

# CS 473: Fundamental Algorithms

Sariel Har-Peled  
sariel@illinois.edu  
3306 SC

University of Illinois, Urbana-Champaign

Fall 2011

## CS 473: Fundamental Algorithms, Fall 2011

# Administrivia, Primality/Factoring, Problems, Intro to Recursion

Lecture 1  
August 23, 2011

# The word “algorithm” comes from...

Muhammad ibn Musa al-Khwarizmi

780-850 AD

The word “algebra” is taken from the title of one of his books.

## Part I

### Administrivia

# Instructional Staff

- **Instructor:** Sariel Har-Peled (sariel)
- **Teaching Assistants:**
  - JohnMark Lau (johnlau) - TA
  - Smit Shah (ssshah5) - TA
  - Curtis Wang (wang505) - TA
- **Office hours:** See course webpage
- **Email:** See course webpage

## Electronic Bulletin Boards

- **Webpage:** [www.cs.illinois.edu/class/fa11/cs473](http://www.cs.illinois.edu/class/fa11/cs473)
- **Newsgroup:** class.fa11.cs473 on server news.cs.illinois.edu
- **Announcement newsgroup:** class.fa11.cs473.announce on server news.cs.illinois.edu

- **Prerequisites:** CS 173 (discrete math), CS 225 (data structures) and CS 373 (theory of computation)
- **Recommended Textbook:** Algorithms by Dasgupta, Papadimitriou, Vazirani.  
Available online for free!
- **Not-recommended Textbook:** Algorithm Design, Kleinberg & Tardos
- **Lecture Notes:** Available on the web-page after every class
- **Additional References**
  - Algorithms: Dasgupta, Papadimitriou, and Vazirani.
  - Previous class notes of Jeff Erickson and the instructor.
  - Introduction to Algorithms: Cormen, Leiserson, Rivest, Stein.
  - Computers and Intractability: Garey and Johnson.

## Prerequisites

- **Asymptotic notation:**  $O()$ ,  $\Omega()$ ,  $o()$ .
- **Discrete Structures:** sets, functions, relations, equivalence classes, partial orders, trees, graphs
- **Logic:** predicate logic, boolean algebra
- **Proofs:** **by induction**, by contradiction
- **Basic sums and recurrences:** sum of a geometric series, unrolling of recurrences, basic calculus
- **Data Structures:** arrays, multi-dimensional arrays, linked lists, trees, balanced search trees, heaps
- **Abstract Data Types:** lists, stacks, queues, dictionaries, priority queues
- **Algorithms:** sorting (merge, quick, insertion), pre/post/in order traversal of trees, depth/breadth first search of trees (maybe graphs)
- **Basic analysis of algorithms:** loops and nested loops, deriving recurrences from a recursive program
- **Concepts from Theory of Computation:** languages, automata, Turing machine, undecidability, non-determinism
- **Programming:** in some general purpose language
- **Elementary Discrete Probability:** event, random variable, independence
- **Mathematical maturity**

# Grading Policy: Overview

- **Quizzes:** 5%
- **Homeworks:** 20%
- **Midterms:** 40% (**2 × 20**)
- **Finals:** 35% (covers the full course content)

## Homeworks

- One quiz every week: Due by midnight on Sunday.
- One homework every week: Assigned on Monday and due the following Monday by midnight.
- Submit online only!
- Homeworks can be worked on in groups of up to 3 and each group submits *one* written solution (except Homework 0).
  - Short quiz-style questions to be answered individually in Compass.
- Groups can be changed a *few* times only
- Unlike previous years no *oral* homework this semester due to large enrollment and insufficient TA support

# More on Homeworks

- No extensions or late homeworks accepted.
- To compensate, the homework with the least score will be dropped in calculating the homework average.
- **Important:** Read homework faq/instructions on website.

## Discussion Sessions

- 50min problem solving session led by TAs
- **Attendance is required.**
- Five sections all in SC 1214.
  - Tuesday
    - 4–4:50pm,
    - 5–5:50pm.
  - Wednesday
    - 2–2:50pm,
    - 3–3:50pm,
    - 6–6:50pm.

- Attend lectures, please ask plenty of questions.
- Attend discussion sessions.
- Don't skip homework and don't copy homework solutions.
- Study regularly and keep up with the course.
- Ask for help promptly. Make use of office hours.

## Homeworks

- HW 0 is posted on the class website.
- due on Monday midnight.
- Online submission.
- HW 0 to be turned in individually.

# Part II

## Course Goals and Overview

## Topics

- Some useful basic algorithms and data structures
- Broadly applicable techniques in algorithm design
  - Understanding problem structure
  - Brute force enumeration and backtrack search
  - Reductions
  - Recursion
    - Divide and Conquer
    - Dynamic Programming
  - Greedy methods
  - Network Flows and Linear/Integer Programming (optional)
- Analysis techniques
  - Correctness of algorithms via induction and other methods
  - Recurrences
  - Amortization and elementary potential functions
- Polynomial-time Reductions, NP-Completeness, Heuristics



- 
- Appreciate the importance of algorithms in computer science and beyond (engineering, mathematics, natural sciences, social sciences, ...)
- Understand/appreciate limits of computation (intractability)
- Learn/remember some basic tricks, algorithms, problems, ideas
- Have fun!!!

