

Data Structures

Path Compression Proof and Graph Fundamentals

CS 225

March 27, 2026

Brad Solomon



UNIVERSITY OF
ILLINOIS
URBANA - CHAMPAIGN

Department of Computer Science

Learning Objectives

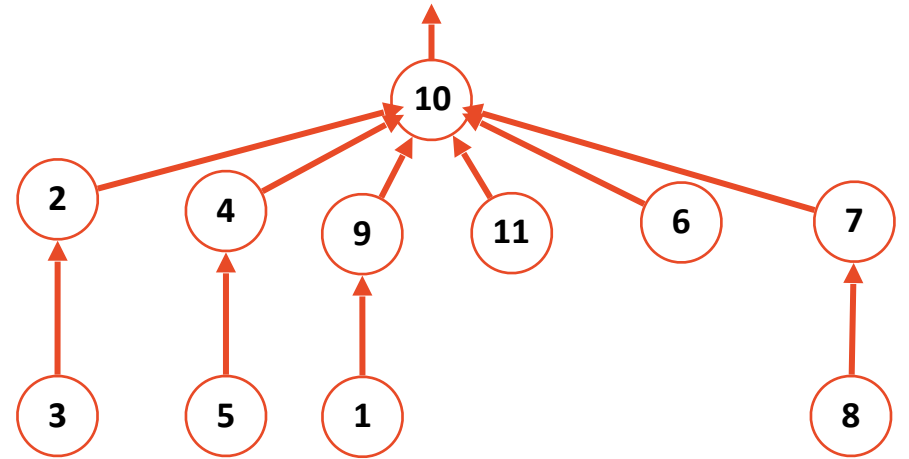
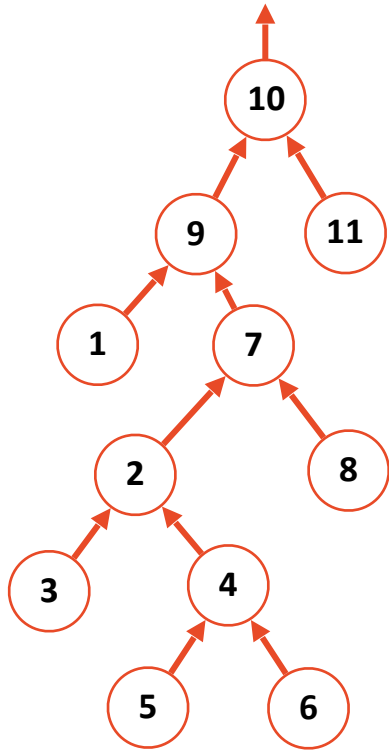
See a high level summary of path compression proof

Define graph vocabulary

Discuss graph implementation and storage strategies

Path Compression

Find(6)



This seems good — but how good in theory?

Path Compression Analysis

Two major problems here:

- 1) Our efficiency changes ***over repeated calls to find()***
- 2) Our height changes so we cant use union by height

Amortized Time Review

We have **n items**. We make **n insert()** calls.

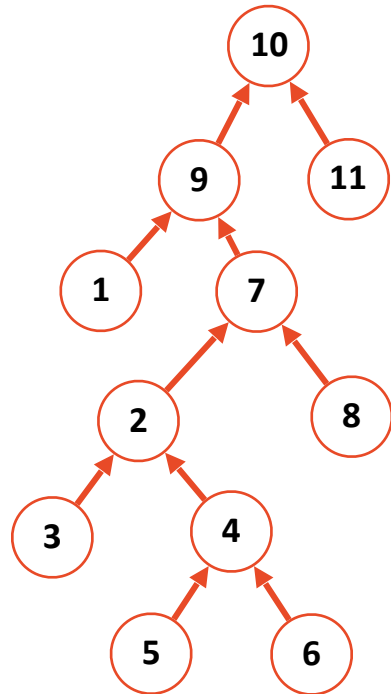
We are interested in the **worst case work** possible **over n calls**.



Amortized Time (Path Compression)

We have **n items** in an Uptree. We make **m find()** calls.

We are interested in the **worst case work** possible **over m calls**.



Union by Rank (Not Height)

Once I do path compression, I change the height of tree!

So we need a new way of approximating height.

Rank is a way of remembering what our height was before P.C.

Union by Rank (Not Height)

New UpTrees have rank = 0

Let A, B be two sets being unioned. If:

rank(A) == rank(B): The merged UpTree has rank + 1

rank(A) > rank(B): The merged UpTree has rank(A)

rank(B) > rank(A): The merged UpTree has rank(B)

Key Properties of UpTree by rank w/ PC

The parent of a node is always higher rank than the node.

There are at least $\geq 2^r$ nodes in a root of rank r .

For any integer r , there are at most $\frac{n}{2^r}$ nodes of rank r .

Key Properties of UpTree by rank w/ PC

The parent of a node is always higher rank than the node.

This comes from how we set up rank union

(Take larger of two rank or add one if tied)

There are at least $\geq 2^r$ nodes in a root of rank r .

Proof by Induction: To create rank r set, we merge two $r - 1$ sets

By IH (not shown), those sets have $2^{r-1} + 2^{r-1} = 2^r$ nodes

For any integer r , there are at most $\frac{n}{2^r}$ nodes of rank r .

A rewrite of the above logic given n nodes

Amortized Time (Rank w/ Path Compression)

Put every non-root node in a bucket by rank!

Structure buckets to store ranks $[r, 2^r - 1]$

Where did number range come from?

Ranks	Bucket
0	0
1	1
2 - 3	2
4 - 15	3
16 - 65535	4
$65536 - 2^{\{65536\}} - 1$	5

Iterated Logarithm Function ($\log^* n$)

The number of times you can take a log of a number

$$\log^*(n) = \begin{cases} 0 & , n \leq 1 \\ 1 + \log^*(\log(n)) & , n > 1 \end{cases}$$

$$\log^*(2^{65536}) = 5$$

$$2^{65536}$$

$$2^{16} = 65536$$

$$2^4 = 16$$

$$2^2 = 4$$

$$2^1 = 2$$

$$2^0 = 1$$

Amortized Time (Rank w/ Path Compression)

The work of **find(x)** are the steps taken on the path from a node x to the root (or immediate child of the root) of the UpTree containing x

We can split this into two cases:

Case 1: We take a step from one bucket to another bucket.

Case 2: We take a step from one item to another inside the same bucket.

Amortized Time (Rank w/ Path Compression)

The work of **find(x)** are the steps taken on the path from a node x to the root (or immediate child of the root) of the UpTree containing x

We can split this into two cases:

Case 1: We take a step from one bucket to another bucket.

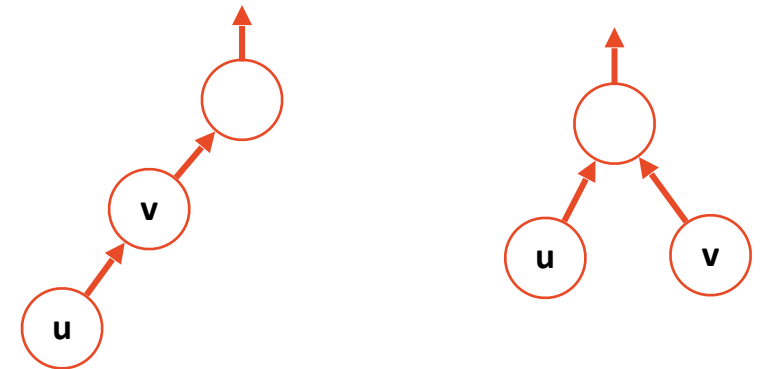
We have at most $\log^*(n)$ buckets so for m finds, this is $O(m \log^* n)$

Case 2: We take a step from one item to another inside the same bucket.

Let's call this the step from u to v .

Every time we do this, we do path compression:

We set $\text{parent}(u)$ a little closer to root



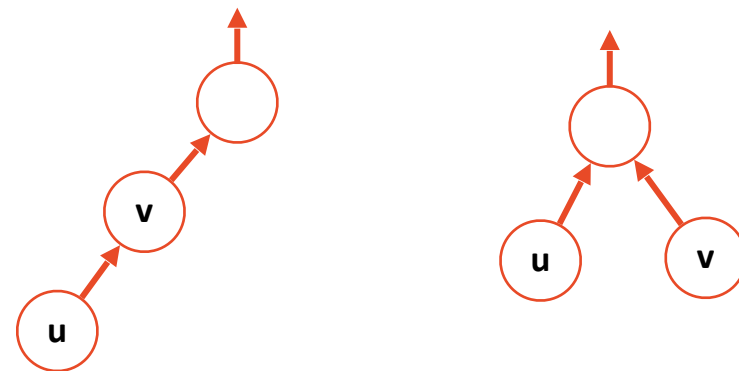
Amortized Time (Rank w/ Path Compression)

Case 2: We take a step from one item to another *inside* the same bucket.

Let's call this the step from **u** to **v**.

Every time we do this, we do path compression:

We set $\text{parent}(u)$ a little closer to root



How many total times can I do this for each **u** in a bucket?

By definition of our bucket ranges $\sim 2^r$

How many nodes are in bucket **r**?

By definition of how we set up rank: $\frac{n}{2^r}$

Given we have $\log^*(n)$ buckets:

Case 2 work is $n \log^*(n)$

Final Result



We have **n items** in an Uptree. We make **m find()** calls. Total work is:

Amortized $(n + m) \log^* (n)$

In terms of real world data, this is practically a constant.

