


**MA/CSSE 474**  
**Theory of Computation**

Computational Complexity



**Announcements**

- Don't forget the course evaluations on Banner Web.
- "HW 16" practice problems will be available today.
- Solutions available by Monday
- Final Exam Tuesday 1:00 O-201
  - Covers whole course
  - More emphasis on later stuff



## Undecidable Problems About CFLs

1. Given a CFL  $L$  and a string  $s$ , is  $s \in L$ ? (decidable)
2. Given a CFL  $L$ , is  $L = \emptyset$ ?
3. Given a CFL  $L$ , is  $L = \Sigma^*$ ?
4. Given CFLs  $L_1$  and  $L_2$ , is  $L_1 = L_2$ ?
5. Given CFLs  $L_1$  and  $L_2$ , is  $L_1 \subseteq L_2$ ?
6. Given a CFL  $L$ , is  $\neg L$  context-free?
7. Given a CFL  $L$ , is  $L$  regular?
8. Given two CFLs  $L_1$  and  $L_2$ , is  $L_1 \cap L_2 = \emptyset$ ?
9. Given a CFL  $L$ , is  $L$  inherently ambiguous?
10. Given PDAs  $M_1$  and  $M_2$ , is  $M_2$  a minimization of  $M_1$ ?
11. Given a CFG  $G$ , is  $G$  ambiguous?



## Time Complexity Classes



## Asymptotic Analysis Review

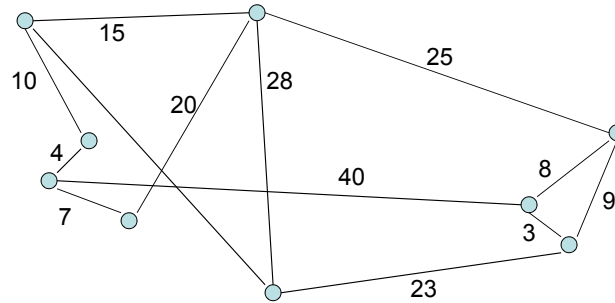
in case it's been a while ...



## Are All Decidable Languages Equal?

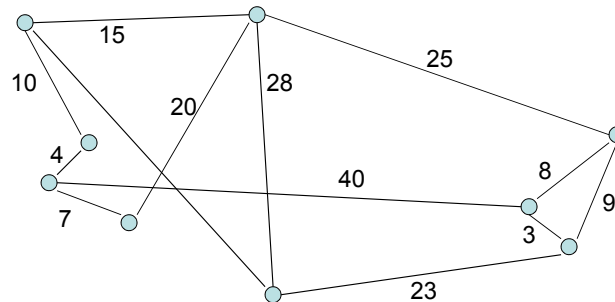
- $(ab)^*$
- $WW^R = \{ww^R : w \in \{a, b\}^*\}$
- $WW = \{ww : w \in \{a, b\}^*\}$
- $SAT = \{w : w \text{ is a wff in Boolean logic and } w \text{ is satisfiable}\}$

## The Traveling Salesman Problem



Given  $n$  cities and the distances between each pair of them, find the shortest tour that returns to its starting point and visits each other city exactly once along the way.

## The Traveling Salesman Problem



Given  $n$  cities:

Choose a first city

$n$

Choose a second

$n-1$

Choose a third

$\frac{n-2}{n!}$

...

$n!$

## The Traveling Salesman Problem

Can we do better than  $n!$

- First city doesn't matter.
- Order doesn't matter.

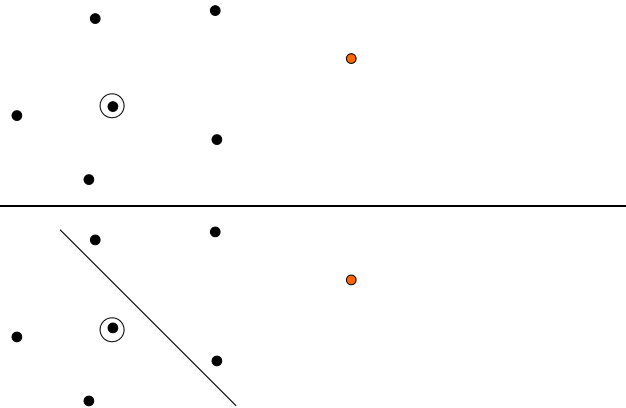
So we get  $(n-1)!/2$ .

## The Growth Rate of $n!$

<b>2</b>	2	<b>11</b>	479001600
<b>3</b>	6	<b>12</b>	6227020800
<b>4</b>	24	<b>13</b>	87178291200
<b>5</b>	120	<b>14</b>	1307674368000
<b>6</b>	720	<b>15</b>	20922789888000
<b>7</b>	5040	<b>16</b>	355687428096000
<b>8</b>	40320	<b>17</b>	6402373705728000
<b>9</b>	362880	<b>18</b>	121645100408832000
<b>10</b>	3628800	<b>19</b>	2432902008176640000
<b>11</b>	39916800	<b>36</b>	$3.6 \cdot 10^{41}$

## Tackling Hard Problems

1. Use a technique that is guaranteed to find an optimal solution and likely to do so quickly. Linear programming:



The [Concorde TSP Solver](#) found an optimal route that visits 24,978 cities in Sweden.

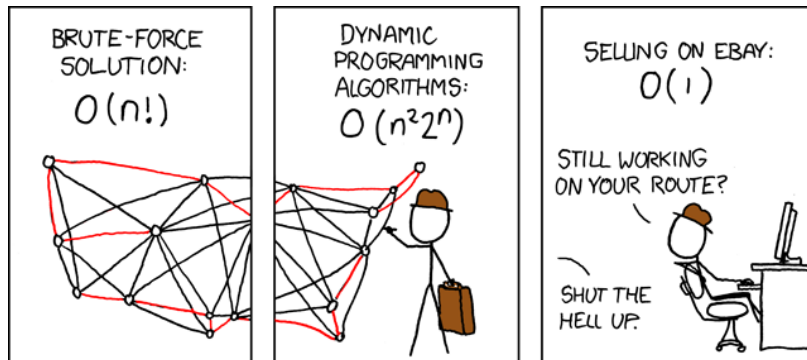
## Tackling Hard Problems

2. Use a technique that is guaranteed to run quickly and find a “good” solution.

- [The World Tour Problem](#)

Does it make sense to insist on true optimality if the description of the original problem was approximate?

## Modern TSP



From: <http://xkcd.com/399/>

## The Complexity Zoo

The attempt to characterize the decidable languages by their complexity:

[http://qwiki.stanford.edu/wiki/Complexity\\_Zoo](http://qwiki.stanford.edu/wiki/Complexity_Zoo)

See especially the [Petting Zoo](#) page.

## All Problems Are Decision Problems

### The Towers of Hanoi

Requires at least enough time to write the solution.

