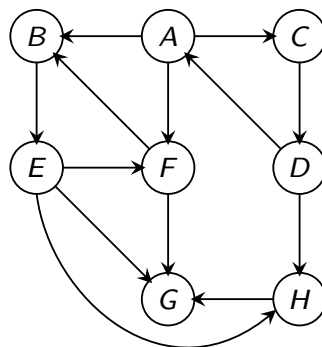


# DFS in Directed Graphs, Strong Connected Components, and DAGs

## Lecture 2

January 20, 2011

## Strong Connected Components (SCCs)



### Algorithmic Problem

Find all **SCCs** of a given directed graph.

Previous lecture: saw an  $O(n \cdot (n + m))$  time algorithm.

This lecture:  $O(n + m)$  time algorithm.

# Graph of SCCs

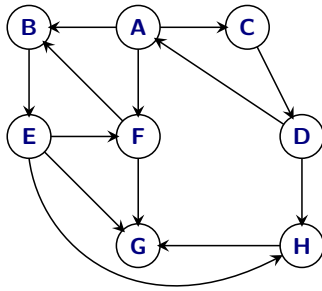


Figure: Graph  $G$

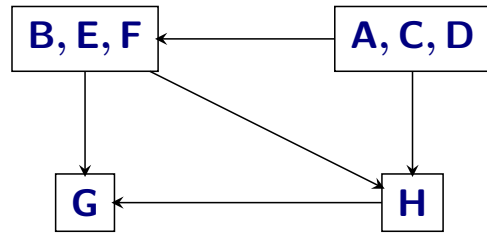


Figure: Graph of SCCs  $G$

## Meta-graph of SCCs

Let  $S_1, S_2, \dots, S_k$  be the SCCs of  $G$ . The graph of SCCs is  $G^{SCC}$

- Vertices are  $S_1, S_2, \dots, S_k$
- There is an edge  $(S_i, S_j)$  if there is some  $u \in S_i$  and  $v \in S_j$  such that  $(u, v)$  is an edge in  $G$ .

## Reversal and SCCs

### Proposition

For any graph  $G$ , the graph of SCCs of  $G^{\text{rev}}$  is the same as the reversal of  $G^{SCC}$ .

### Proof.

Exercise. □

# SCCs and DAGs

## Proposition

For any graph  $G$ , the graph  $G^{\text{SCC}}$  has no directed cycle.

## Proof.

If  $G^{\text{SCC}}$  has a cycle  $S_1, S_2, \dots, S_k$  then  $S_1 \cup S_2 \cup \dots \cup S_k$  is an SCC in  $G$ . Formal details: exercise.  $\square$

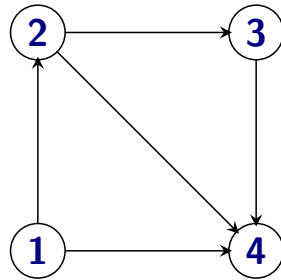
## Part I

# Directed Acyclic Graphs

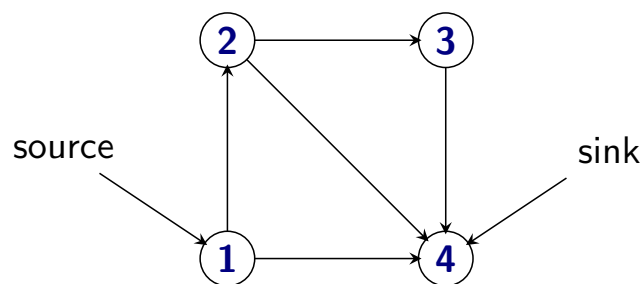
# Directed Acyclic Graphs

## Definition

A directed graph  $G$  is a **directed acyclic graph (DAG)** if there is no directed cycle in  $G$ .



## Sources and Sinks



## Definition

- A vertex  $u$  is a **source** if it has no in-coming edges.
- A vertex  $u$  is a **sink** if it has no out-going edges.

# Simple DAG Properties

- Every **DAG**  $G$  has at least one source and at least one sink.
- If  $G$  is a **DAG** if and only if  $G^{\text{rev}}$  is a **DAG**.
- $G$  is a **DAG** if and only if each node is in its own strong connected component.

