



CSE373: Data Structures & Algorithms

Lecture 10: Amortized Analysis and Disjoint Sets

Aaron Bauer

Winter 2014

Amortized

- Recall our plain-old stack implemented as an array that doubles its size if it runs out of room
 - How can we claim **push** is $O(1)$ time if resizing is $O(n)$ time?
 - We *can't*, but we *can* claim it's an $O(1)$ **amortized operation**
- What does amortized mean?
- When are amortized bounds good enough?
- How can we prove an amortized bound?

Will just do two simple examples

- Text has more sophisticated examples and proof techniques
- *Idea* of how amortized describes average cost is essential

Amortized Complexity

If a sequence of M operations takes $O(M f(n))$ time,
we say the amortized runtime is $O(f(n))$

Amortized bound: worst-case guarantee over sequences of operations

- Example: If any n operations take $O(n)$, then amortized $O(1)$
- Example: If any n operations take $O(n^3)$, then amortized $O(n^2)$

- The worst case time per operation can be larger than $f(n)$
 - As long as the worst case is *always* “rare enough” in *any* sequence of operations

Amortized guarantee ensures the average time per operation for any sequence is $O(f(n))$

“Building Up Credit”

- Can think of preceding “cheap” operations as building up “credit” that can be used to “pay for” later “expensive” operations
- Because any sequence of operations must be under the bound, enough “cheap” operations must come *first*
 - Else a prefix of the sequence, which is also a sequence, would violate the bound

Example #1: Resizing stack

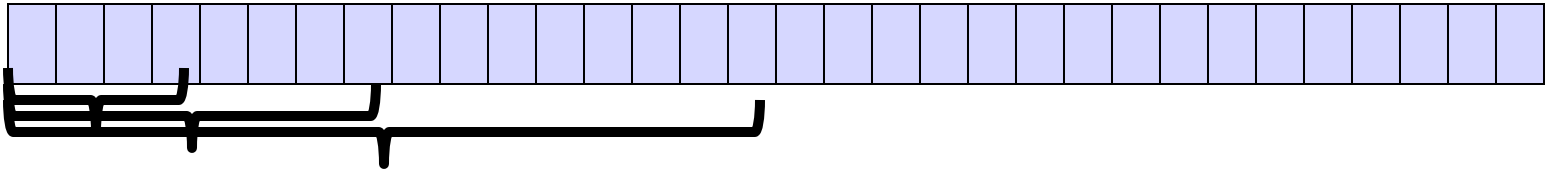
A stack implemented with an array where we double the size of the array if it becomes full

Claim: Any sequence of `push/pop/isEmpty` is amortized $O(1)$

Need to show any sequence of M operations takes time $O(M)$

- Recall the non-resizing work is $O(M)$ (i.e., $M * O(1)$)
- The resizing work is proportional to the total number of element copies we do for the resizing
- So it suffices to show that:
 - After M operations, we have done $< 2M$ total element copies
(So average number of copies per operation is bounded by a constant)

Amount of copying



After M operations, we have done $< 2M$ total element copies

Let n be the size of the array after M operations

- Then we have done a total of:

$$n/2 + n/4 + n/8 + \dots \text{INITIAL_SIZE} < n$$

element copies

- Because we must have done at least enough **push** operations to cause resizing up to size n :

$$M \geq n/2$$

- So

$$2M \geq n > \text{number of element copies}$$

Other approaches

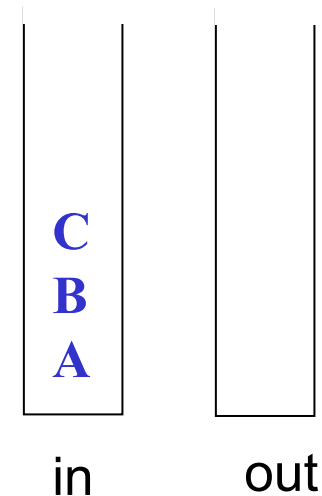
- If array grows by a constant amount (say 1000), operations are **not** amortized $O(1)$
 - After $O(M)$ operations, you may have done $\Theta(M^2)$ copies
- If array shrinks when 1/2 empty, operations are **not** amortized $O(1)$
 - Terrible case: **pop** once and shrink, **push** once and grow, **pop** once and shrink, ...
- If array shrinks when 3/4 empty, it **is** amortized $O(1)$
 - Proof is more complicated, but basic idea remains: by the time an expensive operation occurs, many cheap ones occurred