By restricting our attention to decision problems, the length of the answer is not a factor.

## Encoding Types Other Than Strings

The length of the encoding matters.

**Integers:** use any base other than 1.

11111111111	vs	1100
11111111111111111111111111111111	vs	11110

$$\log_a x = \log_a b \log_b x$$

- PRIMES = { $w$  :  $w$  is the binary encoding of a prime number}



## Encoding Types Other Than Strings

**Graphs:** use an adjacency matrix:

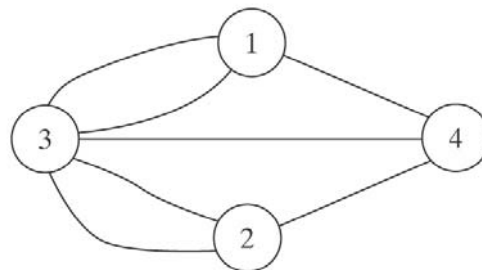
	1	2	3	4	5	6	7
1			•				
2				•			
3		•					
4					•		
5							
6							
7							

Or a list of edges:

101/1/11/11/10/10/100/100/101

## Graph Languages

- CONNECTED =  $\{ \langle G \rangle : G \text{ is an undirected graph and } G \text{ is connected} \}$ .
- HAMILTONIANCIRCUIT =  $\{ \langle G \rangle : G \text{ is an undirected graph that contains a } \mathbf{Hamiltonian circuit} \}$ .





## Characterizing Optimization Problems as Languages

- TSP-DECIDE =  $\{ \langle G, cost \rangle : \langle G \rangle \text{ encodes an undirected graph with a positive distance attached to each of its edges and } G \text{ contains a Hamiltonian circuit whose total cost is less than } \langle cost \rangle \}$ .



## Choosing A Model of Computation

We'll use Turing machines:

- Tape alphabet size?
- How many tapes?
- Deterministic vs. nondeterministic?

## Measuring Time and Space Requirements

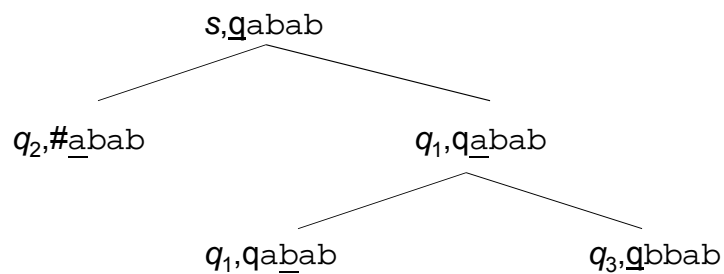
$timereq(M)$  is a function of  $n$ :

- If  $M$  is a **deterministic** Turing machine that halts on all inputs, then:

$timereq(M) = f(n)$  = the maximum number of steps that  $M$  executes on any input of length  $n$ .

## Measuring Time and Space Requirements

- If  $M$  is a **nondeterministic** Turing machine all of whose computational paths halt on all inputs, then:



$timereq(M) = f(n)$  = the number of steps on the longest path that  $M$  executes on any input of length  $n$ .

## Measuring Time and Space Requirements

$spacereq(M)$  is a function of  $n$ :

- If  $M$  is a **deterministic** Turing machine that halts on all inputs, then:

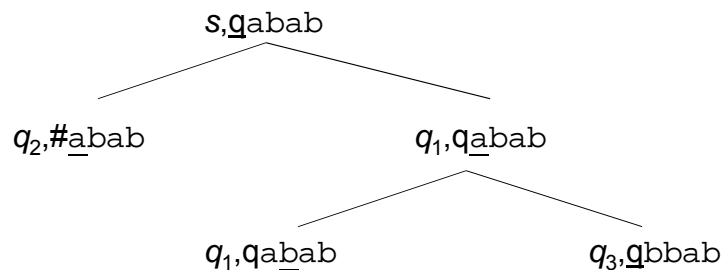
$spacereq(M) = f(n)$  = the maximum number of tape squares that  $M$  reads on any input of length  $n$ .

- If  $M$  is a **nondeterministic** Turing machine all of whose computational paths halt on all inputs, then:

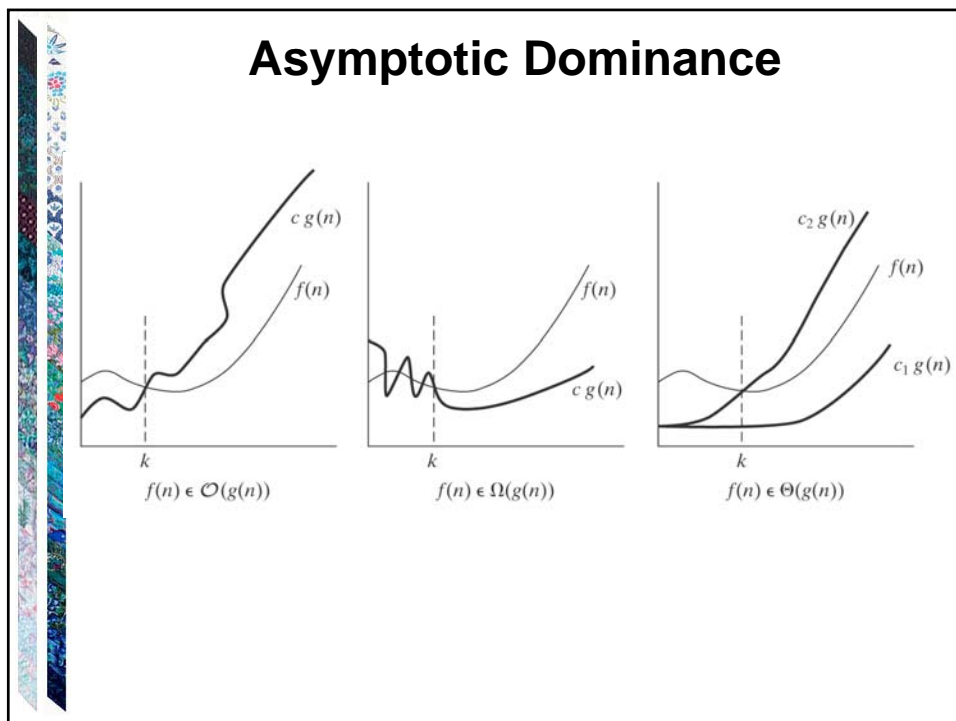
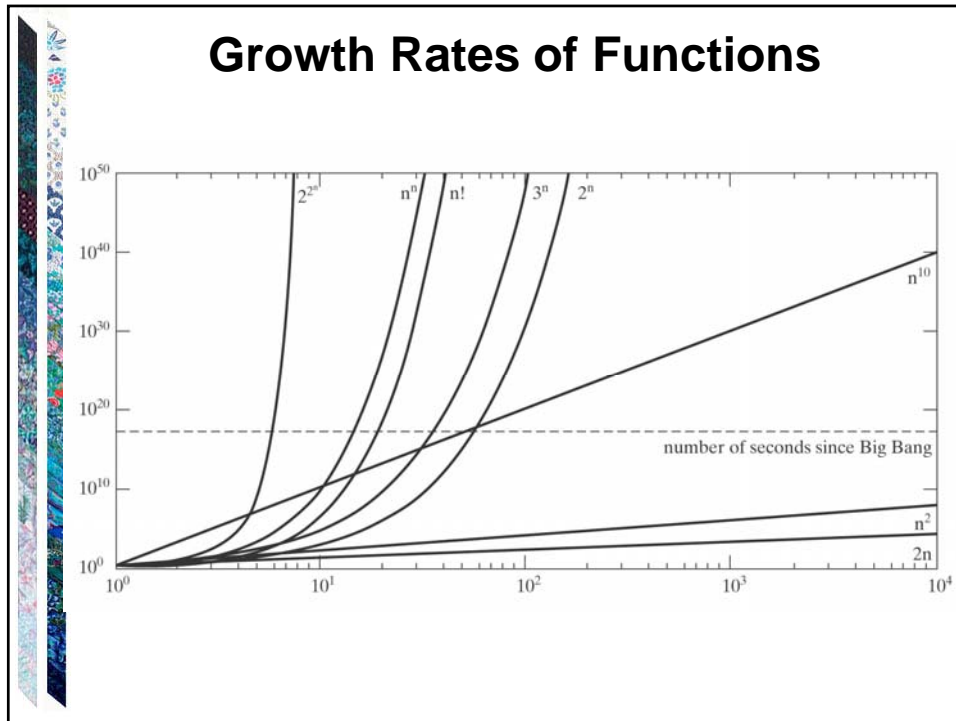
$spacereq(M) = f(n)$  = the maximum number of tape squares that  $M$  reads on any path that it executes on any input of length  $n$ .

