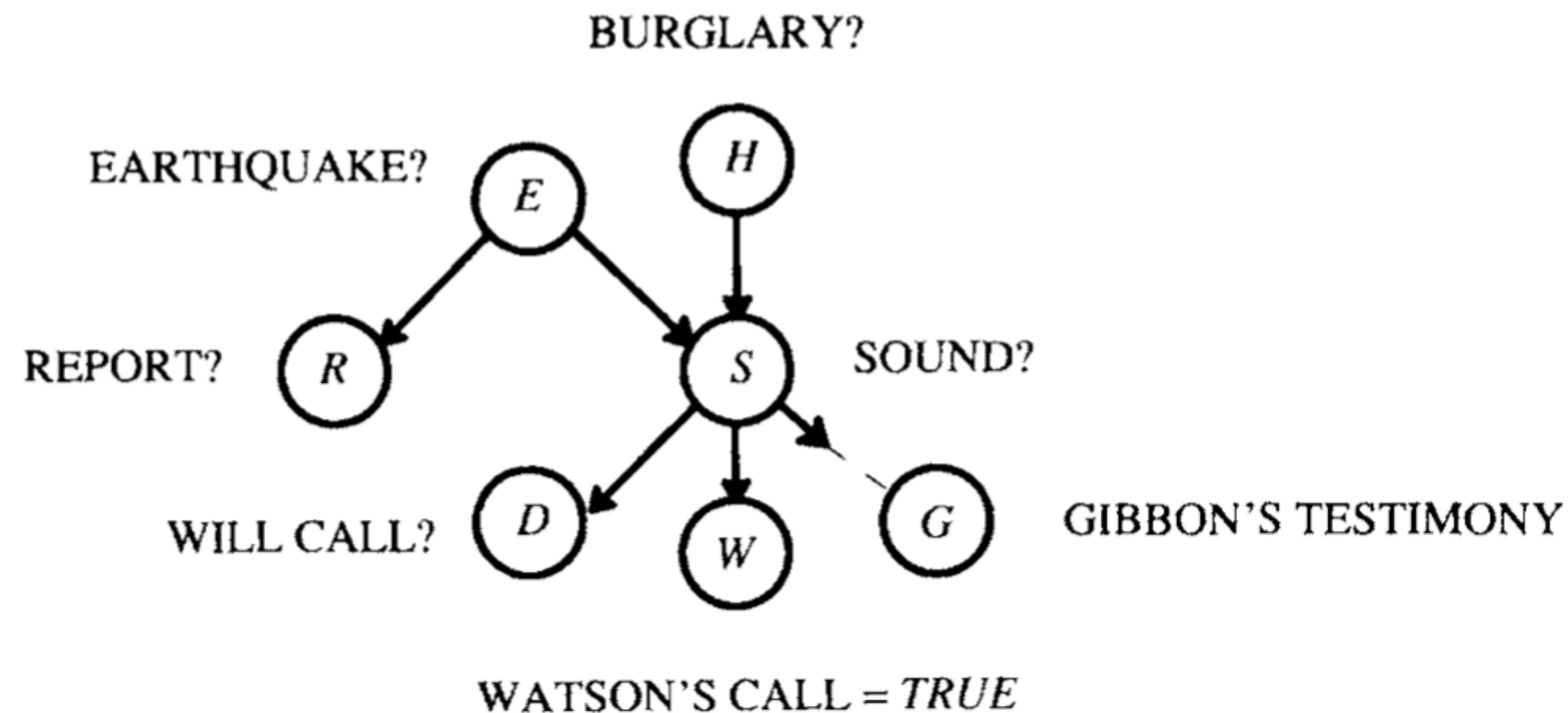


# Graph Neural Networks

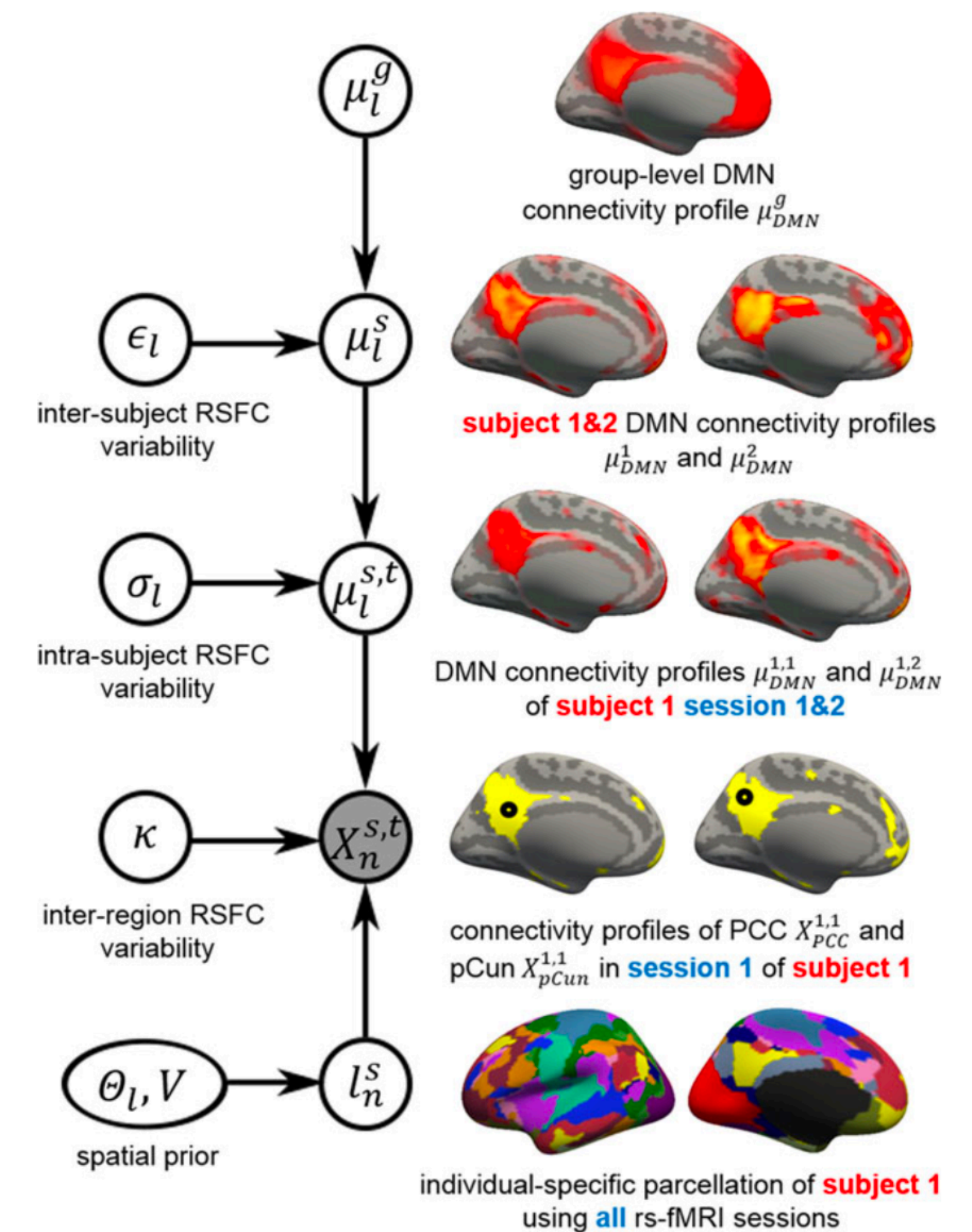
Demian Wassermann, Inria

Graphical Models: Discrete Inference and Learning

# Introduction to DAG and their relationship with Probability Functions (Pearl)



[Pearl 1987]



[Kong et al 2019]



# And the Usual Graph Slide



Image credit: [Medium](#)

**Social Networks**

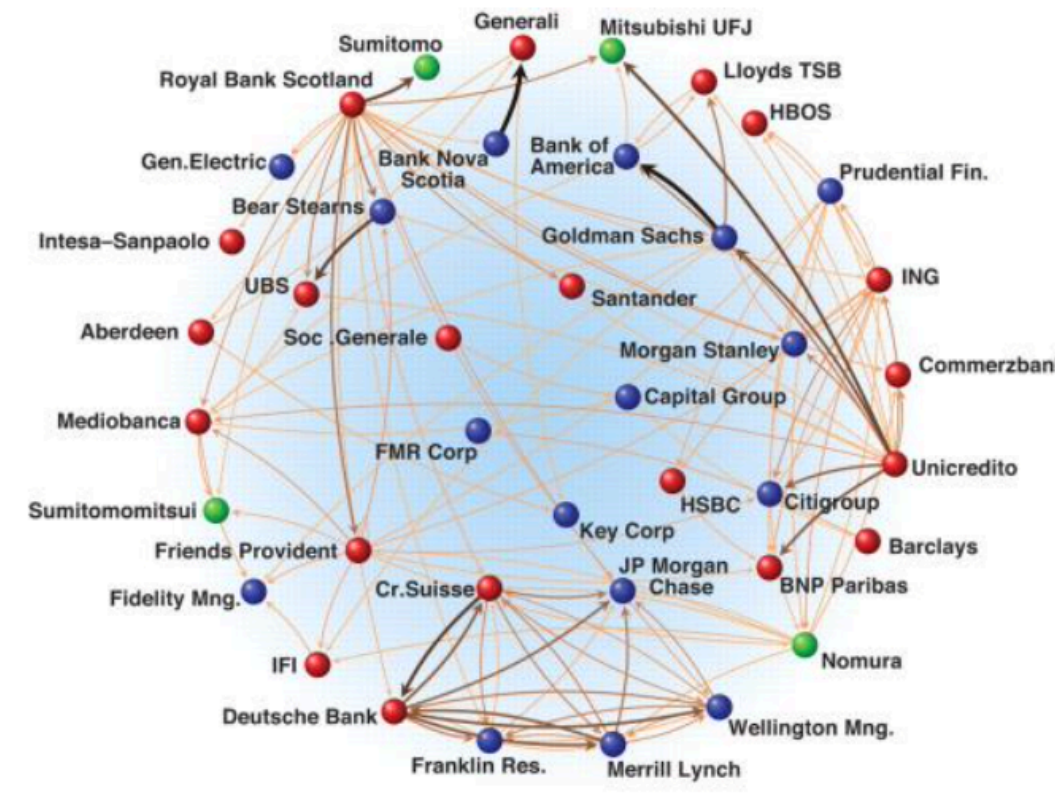


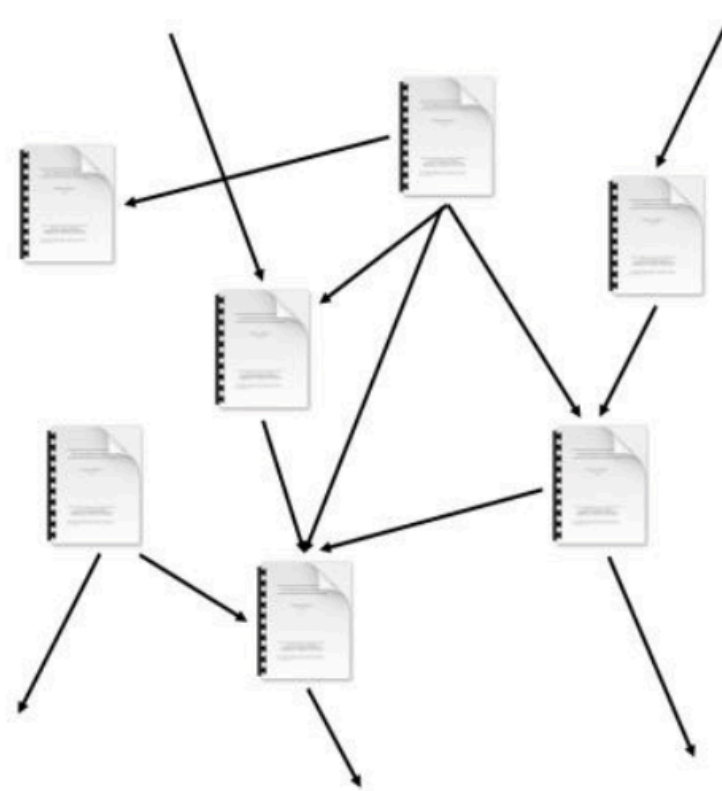
Image credit: [Science](#)

**Economic Networks**



Image credit: [Lumen Learning](#)

**Communication Networks**



**Citation Networks**



Image credit: [Missoula Current News](#)

**Internet**

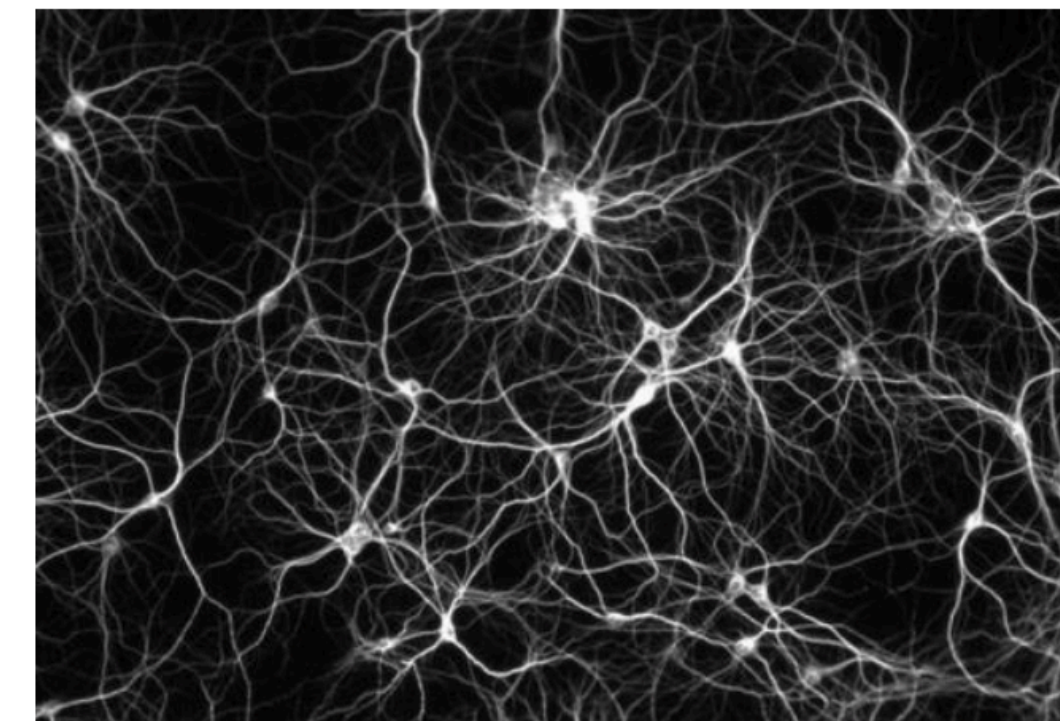
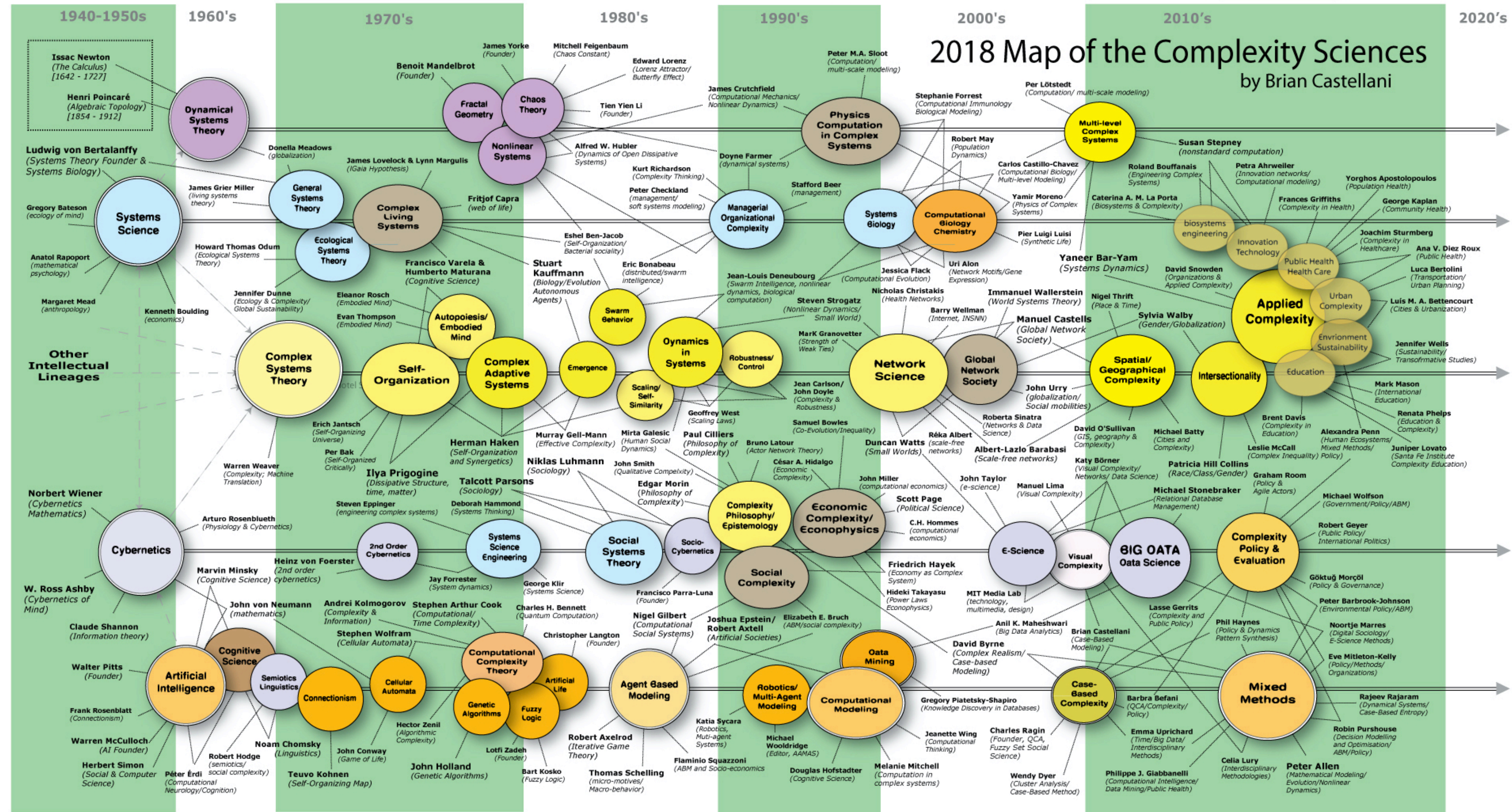


Image credit: [The Conversation](#)

**Networks of Neurons**



# Complex Systems to Understand the World





# Main Epistemological Angles on Graphs and Knowledge

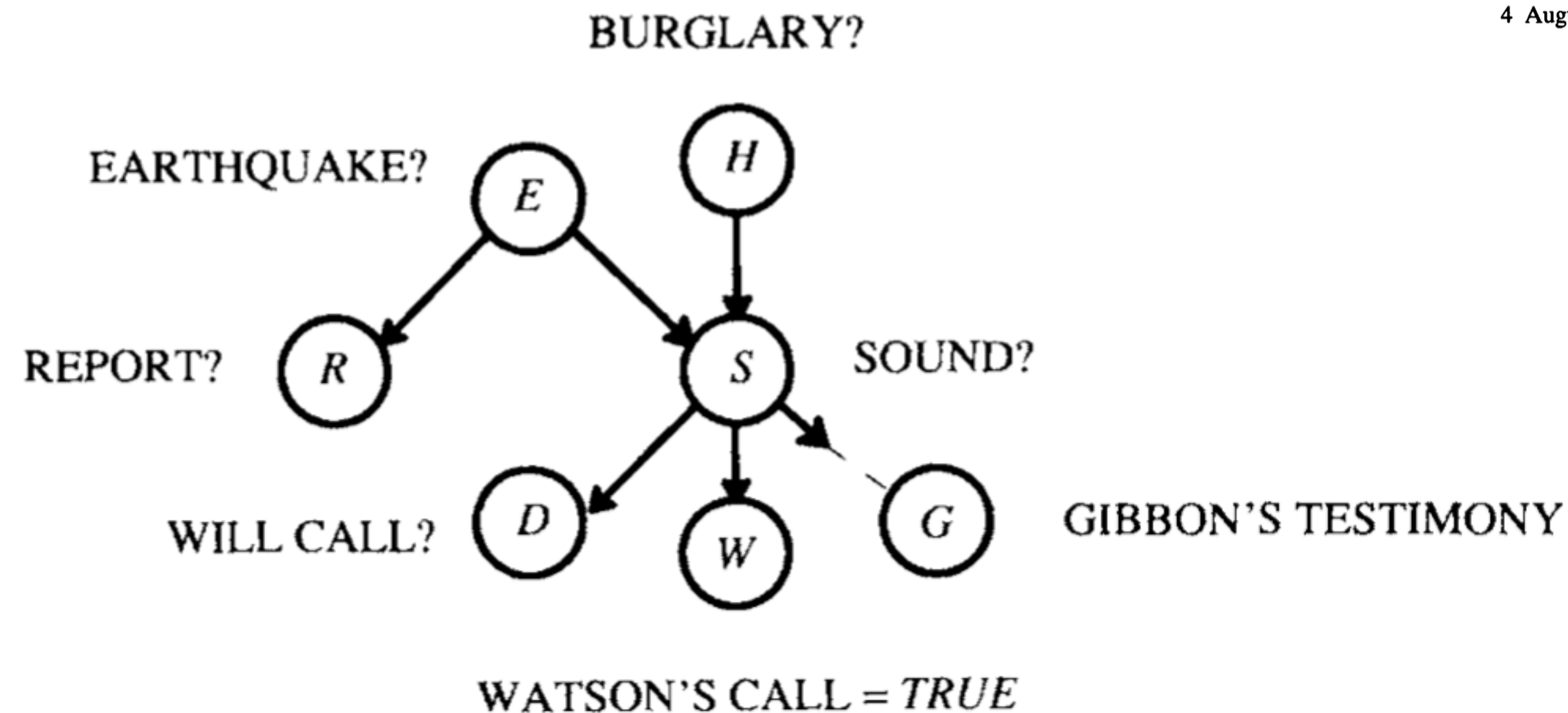
4 August 1972, Volume 177, Number 4047

**SCIENCE**

## More Is Different

Broken symmetry and the nature of the hierarchical structure of science.

