

Complex Networks Comparison using Graphlets

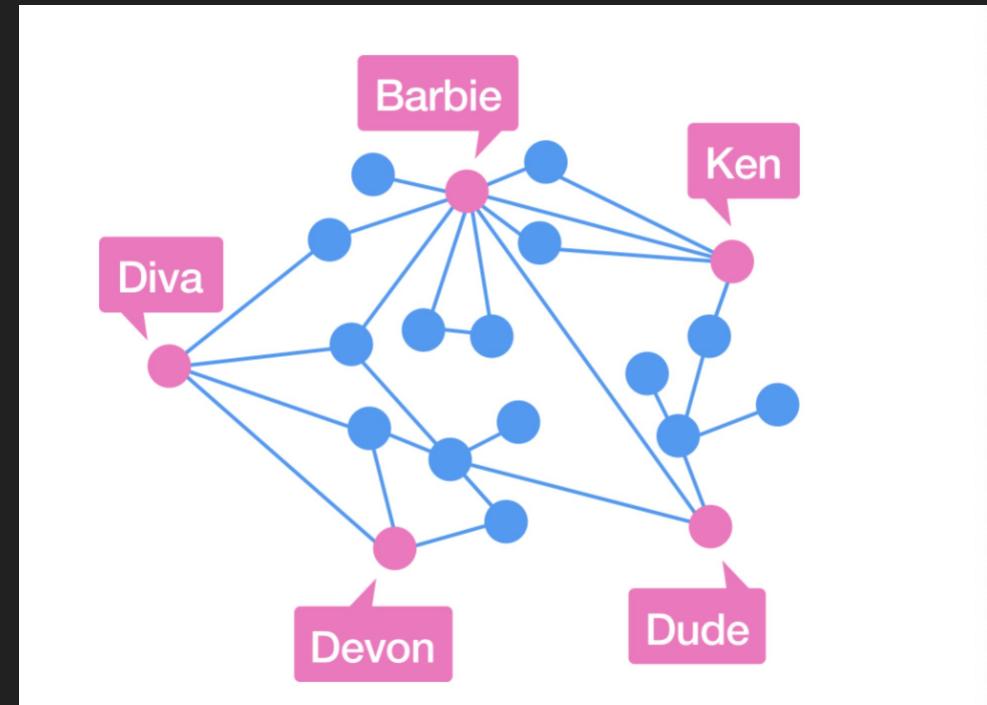
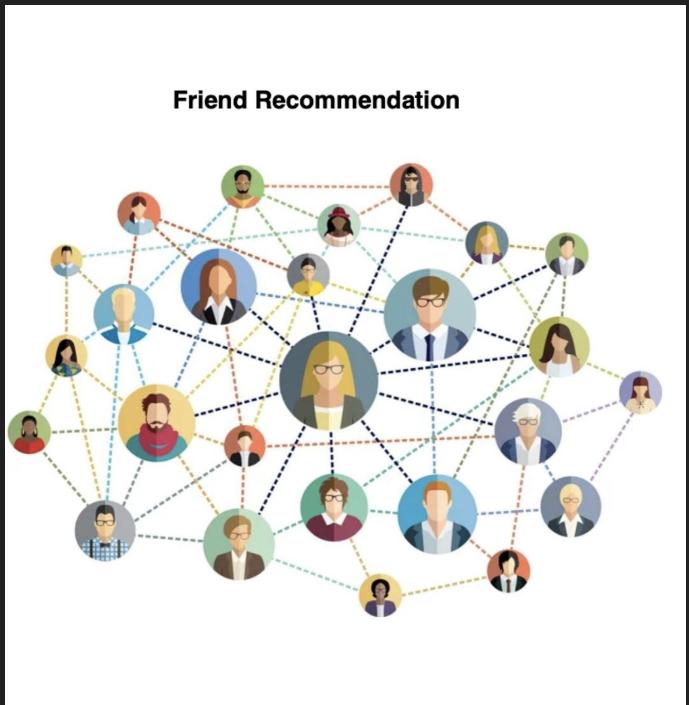
- Jiarong Li

Networks

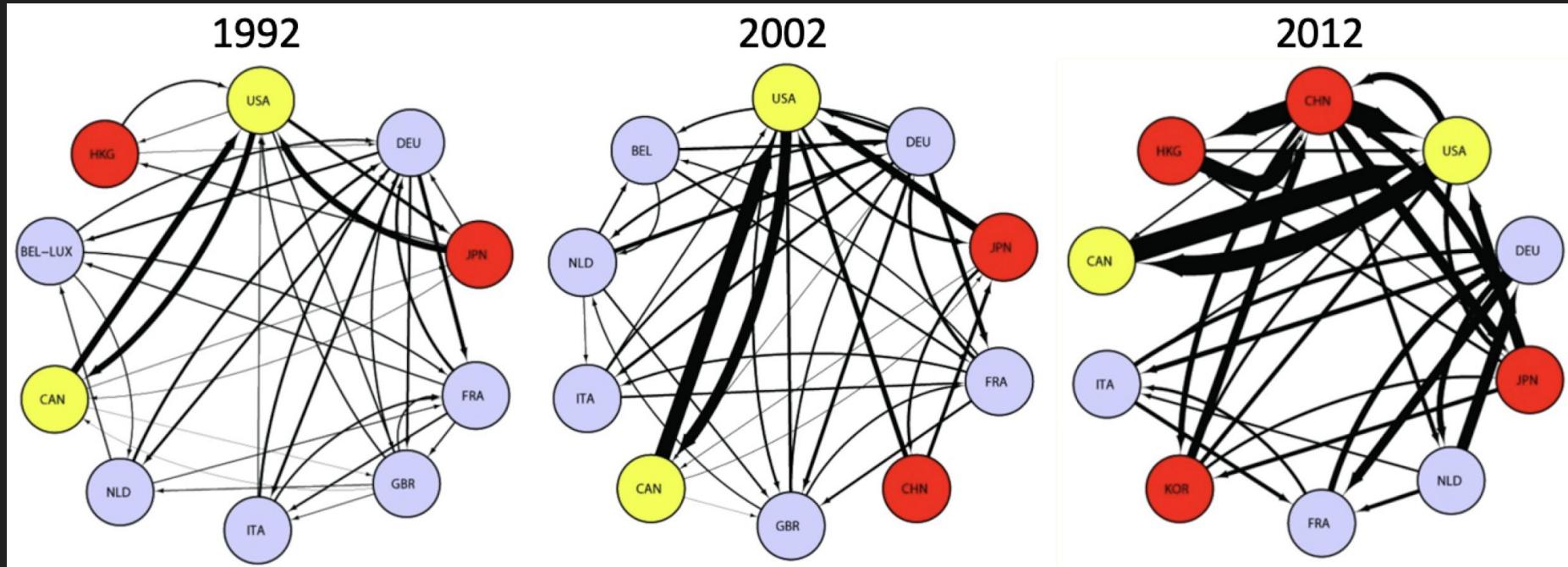
A very flexible and powerful way of modeling many real-world systems

It represents the relationships among individuals or among different entities in some collection.