Example #2: Queue with two stacks

A clever and simple queue implementation using only stacks

```
class Queue<E> {
    Stack<E> in = new Stack<E>();
    Stack<E> out = new Stack<E>();
    void enqueue(E x) { in.push(x); }
    E dequeue() {
        if(out.isEmpty()) {
            while(!in.isEmpty()) {
                out.push(in.pop());
            }
        }
        return out.pop();
    }
}
```

enqueue: A, B, C

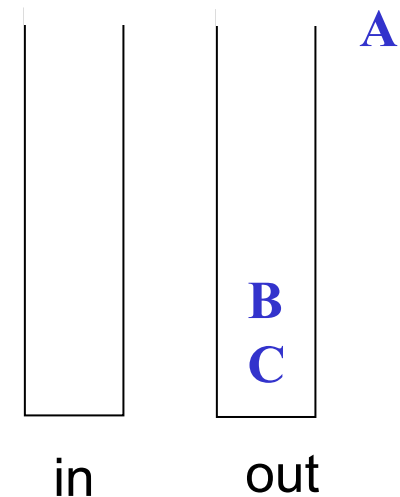


Example #2: Queue with two stacks

A clever and simple queue implementation using only stacks

```
class Queue<E> {  
    Stack<E> in = new Stack<E>();  
    Stack<E> out = new Stack<E>();  
    void enqueue(E x) { in.push(x); }  
    E dequeue() {  
        if (out.isEmpty()) {  
            while (!in.isEmpty()) {  
                out.push(in.pop());  
            }  
        }  
        return out.pop();  
    }  
}
```

dequeue

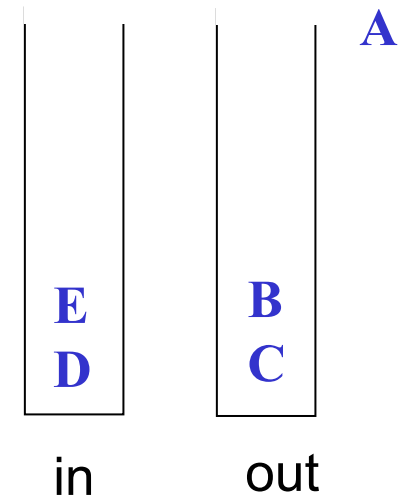


Example #2: Queue with two stacks

A clever and simple queue implementation using only stacks

```
class Queue<E> {
    Stack<E> in = new Stack<E>();
    Stack<E> out = new Stack<E>();
    void enqueue(E x) { in.push(x); }
    E dequeue() {
        if(out.isEmpty()) {
            while(!in.isEmpty()) {
                out.push(in.pop());
            }
        }
        return out.pop();
    }
}
```