## Measuring Time and Space Requirements

- If  $M$  is a **nondeterministic** Turing machine all of whose computational paths halt on all inputs, then:



$timereq(M) = f(n)$  = the number of steps on the longest path that  $M$  executes on any input of length  $n$ .



## Asymptotic Dominance - $\mathcal{O}$

$f(n) \in \mathcal{O}(g(n))$  iff there exists a positive integer  $k$  and a positive constant  $c$  such that:

$$\forall n \geq k (f(n) \leq c g(n)).$$

Alternatively, if the limit exists:

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} < \infty$$

Or,  $g$  grows at least as fast as  $f$  does.

## Summarizing $\mathcal{O}$

$$\mathcal{O}(c) \subseteq \mathcal{O}(\log_a n) \subseteq \mathcal{O}(n^b) \subseteq \mathcal{O}(d^n) \subseteq \mathcal{O}(n!) \subseteq \mathcal{O}(n^n)$$

## $\mathcal{O}$ (little oh)

**Asymptotic strong upper bound:**  $f(n) \in \mathcal{O}(g(n))$  iff, for every positive  $c$ , there exists a positive integer  $k$  such that:

$$\forall n \geq k (f(n) < c g(n)).$$

Alternatively, if the limit exists:

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0$$

In this case, we'll say that  $f$  is "little-oh" of  $g$  or that  $g$  grows strictly faster than  $f$  does.

## $\Omega$

• **Asymptotic lower bound:**  $f(n) \in \Omega(g(n))$  iff there exists a positive integer  $k$  and a positive constant  $c$  such that:

$$\forall n \geq k (f(n) \geq c g(n)).$$

In other words, ignoring some number of small cases (all those of size less than  $k$ ), and ignoring some constant factor  $c$ ,  $f(n)$  is bounded from below by  $g(n)$ .

Alternatively, if the limit exists:

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} > 0$$

In this case, we'll say that  $f$  is "big-Omega" of  $g$  or that  $g$  grows no faster than  $f$ .

## Ω

- **Asymptotic strong lower bound:**  $f(n) \in \omega(g(n))$   
iff, for every positive  $c$ , there exists a positive integer  $k$  such that:

$$\forall n \geq k (f(n) > c g(n)).$$

Alternatively, if the required limit exists:

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = \infty$$

In this case, we'll say that  $f$  is "little-omega" of  $g$  or that  $g$  grows strictly slower than  $f$  does.

## Θ

$f(n) \in \Theta(g(n))$  iff there exists a positive integer  $k$  and positive constants  $c_1$ , and  $c_2$  such that:

$$\forall n \geq k (c_1 g(n) \leq f(n) \leq c_2 g(n))$$

Or:

$$f(n) \in \Theta(g(n)) \text{ iff:}$$

$$f(n) \in \mathcal{O}(g(n)), \text{ and}$$

$$g(n) \in \mathcal{O}(f(n)).$$

Or:

$$f(n) \in \Theta(g(n)) \text{ iff:}$$

$$f(n) \in \mathcal{O}(g(n)), \text{ and}$$

$$f(n) \in \Omega(g(n)).$$

Is  $n^3 \in \Theta(n^3)$ ?

Is  $n^3 \in \Theta(n^4)$ ?

Is  $n^3 \in \Theta(n^5)$ ?



## Algorithmic Gaps

We'd like to show:

1. Upper bound: There exists an algorithm that decides  $L$  and that has complexity  $C_1$ .
2. Lower bound: Any algorithm that decides  $L$  must have complexity at least  $C_2$ .
3.  $C_1 = C_2$ .

If  $C_1 = C_2$ , we are done. Often, we're not done.

## Algorithmic Gaps

Example: TSP

- Upper bound:  $timereq \in \mathcal{O}(2^{(n^k)})$ .
- Don't have a lower bound that says polynomial isn't possible.

We group languages by what we know. And then we ask:  
"Is class  $CL_1$  equal to class  $CL_2$ ?"

## A Simple Example of Polynomial Speedup

Given a list of  $n$  numbers, find the minimum and the maximum elements in the list. Or, as a language recognition problem:

$L = \langle \text{list of numbers}, \text{number}_1, \text{number}_2 \rangle$ :  
 $\text{number}_1$  is the minimum element of the list and  
 $\text{number}_2$  is the maximum element}.

$(23, 45, 73, 12, 45, 197; 12; 197) \in L$ .

## A Simple Example of Polynomial Speedup

The straightforward approach:

$\text{simplecompare}(\text{list: list of numbers}) =$

$\text{max} = \text{list}[1].$

$\text{min} = \text{list}[1].$

For  $i = 2$  to  $\text{length}(\text{list})$  do:

  If  $\text{list}[i] < \text{min}$  then  $\text{min} = \text{list}[i].$

  If  $\text{list}[i] > \text{max}$  then  $\text{max} = \text{list}[i].$

Requires  $2(n-1)$  comparisons. So  $\text{simplecompare}$  is  $\mathcal{O}(n)$ .