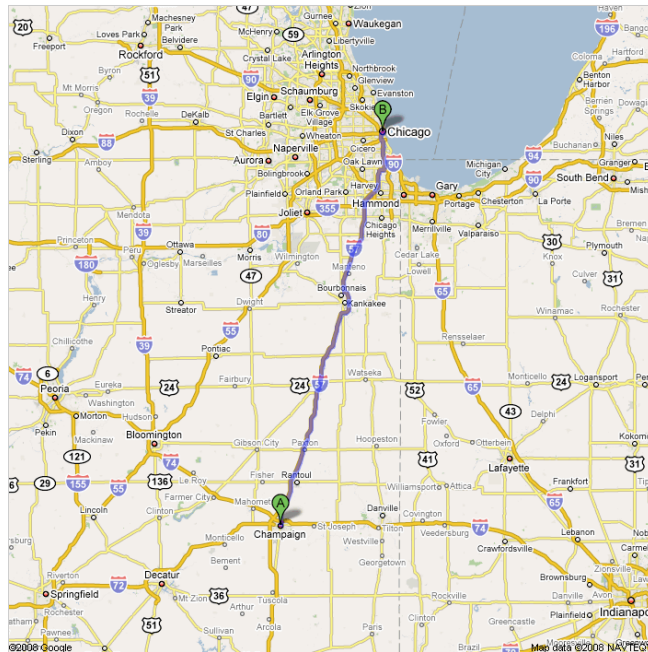
## Part III

# Algorithmic Problems in the Real World

# Shortest Paths

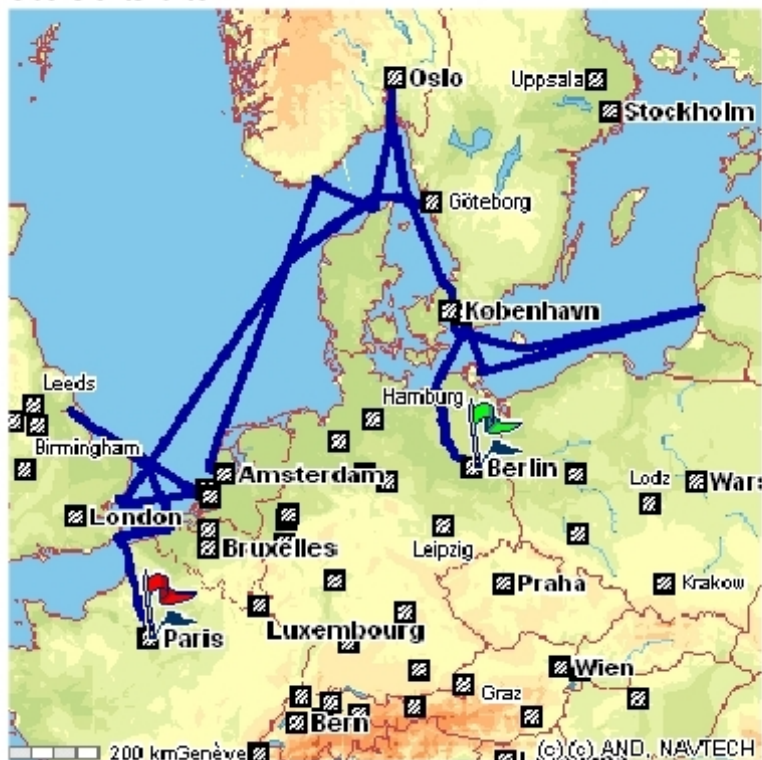


Directions to Chicago, IL  
136 mi – about 2 hours 20 mins



# Shortest Paths - Paris to Berlin

Übersichtskarte



Lossless compression:

Lossy compression: music, images, video (JPEG, MPEG standards)

## Search and Indexing: String Matching and Link Analysis

- Web search: Google, Yahoo!, Microsoft, Ask, ...
- Text search: Text editors (Emacs, Word, Browsers, ...)
- Regular expression search: grep, egrep, emacs, Perl, Awk, compilers

# Public-Key Cryptography

Foundation of Electronic Commerce

RSA Crypto-system: generate key  $n = pq$  where  $p, q$  are *primes*

**Primality:** Given a number  $N$ , check if  $N$  is a prime or composite.

**Factoring:** Given a composite number  $N$ , find a non-trivial factor

## Programming: Parsing and Debugging

```
[godavari: /temp/test] chekuri % gcc main.c
```

**Parsing:** Is main.c a syntactically valid C program?

**Debugging:** Will main.c go into an infinite loop on some input?

**Easier problem ???** Will main.c halt on the specific input 10?

Find the cheapest of most profitable way to do things

- Airline schedules - AA, Delta, ...
- Vehicle routing - trucking and transportation (UPS, FedEx, Union Pacific, ...)
- Network Design - AT&T, Sprint, Level3 ...

Linear and Integer programming problems

## Part IV

# Algorithm Design

# Important Ingredients in Algorithm Design

- What is the problem (really)?
  - What is the input? How is it represented?
  - What is the output?
- What is the model of computation?
- What basic operations are allowed?
- Algorithm design
- Analysis of correctness, running time, space etc.
- Algorithmic engineering: evaluation an understanding of algorithm's performance in practice, performance tweaks, comparison with other algorithms etc. (Not covered in this course)

## Primality testing

### Problem

Given an integer  $N > 0$ , is  $N$  a prime?

### SimpleAlgorithm:

```
for  $i = 2$  to  $\lfloor \sqrt{N} \rfloor$  do
  if  $i$  divides  $N$  then
    return ``COMPOSITE''
return ``PRIME''
```

Correctness? If  $N$  is composite, at least one factor in  $\{2, \dots, \sqrt{N}\}$   
Running time?  $O(\sqrt{N})$  divisions? Sub-linear in input size! **Wrong!**

# Primality testing

...Polynomial means... in input size

How many bits to represent  $N$  in binary?  $\lceil \log N \rceil$  bits.