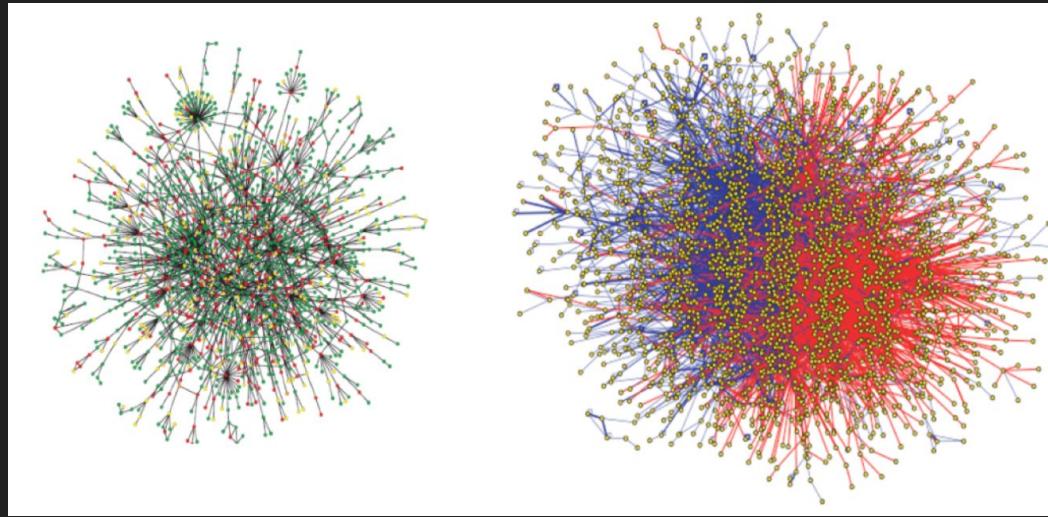
Social Networks



World Trade Networks



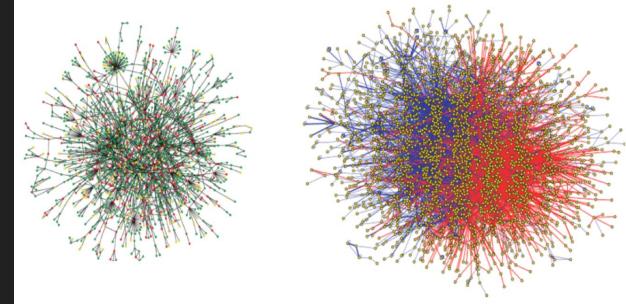
PPI Networks Reconstruction



Yeast protein Interactome

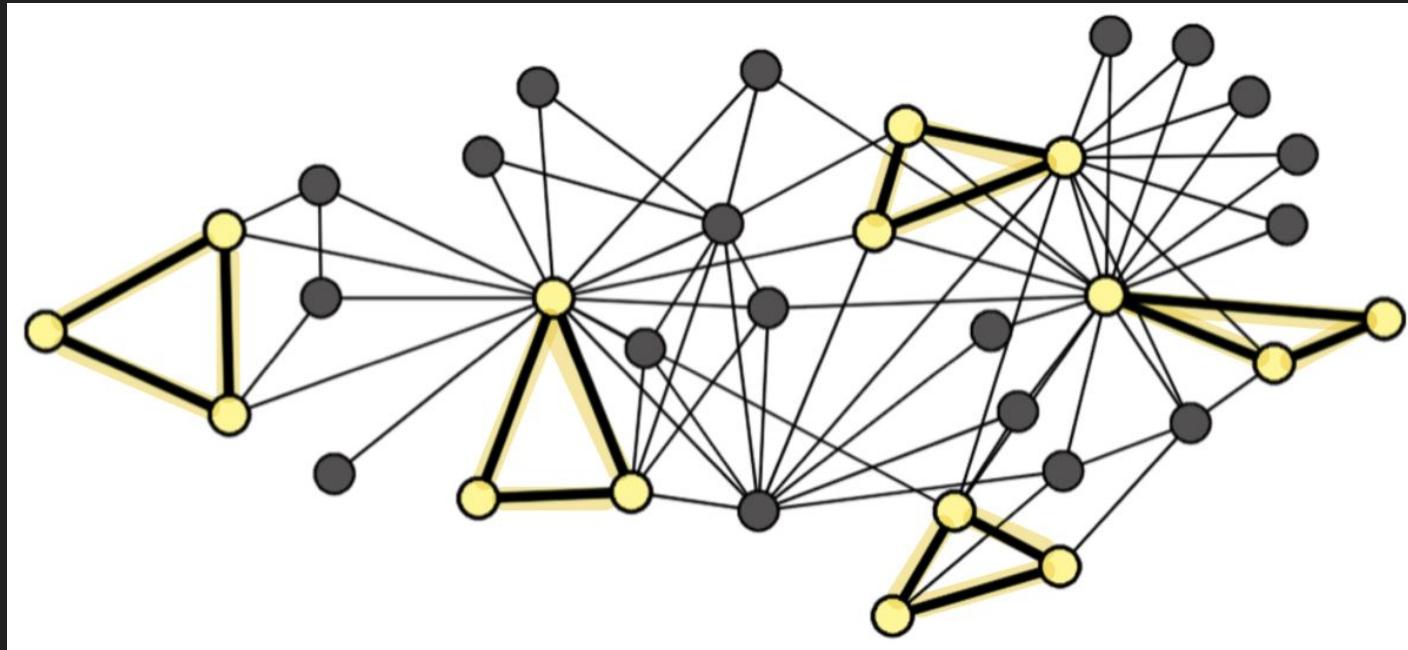
Human Interactome

How do we compare networks?

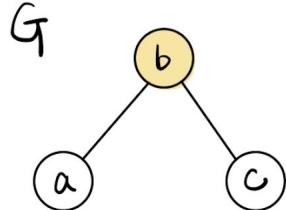


Substructures of a huge network: subgraphs

They give us power to characterize and discriminate networks



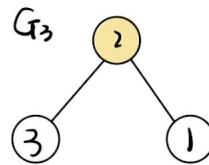
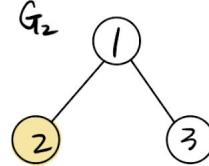
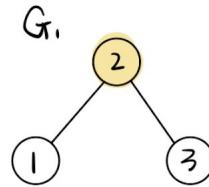
Graph Isomorphism



$$f_1: \begin{array}{l} a \mapsto 1 \\ b \mapsto 2 \\ c \mapsto 3 \end{array}$$

$$f_2: \begin{array}{l} a \mapsto 1 \\ b \mapsto 2 \\ c \mapsto 3 \end{array}$$

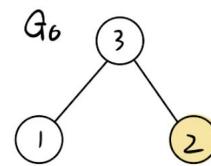
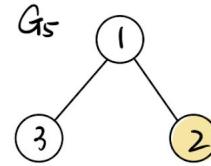
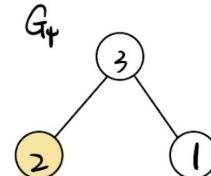
$$f_3: \begin{array}{l} a \mapsto 1 \\ b \mapsto 2 \\ c \mapsto 3 \end{array}$$



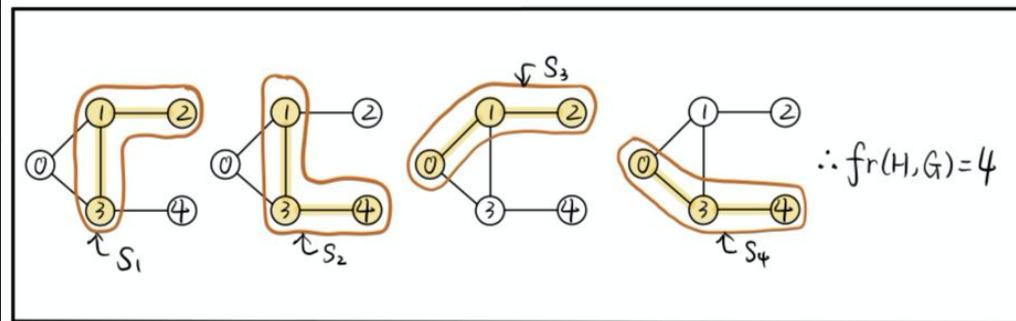
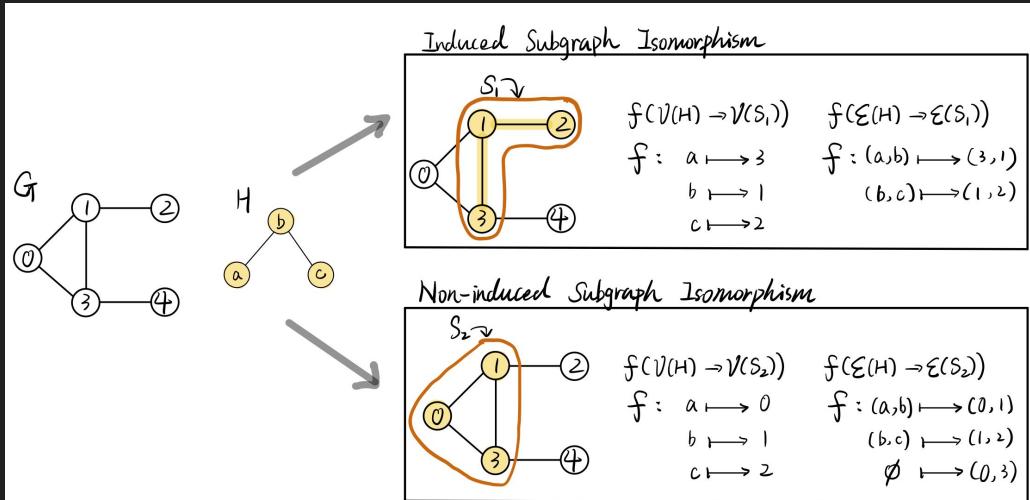
$$f_4: \begin{array}{l} a \mapsto 1 \\ b \mapsto 2 \\ c \mapsto 3 \end{array}$$