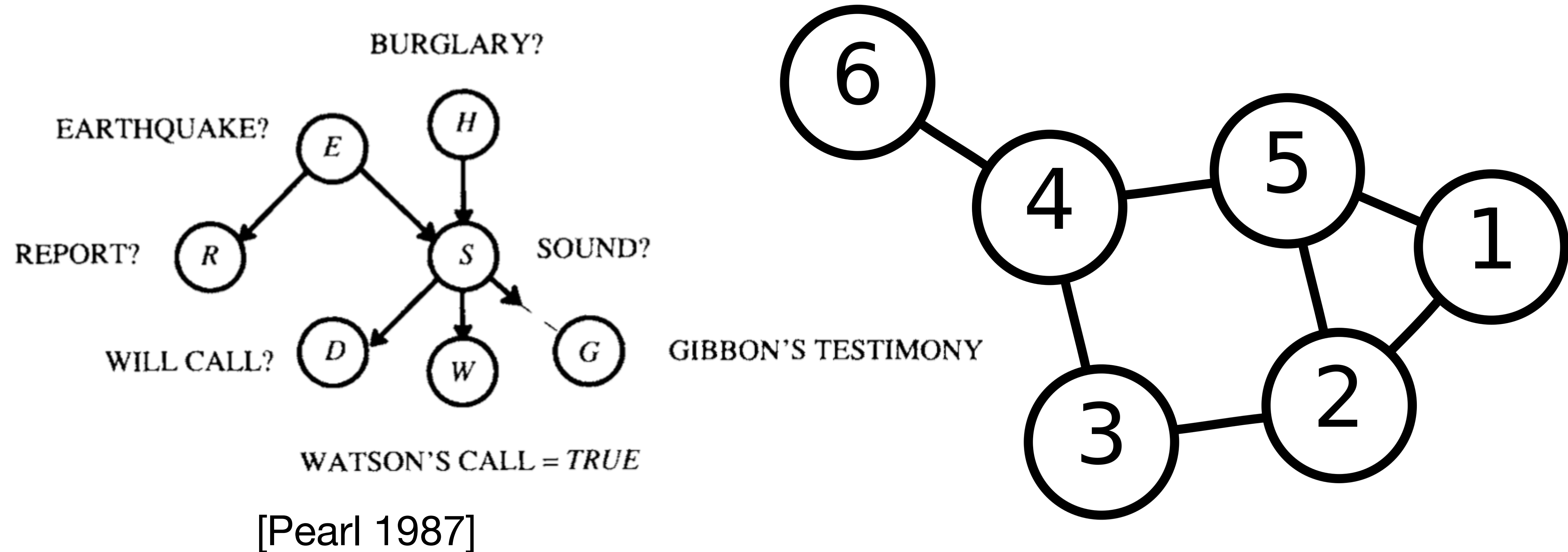
P. W. Anderson



[Pearl 1987]



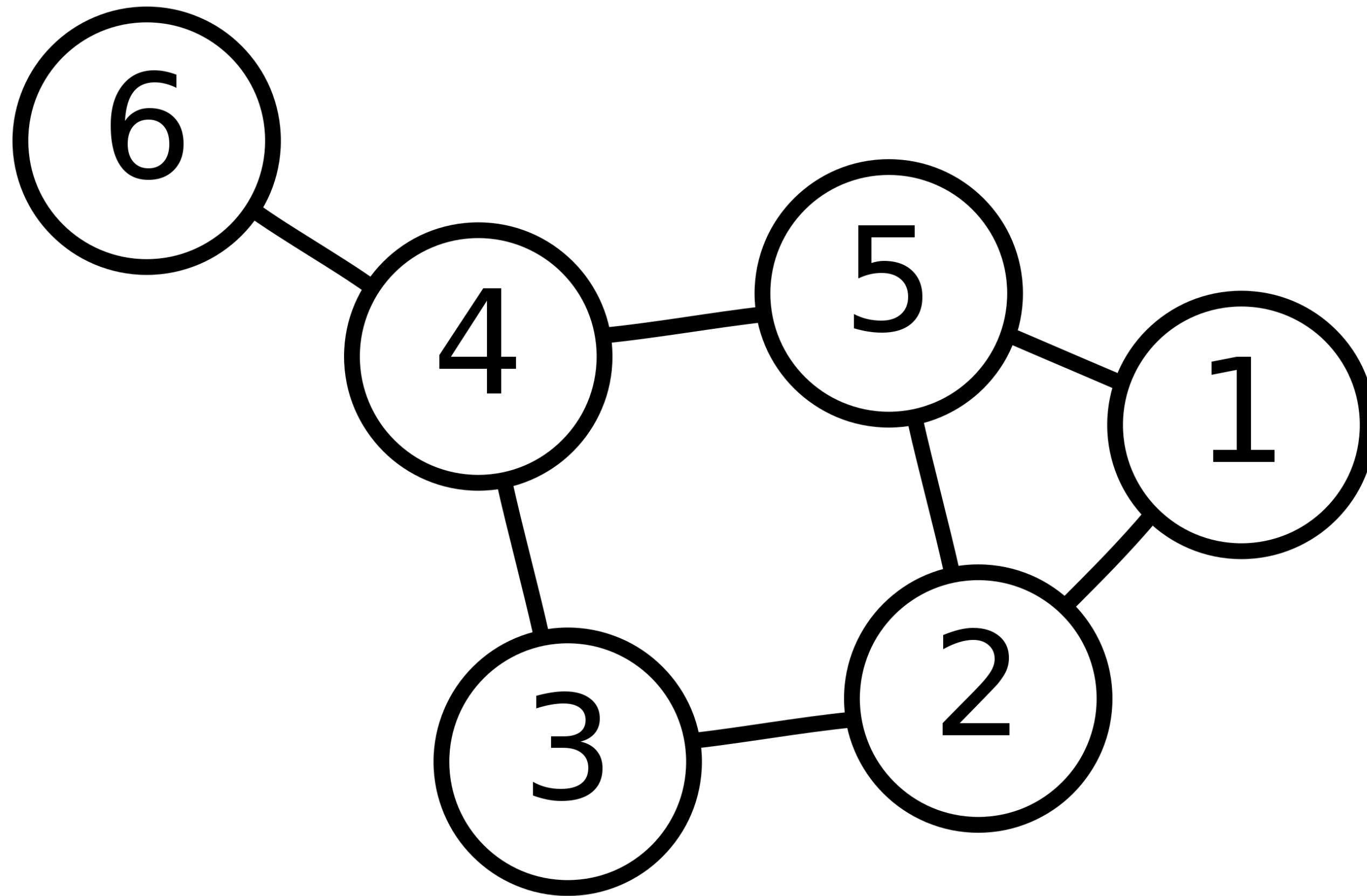
# Main Epistemological Angles on Graphs and Knowledge





# Alternative General Undirected Graphs

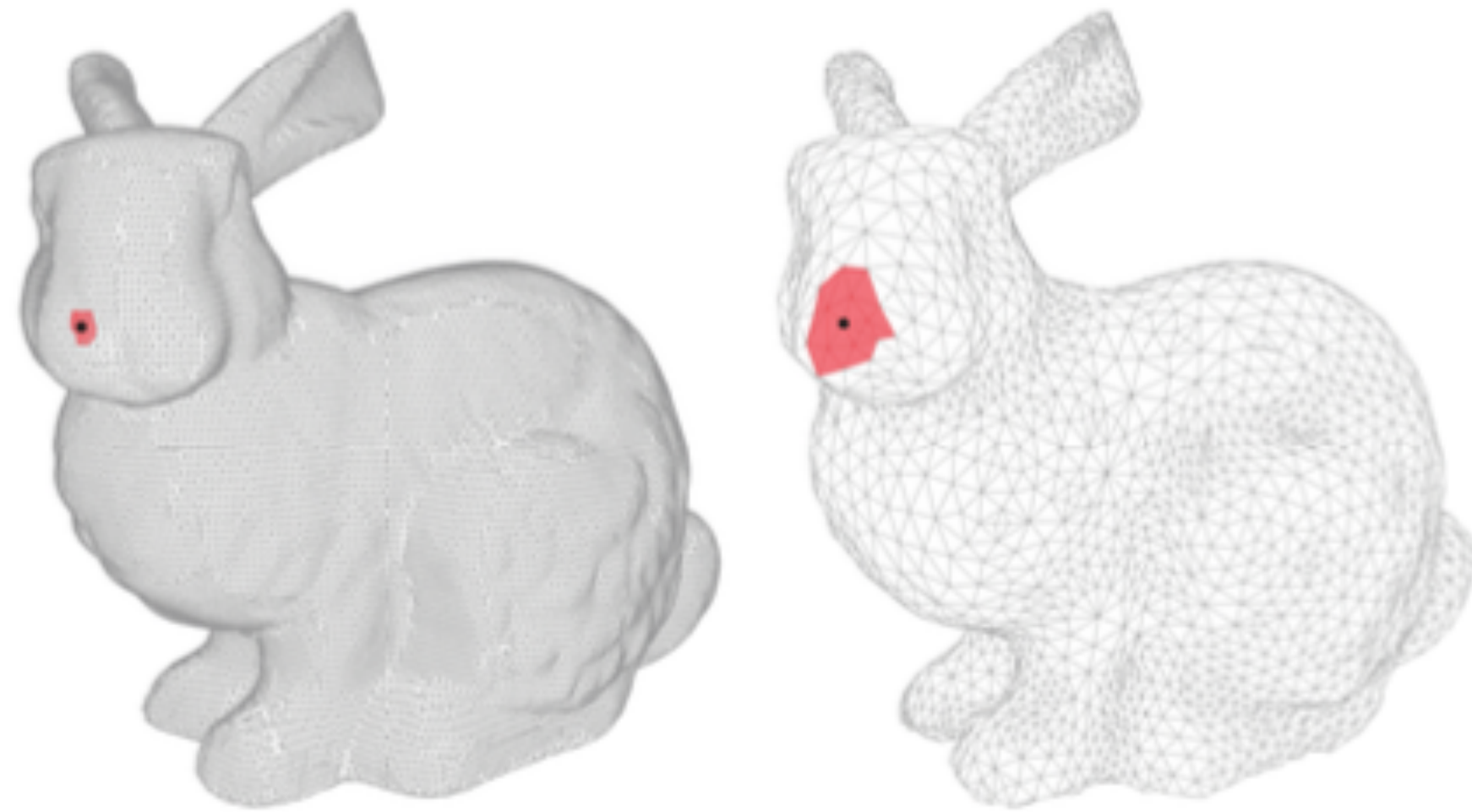
- Graph  $G(V, E, f, g), E \subseteq V \times V, f: V \mapsto F, g: E \mapsto G$



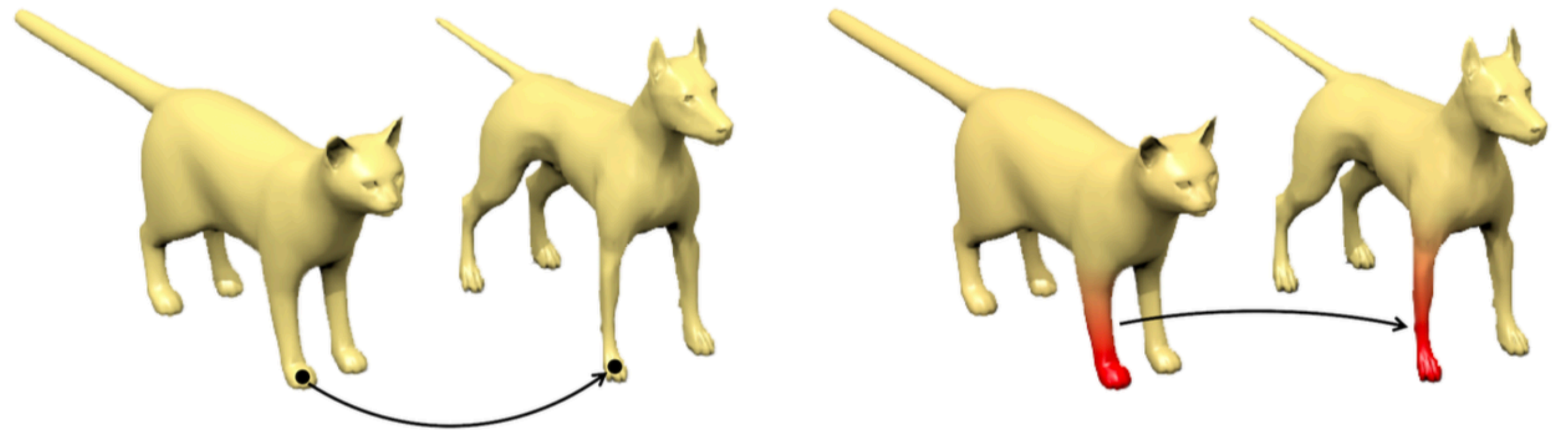
- $V$  are the vertices
- $E$  are the edges
- $f$  is a mapping of features for vertices
- $g$  is a mapping of features for edges



# Structure and Parameterization are Important



Two Different Resolutions



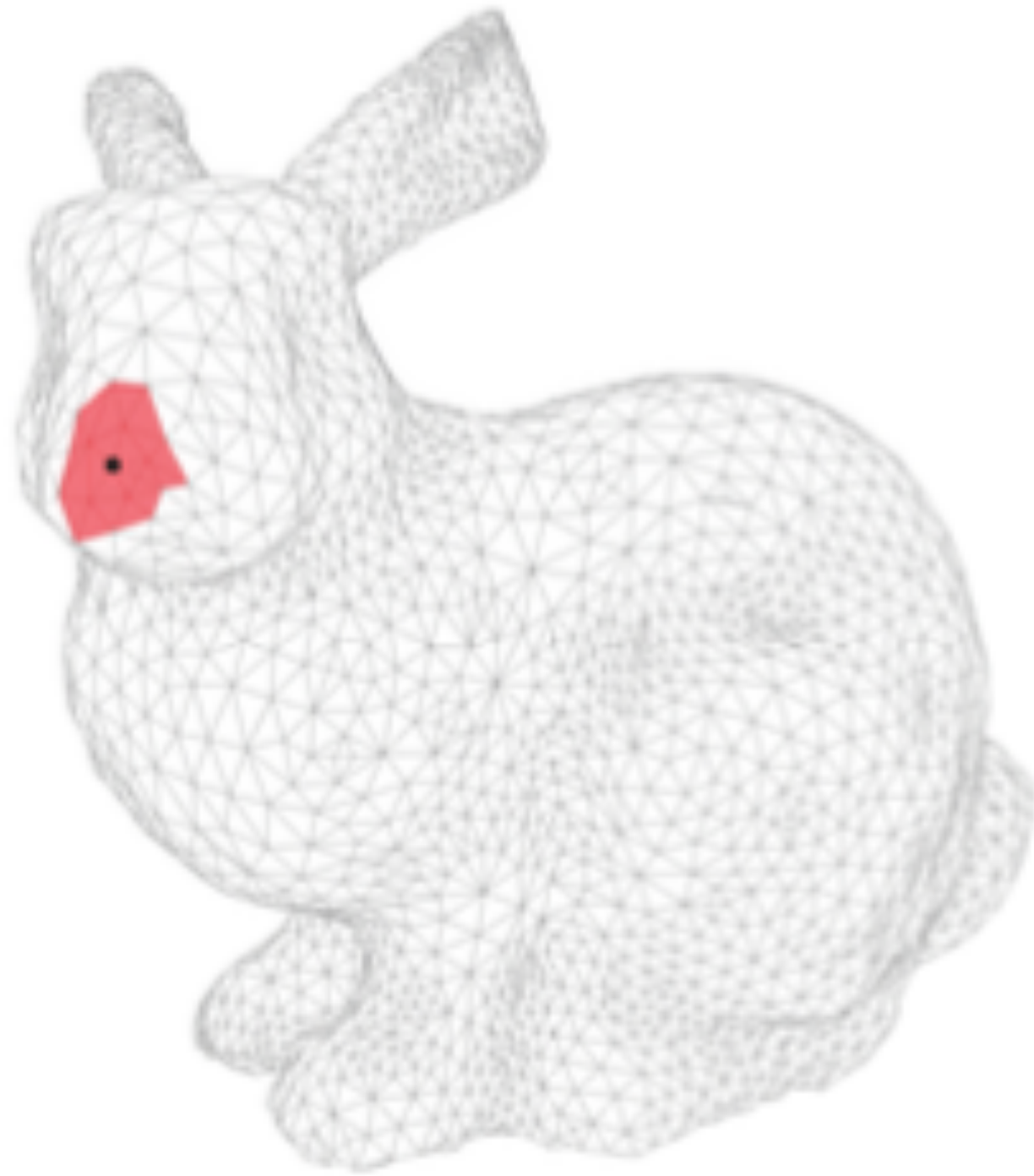
Point-wise vs Functional Equivalence



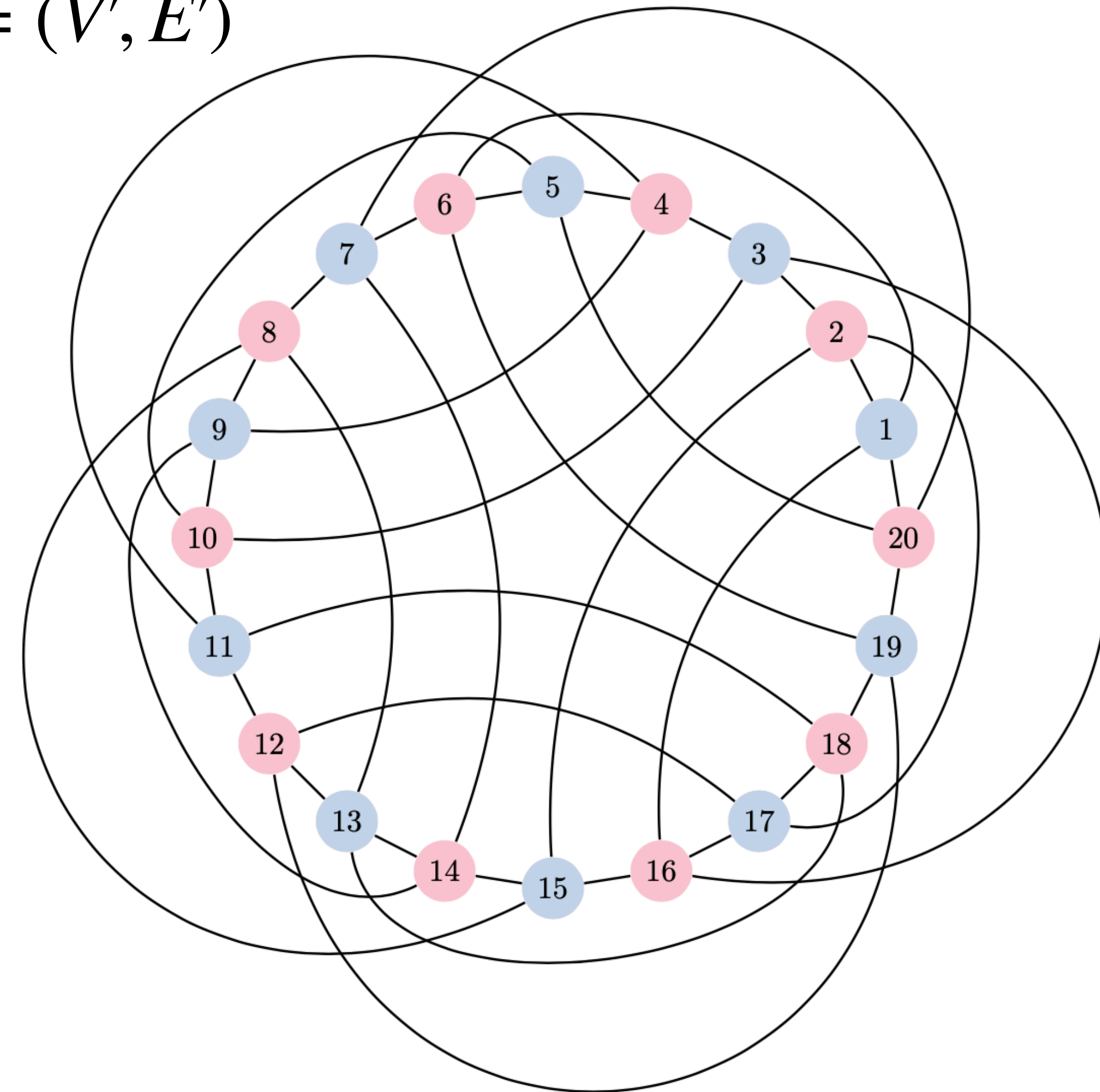
# Automorphisms

$G = (V, E)$  automorphic to  $G' = (V', E')$  iif

$\exists \phi : V \mapsto V'$  s.t.  $(\phi(V), \{ \phi(u) \rightarrow \phi(v) : u \rightarrow v \in E \}) = (V', E')$