But we can solve this problem in  $(3/2)(n-1)$  comparisons.

How?

## A Simple Example of Polynomial Speedup

*efficientcompare(list: list of numbers) =*

*max = list[1].*

*min = list[1].*

For  $i = 3$  to  $\text{length}(\text{list})$  by 2 do:

  If  $\text{list}[i] < \text{list}[i-1]$  then:

    If  $\text{list}[i] < \text{min}$  then  $\text{min} = \text{list}[i]$ .

    If  $\text{list}[i-1] > \text{max}$  then  $\text{max} = \text{list}[i-1]$ .

  Else:

    If  $\text{list}[i-1] < \text{min}$  then  $\text{min} = \text{list}[i-1]$ .

    If  $\text{list}[i] > \text{max}$  then  $\text{max} = \text{list}[i]$ .

If  $\text{length}(\text{list})$  is even then check the last element.

Requires  $3/2(n-1)$  comparisons.

## String Search

*t:*     a b c a b a b c a b d

*p:*     a b c d

       a b c d

          a b c d

                  . . .

## String Search

```

simple-string-search( $t, p$ : strings) =
   $i = 0.$ 
   $j = 0.$ 
  While  $i \leq |t| - |p|$  do:
    While  $j < |p|$  do:
      If  $t[i+j] = p[j]$  then  $j = j + 1.$ 
      Else exit this loop.
      If  $j = |p|$  then halt and accept.
      Else:
         $i = i + 1.$ 
         $j = 0.$ 
  Halt and reject.

```

Let  $n$  be  $|t|$  and let  $m$  be  $|p|$ . In the worst case (in which it doesn't find an early match), *simple-string-search* will go through its outer loop almost  $n$  times and, for each of those iterations, it will go through its inner loop  $m$  times.

So  $\text{timereq}(\text{simple-string-search}) \in \mathcal{O}(nm)$ .

K-M-P algorithm is  
 $\mathcal{O}(n+m)$

## Replacing an Exponential Algorithm with a Polynomial One

- Context-free parsing can be done in  $\mathcal{O}(n^3)$  time instead of  $\mathcal{O}(2^n)$  time. (CYK algorithm)
- Finding the greatest common divisor of two integers can be done in  $\mathcal{O}(\log_2(\max(n, m)))$  time instead of exponential time.

## The Language Class P

$L \in P$  iff

- there exists some deterministic Turing machine  $M$  that decides  $L$ , and
- $\text{timereq}(M) \in \mathcal{O}(n^k)$  for some  $k$ .

We'll say that  $L$  is **tractable** iff it is in P.

## Closure under Complement

**Theorem:** The class P is closed under complement.

**Proof:** If  $M$  accepts  $L$  in polynomial time, swap accepting and non accepting states to accept  $\neg L$  in polynomial time.

## Defining Complement

- $\text{CONNECTED} = \{ \langle G \rangle : G \text{ is an undirected graph and } G \text{ is connected} \}$  is in P.
- $\text{NOTCONNECTED} = \{ \langle G \rangle : G \text{ is an undirected graph and } G \text{ is not connected} \}$ .
- $\neg\text{CONNECTED} = \text{NOTCONNECTED} \cup \{ \text{strings that are not syntactically legal descriptions of undirected graphs} \}$ .

$\neg\text{CONNECTED}$  is in P by the closure theorem. What about NOTCONNECTED?

If we can check for legal syntax in polynomial time, then we can consider the universe of strings whose syntax is legal. Then we can conclude that NOTCONNECTED is in P if CONNECTED is.

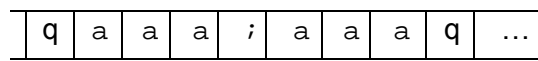
## Languages That Are in P

- Every regular language.
- Every context-free language since there exist context-free parsing algorithms that run in  $\mathcal{O}(n^3)$  time.
- Others:
  - $A^n B^n C^n$
  - Nim

## To Show That a Language Is In P

- Describe a one-tape, deterministic Turing machine.
- It may use multiple tapes. Price:
- State an algorithm that runs on a conventional computer.  
Price:

How long does it take to compare two strings?



Bottom line: If ignoring polynomial factors, then just describe a deterministic algorithm.

## Regular Languages

**Theorem:** Every regular language can be decided in linear time. So every regular language is in P.

**Proof:** If  $L$  is regular, there exists some DFSM  $M$  that decides it. Construct a deterministic TM  $M'$  that simulates  $M$ , moving its read/write head one square to the right at each step. When  $M'$  reads a  $q$ , it halts. If it is in an accepting state, it accepts; otherwise it rejects.

On any input of length  $n$ ,  $M'$  will execute  $n + 2$  steps.

So  $\text{timereq}(M') \in \mathcal{O}(n)$ .

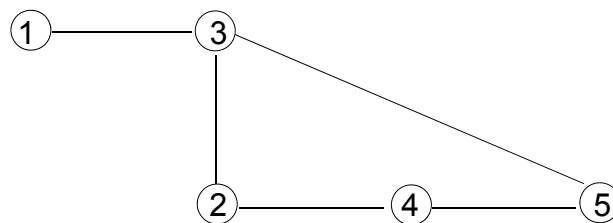
## Context-Free Languages

**Theorem:** Every context-free language can be decided in  $\mathcal{O}(n^{18})$  time. So every context-free language is in P.

**Proof:** The Cocke-Kasami-Younger (CKY) algorithm can parse any context-free language in time that is  $\mathcal{O}(n^3)$  if we count operations on a conventional computer. That algorithm can be simulated on a standard, one-tape Turing machine in  $\mathcal{O}(n^{18})$  steps.

## Graph Languages

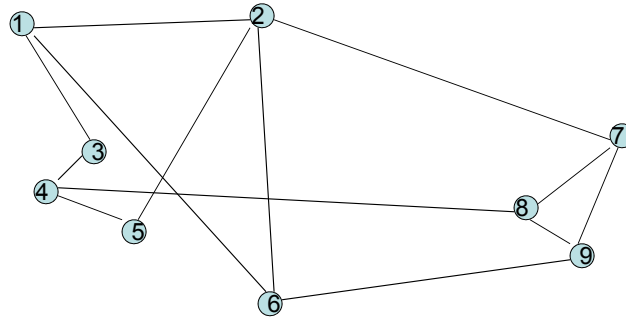
Represent a graph  $G = (V, E)$  as a list of edges:



101/1/11/11/10/10/100/100/101/11/101



## Graph Languages



CONNECTED =  
 $\{ \langle G \rangle : G \text{ is an undirected graph and } G \text{ is connected} \}$ .

Is CONNECTED in P?

