

# Lecture 13: Basic Graph Algorithms

Michael Dinitz

October 12, 2021

601.433/633 Introduction to Algorithms

# Introduction

Next 3-4 weeks: graphs!

- ▶ Super important abstractions, used all over the place in CS
- ▶ Most of my research is in graph algorithms (particularly when graph represents computer/communication network)
- ▶ Great course on Graph Theory in AMS

Today: review of basic graph algorithms from Data Structures, one or two new

- ▶ Going to move pretty quickly, since much review: see CLRS for details!

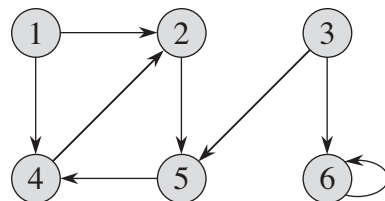
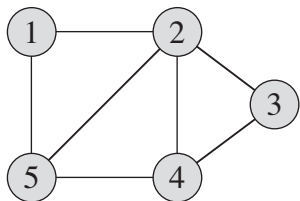
# Basic Definitions

## Definition

A graph  $\mathbf{G} = (\mathbf{V}, \mathbf{E})$  is a pair where  $\mathbf{V}$  is a set and  $\mathbf{E} \subseteq \binom{\mathbf{V}}{2}$  (unordered pairs) or  $\mathbf{E} \subseteq \mathbf{V} \times \mathbf{V}$  (ordered pairs).

## Notation:

- ▶ Elements of  $\mathbf{V}$  are called *vertices* or *nodes*
- ▶ Elements of  $\mathbf{E}$  are called *edges* or *arcs*.
- ▶ If  $\mathbf{E} \subseteq \binom{\mathbf{V}}{2}$  then graph is *undirected*, if  $\mathbf{E} \subseteq \mathbf{V} \times \mathbf{V}$  graph is *directed*
- ▶  $|\mathbf{V}| = \mathbf{n}$  and  $|\mathbf{E}| = \mathbf{m}$  (usually)
- ▶ So “size of input” =  $\mathbf{n} + \mathbf{m}$



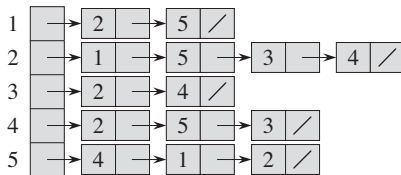
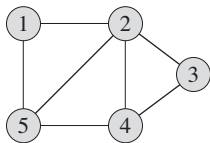
# Representations

## Adjacency List:

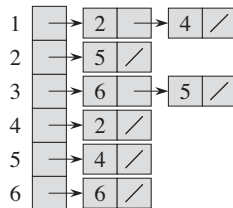
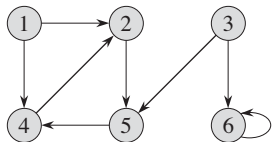
- ▶ Array **A** of length **n**
- ▶ **A[v]** is linked list of vertices *adjacent* to **v** (edge from **u** to **v**)

## Adjacency Matrix:

- ▶  $\mathbf{A} \in \{0, 1\}^{n \times n}$
- ▶  $A_{ij} = \begin{cases} 1 & \text{if } (i, j) \in E \\ 0 & \text{otherwise} \end{cases}$



	1	2	3	4	5
1	0	1	0	0	1
2	1	0	1	1	1
3	0	1	0	1	0
4	0	1	1	0	1
5	1	1	0	1	0



	1	2	3	4	5	6
1	0	1	0	1	0	0
2	0	0	0	0	1	0
3	0	0	0	0	1	1
4	0	1	0	0	0	0
5	0	0	0	1	0	0
6	0	0	0	0	0	1

# Representations (cont'd)

## Adjacency List:

- ▶ Pros:
  - ▶  $O(n + m)$  space
  - ▶ Can iterate through edges adjacent to  $v$  very efficiently
- ▶ Cons:
  - ▶ Hard to check if an edge exists:  $O(d(u))$  or  $O(d(v))$  (where  $d(v)$  is the degree of  $v$ : # edges with  $v$  as endpoint)

## Adjacency Matrix:

- ▶ Pros:
  - ▶ Can check if  $e = (u, v)$  an edge in  $O(1)$  time
- ▶ Cons:
  - ▶ Takes  $\Theta(n^2)$  space: if  $m$  small, lots wasted!
  - ▶ Iterating through edges incident on  $v$  takes time  $\Theta(n)$ , even if  $d(v)$  small.

# Representations (cont'd)

## Adjacency List:

- ▶ Pros:
  - ▶  $O(n + m)$  space
  - ▶ Can iterate through edges adjacent to  $\mathbf{v}$  very efficiently
- ▶ Cons:
  - ▶ Hard to check if an edge exists:  
 $O(d(\mathbf{u}))$  or  $O(d(\mathbf{v}))$  (where  $d(\mathbf{v})$  is the degree of  $\mathbf{v}$ : # edges with  $\mathbf{v}$  as endpoint)

This class: adjacency list unless otherwise specified.

## Adjacency Matrix:

- ▶ Pros:
  - ▶ Can check if  $\mathbf{e} = (\mathbf{u}, \mathbf{v})$  an edge in  $O(1)$  time
- ▶ Cons:
  - ▶ Takes  $\Theta(n^2)$  space: if  $\mathbf{m}$  small, lots wasted!
  - ▶ Iterating through edges incident on  $\mathbf{v}$  takes time  $\Theta(n)$ , even if  $d(\mathbf{v})$  small.

# Representations (cont'd)

## Adjacency List:

- ▶ Pros:
  - ▶  $O(n + m)$  space
  - ▶ Can iterate through edges adjacent to  $v$  very efficiently
- ▶ Cons:
  - ▶ Hard to check if an edge exists:  $O(d(u))$  or  $O(d(v))$  (where  $d(v)$  is the degree of  $v$ : # edges with  $v$  as endpoint)

This class: adjacency list unless otherwise specified.

Any way to improve these?

## Adjacency Matrix:

- ▶ Pros:
  - ▶ Can check if  $e = (u, v)$  an edge in  $O(1)$  time
- ▶ Cons:
  - ▶ Takes  $\Theta(n^2)$  space: if  $m$  small, lots wasted!
  - ▶ Iterating through edges incident on  $v$  takes time  $\Theta(n)$ , even if  $d(v)$  small.

# Representations (cont'd)

## Adjacency List:

- ▶ Pros:
  - ▶  $O(n + m)$  space
  - ▶ Can iterate through edges adjacent to  $v$  very efficiently
- ▶ Cons:
  - ▶ Hard to check if an edge exists:  $O(d(u))$  or  $O(d(v))$  (where  $d(v)$  is the degree of  $v$ : # edges with  $v$  as endpoint)

This class: adjacency list unless otherwise specified.

Any way to improve these?

- ▶ Replace adjacency *list* with adjacency *structure*: Red-black tree, hash table, etc.
- ▶ Not traditional, doesn't gain us much, and more complicated. But better!

## Adjacency Matrix:

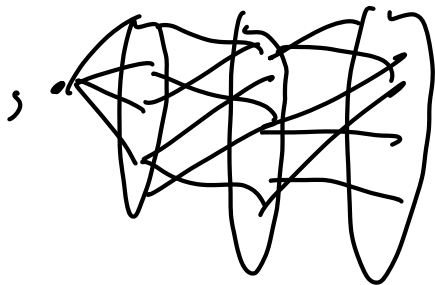
- ▶ Pros:
  - ▶ Can check if  $e = (u, v)$  an edge in  $O(1)$  time
- ▶ Cons:
  - ▶ Takes  $\Theta(n^2)$  space: if  $m$  small, lots wasted!
  - ▶ Iterating through edges incident on  $v$  takes time  $\Theta(n)$ , even if  $d(v)$  small.