Folkman graph  
3840 automorphisms

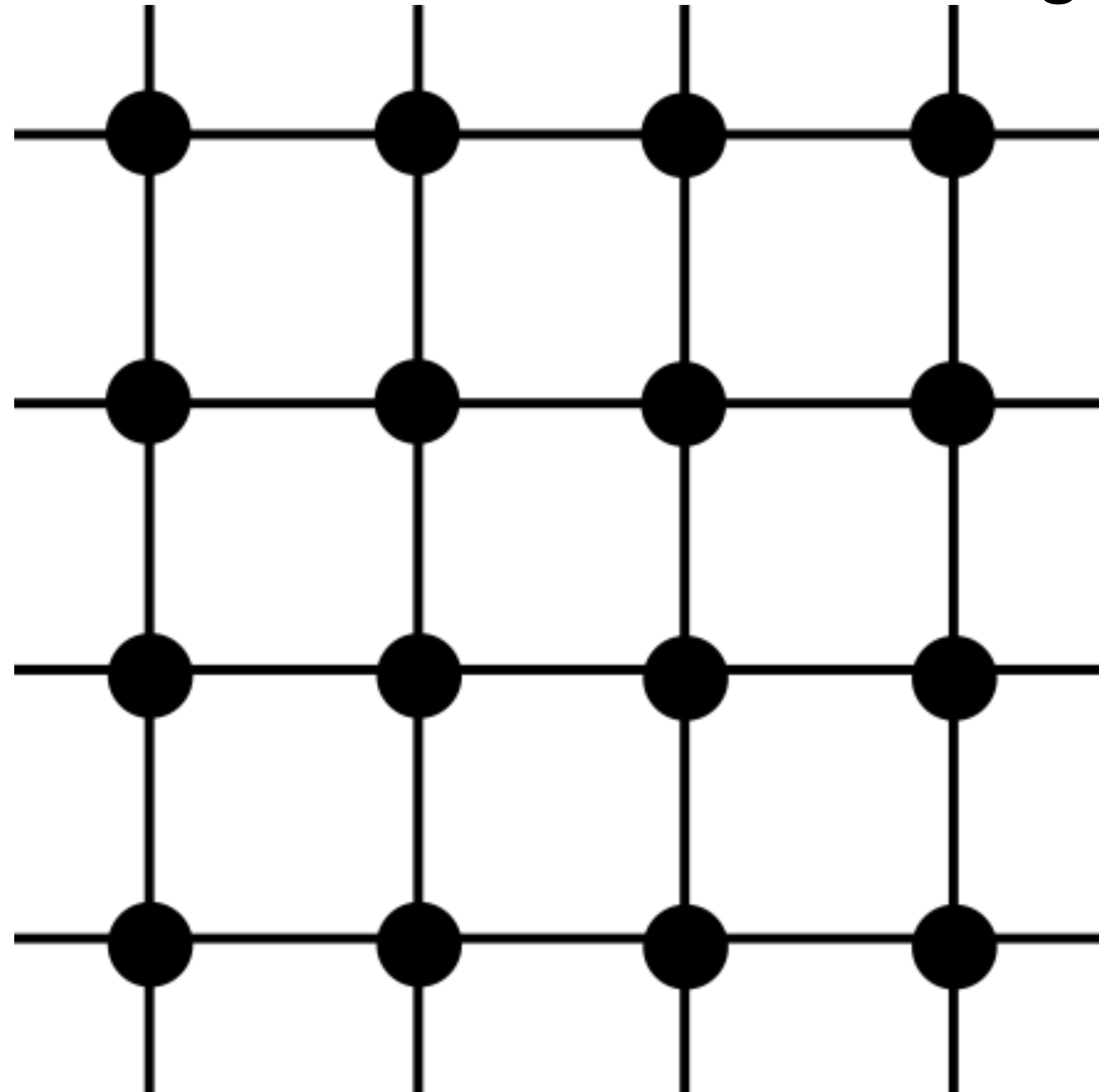




# Alternative General Undirected Graphs

- Graph  $G(V, E, f, g), E \subseteq V \times V, f: V \mapsto F, g: E \mapsto G$

The Lattice Case, or an Image



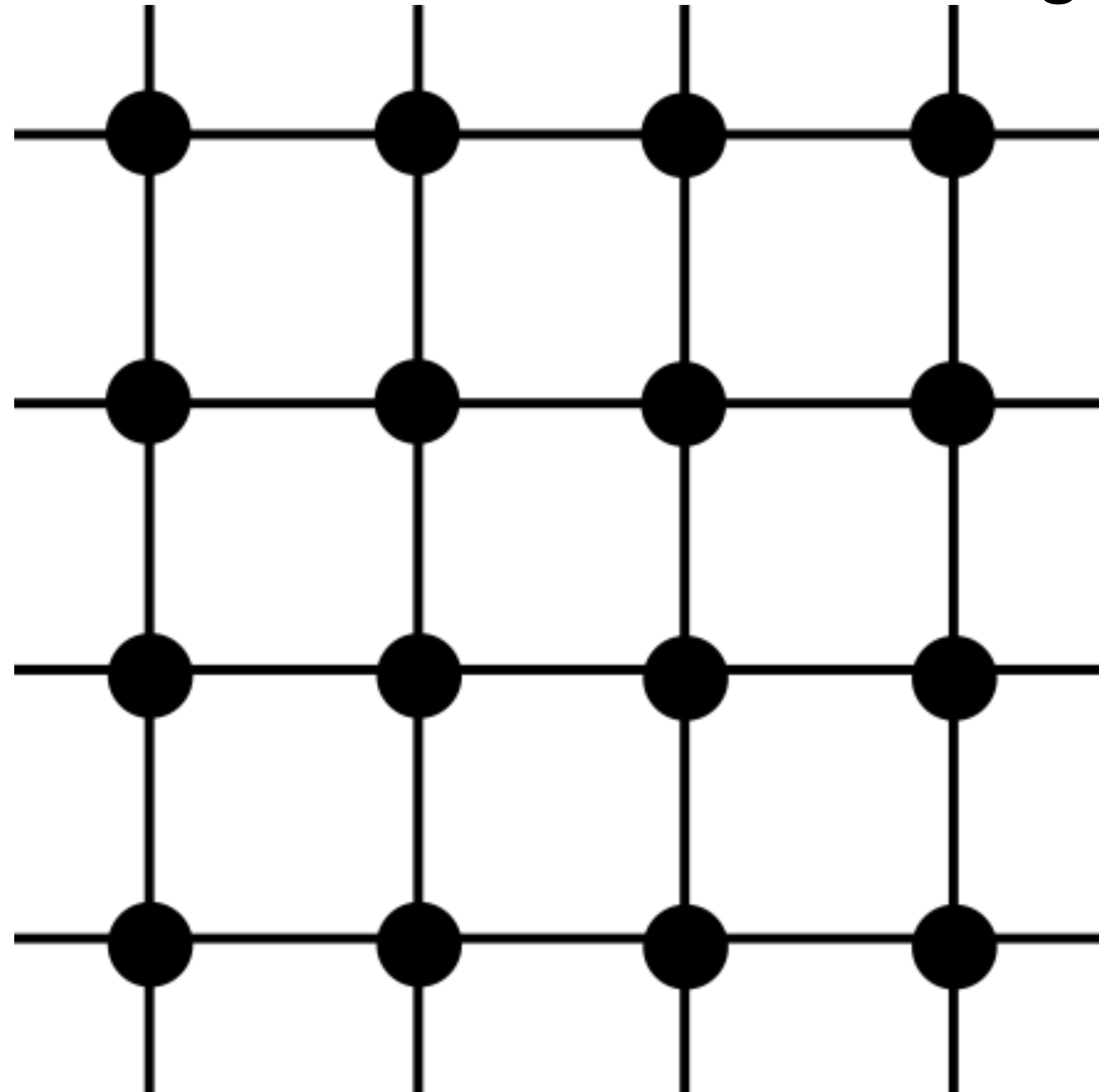
- $V$  are the vertices
- $E$  are the edges
- $f$  is a mapping of features for vertices
- $g$  is a mapping of features for edges



# Alternative General Undirected Graphs

- Graph  $G(V, E, f, g), E \subseteq V \times V, f: V \mapsto F, g: E \mapsto G$

The Lattice Case, or an Image



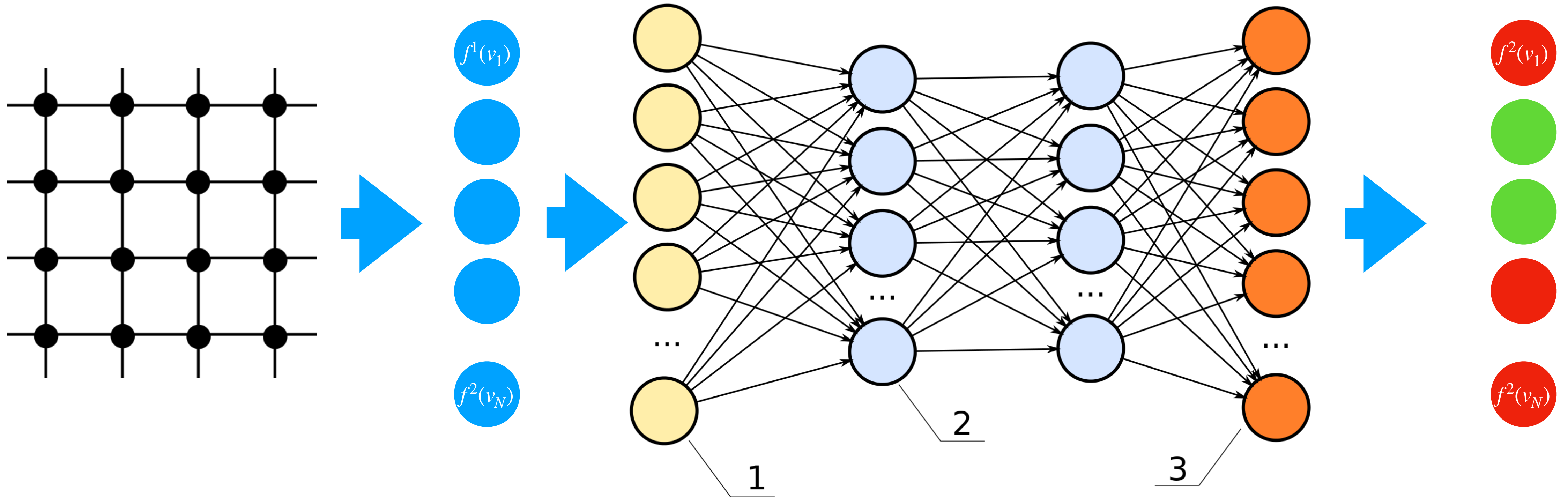
- $V$  are the vertices
- $E$  are the edges
- $f$  is a mapping of features for vertices
- Learn  
$$\arg \min_{\theta} \mathcal{L}([f(v_i)]_i, [\phi_{\theta}(v_i)]_i)$$
  
for a fixed ordering  $i$



# Alternative General Undirected Graphs

- Learn  $\arg \min_{\theta} \mathcal{L}([f^2(v_i)]_i, [\phi_{\theta}(f^1(v_i))]_i)$  for a fixed ordering  $i$

The Lattice Case, or an Image

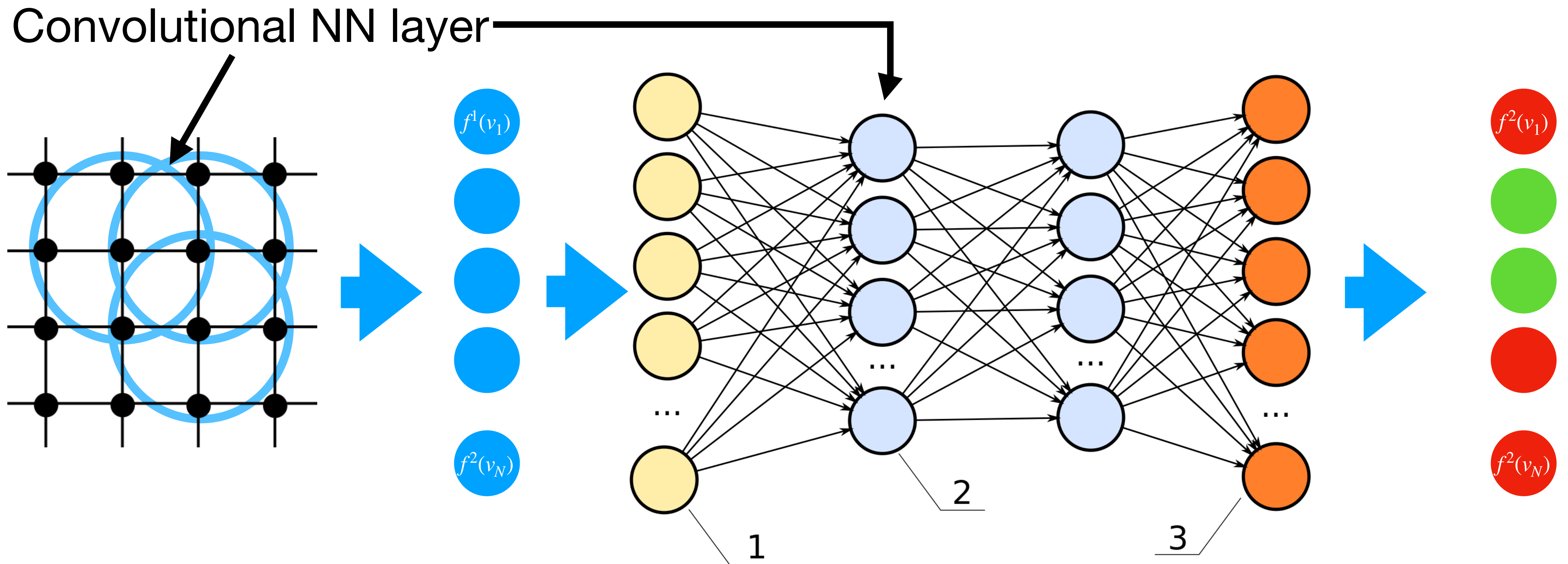




# Alternative General Undirected Graphs

- Learn  $\arg \min_{\theta} \mathcal{L}([f^2(v_i)]_i, [\phi_{\theta}(f^1(v_i))]_i)$  for a fixed ordering  $i$

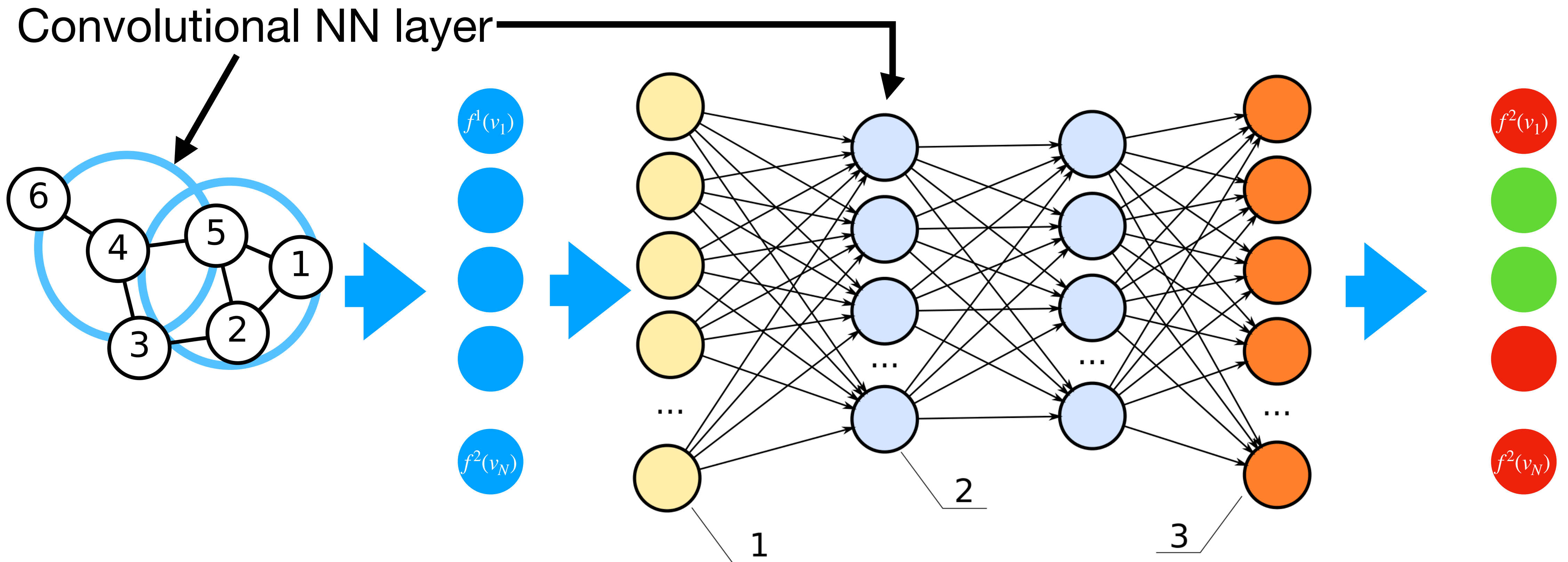
The Lattice Case, or an Image



# Alternative General Undirected Graphs

- Learn  $\arg \min_{\theta} \mathcal{L}([f^2(v_i)]_i, [\phi_{\theta}(f^1(v_i))]_i)$  for a fixed **structure**

The Graph Case, and now?