Alternative Not-Actually-A-Proof

Unproven Claim: A disjoint set implemented with smart union and path compression with **m** find calls and **n** items has a worst case running time of **inverse Ackerman**. $[O(m \alpha(n))]$

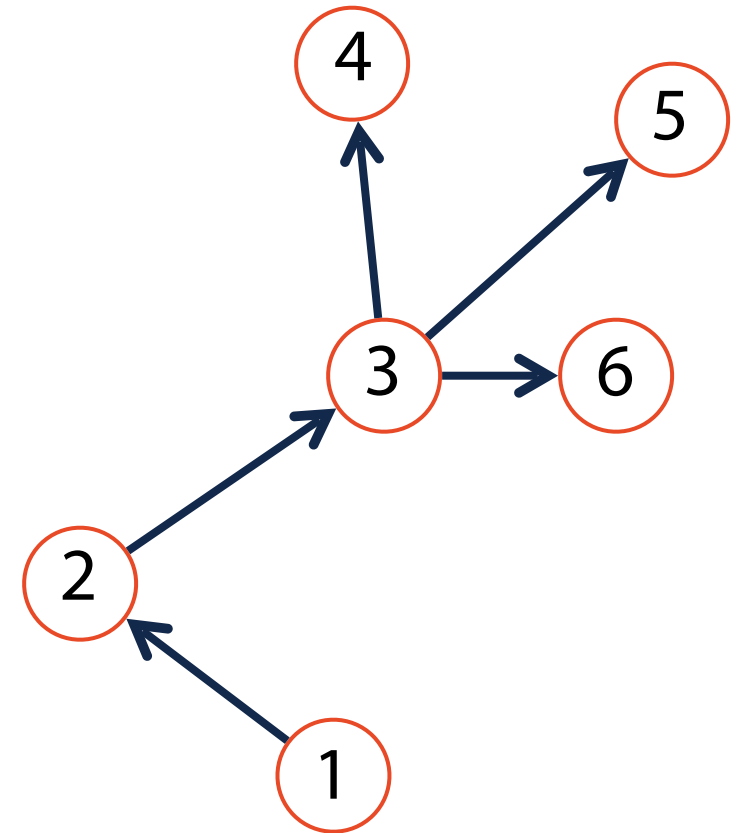
This grows *very* slowly to the point of being treated a constant in CS.

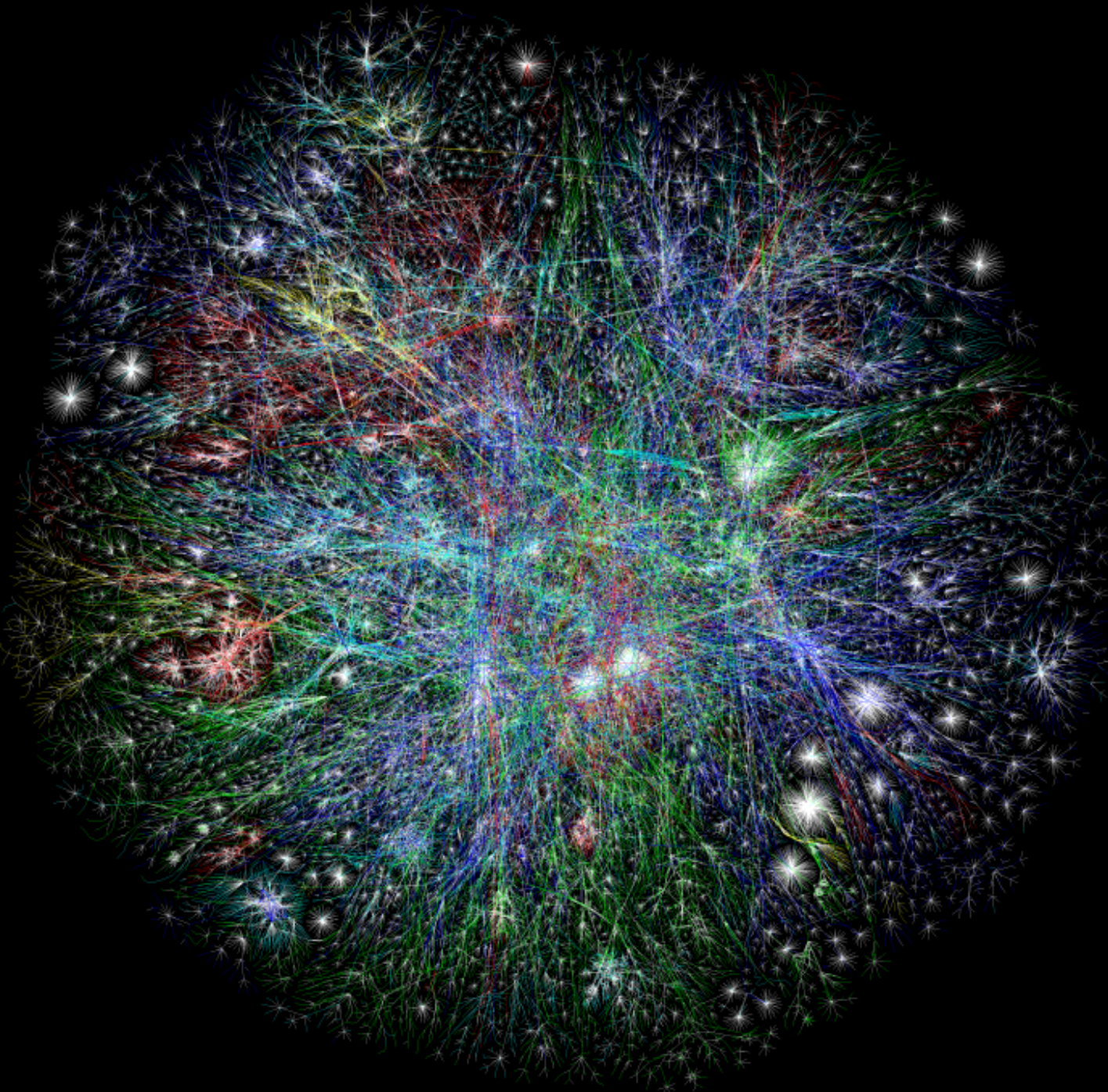
Whats next?

A non-linear data structure defined recursively as a collection of nodes where each node contains a value and zero or more connected nodes.

(In CS 225) a tree is also:

- 1) Acyclic — contains no cycles
- 2) Rooted — root node connected to all nodes



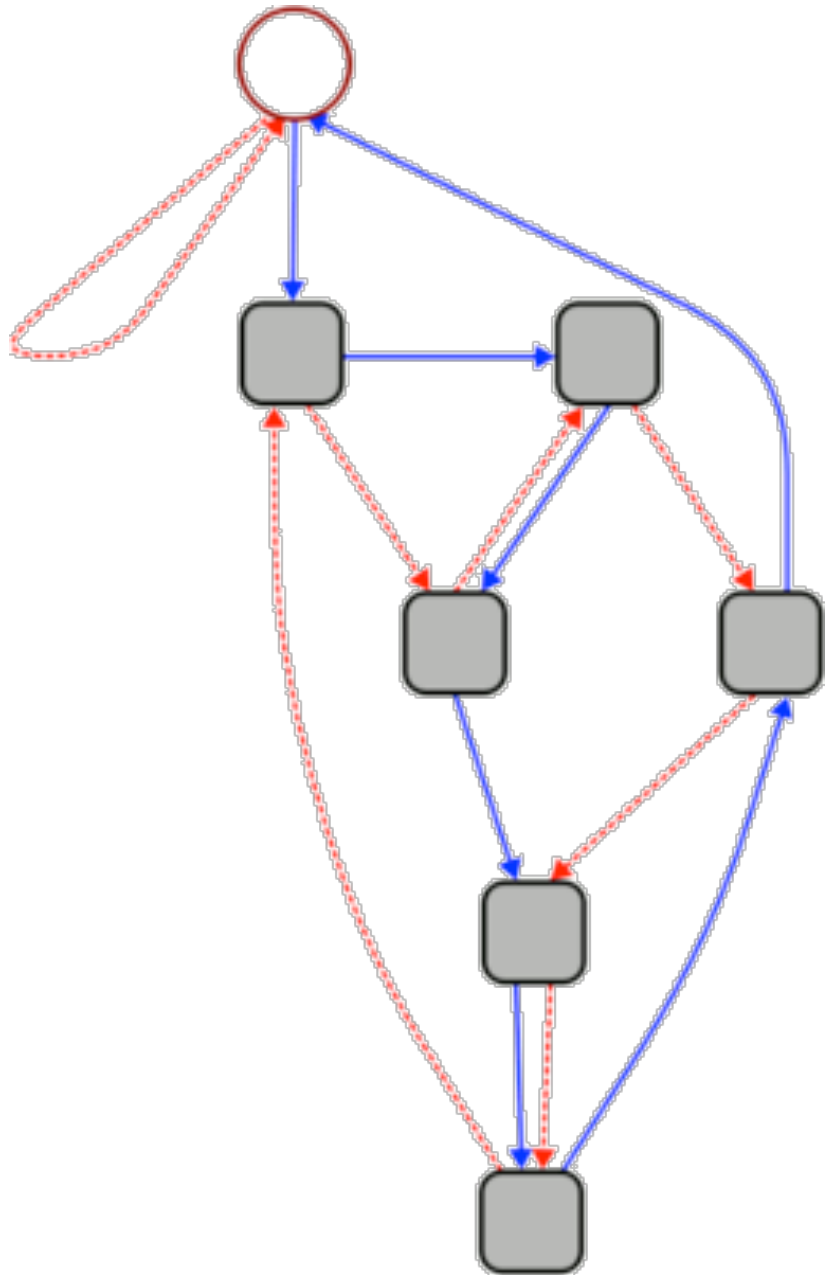


Nodes: Routers and servers

Edges: Connections

The Internet 2003

[The OPTE Project](#) (2003)



This graph can be used to quickly calculate whether a given number is divisible by 7.