# Breadth-First Search (BFS)

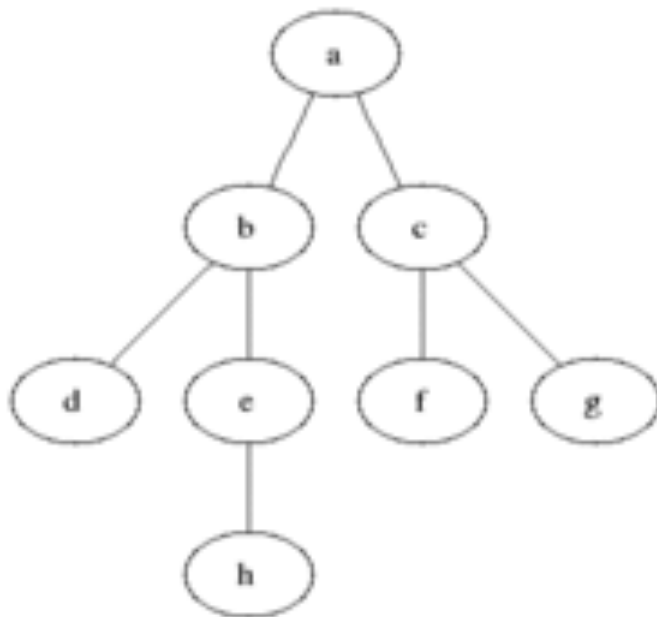
# BFS Definition

Idea: explore graph in *levels* or *layers* from source  $s$



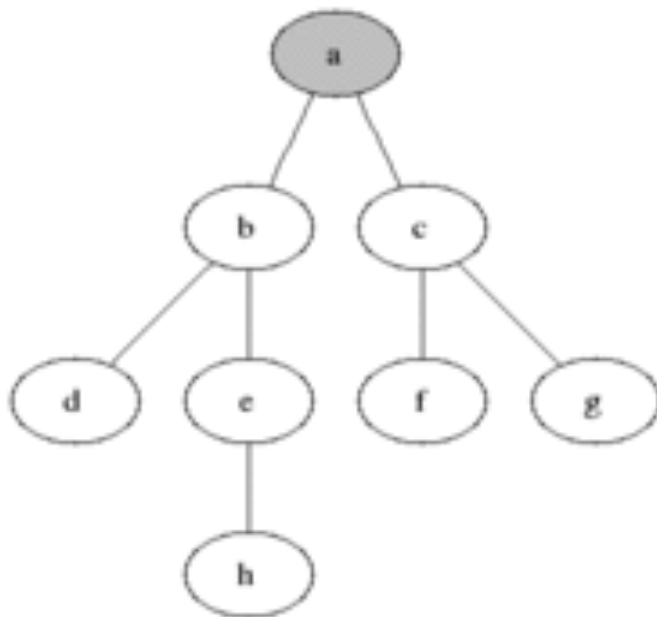
# BFS Definition

Idea: explore graph in *levels* or *layers* from source **s**



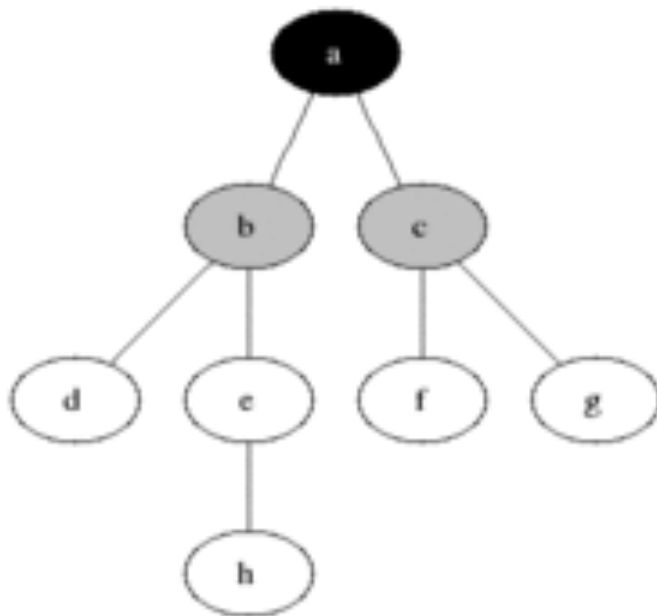
# BFS Definition

Idea: explore graph in *levels* or *layers* from source  $s$



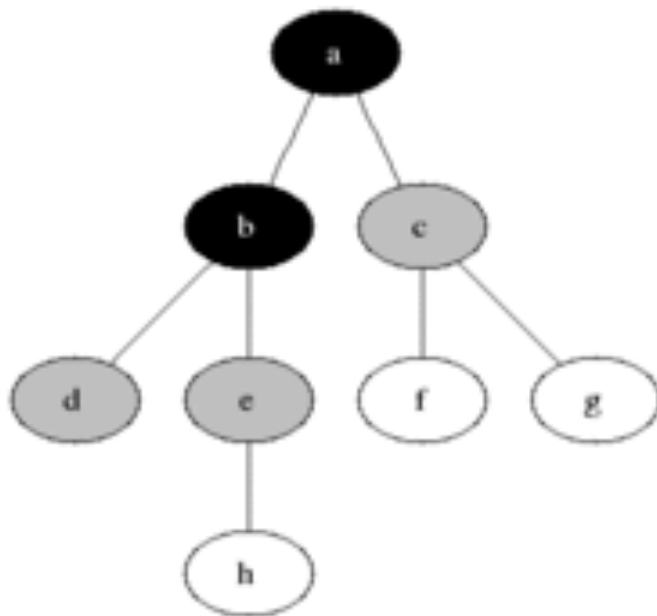
# BFS Definition

Idea: explore graph in *levels* or *layers* from source  $s$



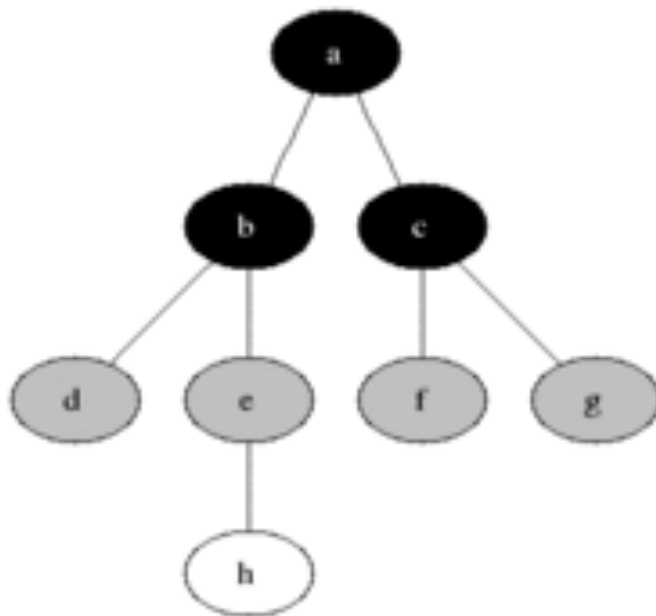
# BFS Definition

Idea: explore graph in *levels* or *layers* from source  $s$



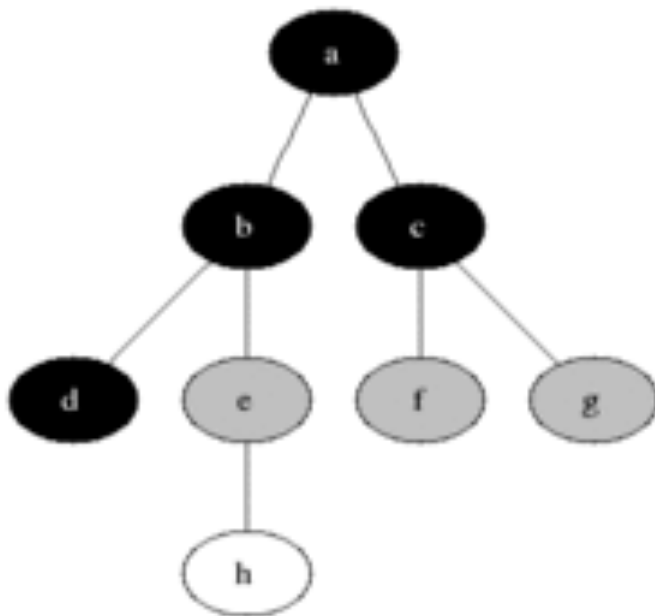
# BFS Definition

Idea: explore graph in *levels* or *layers* from source  $s$



# BFS Definition

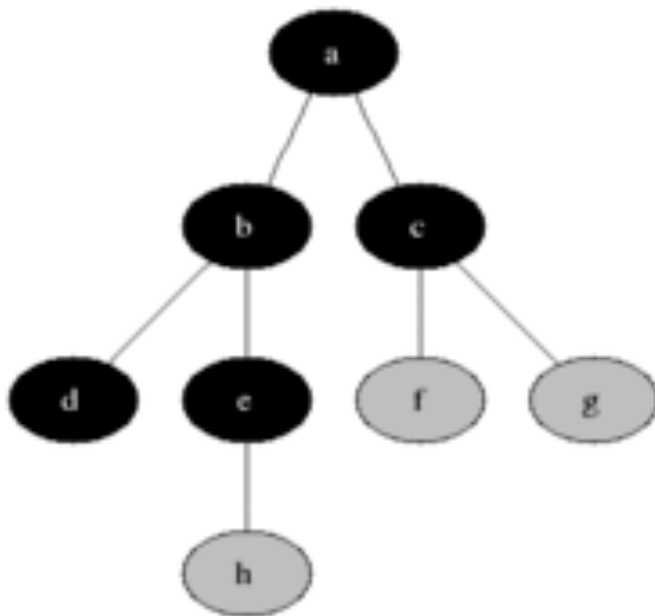
Idea: explore graph in *levels* or *layers* from source  $s$