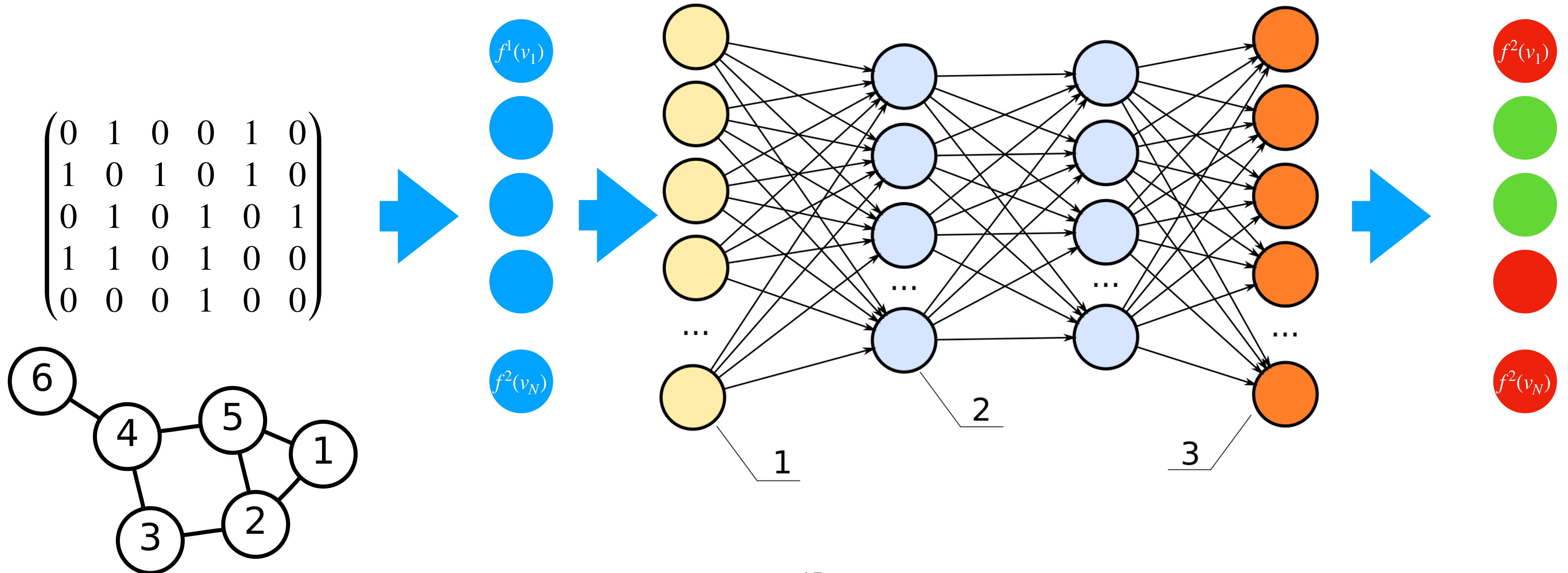




# Alternative General Undirected Graphs

- Learn  $\arg \min_{\theta} \mathcal{L}([f^2(v_i)]_i, [\phi_{\theta}(f^1(v_i))]_i)$  for a fixed **structure**

The Graph Case, use the affinity matrix



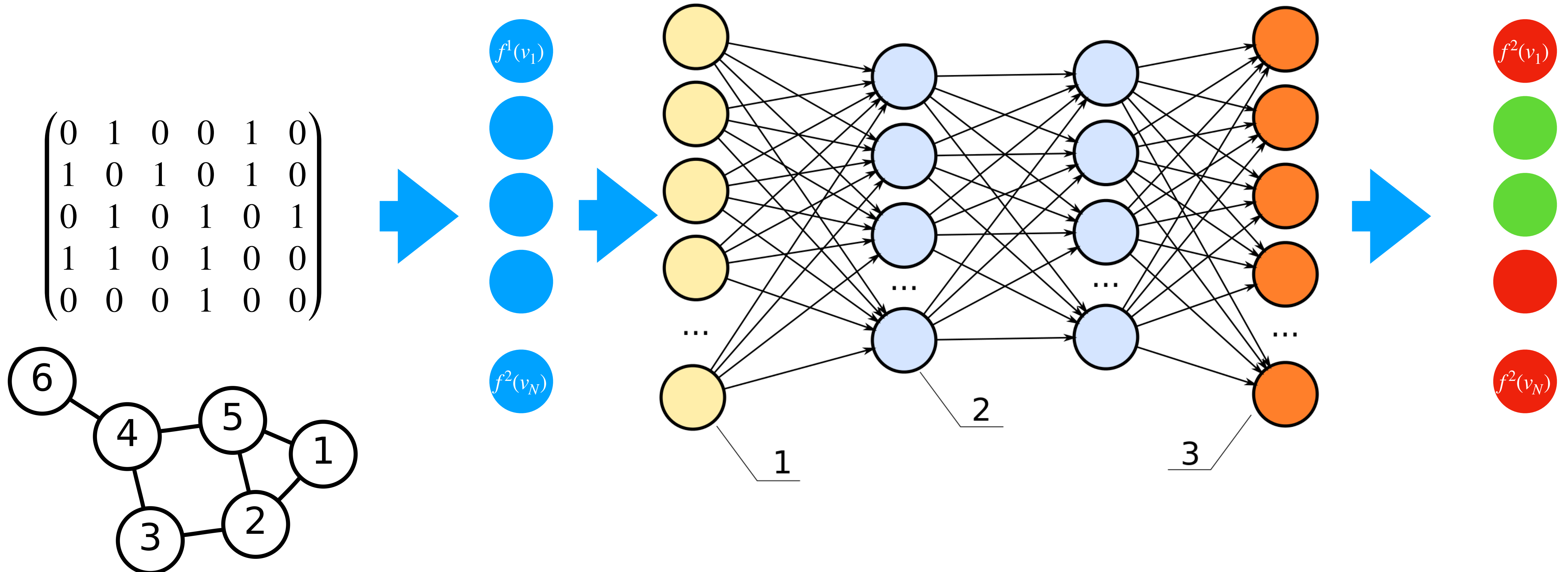
# Alternative General Undirected Graphs

- Learn  $\arg \min_{\theta} \mathcal{L}([f^2(v_i)]_i, [\phi_{\theta}(f^1(v_i))]_i)$  for a fixed **structure**

The Graph Case, use the affinity matrix

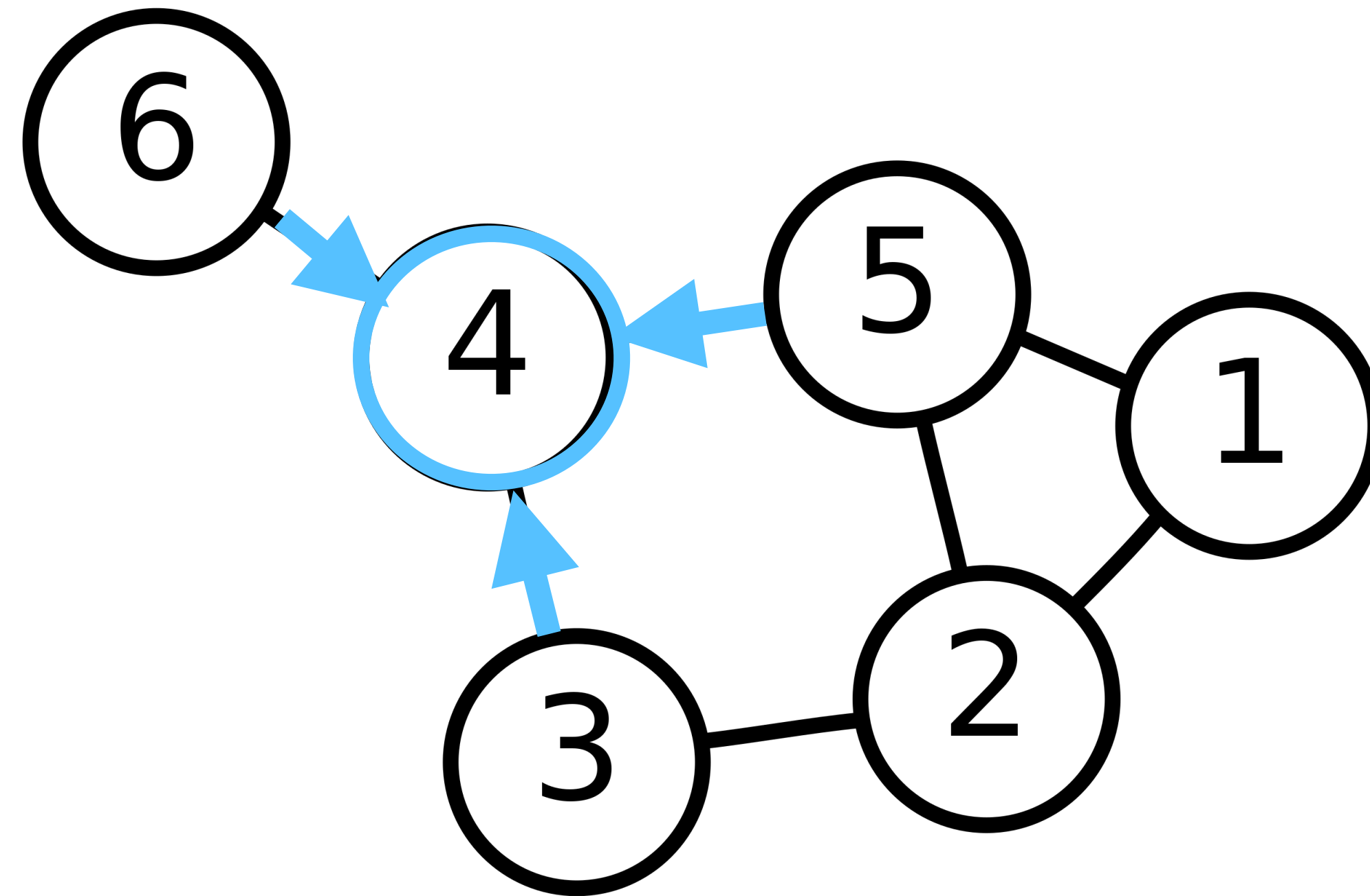
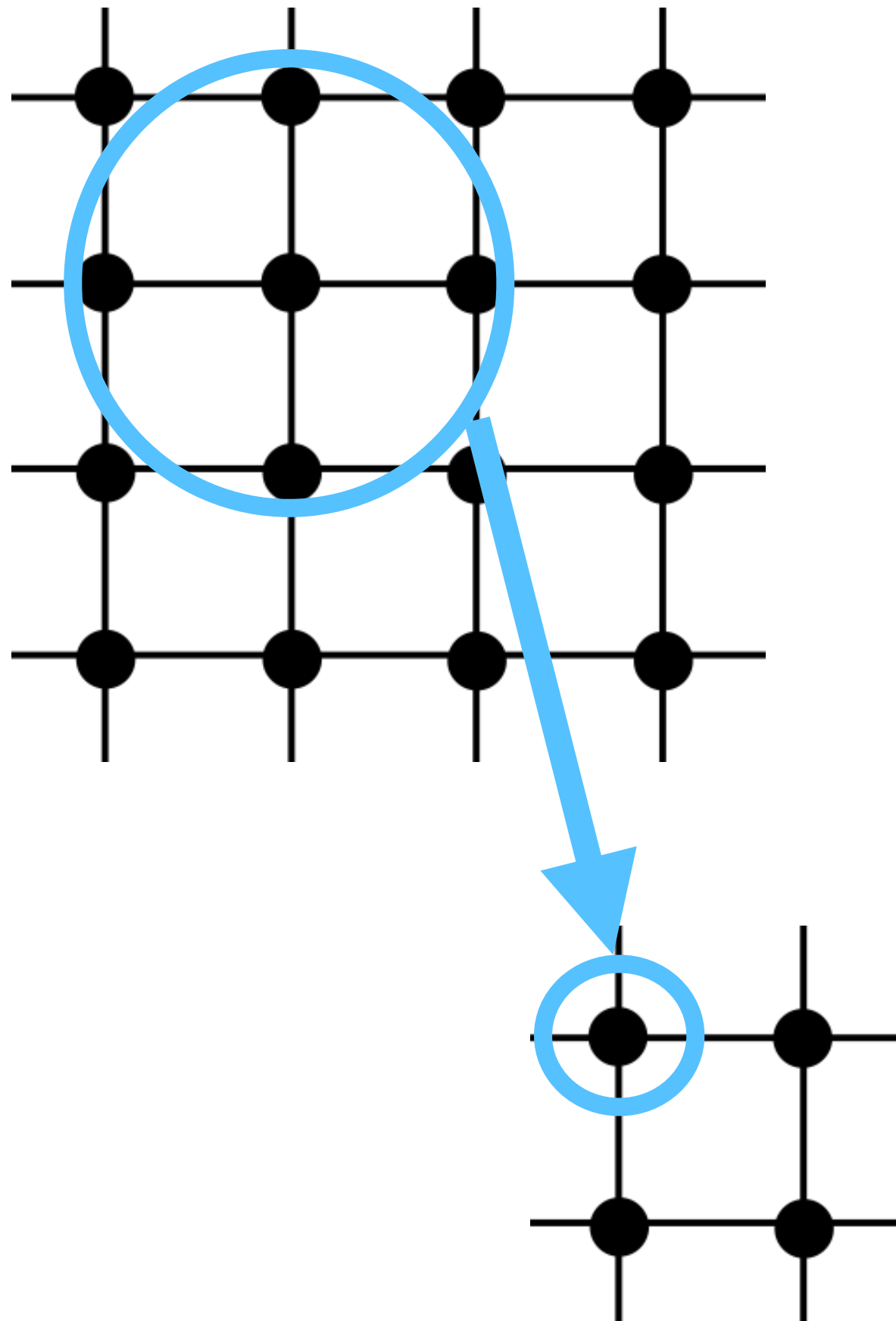
## Issues

- $O(|V|)$  Parameters
- Not applicable different sizes
- Sensitive to node ordering