$$f_5: \begin{array}{l} a \mapsto 1 \\ b \mapsto 2 \\ c \mapsto 3 \end{array}$$

$$f_6: \begin{array}{l} a \mapsto 1 \\ b \mapsto 2 \\ c \mapsto 3 \end{array}$$

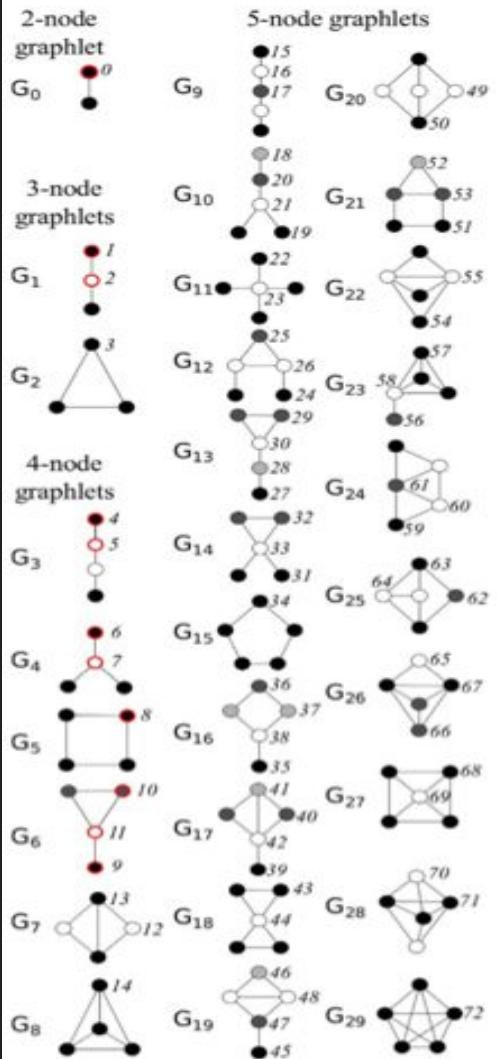


Subgraph Counting



Which subgraphs do we use to measure a huge network?

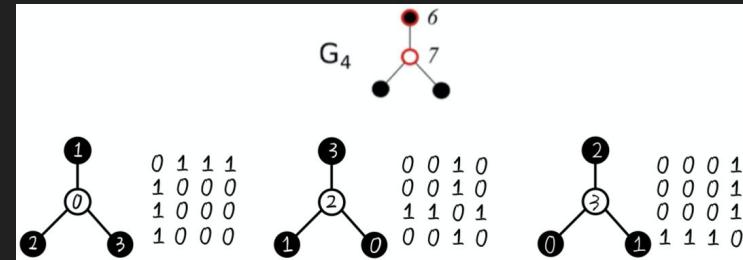
- Graphlets and Orbits



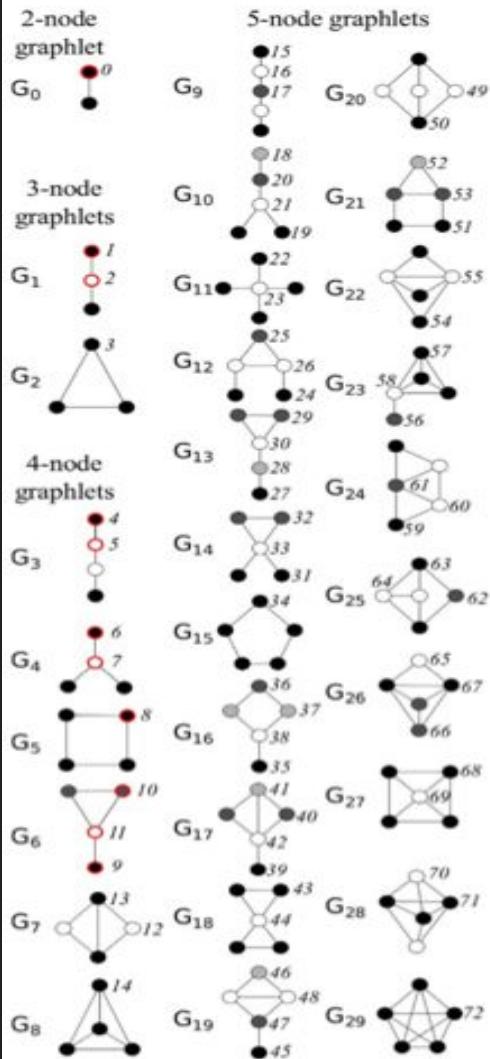
What are graphlets?

- Small induced non-isomorphic subgraphs of a large network that differentiate nodes according to their subgraph positions
- There are 30 graphlets of up to 5 nodes (G₀...G₂₉)

Matrix ID for each Graphlet



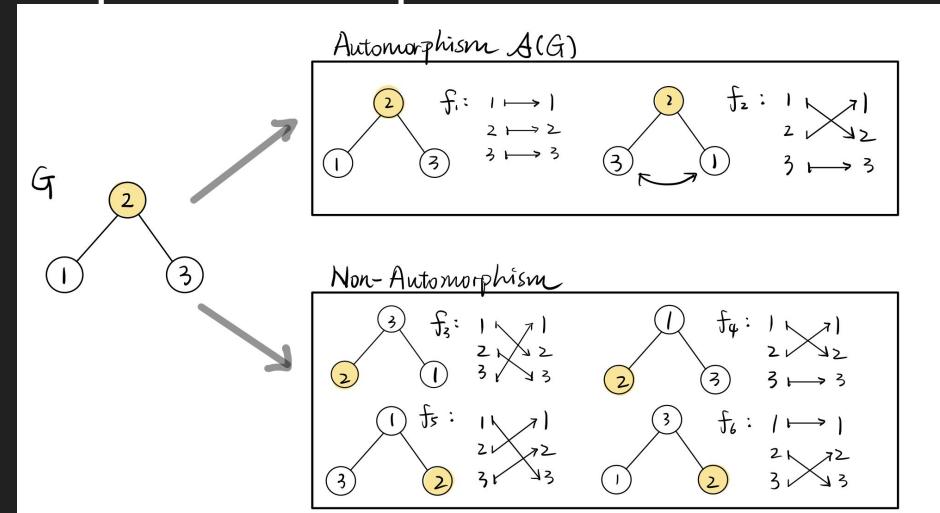
<i>Graphlets</i>	G ₀ 	G ₁ 	G ₂ 		
<i>Adjacency Matrix</i>	$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$	$\begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix}$	$\begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}$		
<i>Graphlets</i>	G ₃ 	G ₄ 	G ₅ 	G ₆ 	G ₇
<i>Adjacency Matrix</i>	$\begin{pmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}$	$\begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}$	$\begin{pmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}$	$\begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \end{pmatrix}$	$\begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$



What are automorphism orbits?

- Symmetry groups of nodes within graphlets
- A bijection of nodes that preserves node adjacency
- There are 73 automorphism orbits of graphlets of up to 5 nodes (red circles)

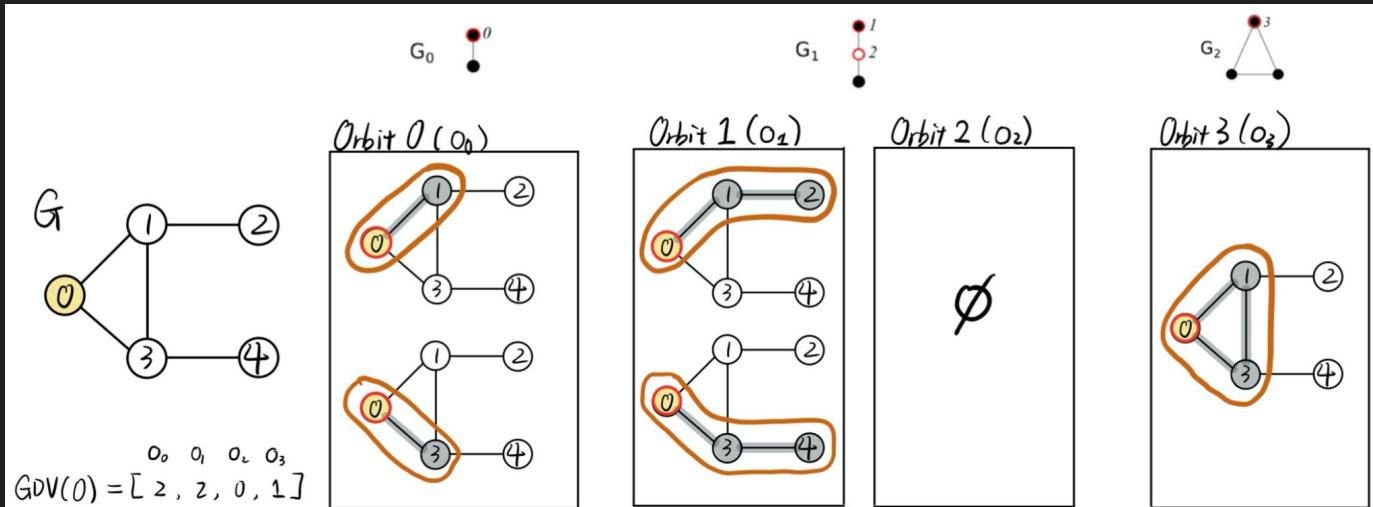
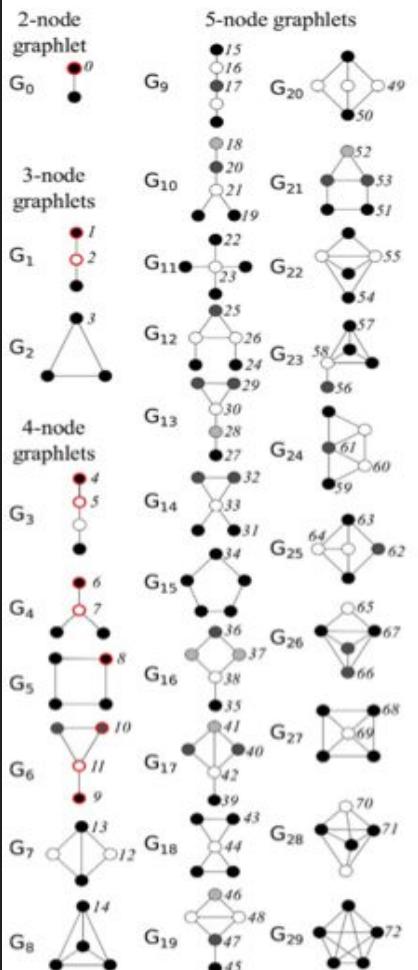
Graph Automorphism