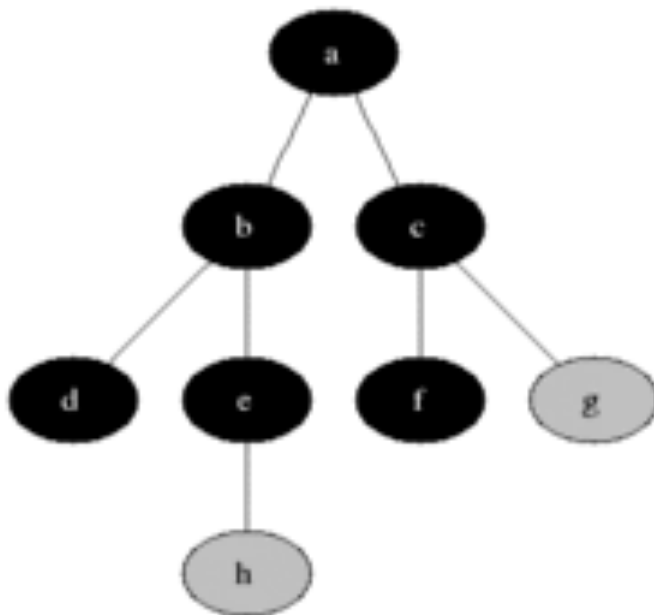
# BFS Definition

Idea: explore graph in *levels* or *layers* from source  $s$



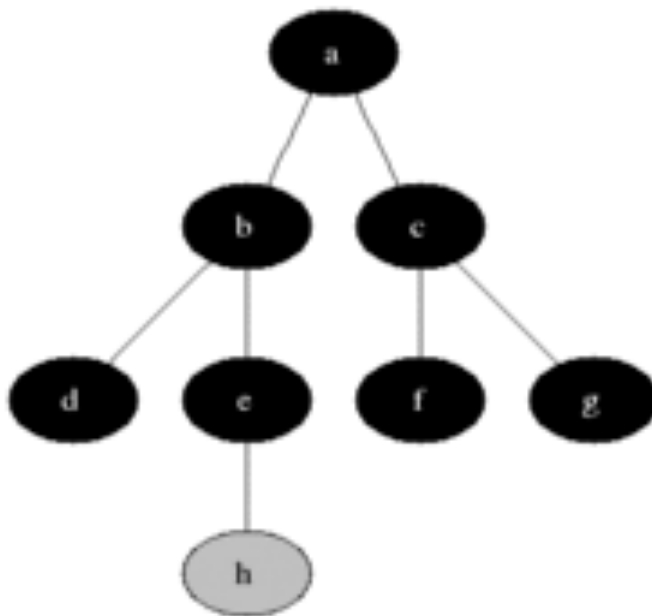
# BFS Definition

Idea: explore graph in *levels* or *layers* from source  $s$



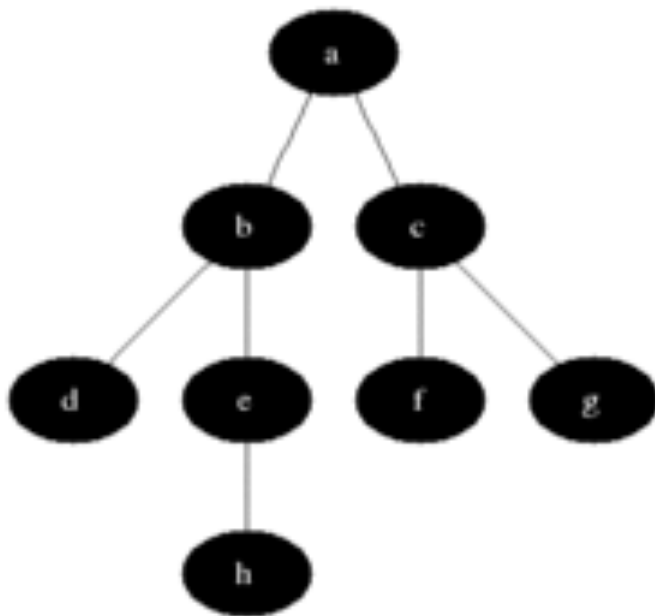
# BFS Definition

Idea: explore graph in *levels* or *layers* from source  $s$



# BFS Definition

Idea: explore graph in *levels* or *layers* from source  $s$



# BFS Pseudocode

Idea: explore with a queue (LIFO)

```
BFS(G = (V, E), s) {  
  Set mark(s) = True;  
  Set mark(v) = False for all v  $\in$  V  $\setminus$  {s};  
  Enqueue(s);  
  while(queue not empty) {  
    v = Dequeue();  
    forall neighbors u of v {  
      if(mark(u) == False) {  
        mark(u) = True;  
        Enqueue(u);  
      }  
    }  
  }  
}
```

# BFS Pseudocode

Idea: explore with a queue (LIFO)

```
BFS(G = (V, E), s) {  
  Set mark(s) = True;  
  Set mark(v) = False for all v  $\in$  V  $\setminus$  {s};  
  Enqueue(s);  
  while(queue not empty) {  
    v = Dequeue();  
    forall neighbors u of v {  
      if(mark(u) == False) {  
        mark(u) = True;  
        Enqueue(u);  
      }  
    }  
  }  
}
```

**Running Time:**

# BFS Pseudocode

Idea: explore with a queue (LIFO)

```
BFS(G = (V, E), s) {  
  Set mark(s) = True;  
  Set mark(v) = False for all v  $\in$  V  $\setminus$  {s};  
  Enqueue(s);  
  while(queue not empty) {  
    v = Dequeue();  
    forall neighbors u of v {  
      if(mark(u) == False) {  
        mark(u) = True;  
        Enqueue(u);  
      }  
    }  
  }  
}
```

**Running Time:  $O(n + m)$**

# BFS Pseudocode

Idea: explore with a queue (LIFO)

```
BFS( $\mathbf{G} = (\mathbf{V}, \mathbf{E}), s$ ) {  
  Set  $\mathbf{mark}(s) = \mathbf{True}$ ;  
  Set  $\mathbf{mark}(v) = \mathbf{False}$  for all  $v \in \mathbf{V} \setminus \{s\}$ ;  
  Enqueue( $s$ );  
  while(queue not empty) {  
     $v = \mathbf{Dequeue}()$ ;  
    forall neighbors  $u$  of  $v$  {  
      if( $\mathbf{mark}(u) == \mathbf{False}$ ) {  
         $\mathbf{mark}(u) = \mathbf{True}$ ;  
        Enqueue( $u$ );  
      }  
    }  
  }  
}
```

*Handwritten red annotations:*  
A red curly brace on the right side of the pseudocode groups the initialization steps (Set mark(s) = True; Set mark(v) = False for all v in V \ {s}; Enqueue(s);) and is labeled  $O(n)$ .  
A larger red curly brace on the right side groups the main while loop (while(queue not empty) { ... }) and is labeled  $O(m)$ .

**Running Time:  $O(n + m)$**

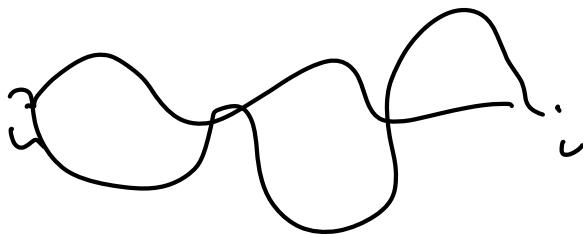
- ▶  $O(n)$  for initialization
- ▶  $O(m)$  for main while loop
  - ▶ Examine every edge twice:  
when each endpoint dequeued
  - ▶ Or (equivalent): Adjacency list  
scanned only when vertex  
dequeued



## Correctness / Shortest Paths

**Definition:** Distance  $d(u, v)$  from  $u$  to  $v$  is min # edges in any path from  $u$  to  $v$

**Theorem (informal):** BFS( $s$ ) gives shortest paths from  $s$  to all other nodes



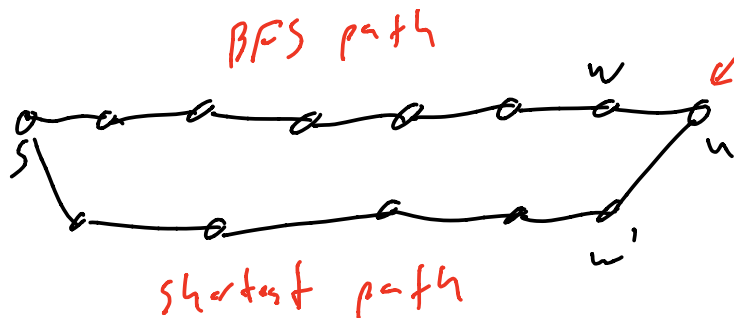
# Correctness / Shortest Paths

**Definition:** Distance  $d(u, v)$  from  $u$  to  $v$  is min # edges in any path from  $u$  to  $v$

**Theorem (informal):** BFS( $s$ ) gives shortest paths from  $s$  to all other nodes

## Proof Sketch:

Assume false for contradiction, let  $u$  be closest node to  $s$  where BFS( $s$ ) doesn't give shortest path



$$d(s, w') < d(s, w)$$

- $\Rightarrow w'$  dequeued before  $w$  (since  $w'$  has correct distance by def of  $u$ )
- $\Rightarrow u$  will be enqueued from  $w'$ , not  $w$ . Contradiction.