Formal proofs: exercise.

# Topological Ordering/Sorting

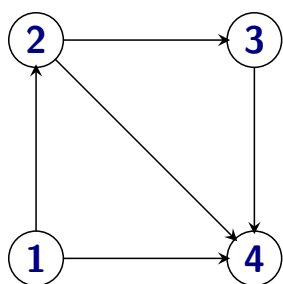


Figure: Graph  $G$

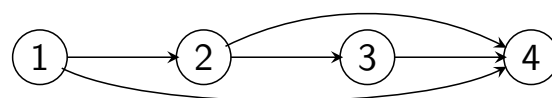


Figure: Topological Ordering of  $G$

## Definition

A **topological ordering/topological sorting** of  $G = (V, E)$  is an ordering  $<$  on  $V$  such that if  $(u, v) \in E$  then  $u < v$ .

# DAGs and Topological Sort

## Lemma

A directed graph  $G$  can be topologically ordered iff it is a DAG.

## Proof.

Only if: Suppose  $G$  is not a DAG and has a topological ordering  $<$ .  
 $G$  has a cycle  $C = u_1, u_2, \dots, u_k, u_1$ .

Then  $u_1 < u_2 < \dots < u_k < u_1$ ! A contradiction.  $\square$

## Proof.

If: Consider the following algorithm:

- Pick a source  $u$ , output it.
- Remove  $u$  and all edges out of  $u$ .
- Repeat until graph is empty.

- Exercise: prove this gives an ordering.  $\square$

Sariel (UIUC)

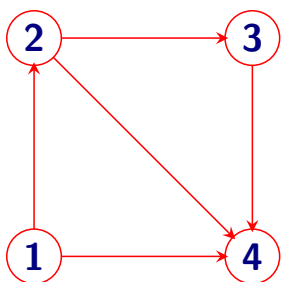
CS473

11

Spring 2011

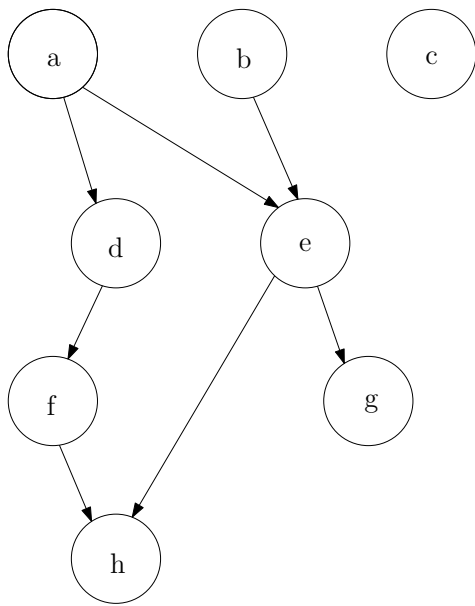
11 / 49

## Topological Sort: An Example



Output: 1 2 3 4

# Topological Sort: Another Example



## DAGs and Topological Sort

**Note:** A **DAG** **G** may have many different topological sorts.

**Question:** What is a **DAG** with the most number of distinct topological sorts for a given number **n** of vertices?

**Question:** What is a **DAG** with the least number of distinct topological sorts for a given number **n** of vertices?

# Using DFS...

... to check for Acyclicity and compute Topological Ordering

## Question

Given  $G$ , is it a **DAG**? If it is, generate a topological sort.

**DFS** based algorithm:

- Compute **DFS**( $G$ )
- If there is a back edge then  $G$  is not a **DAG**.
- Otherwise output nodes in decreasing post-visit order.

Correctness relies on the following:

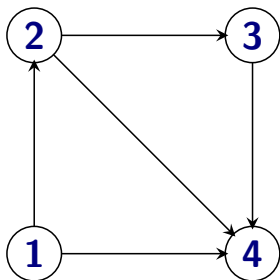
## Proposition

$G$  is a **DAG** iff there is no back-edge in **DFS**( $G$ ).

## Proposition

If  $G$  is a **DAG** and  $post(v) > post(u)$ , then  $(u, v)$  is not in  $G$ .