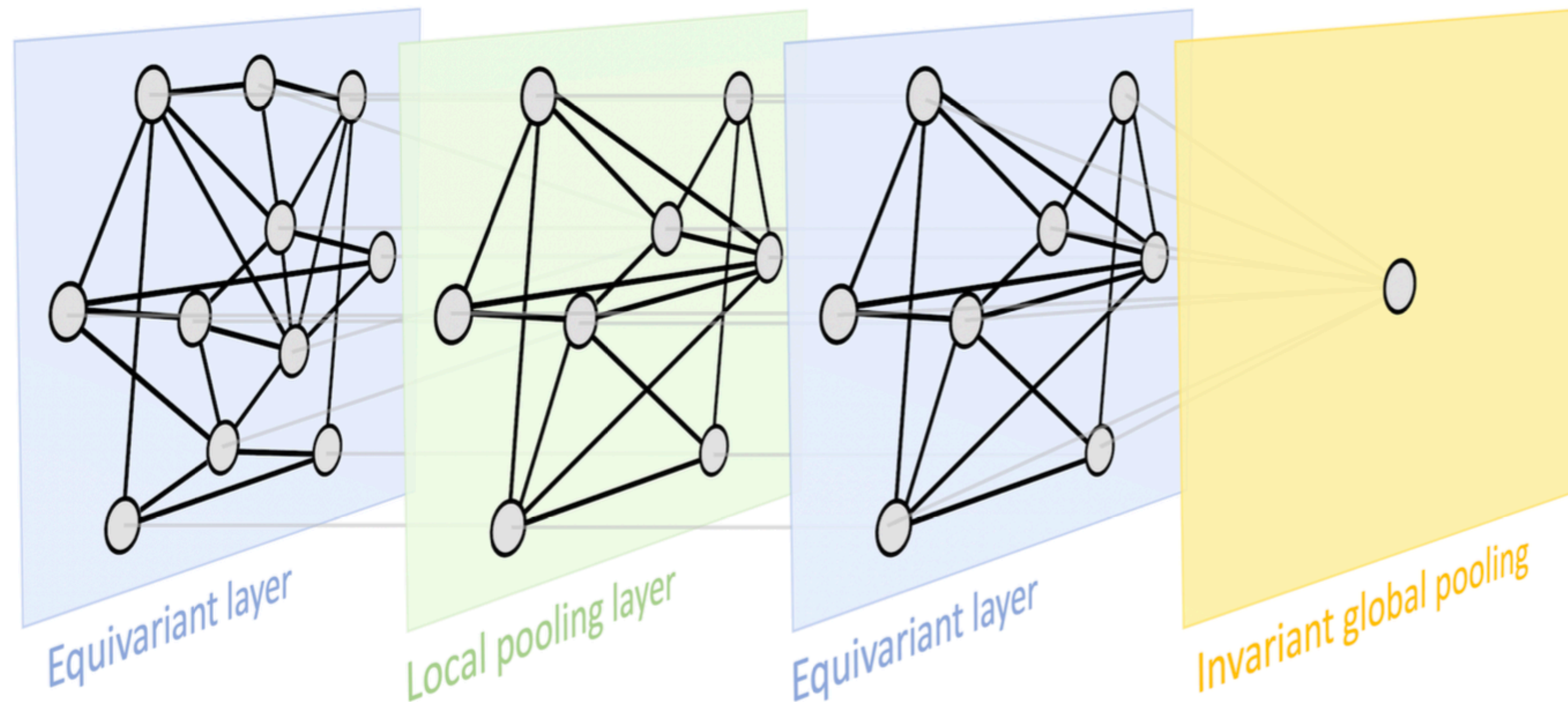




# Convolutions on Graphs: Message Passing

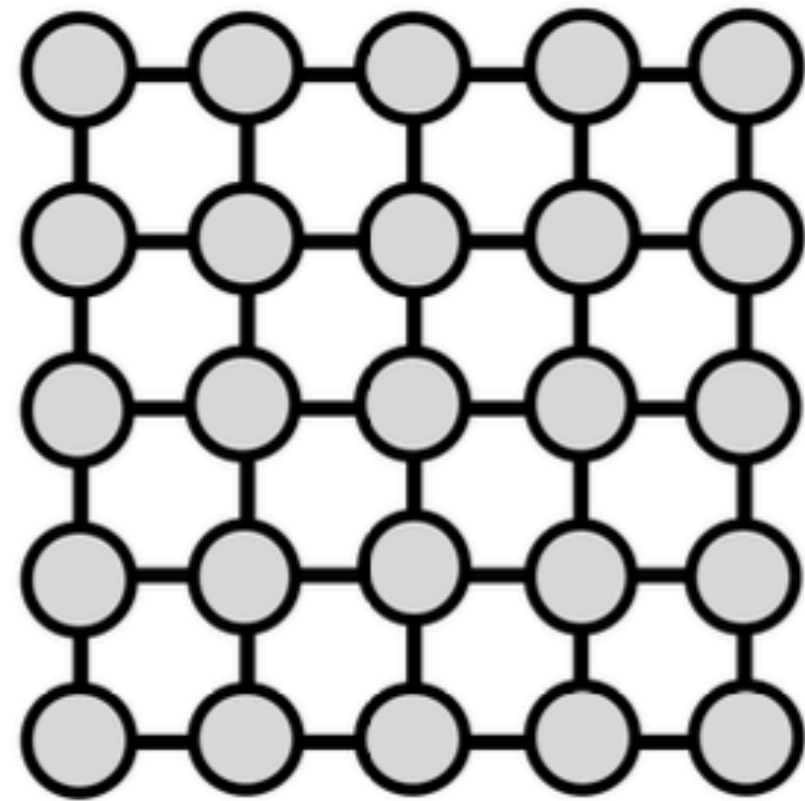


# The Blueprint (Bronstein et al 2023)

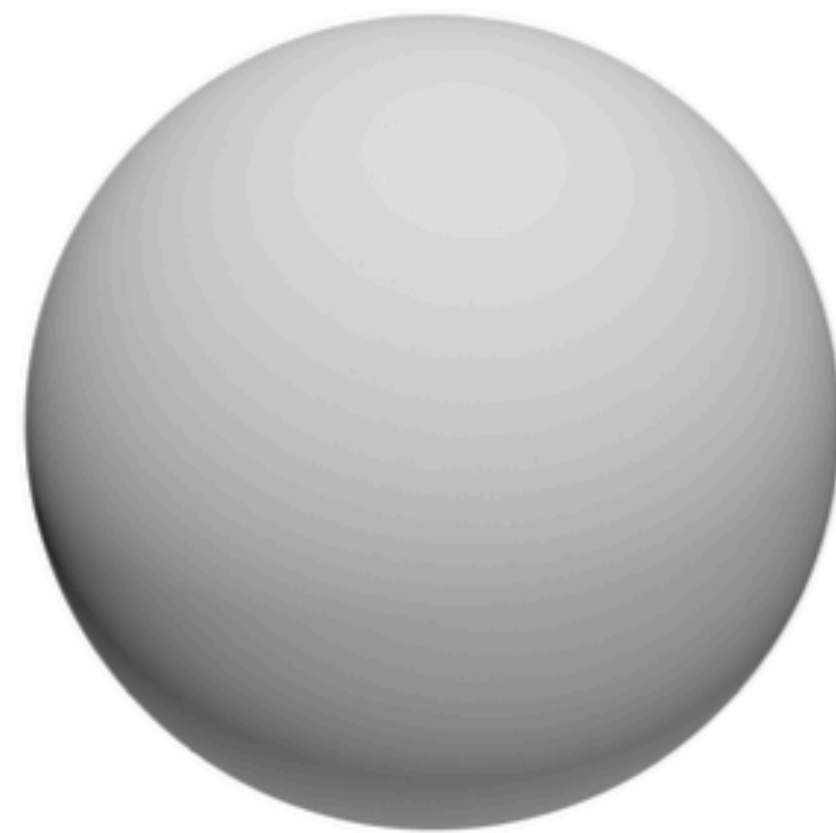




# The 5G of Geometric/Graph DL



**Grids**



**Groups**

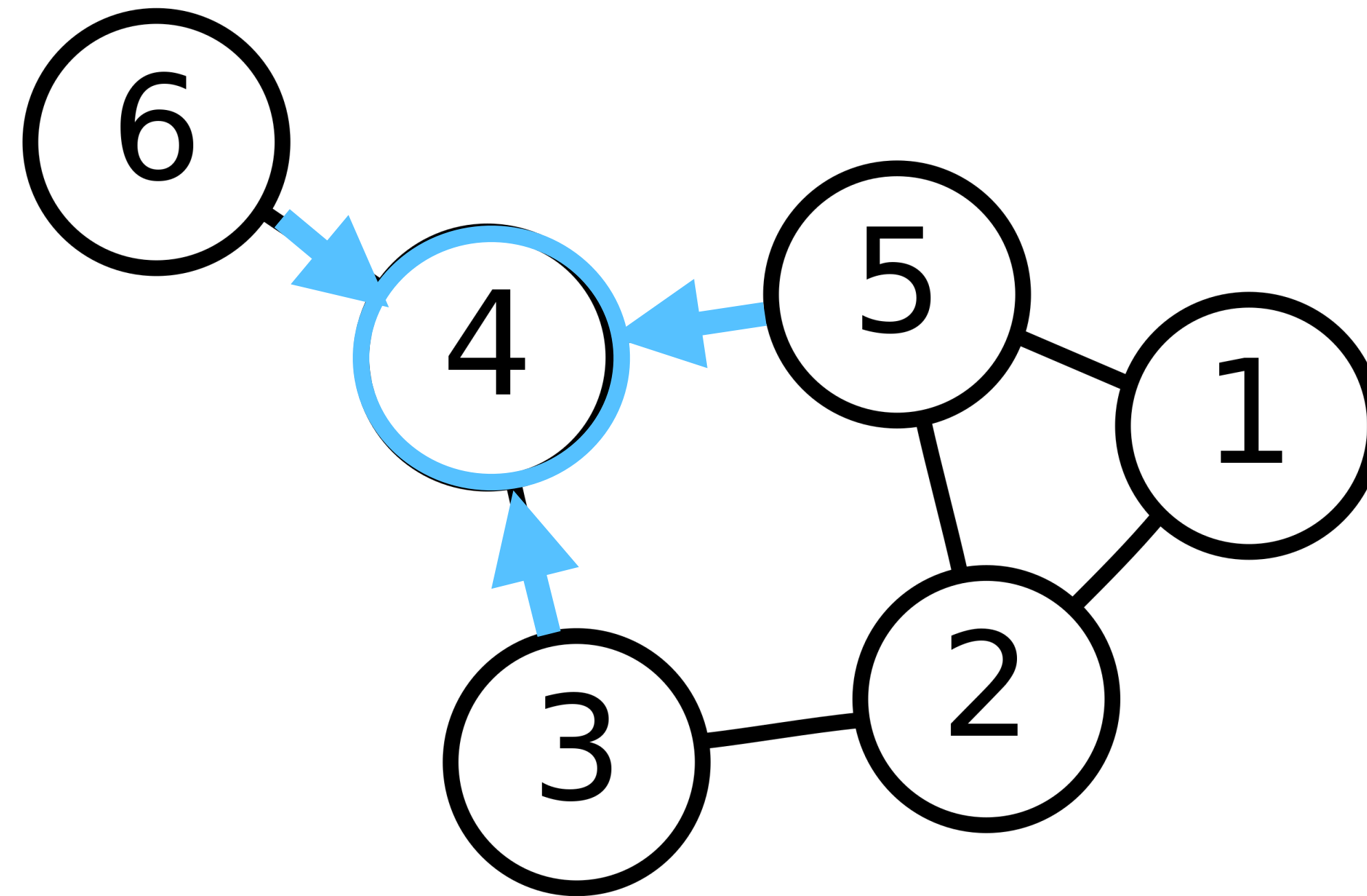
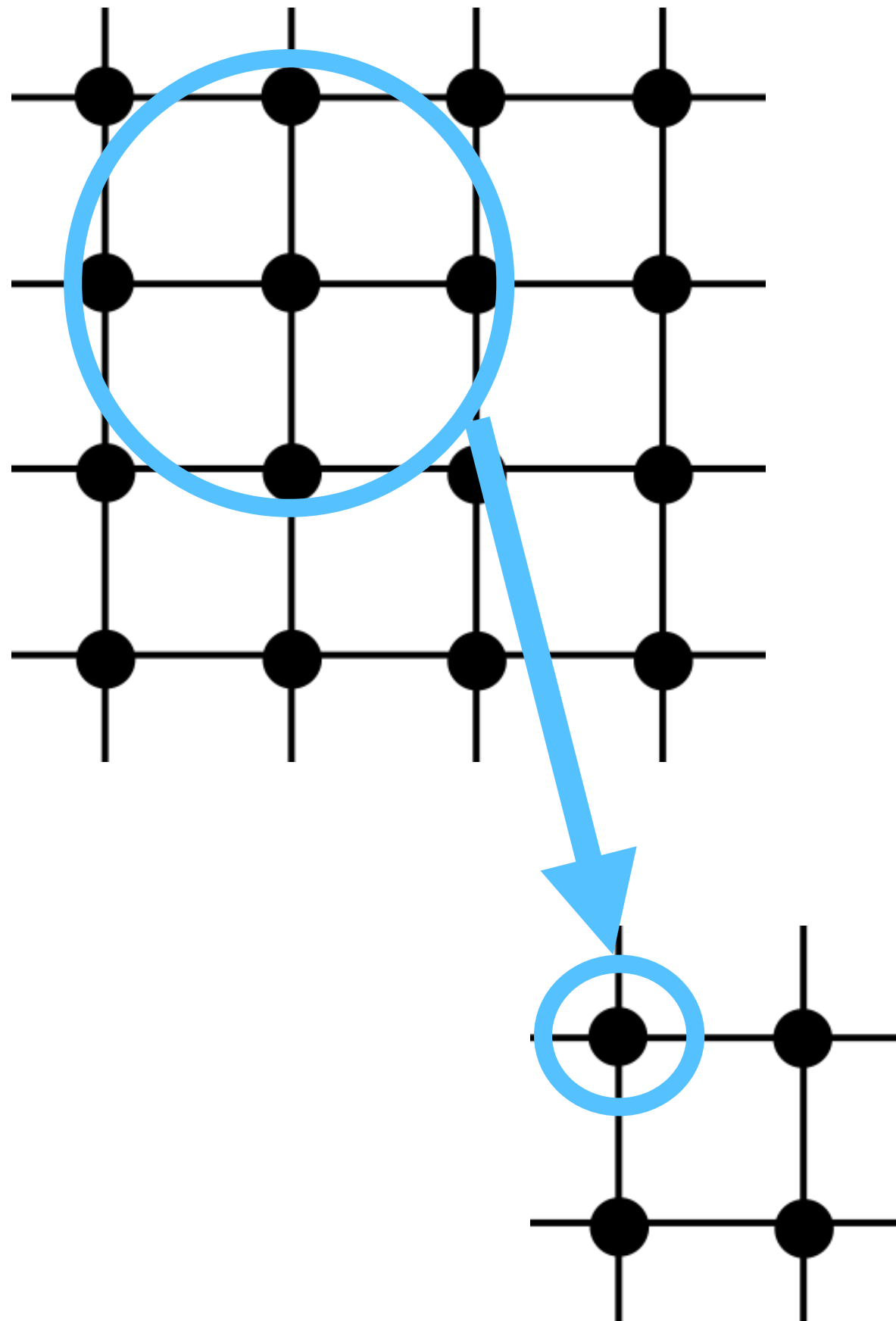


**Graphs**



**Geodesics & Gauges**

# Convolutions on Graphs: Message Passing





# Convolutions on Graphs: Message Passing

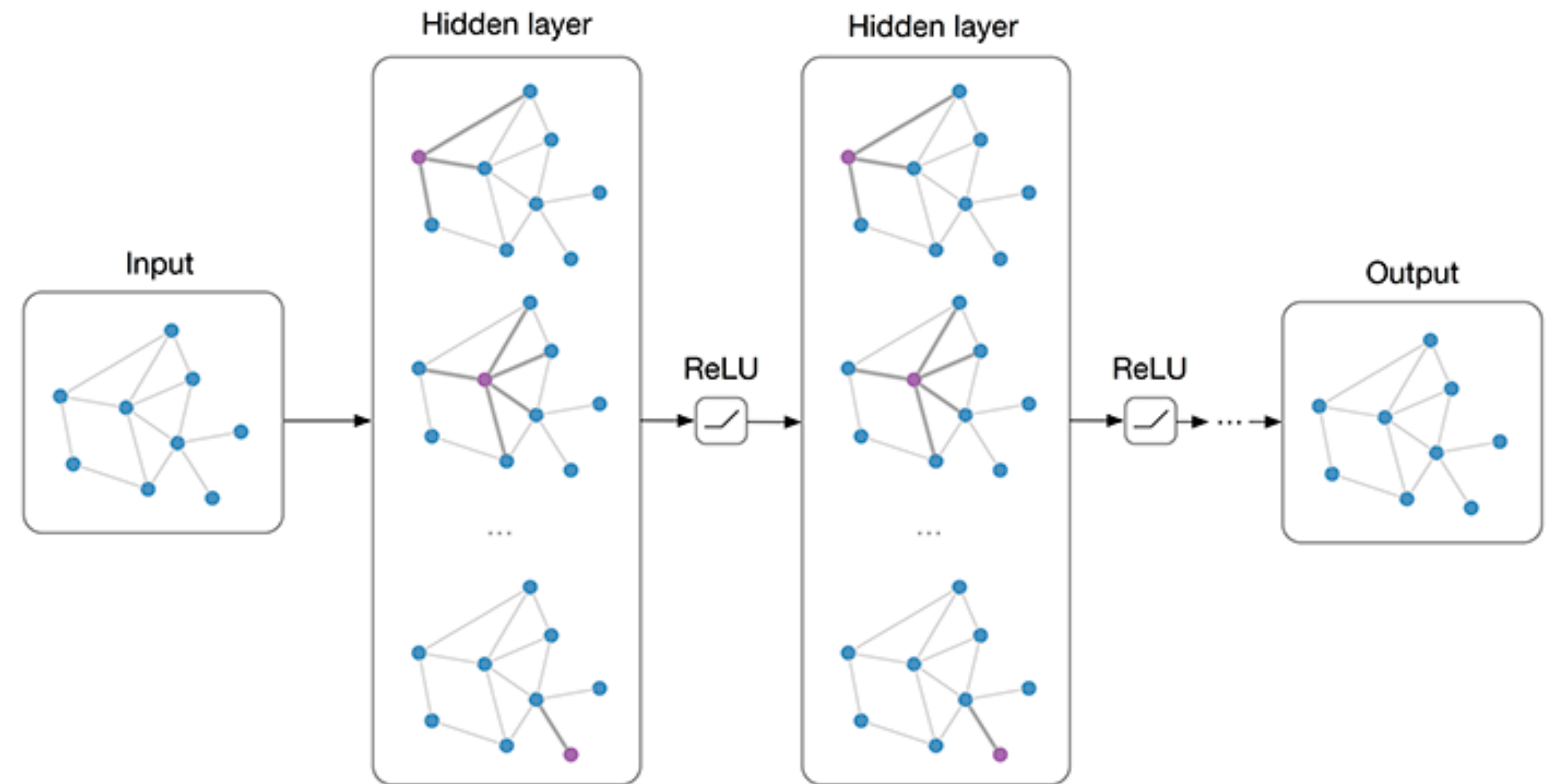
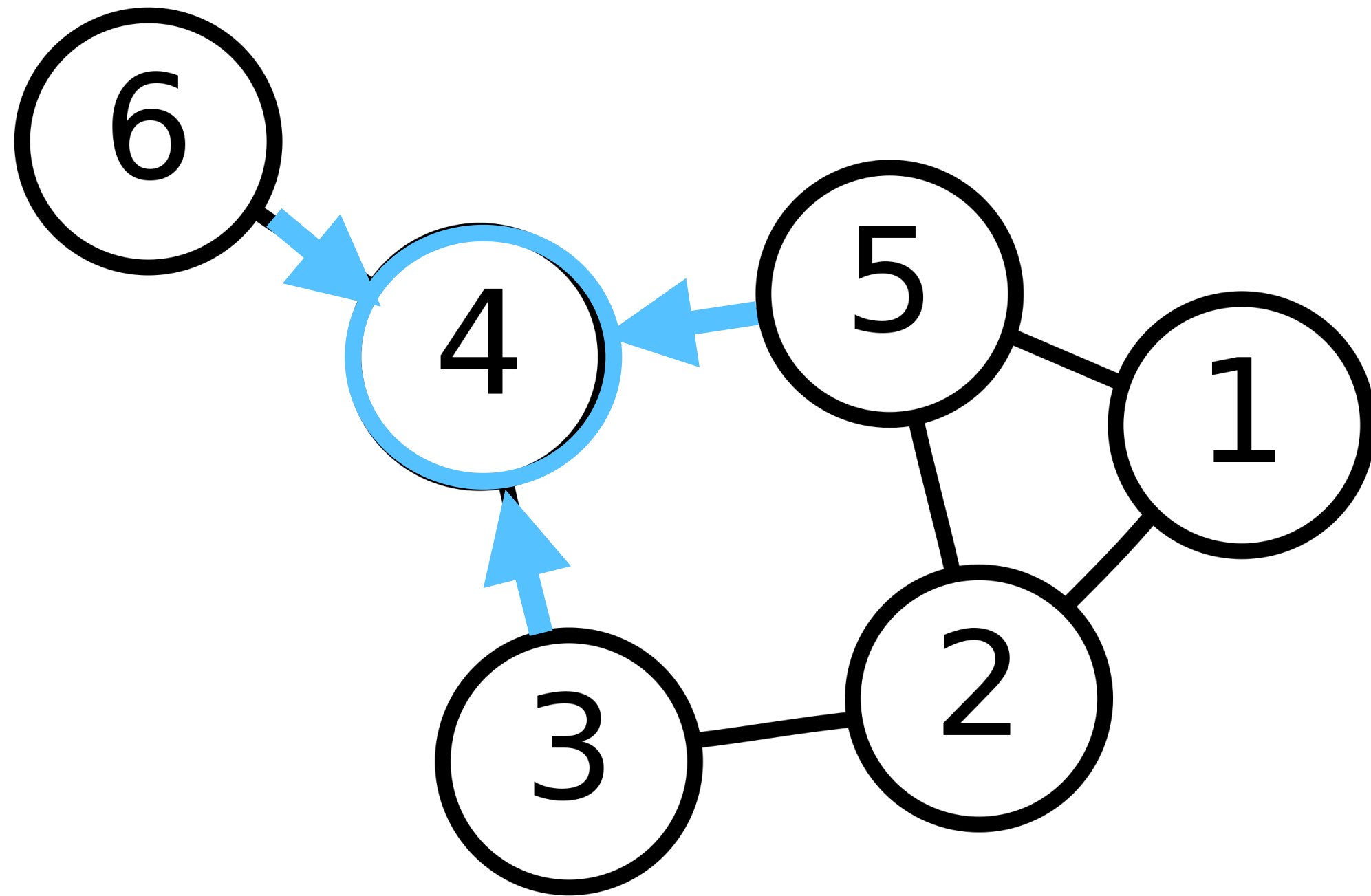


Image from T. Kipf's blog. Kipf et al. ICLR 2017

# Convolutions on Graphs: Message Passing

$$h_u^{(k+1)} = \text{update}^{(k)}(h_u^{(k)}, \text{agg}_{v \in \mathcal{N}(u)}^{(k)} h_v^{(k)})$$

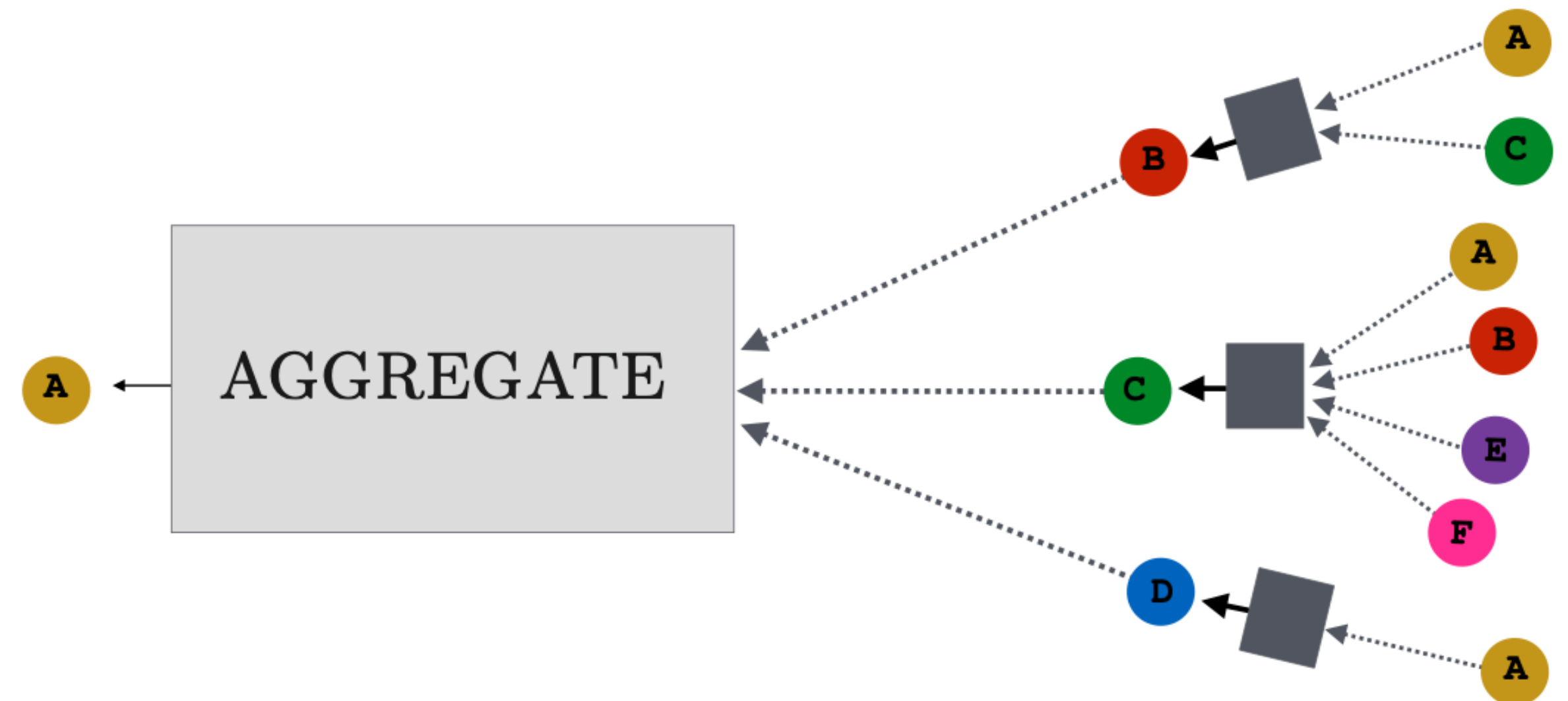
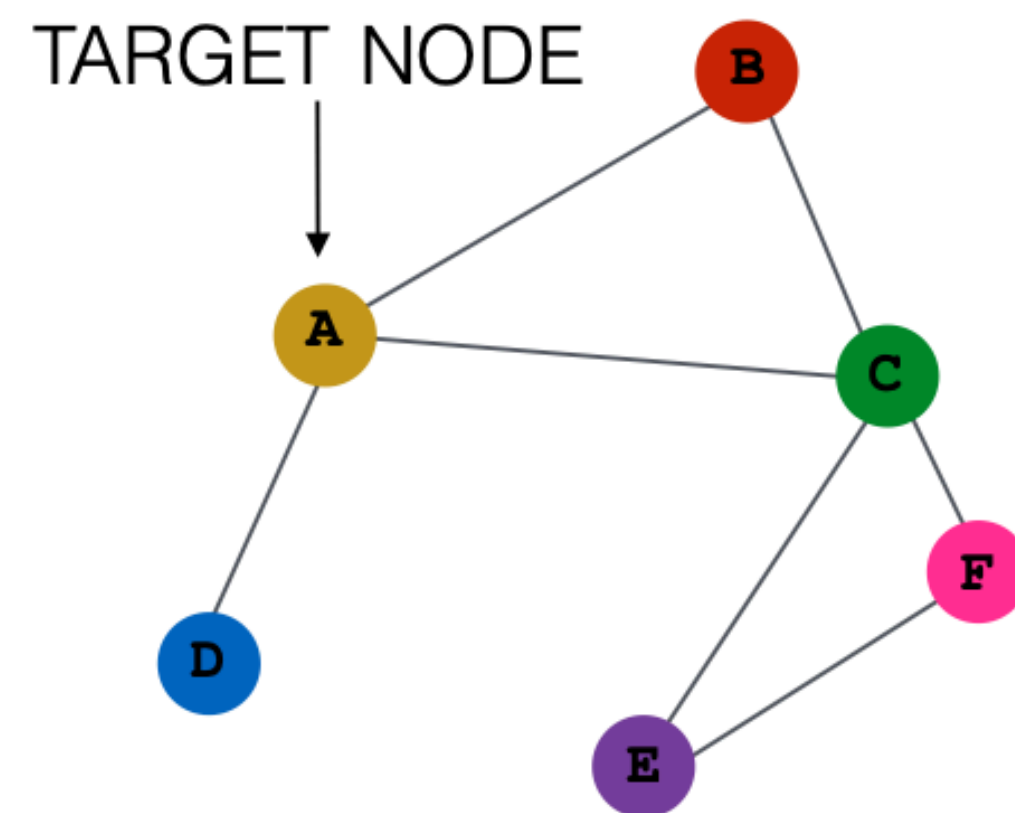
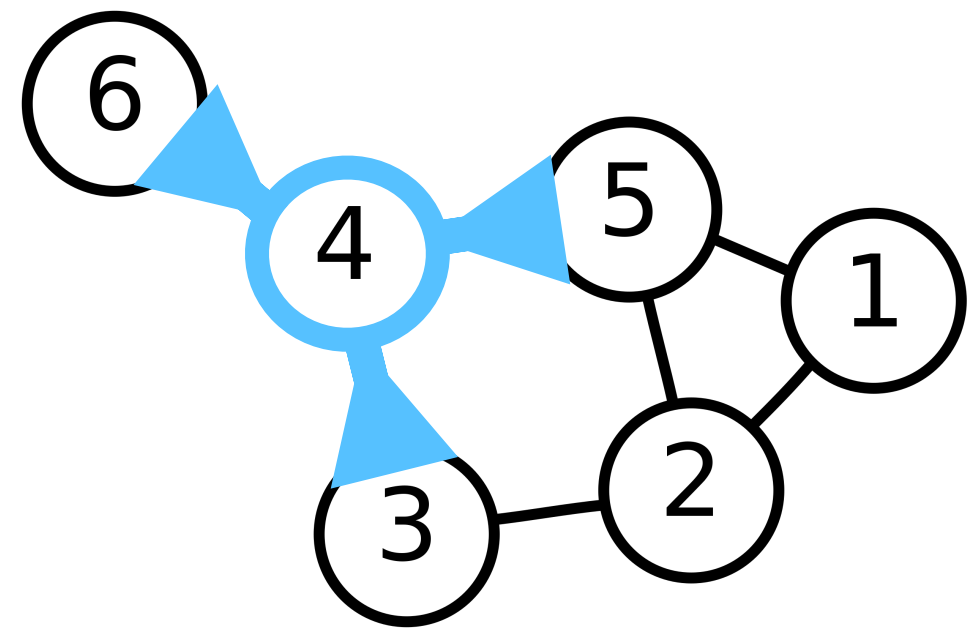