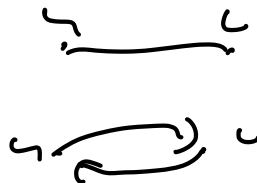
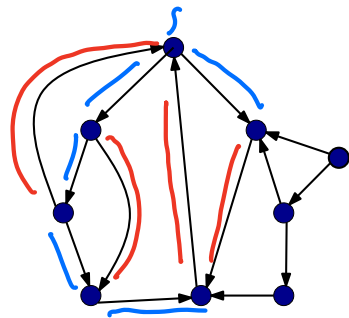
# Depth-First Search (DFS)

# DFS: Definition

Intuition: Instead of exploring wide (breadth), explore far (deep): just keep walking until see a node we've already seen, then backtrack!

```
Init: for each  $v \in V$ ,  $\text{mark}(v) = \text{False}$ ;
```

```
DFS( $v$ ) {  
   $\text{mark}(v) = \text{True}$ ;  
  for each edge  $(v, u) \in A[v]$  {  
    if  $\text{mark}(u) == \text{False}$  then DFS( $u$ );  
  }  
}
```

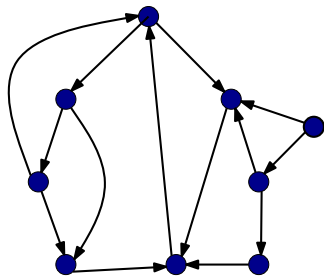


# DFS: Definition

Intuition: Instead of exploring wide (breadth), explore far (deep): just keep walking until see a node we've already seen, then backtrack!

```
Init: for each  $v \in V$ ,  $\text{mark}(v) = \text{False}$ ;
```

```
DFS( $v$ ) {  
   $\text{mark}(v) = \text{True}$ ;  
  for each edge  $(v, u) \in A[v]$  {  
    if  $\text{mark}(u) == \text{False}$  then DFS( $u$ );  
  }  
}
```



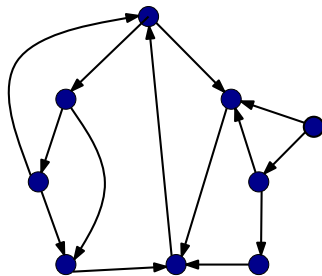
Running time:

## DFS: Definition

Intuition: Instead of exploring wide (breadth), explore far (deep): just keep walking until see a node we've already seen, then backtrack!

```
Init: for each  $v \in V$ ,  $\text{mark}(v) = \text{False}$ ;
```

```
DFS( $v$ ) {  
   $\text{mark}(v) = \text{True}$ ;  
  for each edge  $(v, u) \in A[v]$  {  
    if  $\text{mark}(u) == \text{False}$  then DFS( $u$ );  
  }  
}
```



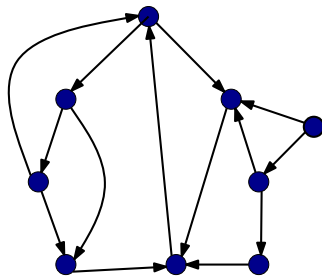
Running time:  $O(m + n)$

## DFS: Definition

Intuition: Instead of exploring wide (breadth), explore far (deep): just keep walking until see a node we've already seen, then backtrack!

```
Init: for each  $v \in V$ ,  $\text{mark}(v) = \text{False}$ ;
```

```
DFS( $v$ ) {  
     $\text{mark}(v) = \text{True}$ ;  
    for each edge  $(v, u) \in A[v]$  {  
        if  $\text{mark}(u) == \text{False}$  then DFS( $u$ );  
    }  
}
```



Running time:  $O(m + n)$

- ▶  $O(n)$  initialization
- ▶ Every edge considered at most twice

## DFS: Correctness

**Definition:**  $u$  is *reachable* from  $v$  if there is a path  $v = v_0, v_1, \dots, v_k = u$  such that  $(v_i, v_{i+1}) \in E$  for all  $i \in \{0, 1, \dots, k-1\}$ .

### Theorem

*When  $DFS(v)$  terminates, it has visited (marked) all nodes that are reachable from  $v$ .*

### Proof.

Suppose  $u$  reachable from  $v$  but not marked when  $DFS(v)$  terminates.



## DFS: Correctness

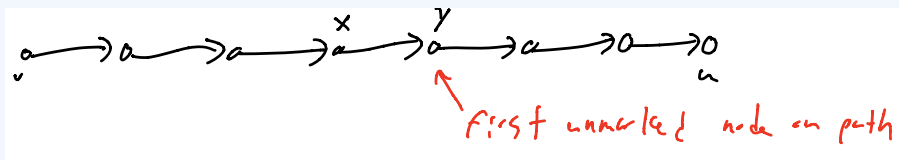
**Definition:**  $u$  is *reachable* from  $v$  if there is a path  $v = v_0, v_1, \dots, v_k = u$  such that  $(v_i, v_{i+1}) \in E$  for all  $i \in \{0, 1, \dots, k-1\}$ .

### Theorem

When  $DFS(v)$  terminates, it has visited (marked) all nodes that are reachable from  $v$ .

### Proof.

Suppose  $u$  reachable from  $v$  but not marked when  $DFS(v)$  terminates.



## DFS: Correctness

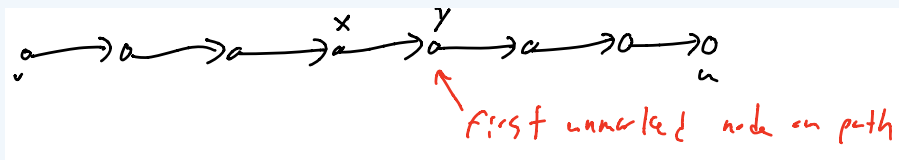
**Definition:**  $u$  is *reachable* from  $v$  if there is a path  $v = v_0, v_1, \dots, v_k = u$  such that  $(v_i, v_{i+1}) \in E$  for all  $i \in \{0, 1, \dots, k-1\}$ .

### Theorem

When  $DFS(v)$  terminates, it has visited (marked) all nodes that are reachable from  $v$ .

### Proof.

Suppose  $u$  reachable from  $v$  but not marked when  $DFS(v)$  terminates.



$x$  is marked so  $DFS(x)$  must have been called

## DFS: Correctness

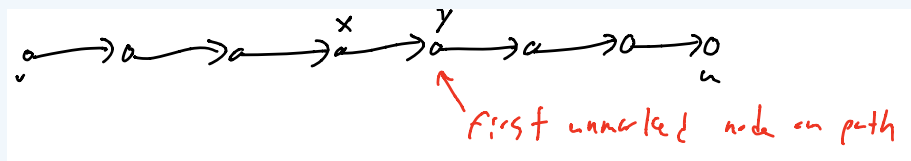
**Definition:**  $u$  is *reachable* from  $v$  if there is a path  $v = v_0, v_1, \dots, v_k = u$  such that  $(v_i, v_{i+1}) \in E$  for all  $i \in \{0, 1, \dots, k-1\}$ .

### Theorem

When  $DFS(v)$  terminates, it has visited (marked) all nodes that are reachable from  $v$ .

### Proof.

Suppose  $u$  reachable from  $v$  but not marked when  $DFS(v)$  terminates.



$x$  is marked so  $DFS(x)$  must have been called

$\implies y$  was either marked or  $DFS(y)$  called and it became marked.

## DFS: Correctness

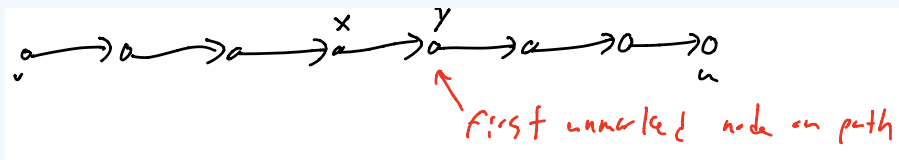
**Definition:**  $u$  is *reachable* from  $v$  if there is a path  $v = v_0, v_1, \dots, v_k = u$  such that  $(v_i, v_{i+1}) \in E$  for all  $i \in \{0, 1, \dots, k-1\}$ .

### Theorem

When  $DFS(v)$  terminates, it has visited (marked) all nodes that are reachable from  $v$ .

### Proof.

Suppose  $u$  reachable from  $v$  but not marked when  $DFS(v)$  terminates.



$x$  is marked so  $DFS(x)$  must have been called

$\implies y$  was either marked or  $DFS(y)$  called and it became marked.

Contradiction. □

## Graph variant

After  $\text{DFS}(\mathbf{v})$ , node marked if and only if reachable from  $\mathbf{v}$ .