## Example





# Back edge and Cycles

## Proposition

**G** has a cycle iff there is a back-edge in **DFS(G)**.

## Proof.

If:  $(u, v)$  is a back edge implies there is a cycle **C** consisting of the path from **v** to **u** in **DFS** search tree and the edge  $(u, v)$ .

Only if: Suppose there is a cycle  $C = v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_k \rightarrow v_1$ . Let  $v_i$  be first node in **C** visited in **DFS**.

All other nodes in **C** are descendants of  $v_i$  since they are reachable from  $v_i$ .

Therefore,  $(v_{i-1}, v_i)$  (or  $(v_k, v_1)$  if  $i = 1$ ) is a back edge.  $\square$

# DAGs and Partial Orders

## Definition

A **partially ordered set** is a set **S** along with a binary relation  $\preceq$  such that  $\preceq$  is

- 1 reflexive ( $a \preceq a$  for all  $a \in V$ ),
- 2 anti-symmetric ( $a \preceq b$  and  $a \neq b$  implies  $b \not\preceq a$ ), and
- 3 transitive ( $a \preceq b$  and  $b \preceq c$  implies  $a \preceq c$ ).

**Example:** For numbers in the plane define  $(x, y) \preceq (x', y')$  iff  $x \leq x'$  and  $y \leq y'$ .

**Observation:** A *finite* partially ordered set is equivalent to a **DAG**.

**Observation:** A topological sort of a **DAG** corresponds to a complete (or total) ordering of the underlying partial order.

## Part II

Linear time algorithm for finding all strong connected components of a directed graph

## Finding all SCCs of a Directed Graph

### Problem

Given a directed graph  $\mathbf{G} = (\mathbf{V}, \mathbf{E})$ , output *all* its strong connected components.

Straightforward algorithm:

For each vertex  $\mathbf{u} \in \mathbf{V}$  do

find  $\mathbf{SCC}(\mathbf{G}, \mathbf{u})$  the strong component containing  $\mathbf{u}$  as follows:

Obtain  $\mathbf{rch}(\mathbf{G}, \mathbf{u})$  using  $\mathbf{DFS}(\mathbf{G}, \mathbf{u})$

Obtain  $\mathbf{rch}(\mathbf{G}^{\text{rev}}, \mathbf{u})$  using  $\mathbf{DFS}(\mathbf{G}^{\text{rev}}, \mathbf{u})$

Output  $\mathbf{SCC}(\mathbf{G}, \mathbf{u}) = \mathbf{rch}(\mathbf{G}, \mathbf{u}) \cap \mathbf{rch}(\mathbf{G}^{\text{rev}}, \mathbf{u})$

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

Is there an  $\mathbf{O}(n + m)$  time algorithm?

# Structure of a Directed Graph

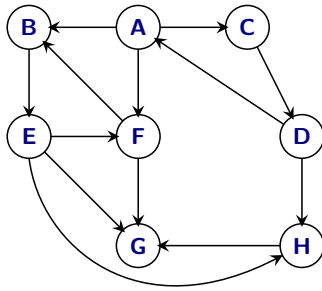


Figure: Graph  $G$

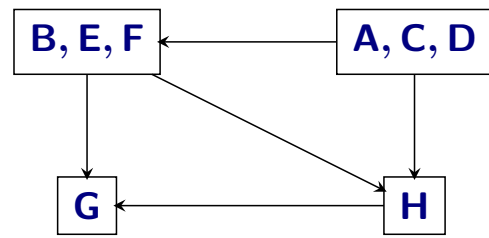


Figure: Graph of SCCs  $G^{SCC}$

## Proposition

For a directed graph  $G$ , its meta-graph  $G^{SCC}$  is a DAG.

## Linear-time Algorithm for SCCs: Ideas

Exploit structure of meta-graph.

### Algorithm

- Let  $u$  be a vertex in a sink SCC of  $G^{SCC}$
- Do  $DFS(u)$  to compute  $SCC(u)$
- Remove  $SCC(u)$  and repeat

### Justification

- $DFS(u)$  only visits vertices (and edges) in  $SCC(u)$
- $DFS(u)$  takes time proportional to size of  $SCC(u)$
- Therefore, total time  $O(n + m)$ !