Simple Algorithm takes  $\sqrt{N} = 2^{(\log N)/2}$  time.

*Exponential* in the input size  $n = \log N$ .

- Modern cryptography: binary numbers with 128, 256, 512 bits.
- Simple Algorithm will take  $2^{64}$ ,  $2^{128}$ ,  $2^{256}$  steps!
- Fastest computer today about 3 petaFlops/sec:  $3 \times 2^{50}$  floating point ops/sec.

## Lesson:

Pay attention to representation size in analyzing efficiency of algorithms. Especially in *number* problems.

## Efficient algorithms

So, is there an *efficient/good/effective* algorithm for primality?

## Question:

What does efficiency mean?

In this class *efficiency* is broadly equated to *polynomial time*.

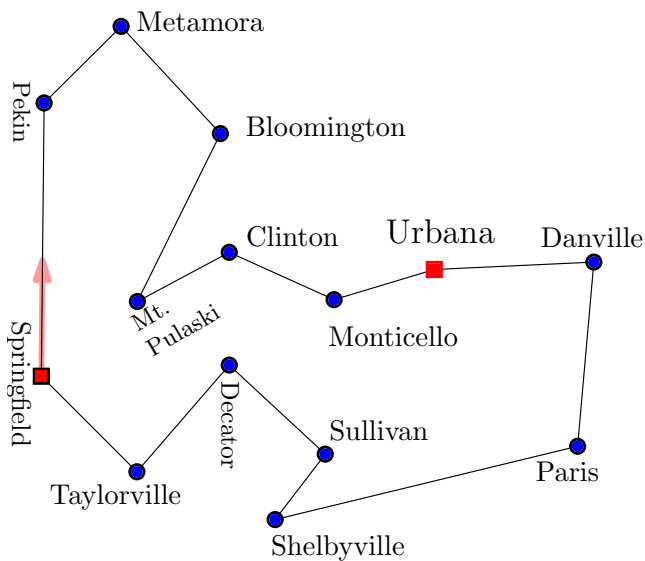
$O(n)$ ,  $O(n \log n)$ ,  $O(n^2)$ ,  $O(n^3)$ ,  $O(n^{100})$ , ... where  $n$  is size of the input.

Why? Is  $n^{100}$  really efficient/practical? Etc.

Short answer: polynomial time is a robust, mathematically sound way to define efficiency. Has been useful for several decades.

# TSP problem

## Lincoln's tour



- Circuit court - ride through counties staying a few days in each town.
- Lincoln was a lawyer traveling with the Eighth Judicial Circuit.
- Picture: travel during 1850.
  - Very close to optimal tour.
  - Might have been optimal at the time..

## Solving TSP by a Computer

### Is it hard?

- $n$  = number of cities.
- $n^2$ : size of input.
- Number of possible solutions is

$$n * (n - 1) * (n - 2) * \dots * 2 * 1 = n!.$$

- $n!$  grows very quickly as  $n$  grows.
  - $n = 10$ :  $n! \approx 3628800$
  - $n = 50$ :  $n! \approx 3 * 10^{64}$
  - $n = 100$ :  $n! \approx 9 * 10^{157}$



# Solving TSP by a Computer

Fastest computer...

- Fastest super computer can do (roughly)

$$2.5 * 10^{15}$$

operations a second.

- Assume: computer checks  $2.5 * 10^{15}$  solutions every second, then...
  - $n = 20 \implies$  2 hours.
  - $n = 25 \implies$  200 years.
  - $n = 37 \implies 2 * 10^{20}$  years!!!

## What is a good algorithm?

Running time...

Input size	$n^2$ ops	$n^3$ ops	$n^4$ ops	$n!$ ops
5	0 secs	0 secs	0 secs	0 secs
20	0 secs	0 secs	0 secs	16 mins
30	0 secs	0 secs	0 secs	$3 \cdot 10^9$ years
100	0 secs	0 secs	0 secs	never
8000	0 secs	0 secs	1 secs	never
16000	0 secs	0 secs	26 secs	never
32000	0 secs	0 secs	6 mins	never
64000	0 secs	0 secs	111 mins	never
200,000	0 secs	3 secs	7 days	never
2,000,000	0 secs	53 mins	202.943 years	never
$10^8$	4 secs	12.6839 years	$10^9$ years	never
$10^9$	6 mins	12683.9 years	$10^{13}$ years	never

# What is a good algorithm?

Running time...

ALL RIGHTS RESERVED  
<http://www.cartoonbank.com>



## Primes is in **P**!

### Theorem (Agrawal-Kayal-Saxena'02)

*There is a polynomial time algorithm for primality.*

First polynomial time algorithm for testing primality. Running time is  $O(\log^{12} N)$  further improved to about  $O(\log^6 N)$  by others. In terms of input size  $n = \log N$ , time is  $O(n^6)$ .

Breakthrough announced in August 2002. Three days later announced in New York Times. Only 9 pages!

Neeraj Kayal and Nitin Saxena were undergraduates at IIT-Kanpur!

# What about before 2002?

Primality testing a key part of cryptography. What was the algorithm being used before 2002?

Miller-Rabin *randomized* algorithm:

- runs in polynomial time:  $O(\log^3 N)$  time
- if  $N$  is prime correctly says “yes”.
- if  $N$  is composite it says “yes” with probability at most  $1/2^{100}$  (can be reduced further at the expense of more running time).

Based on Fermat’s little theorem and some basic number theory.