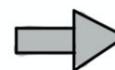
Quantifying networks' structural information

- Graphlets Degree Distribution (GDD)

Example of GDD calculation



$$Fr_q = \begin{bmatrix} O_0 & O_1 & O_2 & O_3 \\ 0 & 2 & 2 & 1 \\ 1 & 3 & 2 & 1 \\ 2 & 3 & 2 & 1 \\ 3 & 1 & 2 & 0 \\ 4 & 1 & 2 & 0 \end{bmatrix}$$



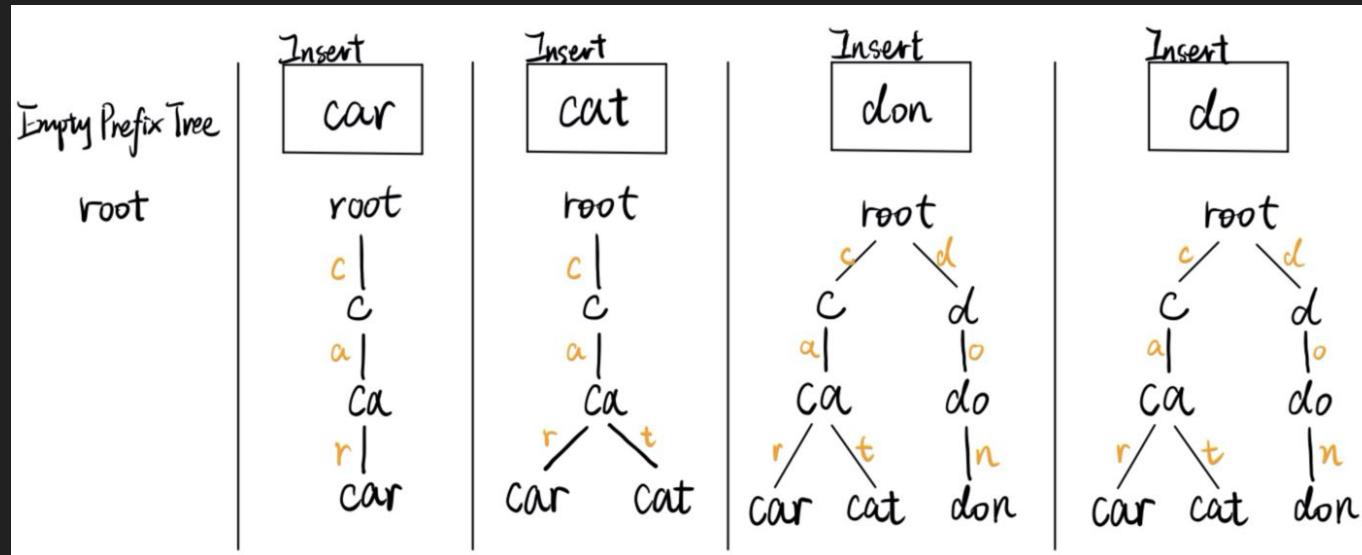
$$GDD_q = \begin{bmatrix} 1 & 2 & 3 \\ O_0 & 2 & 1 & 2 \\ O_1 & 0 & 5 & 0 \\ O_2 & 2 & 0 & 0 \\ O_3 & 3 & 0 & 0 \end{bmatrix}$$

How do we count graphlets efficiently?

Trie (Prefix Tree)

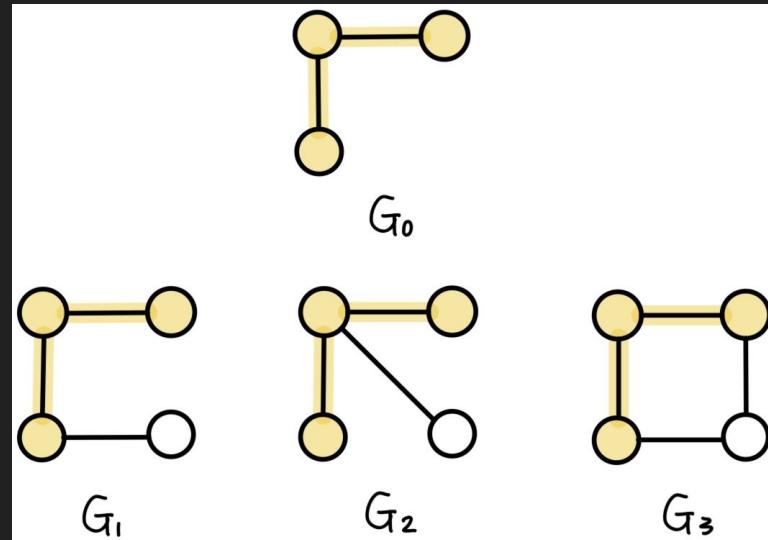
A tree data structure used for locating specific keys within a set.

E.g. Looking for a set of substrings in a given long string

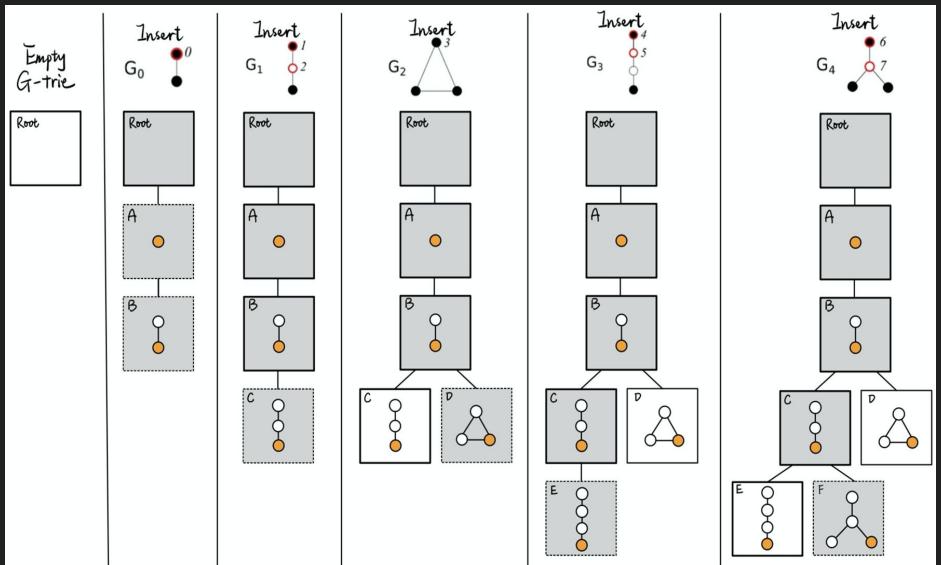


G-Trie (Prefix Tree for Graphs)

