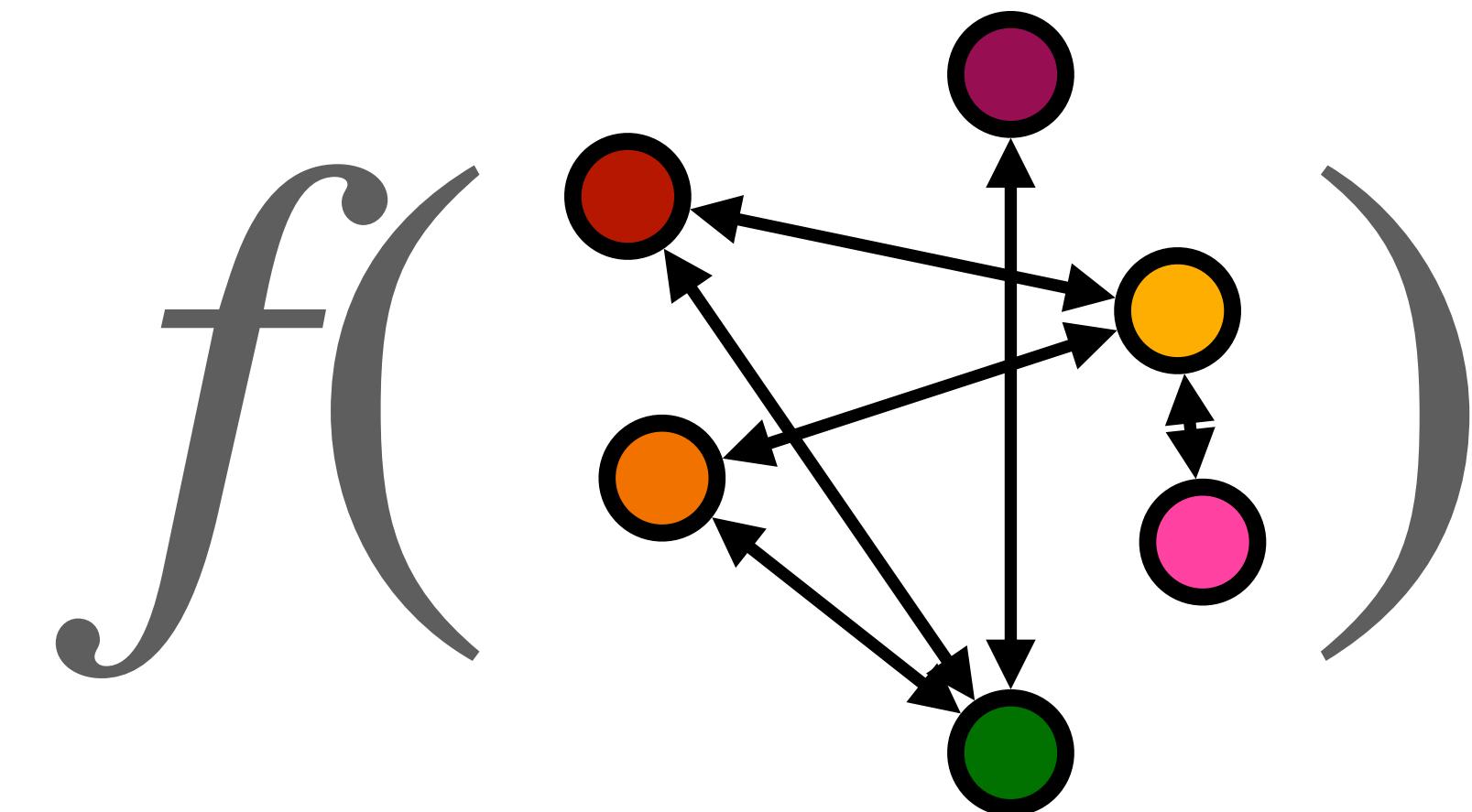


Function Symmetries

Graph-level functions should
be permutation-invariant



order should not matter

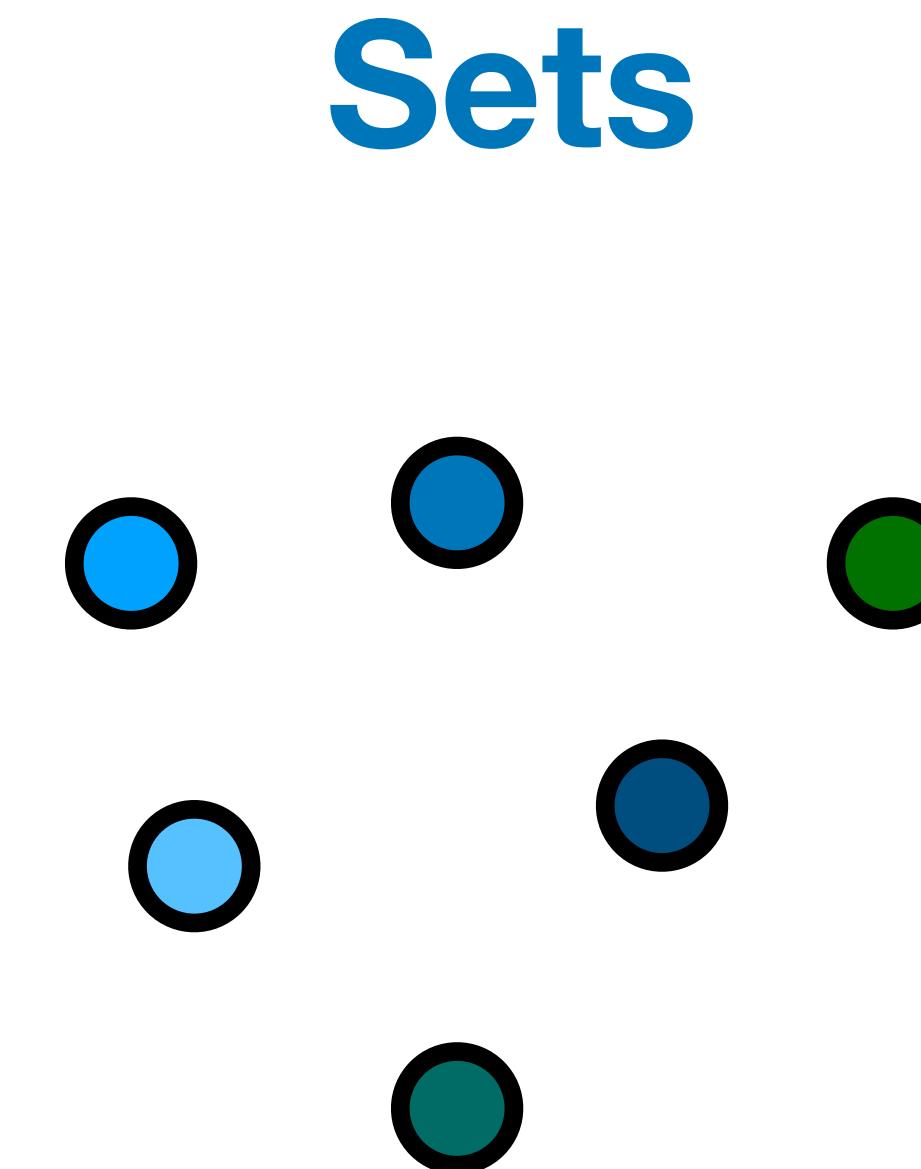


purely a function of the graph

Graphs w/o Edges

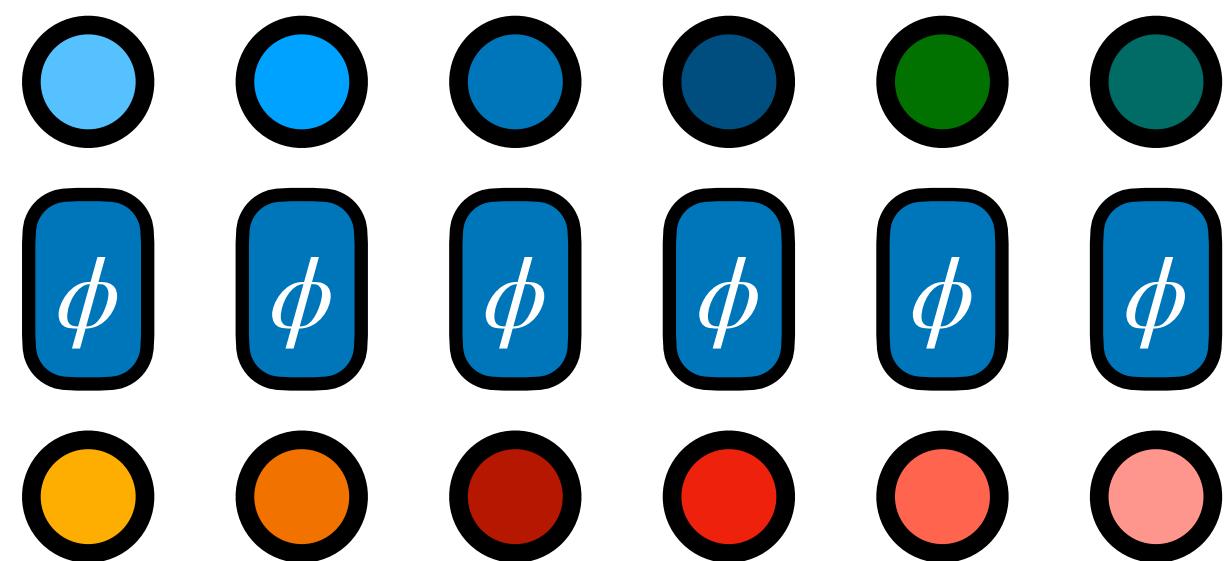
Function on Sets

We can get a glimpse how to get permutation invariant functions by looking at **graphs without connectivity**:



Function on Sets

Without any connectivity, there is no neighborhood, the only way to transform data is element-wise

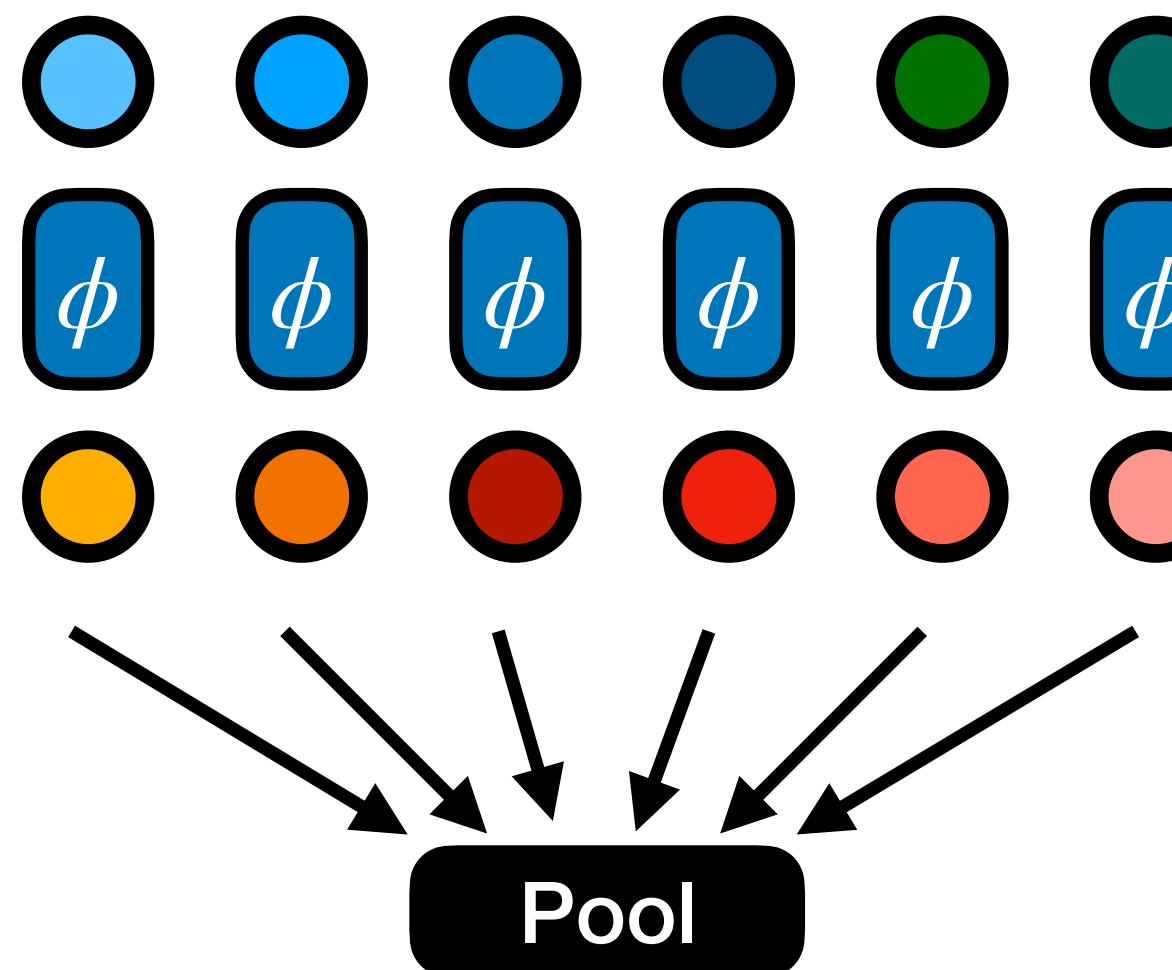


$$x_i \rightarrow y_i = \phi(x_i)$$

No connections = trivially permutation equivariant

Permutation-Invariance

To get combine the outputs we can add a **permutation-invariant pooling function**