## Factoring

- Modern public-key cryptography based on RSA (Rivest-Shamir-Adelman) system.
- Relies on the difficulty of factoring a composite number into its prime factors.
- There is a polynomial time algorithm that decides whether a given number  $N$  is prime or not (hence composite or not) but no known polynomial time algorithm to factor a given number.

### Lesson

Intractability can be useful!

Three variants of problems.

- **Decision problem:** answer is yes or no.  
**Example:** Given integer  $N$ , is it a composite number?
- **Search problem:** answer is a feasible solution if it exists.  
**Example:** Given integer  $N$ , if  $N$  is composite output a non-trivial factor  $p$  of  $N$ .
- **Optimization problem:** answer is the *best* feasible solution (if one exists).  
**Example:** Given integer  $N$ , if  $N$  is composite output the *smallest* non-trivial factor  $p$  of  $N$ .

For a given underlying problem:

$$\text{Optimization} \geq \text{Search} \geq \text{Decision}$$

## Quantum Computing

### Theorem (Shor'1994)

*There is a polynomial time algorithm for factoring on a quantum computer.*

RSA and current commercial cryptographic systems can be broken if a quantum computer can be built!

### Lesson

Pay attention to the model of computation.

Many many different problems.

- Adding two numbers: efficient and simple algorithm
- Sorting: efficient and not too difficult to design algorithm
- Primality testing: simple and basic problem, took a long time to find efficient algorithm
- Factoring: no efficient algorithm known.
- Halting problem: important problem in practice, undecidable!

## Multiplying Numbers

**Problem** Given two  $n$ -digit numbers  $x$  and  $y$ , compute their product.

### Grade School Multiplication

Compute “partial product” by multiplying each digit of  $y$  with  $x$  and adding the partial products.

$$\begin{array}{r} 3141 \\ \times 2718 \\ \hline 25128 \\ 3141 \\ 21987 \\ 6282 \\ \hline 8537238 \end{array}$$

# Time analysis of grade school multiplication

- Each partial product:  $\Theta(n)$  time
- Number of partial products:  $\leq n$
- Adding partial products:  $n$  additions each  $\Theta(n)$  (Why?)
- Total time:  $\Theta(n^2)$
- Is there a faster way?

## Fast Multiplication

Best known algorithm:  $O(n \log n \cdot 2^{O(\log^* n)})$  time [Furer 2008]

Previous best time:  $O(n \log n \log \log n)$  [Schönhage-Strassen 1971]

**Conjecture:** there exists an  $O(n \log n)$  time algorithm

We don't fully understand multiplication!

Computation and algorithm design is non-trivial!

# Course Approach

Algorithm design requires a mix of skill, experience, mathematical background/maturity and ingenuity.

Approach in this class and many others:

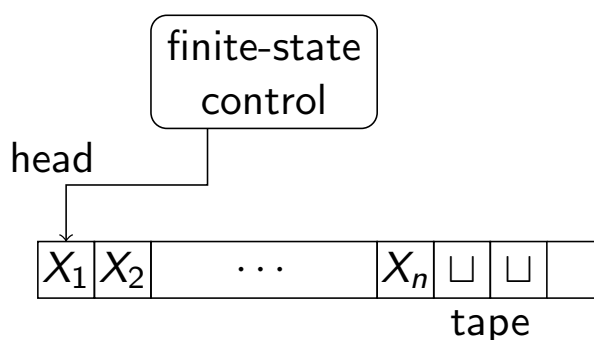
- Improve skills by showing various tools in the abstract and with concrete examples
- Improve experience by giving **many** problems to solve
- Motivate and inspire
- Creativity: you are on your own!

## What model of computation do we use?

Turing Machine?

# Turing Machines: Recap

- Infinite tape
- Finite state control
- Input at beginning of tape
- Special tape letter “blank”  
□
- Head can move only one cell to left or right



# Turing Machines

- Basic unit of data is a bit (or a single character from a finite alphabet)
- Algorithm is the finite control
- Time is number of steps/head moves

## Pros and Cons:

- theoretically sound, robust and simple model that underpins computational complexity.
- polynomial time equivalent to any reasonable “real” computer: Church-Turing thesis
- too low-level and cumbersome, does not model actual computers for many realistic settings



# “Real” Computers vs Turing Machines

How do computers in use differ from TMs?

- random access to memory
- pointers
- arithmetic operations (addition, subtraction, multiplication, division) in constant time

How do they do it?

- basic data type is a word: currently 64 bits
- arithmetic on words are basic instructions of computer
- memory requirements assumed to be  $\leq 2^{64}$  which allows for pointers and indirect addressing as well as random access

## Unit-Cost RAM Model

Informal description:

- Basic data type is an integer/floating point number
- Numbers in input fit in a word
- Arithmetic/comparison operations on words take constant time
- Arrays allow random access (constant time to access  $A[i]$ )
- Pointer based data structures via storing addresses in a word

# Example

Sorting: input is an array of  $n$  numbers

- input size is  $n$  (ignore the bits in each number),
- comparing two numbers takes  $O(1)$  time,
- random access to array elements,
- addition of indices takes constant time,
- basic arithmetic operations take constant time,
- reading/writing one word from/to memory takes constant time.

We will usually not allow (or be careful about allowing):

- bitwise operations (and, or, xor, shift, etc).
- floor function.
- limit word size (usually assume unbounded word size).

## Caveats of RAM Model

Unit-Cost RAM model is applicable in wide variety of settings in practice. However it is not a proper model in several important situations so one has to be careful.

- For some problems such as basic arithmetic computation, unit-cost model makes no sense. Examples: multiplication of two  $n$ -digit numbers, primality etc.
- Input data is very large and does not satisfy the assumptions that individual numbers fit into a word or that total memory is bounded by  $2^k$  where  $k$  is word length.
- Assumptions valid only for certain type of algorithms that do not create large numbers from initial data. For example, exponentiation creates very big numbers from initial numbers.