ca
cat car



G-Tries - A data-structure to store and enumerate subgraphs



Algorithm 1 Populate a g-trie T with subgraphs $G \in GSet$

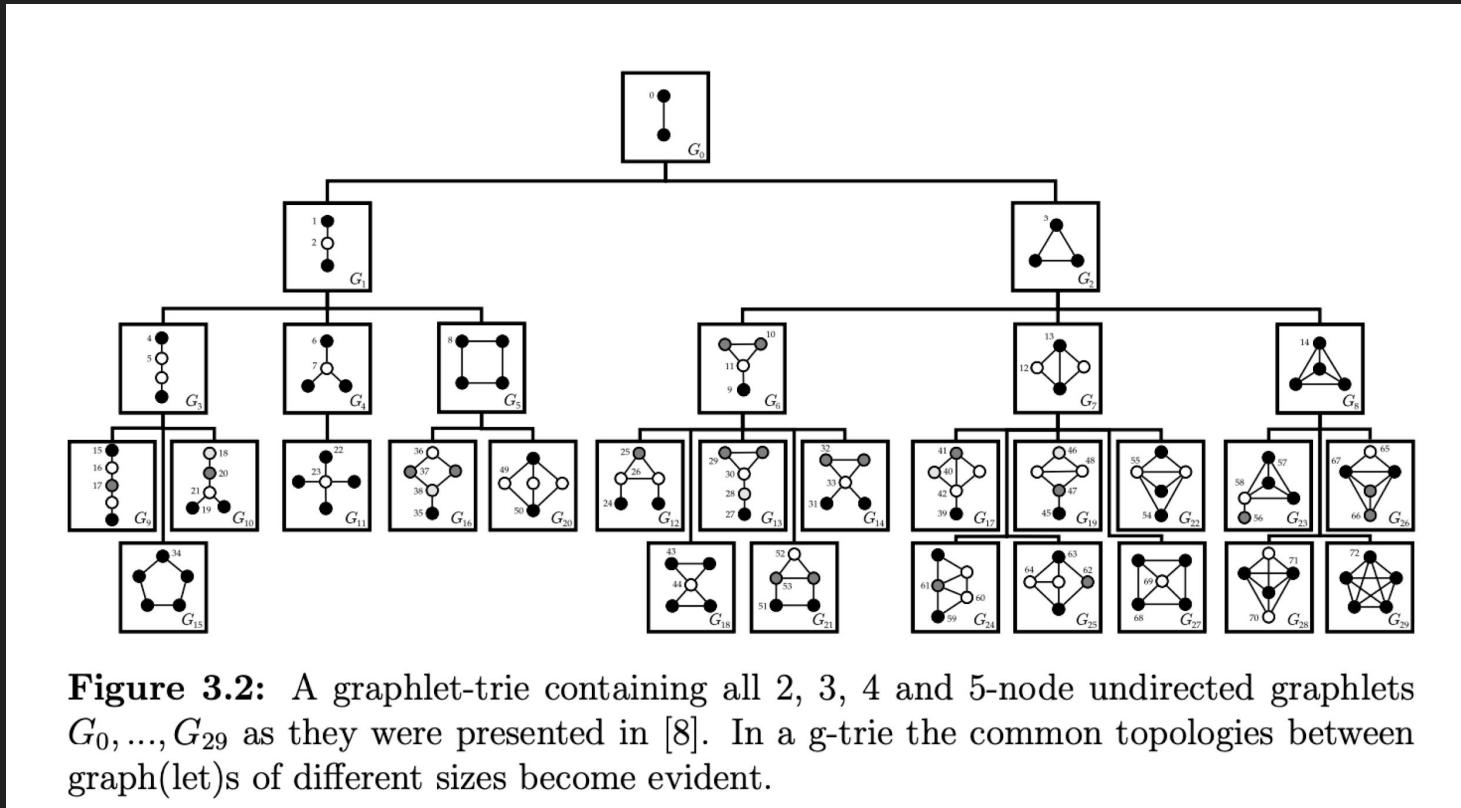
```

1: procedure CREATEGTRIE( $GSet$ )
2:    $T \leftarrow \text{EmptyGtrie}$ 
3:   for all  $G \in GSet$  do
4:     INSERT( $T.\text{root}, G, 0$ )
5:   return  $T$ 
6: procedure INSERT( $V, G, depth$ )
7:   if  $k < \text{numberRows}(G.\text{MatrixID})$  then
8:     for all  $C \in V.\text{children}$  do
9:       if SHARECOMMONTOPOLOGY( $C, G, depth$ ) then
10:         INSERT( $C, G, depth + 1$ )
11:       return
12:      $NewChild \leftarrow N.\text{ADDCHILD}(G, depth)$ 
13:     INSERT( $NewChild, G, depth + 1$ )

```

Starts by looking for occurrences of the smaller common subgraph to save searching time.

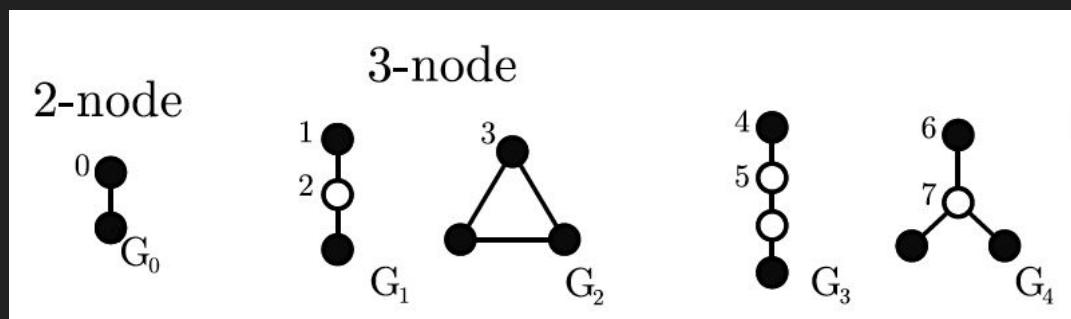
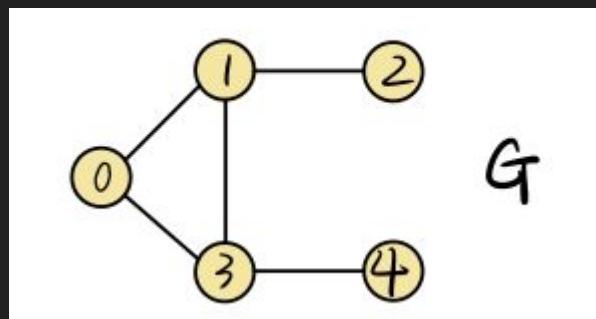
Graphlets-tries - g-tries for graphlets



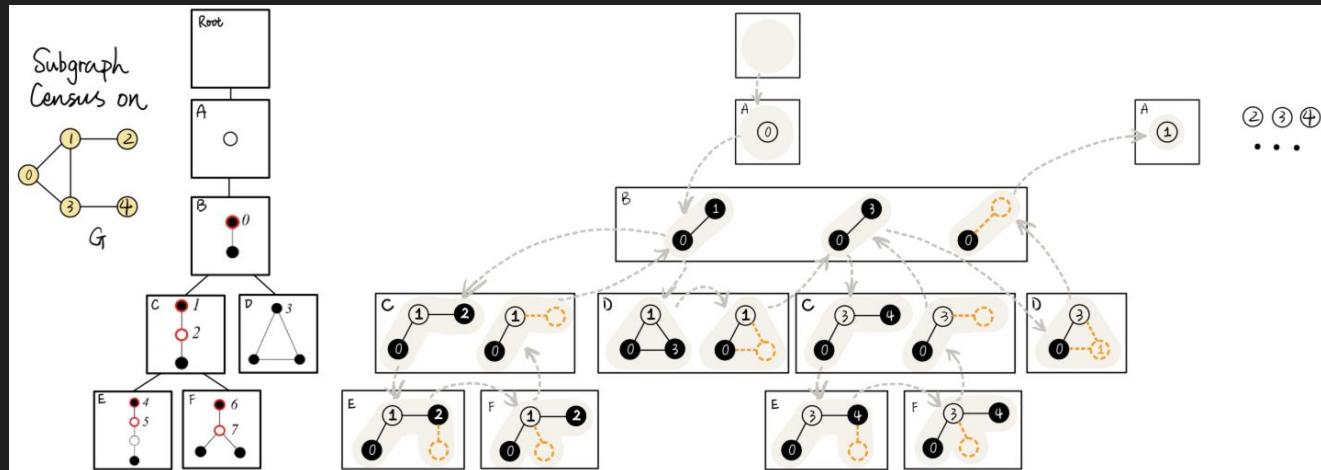
Graphlet Census using Graphlet-trie (GT-Scanner)

Let's count the frequencies of the first 5 graphlets in graph G.

We will use graphlet-trie as the data structure to perform frequencies counting...