enqueue D, E

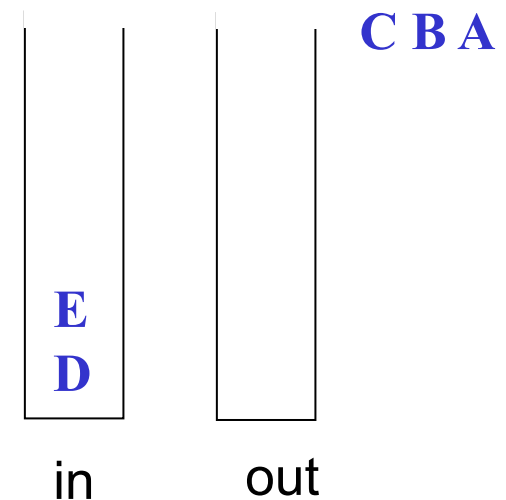


Example #2: Queue with two stacks

A clever and simple queue implementation using only stacks

```
class Queue<E> {  
    Stack<E> in = new Stack<E>();  
    Stack<E> out = new Stack<E>();  
    void enqueue(E x) { in.push(x); }  
    E dequeue() {  
        if (out.isEmpty()) {  
            while (!in.isEmpty()) {  
                out.push(in.pop());  
            }  
        }  
        return out.pop();  
    }  
}
```

dequeue twice

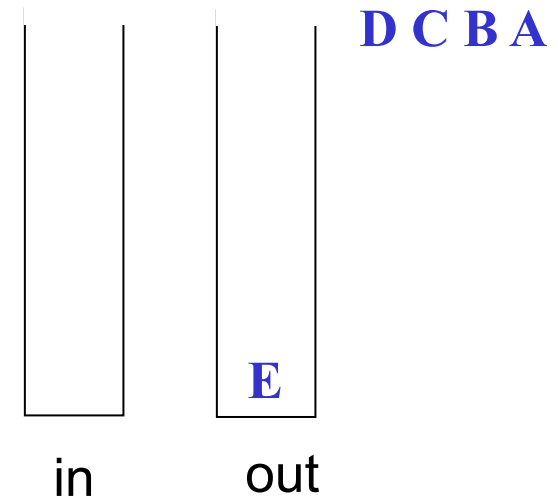


Example #2: Queue with two stacks

A clever and simple queue implementation using only stacks

```
class Queue<E> {
    Stack<E> in = new Stack<E>();
    Stack<E> out = new Stack<E>();
    void enqueue(E x) { in.push(x); }
    E dequeue() {
        if(out.isEmpty()) {
            while(!in.isEmpty()) {
                out.push(in.pop());
            }
        }
        return out.pop();
    }
}
```

dequeue again



Correctness and usefulness

- If x is enqueued before y , then x will be popped from **in** later than y and therefore popped from **out** sooner than y
 - So it is a queue
- Example:
 - Wouldn't it be nice to have a queue of t-shirts to wear instead of a stack (like in your dresser)?
 - So have two stacks
 - *in*: stack of t-shirts go after you wash them
 - *out*: stack of t-shirts to wear
 - if *out* is empty, reverse *in* into *out*

Analysis

- **dequeue** is not $O(1)$ worst-case because **out** might be empty and **in** may have lots of items
- But if the stack operations are (amortized) $O(1)$, then any sequence of queue operations is amortized $O(1)$
 - The total amount of work done per element is 1 **push** onto **in**, 1 **pop** off of **in**, 1 **push** onto **out**, 1 **pop** off of **out**
 - When you reverse n elements, there were n earlier $O(1)$ **enqueue** operations to average with

Amortized useful?

- When the average per operation is all we care about (i.e., sum over all operations), amortized is perfectly fine
 - This is the usual situation
- If we need every operation to finish quickly (e.g., in a web server), amortized bounds may be too weak
- While amortized analysis is about averages, we are averaging cost-per-operation on worst-case input
 - Contrast: Average-case analysis is about averages across possible inputs. Example: if all initial permutations of an array are equally likely, then quicksort is $O(n \log n)$ on average even though on some inputs it is $O(n^2)$

Not always so simple

- Proofs for amortized bounds can be much more complicated
- Example: Splay trees are dictionaries with amortized $O(\log n)$ operations
 - No extra height field like AVL trees
 - See Chapter 4.5 if curious
- For more complicated examples, the proofs need much more sophisticated invariants and “potential functions” to describe how earlier cheap operations build up “energy” or “money” to “pay for” later expensive operations
 - See Chapter 11 if curious
- But complicated *proofs* have nothing to do with the code!