## CONNECTED is in P

$connected(\langle G = (V, E) \rangle =$

1. Set all vertices to be unmarked.
2. Mark vertex 1.
3. Initialize  $L$  to  $\{1\}$ .
4. Initialize *marked-vertices-counter* to 1.
5. Until  $L$  is empty do:
  - 5.1. Remove the first element from  $L$ . Call it *current-vertex*.
  - 5.2. For each edge  $e$  that has *current-vertex* as an endpoint do:
 

Call the other endpoint of  $e$  *next-vertex*. If *next-vertex* is not already marked then do:

Mark *next-vertex*.

Add *next-vertex* to  $L$ .

Increment *marked-vertices-counter* by 1.
6. If *marked-vertices-counter* =  $|V|$  accept. Else reject.

## Analyzing *connected*

- Step 1 takes time that is  $\mathcal{O}(|V|)$ .
- Steps 2, 3, and 4 each take constant time.
- The loop of step 5 can be executed at most  $|V|$  times.
  - Step 5.1 takes constant time.
  - Step 5.2 can be executed at most  $|E|$  times. Each time, it requires at most  $\mathcal{O}(|V|)$  time.
- Step 6 takes constant time.

So *timereq(connected)* is:

$$|V| \cdot \mathcal{O}(|E|) \cdot \mathcal{O}(|V|) = \mathcal{O}(|V|^2|E|).$$

But  $|E| \leq |V|^2$ . So *timereq(connected)* is:

$$\mathcal{O}(|V|^4).$$

## Primality Testing

RELATIVELY-PRIME =

$\{ \langle n, m \rangle : n \text{ and } m \text{ are integers that are relatively prime} \}$ .

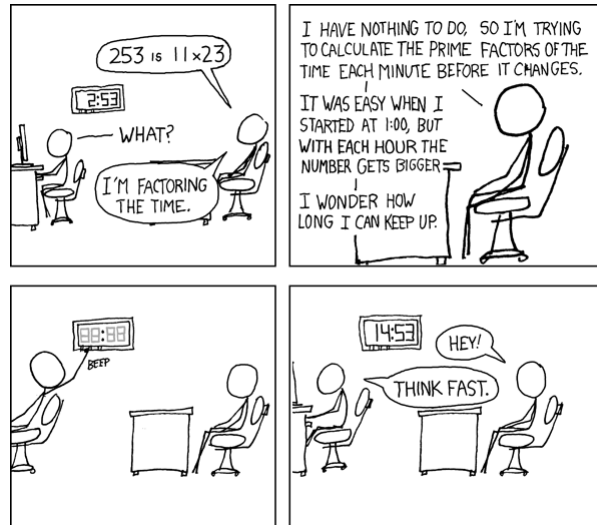
PRIMES =

$\{ w : w \text{ is the binary encoding of a prime number} \}$

COMPOSITES =

$\{ w : w \text{ is the binary encoding of a nonprime number} \}$

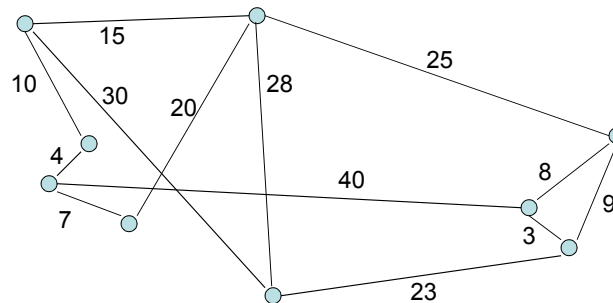
## But Finding Factors Remains Hard



<http://xkcd.com/247/>

## Returning to TSP

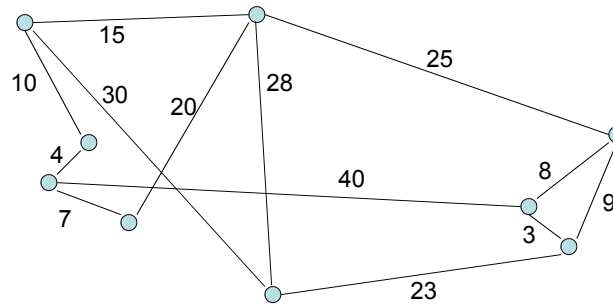
TSP-DECIDE =  $\{ \langle G, cost \rangle : \langle G \rangle \text{ encodes an undirected graph with a positive distance attached to each of its edges and } G \text{ contains a Hamiltonian circuit whose total cost is less than } \langle cost \rangle \}$ .



An NDTM to decide TSP-DECIDE:

## Returning to TSP

An NDTM to decide TSP-DECIDE:



1. For  $i = 1$  to  $|V|$  do:
  - Choose a vertex that hasn't yet been chosen.
2. Check that the path defined by the chosen sequence of vertices is a Hamiltonian circuit through  $G$  with distance less than cost.

## TSP and Other Problems Like It

TSP-DECIDE, and other problems like it, share three properties:

1. The problem can be solved by searching through a space of partial solutions (such as routes). The size of this space grows exponentially with the size of the problem.
2. No better (i.e., not based on search) technique for finding an exact solution is known.
3. But, if a proposed solution were suddenly to appear, it could be checked for correctness very efficiently.

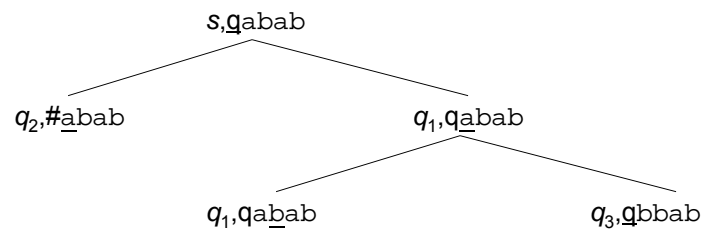
## The Language Class NP

Nondeterministic deciding:

$L \in \text{NP}$  iff:

- there is some NDTM  $M$  that decides  $L$ , and
- $\text{timereq}(M) \in \mathcal{O}(n^k)$  for some  $k$ .

NDTM deciders:



## TSP Again

TSP-DECIDE =  $\{ \langle G, cost \rangle : \langle G \rangle$  encodes an undirected graph with a positive distance attached to each of its edges and  $G$  contains a Hamiltonian circuit whose total cost is less than  $\langle cost \rangle$   $\}$ .

Suppose some Oracle presented a candidate path  $c$ :

$\langle G, cost, v_1, v_7, v_4, v_3, v_8, v_5, v_2, v_6, v_1 \rangle$

How long would it take to verify that  $c$  proves that:

$\langle G, cost \rangle$  is in TSP-DECIDE?

## Deterministic Verifying

A Turing machine  $V$  is a **verifier** for a language  $L$  iff:

$w \in L$  iff  $\exists c \langle w, c \rangle \in L(V)$ .

We'll call  $c$  a **certificate**.