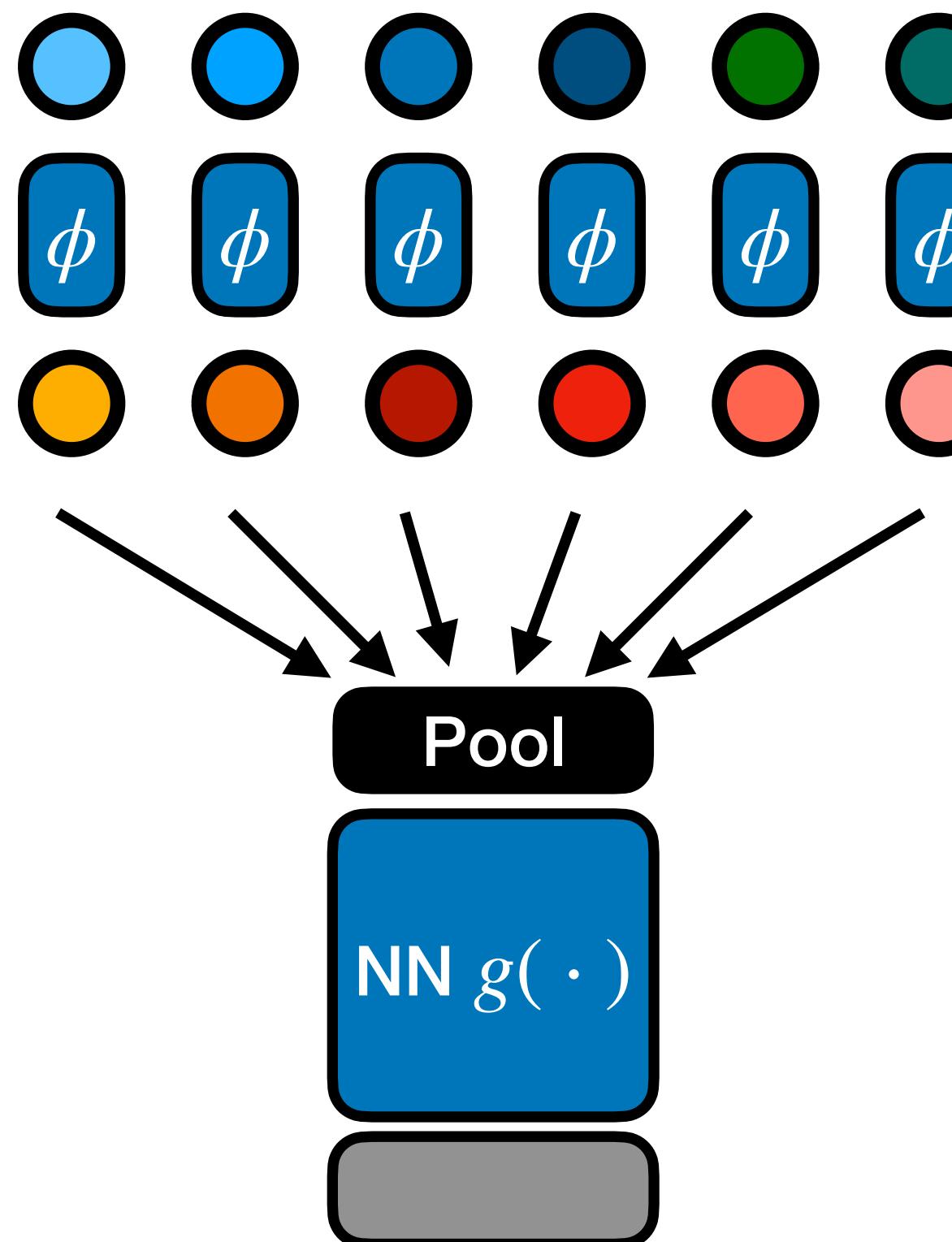


Example:
The sum is permutation-invariant

$$f(x_1, \dots, x_n) = \sum_i y_i = \sum_i f(x_i)$$

Permutation-Invariance

The pooled output can then be processed by any given “head” function to produce the final output



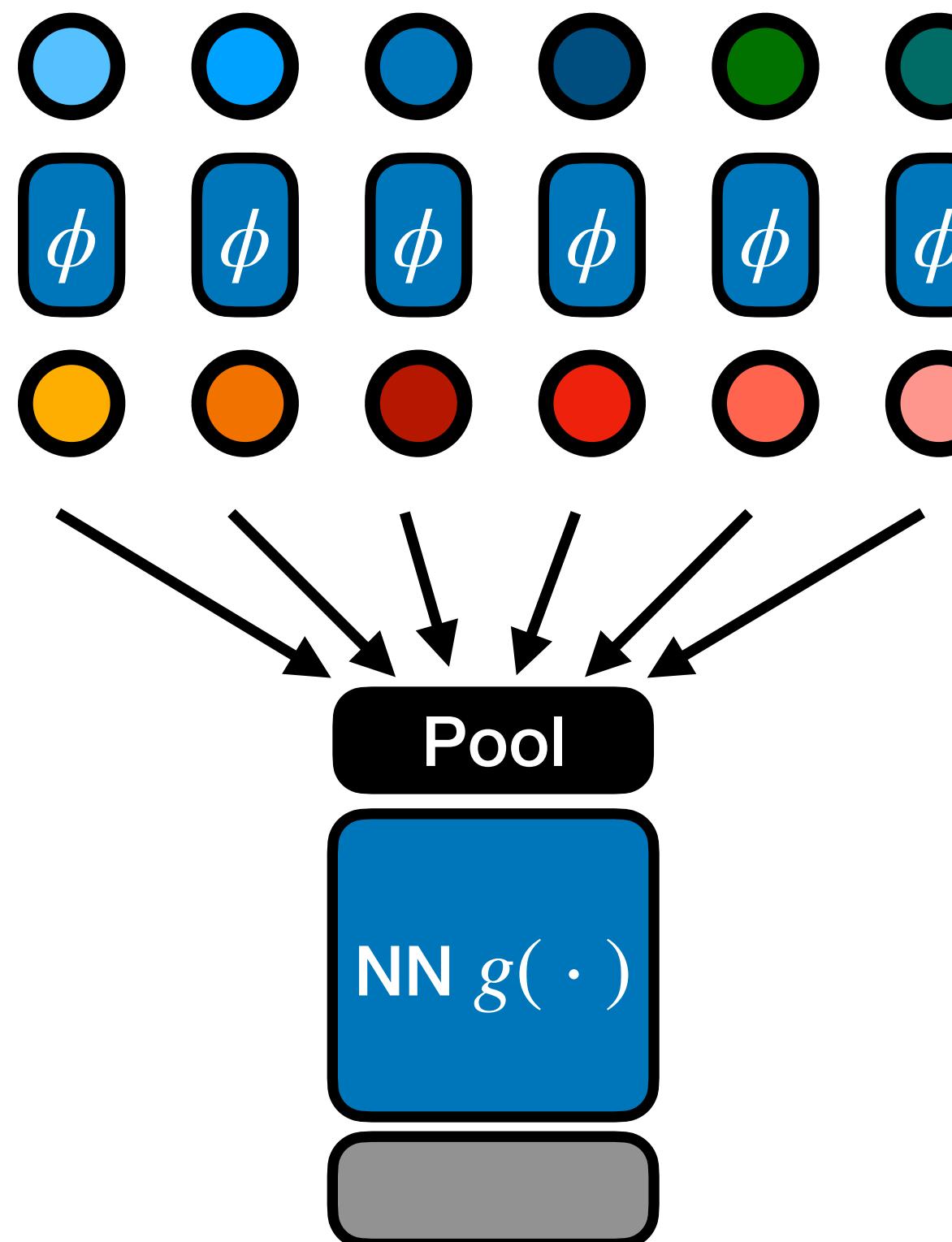
$$f(X) = g(\phi(x_1, \dots, x_n))$$

**the resulting function is expressive via
 $\phi(\cdot), g(\cdot)$ but permutation invariant**

$$f(X) = f(PX)$$

Permutation-Invariance

The pooled output can then be processed by any given “head” function to produce the final output



$$f(X) = g(\phi(x_1, \dots, x_n))$$

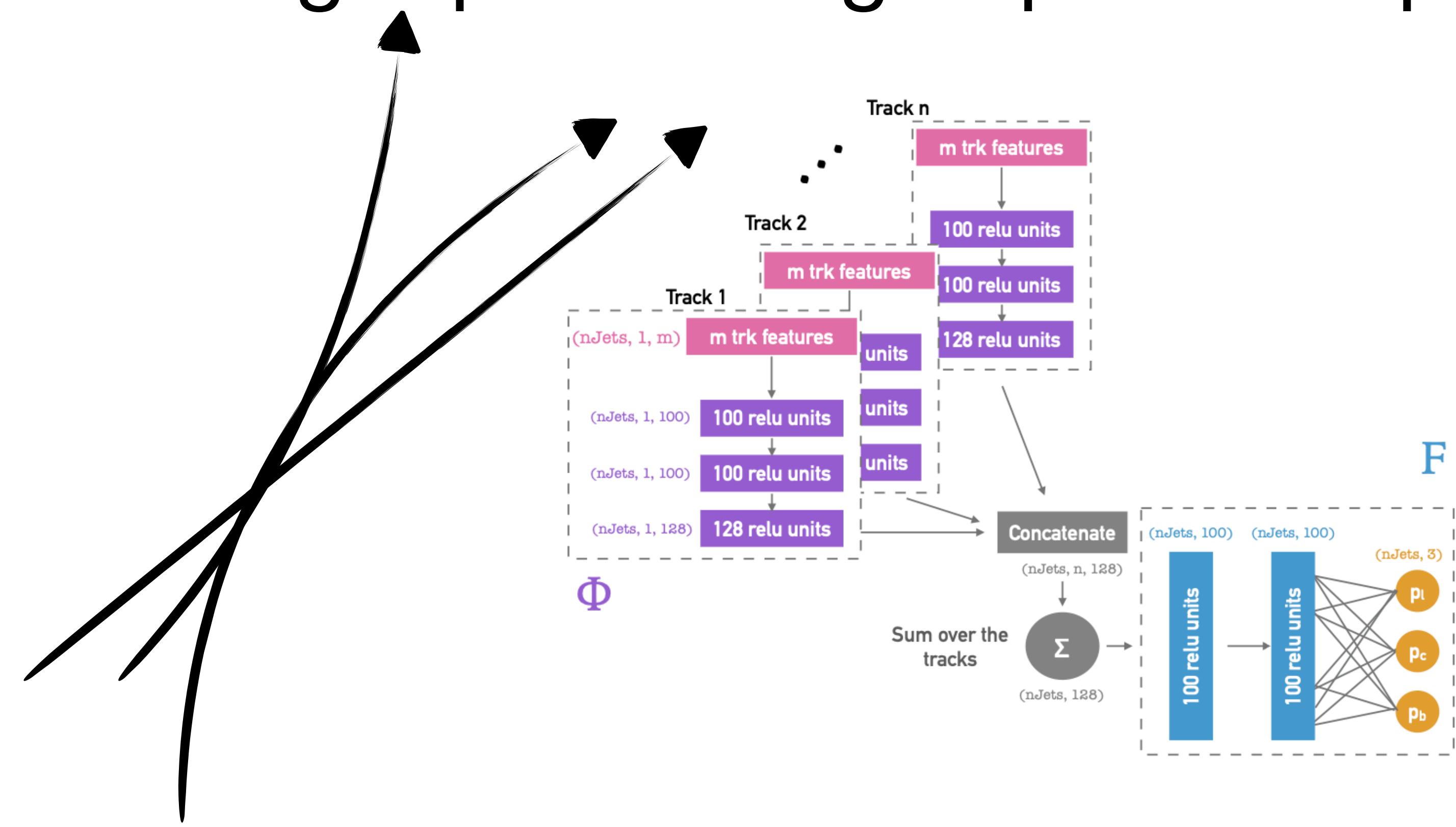
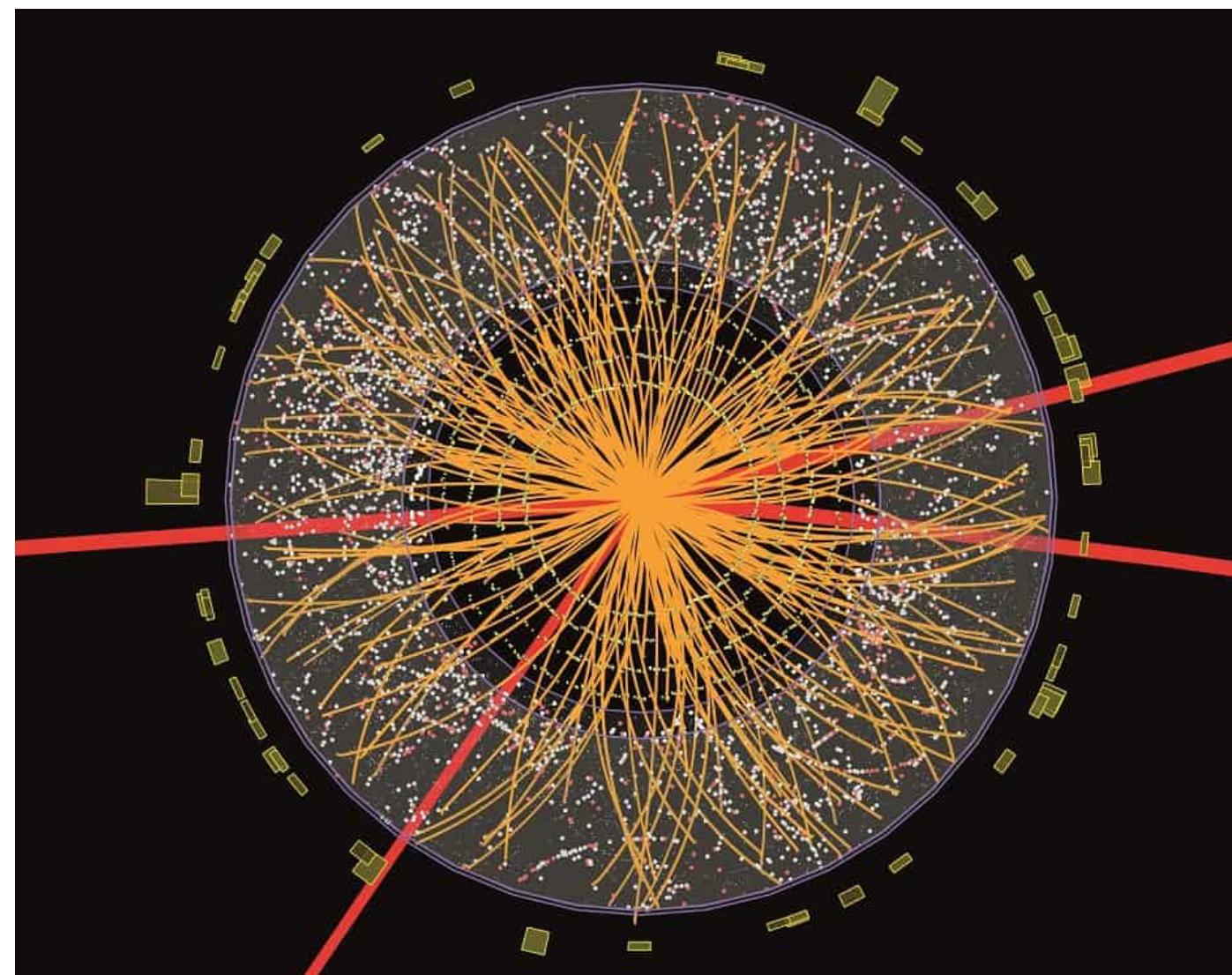
These networks are called:
Deep Sets

they encode a set of object into
some representation

Example

Particle Physics often deals w/ set of objects w/ no order

Here: per-tracks embedding + processing of pooled rep.