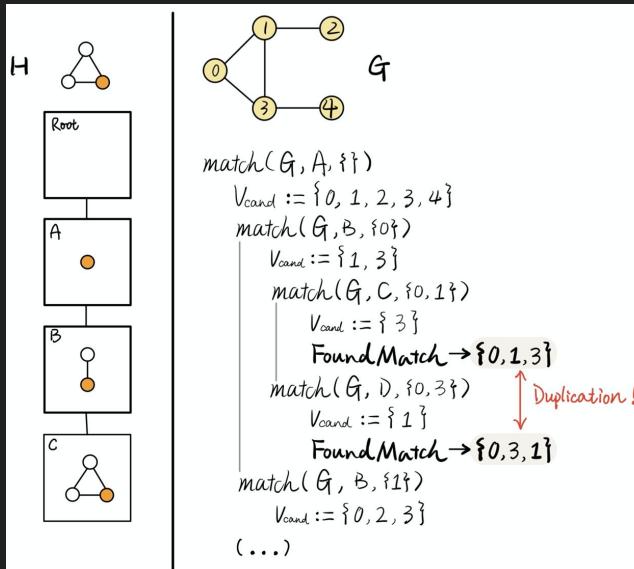


Graphlet Census using Graphlet-trie (GT-Scanner)



GT-Scanner's Symmetry Breaking Condition:

Vertices that appear later in the same orbit are only valid if they have a bigger index than the previous vertices of the same orbit.



Algorithm 2 Graphlets Census using g-trie T in network G

```

1: procedure CENSUS( $G, T$ )
2:   for all children  $c$  of  $T.root$  do
3:     MATCH( $G, c, \emptyset$ )
4: procedure MATCH( $G, T, V_{used}$ )
5:    $V_{cand} \leftarrow \text{GETVALIDCANDIDATES}(G, T, V_{used})$ 
6:   for all vertex  $v \in V_{cand}$  do
7:     if  $T.\text{FOUNDMATCH}(V_{used} \cup \{v\})$  then
8:       output  $(V_{used} \cup \{v\})$ 
9:   for all children  $c$  of  $T$  do
10:    MATCH( $G, c, V_{used} \cup \{v\}$ )
11: function GETVALIDCANDIDATES( $G, T, V_{used}$ )
12:   if  $V_{used} = \emptyset$  then  $V_{cand} \leftarrow \mathcal{V}(G)$ 
13:   else
14:      $V_{conn} \leftarrow$  vertices connected to the vertex being added
15:      $m \leftarrow$  vertex of  $V_{conn}$  with smallest neighborhood
16:      $V_{cand} \leftarrow$  neighbors of  $m$  that respect both
17:       connections to ancestors and
18:       symmetry breaking condition
19:   return  $V_{cand}$ 

```

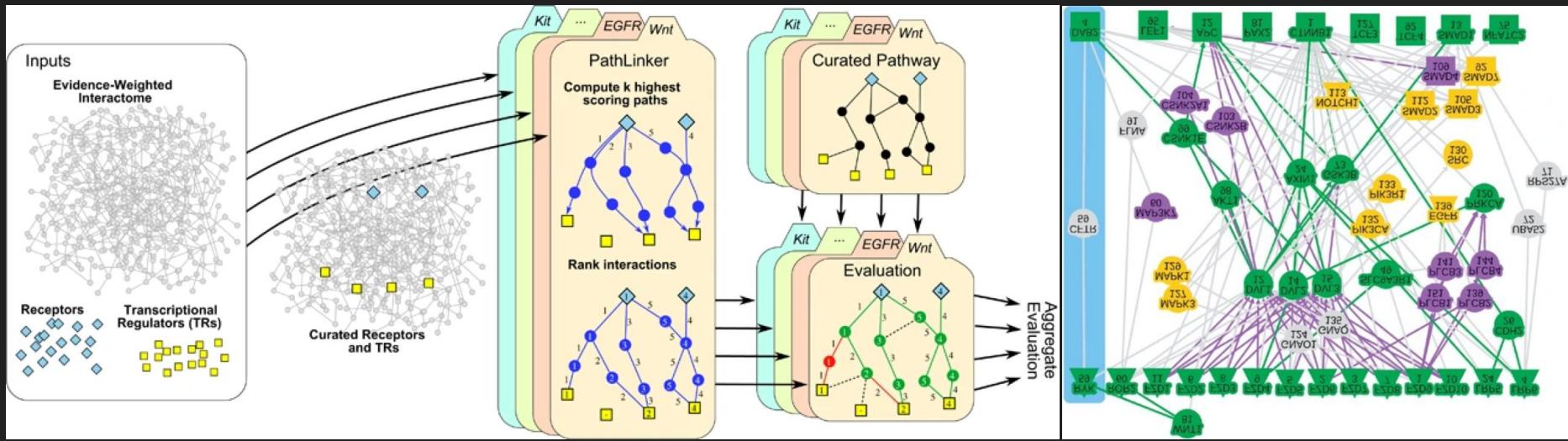
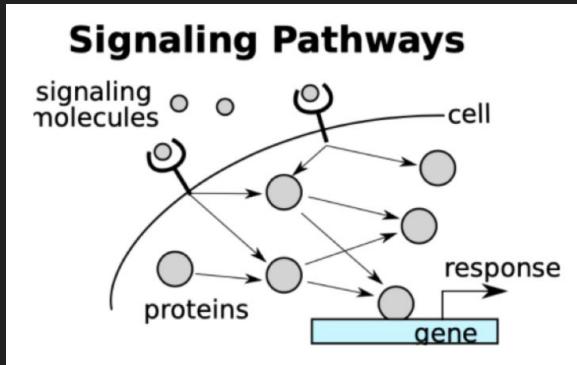
Experimental Comparisons

- Wnt Signaling Pathway
- Manually Curated: 106 vertices and 220 edges (gained from NetPath database)
- Reconstruct Wnt pathway from Interactome

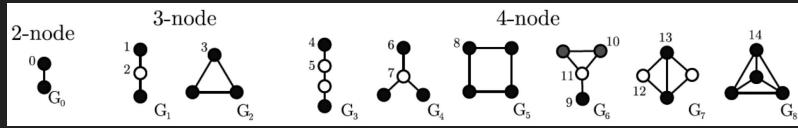
Curated vs PathLinker

Glance of PathLinker

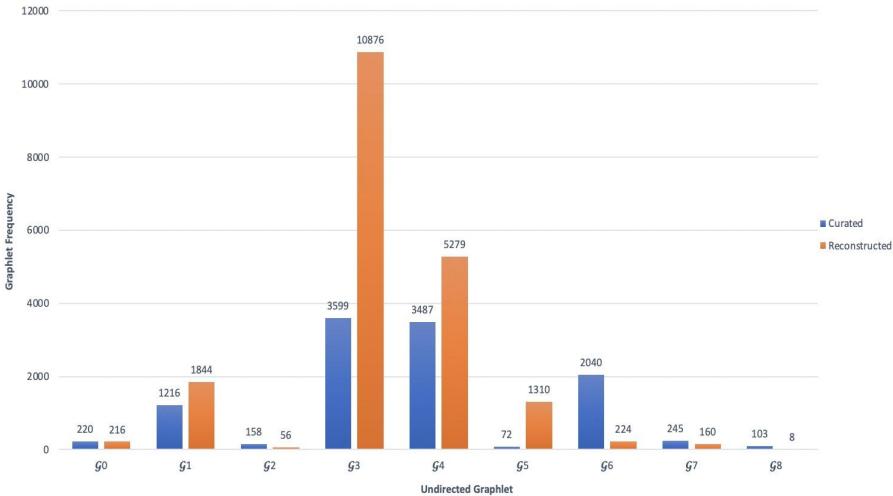
Given as input a set of sources and a set of targets for a particular curated pathway, PathLinker reconstructs the pathway by finding the k shortest paths from any source to any target within the interactome.