# Models used in class

In this course:

- Assume unit-cost **RAM** by default.
- We will explicitly point out where unit-cost RAM is not applicable for the problem at hand.

## Part V

### Graph Basics

# Why Graphs?

- Graphs help model networks which are ubiquitous: transportation networks (rail, roads, airways), social networks (interpersonal relationships), information networks (web page links) etc etc.
- Fundamental objects in Computer Science, Optimization, Combinatorics
- Many important and useful optimization problems are graph problems
- Graph theory: elegant, fun and deep mathematics

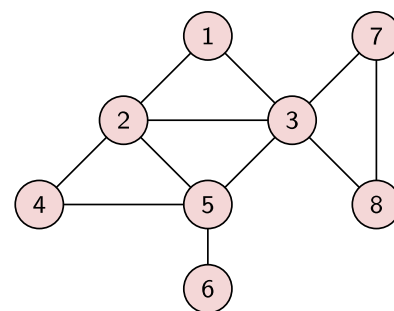
## Graph

### Definition

An undirected (simple) graph

$G = (V, E)$  is a 2-tuple:

- $V$  is a set of vertices (also referred to as nodes/points)
- $E$  is a set of edges where each edge  $e \in E$  is a set of the form  $\{u, v\}$  with  $u, v \in V$  and  $u \neq v$ .



### Example

In figure,  $G = (V, E)$  where  $V = \{1, 2, 3, 4, 5, 6, 7, 8\}$  and  $E = \{\{1, 2\}, \{1, 3\}, \{2, 3\}, \{2, 4\}, \{2, 5\}, \{3, 5\}, \{3, 7\}, \{3, 8\}, \{4, 5\}, \{5, 6\}, \{7, 8\}\}$ .

## Notation

An edge in an undirected graphs is an *unordered* pair of nodes and hence it is a set. Conventionally we use  $(u, v)$  for  $\{u, v\}$  when it is clear from the context that the graph is undirected.

- $u$  and  $v$  are the **end points** of an edge  $\{u, v\}$
- **Multi-graphs** allow
  - *loops* which are edges with the same node appearing as both end points
  - *multi-edges*: different edges between same pairs of nodes
- In this class we will assume that a graph is a simple graph unless explicitly stated otherwise.

## Graph Representation I

### Adjacency Matrix

Represent  $G = (V, E)$  with  $n$  vertices and  $m$  edges using a  $n \times n$  adjacency matrix  $A$  where

- $A[i, j] = A[j, i] = 1$  if  $\{i, j\} \in E$  and  $A[i, j] = A[j, i] = 0$  if  $\{i, j\} \notin E$ .
- Advantage: can check if  $\{i, j\} \in E$  in  $O(1)$  time
- Disadvantage: needs  $\Omega(n^2)$  space even when  $m \ll n^2$

## Adjacency Lists

Represent  $G = (V, E)$  with  $n$  vertices and  $m$  edges using adjacency lists:

- For each  $u \in V$ ,  $\text{Adj}(u) = \{v \mid \{u, v\} \in E\}$ , that is neighbors of  $u$ . Sometimes  $\text{Adj}(u)$  is the list of edges incident to  $u$ .
- Advantage: space is  $O(m + n)$
- Disadvantage: cannot “easily” determine in  $O(1)$  time whether  $\{i, j\} \in E$ 
  - By sorting each list, one can achieve  $O(\log n)$  time
  - By hashing “appropriately”, one can achieve  $O(1)$  time

**Note:** In this class we will assume that by default, graphs are represented using plain vanilla (unsorted) adjacency lists.

## Connectivity

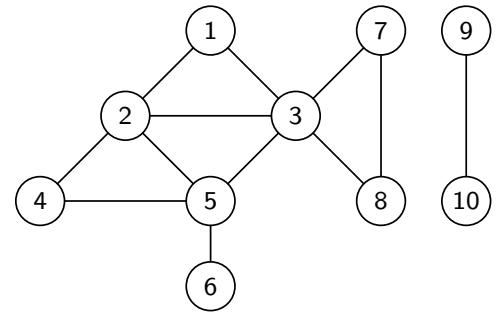
Given a graph  $G = (V, E)$ :

- A **path** is a sequence of *distinct* vertices  $v_1, v_2, \dots, v_k$  such that  $\{v_i, v_{i+1}\} \in E$  for  $1 \leq i \leq k - 1$ . The length of the path is  $k - 1$  and the path is from  $v_1$  to  $v_k$
- A **cycle** is a sequence of *distinct* vertices  $v_1, v_2, \dots, v_k$  such that  $\{v_i, v_{i+1}\} \in E$  for  $1 \leq i \leq k - 1$  and  $\{v_1, v_k\} \in E$ .
- A vertex  $u$  is **connected** to  $v$  if there is a path from  $u$  to  $v$ .
- The **connected component** of  $u$ ,  $\text{con}(u)$ , is the set of all vertices connected to  $u$ .

# Connectivity contd

Define a relation  $C$  on  $V \times V$  as  $uCv$  if  $u$  is connected to  $v$

- In undirected graphs, connectivity is a reflexive, symmetric, and transitive relation. Connected components are the equivalence classes.
- Graph is **connected** if only one connected component



## Connectivity Problems

### Algorithmic Problems

- Given graph  $G$  and nodes  $u$  and  $v$ , is  $u$  connected to  $v$ ?
- Given  $G$  and node  $u$ , find all nodes that are connected to  $u$ .
- Find all connected components of  $G$ .