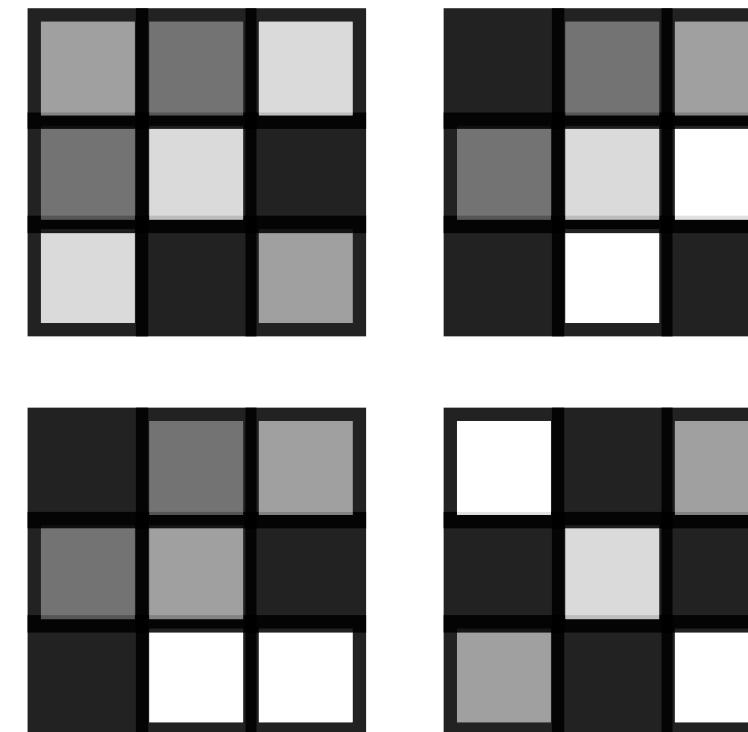
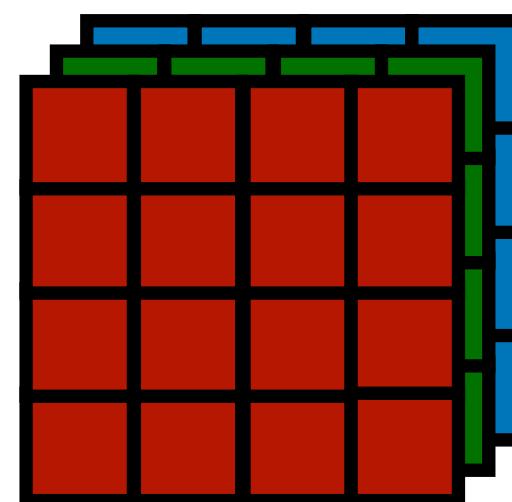
Graphs w/ Edges

Graph Convolutions

In CNNs we had the idea local processing

- underlying concept for scalability of CNN to arbitrary sizes

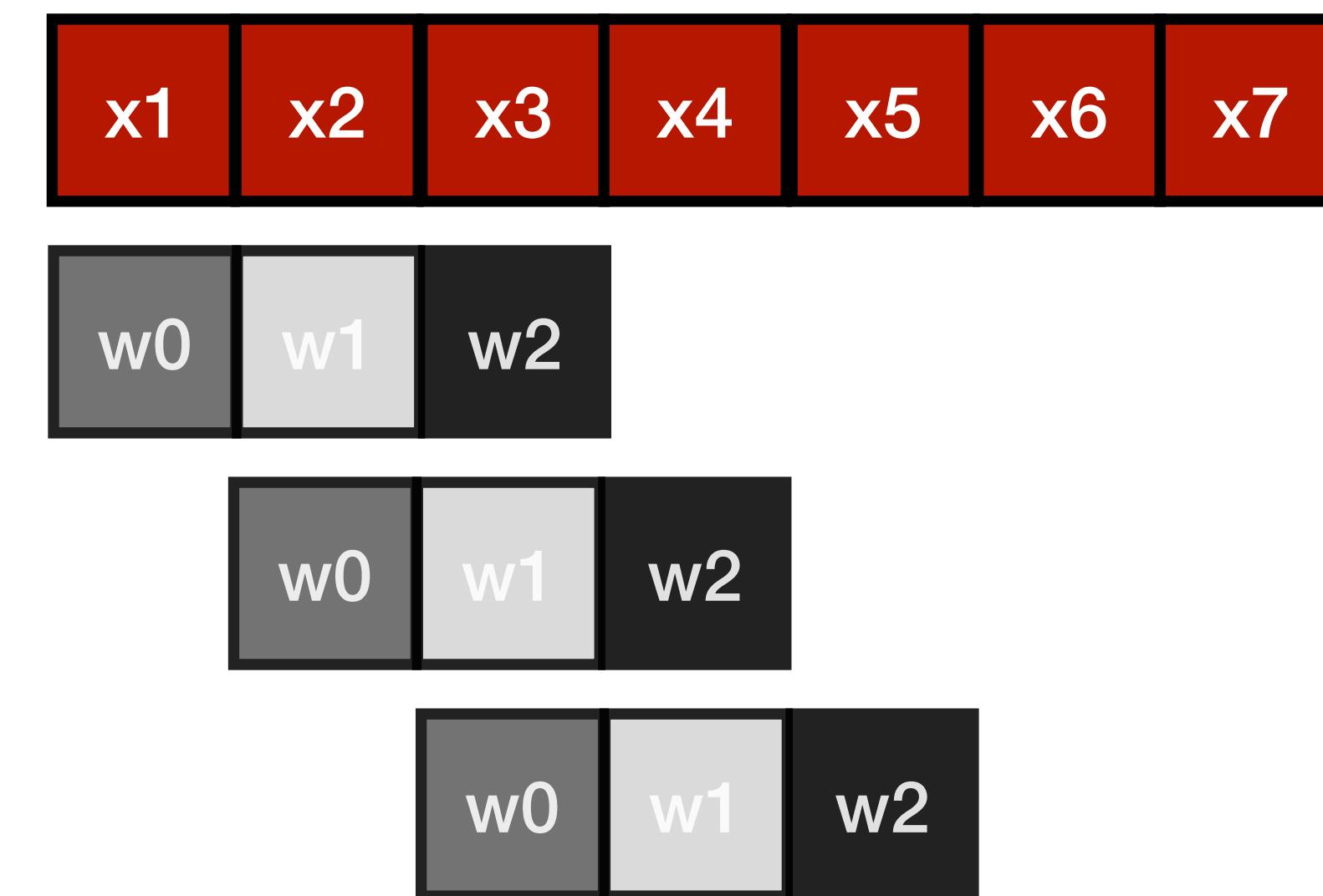
How can we extend this to graphs?



Graph Convolutions

We can go back to 1-D convolutions

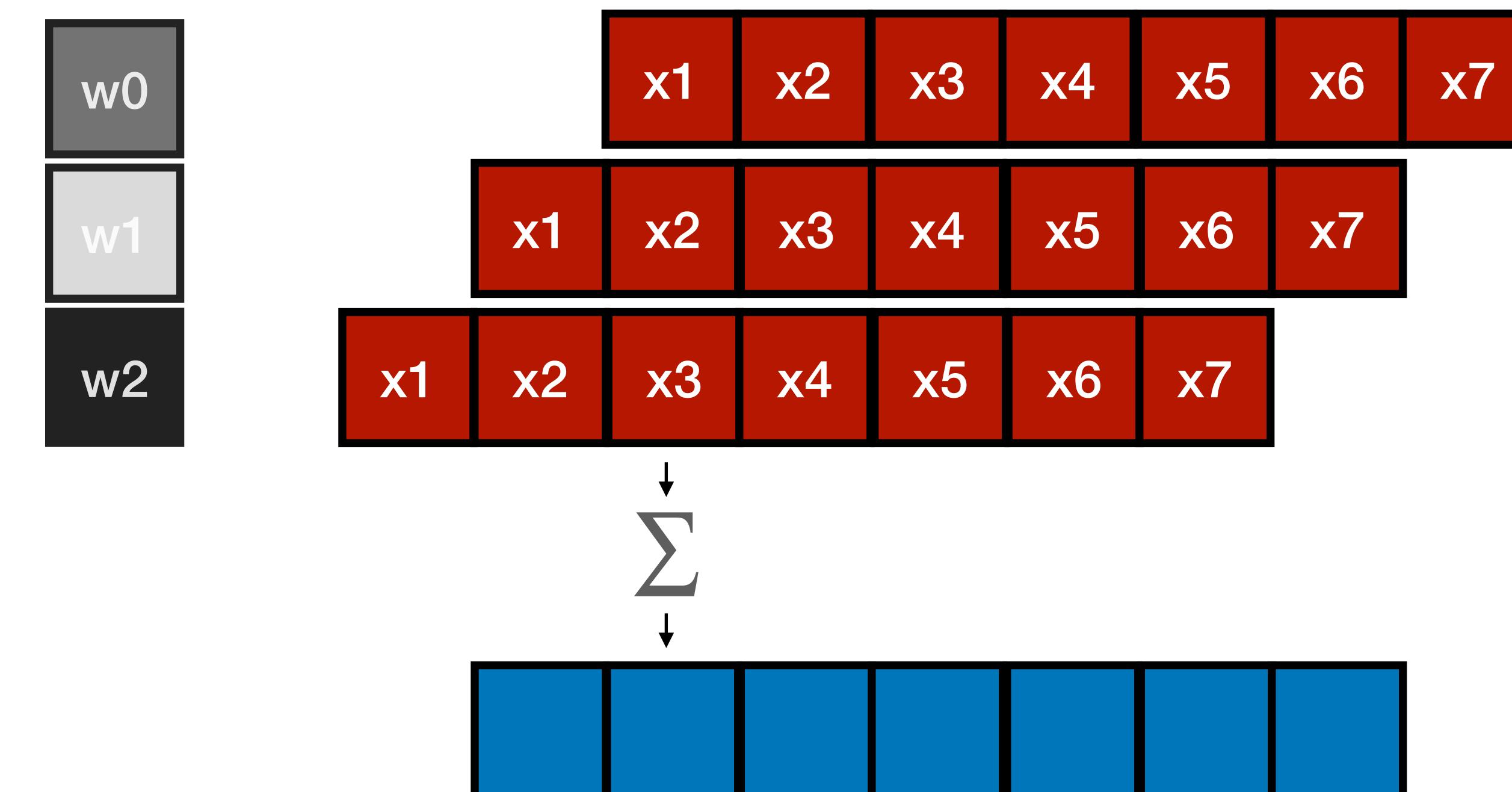
- normally we think of it at sliding the filter across



Graph Convolutions

We can go back to 1-D convolutions

- but we can also look at shifting the data



$$y_i = \sum x_{i+o} w_o$$

Graph Convolutions

We can go back to 1-D convolutions

- can rewrite x_{i+o} as the result of applying a shift operation
- convolution becomes formally like a polynomial

$$x_{i+1} = S(x_i)$$

$$x_{i+2} = S(S(x_i)) = S^2(x_i)$$

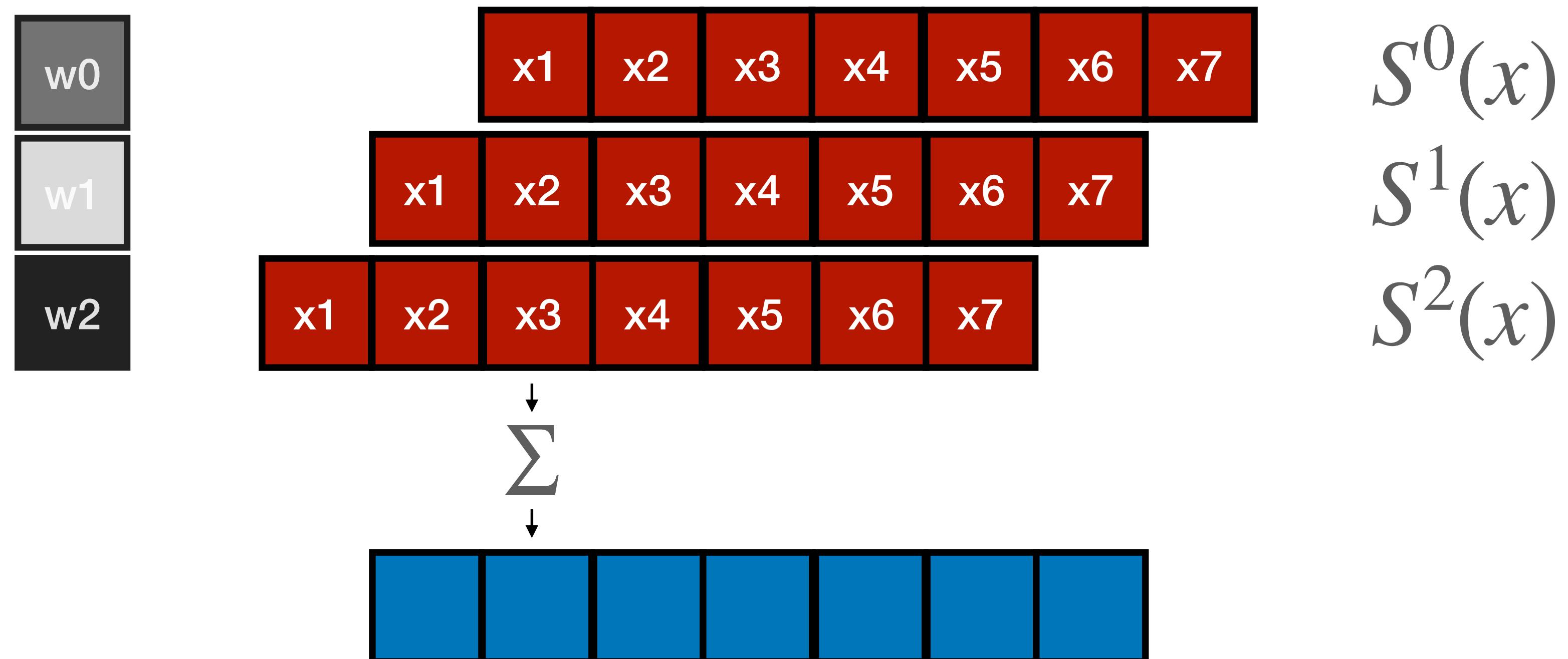
...

$$y_i = \sum x_{i+k} w_k$$

$$y_i = \sum_k S^k(x) w_k$$

higher orders in polynomial = contributions from further away

Graph Convolutions



$$y_i = \sum_k S^k(x) w_k$$

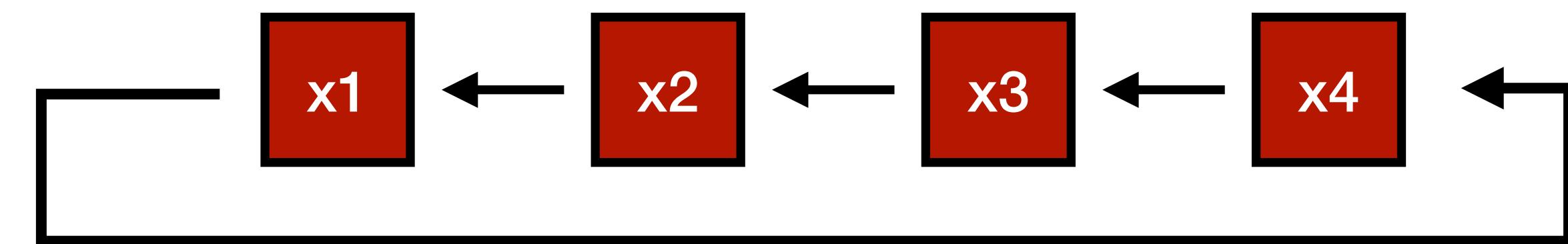
Graph Convolutions

To connect it back to the graph, picture, we can view a single line as a directed graph

- adjacency matrix = diagonal offset by 1

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$



Graph Convolutions

The shift operation $S(x)$ is simply a matrix multiplication by the adjacency matrix of the graph

$$A = \begin{matrix} & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \end{matrix} \quad X = \begin{matrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{matrix} \quad = \quad \begin{matrix} x_2 \\ x_3 \\ x_4 \\ x_1 \end{matrix} = \begin{matrix} & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \end{matrix} \quad \begin{matrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{matrix}$$

$$S(x) = AX$$

Recovering Convolution

$$y_i = \sum_k S^k(x) w_k = \text{convolution} \quad S^{k=-1} \quad w_{-1} + S^{k=0} \quad w_0 + S^{k=1} \quad w_1$$

The diagram illustrates the decomposition of a convolution operation. At the top, a 5x5 input matrix is shown being equal to the sum of three 5x5 weight matrices and a 5x5 bias matrix. The input matrix has colored blocks (red, orange, teal) at specific positions. The weight matrices w_{-1} , w_0 , and w_1 also have colored blocks at specific positions. The bias matrix $S^{k=0}$ has colored blocks at specific positions. Below this, another 5x5 input matrix is shown being equal to the sum of three 5x5 weight matrices and a 5x5 bias matrix. The input matrix has black blocks at specific positions. The weight matrices $S^{k=-1}$, $S^{k=0}$, and $S^{k=1}$ also have black blocks at specific positions. The bias matrix w_{-1} , w_0 , and w_1 also have black blocks at specific positions.