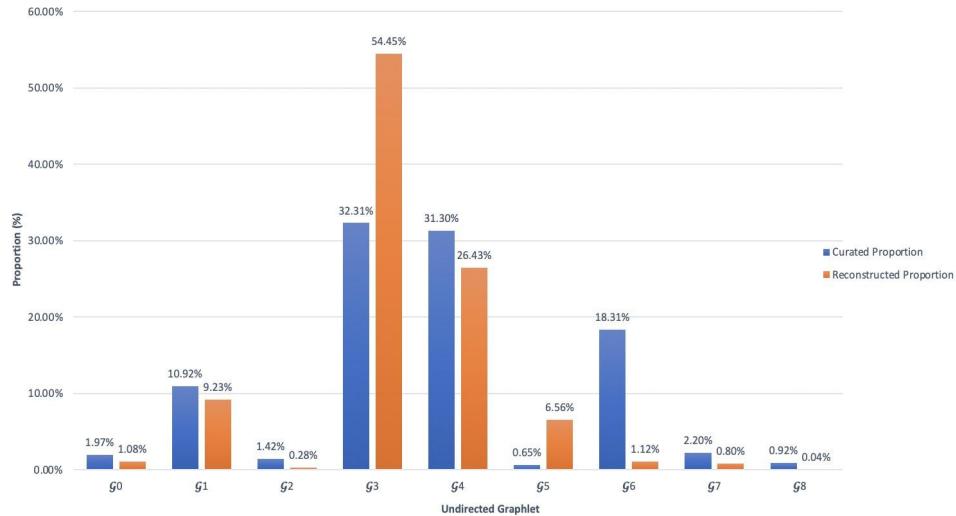
Curated Pathway vs PathLinker's Reconstruction (Wnt, k = 200)



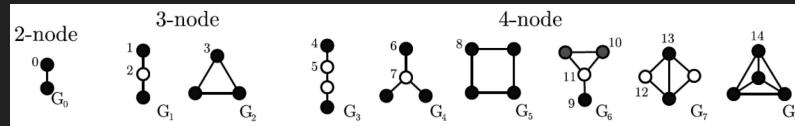
Graphlet Frequencies Comparison (k = 200, Wnt)



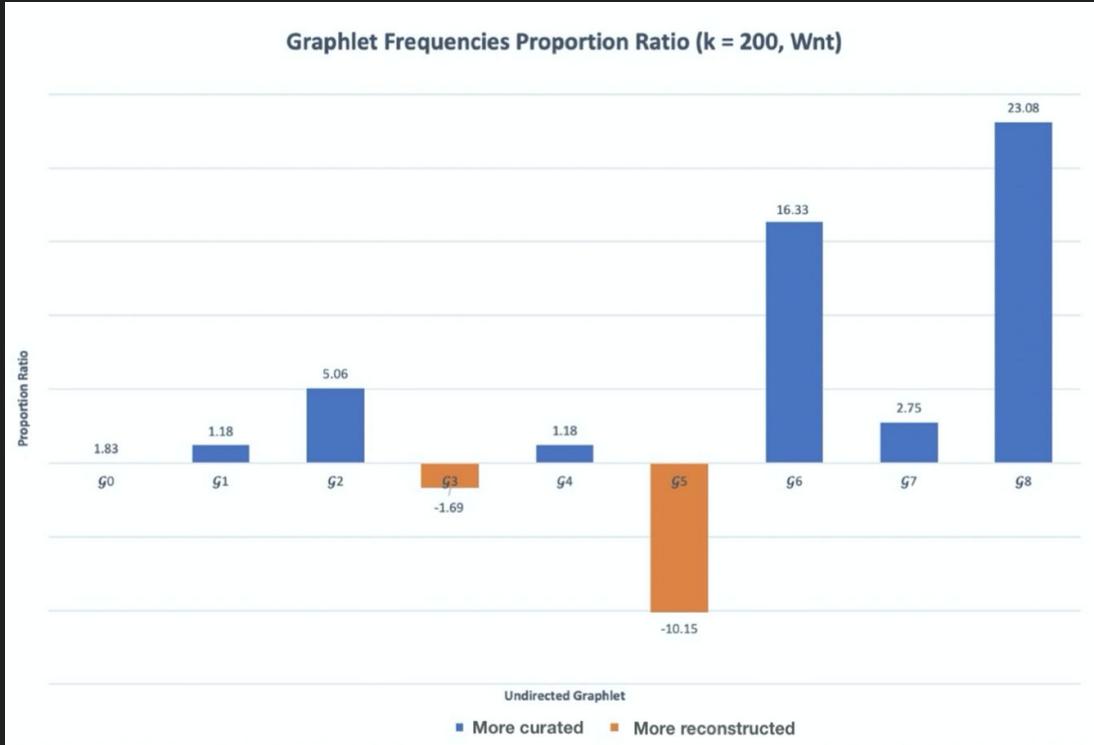
Graphlet Frequencies Proportion Comparison (k = 200, Wnt)



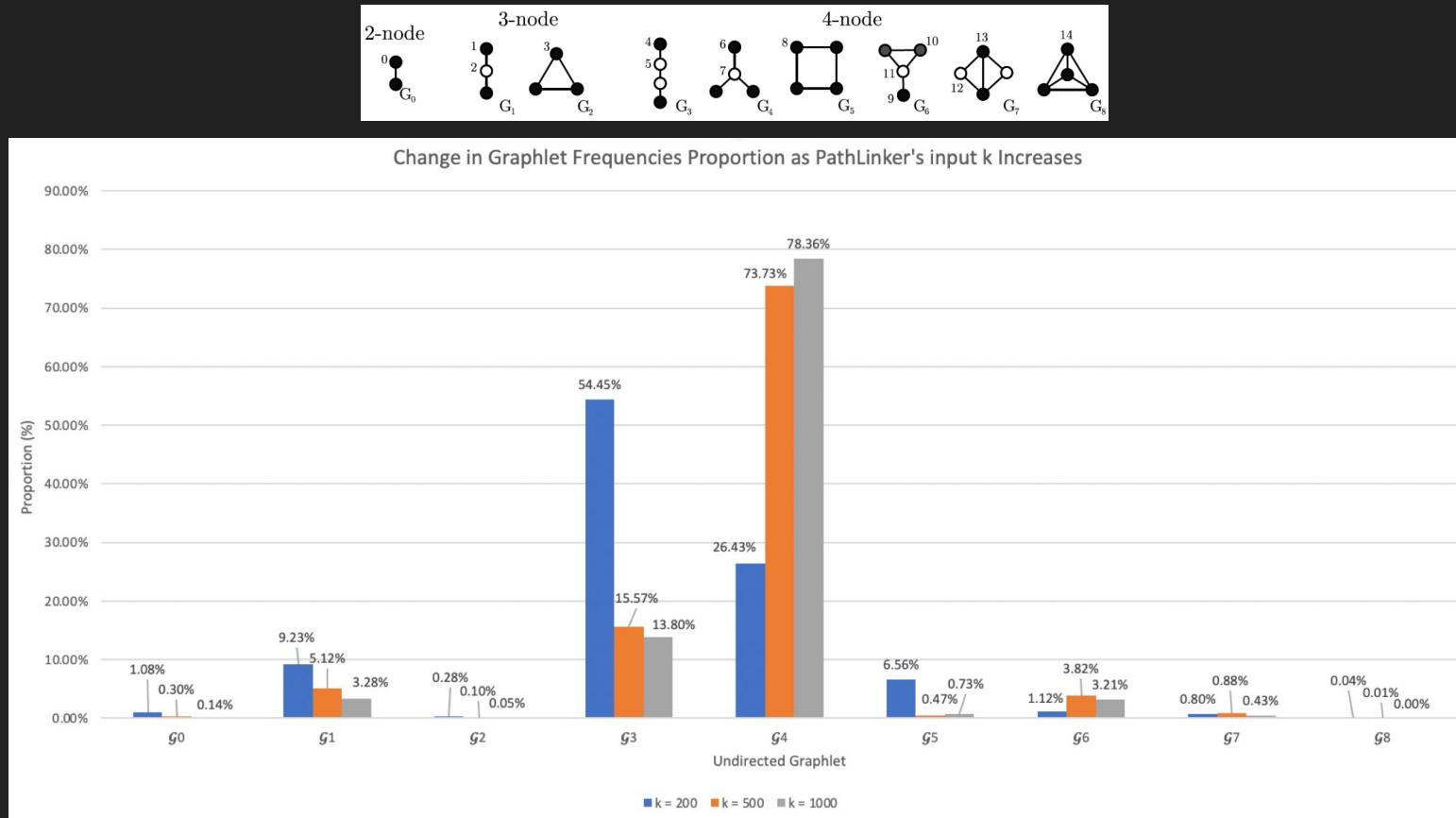
Curated Pathway vs PathLinker's Reconstruction (Wnt, k = 200)