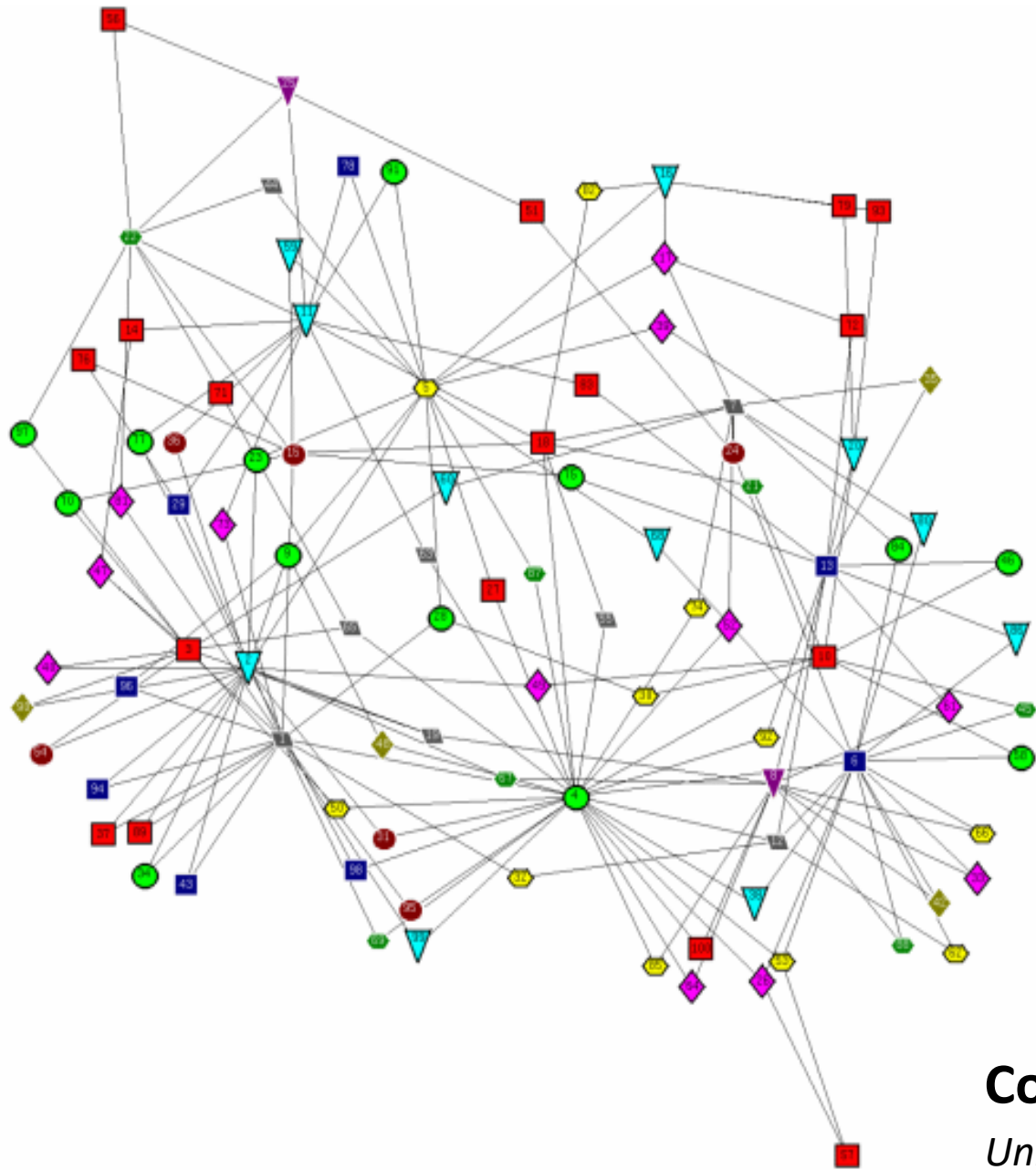
1. Start at the circle node at the top.
2. For each digit **d** in the given number, follow **d blue (solid) edges** in succession. As you move from one digit to the next, follow **1 red (dashed) edge**.
3. If you end up back at the circle node, your number is divisible by 7.

3703

“Rule of 7”

Unknown Source

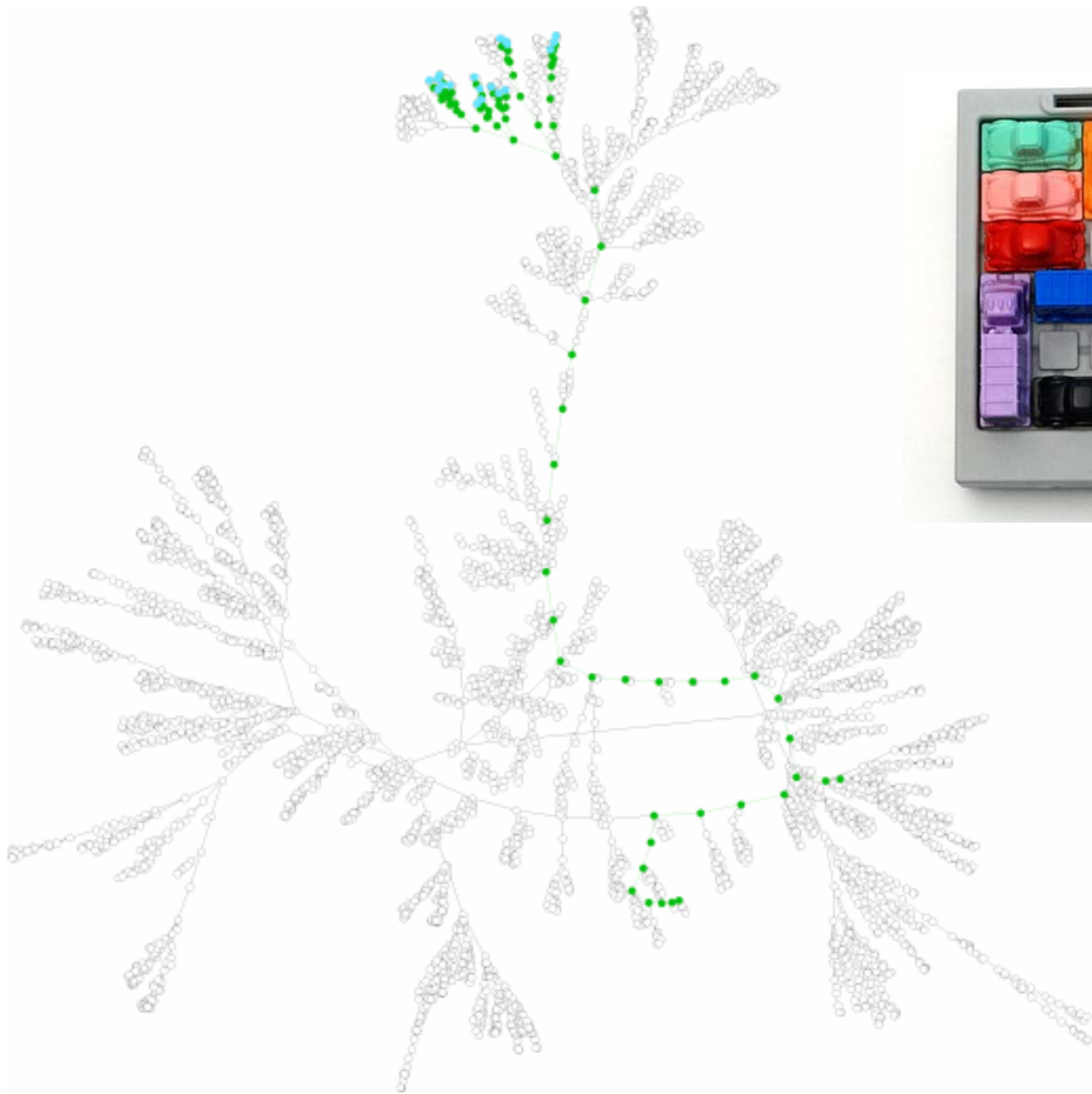
Presented by Cinda Heeren, 2016



Conflict-Free Final Exam Scheduling Graph

Unknown Source

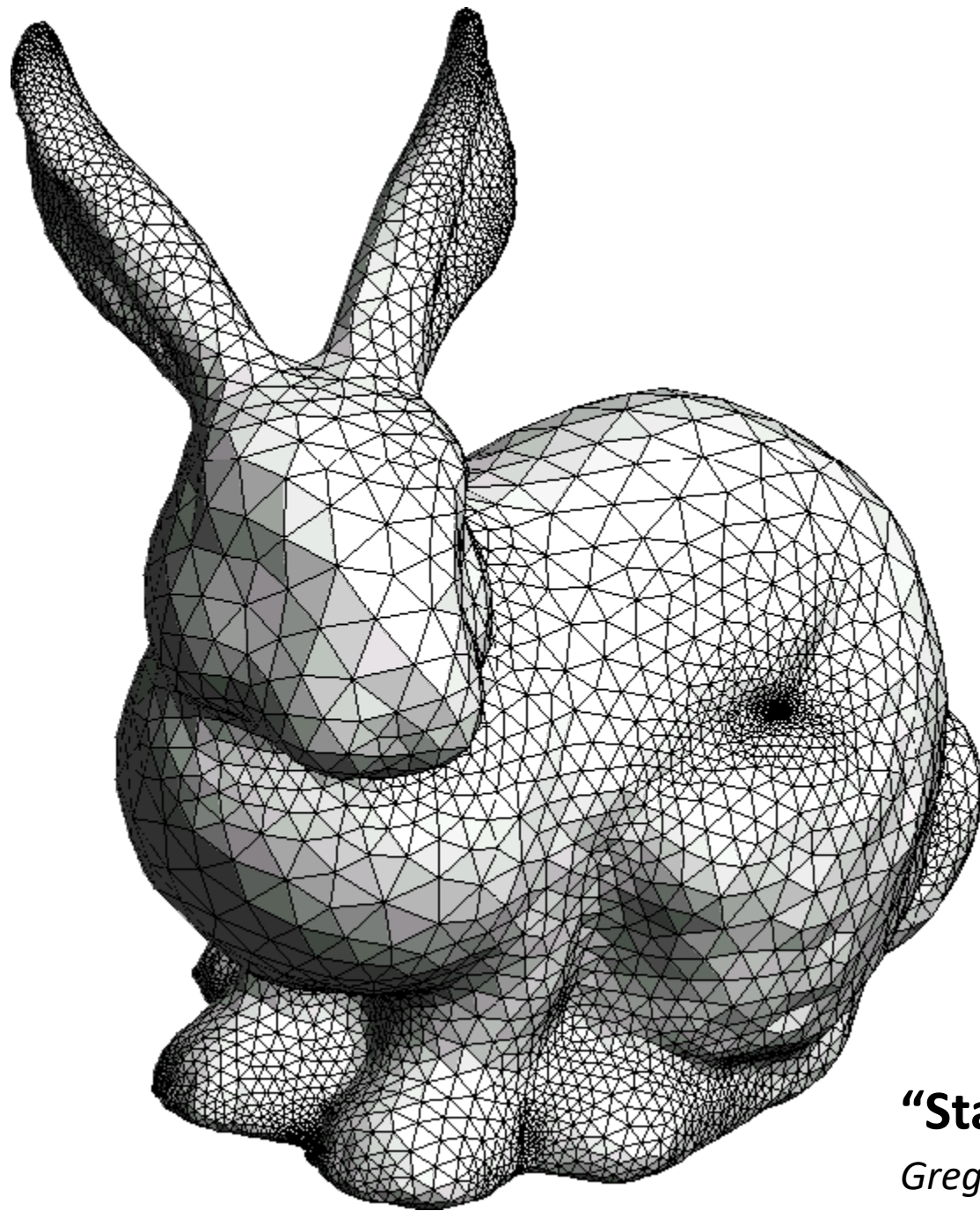
Presented by Cinda Heeren, 2016



“Rush Hour” Solution

Unknown Source

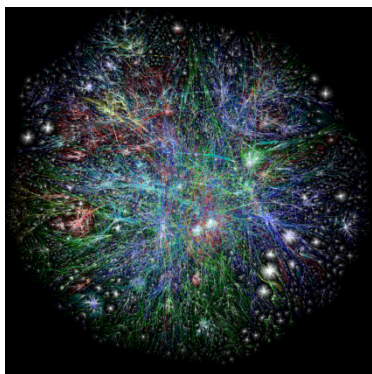
Presented by Cinda Heeren, 2016



“Stanford Bunny”