Graph Convolutions

With this, we can generalize the concept of a convolution to graphs, using a generalized shift operation on graphs

Graph Filter

$$y_i = \sum_k S^k(x)w_k$$

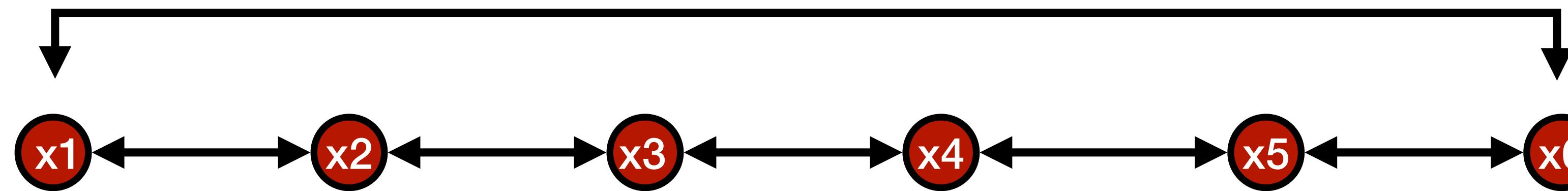
convolution

$$Y = \sum_k S^k X W_k$$

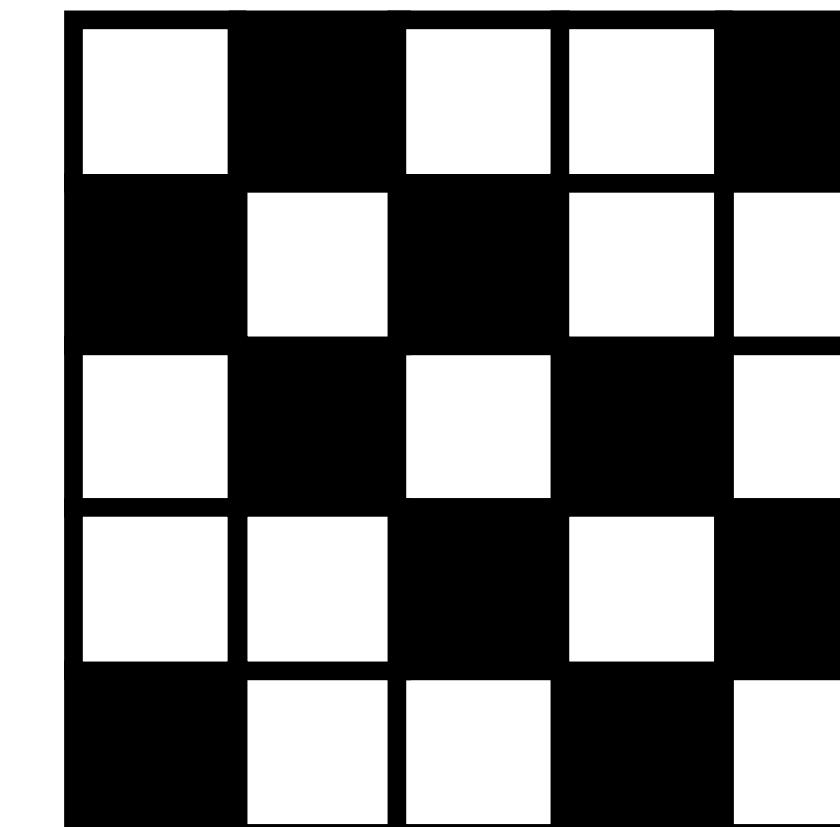
graph convolution

Graph Convolutions

Let's look at what this means for undirected linear graphs



| |
|----------------|
| w(x_2+x_5) |
| w(x_1+x_3) |
| w(x_2+x_4) |
| w(x_3+x_5) |
| w(x_4+x_1) |



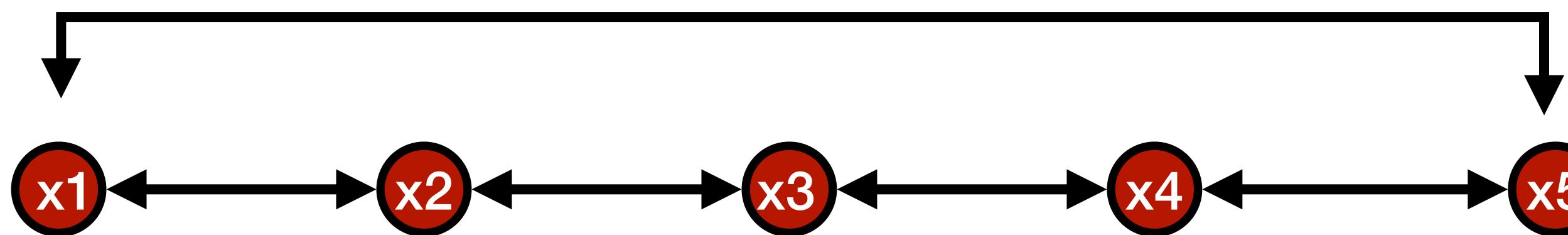
| |
|----|
| x1 |
| x2 |
| x3 |
| x4 |
| x5 |

W

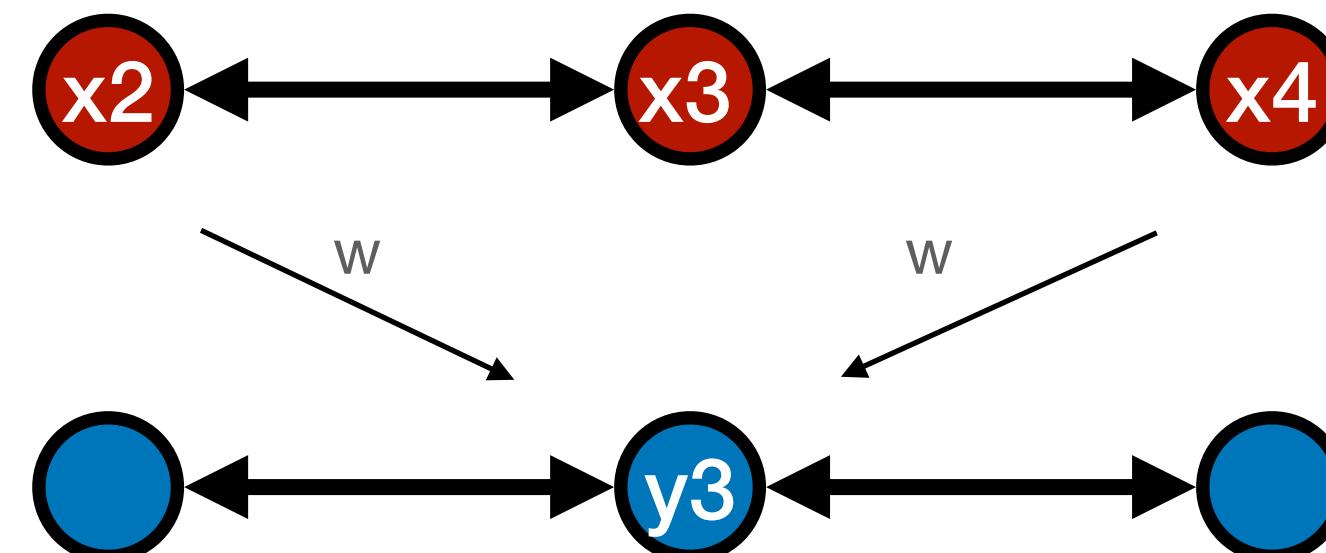
SXW_1

Graph Convolutions

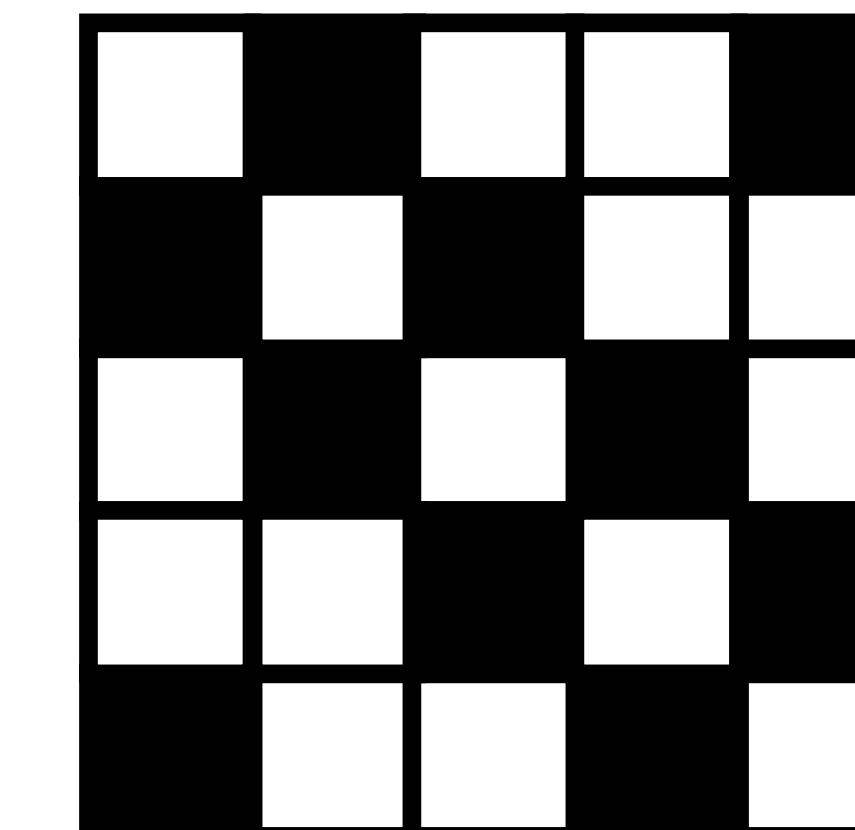
Let's look at what this means for undirected linear graphs



*Each node, gets a **weighted signal** from its local neighborhood, but signal is symmetric (all share the same weight)*



| |
|------------------------------------|
| w(x ₂ +x ₅) |
| w(x ₁ +x ₃) |
| w(x ₂ +x ₄) |
| w(x ₃ +x ₅) |
| w(x ₄ +x ₁) |



| |
|----------------|
| x ₁ |
| x ₂ |
| x ₃ |
| x ₄ |
| x ₅ |

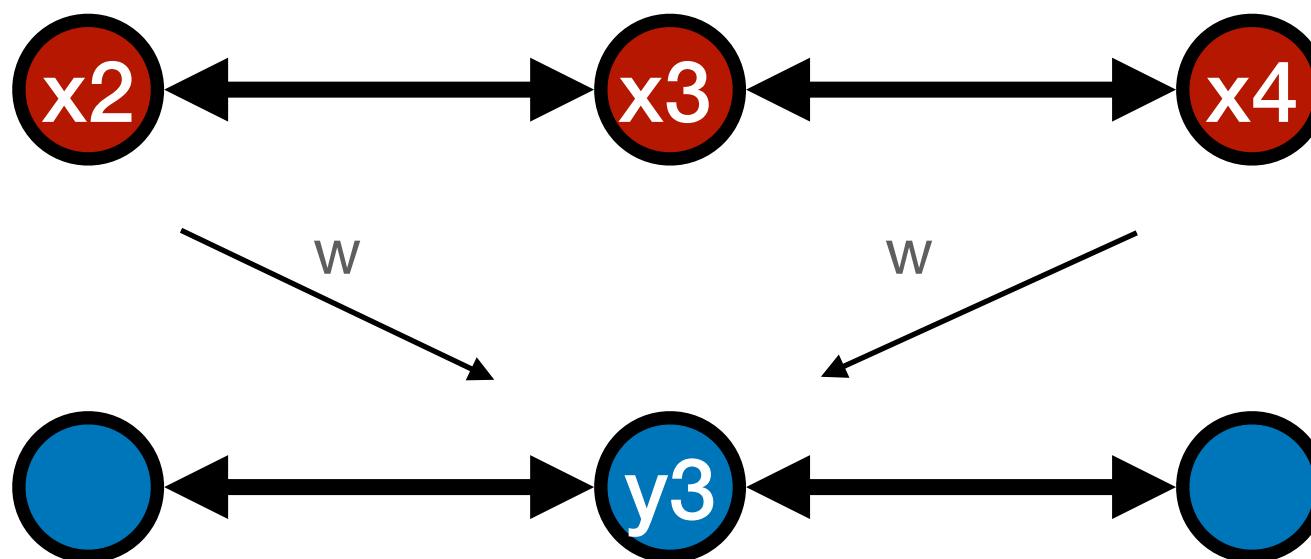
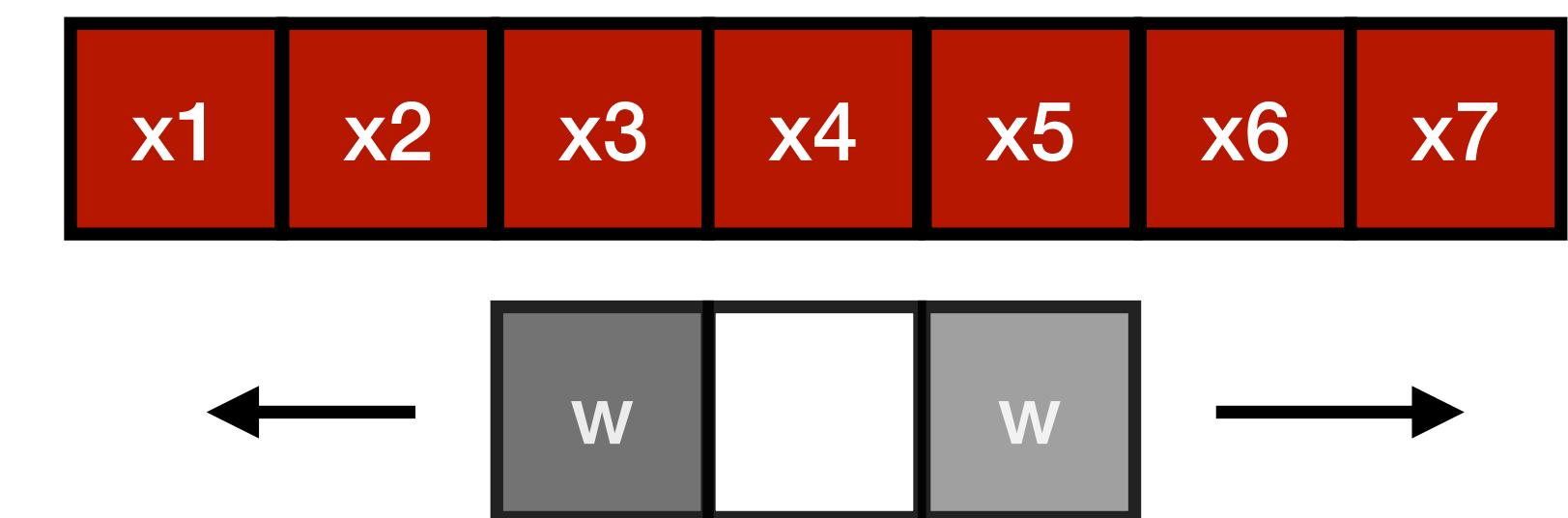
W