Might want to continue until all nodes marked.

```
DFS(G) {  
  for all  $\mathbf{v} \in \mathbf{V}$ , set  $\text{mark}(\mathbf{v}) = \text{False}$ ;  
  while there exists an unmarked node  $\mathbf{v}$  {  
    DFS( $\mathbf{v}$ );  
  }  
}
```

# Timestamps

Explicitly keep track of “start” and “finishing” times

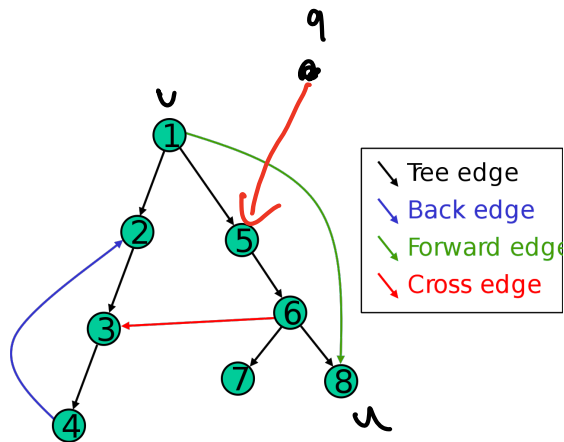
- ▶ Replaces **mark**

```
DFS(G) {  
    t = 0;  
    for all v ∈ V {  
        start(v) = 0;  
        finish(v) = 0;  
    }  
    while ∃v ∈ V with start(v) = 0 {  
        DFS(v);  
    }  
}
```

```
DFS(v) {  
    t = t + 1;  
    start(v) = t;  
    for each edge (v, u) ∈ A[v] {  
        if start(u) == 0 then DFS(u);  
    }  
    t = t + 1;  
    finish(v) = t;  
}
```

# Edge Types

DFS naturally gives a spanning forest: edge  $(v, u)$  if  $\text{DFS}(v)$  calls  $\text{DFS}(u)$



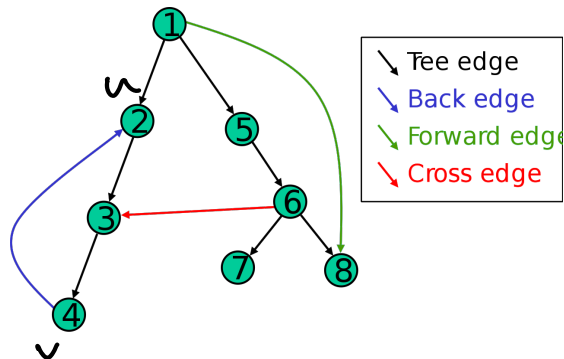
**Forward Edges:**  $(v, u)$  such that  $u$  descendent of  $v$   
(includes tree edges)

**Back Edges:**  $(v, u)$  such that  $u$  an ancestor of  $v$

**Cross Edges:**  $(v, u)$  such that  $u$  neither a descendent nor an ancestor of  $v$

# Edge Types

DFS naturally gives a spanning forest: edge  $(v, u)$  if  $\text{DFS}(v)$  calls  $\text{DFS}(u)$



**Forward Edges:**  $(v, u)$  such that  $u$  descendent of  $v$   
(includes tree edges)

$\text{start}(v) < \text{start}(u) < \text{finish}(u) < \text{finish}(v)$

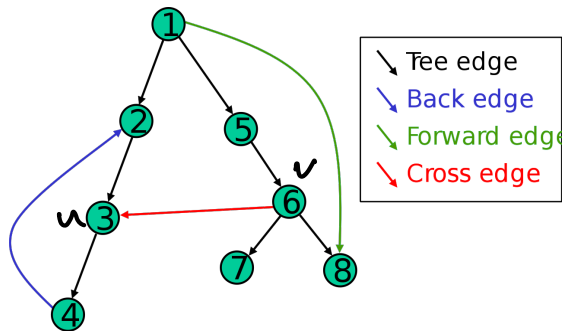
**Back Edges:**  $(v, u)$  such that  $u$  an ancestor of  $v$

**Cross Edges:**  $(v, u)$  such that  $u$  neither a  
descendent nor an ancestor of  $v$



# Edge Types

DFS naturally gives a spanning forest: edge  $(v, u)$  if  $\text{DFS}(v)$  calls  $\text{DFS}(u)$



**Forward Edges:**  $(v, u)$  such that  $u$  descendent of  $v$   
(includes tree edges)

$\text{start}(v) < \text{start}(u) < \text{finish}(u) < \text{finish}(v)$

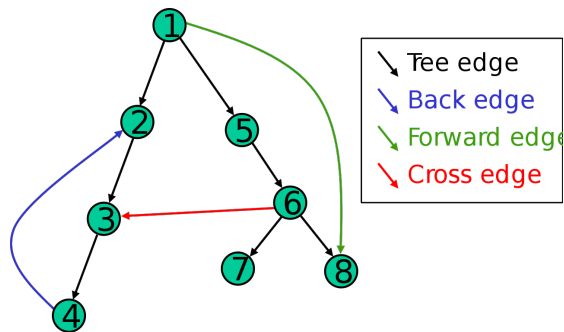
**Back Edges:**  $(v, u)$  such that  $u$  an ancestor of  $v$

$\text{start}(u) < \text{start}(v) < \text{finish}(v) < \text{finish}(u)$

**Cross Edges:**  $(v, u)$  such that  $u$  neither a  
descendent nor an ancestor of  $v$

# Edge Types

DFS naturally gives a spanning forest: edge  $(v, u)$  if  $\text{DFS}(v)$  calls  $\text{DFS}(u)$



**Forward Edges:**  $(v, u)$  such that  $u$  descendent of  $v$   
(includes tree edges)

$\text{start}(v) < \text{start}(u) < \text{finish}(u) < \text{finish}(v)$

**Back Edges:**  $(v, u)$  such that  $u$  an ancestor of  $v$

$\text{start}(u) < \text{start}(v) < \text{finish}(v) < \text{finish}(u)$

**Cross Edges:**  $(v, u)$  such that  $u$  neither a  
descendent nor an ancestor of  $v$

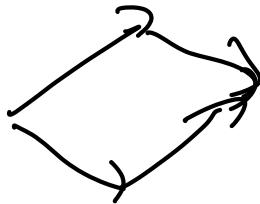
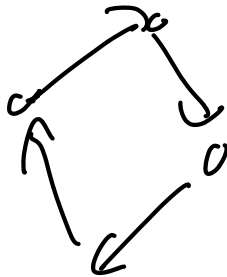
$\text{start}(u) < \text{finish}(u) < \text{start}(v) < \text{finish}(v)$

# Topological Sort

# Definitions

## Definition

A directed graph  $\mathbf{G}$  is a *Directed Acyclic Graph (DAG)* if it has no directed cycles.



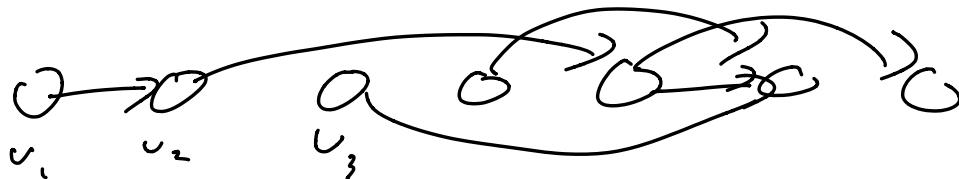
# Definitions

## Definition

A directed graph  $\mathbf{G}$  is a *Directed Acyclic Graph (DAG)* if it has no directed cycles.

## Definition

A *topological sort*  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$  of a DAG is an ordering of the vertices such that all edges are of the form  $(\mathbf{v}_i, \mathbf{v}_j)$  with  $i < j$ .



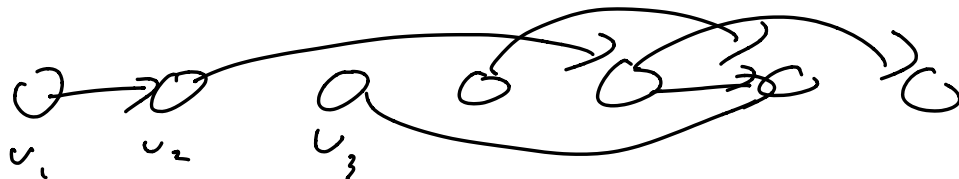
# Definitions

## Definition

A directed graph  $\mathbf{G}$  is a *Directed Acyclic Graph (DAG)* if it has no directed cycles.

## Definition

A *topological sort*  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$  of a DAG is an ordering of the vertices such that all edges are of the form  $(\mathbf{v}_i, \mathbf{v}_j)$  with  $i < j$ .



Can use DFS to characterize DAGs and compute topological sort!

# Characterizing DAGs

## Theorem

*A directed graph  $\mathbf{G}$  is a DAG if and only if  $DFS(\mathbf{G})$  has no back edges.*

# Characterizing DAGs

## Theorem

*A directed graph  $\mathbf{G}$  is a DAG if and only if  $DFS(\mathbf{G})$  has no back edges.*

## Proof.

Only if: contrapositive. If  $\mathbf{G}$  has a back edge:



# Characterizing DAGs

## Theorem

A directed graph  $\mathbf{G}$  is a DAG if and only if  $\text{DFS}(\mathbf{G})$  has no back edges.

## Proof.

Only if: contrapositive. If  $\mathbf{G}$  has a back edge: Directed cycle! Not a DAG.