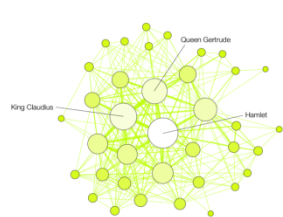
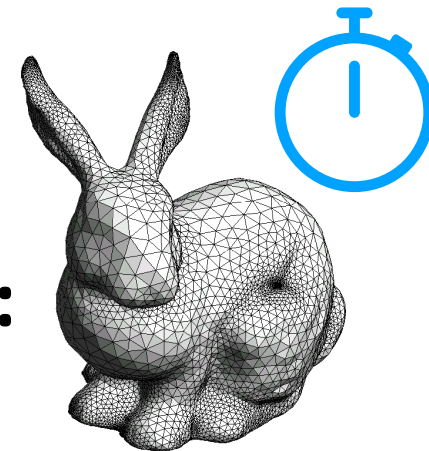
Greg Turk and Mark Levoy (1994)

Graphs



To study all of these structures:

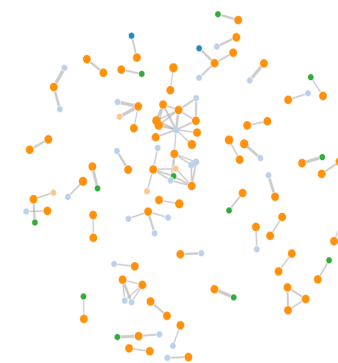
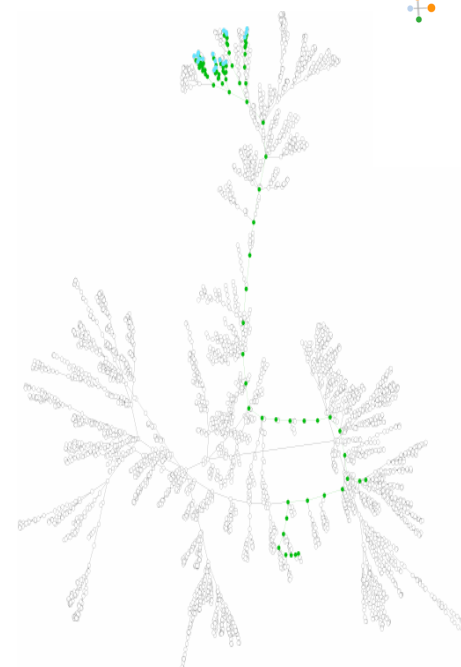
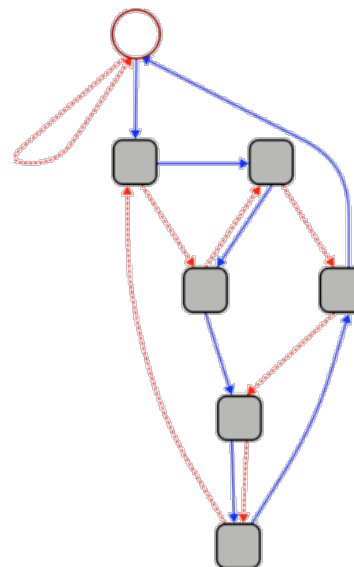
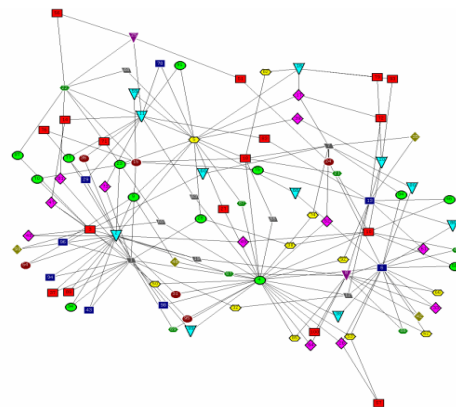
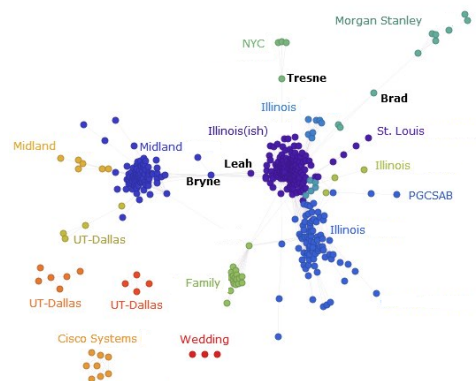
1. A common vocabulary
2. Graph implementations
3. Graph traversals
4. Graph algorithms