The plan

- What are *disjoint sets*
 - And how are they “the same thing” as *equivalence relations*
- The union-find ADT for disjoint sets
- Applications of union-find

Next lecture:

- Basic implementation of the ADT with “up trees”
- Optimizations that make the implementation much faster

Disjoint sets

- A **set** is a collection of elements (no-repeats)
- Two sets are **disjoint** if they have no elements in common
 - $S_1 \cap S_2 = \emptyset$
- Example: {a, e, c} and {d, b} are disjoint
- Example: {x, y, z} and {t, u, x} are not disjoint

Partitions

A **partition** P of a set S is a set of sets $\{S_1, S_2, \dots, S_n\}$ such that every element of S is in **exactly one** S_i

Put another way:

- $S_1 \cup S_2 \cup \dots \cup S_k = S$
- $i \neq j$ implies $S_i \cap S_j = \emptyset$ (sets are disjoint with each other)

Example:

- Let S be $\{a, b, c, d, e\}$
- One partition: $\{a\}, \{d, e\}, \{b, c\}$
- Another partition: $\{a, b, c\}, \emptyset, \{d\}, \{e\}$
- A third: $\{a, b, c, d, e\}$
- Not a partition: $\{a, b, d\}, \{c, d, e\}$
- Not a partition of S : $\{a, b\}, \{e, c\}$

Binary relations

- $S \times S$ is the set of all pairs of elements of S
 - Example: If $S = \{a,b,c\}$
then $S \times S = \{(a,a),(a,b),(a,c),(b,a),(b,b),(b,c), (c,a),(c,b),(c,c)\}$
- A **binary relation** R on a set S is any subset of $S \times S$
 - Write $R(x,y)$ to mean (x,y) is “in the relation”
 - (Unary, ternary, quaternary, ... relations defined similarly)
- Examples for $S = \text{people-in-this-room}$
 - Sitting-next-to-each-other relation
 - First-sitting-right-of-second relation
 - Went-to-same-high-school relation
 - Same-gender-relation
 - First-is-younger-than-second relation

Properties of binary relations

- A binary relation R over set S is **reflexive** means
 $R(a,a)$ for *all* a in S
- A binary relation R over set S is **symmetric** means
 $R(a,b)$ if and only if $R(b,a)$ for *all* a,b in S
- A binary relation R over set S is **transitive** means
If $R(a,b)$ and $R(b,c)$ then $R(a,c)$ for *all* a,b,c in S
- Examples for $S =$ people-in-this-room
 - Sitting-next-to-each-other relation
 - First-sitting-right-of-second relation
 - Went-to-same-high-school relation
 - Same-gender-relation
 - First-is-younger-than-second relation

Equivalence relations

- A binary relation R is an **equivalence relation** if R is reflexive, symmetric, and transitive
- Examples
 - Same gender
 - Connected roads in the world
 - *Graduated* from same high school?
 - ...

Punch-line

- Every partition induces an *equivalence relation*
- Every equivalence relation *induces* a partition

- Suppose $P = \{S_1, S_2, \dots, S_n\}$ be a partition
 - Define $R(x, y)$ to mean x and y are in the same S_i
 - R is an equivalence relation

- Suppose R is an equivalence relation over S
 - Consider a set of sets S_1, S_2, \dots, S_n where
 - (1) x and y are in the same S_i if and only if $R(x, y)$
 - (2) Every x is in some S_i
 - This set of sets is a partition

Example

- Let S be $\{a,b,c,d,e\}$
- One partition: $\{a,b,c\}, \{d\}, \{e\}$
- The corresponding equivalence relation:
 $(a,a), (b,b), (c,c), (a,b), (b,a), (a,c), (c,a), (b,c), (c,b), (d,d), (e,e)$