SXW_1

Graph Convolutions

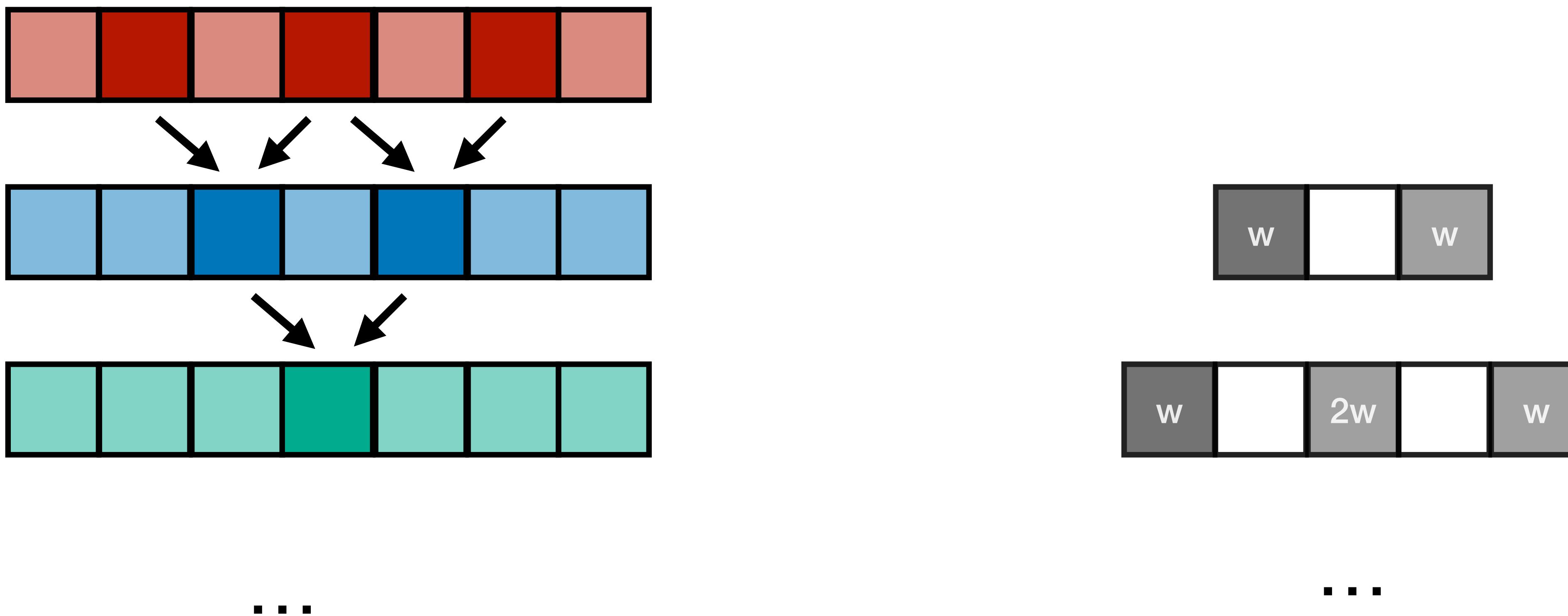
A “size-1” kernel on a symmetric graph is like a symmetric convolution filter extending to the nearest neighbors

Each node, gets a weighted signal from its local neighborhood, but signal is symmetric (all share the same weight)

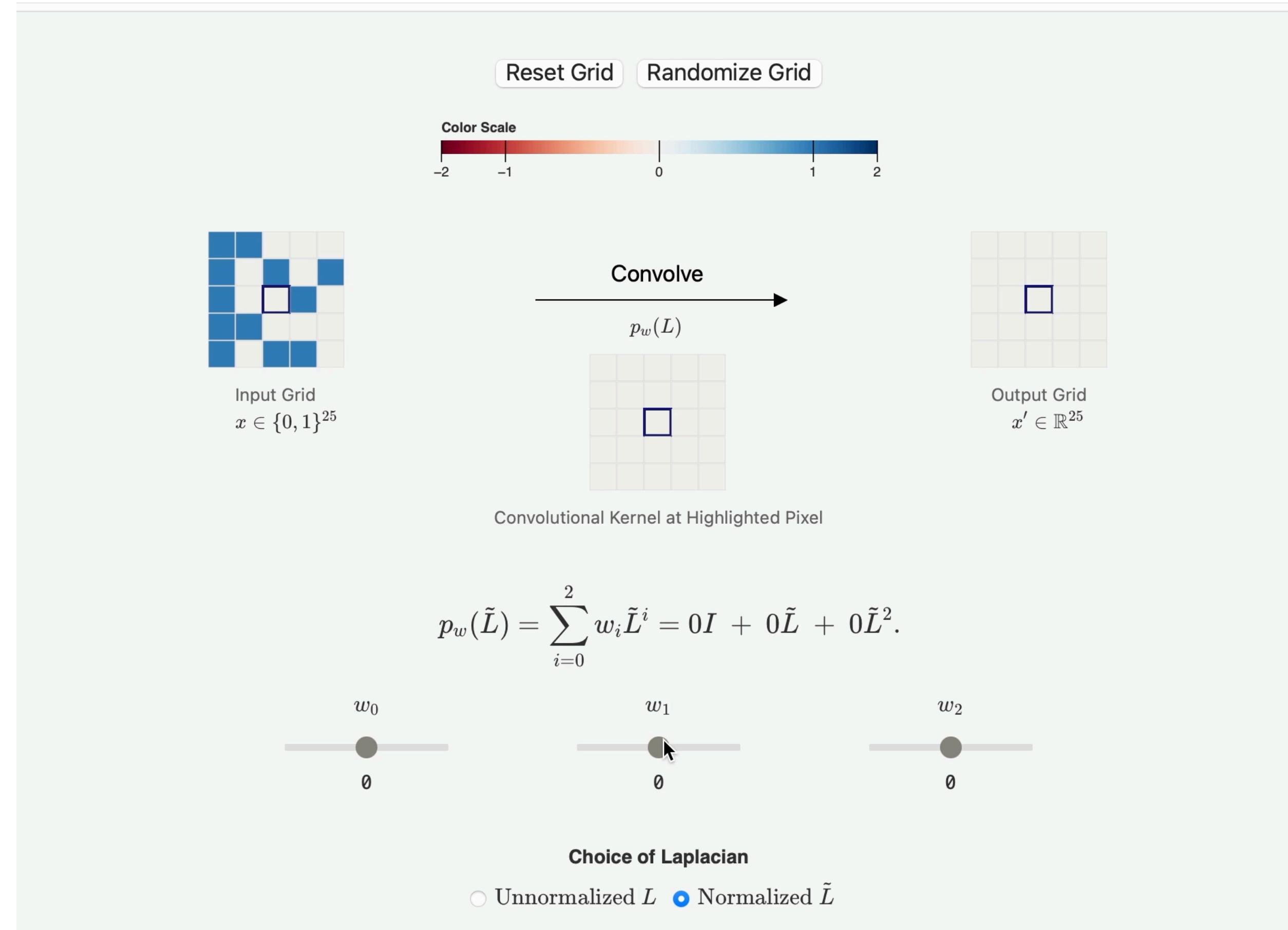


Graph Convolutions

As in the directed convolutions, higher k extend the reach of the neighborhood. Produces “radially symmetric” filters



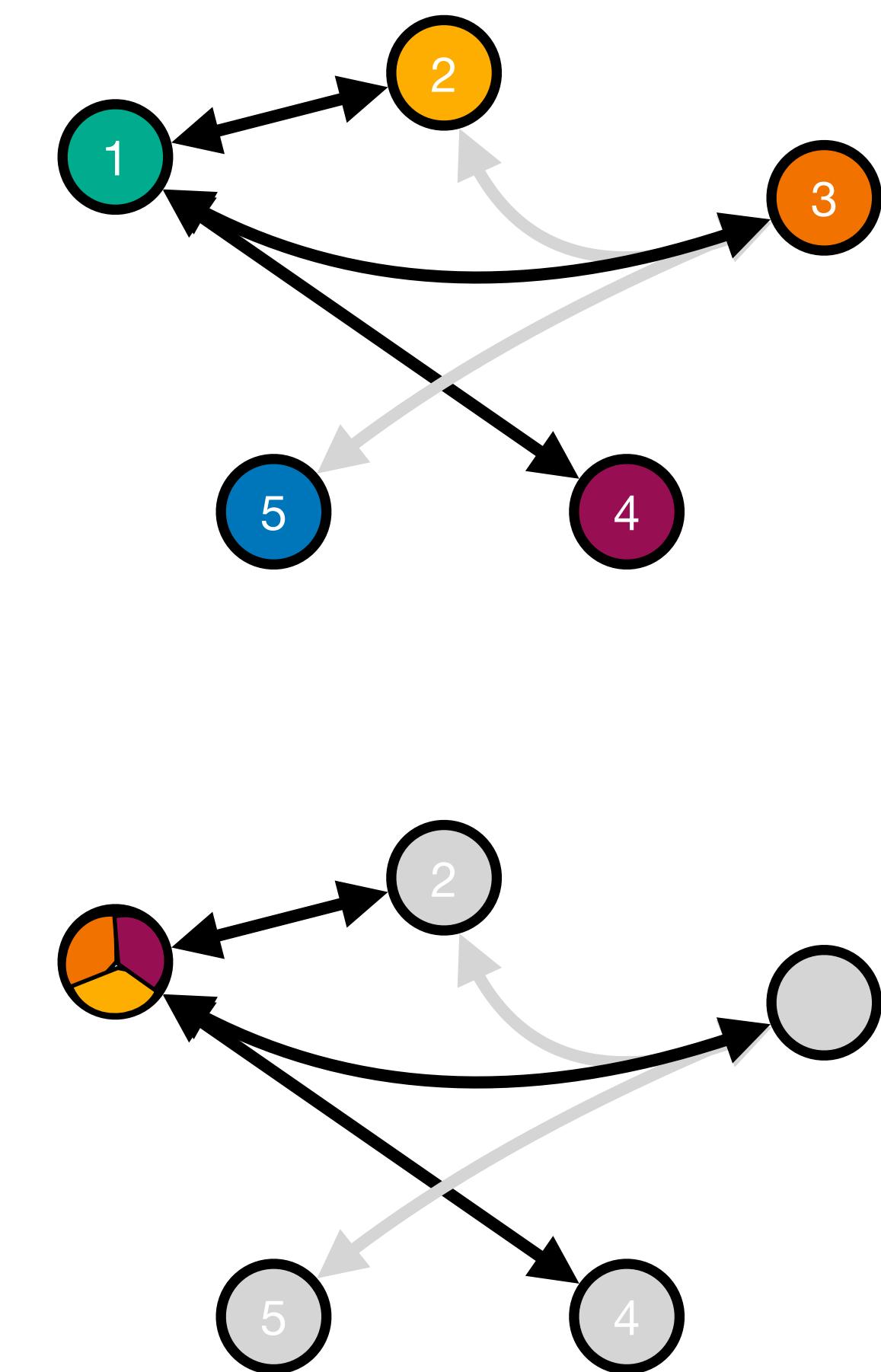
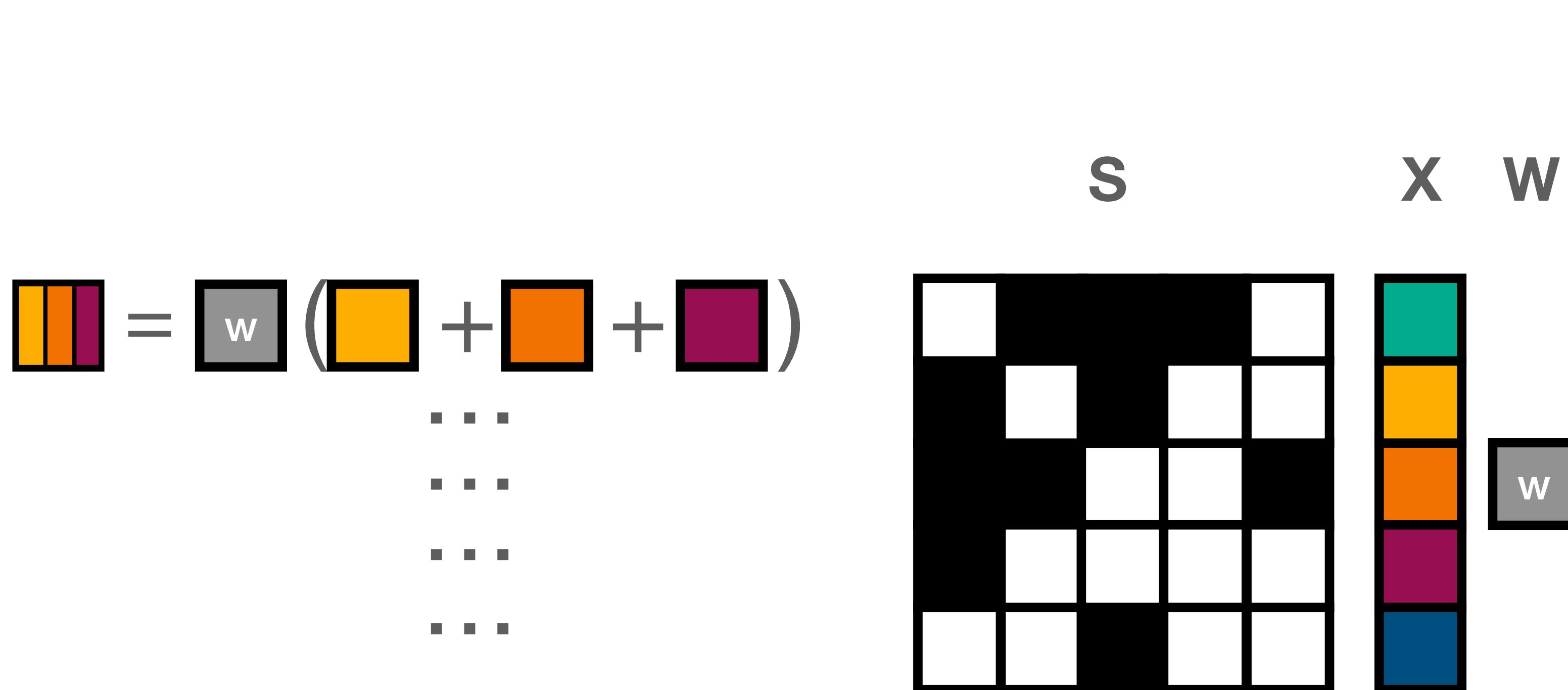
Try it yourself!