## Deterministic Verifying

An alternative definition for the class NP:

$L \in \text{NP}$  iff there exists a deterministic TM  $V$  such that:

- $V$  is a verifier for  $L$ , and
- $\text{timereq}(V) \in \mathcal{O}(n^k)$  for some  $k$ .

## ND Deciding and D Verifying

**Theorem:** These two definitions are equivalent:

- (1)  $L \in \text{NP}$  iff there exists a nondeterministic, polynomial-time TM that decides it.
- (2)  $L \in \text{NP}$  iff there exists a deterministic, polynomial-time verifier for it.

**Proof:** WE skip it

## Proving That a Language is in NP

- Exhibit an NDTM to decide it.
- Exhibit a DTM to verify it.

## Example

- $SAT = \{w : w \text{ is a Boolean wff and } w \text{ is satisfiable}\}$  is in NP.

$$F_1 = P \wedge Q \wedge \neg R ?$$

$$F_2 = P \wedge Q \wedge R ?$$

$$F_3 = P \wedge \neg P ?$$

$$F_4 = P \wedge (Q \vee \neg R) \wedge \neg Q ?$$

$SAT-decide(F_4) =$

$SAT-verify(\langle F_4, (P = True, Q = False, R = False) \rangle) =$

## 3-SAT

- A **literal** is either a variable or a variable preceded by a single negation symbol.
- A **clause** is either a single literal or the disjunction of two or more literals.
- A wff is in **conjunctive normal form** (or CNF) iff it is either a single clause or the conjunction of two or more clauses.
- A wff is in **3-conjunctive normal form** (or 3-CNF) iff it is in conjunctive normal form and each clause contains exactly three literals.



## 3-SAT

	3-CNF	CNF
$(P \vee \neg Q \vee R)$	•	•
$(P \vee \neg Q \vee R) \wedge (\neg P \vee Q \vee \neg R)$	•	•
$P$		•
$(P \vee \neg Q \vee R \vee S) \wedge (\neg P \vee \neg R)$		•
$P \rightarrow Q$		
$(P \wedge \neg Q \wedge R \wedge S) \vee (\neg P \wedge \neg R)$		
$\neg(P \vee Q \vee R)$		

Every wff can be converted to an equivalent wff in CNF.

- 3-SAT = {  $w$  :  $w$  is a wff in Boolean logic,  $w$  is in 3-conjunctive normal form, and  $w$  is satisfiable}.

Is 3-SAT in NP?

## The Relationship Between P and NP

Is  $P = NP$ ?

Here are some things we know:

$P \subseteq NP \subseteq PSPACE \subseteq EXPTIME$

$P \neq EXPTIME$

[The Millenium Prize](#)

## Using Reduction in Complexity Proofs

A **mapping reduction**  $R$  from  $L_1$  to  $L_2$  is a

- Turing machine that
- implements some computable function  $f$  with the property that:

$$\forall x (x \in L_1 \leftrightarrow f(x) \in L_2).$$

If  $L_1 \leq L_2$  and  $M$  decides  $L_2$ , then:

$$C(x) = M(R(x)) \text{ will decide } L_1.$$

## Using Reduction in Complexity Proofs

If  $R$  is deterministic polynomial then:

$$L_1 \leq_P L_2.$$

And, whenever such an  $R$  exists:

- $L_1$  must be in P if  $L_2$  is: if  $L_2$  is in P then there exists some deterministic, polynomial-time Turing machine  $M$  that decides it. So  $M(R(x))$  is also a deterministic, polynomial-time Turing machine and it decides  $L_1$ .
- $L_1$  must be in NP if  $L_2$  is: if  $L_2$  is in NP then there exists some nondeterministic, polynomial-time Turing machine  $M$  that decides it. So  $M(R(x))$  is also a nondeterministic, polynomial-time Turing machine and it decides  $L_1$ .

## Why Use Reduction?

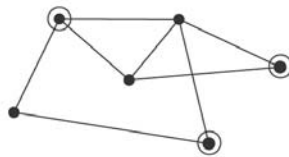
Given  $L_1 \leq_P L_2$ , we can use reduction to:

- Prove that  $L_1$  is in P or in NP because we **already know** that  $L_2$  is.
- Prove that  $L_1$  would be in P or in NP if we **could somehow show** that  $L_2$  is. When we do this, we cluster languages of similar complexity (even if we're not yet sure what that complexity is). In other words,  $L_1$  is no harder than  $L_2$  is.

## INDEPENDENT-SET

- INDEPENDENT-SET =  $\{ \langle G, k \rangle : G \text{ is an undirected graph and } G \text{ contains an independent set of at least } k \text{ vertices} \}$ .

An **independent set** is a set of vertices no two of which are adjacent (i.e., connected by a single edge). So, in the following graph, the circled vertices form an independent set:



In a scheduling program the vertices represent tasks and are connected by an edge if their corresponding tasks conflict. We can find the largest number of tasks that can be scheduled at the same time by finding the largest independent set in the task graph.

## 3-SAT and INDEPENDENT-SET

3-SAT  $\leq_p$  INDEPENDENT-SET.

Strings in 3-SAT describe formulas that contain literals and clauses.

$$(P \vee Q \vee \neg R) \wedge (R \vee \neg S \vee Q)$$

Strings in INDEPENDENT-SET describe graphs that contain vertices and edges.

101/1/11/11/10/10/100/100/101/11/101

## Gadgets

A **gadget** is a structure in the target language that mimics the role of a corresponding structure in the source language.

Example: 3-SAT  $\leq_p$  INDEPENDENT-SET.

$$(P \vee Q \vee \neg R) \wedge (R \vee \neg S \vee Q)$$



(approximately)

101/1/11/11/10/10/100/100/101/11/101

So we need:

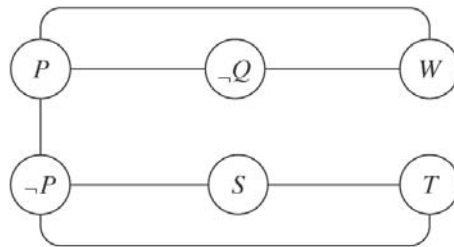
- a gadget that looks like a graph but that mimics a literal, and
- a gadget that looks like a graph but that mimics a clause.

## 3-SAT $\leq_p$ INDEPENDENT-SET

$R(\langle f \rangle)$ : Boolean formula with  $k$  clauses =

1. Build a graph  $G$  by doing the following:
  - 1.1. Create one vertex for each instance of each literal in  $f$ .
  - 1.2. Create an edge between each pair of vertices for symbols in the same clause.
  - 1.3. Create an edge between each pair of vertices for complementary literals.
2. Return  $\langle G, k \rangle$ .