The plan

- What are *disjoint sets*
 - And how are they “the same thing” as *equivalence relations*
- The union-find ADT for disjoint sets
- Applications of union-find

Next lecture:

- Basic implementation of the ADT with “up trees”
- Optimizations that make the implementation much faster

The operations

- Given an unchanging set S , **create** an initial partition of a set
 - Typically each item in its own subset: $\{a\}$, $\{b\}$, $\{c\}$, ...
 - Give each subset a “name” by choosing a *representative element*
- Operation **find** takes an element of S and returns the representative element of the subset it is in
- Operation **union** takes two subsets and (permanently) makes one larger subset
 - A different partition with one fewer set
 - Affects result of subsequent **find** operations
 - Choice of representative element up to implementation

Example

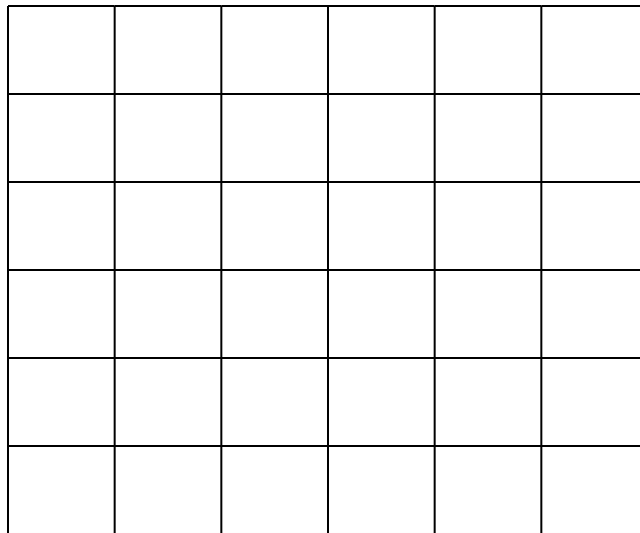
- Let $S = \{1,2,3,4,5,6,7,8,9\}$
- Let initial partition be (will highlight representative elements red)
 $\{\underline{1}\}, \{\underline{2}\}, \{\underline{3}\}, \{\underline{4}\}, \{\underline{5}\}, \{\underline{6}\}, \{\underline{7}\}, \{\underline{8}\}, \{\underline{9}\}$
- `union(2,5)`:
 $\{\underline{1}\}, \{\underline{2}, 5\}, \{\underline{3}\}, \{\underline{4}\}, \{\underline{6}\}, \{\underline{7}\}, \{\underline{8}\}, \{\underline{9}\}$
- `find(4) = 4`, `find(2) = 2`, `find(5) = 2`
- `union(4,6)`, `union(2,7)`
 $\{\underline{1}\}, \{\underline{2}, 5, 7\}, \{\underline{3}\}, \{4, \underline{6}\}, \{\underline{8}\}, \{\underline{9}\}$
- `find(4) = 6`, `find(2) = 2`, `find(5) = 2`
- `union(2,6)`
 $\{\underline{1}\}, \{\underline{2}, 4, 5, 6, 7\}, \{\underline{3}\}, \{\underline{8}\}, \{\underline{9}\}$

No other operations

- All that can “happen” is sets get unioned
 - No “un-union” or “create new set” or ...
- As always: trade-offs – implementations will exploit this small ADT
- Surprisingly useful ADT: list of applications after one example surprising one
 - But not as common as dictionaries or priority queues