# Big Challenge(s)

How do we find a vertex in the sink SCC of  $\mathbf{G}^{\text{SCC}}$ ?

Can we obtain an *implicit* topological sort of  $\mathbf{G}^{\text{SCC}}$  without computing  $\mathbf{G}^{\text{SCC}}$ ?

**Answer:**  $\text{DFS}(\mathbf{G})$  gives some information!

## Post-visit times of SCCs

### Definition

Given  $\mathbf{G}$  and a SCC  $\mathbf{S}$  of  $\mathbf{G}$ , define  $\text{post}(\mathbf{S}) = \max_{u \in \mathbf{S}} \text{post}(u)$  where  $\text{post}$  numbers are with respect to some  $\text{DFS}(\mathbf{G})$ .

# An Example

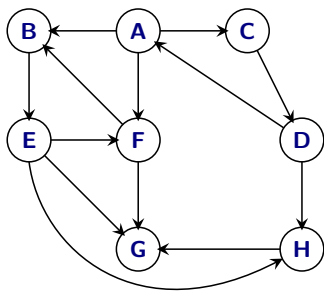


Figure: Graph **G**

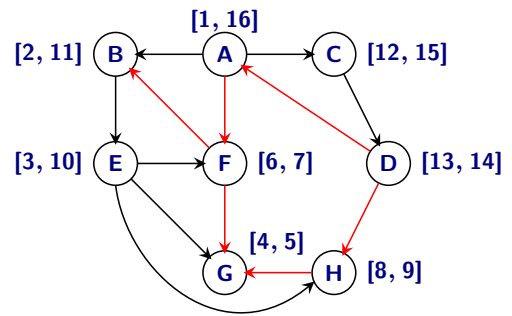


Figure: Graph with pre-post times for **DFS(A)**; black edges in tree

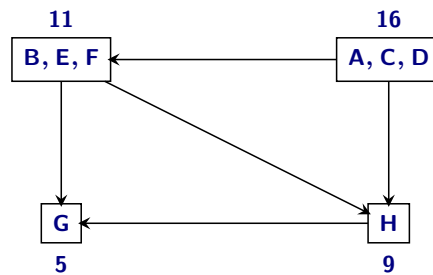


Figure:  $G^{SCC}$  with post times

## $G^{SCC}$ and post-visit times

### Proposition

If **S** and **S'** are **SCCs** in **G** and **(S, S')** is an edge in  $G^{SCC}$  then  $post(S) > post(S')$ .

### Proof.

Let **u** be first vertex in  $S \cup S'$  that is visited.

- If  $u \in S$  then all of **S'** will be explored before **DFS(u)** completes.
- If  $u \in S'$  then all of **S'** will be explored before any of **S**.

□

**A False Statement:** If **S** and **S'** are **SCCs** in **G** and **(S, S')** is an edge in  $G^{SCC}$  then for every  $u \in S$  and  $u' \in S'$ ,  $post(u) > post(u')$ .

# Topological ordering of $G^{SCC}$

## Corollary

Ordering SCCs in decreasing order of  $\text{post}(\mathbf{S})$  gives a topological ordering of  $G^{SCC}$

**Recall:** for a DAG, ordering nodes in decreasing post-visit order gives a topological sort.

So...

**DFS(G)** gives some information on topological ordering of  $G^{SCC}$ !

## Finding Sources

### Proposition

The vertex  $\mathbf{u}$  with the highest post visit time belongs to a source SCC in  $G^{SCC}$

### Proof.

- $\text{post}(\text{SCC}(\mathbf{u})) = \text{post}(\mathbf{u})$
- Thus,  $\text{post}(\text{SCC}(\mathbf{u}))$  is highest and will be output first in topological ordering of  $G^{SCC}$ .



# Finding Sinks

## Proposition

The vertex  $u$  with highest post visit time in  $\text{DFS}(\mathbf{G}^{\text{rev}})$  belongs to a sink SCC of  $\mathbf{G}$ .