# Characterizing DAGs

## Theorem

*A directed graph  $\mathbf{G}$  is a DAG if and only if  $\text{DFS}(\mathbf{G})$  has no back edges.*

## Proof.

Only if: contrapositive. If  $\mathbf{G}$  has a back edge: Directed cycle! Not a DAG.

If: contrapositive. If  $\mathbf{G}$  has a directed cycle  $\mathbf{C}$ :

# Characterizing DAGs

## Theorem

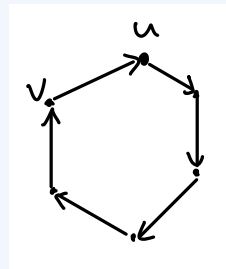
A directed graph  $\mathbf{G}$  is a DAG if and only if  $\text{DFS}(\mathbf{G})$  has no back edges.

## Proof.

Only if: contrapositive. If  $\mathbf{G}$  has a back edge: Directed cycle! Not a DAG.

If: contrapositive. If  $\mathbf{G}$  has a directed cycle  $\mathbf{C}$ :

- ▶ Let  $\mathbf{u} \in \mathbf{C}$  with minimum start value,  $\mathbf{v}$  predecessor in cycle
- ▶ All nodes in  $\mathbf{C}$  reachable from  $\mathbf{u} \implies$  all nodes in  $\mathbf{C}$  descendants of  $\mathbf{u}$
- ▶  $(\mathbf{v}, \mathbf{u})$  a back edge



# Topological Sort

- ▶ Run DFS(**G**)
  - ▶ When DFS(**v**) returns, put **v** at beginning of list

# Topological Sort

- ▶ Run DFS(**G**)
  - ▶ When DFS(**v**) returns, put **v** at beginning of list

**Correctness:** Since **G** a DAG, never see back edge

⇒ Every edge (**v**, **u**) out of **v** a forward or cross edge

⇒ **finish(u)** < **finish(v)**

⇒ **u** already in list



# Topological Sort

- ▶ Run DFS(**G**)
  - ▶ When DFS(**v**) returns, put **v** at beginning of list

**Correctness:** Since **G** a DAG, never see back edge

⇒ Every edge (**v**, **u**) out of **v** a forward or cross edge

⇒ **finish(u)** < **finish(v)**

⇒ **u** already in list

**Running Time:**  $O(m + n)$

# Strongly Connected Components (SCC): Sketch

## Definitions

Another application of DFS. “Kosaraju’s Algorithm”: Developed by Rao Kosaraju, professor emeritus at JHU CS!

$\mathbf{G} = (\mathbf{V}, \mathbf{E})$  a directed graph.

### Definition

$\mathbf{C} \subseteq \mathbf{V}$  is a *strongly connected component (SCC)* if it is a *maximal* subset such that for all  $\mathbf{u}, \mathbf{v} \in \mathbf{C}$ ,  $\mathbf{u}$  can reach  $\mathbf{v}$  and vice versa.



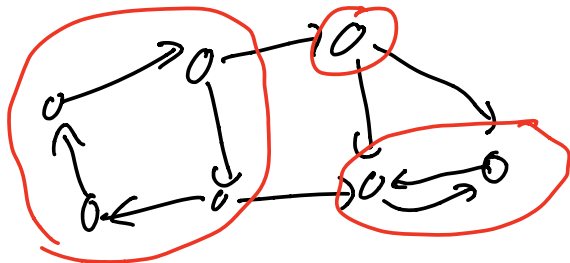
# Definitions

Another application of DFS. “Kosaraju’s Algorithm”: Developed by Rao Kosaraju, professor emeritus at JHU CS!

$G = (V, E)$  a directed graph.

## Definition

$C \subseteq V$  is a *strongly connected component (SCC)* if it is a *maximal* subset such that for all  $u, v \in C$ ,  $u$  can reach  $v$  and vice versa.



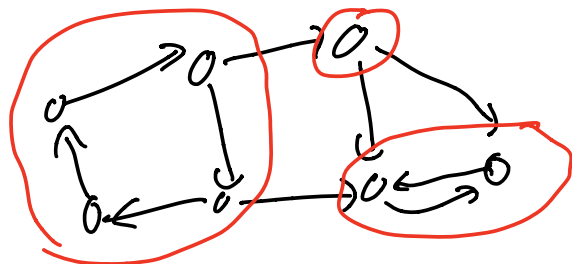
# Definitions

Another application of DFS. “Kosaraju’s Algorithm”: Developed by Rao Kosaraju, professor emeritus at JHU CS!

$\mathbf{G} = (\mathbf{V}, \mathbf{E})$  a directed graph.

## Definition

$\mathbf{C} \subseteq \mathbf{V}$  is a *strongly connected component (SCC)* if it is a *maximal* subset such that for all  $\mathbf{u}, \mathbf{v} \in \mathbf{C}$ ,  $\mathbf{u}$  can reach  $\mathbf{v}$  and vice versa.



**Fact:** There is a *unique* partition of  $\mathbf{V}$  into SCCs

**Proof:** Bireachability is an equivalence relation

# SCC Problem

**Problem:** Given  $\mathbf{G}$  compute SCCs (partition  $\mathbf{V}$  into the SCCs)

# SCC Problem

**Problem:** Given  $\mathbf{G}$  compute SCCs (partition  $\mathbf{V}$  into the SCCs)

**Trivial Algorithm:**

# SCC Problem

**Problem:** Given  $\mathbf{G}$  compute SCCs (partition  $\mathbf{V}$  into the SCCs)

**Trivial Algorithm:** DFS/BFS from every node, keep track of what's reachable from where

# SCC Problem

**Problem:** Given  $\mathbf{G}$  compute SCCs (partition  $\mathbf{V}$  into the SCCs)

**Trivial Algorithm:** DFS/BFS from every node, keep track of what's reachable from where

- ▶ Running time:  $\mathbf{O(n(m + n))}$

# SCC Problem

**Problem:** Given  $\mathbf{G}$  compute SCCs (partition  $\mathbf{V}$  into the SCCs)

**Trivial Algorithm:** DFS/BFS from every node, keep track of what's reachable from where

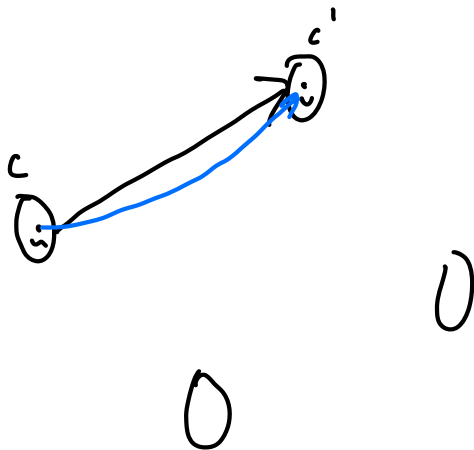
- ▶ Running time:  $\mathbf{O(n(m + n))}$

Can we do better?  $\mathbf{O(m + n)}$ ?

## Graph of SCCs

**Definition:** Let  $\hat{\mathbf{G}}$  be graph of SCCs:

- ▶ Vertex  $\mathbf{v}(\mathbf{C})$  for each SCC  $\mathbf{C}$
- ▶ Edge  $(\mathbf{v}(\mathbf{C}), \mathbf{v}(\mathbf{C}'))$  if  $\exists \mathbf{u} \in \mathbf{C}, \mathbf{v} \in \mathbf{C}'$  such that  $(\mathbf{u}, \mathbf{v}) \in \mathbf{E}$





## Graph of SCCs

**Definition:** Let  $\hat{\mathbf{G}}$  be graph of SCCs:

- ▶ Vertex  $\mathbf{v}(\mathbf{C})$  for each SCC  $\mathbf{C}$
- ▶ Edge  $(\mathbf{v}(\mathbf{C}), \mathbf{v}(\mathbf{C}'))$  if  $\exists \mathbf{u} \in \mathbf{C}, \mathbf{v} \in \mathbf{C}'$  such that  $(\mathbf{u}, \mathbf{v}) \in \mathbf{E}$

### Theorem

$\hat{\mathbf{G}}$  is a DAG.

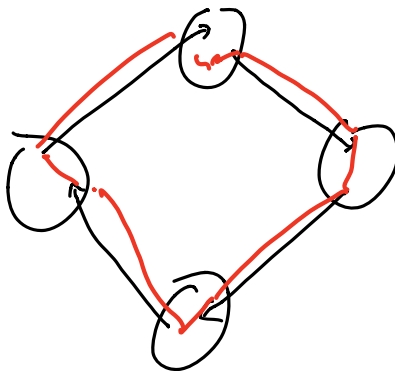
# Graph of SCCs

**Definition:** Let  $\hat{\mathbf{G}}$  be graph of SCCs:

- ▶ Vertex  $\mathbf{v}(\mathbf{C})$  for each SCC  $\mathbf{C}$
- ▶ Edge  $(\mathbf{v}(\mathbf{C}), \mathbf{v}(\mathbf{C}'))$  if  $\exists \mathbf{u} \in \mathbf{C}, \mathbf{v} \in \mathbf{C}'$  such that  $(\mathbf{u}, \mathbf{v}) \in \mathbf{E}$

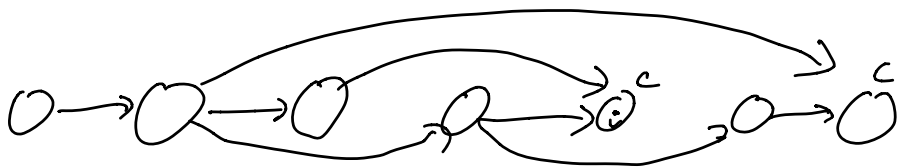
## Theorem

$\hat{\mathbf{G}}$  is a DAG.