Image from Hamilton, "Graph Representation Learning Book"



# Convolutions on Graphs: Message Passing

$$h_u^{(k+1)} = \text{update}^{(k)}(h_u^{(k)}, \text{agg}_{v \in \mathcal{N}(u)}^{(k)} h_v^{(k)})$$

## Main Points

- Each gray block is a trainable network
- Each node has a different architecture
- Some will have the same
- We can learn over different architectures

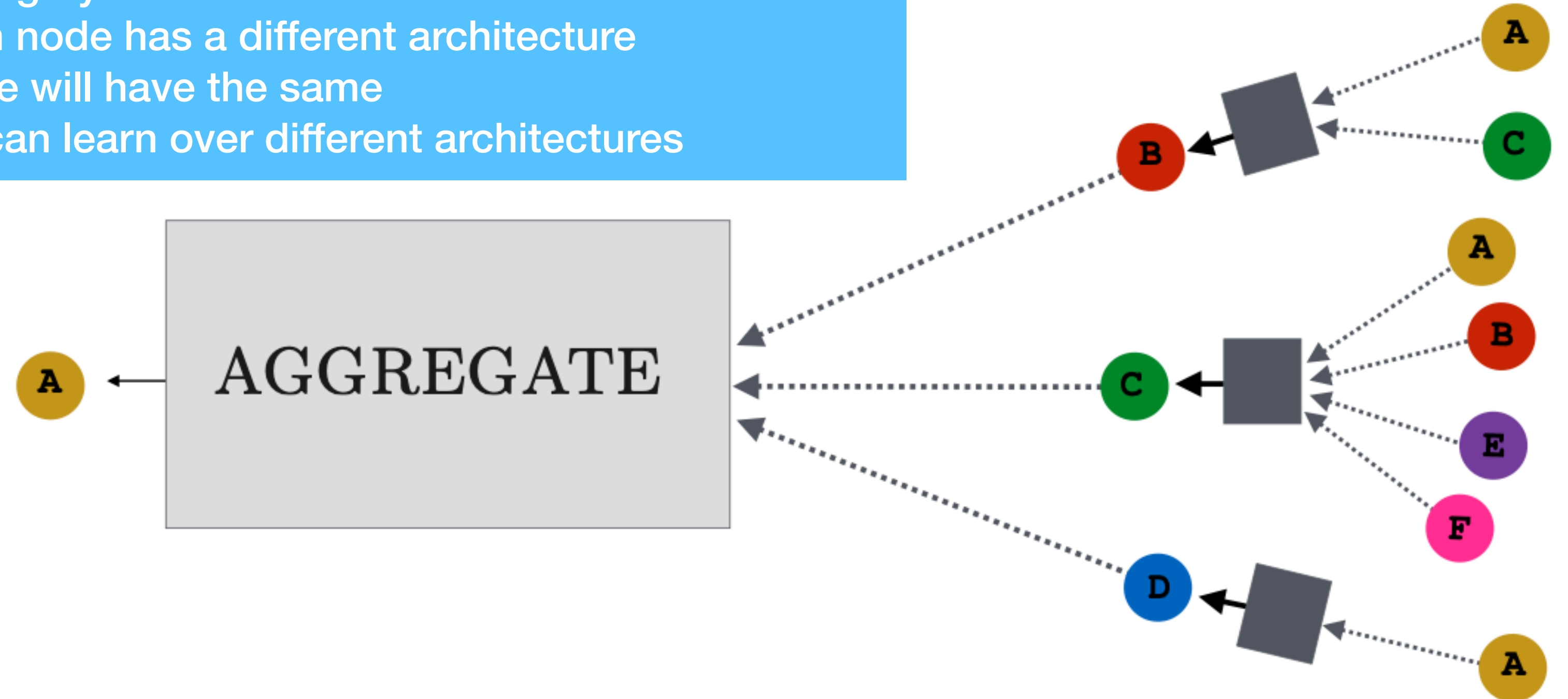
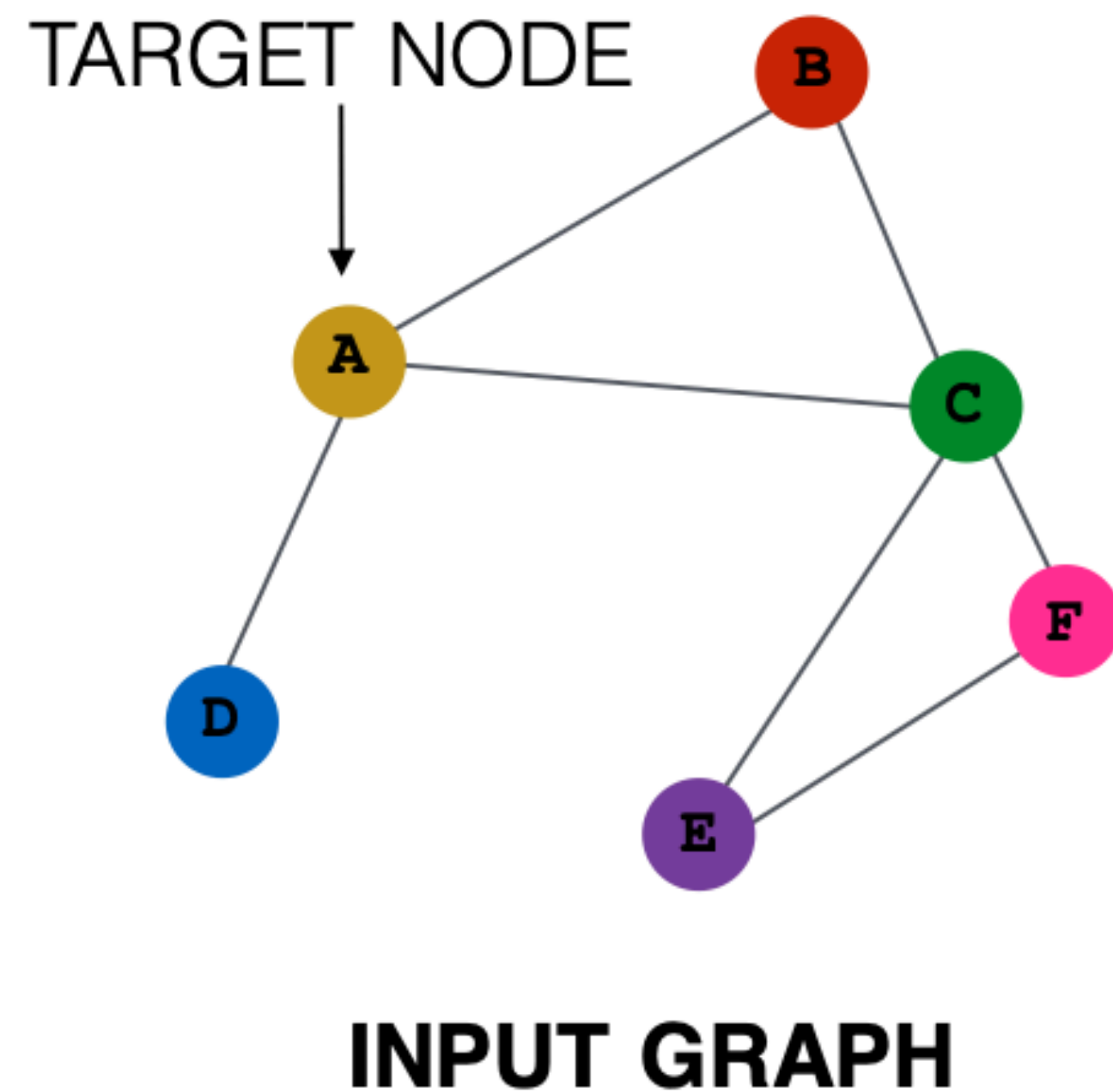


Image from Hamilton, "Graph Representation Learning Book"

# Convolutions on Graphs: Message Passing

$$h_u^{(k+1)} = \text{update}^{(k)}(h_u^{(k)}, \text{agg}_{v \in \mathcal{N}(u)}^{(k)} h_v^{(k)})$$

## Main Points

- Each layer incorporates information from nodes k-hops away
- Neighbours need to be aggregated need to be permutation invariant
- Different aggregations define different networks

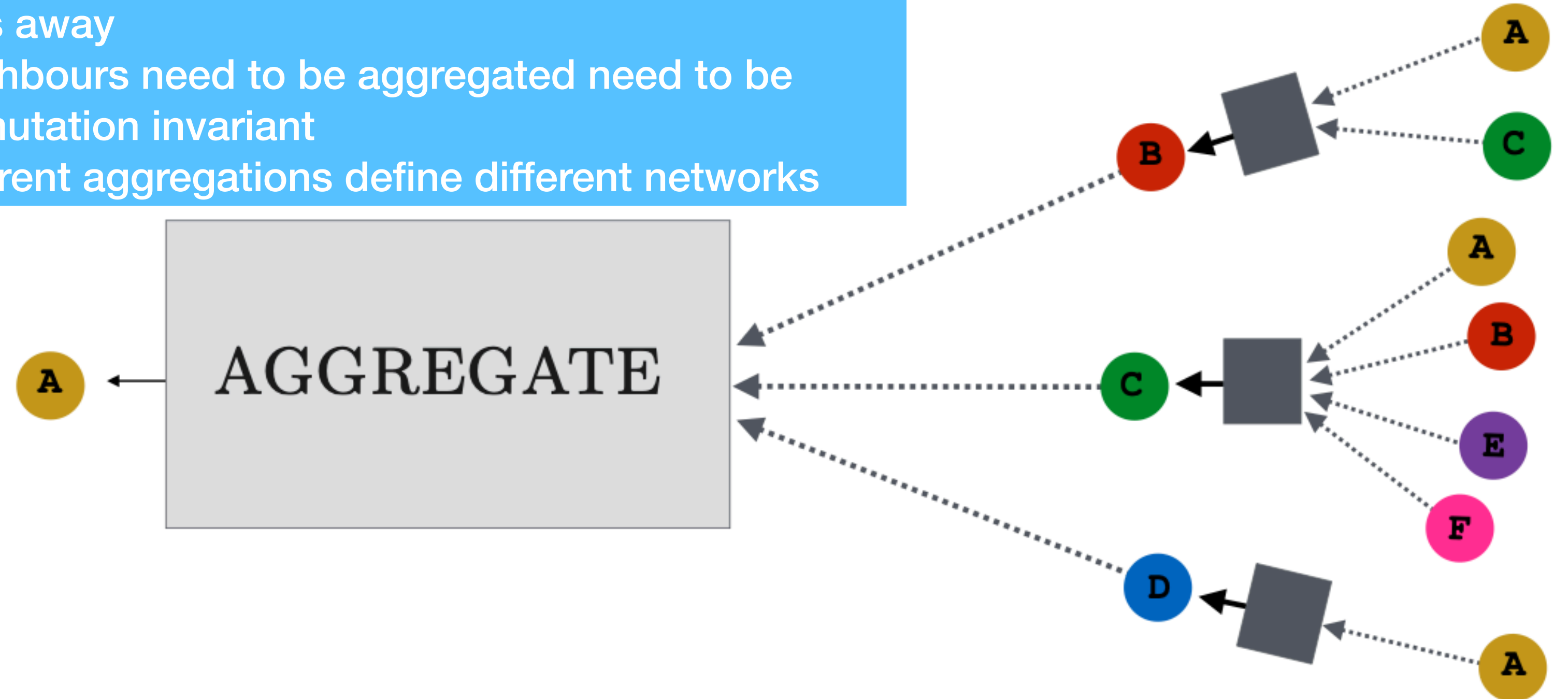
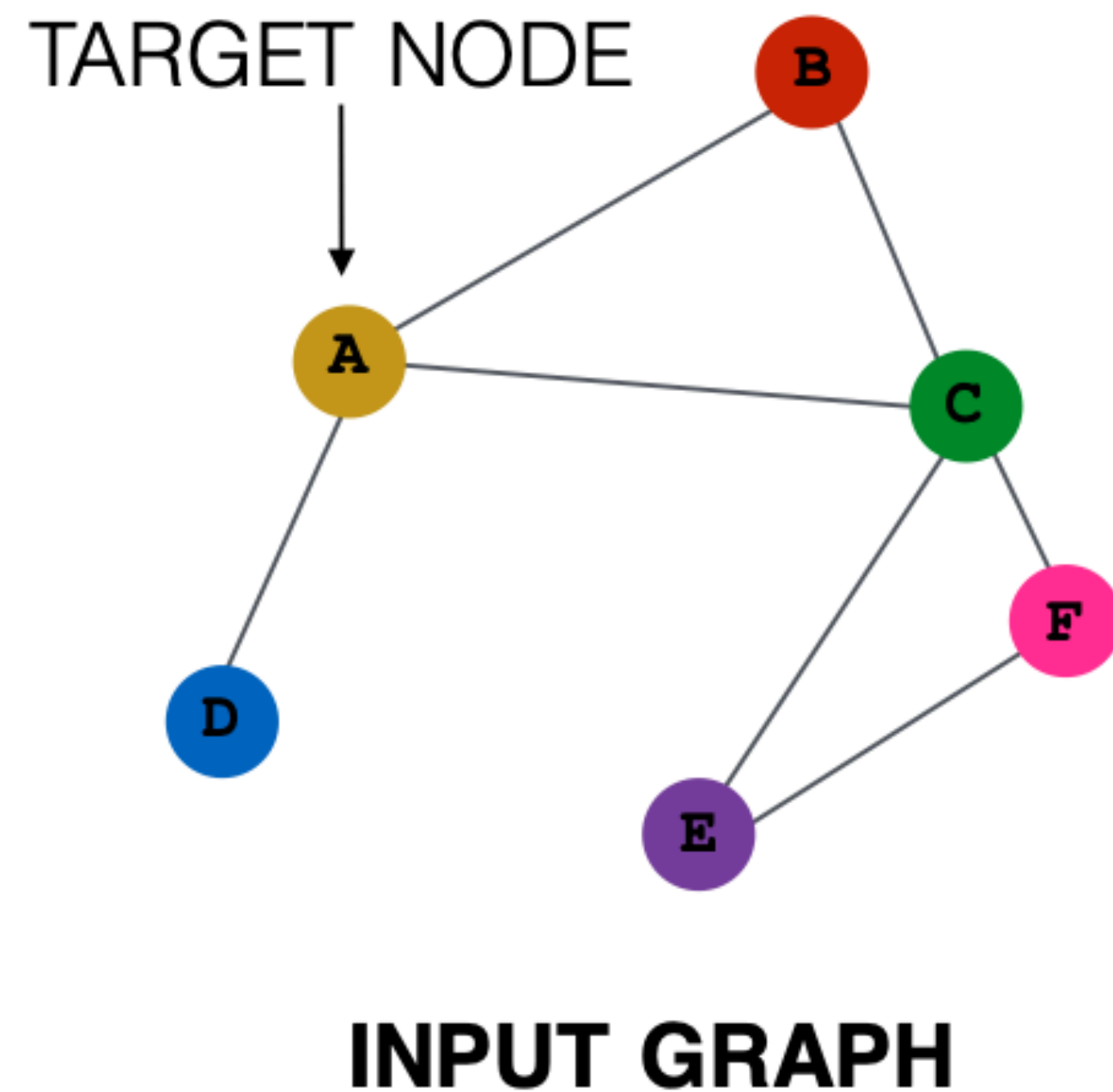


Image from Hamilton, "Graph Representation Learning Book"



# Convolutions on Graphs: Message Passing

$$h_u^{(k+1)} = \sigma \left( W^{(k)} \sum_{v \in \mathcal{N}(u) \cup \{u\}} \frac{h_v}{\sqrt{|\mathcal{N}(u)| |\mathcal{N}(v)|}} \right)$$

## Main Points

- Each layer incorporates information from nodes k-hops away
- Neighbours need to be aggregated need to be permutation invariant
- Different aggregations define different networks