## Proof.

- $u$  belongs to source SCC of  $\mathbf{G}^{\text{rev}}$
- Since graph of SCCs of  $\mathbf{G}^{\text{rev}}$  is the reverse of  $\mathbf{G}^{\text{SCC}}$ ,  $\text{SCC}(u)$  is sink SCC of  $\mathbf{G}$ .  $\square$

# Linear Time Algorithm

Do  $\text{DFS}(\mathbf{G}^{\text{rev}})$  and sort vertices in decreasing post order.

Mark all nodes as unvisited

for each  $u$  in the computed order do

  if  $u$  is not visited then

$\text{DFS}(u)$

    Let  $S_u$  be the nodes reached by  $u$

    Output  $S_u$  as a strong connected component

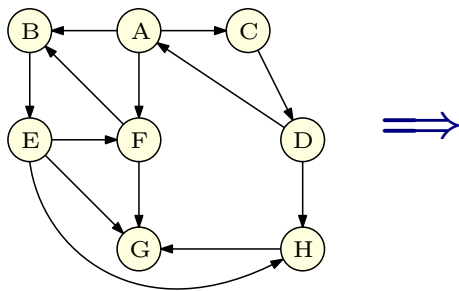
    Remove  $S_u$  from  $\mathbf{G}$

## Analysis

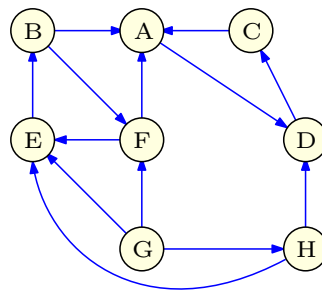
Running time is  $O(n + m)$ . (Exercise)

# Linear Time Algorithm: An Example - Initial steps

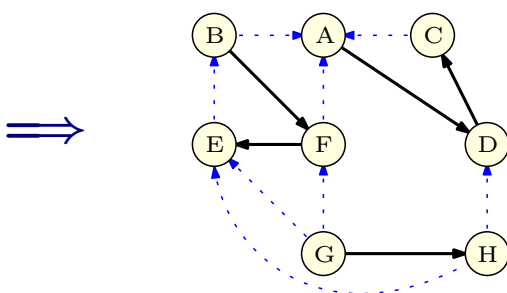
Graph **G**:



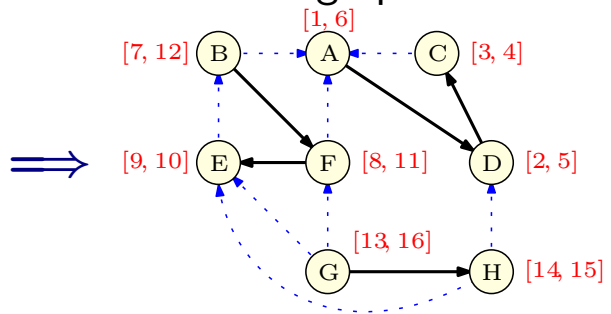
Reverse graph **G<sup>rev</sup>**:



**DFS** of reverse graph:



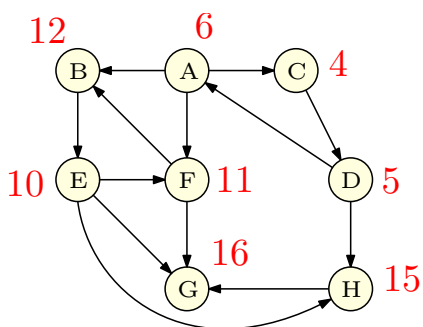
Pre/Post **DFS** numbering of reverse graph:



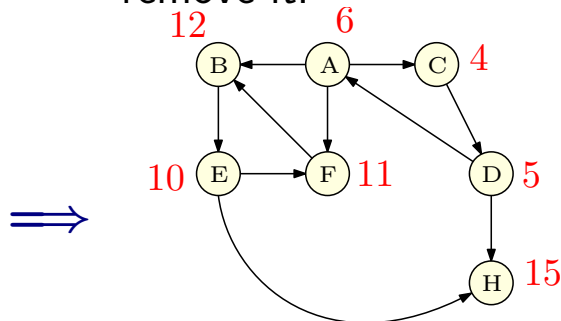
# Linear Time Algorithm: An Example

Removing connected components: 1

Original graph **G** with rev post numbers:



Do **DFS** from vertex **G** remove it.