$(P \vee \neg Q \vee W) \wedge (\neg P \vee S \vee T)$ :

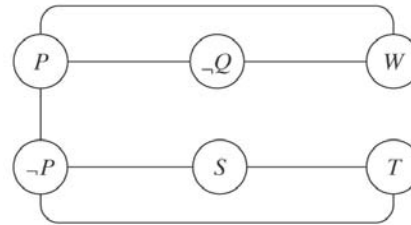


## $R$ is Correct

Show:  $f \in 3\text{-SAT}$  iff  $R(\langle f \rangle) \in \text{INDEPENDENT-SET}$   
by showing:

- $f \in 3\text{-SAT} \rightarrow R(\langle f \rangle) \in \text{INDEPENDENT-SET}$
- $R(\langle f \rangle) \in \text{INDEPENDENT-SET} \rightarrow f \in 3\text{-SAT}$

## One Direction



$f \in 3\text{-SAT} \rightarrow R(\langle f \rangle) \in \text{INDEPENDENT-SET}$ :

There is a satisfying assignment  $A$  to the symbols in  $f$ .

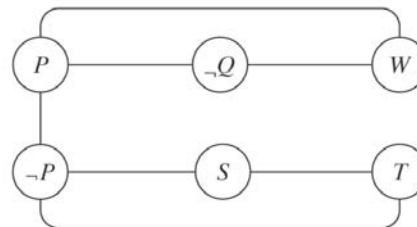
So  $G$  contains an independent set  $S$  of size  $k$ , built by:

1. From each clause gadget choose one literal that is made positive by  $A$ .
2. Add the vertex corresponding to that literal to  $S$ .

$S$  will contain exactly  $k$  vertices and is an independent set:

- No two vertices come from the same clause so step 1.2 could not have created an edge between them.
- No two vertices correspond to complementary literals so step 1.3 could not have created an edge between them.

## The Other Direction



- $R(\langle f \rangle) \in \text{INDEPENDENT-SET}$ .
- So the graph  $G$  that  $R$  builds contains an independent set  $S$  of size  $k$ .
- We prove that there is some satisfying assignment  $A$  for  $f$ .

No two vertices in  $S$  come from the same clause gadget. Since  $S$  contains at least  $k$  vertices, no two are from the same clause, and  $f$  contains  $k$  clauses,  $S$  must contain one vertex from each clause.

Build  $A$  as follows:

1. Assign *True* to each literal that corresponds to a vertex in  $S$ .
2. Assign arbitrary values to all other literals.

Since each clause will contain at least one literal whose value is *True*, the value of  $f$  will be *True*.

## NP-Completeness

A language  $L$  might have these properties:

1.  $L$  is in NP.
  2. Every language in NP is deterministic, polynomial-time reducible to  $L$ .
- $L$  is **NP-hard** iff it possesses property 2.  
An NP-hard language is at least as hard as any other language in NP.
  - $L$  is **NP-complete** iff it possesses *both* property 1 and property 2.  
All NP-complete languages can be viewed as being equivalently hard.

## NP-Hard vs. NP-Complete

An example: puzzles vs. games (Appendix N).

To use this theory to analyze a game like chess, we must generalize it so that we can talk about solution time as a function of problem size:

CHESSESS =  $\{ \langle b \rangle : b \text{ is a configuration of an } n \times n \text{ chess board and there is a guaranteed win for the current player} \}$ .

## Sudoku

- **SUDOKU** =  $\{ \langle b \rangle : b \text{ is a configuration of an } n \times n \text{ grid and } b \text{ has a solution under the rules of Sudoku} \}$ .

	5			9	4	2	1	
4							8	
	3		7					
				2				4
2			4		6			8
6				3				
					8		6	
	7							3
	6	8	9	5			4	

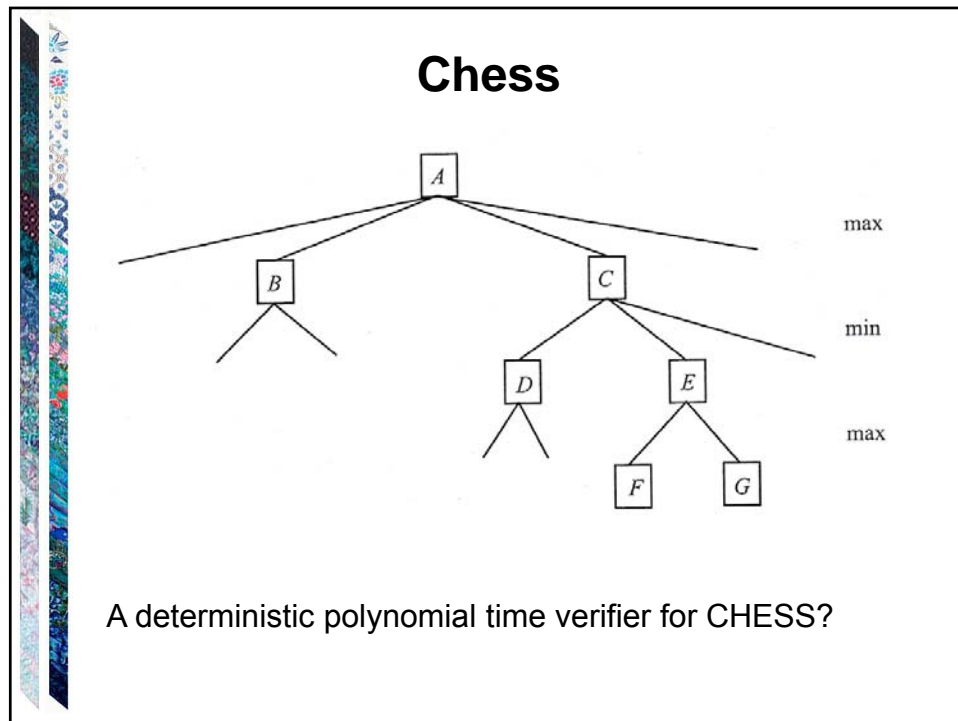
## Sudoku

- **SUDOKU** =  $\{ \langle b \rangle : b \text{ is a configuration of an } n \times n \text{ grid and } b \text{ has a solution under the rules of Sudoku} \}$ .

A deterministic, polynomial time verifier for SUDOKU, on input:

$\langle b, (1,1,1), (1,2,5), (1,3,4), \dots \rangle$





## NP-Hard vs. NP-Complete

SUDOKU =  $\{ \langle b \rangle : b \text{ is a configuration of an } n \times n \text{ grid and } b \text{ has a solution under the rules of Sudoku} \}$ .

*NP-complete.*

CHESS =  $\{ \langle b \rangle : b \text{ is a configuration of an } n \times n \text{ chess board and there is a guaranteed win for the current player} \}$ .