Can be accomplished in  $O(m + n)$  time using **BFS** or **DFS**.

# Basic Graph Search

Given  $G = (V, E)$  and vertex  $u \in V$ :

**Explore**( $u$ ):

Initialize  $S = \{u\}$

**while** there is an edge  $(x, y)$  with  $x \in S$  and  $y \notin S$  **do**  
    add  $y$  to  $S$

## Proposition

**Explore**( $u$ ) terminates with  $S = \text{con}(u)$ .

Running time: depends on implementation

- Breadth First Search (**BFS**): use **queue** data structure
- Depth First Search (**DFS**): use **stack** data structure
- Review CS 225 material!

## Part VI

### DFS



# Depth First Search

**DFS** is a very versatile graph exploration strategy. Hopcroft and Tarjan (Turing Award winners) demonstrated the power of **DFS** to understand graph structure. **DFS** can be used to obtain linear time ( $O(m + n)$ ) time algorithms for

- Finding cut-edges and cut-vertices of undirected graphs
- Finding strong connected components of directed graphs
- Linear time algorithm for testing whether a graph is planar

## DFS in Undirected Graphs

Recursive version.

**DFS(G)**

Mark all nodes  $u$  as unvisited

While there is an unvisited node  $u$  do

**DFS( $u$ )**

**DFS( $u$ )**

Mark  $u$  as visited

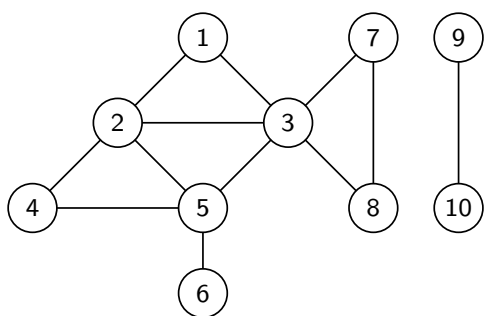
for each edge  $(u,v)$  in  $A_{jd}(u)$  do

if  $v$  is not marked

**DFS( $v$ )**

Global array Mark for all recursive calls.

# Example



## DFS Tree/Forest

### DFS( $G$ )

Mark all nodes  $u$  as unvisited

$T$  is set to  $\emptyset$

While there is an unvisited node  $u$  do

    DFS( $u$ )

Output  $T$

### DFS( $u$ )

Mark  $u$  as visited

for each edge  $(u, v)$  in  $Ajd(u)$  do

    if  $v$  is not marked

        add edge  $(u, v)$  to  $T$

        DFS( $v$ )

Edges classified into two types:  $(u, v) \in E$  is a

- **tree edge**: belongs to  $T$
- **non-tree edge**: does not belong to  $T$

## Proposition

- $T$  is a forest and connected components of  $T$  are same as those of  $G$ .
- If  $(u, v) \in E$  is a non-tree edge then, in  $T$ , either  $u$  is an ancestor of  $v$  or  $v$  is an ancestor of  $u$ .

**Question:** Why are there no cross-edges?

## DFS with Visit Times

Keep track of when nodes are visited.

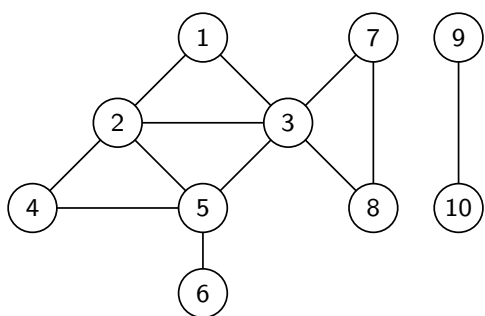
### DFS( $G$ )

```
Mark all nodes  $u \in V(G)$  as unvisited
 $T$  is set to  $\emptyset$ 
 $time = 0$ 
while there is an unvisited node  $u$  do
    DFS( $u$ )
Output  $T$ 
```

### DFS( $u$ )

```
Mark  $u$  as visited
pre( $u$ ) = ++time
for each edge  $(u, v)$  in Out( $u$ ) do
    if  $v$  is not marked then
        add edge  $(u, v)$  to  $T$ 
        DFS( $v$ )
post( $u$ ) = ++time
```

# Example



## pre and post numbers

Node  $u$  is **active** in time interval  $[\text{pre}(u), \text{post}(u)]$

### Proposition

For any two nodes  $u$  and  $v$ , the two intervals  $[\text{pre}(u), \text{post}(u)]$  and  $[\text{pre}(v), \text{post}(v)]$  are disjoint or one is contained in the other.

### Proof.

- Assume without loss of generality that  $\text{pre}(u) < \text{pre}(v)$ . Then  $v$  visited after  $u$ .
- If  $\text{DFS}(v)$  invoked before  $\text{DFS}(u)$  finished,  $\text{post}(u) > \text{post}(v)$ .
- If  $\text{DFS}(v)$  invoked after  $\text{DFS}(u)$  finished,  $\text{pre}(v) > \text{post}(u)$ .

□

pre and post numbers useful in several applications of **DFS**- soon!

# Part VII

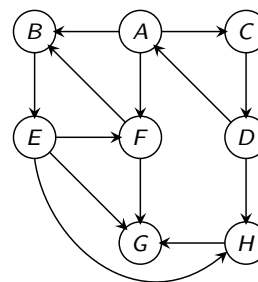
## Directed Graphs and Decomposition

## Directed Graphs

### Definition

A directed graph  $G = (V, E)$  consists of

- set of vertices/nodes  $V$  and
- a set of edges/arcs  $E \subseteq V \times V$ .