HAMLET



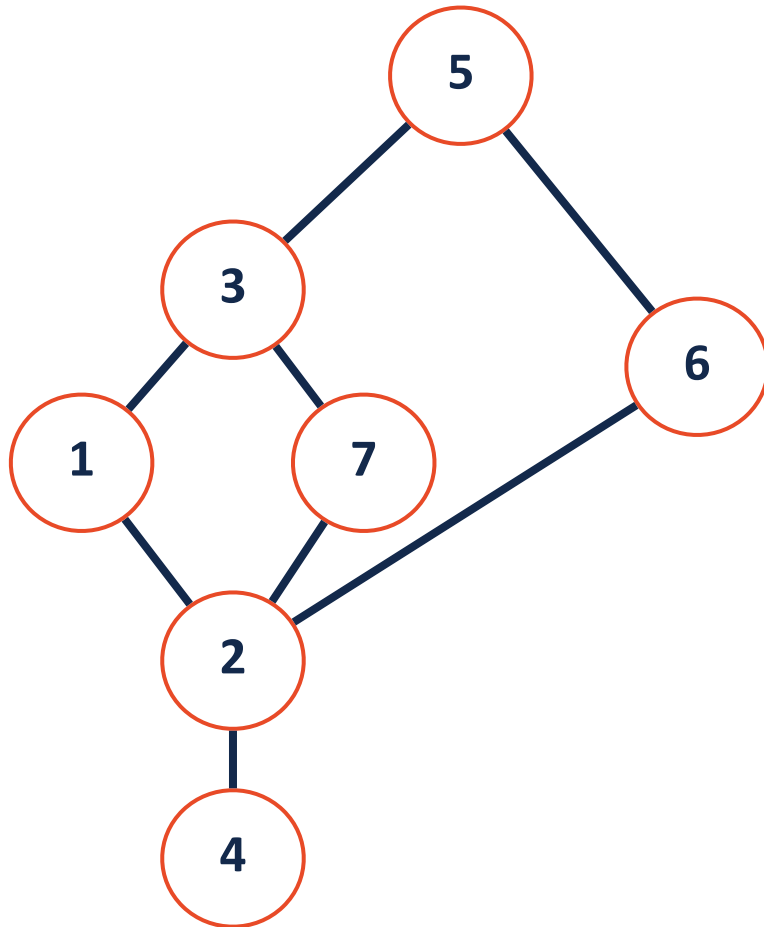
TROILUS AND CRESSIDA



Graph Vocabulary

$$G = (V, E)$$

A **graph** is a data structure containing a set of vertices and a set of edges

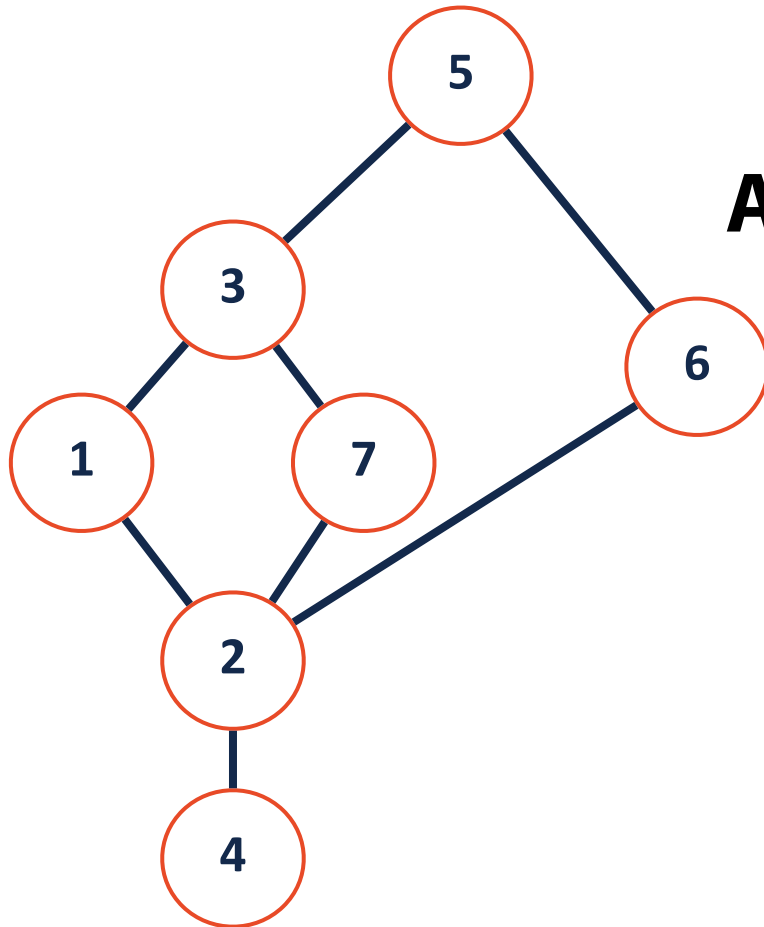


Vertex: Nodes of the graph

Edges: The connections between nodes
Defined by two endpoints

Graph Vocabulary

Degree: # of edges touching a vertex

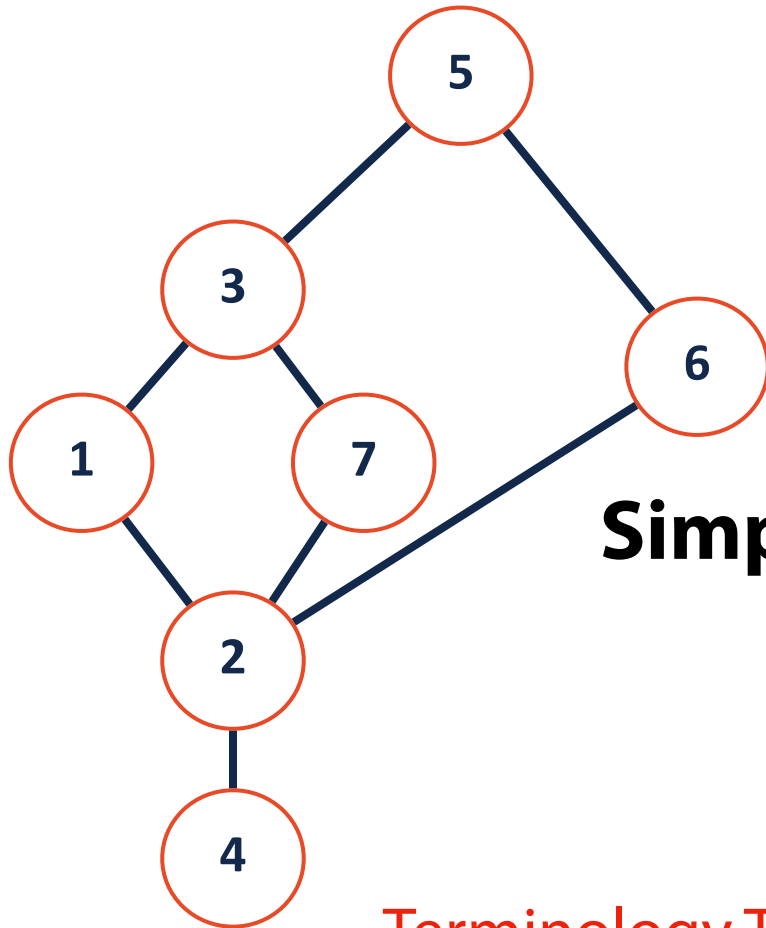


Adjacency: Two vertices are adjacent if they are connected by an edge

Path: A sequence of vertices (or edges) between two nodes

Graph Vocabulary

A graph has **no root** and **may contain cycles**



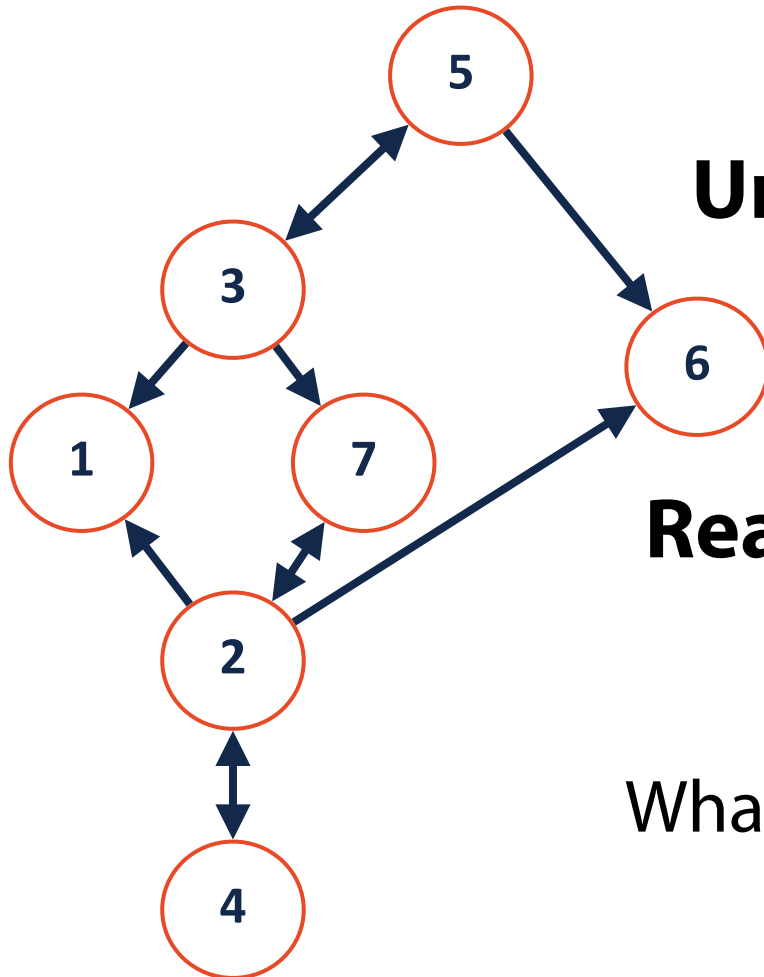
Cycle: A path from a node to itself

Simple Graph: No self-loops or multi-edges

Terminology Trivia: Every tree is a graph but not every graph is a tree

Graph Vocabulary

A graph may be **directed** or **undirected**



Directed: Edges are one way connections

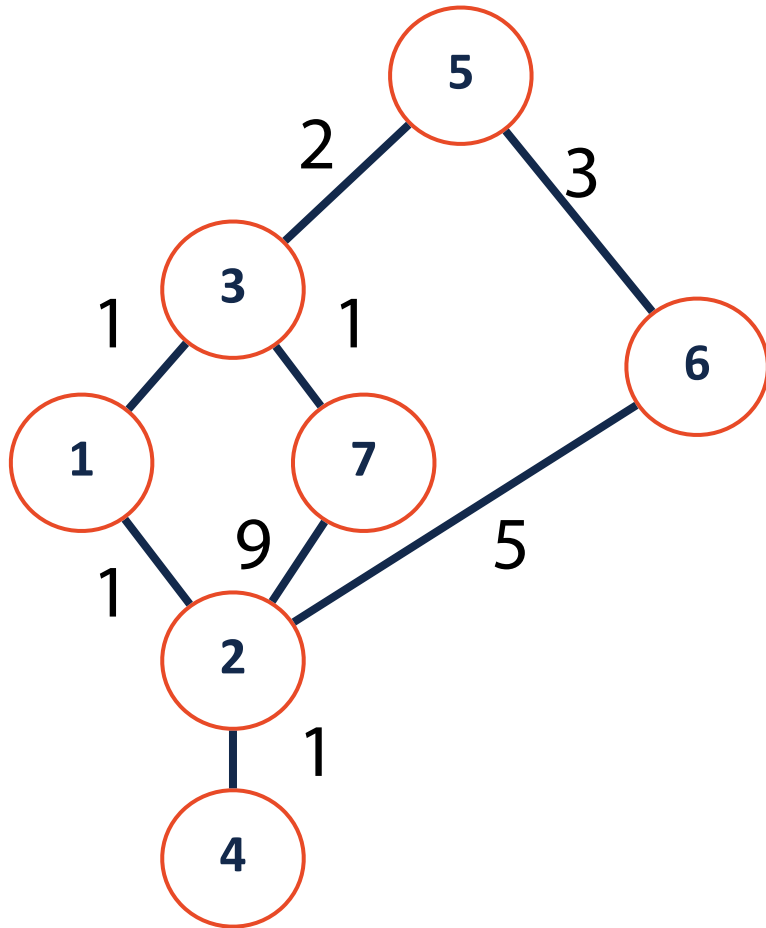
Undirected: Traversable in either direction

Reachability: v_2 is reachable from v_1 if there is a path from v_1 to v_2

What nodes are not reachable from 4?

Graph Vocabulary

A graph may be **weighted** or **unweighted**



Weights: A value associated with an edge

What is the shortest path from 4 to 5?

Graph Vocabulary

$$G = (V, E)$$

$$|V| = n$$

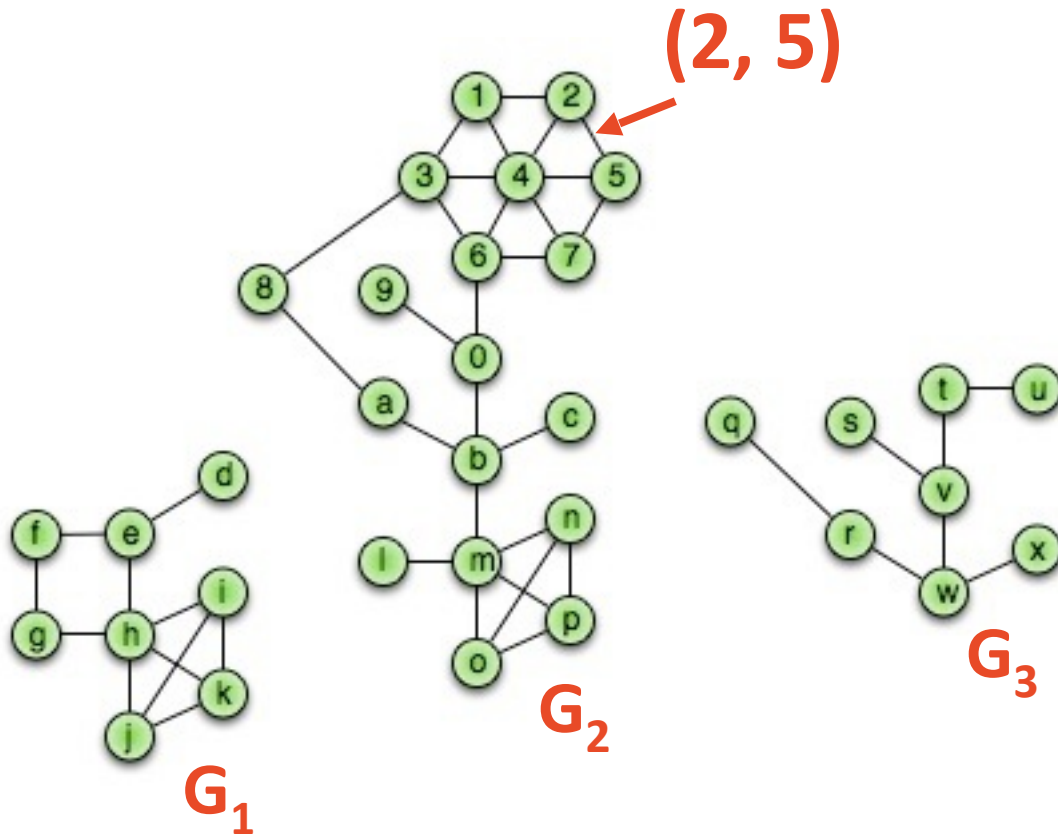
$$|E| = m$$

Subgraph(G):

$$G' = (V', E'):$$

$$V' \subseteq V, E' \subseteq E, \text{ and}$$

$$(u, v) \in E' \rightarrow u \in V', v \in V'$$



Graph Vocabulary

$$G = (V, E)$$

$$|V| = n$$

$$|E| = m$$

Subgraph(G):