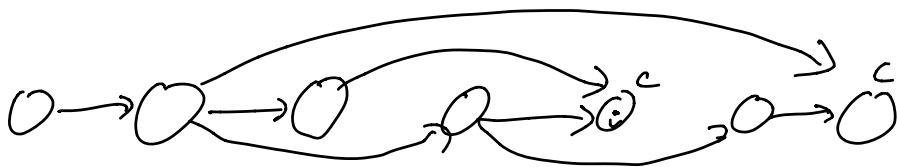
# Sink SCC

Since  $\hat{G}$  a DAG, has a topological sort



# Sink SCC

Since  $\hat{G}$  a DAG, has a topological sort

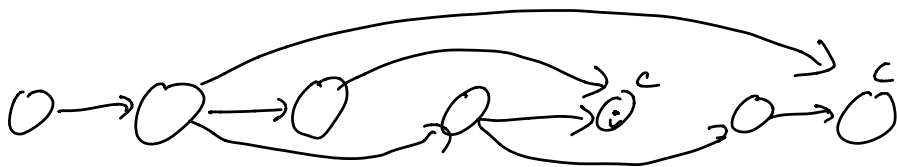


**Definition:** SCC  $C$  is a *sink* SCC if no outgoing edges

- ▶ At least one sink SCC exists

# Sink SCC

Since  $\hat{\mathbf{G}}$  a DAG, has a topological sort



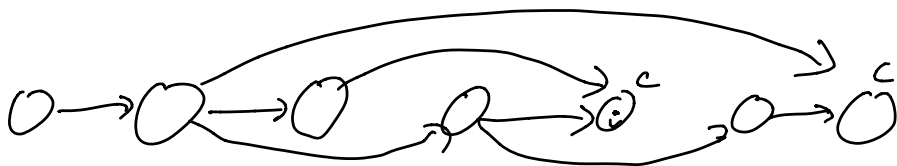
**Definition:** SCC  $\mathbf{C}$  is a *sink* SCC if no outgoing edges

- ▶ At least one sink SCC exists

What happens if we run  $\text{DFS}(\mathbf{v})$  where  $\mathbf{v}$  in a sink SCC?

# Sink SCC

Since  $\hat{\mathbf{G}}$  a DAG, has a topological sort



**Definition:** SCC  $\mathbf{C}$  is a *sink* SCC if no outgoing edges

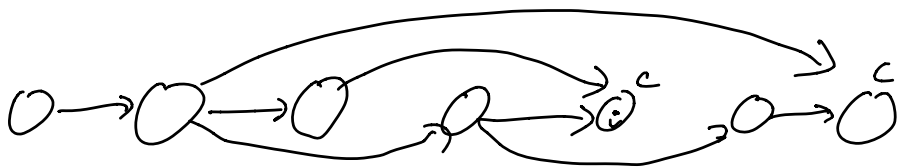
- ▶ At least one sink SCC exists

What happens if we run  $\text{DFS}(\mathbf{v})$  where  $\mathbf{v}$  in a sink SCC?

- ▶ See exactly nodes in  $\mathbf{C}$ !

# Sink SCC

Since  $\hat{\mathbf{G}}$  a DAG, has a topological sort



**Definition:** SCC  $\mathbf{C}$  is a *sink* SCC if no outgoing edges

- ▶ At least one sink SCC exists

What happens if we run  $\text{DFS}(\mathbf{v})$  where  $\mathbf{v}$  in a sink SCC?

- ▶ See exactly nodes in  $\mathbf{C}$ !

Strategy: find node in sink SCC, run DFS, remove nodes found, repeat

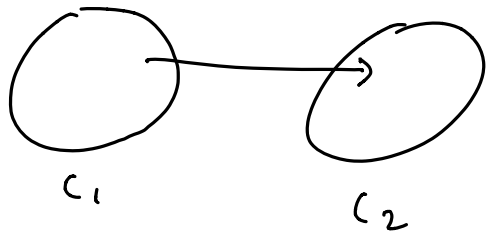
# SCCs and DFS

Run DFS( $\mathbf{G}$ ), and let  $\mathbf{finish}(\mathbf{C}) = \max_{v \in \mathbf{C}} \mathbf{finish}(v)$

## Lemma

Let  $\mathbf{C}_1, \mathbf{C}_2$  distinct SCCs s.t.  $(v(\mathbf{C}_1), v(\mathbf{C}_2)) \in \mathbf{E}(\hat{\mathbf{G}})$ . Then  $\mathbf{finish}(\mathbf{C}_1) > \mathbf{finish}(\mathbf{C}_2)$ .

Let  $x \in \mathbf{C}_1 \cup \mathbf{C}_2$  be first node encountered by DFS



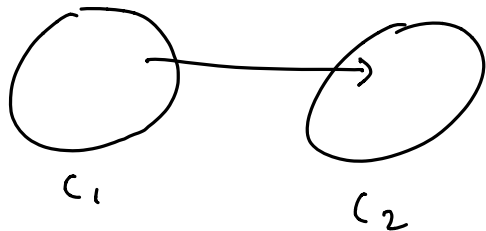


# SCCs and DFS

Run DFS( $\mathbf{G}$ ), and let  $\mathbf{finish}(\mathbf{C}) = \max_{v \in \mathbf{C}} \mathbf{finish}(v)$

## Lemma

Let  $\mathbf{C}_1, \mathbf{C}_2$  distinct SCCs s.t.  $(v(\mathbf{C}_1), v(\mathbf{C}_2)) \in \mathbf{E}(\hat{\mathbf{G}})$ . Then  $\mathbf{finish}(\mathbf{C}_1) > \mathbf{finish}(\mathbf{C}_2)$ .



Let  $x \in \mathbf{C}_1 \cup \mathbf{C}_2$  be first node encountered by DFS

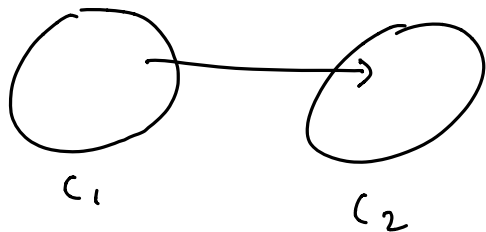
- ▶ If  $x \in \mathbf{C}_1$ :

# SCCs and DFS

Run DFS( $\mathbf{G}$ ), and let  $\mathbf{finish}(\mathbf{C}) = \max_{v \in \mathbf{C}} \mathbf{finish}(v)$

## Lemma

Let  $\mathbf{C}_1, \mathbf{C}_2$  distinct SCCs s.t.  $(v(\mathbf{C}_1), v(\mathbf{C}_2)) \in \mathbf{E}(\hat{\mathbf{G}})$ . Then  $\mathbf{finish}(\mathbf{C}_1) > \mathbf{finish}(\mathbf{C}_2)$ .



Let  $\mathbf{x} \in \mathbf{C}_1 \cup \mathbf{C}_2$  be first node encountered by DFS

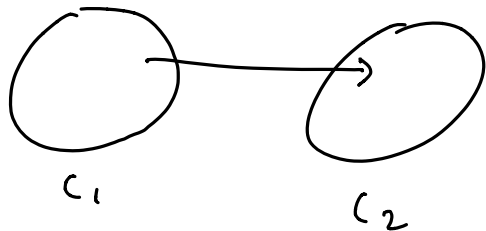
- ▶ If  $\mathbf{x} \in \mathbf{C}_1$ : all of  $\mathbf{C}_2$  reachable from  $\mathbf{x}$ , so DFS( $\mathbf{x}$ ) does not complete until all of  $\mathbf{C}_2$  finished

# SCCs and DFS

Run DFS( $\mathbf{G}$ ), and let  $\mathbf{finish}(\mathbf{C}) = \max_{v \in \mathbf{C}} \mathbf{finish}(v)$

## Lemma

Let  $\mathbf{C}_1, \mathbf{C}_2$  distinct SCCs s.t.  $(v(\mathbf{C}_1), v(\mathbf{C}_2)) \in \mathbf{E}(\hat{\mathbf{G}})$ . Then  $\mathbf{finish}(\mathbf{C}_1) > \mathbf{finish}(\mathbf{C}_2)$ .