An edge is an *ordered* pair of vertices.  $(u, v)$  different from  $(v, u)$ .

# Examples of Directed Graphs

In many situations relationship between vertices is asymmetric:

- Road networks with one-way streets.
- Web-link graph: vertices are web-pages and there is an edge from page  $p$  to page  $p'$  if  $p$  has a link to  $p'$ . Web graphs used by Google with PageRank algorithm to rank pages.
- Dependency graphs in variety of applications: link from  $x$  to  $y$  if  $y$  depends on  $x$ . Make files for compiling programs.
- Program Analysis: functions/procedures are vertices and there is an edge from  $x$  to  $y$  if  $x$  calls  $y$ .

## Representation

Graph  $G = (V, E)$  with  $n$  vertices and  $m$  edges:

**Matrix**  $n \times n$  asymmetric matrix  $A$ .  $A[u, v] = 1$  if  $(u, v) \in E$  and  $A[u, v] = 0$  if  $(u, v) \notin E$ .  $A[u, v]$  is not same as  $A[v, u]$ .

**Adjacency Lists** for each node  $u$ ,  $Out(u)$  (also referred to as  $Adj(u)$ ) and  $In(u)$  store out-going edges and in-coming edges from  $u$ .

Default representation is adjacency lists.

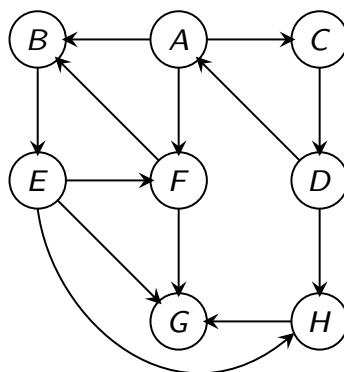
# Directed Connectivity

Given a graph  $G = (V, E)$ :

- A **(directed) path** is a sequence of *distinct* vertices  $v_1, v_2, \dots, v_k$  such that  $(v_i, v_{i+1}) \in E$  for  $1 \leq i \leq k - 1$ . The length of the path is  $k - 1$  and the path is from  $v_1$  to  $v_k$
- A **cycle** is a sequence of *distinct* vertices  $v_1, v_2, \dots, v_k$  such that  $(v_i, v_{i+1}) \in E$  for  $1 \leq i \leq k - 1$  and  $(v_k, v_1) \in E$ .
- A vertex  $u$  can **reach**  $v$  if there is a path from  $u$  to  $v$ .  
Alternatively  $v$  can be reached from  $u$
- Let  $\text{rch}(u)$  be the set of all vertices reachable from  $u$ .

## Connectivity contd

**Asymmetry:**  $A$  can reach  $B$  but  $B$  cannot reach  $A$



### Questions:

- Is there a notion of connected components?
- How do we understand connectivity in directed graphs?

# Connectivity and Strong Connected Components

## Definition

Given a directed graph  $G$ ,  $u$  is strongly connected to  $v$  if  $u$  can reach  $v$  and  $v$  can reach  $u$ . In other words  $v \in \text{rch}(u)$  and  $u \in \text{rch}(v)$ .

Define relation  $C$  where  $uCv$  if  $u$  is (strongly) connected to  $v$ .

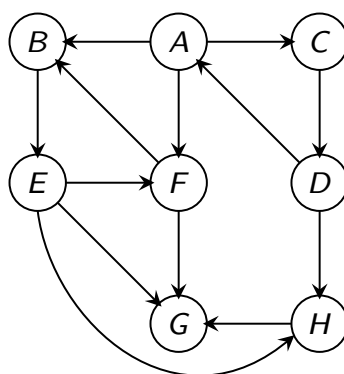
## Proposition

$C$  is an equivalence relation, that is reflexive, symmetric and transitive.

Equivalence classes of  $C$ : *strong connected components* of  $G$ .  
They *partition* the vertices of  $G$ .

$\text{SCC}(u)$ : strongly connected component containing  $u$ .

## Strongly Connected Components: Example





# Directed Graph Connectivity Problems

- Given  $G$  and nodes  $u$  and  $v$ , can  $u$  reach  $v$ ?
- Given  $G$  and  $u$ , compute  $\text{rch}(u)$ .
- Given  $G$  and  $u$ , compute all  $v$  that can reach  $u$ , that is all  $v$  such that  $u \in \text{rch}(v)$ .
- Find the strongly connected component containing node  $u$ , that is  $\text{SCC}(u)$ .
- Is  $G$  strongly connected (a single strong component)?
- Compute *all* strongly connected components of  $G$ .

First four problems can be solve in  $O(n + m)$  time by adapting **BFS/DFS** to directed graphs. The last one requires a clever **DFS** based algorithm.

## DFS in Directed Graphs

### DFS( $G$ )

```
Mark all nodes  $u$  as unvisited
 $T$  is set to  $\emptyset$ 
 $time = 0$ 
while there is an unvisited node  $u$  do
    DFS( $u$ )
```

Output  $T$

### DFS( $u$ )

```
Mark  $u$  as visited
 $pre(u) = ++time$ 
for each edge  $(u, v)$  in  $Out(u)$  do
    if  $v$  is not marked
        add edge  $(u, v)$  to  $T$ 
        DFS( $v$ )
 $post(u) = ++time$ 
```

Generalizing ideas from undirected graphs:

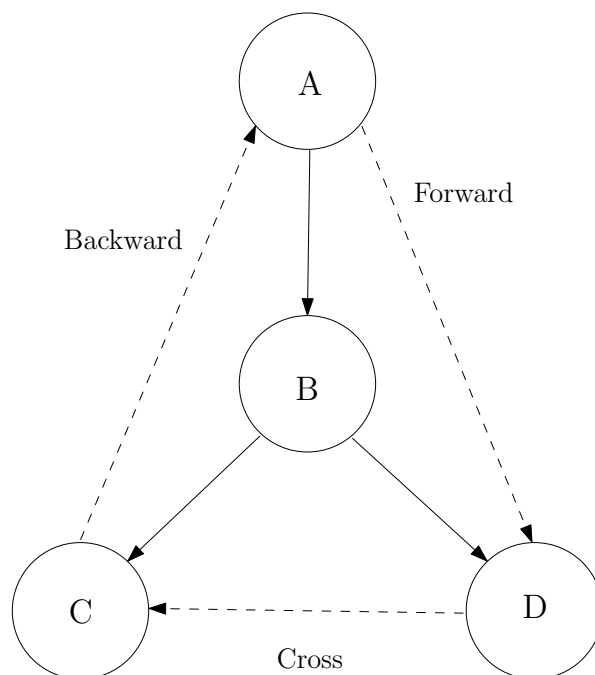
- **DFS**( $u$ ) outputs a directed out-tree  $T$  rooted at  $u$
- A vertex  $v$  is in  $T$  if and only if  $v \in \text{rch}(u)$
- For any two vertices  $x, y$  the intervals  $[\text{pre}(x), \text{post}(x)]$  and  $[\text{pre}(y), \text{post}(y)]$  are either disjoint or one is contained in the other.
- The running time of **DFS**( $u$ ) is  $O(k)$  where  $k = \sum_{v \in \text{rch}(u)} |\text{Adj}(v)|$  plus the time to initialize the Mark array.
- **DFS**( $G$ ) takes  $O(m + n)$  time. Edges in  $T$  form a disjoint collection of out-trees. Output of **DFS**( $G$ ) depends on the order in which vertices are considered.

## DFS Tree

Edges of  $G$  can be classified with respect to the **DFS** tree  $T$  as:

- **Tree edges** that belong to  $T$
- A **forward edge** is a non-tree edge  $(x, y)$  such that  $\text{pre}(x) < \text{pre}(y) < \text{post}(y) < \text{post}(x)$ .
- A **backward edge** is a non-tree edge  $(x, y)$  such that  $\text{pre}(y) < \text{pre}(x) < \text{post}(x) < \text{post}(y)$ .
- A **cross edge** is a non-tree edge  $(x, y)$  such that the intervals  $[\text{pre}(x), \text{post}(x)]$  and  $[\text{pre}(y), \text{post}(y)]$  are disjoint.

# Types of Edges



## Directed Graph Connectivity Problems

- Given  $G$  and nodes  $u$  and  $v$ , can  $u$  reach  $v$ ?
- Given  $G$  and  $u$ , compute  $\text{rch}(u)$ .
- Given  $G$  and  $u$ , compute all  $v$  that can reach  $u$ , that is all  $v$  such that  $u \in \text{rch}(v)$ .
- Find the strongly connected component containing node  $u$ , that is  $\text{SCC}(u)$ .
- Is  $G$  strongly connected (a single strong component)?
- Compute *all* strongly connected components of  $G$ .

# Algorithms via DFS- I

- Given  $G$  and nodes  $u$  and  $v$ , can  $u$  reach  $v$ ?
- Given  $G$  and  $u$ , compute  $\text{rch}(u)$ .

Use  $\text{DFS}(G, u)$  to compute  $\text{rch}(u)$  in  $O(n + m)$  time.

# Algorithms via DFS- II

- Given  $G$  and  $u$ , compute all  $v$  that can reach  $u$ , that is all  $v$  such that  $u \in \text{rch}(v)$ .

Given  $G = (V, E)$ ,  $G^{\text{rev}}$  is the graph with edge directions reversed  
 $G^{\text{rev}} = (V, E')$  where  $E' = \{(y, x) \mid (x, y) \in E\}$

Compute  $\text{rch}(u)$  in  $G^{\text{rev}}$ !

- **Correctness:** exercise
- **Running time:**  $O(n + m)$  to obtain  $G^{\text{rev}}$  from  $G$  and  $O(n + m)$  time to compute  $\text{rch}(u)$  via **DFS**. If both  $\text{Out}(v)$  and  $\text{In}(v)$  are available at each  $v$  then no need to explicitly compute  $G^{\text{rev}}$ . Can do it  $\text{DFS}(u)$  in  $G^{\text{rev}}$  implicitly.

## Algorithms via DFS- III

$$SC(\mathbf{G}, u) = \{v \mid u \text{ is strongly connected to } v\}$$

- Find the strongly connected component containing node  $u$ . That is, compute  $SCC(\mathbf{G}, u)$ .

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

Hence,  $SCC(\mathbf{G}, u)$  can be computed with two DFSes, one in  $\mathbf{G}$  and the other in  $\mathbf{G}^{rev}$ . Total  $O(n + m)$  time.

## Algorithms via DFS- IV

- Is  $\mathbf{G}$  strongly connected?

Pick arbitrary vertex  $u$ . Check if  $SC(\mathbf{G}, u) = V$ .

# Algorithms via DFS- V

- Find *all* strongly connected components of  $G$ .

For each vertex  $u \in V$  do  
  find  $SC(G, u)$

Running time:  $O(n(n + m))$ .

Can we do it in  $O(n + m)$  time?

## Reading Assignment and Homework 0

Chapters 1-3 from the textbook. At least half is prereq material.  
Proving algorithms correct - Jeff Erickson's notes (see link on website)