*NP-hard, not thought to be in NP.*

*If fixed number of pieces: PSPACE-complete.*

*If variable number of pieces: EXPTIME-complete.*

## Showing that $L$ is NP-Complete

How about: Take a list of known NP languages and crank out the reductions?

$$NPL_1 \geq L$$

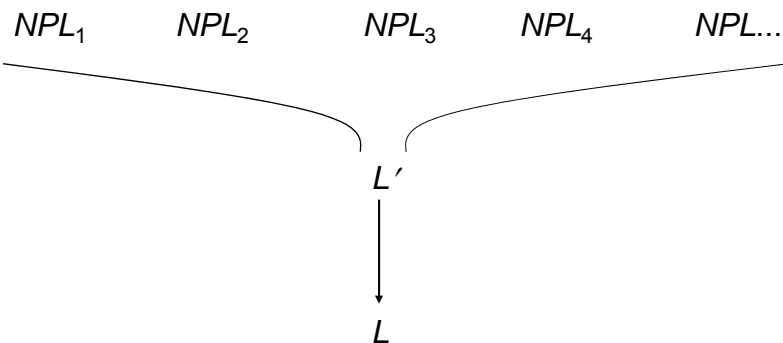
$$NPL_2 \geq L$$

$$NPL_3 \geq L$$

...

## Showing that $L$ is NP-Complete

Suppose we had one NP-complete language  $L'$ :



## Finding an $L'$

- The key property that every NP language has is that it can be decided by a polynomial time NDTM.
- So we need a language in which we can describe computations of NDTMs.

## The Cook-Levin Theorem

**Define:**  $SAT = \{w : w \text{ is a wff in Boolean logic and } w \text{ is satisfiable}\}$

**Theorem:** SAT is NP-complete.

**Proof:**

- SAT is in NP.
- SAT is NP-hard.

## NP-Complete Languages

- SUBSET-SUM =  $\{ \langle S, k \rangle : S \text{ is a multiset of integers, } k \text{ is an integer, and there exists some subset of } S \text{ whose elements sum to } k \}$ .
- SET-PARTITION =  $\{ \langle S \rangle : S \text{ is a multiset of objects each of which has an associated cost and there exists a way to divide } S \text{ into two subsets, } A \text{ and } S - A, \text{ such that the sum of the costs of the elements in } A \text{ equals the sum of the costs of the elements in } S - A \}$ .
- KNAPSACK =  $\{ \langle S, v, c \rangle : S \text{ is a set of objects each of which has an associated cost and an associated value, } v \text{ and } c \text{ are integers, and there exists some way of choosing elements of } S \text{ (duplicates allowed) such that the total cost of the chosen objects is at most } c \text{ and their total value is at least } v \}$ .

## NP-Complete Languages

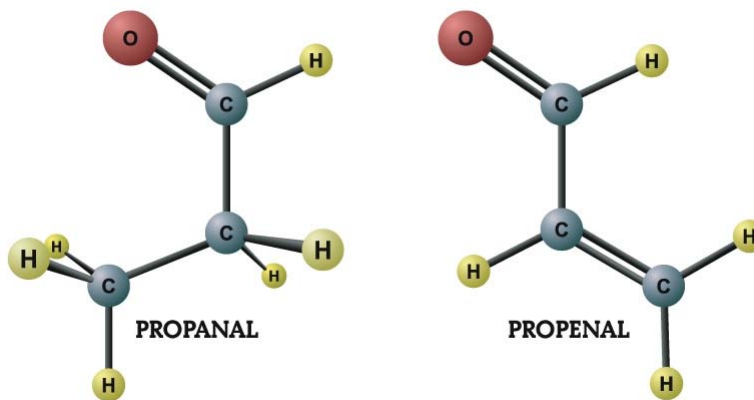
- TSP-DECIDE.
- HAMILTONIAN-PATH =  $\{ \langle G \rangle : G \text{ is an undirected graph and } G \text{ contains a Hamiltonian path} \}$ .
- HAMILTONIAN-CIRCUIT =  $\{ \langle G \rangle : G \text{ is an undirected graph and } G \text{ contains a Hamiltonian circuit} \}$ .
- CLIQUE =  $\{ \langle G, k \rangle : G \text{ is an undirected graph with vertices } V \text{ and edges } E, k \text{ is an integer, } 1 \leq k \leq |V|, \text{ and } G \text{ contains a } k\text{-clique} \}$ .
- INDEPENDENT-SET =  $\{ \langle G, k \rangle : G \text{ is an undirected graph and } G \text{ contains an independent set of at least } k \text{ vertices} \}$ .

## NP-Complete Languages

- SUBGRAPH-ISOMORPHISM =  $\{ \langle G_1, G_2 \rangle : G_1 \text{ is isomorphic to some subgraph of } G_2 \}$ .

Two graphs  $G$  and  $H$  are **isomorphic** to each other iff there exists a way to rename the vertices of  $G$  so that the result is equal to  $H$ . Another way to think about isomorphism is that two graphs are isomorphic iff their drawings are identical except for the labels on the vertices.

## SUBGRAPH-ISOMORPHISM



## NP-Complete Languages

- BIN-PACKING =  $\{ \langle S, c, k \rangle : S \text{ is a set of objects each of which has an associated size and it is possible to divide the objects so that they fit into } k \text{ bins, each of which has size } c \}$ .

## BIN-PACKING

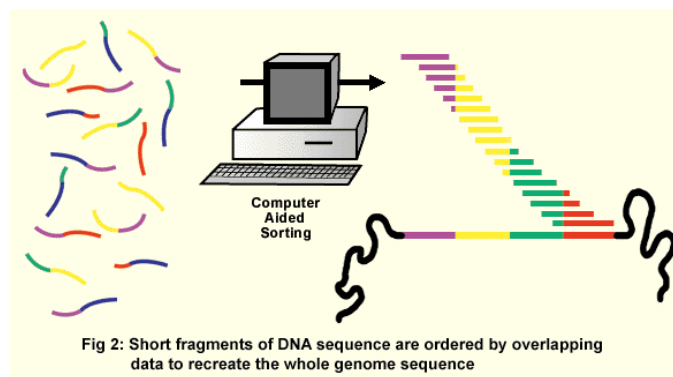
In three dimensions:



## NP-Complete Languages

- SHORTEST-SUPERSTRING =  $\{ \langle S, k \rangle : S \text{ is a set of strings and there exists some superstring } T \text{ such that every element of } S \text{ is a substring of } T \text{ and } T \text{ has length less than or equal to } k \}$ .

## SHORTEST-SUPERSTRING



Source: Wiley: Interactive Concepts in Biology

## Proving that $L$ is NP-Complete

$NPL_1$     $NPL_2$     $NPL_3$     $NPL_4$     $NPL\dots$

$L_1$   
↓  
 $L_2$

**Theorem:**

If:

$L_1$  is NP-complete,

$L_1 \leq_P L_2$ , and

$L_2$  is in NP,

Then  $L_2$  is also NP-complete.