Let  $\mathbf{x} \in \mathbf{C}_1 \cup \mathbf{C}_2$  be first node encountered by DFS

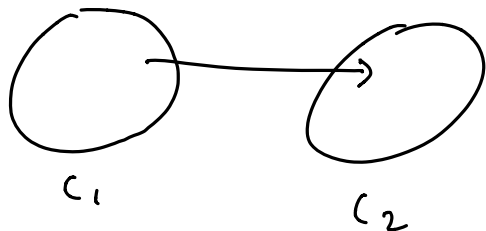
- ▶ If  $\mathbf{x} \in \mathbf{C}_1$ : all of  $\mathbf{C}_2$  reachable from  $\mathbf{x}$ , so DFS( $\mathbf{x}$ ) does not complete until all of  $\mathbf{C}_2$  finished
- ▶ If  $\mathbf{x} \in \mathbf{C}_2$ :

# SCCs and DFS

Run DFS( $\mathbf{G}$ ), and let  $\mathbf{finish}(\mathbf{C}) = \max_{v \in \mathbf{C}} \mathbf{finish}(v)$

## Lemma

Let  $\mathbf{C}_1, \mathbf{C}_2$  distinct SCCs s.t.  $(v(\mathbf{C}_1), v(\mathbf{C}_2)) \in \mathbf{E}(\hat{\mathbf{G}})$ . Then  $\mathbf{finish}(\mathbf{C}_1) > \mathbf{finish}(\mathbf{C}_2)$ .



Let  $\mathbf{x} \in \mathbf{C}_1 \cup \mathbf{C}_2$  be first node encountered by DFS

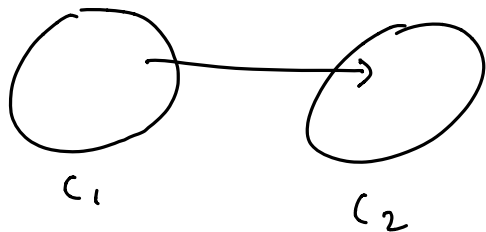
- ▶ If  $\mathbf{x} \in \mathbf{C}_1$ : all of  $\mathbf{C}_2$  reachable from  $\mathbf{x}$ , so DFS( $\mathbf{x}$ ) does not complete until all of  $\mathbf{C}_2$  finished
- ▶ If  $\mathbf{x} \in \mathbf{C}_2$ : all of  $\mathbf{C}_2$  reachable from  $\mathbf{x}$  but nothing from  $\mathbf{C}_1$ , so  $\mathbf{x}$  finishes before any node in  $\mathbf{C}_1$  starts

# SCCs and DFS

Run DFS( $\mathbf{G}$ ), and let  $\mathbf{finish}(\mathbf{C}) = \max_{v \in \mathbf{C}} \mathbf{finish}(v)$

## Lemma

Let  $\mathbf{C}_1, \mathbf{C}_2$  distinct SCCs s.t.  $(v(\mathbf{C}_1), v(\mathbf{C}_2)) \in \mathbf{E}(\hat{\mathbf{G}})$ . Then  $\mathbf{finish}(\mathbf{C}_1) > \mathbf{finish}(\mathbf{C}_2)$ .



Let  $x \in \mathbf{C}_1 \cup \mathbf{C}_2$  be first node encountered by DFS

- ▶ If  $x \in \mathbf{C}_1$ : all of  $\mathbf{C}_2$  reachable from  $x$ , so DFS( $x$ ) does not complete until all of  $\mathbf{C}_2$  finished
- ▶ If  $x \in \mathbf{C}_2$ : all of  $\mathbf{C}_2$  reachable from  $x$  but nothing from  $\mathbf{C}_1$ , so  $x$  finishes before any node in  $\mathbf{C}_1$  starts

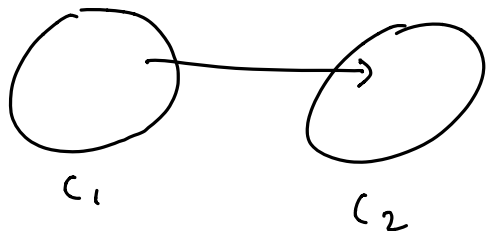
So node with max finish time in a *source* SCC. Want sink.

# SCCs and DFS

Run DFS( $\mathbf{G}$ ), and let  $\mathbf{finish}(\mathbf{C}) = \max_{v \in \mathbf{C}} \mathbf{finish}(v)$

## Lemma

Let  $\mathbf{C}_1, \mathbf{C}_2$  distinct SCCs s.t.  $(v(\mathbf{C}_1), v(\mathbf{C}_2)) \in \mathbf{E}(\hat{\mathbf{G}})$ . Then  $\mathbf{finish}(\mathbf{C}_1) > \mathbf{finish}(\mathbf{C}_2)$ .



Let  $x \in \mathbf{C}_1 \cup \mathbf{C}_2$  be first node encountered by DFS

- ▶ If  $x \in \mathbf{C}_1$ : all of  $\mathbf{C}_2$  reachable from  $x$ , so DFS( $x$ ) does not complete until all of  $\mathbf{C}_2$  finished
- ▶ If  $x \in \mathbf{C}_2$ : all of  $\mathbf{C}_2$  reachable from  $x$  but nothing from  $\mathbf{C}_1$ , so  $x$  finishes before any node in  $\mathbf{C}_1$  starts

So node with max finish time in a *source* SCC. Want sink. Reverse all edges!

# Kosaraju's Algorithm

**Definition:**  $\mathbf{G}^T$  is  $\mathbf{G}$  with all edges reversed.

```
DFS( $\mathbf{G}^T$ ) to get finishing times
while( $\mathbf{G}$  non-empty) {
    Let  $\mathbf{v}$  be vertex in  $\mathbf{G}$  with largest finishing time (from original DFS of  $\mathbf{G}^T$ )
    Run DFS( $\mathbf{v}$ ), let  $\mathbf{C}$  be all nodes found
    Delete  $\mathbf{C}$  from  $\mathbf{G}$ 
}
```

# Kosaraju's Algorithm

**Definition:**  $\mathbf{G}^T$  is  $\mathbf{G}$  with all edges reversed.

```
DFS( $\mathbf{G}^T$ ) to get finishing times
while( $\mathbf{G}$  non-empty) {
    Let  $\mathbf{v}$  be vertex in  $\mathbf{G}$  with largest finishing time (from original DFS of  $\mathbf{G}^T$ )
    Run DFS( $\mathbf{v}$ ), let  $\mathbf{C}$  be all nodes found
    Delete  $\mathbf{C}$  from  $\mathbf{G}$ 
}
```

Some implementation details missing (repeatedly finding max finishing time without using heap): see book



# Kosaraju's Algorithm

**Definition:**  $G^T$  is  $G$  with all edges reversed.

```
DFS( $G^T$ ) to get finishing times
while( $G$  non-empty) {
  Let  $v$  be vertex in  $G$  with largest finishing time (from original DFS of  $G^T$ )
  Run DFS( $v$ ), let  $C$  be all nodes found
  Delete  $C$  from  $G$ 
}
```

Some implementation details missing (repeatedly finding max finishing time without using heap): see book

**Running Time:**  $O(m + n)$

## Correctness Sketch

Let  $C_1, C_2, \dots, C_k$  be set identified by algorithm (in order)

### Theorem

$C_i$  is a sink SCC of  $G \setminus \left( \bigcup_{j=1}^{i-1} C_j \right)$

## Correctness Sketch

Let  $C_1, C_2, \dots, C_k$  be set identified by algorithm (in order)

### Theorem

$C_i$  is a sink SCC of  $G \setminus \left( \bigcup_{j=1}^{i-1} C_j \right)$

Induction on  $i$ .

## Correctness Sketch

Let  $\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_k$  be set identified by algorithm (in order)

### Theorem

$\mathbf{C}_i$  is a sink SCC of  $\mathbf{G} \setminus \left( \bigcup_{j=1}^{i-1} \mathbf{C}_j \right)$

Induction on  $i$ .

**Base case:  $i = 1$ .** By previous argument, largest finishing time in  $\mathbf{G}^T \implies$  in sink SCC of  $\mathbf{G} \implies \mathbf{C}_1$  is sink SCC of  $\mathbf{G}$

## Correctness Sketch

Let  $\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_k$  be set identified by algorithm (in order)

### Theorem

$\mathbf{C}_i$  is a sink SCC of  $\mathbf{G} \setminus \left( \bigcup_{j=1}^{i-1} \mathbf{C}_j \right)$

Induction on  $i$ .

**Base case:**  $i = 1$ . By previous argument, largest finishing time in  $\mathbf{G}^T \implies$  in sink SCC of  $\mathbf{G} \implies \mathbf{C}_1$  is sink SCC of  $\mathbf{G}$

**Inductive case:** Let  $\mathbf{v}$  node remaining with largest finishing time.

- ▶ By induction, current graph is  $\mathbf{G}$  minus  $i - 1$  SCCs of  $\mathbf{G}$
- ▶ Implies  $\mathbf{v}$  must be in sink SCC of remaining graph, so get an SCC of remaining graph when run DFS
- ▶ By induction, also an SCC of original graph