<https://distill.pub/2021/understanding-gnns/>

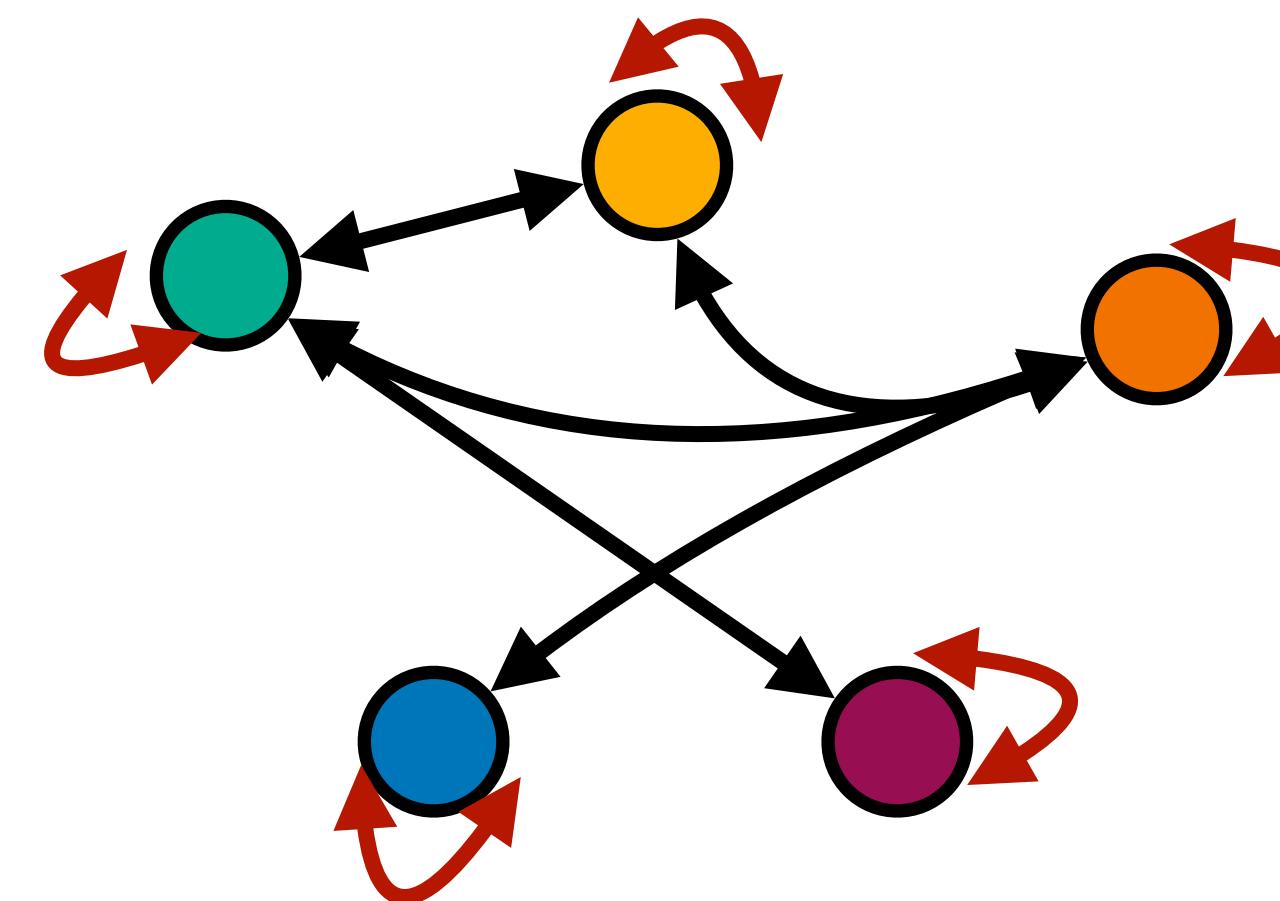
Non-linear Graphs

Our approach extends now easily to arbitrary graphs



Non-linear Graphs

In order to not forget any previous state, it's useful to add self-connections as well

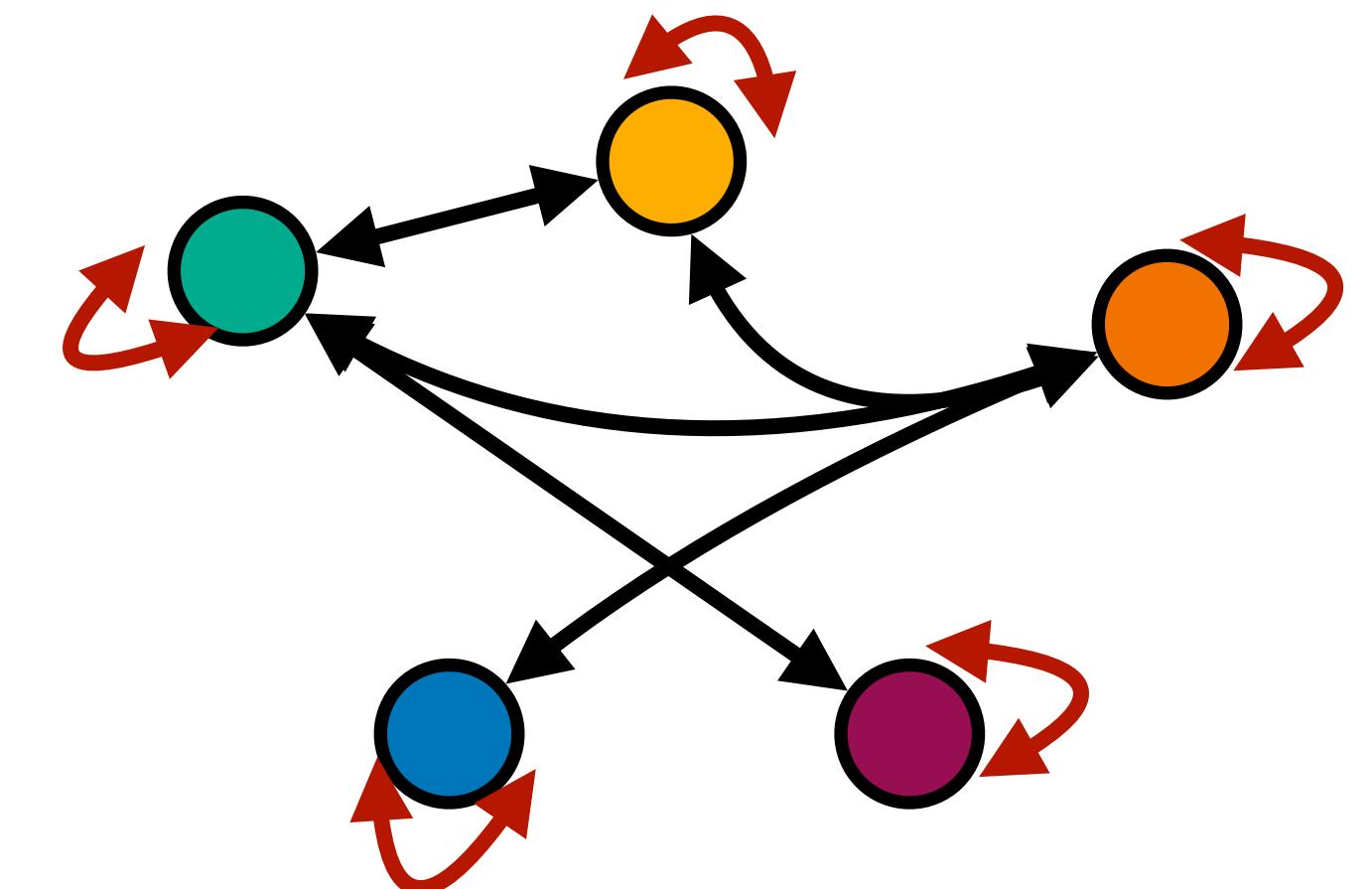
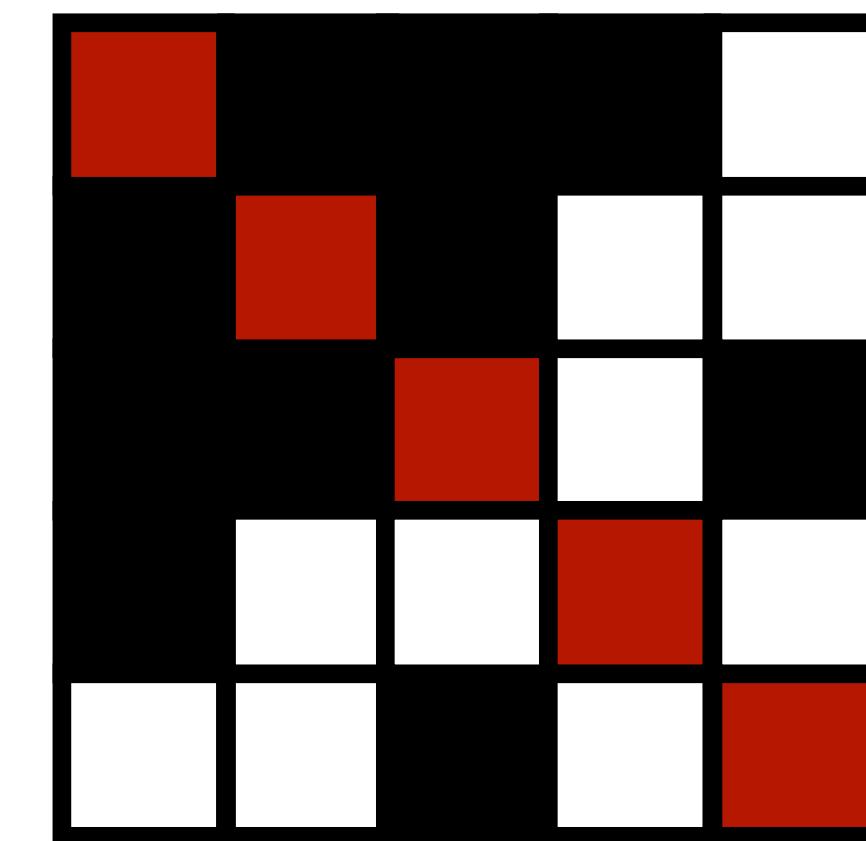


Choosing S

In practice a good choice for S is a combination of the adjacency matrix (neighborhood) and identity (self-links)

- scaled by inverse of number of edges (degree) for stability

$$S = D^{-1/2}(I + A)D^{-1/2}$$



Permutation Equivariance

Under permutation our adjacency and data transform like
 $A \rightarrow PAP^T$ and $X \rightarrow PX$, i.e. $S \rightarrow PSP^T$

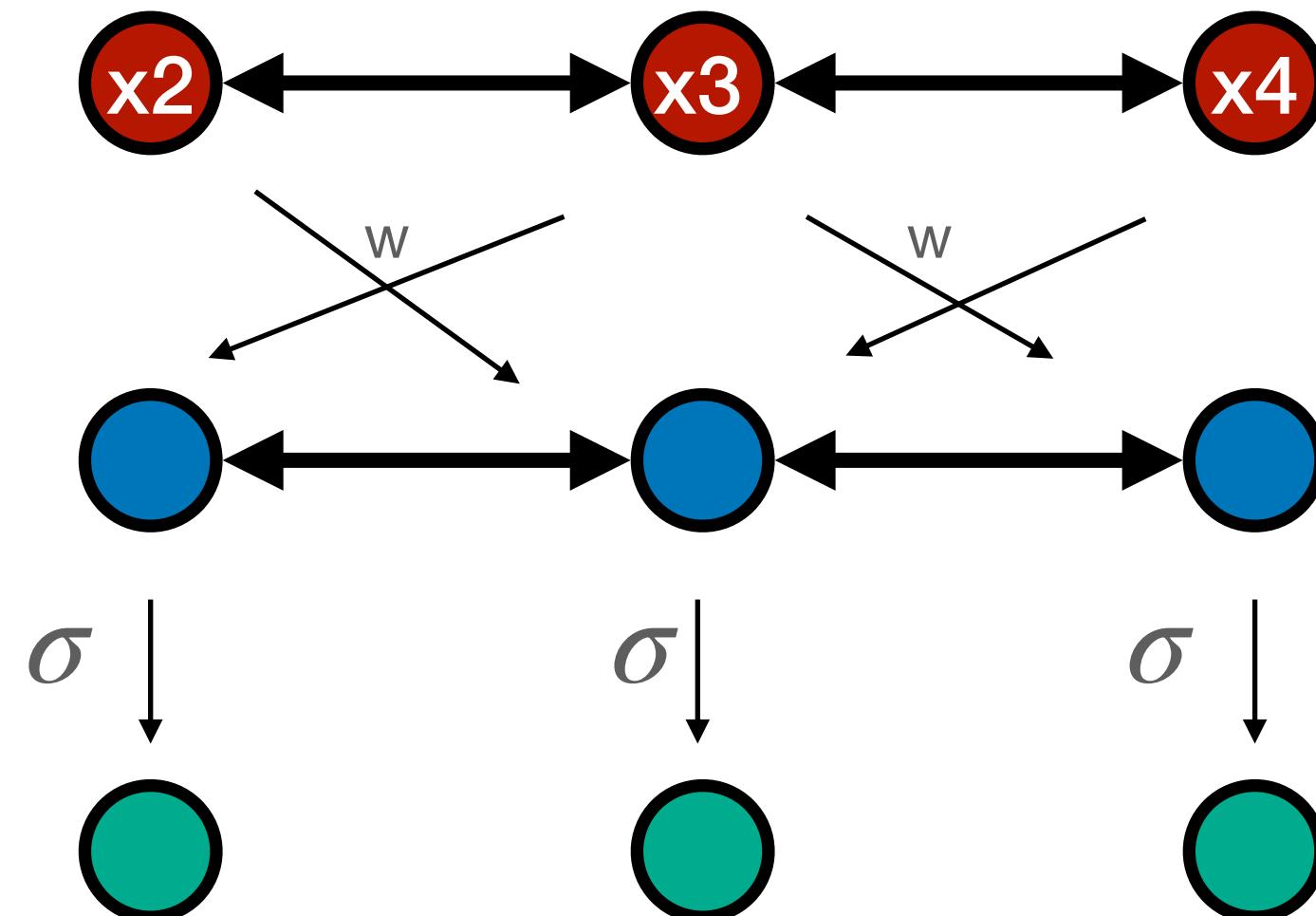
So we can see that:

$$f(X', S') = PSP^T PXW = PSXW = Pf(X, S)$$

- i.e out graph convolution is **permutation equivariant!**

Adding a non-linearity

To complete our “Graph Convolution” operation, we add a element-wise non-linearity :