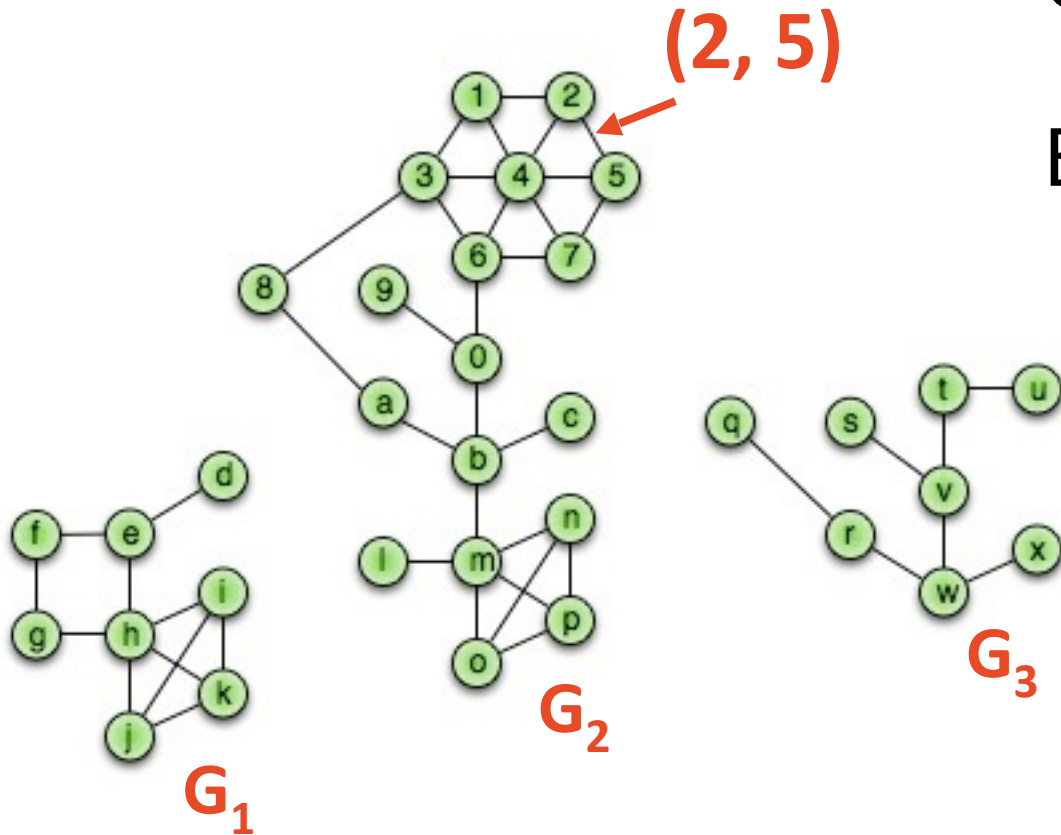
$$G' = (V', E')$$

$V' \subseteq V, E' \subseteq E$, and

$$(u, v) \in E' \rightarrow u \in V', v \in V'$$

Complete Subgraph:

Every pair of vertices are adjacent



Graph Vocabulary

$$G = (V, E)$$

$$|V| = n$$

$$|E| = m$$

Subgraph(G):

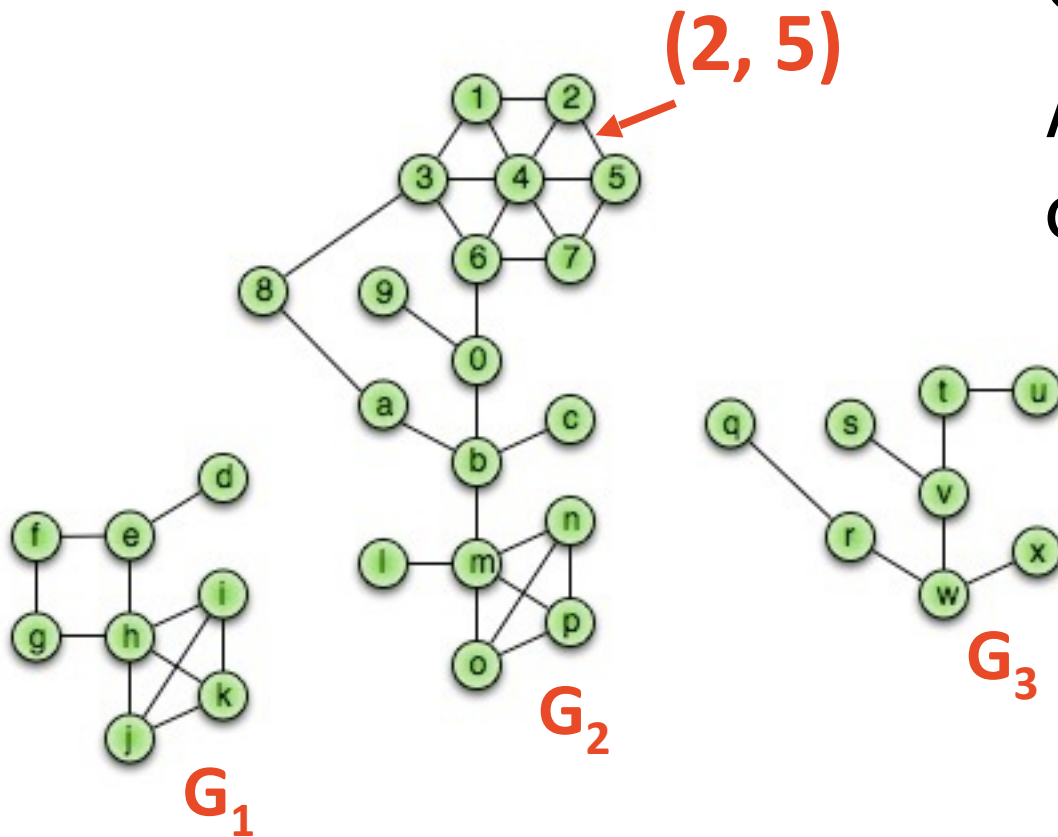
$$G' = (V', E')$$

$V' \subseteq V, E' \subseteq E$, and

$$(u, v) \in E' \rightarrow u \in V', v \in V'$$

Connected Subgraph:

A path exists between every pair of vertices



Graph Vocabulary

$$G = (V, E)$$

$$|V| = n$$

$$|E| = m$$

Subgraph(G):

$$G' = (V', E'):$$

$$V' \subseteq V, E' \subseteq E, \text{ and}$$

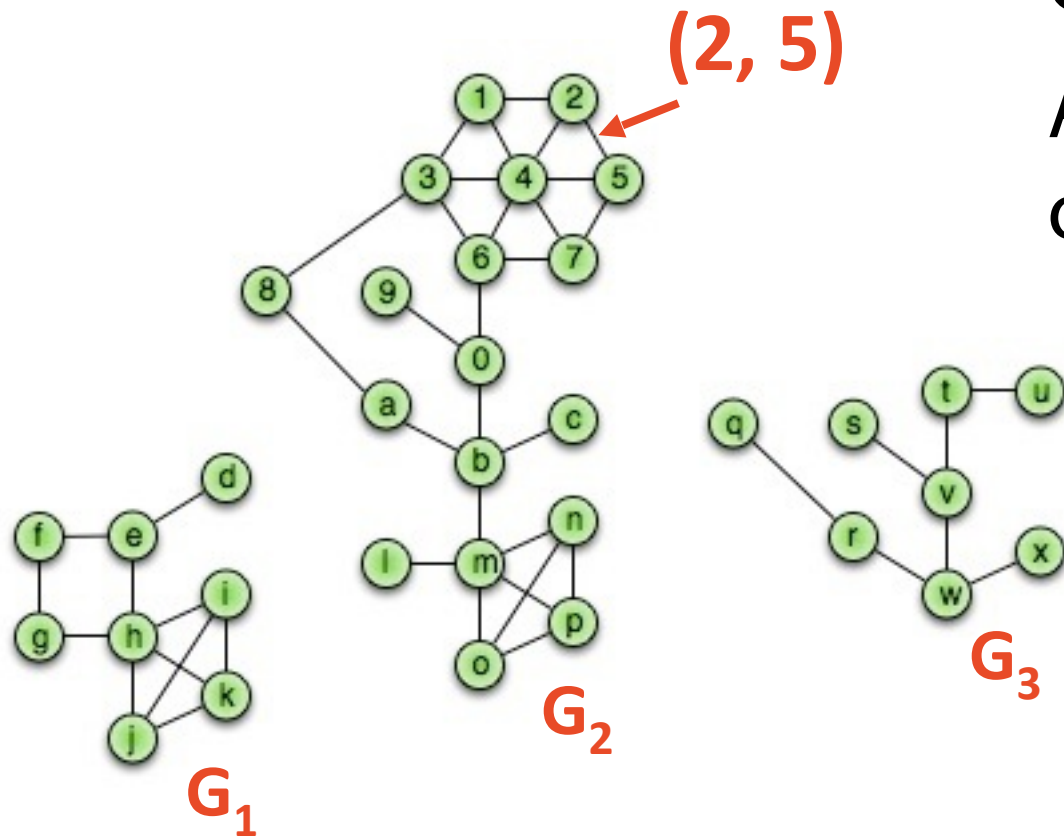
$$(u, v) \in E' \rightarrow u \in V', v \in V'$$