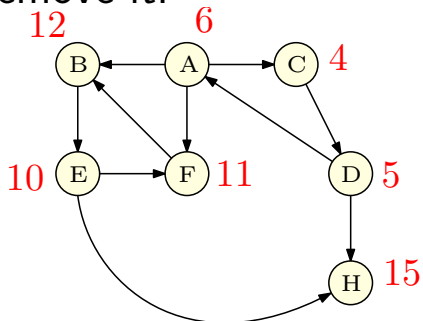
**SCC** computed:  
**{G}**



# Linear Time Algorithm: An Example

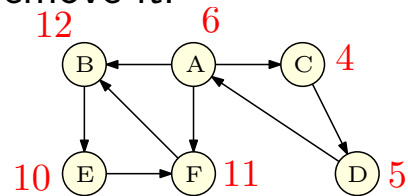
Removing connected components: 2

Do **DFS** from vertex **G**  
remove it.



**SCC** computed:  
**{G}**

Do **DFS** from vertex **H**,  
remove it.

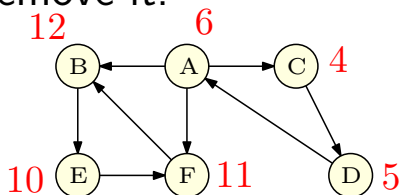


**SCC** computed:  
**{G}, {H}**

# Linear Time Algorithm: An Example

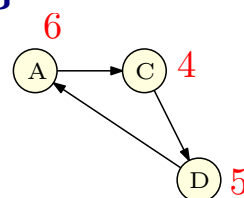
Removing connected components: 3

Do **DFS** from vertex **H**,  
remove it.



**SCC** computed:  
**{G}, {H}**

Do **DFS** from vertex **F**  
Remove visited vertices:  
**{F, B, E}**.

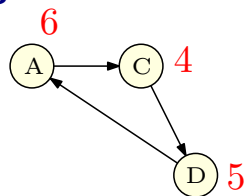


**SCC** computed:  
**{G}, {H}, {F, B, E}**

# Linear Time Algorithm: An Example

Removing connected components: 4

Do **DFS** from vertex **F**  
Remove visited vertices:  
**{F, B, E}**.



**SCC** computed:  
**{G}, {H}, {F, B, E}**

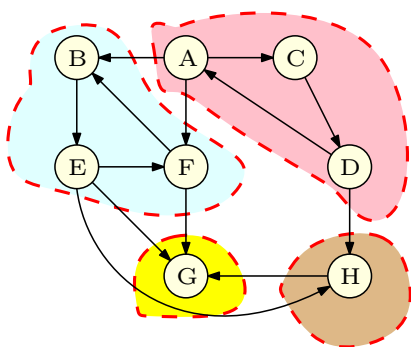
Do **DFS** from vertex **A**  
Remove visited vertices:  
**{A, C, D}**.



**SCC** computed:  
**{G}, {H}, {F, B, E}, {A, C, D}**

# Linear Time Algorithm: An Example

Final result



**SCC** computed:  
**{G}, {H}, {F, B, E}, {A, C, D}**  
Which is the correct answer!

# Obtaining the meta-graph from strong connected components

**Exercise:** Given all the strong connected components of a directed graph  $\mathbf{G} = (\mathbf{V}, \mathbf{E})$  show that the meta-graph  $\mathbf{G}^{\text{SCC}}$  can be obtained in  $O(m + n)$  time.