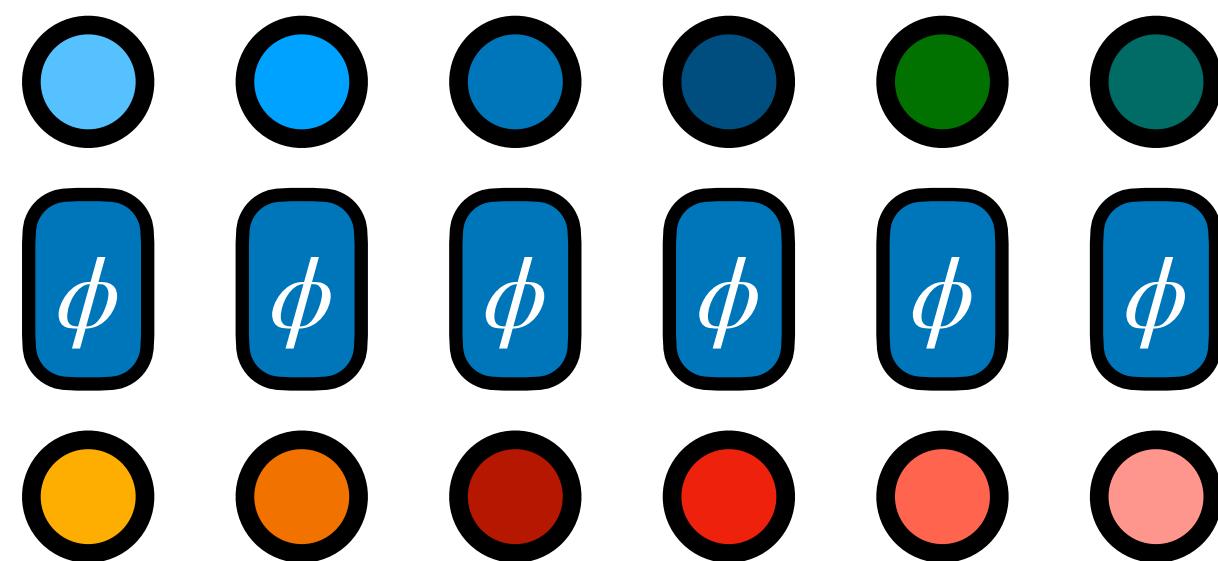


$$f(X) = \sigma(SXW)$$

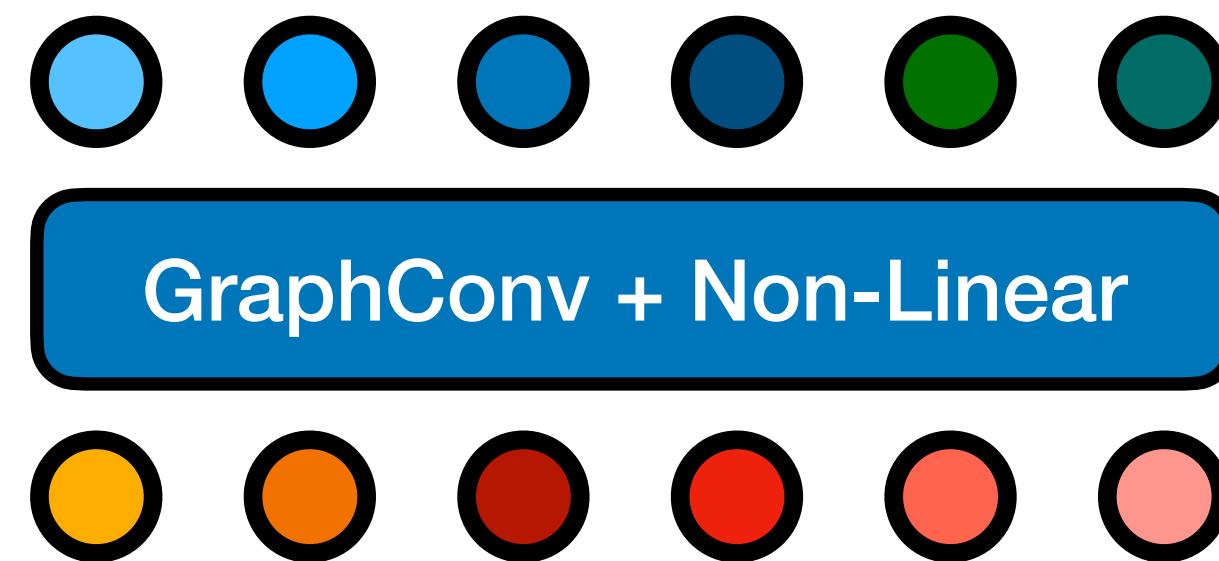
Similar to how we added non-linearities in CNN

Generalizing the Deep Set & CNN

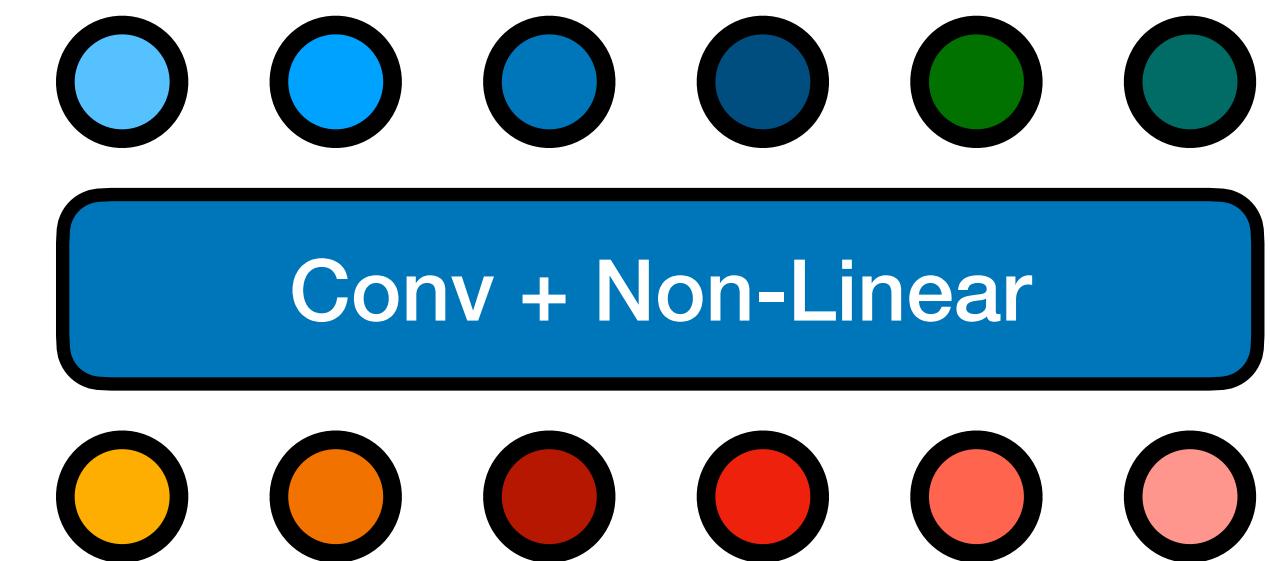
We can now build finally graph neural networks, by combining permutation-equivariant non-linear mappings



*Deep Set Layer
(permutation-invariant)*



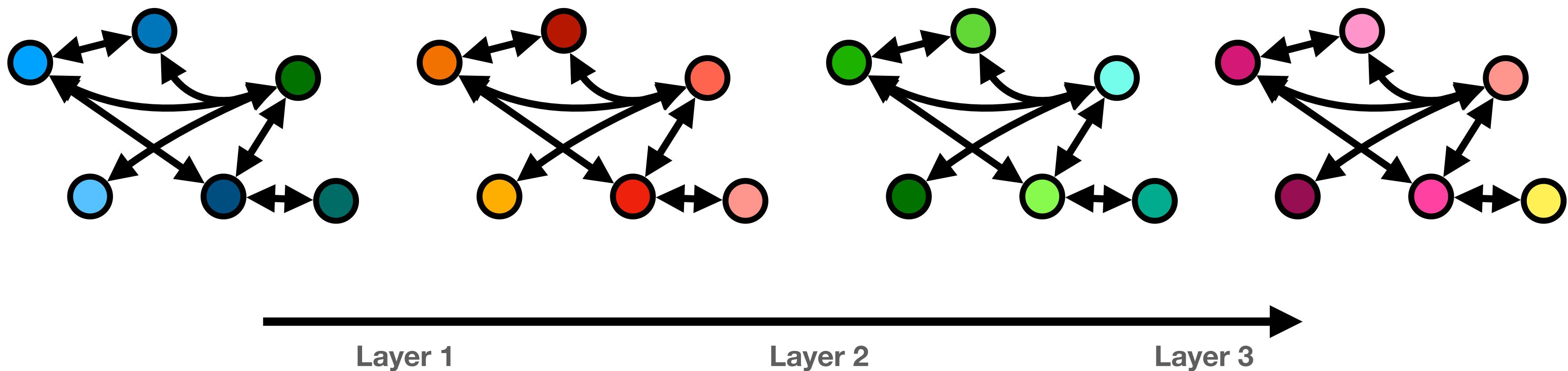
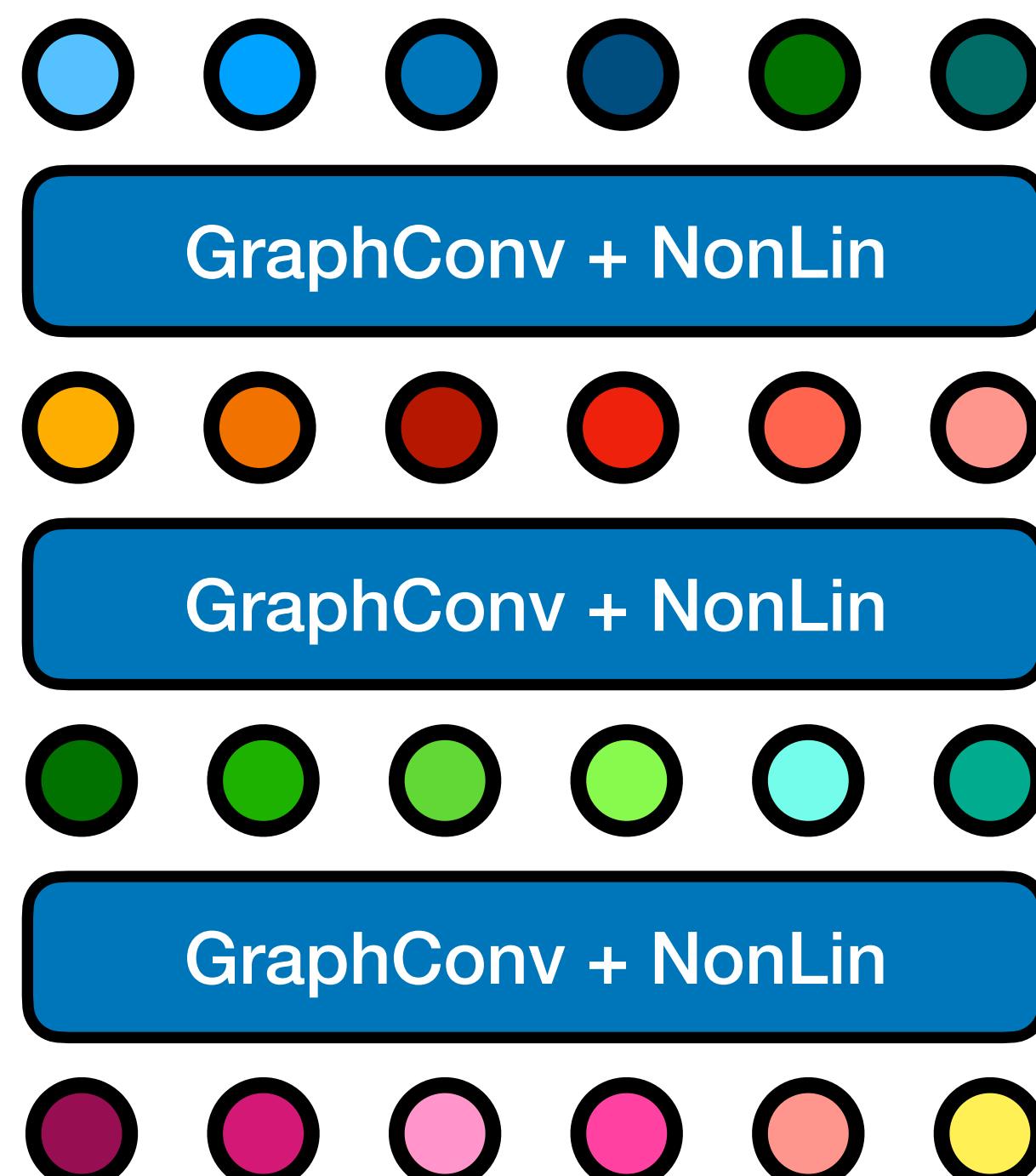
*Graph Conv Layer
(permutation-invariant)*