Connected Subgraph:

A path exists between every pair of vertices

Connected Components:

A connected subgraph that is not part of a larger subgraph



Graph Vocabulary

$$G = (V, E)$$

$$|V| = n$$

$$|E| = m$$

Subgraph(G):

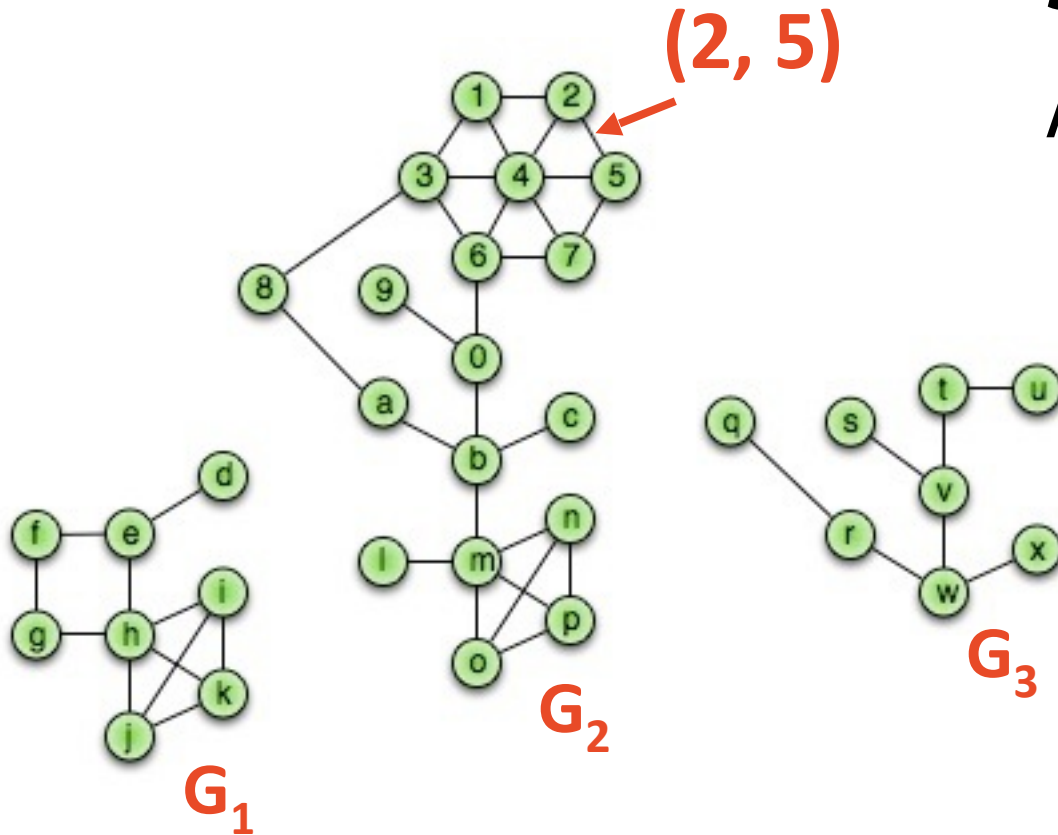
$$G' = (V', E'):$$

$$V' \subseteq V, E' \subseteq E, \text{ and}$$

$$(u, v) \in E' \rightarrow u \in V', v \in V'$$

Spanning Tree:

A connected graph with no cycles



Graph Vocabulary



$$G = (V, E)$$

$$|V| = n$$

$$|E| = m$$

Graph Terminology is very important!

Degree

Weight

Direction

Adjacency

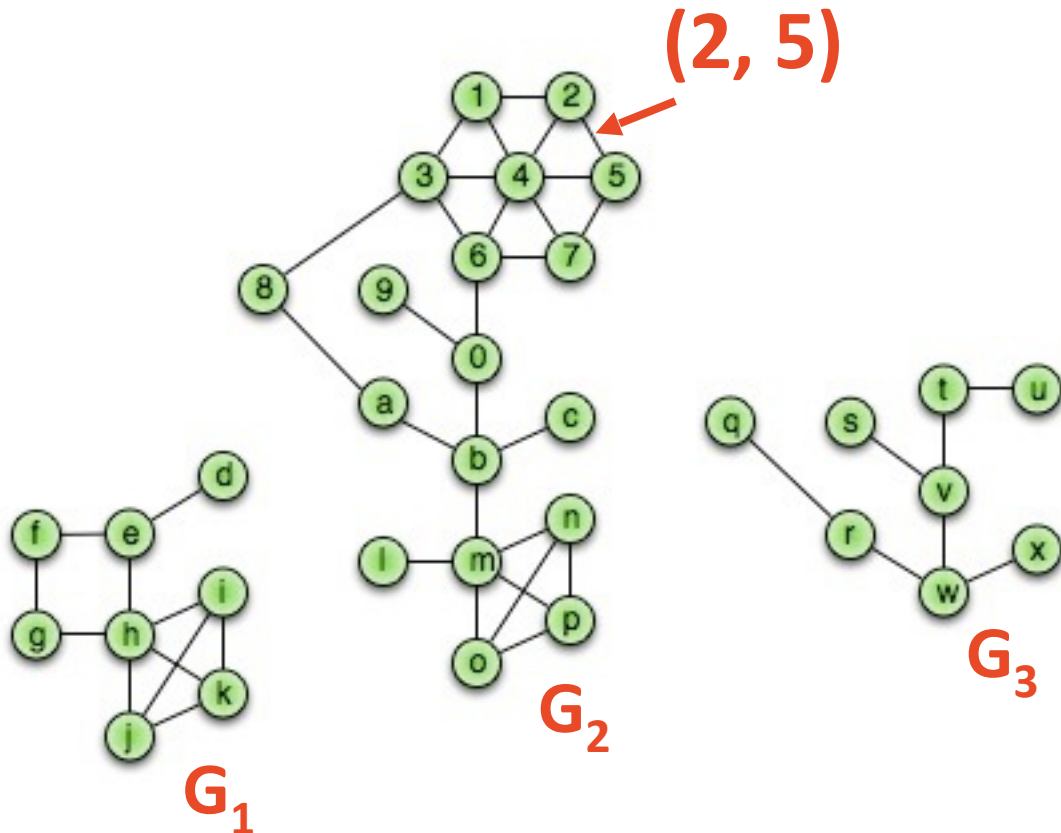
Complete

Connected

Acyclic

Spanning

And more...



Running times are often reported by n , the number of vertices, but often depend on m , the number of edges.

Whats the relationship between n and m ?



Join Code: 225

Minimum Edges:

Unconnected Graph:

Connected (Simple) Graph:

Maximum Edges:

Connected (Simple) Graph:

$$\sum_{v \in V} \deg(v) =$$

Graphs

Given a collection of individual DMs between individuals, you want to build a graph of connections in a social network.

What is a vertex?

What is an edge?

Are the edges directed or undirected?

Are the edges weighted or unweighted?

Graphs

Given a collection of roads between cities in Illinois, you want to build a graph of the transportation infrastructure in the state.

What is a vertex?

What is an edge?

Are the edges directed or undirected?

Are the edges weighted or unweighted?

Graphs

It is important to be able to describe the structure of a graph given input.

Some other common questions:

Does your graph have cycles?

What is the largest / smallest / average degree in your graph?

What is the total number of edges?

...

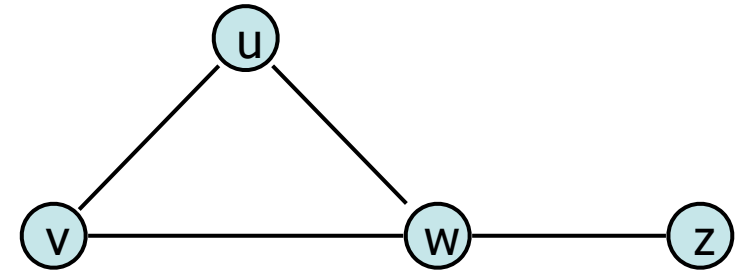
Of course, we also have to understand the graph as a **data structure**

Graph Implementation

What information do we need to store to fully define a graph?

Vertex:

Edge:



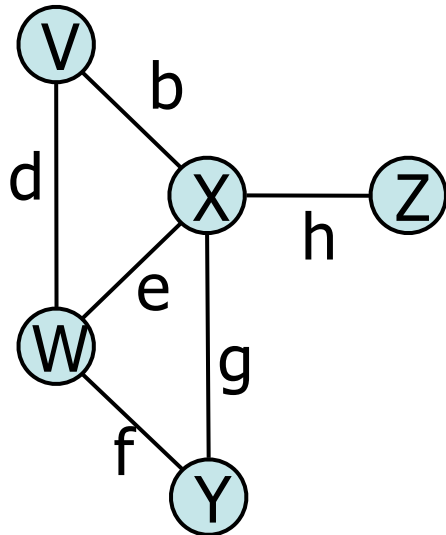
What information do we want to be able to find out quickly?

What operations do we want to prioritize?

Graph ADT

Data:

- Vertices
- Edges
- Some data structure maintaining the structure between vertices and edges.