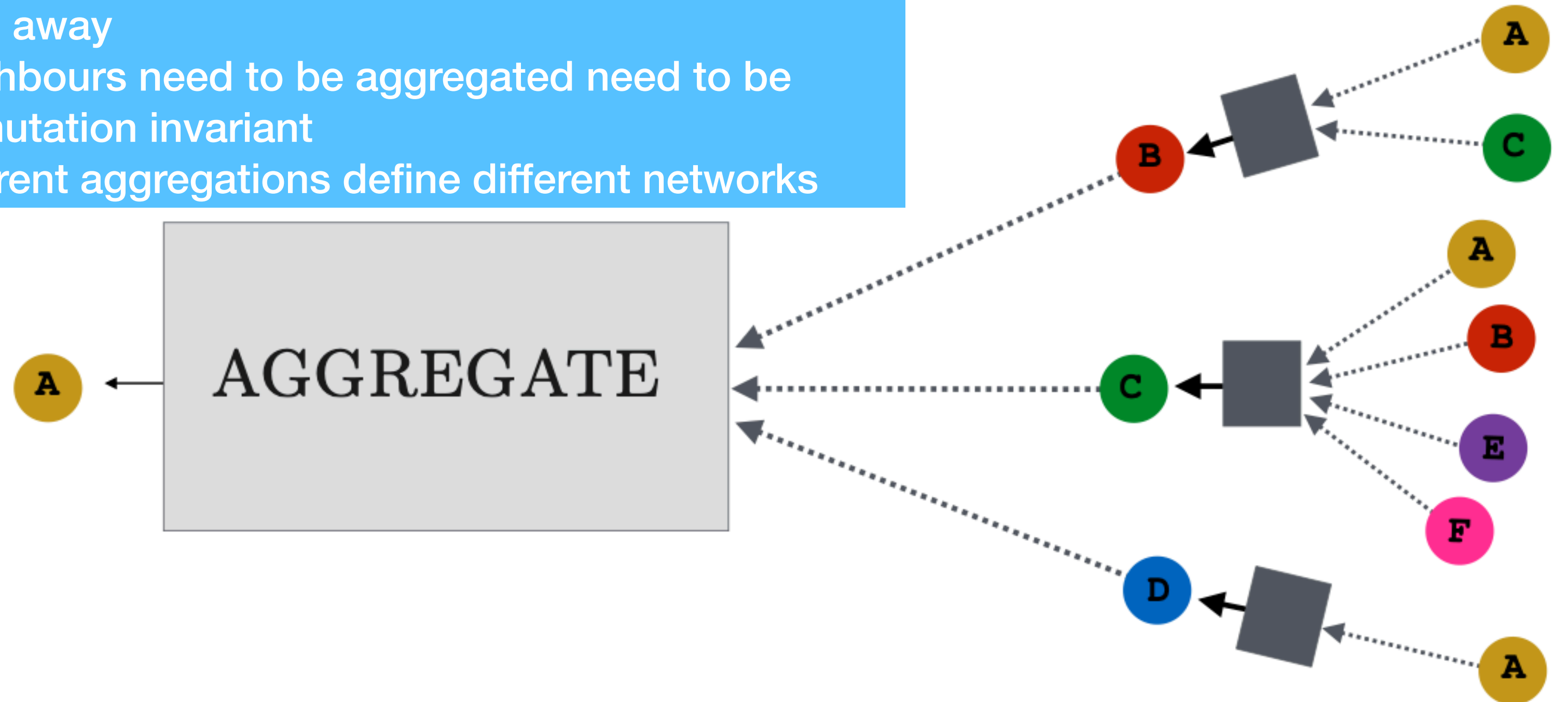
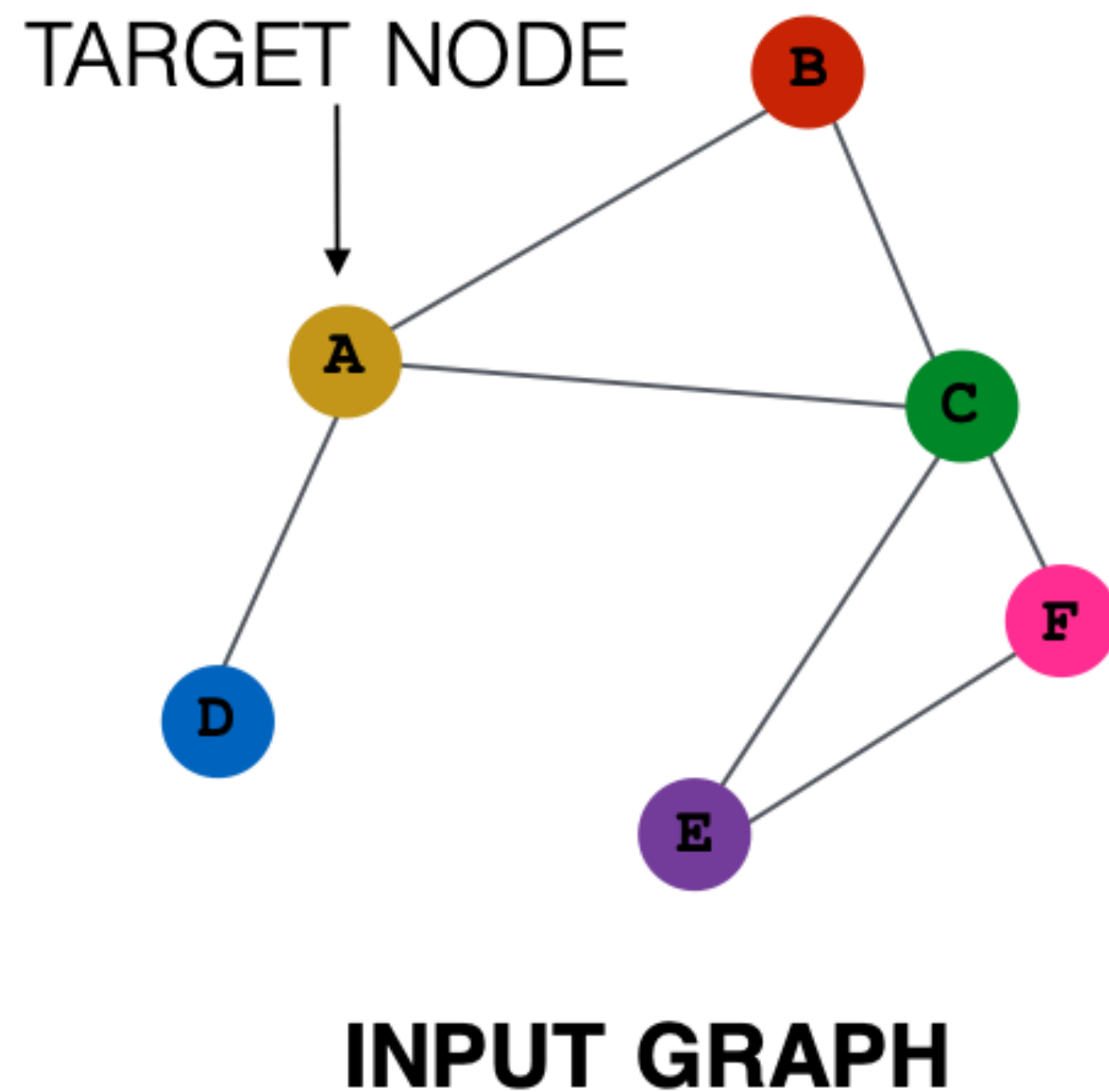


Image from Hamilton, "Graph Representation Learning Book"

# Convolutions on Graphs: Deep Encoder

$$h_u^{(k+1)} = \sigma \left( W^{(k)} \sum_{v \in \mathcal{N}(u)} \frac{h_v^{(k)}}{|\mathcal{N}(u)|} + B^{(k)} h_u^{(k)} \right), h_{(u)}^0 = x_v, h_u^{(L)} = z_v$$

## Main Points

- Each layer incorporates information from nodes k-hops away
- Neighbours need to be aggregated need to be permutation invariant
- Different aggregations define different networks

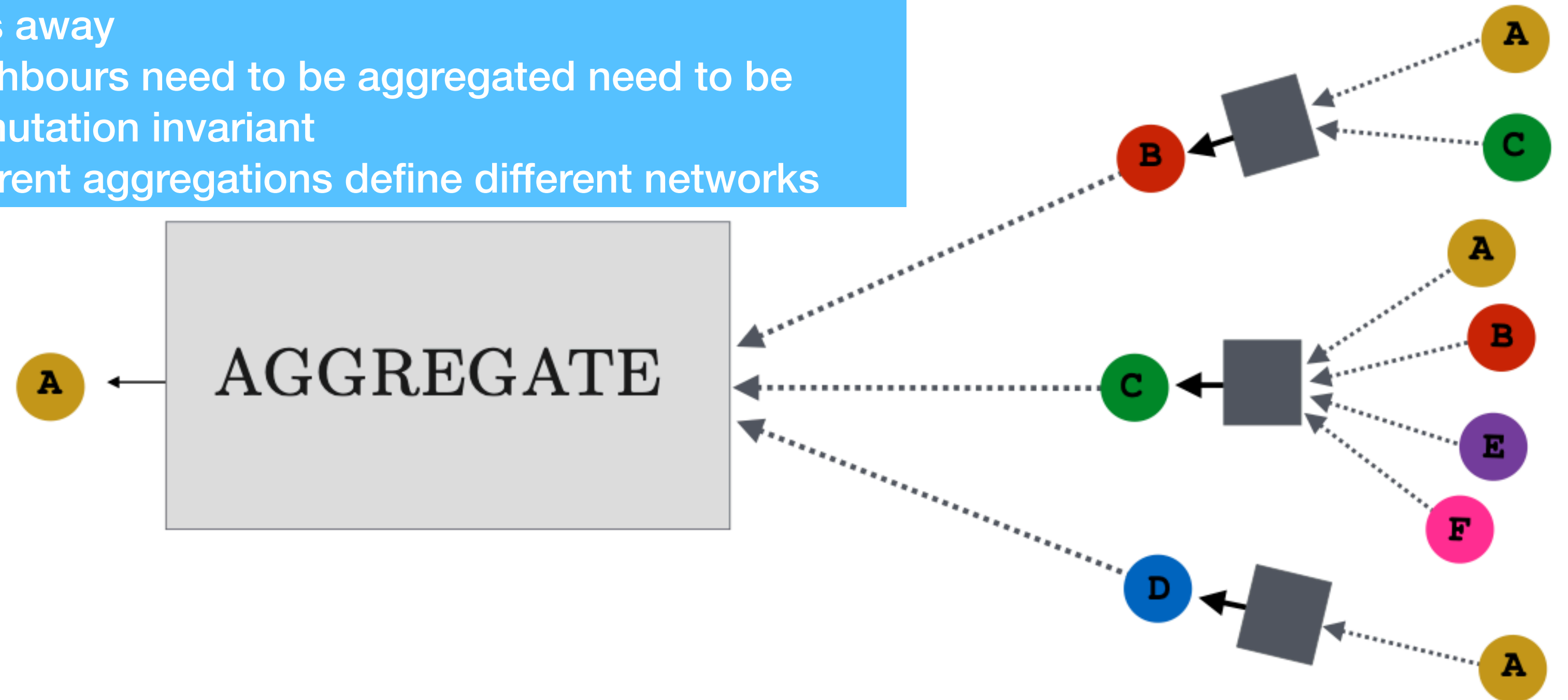
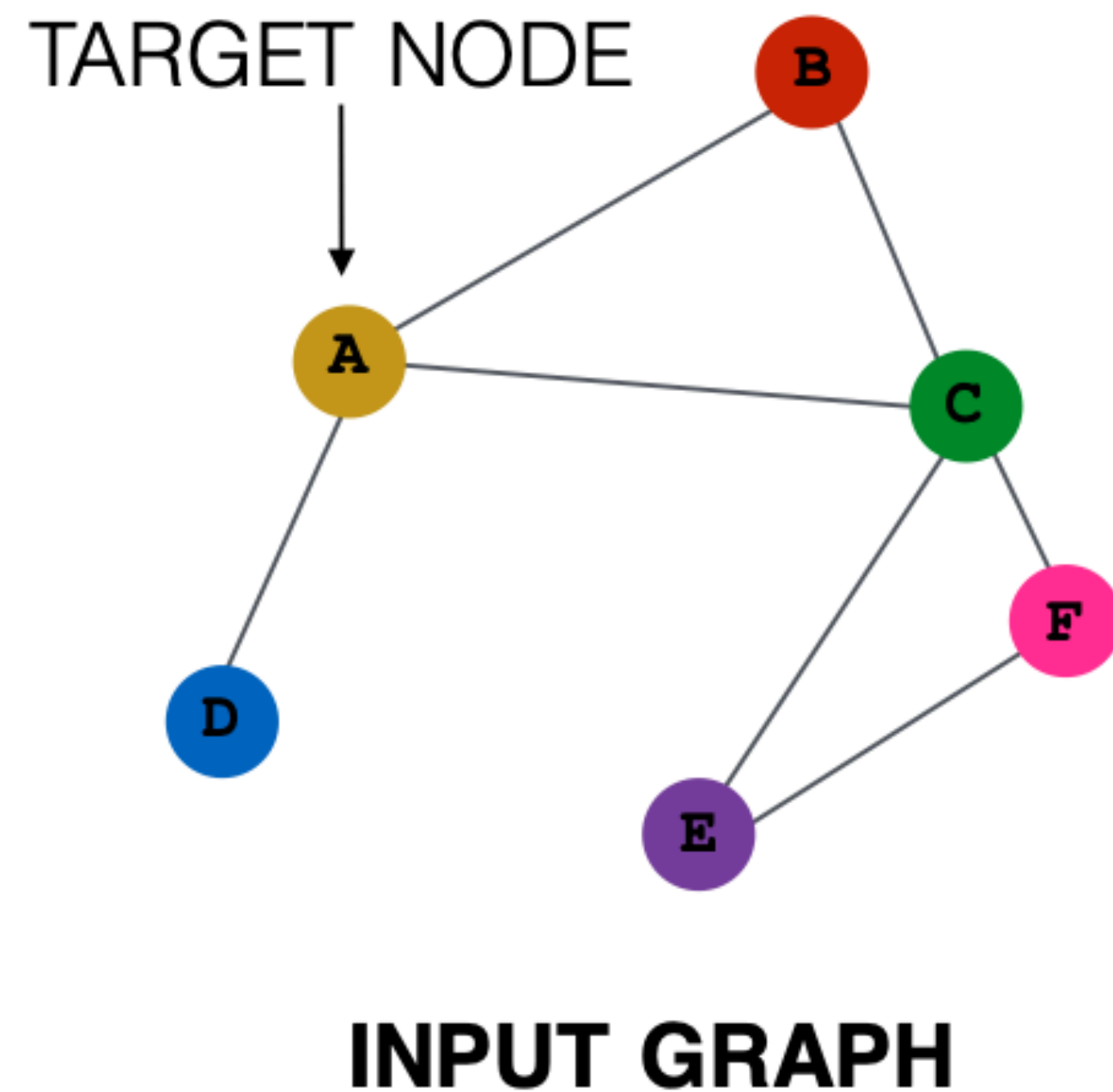


Image from Hamilton, "Graph Representation Learning Book"



# Convolutions on Graphs: Deep Encoder

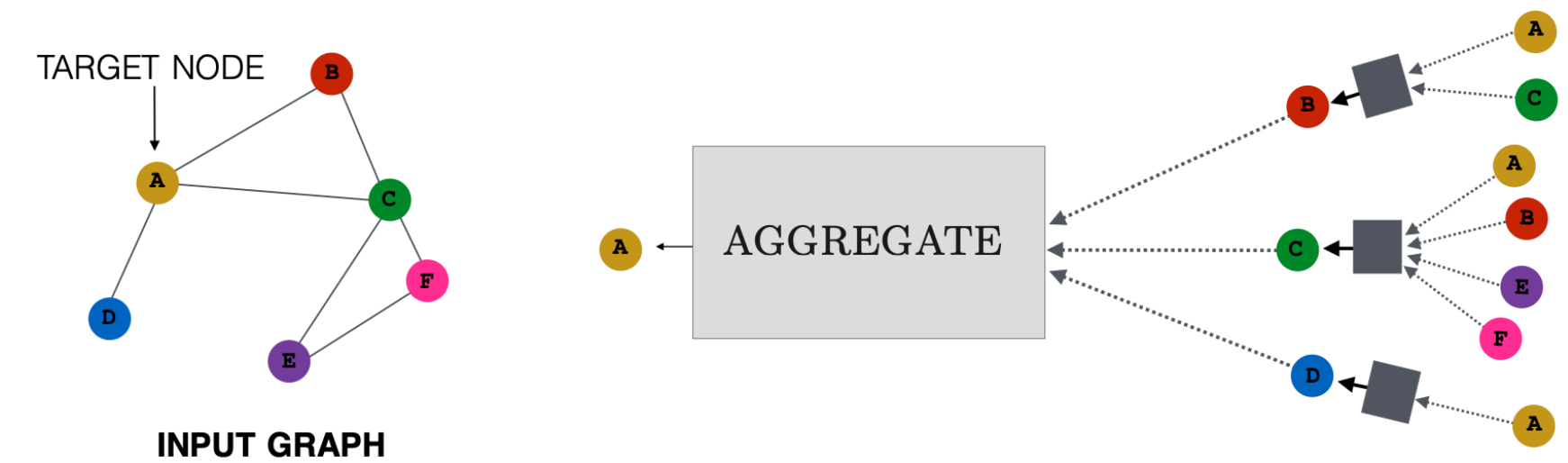
$$h_u^{(k+1)} = \sigma \left( W^{(k)} \sum_{v \in \mathcal{N}(u)} \frac{h_v^{(k)}}{|\mathcal{N}(u)|} + B^{(k)} h_u^{(k)} \right), h_{(u)}^0 = x_v, h_u^{(L)} = z_v$$

To train this model:

- Feed  $W^{(k)}$  and  $B^{(k)}$  to a loss and minimise with Stochastic gradient descent
- The matrices need to be shared across nodes
- In general, we can do this in matrix form  $H^{(k)} = [h_i^{(k)}]_i$ , define the diagonal degree matrix  $D_{u,u} = \text{Deg}(u) = |\mathcal{N}(u)|$
- Then,  $H^{(k+1)} = \sigma(D^{-1}AH^{(k)}W^{kT} + H^{(k)}B^{kT})$
- $D^{-1}A$  is sparse
- If the aggregation function is complex, the matrix formulation does not work

# Convolutions on Graphs: Deep Encoder

$$H^{(k+1)} = \sigma(D^{-1}AH^{(k)}W^{kT} + H^{(k)}B^{kT}), h_u^{(L)} = z_v$$

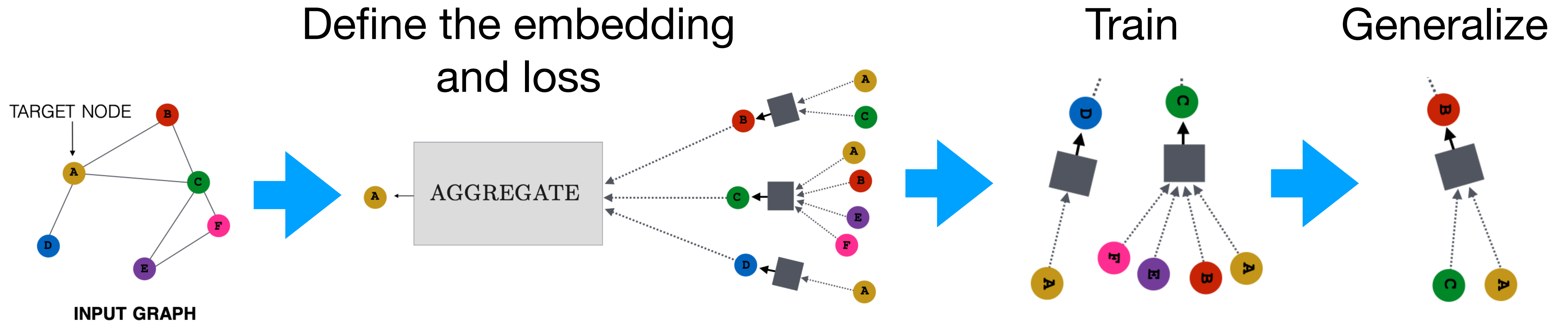


To train this model:

- Supervised learning  $\arg \min_{\theta} \mathcal{L}(Y, f(z_v))$ 
  - Example for classification with cross entropy loss
$$\mathcal{L} = \sum_u y_u \log(\sigma(z_u^T \theta)) + (1 - y_u) \log(1 - \sigma(z_u^T \theta))$$
- For unsupervised learning  $\mathcal{L} = \sum_{u,v} CE(y_{u,v}, DEC(z_u, z_v))$  where  $y_{u,v} = 1$  if nodes  $u$  and  $v$  are similar



# Convolutions on Graphs: Deep Encoder Design



Design The model:

1. Define the embedding: define a neighbourhood aggregation function
2. Define a loss function on the embedding and batch-train
3. Train the model on batch-computed graphs. Which are selected from a node batch + subgraph
4. Generate embeddings for nodes. This is applicable to different graphs and nodes

# Convolutions on Graphs: Results

Published as a conference paper at ICLR 2017

## SEMI-SUPERVISED CLASSIFICATION WITH GRAPH CONVOLUTIONAL NETWORKS

**Thomas N. Kipf**  
University of Amsterdam  
T.N.Kipf@uva.nl

**Max Welling**  
University of Amsterdam  
Canadian Institute for Advanced Research (CIFAR)  
M.Welling@uva.nl

Summary of results in terms of classification accuracy (in percent).