*ConvNet Layer
(translation-invariant)*

Building Graph Neural Nets

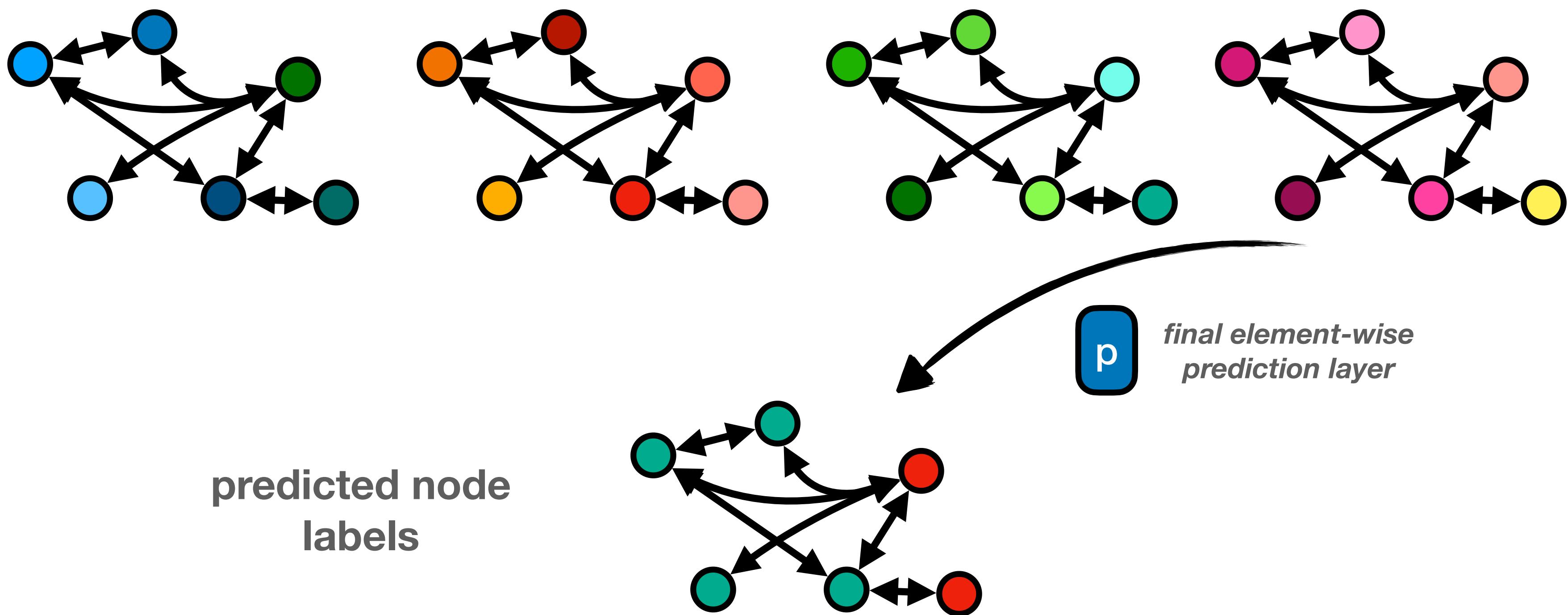
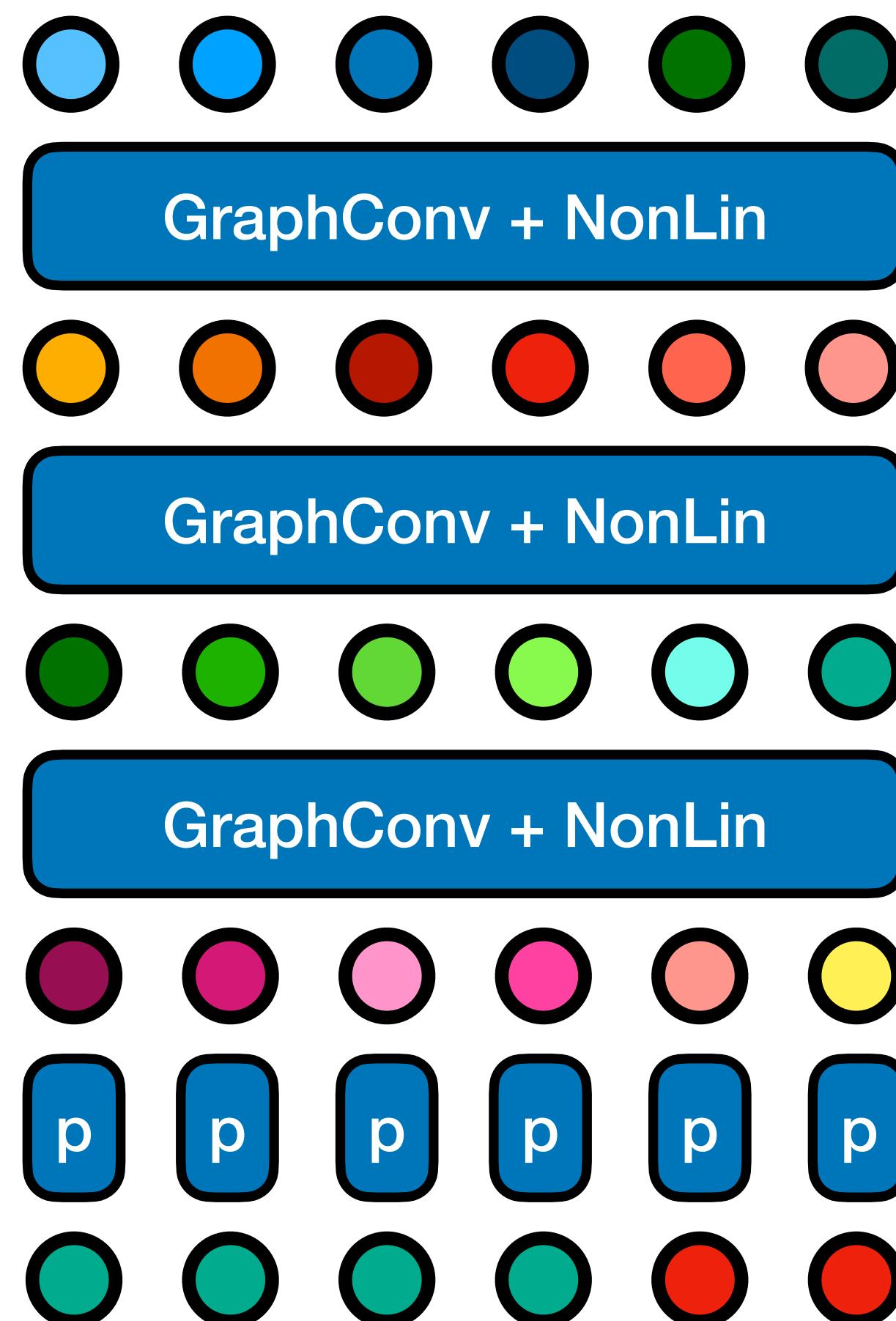
With Graph Convolutions & Non-linearities we can now build deep networks based on these operations



*Note: the graph structure doesn't change
only the node values through multiple rounds
of accumulating signals from the neighborhood*

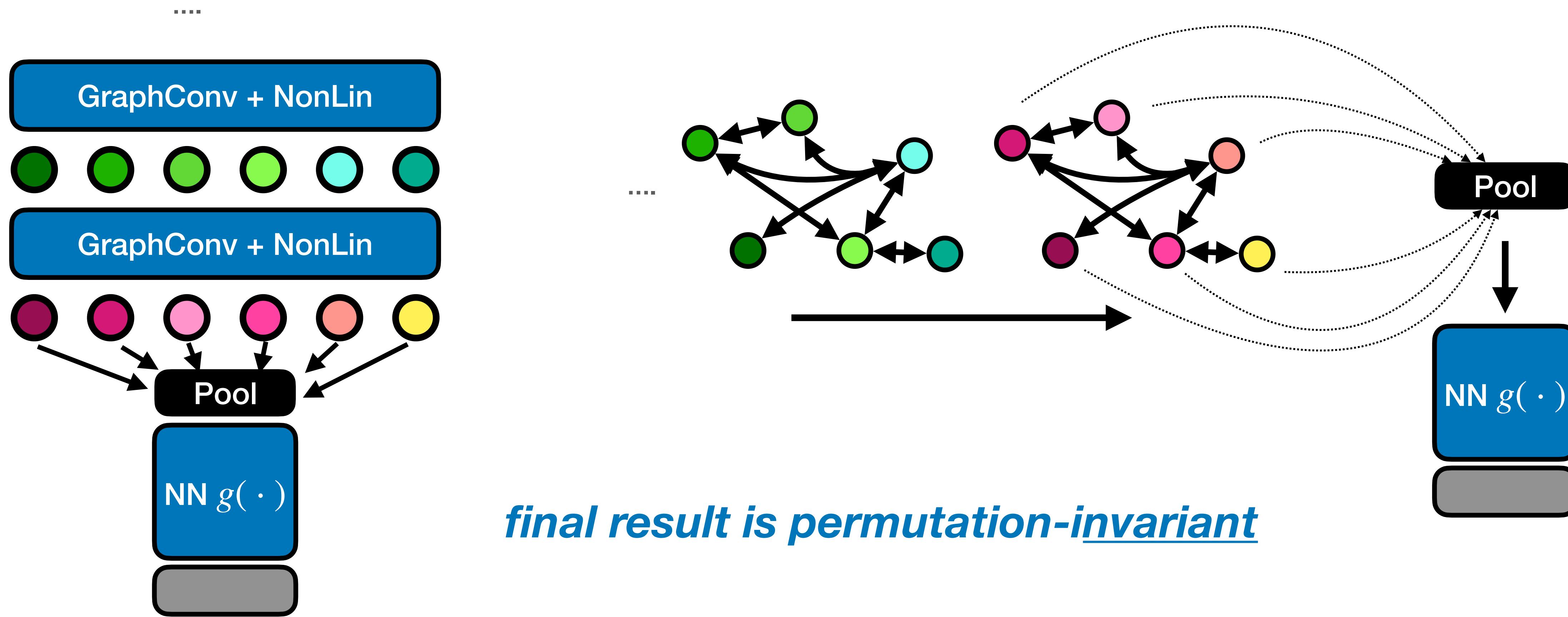
Node Prediction

Once a suitable representation of the data (incl. neighborhood) information is collected, node prediction is elementwise



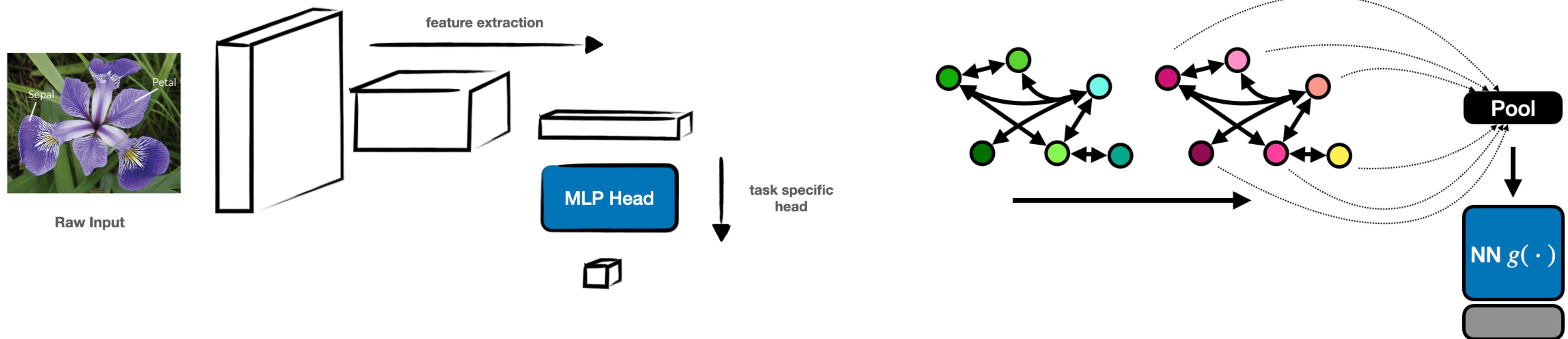
Graph Prediction

Once a suitable representation of the data (incl. neighborhood) information is collected, node prediction is elementwise



CNN v GNN

Same pattern of deep learning: learned data processing (i.e. feature extraction) followed by final high-level task

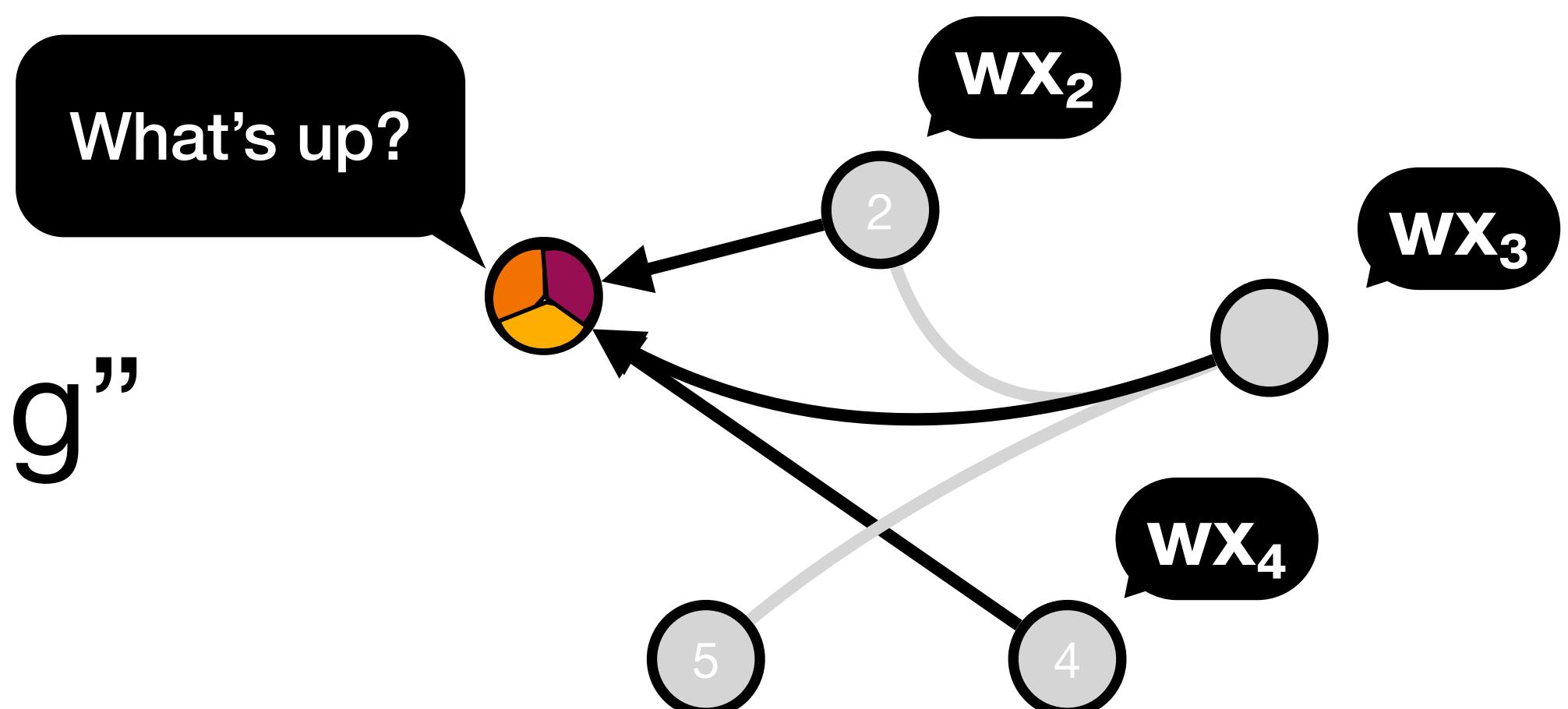


Beyond GCN

We can look at Graph Neural Nets in a more general way

- nodes sending signals (“messages”) to each other
- in graph convolution, messages are fixed by choice of S

GCN are a instance of a more general paradigm of “message passing” graph neural networks



Beyond GCN

General Setup involves three components

**message
creation**

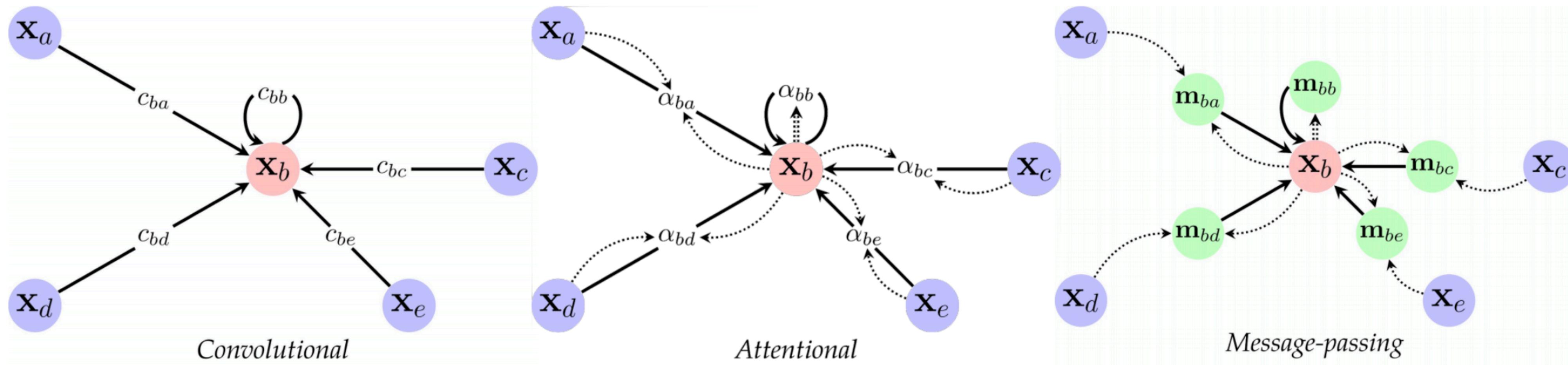
**permutation invariant
pooling of messages**

Non-Linearity

$$\sigma(x_i, \text{Pool}_{j \in N}(m(x_i, x_j)))$$

Beyond GCN

In the literature you'll find entire families of GNN types, depending on choice of message, pooling, non-linearity



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

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

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

Summary

Graph Neural Networks are a very broad generalization of ideas we had in CNN to “structured data”

Processes signals “on top of a graph” which is assumed to be given externally by exchanging signals across edges

Corrolary: the creation of the graph is of utmost importance as it encodes our domain knowledge (“who talks to whom”)

Some GNN do mutate the graph, but often it’s fixed.

End.