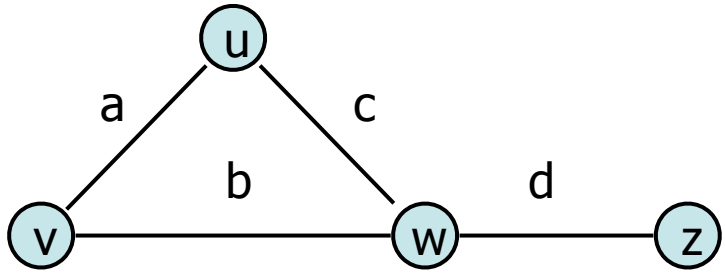


Functions:

- insertVertex(K key);
- insertEdge(Vertex v1, Vertex v2, K key);
- removeVertex(Vertex v);
- removeEdge(Vertex v1, Vertex v2);
- getEdges(Vertex v);
- areAdjacent(Vertex v1, Vertex v2);
- origin(Edge e);
- destination(Edge e);

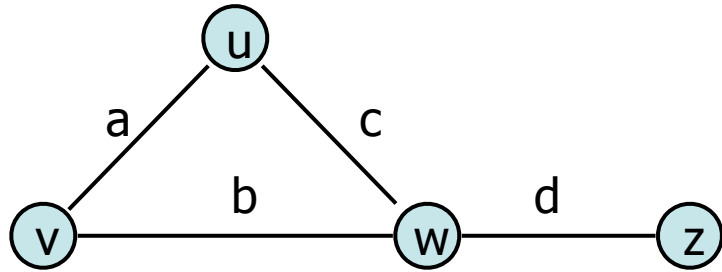


Graph Implementation Idea

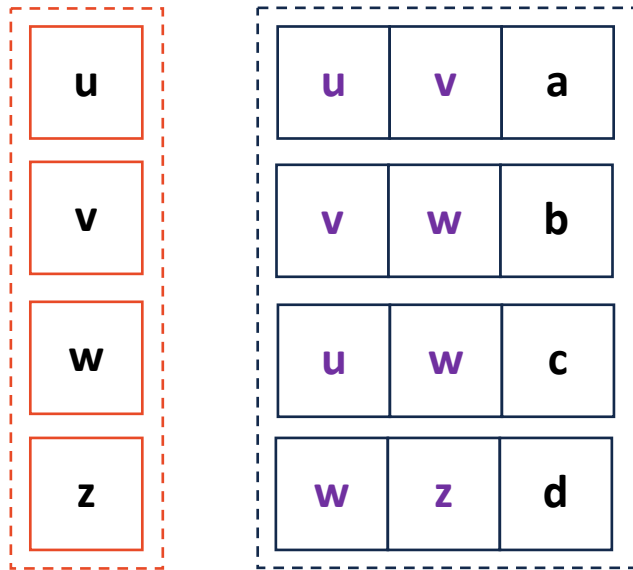


Graph Implementation: Edge List $|V| = n, |E| = m$

The equivalent of an 'unordered' data structure

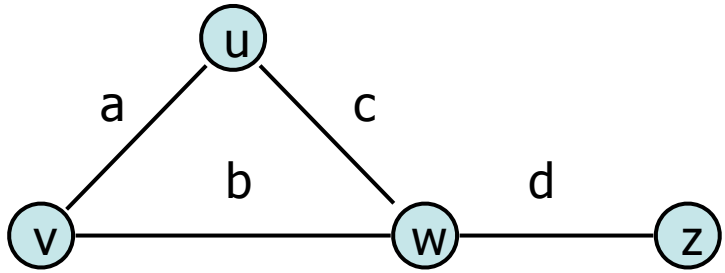


Vertex Storage:

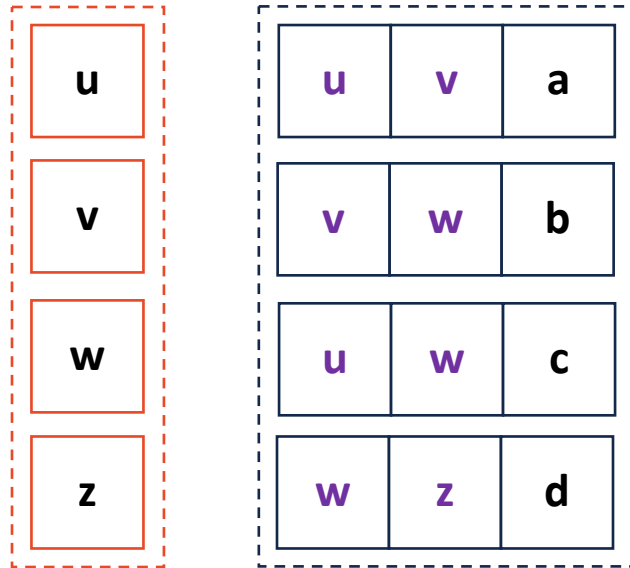


Edge Storage:

Graph Implementation: Edge List $|V| = n, |E| = m$

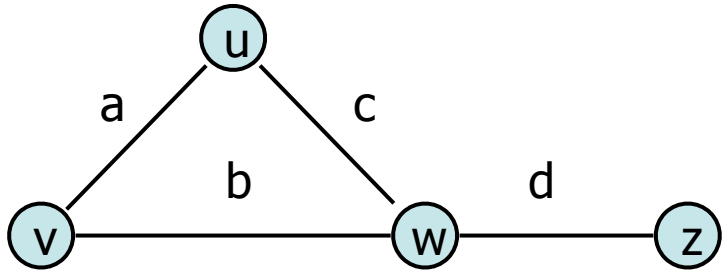


getEdges(Vertex v)

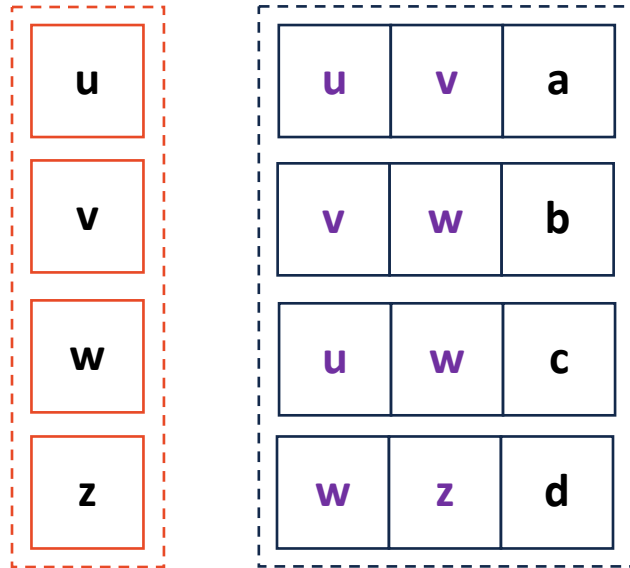


areAdjacent(Vertex v1, Vertex v2)

Graph Implementation: Edge List $|V| = n, |E| = m$

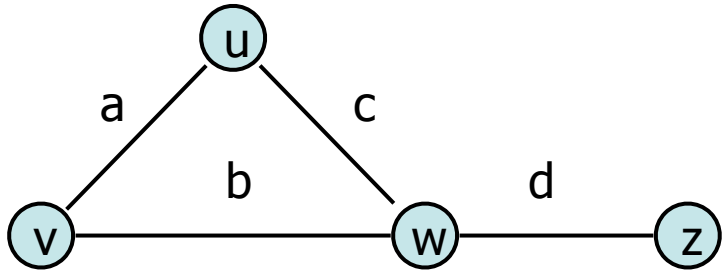


insertVertex(K key)

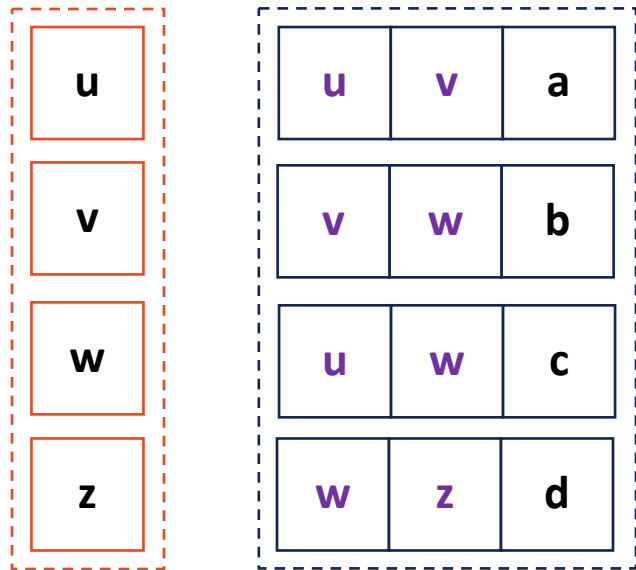


removeVertex(Vertex v)

Graph Implementation: Edge List $|V| = n, |E| = m$



insertEdge(Vertex v1, Vertex v2, K key)



removeEdge(Vertex v1, Vertex v2)

Graph Implementation: Edge List



Pros:

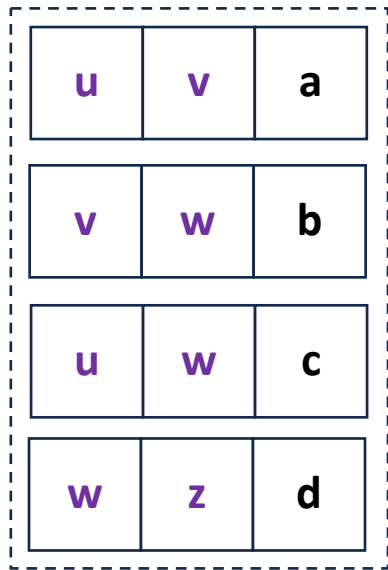
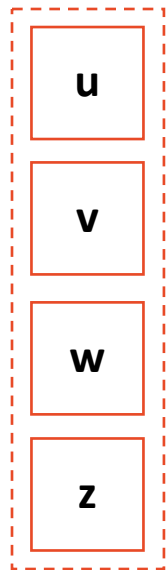
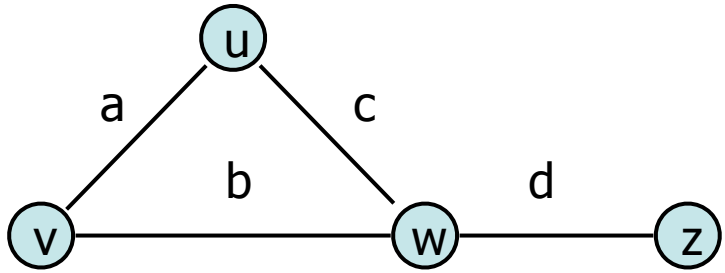
Cons:

Graph Implementation: Brainstorming better

What operations might I want to do very quickly?

What modifications might allow me to do these things faster?

Graph Implementation: Adjacency Matrix



	u	v	w	z
u				
v				
w				
z				