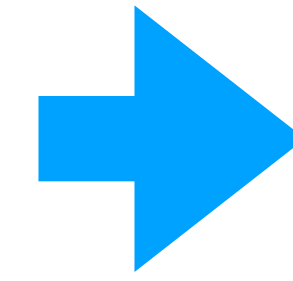
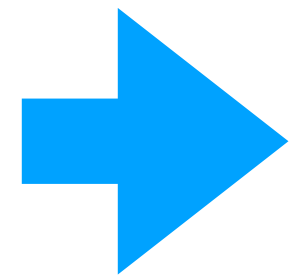
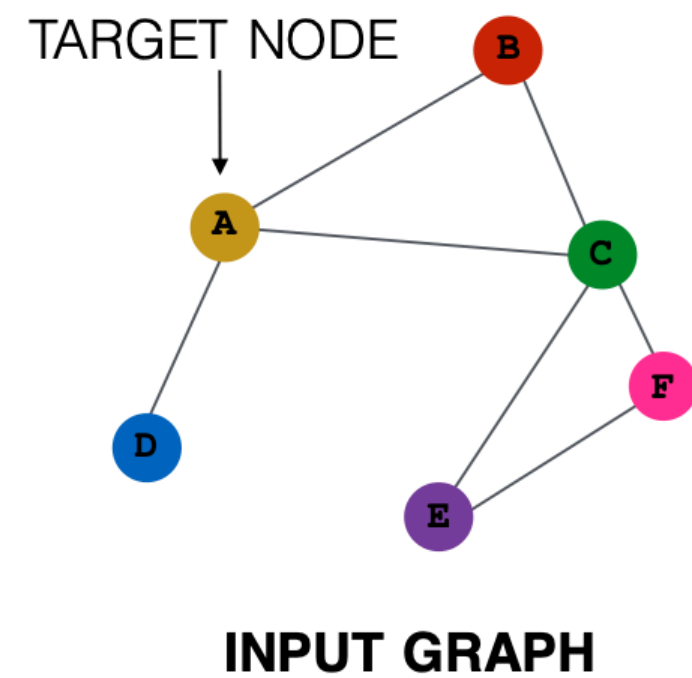
Method	Citeseer	Cora	Pubmed	NELL
ManiReg [3]	60.1	59.5	70.7	21.8
SemiEmb [28]	59.6	59.0	71.1	26.7
LP [32]	45.3	68.0	63.0	26.5
DeepWalk [22]	43.2	67.2	65.3	58.1
ICA [18]	69.1	75.1	73.9	23.1
Planetoid* [29]	64.7 (26s)	75.7 (13s)	77.2 (25s)	61.9 (185s)
<b>GCN (this paper)</b>	<b>70.3</b> (7s)	<b>81.5</b> (4s)	<b>79.0</b> (38s)	<b>66.0</b> (48s)
GCN (rand. splits)	67.9 $\pm$ 0.5	80.1 $\pm$ 0.5	78.9 $\pm$ 0.7	58.4 $\pm$ 1.7



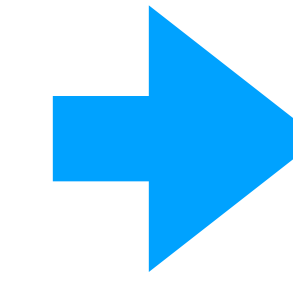
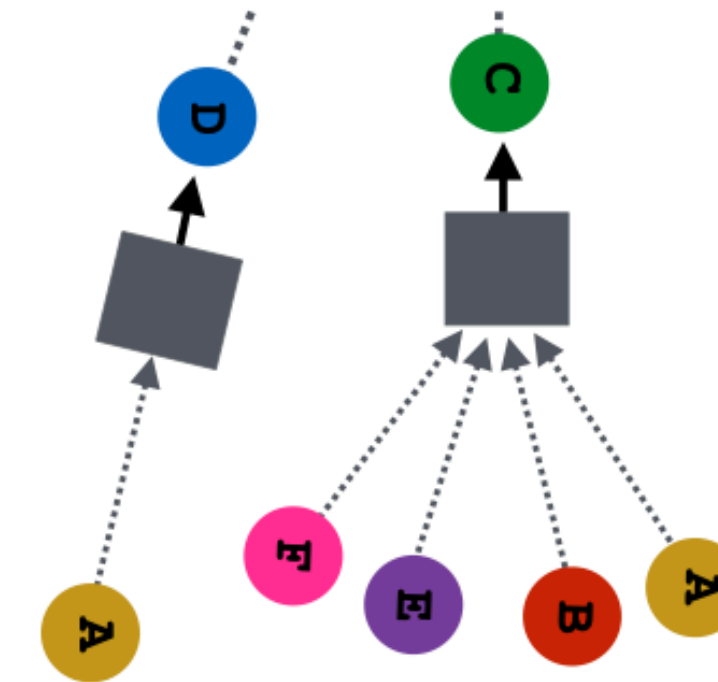
# Convolutions on Graphs: Deep Encoder

## General Case Across Different GNNs

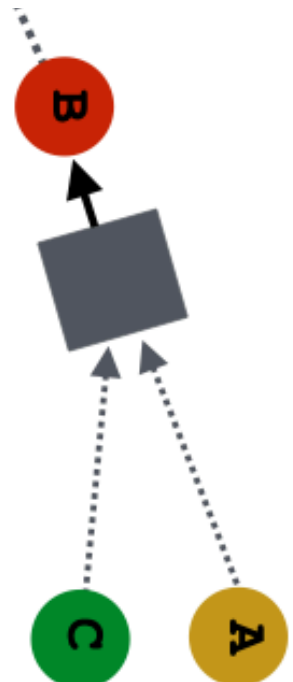
Define the embedding  
and loss



Train



Generalize



$$h_u^{(k+1)} = \sigma \left( \text{Agg}_{v \in \mathcal{N}(u)} (\text{Msg } \phi(h_v^{(k)})), \psi(h_u^{(k)}) \right), h_{(u)}^0 = x_v, h_u^{(L)} = z_v$$

Change over different GNNs

# Convolutions on Graphs: Deep Encoder

## General Case Across Different GNNs

$$h_u^{(k+1)} = \sigma \left( \text{Agg}_{v \in \mathcal{N}(u)} (\text{Msg } \phi(h_v^{(k)})), \psi(h_u^{(k)}) \right), h_{(u)}^0 = x_v, h_u^{(L)} = z_v$$

### GCN

- $\sigma$  is the sigmoid function
- $\text{Msg}$  is the weighting of nodes  $W^{(k+1)}h_v^{(k)}$
- $\text{Agg}$  is the average
- $\phi$  is the identity function
- $\psi$  is null

### GraphSage

- $\sigma$  is the sigmoid function
- $\text{Msg}$  is the weighting of nodes  $W^{(k+1)}h_v^{(k)}$
- $\text{Agg}$  is two-stage
  - Aggregate from networks
  - Different aggregation with the node itself
- $\phi$  is the identity function
- $\psi$  is null

# Convolutions on Graphs: Deep Encoder

## General Case Across Different GNNs

$$h_u^{(k+1)} = \sigma \left( \text{Agg}_{v \in \mathcal{N}(u)} (\text{Msg } \phi(h_v^{(k)})), \psi(h_u^{(k)}) \right), h_{(u)}^0 = x_v, h_u^{(L)} = z_v$$

### GraphSage

- $\sigma$  is the sigmoid function

sometimes  $\sigma(x) = \sigma \left( \frac{x}{\|x\|_2} \right)$

- Msg is the weighting of nodes  $W^{(k+1)}h_v^{(k)}$
- Agg is two-stage
  - Aggregate from networks
  - Different aggregation with the node itself
- $\phi$  is the identity function
- $\psi$  is null

### GraphSage Aggregations

- Mean, like GCN
- Pool, by applying a non-identity  $\phi$
- LSTM aggregations



# GraphSage: Results

## Inductive Representation Learning on Large Graphs

**William L. Hamilton\***

wleif@stanford.edu

**Rex Ying\***

rexying@stanford.edu

**Jure Leskovec**

jure@cs.stanford.edu

Department of Computer Science  
Stanford University  
Stanford, CA, 94305

Table 1: Prediction results for the three datasets (micro-averaged F1 scores). Results for unsupervised and fully supervised GraphSAGE are shown. Analogous trends hold for macro-averaged scores.

Name	Citation		Reddit		PPI	
	Unsup. F1	Sup. F1	Unsup. F1	Sup. F1	Unsup. F1	Sup. F1
Random	0.206	0.206	0.043	0.042	0.396	0.396
Raw features	0.575	0.575	0.585	0.585	0.422	0.422
DeepWalk	0.565	0.565	0.324	0.324	—	—
DeepWalk + features	0.701	0.701	0.691	0.691	—	—
GraphSAGE-GCN	0.742	0.772	<b>0.908</b>	0.930	0.465	0.500
GraphSAGE-mean	0.778	0.820	0.897	0.950	0.486	0.598
GraphSAGE-LSTM	0.788	0.832	<b>0.907</b>	<b>0.954</b>	0.482	<b>0.612</b>
GraphSAGE-pool	<b>0.798</b>	<b>0.839</b>	0.892	0.948	<b>0.502</b>	0.600
% gain over feat.	39%	46%	55%	63%	19%	45%

# Convolutions on Graphs: Deep Encoder

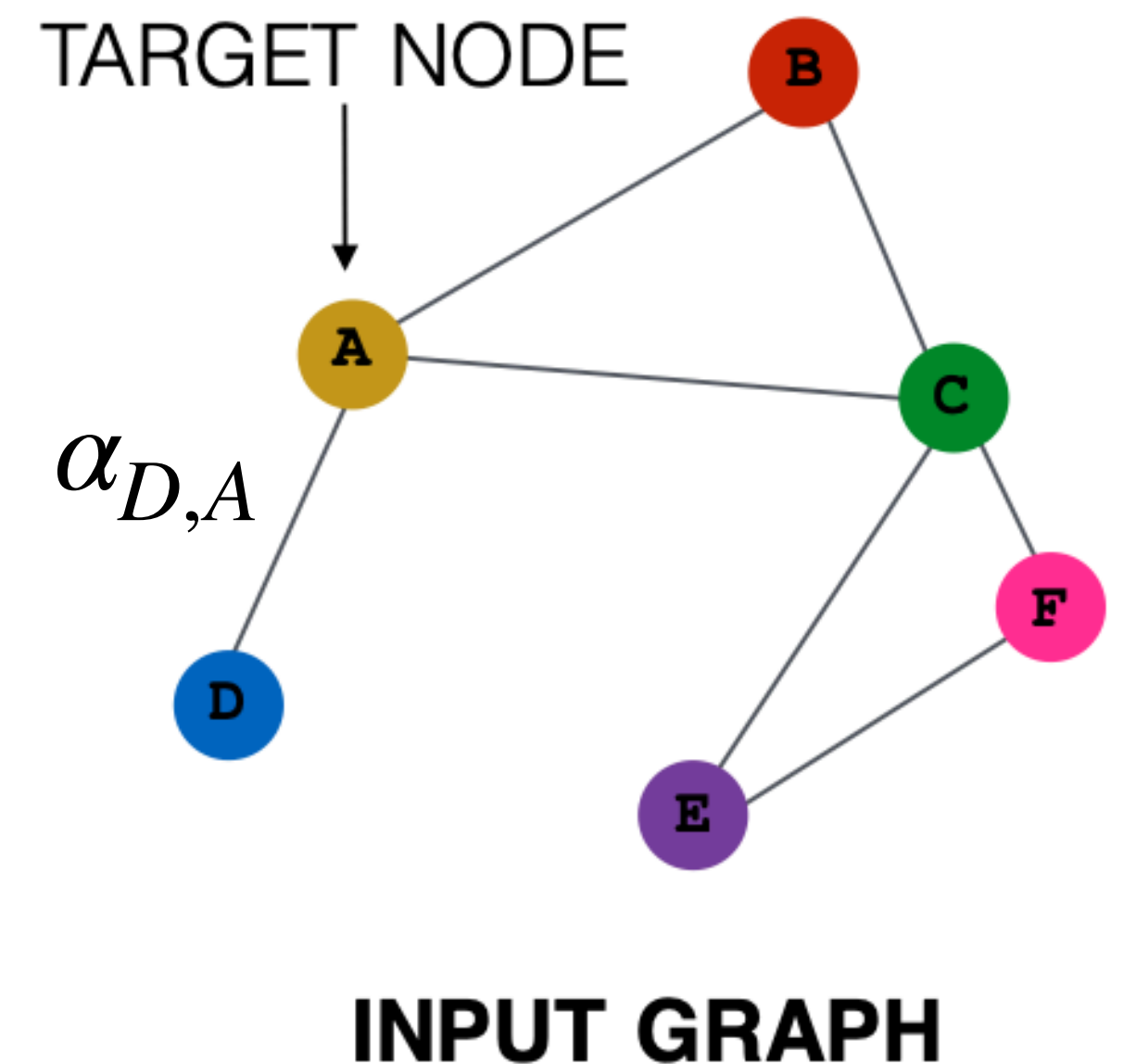
## General Case Across Different GNNs

$$h_u^{(k+1)} = \sigma \left( \text{Agg}_{v \in \mathcal{N}(u)} (\text{Msg } \phi(h_v^{(k)})), \psi(h_u^{(k)}) \right), h_{(u)}^0 = x_v, h_u^{(L)} = z_v$$

### Graph Attention Network (GAT)

- Most parameters are arbitrary
- Msg is the weighting of nodes  $\alpha_{uv} W^{(k+1)} h_v^{(k)}$
- this learnable weighting  $\alpha_{u,v}$  will learn which nodes are more important for any give node embedding.

$$h_u^{(k+1)} = \sigma \left( \sum_{v \in \mathcal{N}(u)} \alpha_{u,v} W^{(k)} \frac{h_v^{(k)}}{|\mathcal{N}(u)|} + B^{(k)} h_u^{(k)} \right), h_{(u)}^0 = x_v, h_u^{(L)} = z_v, \alpha_{u,\cdot} = \sum_{v \in \mathcal{N}(u)} \alpha_{u,v} = 1$$



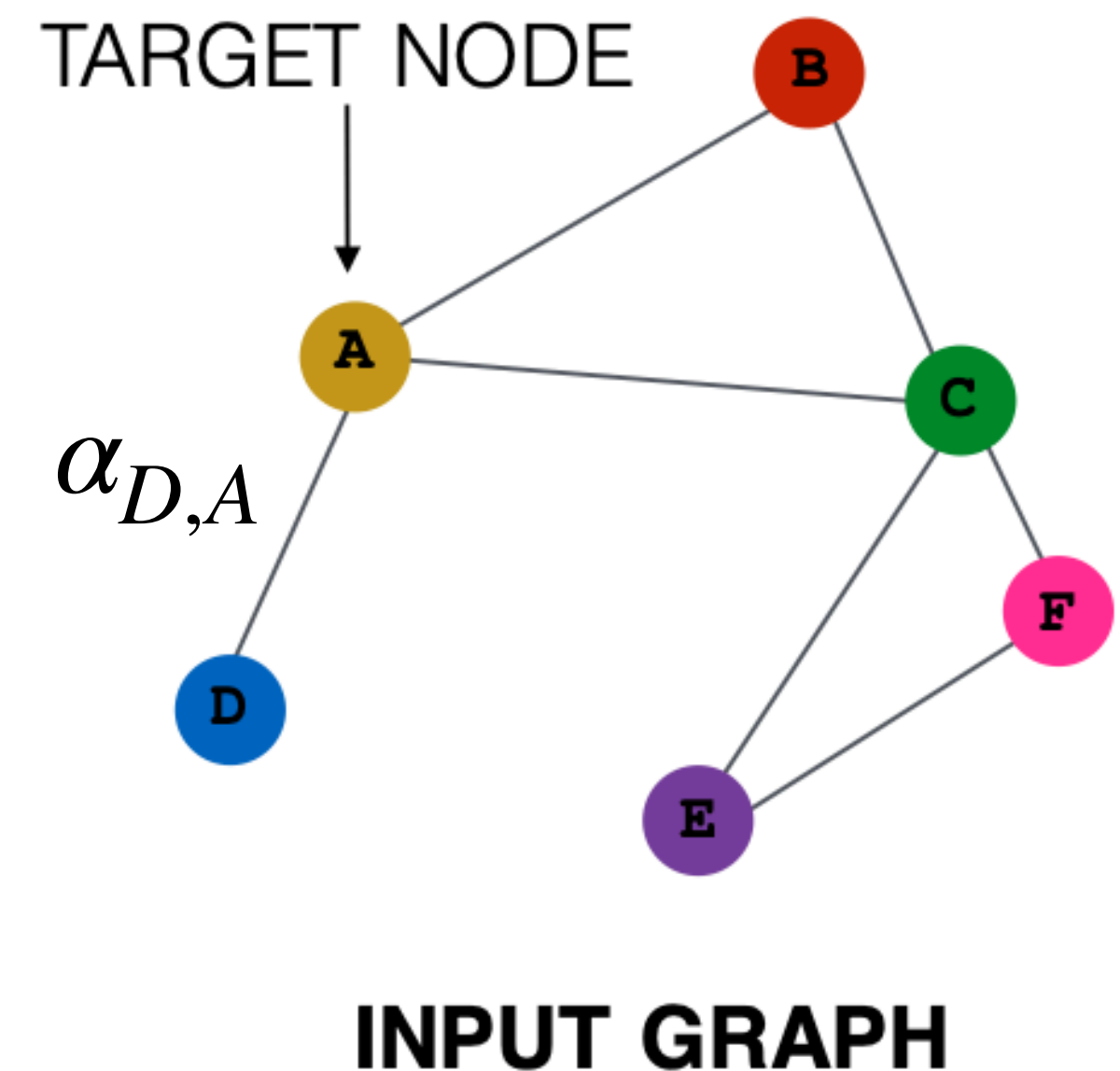
# Convolutions on Graphs: Deep Encoder

## General Case Across Different GNNs

$$h_u^{(k+1)} = \sigma \left( \text{Agg}_{v \in \mathcal{N}(u)} (\text{Msg } \phi(h_v^{(k)})), \psi(h_u^{(k)}) \right), h_{(u)}^0 = x_v, h_u^{(L)} = z_v$$

### Graph Attention Network (GAT)

- Main Benefits:
  - Implicit importance of neighbours
  - Computationally efficient
  - Storage efficient  $O(V+E)$  entries and fixed
  - Localised
  - Inductive, it doesn't depend on the graph structure





# Graphic Attention Network Results

Published as a conference paper at ICLR 2018

## Transductive

Method	Cora	Citeseer	Pubmed
MLP	55.1%	46.5%	71.4%
ManiReg (Belkin et al., 2006)	59.5%	60.1%	70.7%
SemiEmb (Weston et al., 2012)	59.0%	59.6%	71.7%
LP (Zhu et al., 2003)	68.0%	45.3%	63.0%
DeepWalk (Perozzi et al., 2014)	67.2%	43.2%	65.3%
ICA (Lu & Getoor, 2003)	75.1%	69.1%	73.9%
Planetoid (Yang et al., 2016)	75.7%	64.7%	77.2%
Chebyshev (Defferrard et al., 2016)	81.2%	69.8%	74.4%
GCN (Kipf & Welling, 2017)	81.5%	70.3%	<b>79.0%</b>
MoNet (Monti et al., 2016)	81.7 ± 0.5%	—	78.8 ± 0.3%
GCN-64*	81.4 ± 0.5%	70.9 ± 0.5%	<b>79.0 ± 0.3%</b>
<b>GAT</b> (ours)	<b>83.0 ± 0.7%</b>	<b>72.5 ± 0.7%</b>	<b>79.0 ± 0.3%</b>

## GRAPH ATTENTION NETWORKS

<b>Petar Veličković*</b> Department of Computer Science and Technology University of Cambridge petar.velickovic@cst.cam.ac.uk	<b>Guillem Cucurull*</b> Centre de Visió per Computador, UAB gcucurull@gmail.com
<b>Arantxa Casanova*</b> Centre de Visió per Computador, UAB ar.casanova.8@gmail.com	<b>Adriana Romero</b> Montréal Institute for Learning Algorithms adriana.romero.soriano@umontreal.ca
<b>Pietro Liò</b> Department of Computer Science and Technology University of Cambridge pietro.lio@cst.cam.ac.uk	<b>Yoshua Bengio</b> Montréal Institute for Learning Algorithms yoshua.umontreal@gmail.com

## Inductive

Method	PPI
Random	0.396
MLP	0.422
GraphSAGE-GCN (Hamilton et al., 2017)	0.500
GraphSAGE-mean (Hamilton et al., 2017)	0.598
GraphSAGE-LSTM (Hamilton et al., 2017)	0.612
GraphSAGE-pool (Hamilton et al., 2017)	0.600
GraphSAGE*	0.768
Const-GAT (ours)	0.934 ± 0.006
<b>GAT</b> (ours)	<b>0.973 ± 0.002</b>

# Summarising

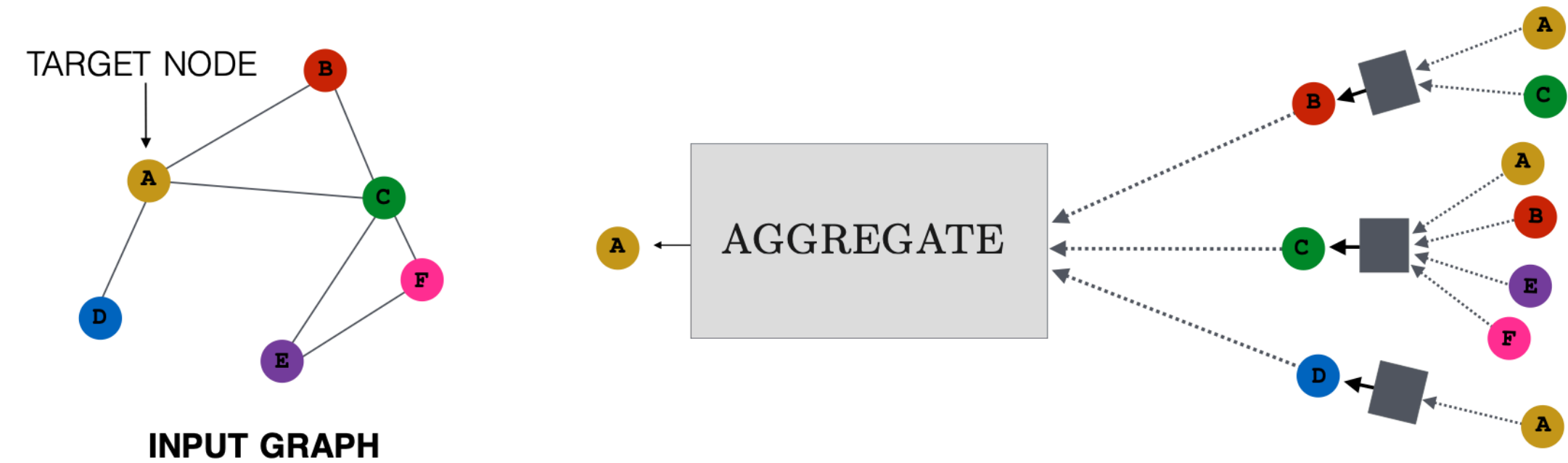


Image from Hamilton, “Graph Representation Learning Book”

- Graphs are good representations of data support and relationships
- Going from grid/lattices to graphs is non-trivial, we need permutation invariance, scale invariance, and sometimes topological invariance
- The main trick, is detect patches or motifs and generalise them to global structure
- There's ample evidence that taking into account heterogeneous structure improves supervised/semi-supervised tasks
- Novel techniques in graph networks, like GCNN, GraphSage, GAT, and others clearly improve on the results that don't take into account structure