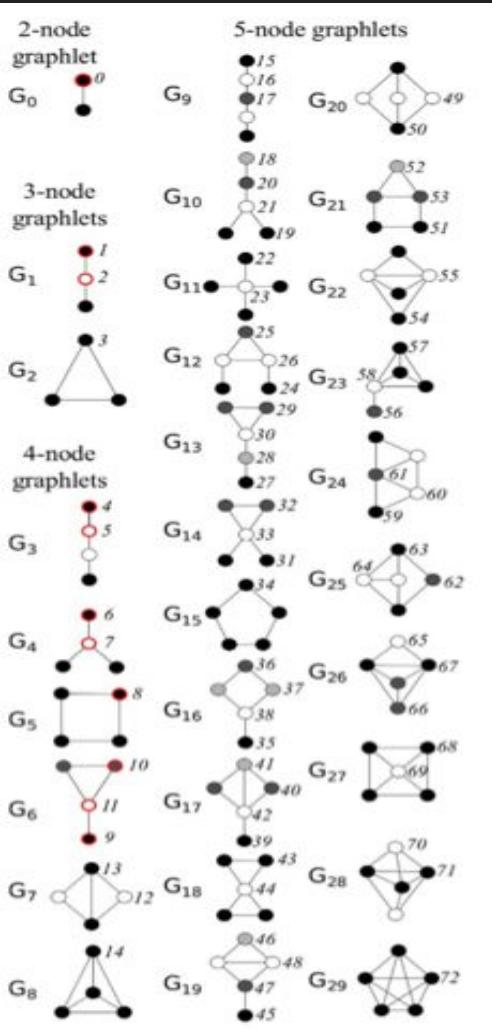
$$\frac{\max(P_c(\mathcal{G}_n), P_r(\mathcal{G}_n))}{\min(P_c(\mathcal{G}_n), P_r(\mathcal{G}_n))} \times \text{sign}(P_c(\mathcal{G}_n) - P_r(\mathcal{G}_n))$$



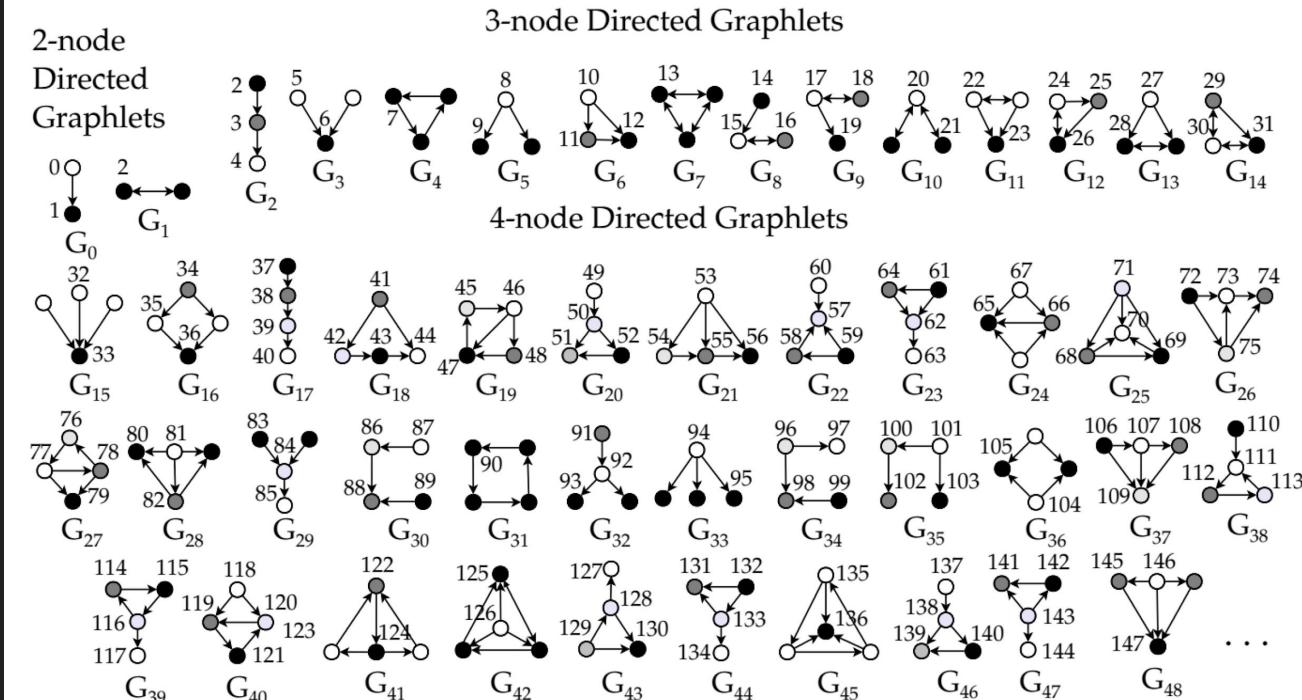
Comparisons among PathLinker's Reconstructions (Wnt, k = 200, 500, 1000)



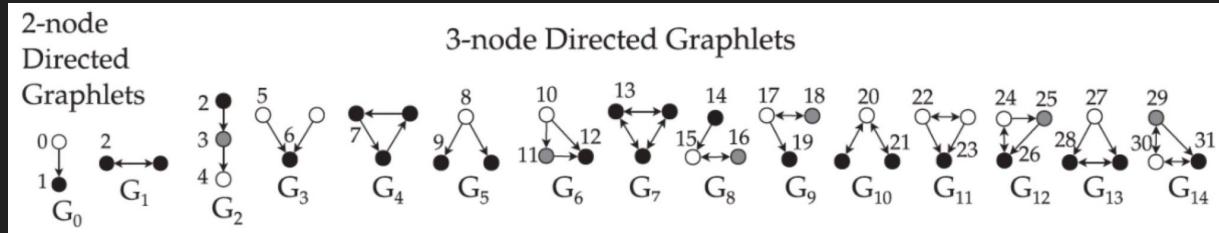
Some other fun facts about graphlets...



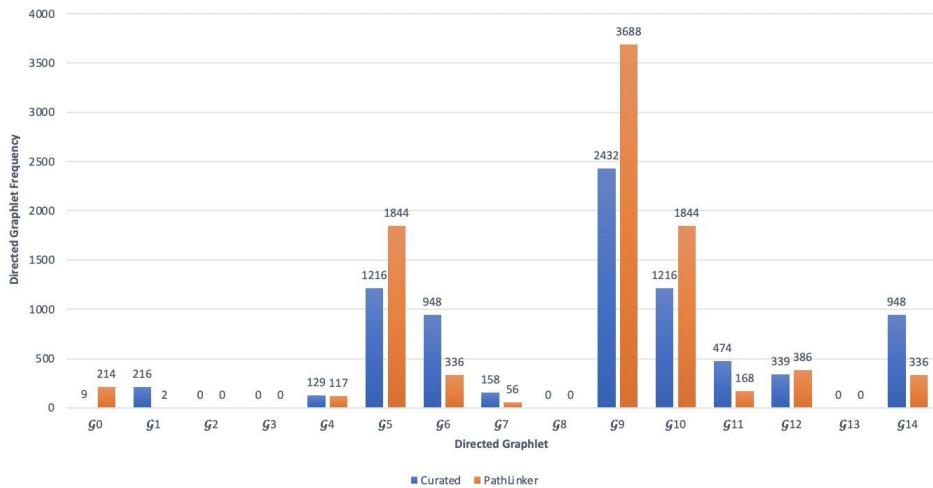
Undirected Graphlets vs Directed Graphlets



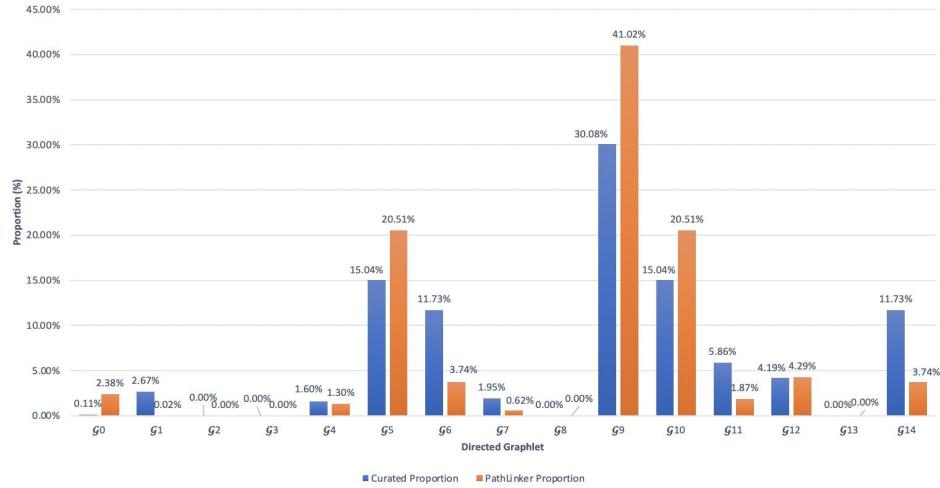
Curated vs PathLiker using Directed Graphlets (Wnt, k = 200)



Directed Graphlet Frequency Comparison (Wnt, k = 200)



Directed Graphlet Frequency Proportion Comparison (Wnt, k = 200)



Curated vs PathLiker using Directed Graphlets (Wnt, k = 200)