## Correctness: more details

- let  $\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_k$  be strong components in  $\mathbf{G}$
- Strong components of  $\mathbf{G}^{\text{rev}}$  and  $\mathbf{G}$  are same and meta-graph of  $\mathbf{G}$  is reverse of meta-graph of  $\mathbf{G}^{\text{rev}}$ .
- consider  $\text{DFS}(\mathbf{G}^{\text{rev}})$  and let  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$  be such that  $\text{post}(\mathbf{u}_i) = \text{post}(\mathbf{S}_i) = \max_{\mathbf{v} \in \mathbf{S}_i} \text{post}(\mathbf{v})$ .
- Assume without loss of generality that  $\text{post}(\mathbf{u}_k) > \text{post}(\mathbf{u}_{k-1}) \geq \dots \geq \text{post}(\mathbf{u}_1)$  (renumber otherwise). Then  $\mathbf{S}_k, \mathbf{S}_{k-1}, \dots, \mathbf{S}_1$  is a topological sort of meta-graph of  $\mathbf{G}^{\text{rev}}$  and hence  $\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_k$  is a topological sort of the meta-graph of  $\mathbf{G}$ .
- $\mathbf{u}_k$  has highest post number and  $\text{DFS}(\mathbf{u}_k)$  will explore all of  $\mathbf{S}_k$  which is a sink component in  $\mathbf{G}$ .
- After  $\mathbf{S}_k$  is removed  $\mathbf{u}_{k-1}$  has highest post number and  $\text{DFS}(\mathbf{u}_{k-1})$  will explore all of  $\mathbf{S}_{k-1}$  which is a sink component in remaining graph  $\mathbf{G} - \mathbf{S}_k$ . Formal proof by induction.

## Part III

### An Application to make

### make Utility [Feldman]

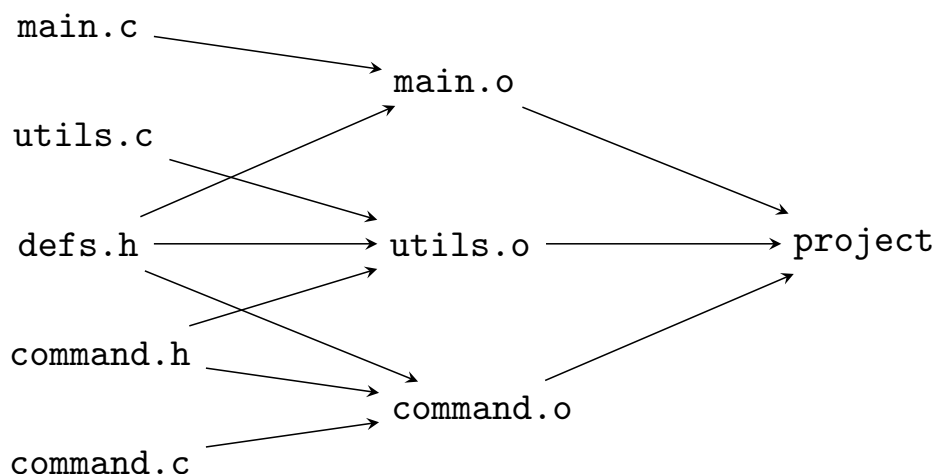
- Unix utility for automatically building large software applications
- A makefile specifies
  - Object files to be created,
  - Source/object files to be used in creation, and
  - How to create them

# An Example makefile

```
project: main.o utils.o command.o
    cc -o project main.o utils.o command.o

main.o: main.c defs.h
    cc -c main.c
utils.o: utils.c defs.h command.h
    cc -c utils.c
command.o: command.c defs.h command.h
    cc -c command.c
```

## makefile as a Digraph



# Computational Problems for make

- Is the `makefile` reasonable?
- If it is reasonable, in what order should the object files be created?
- If it is not reasonable, provide helpful debugging information.
- If some file is modified, find the fewest compilations needed to make application consistent.

# Algorithms for make

- Is the `makefile` reasonable? **Is  $G$  a DAG?**
- If it is reasonable, in what order should the object files be created? **Find a topological sort of a DAG.**
- If it is not reasonable, provide helpful debugging information. **Output a cycle. More generally, output all strong connected components.**
- If some file is modified, find the fewest compilations needed to make application consistent.
  - **Find all vertices reachable (using DFS/BFS) from modified files in directed graph, and recompile them in proper order. Verify that one can find the files to recompile and the ordering in linear time.**

# Take away Points

- Given a directed graph  $G$ , its **SCCs** and the associated acyclic meta-graph  $G^{SCC}$  give a structural decomposition of  $G$  that should be kept in mind.
- There is a **DFS** based linear time algorithm to compute all the **SCCs** and the meta-graph. Properties of **DFS** crucial for the algorithm.
- **DAGs** arise in many application and topological sort is a key property in algorithm design. Linear time algorithms to compute a topological sort (there can be many possible orderings so not unique).