Example application: maze-building

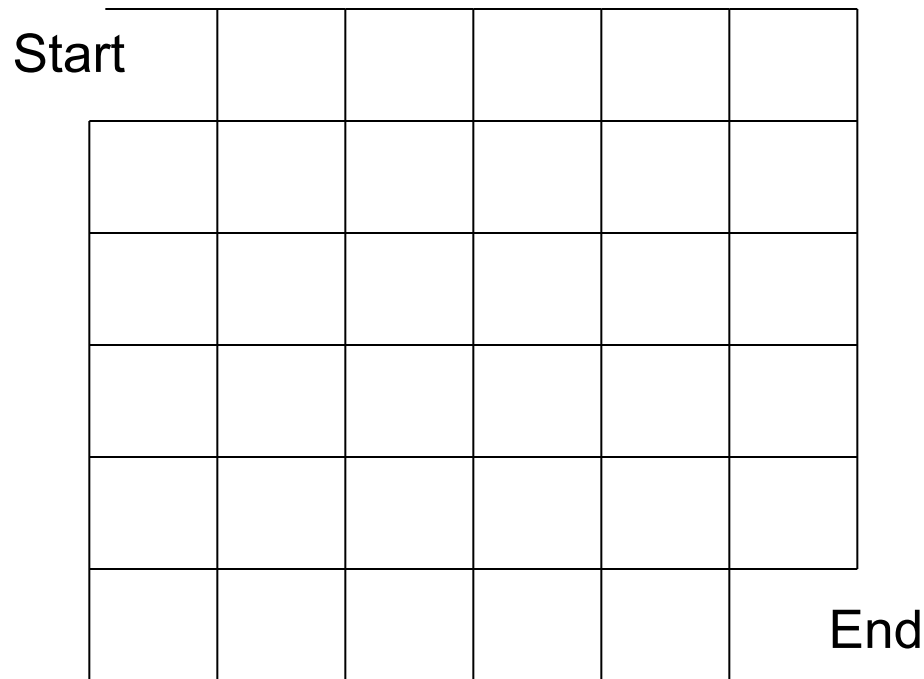
- Build a random maze by erasing edges



- Possible to get from anywhere to anywhere
 - Including “start” to “finish”
- No loops possible without backtracking
 - After a “bad turn” have to “undo”

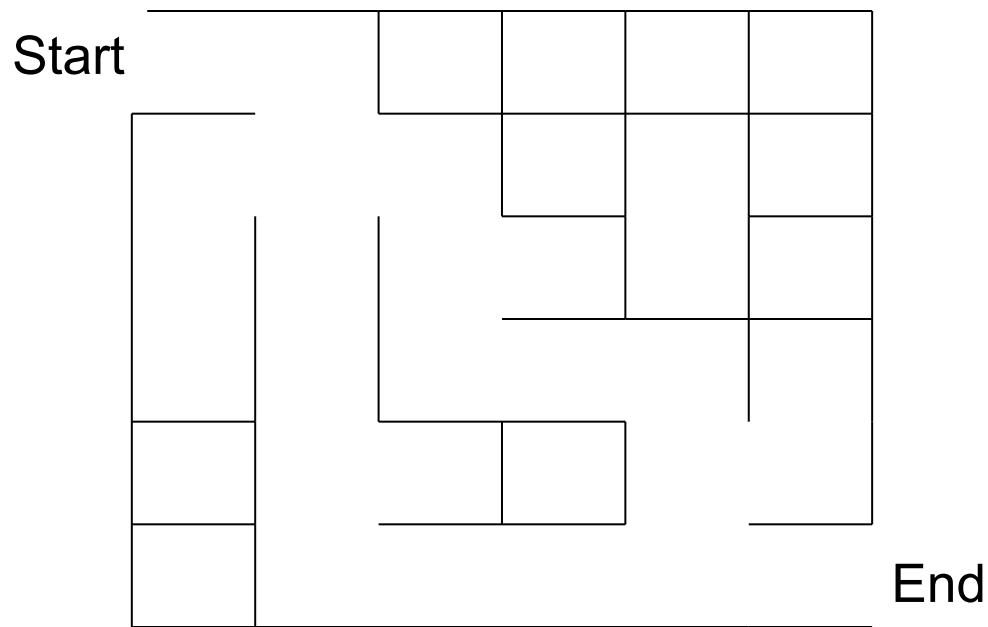
Maze building

Pick start edge and end edge



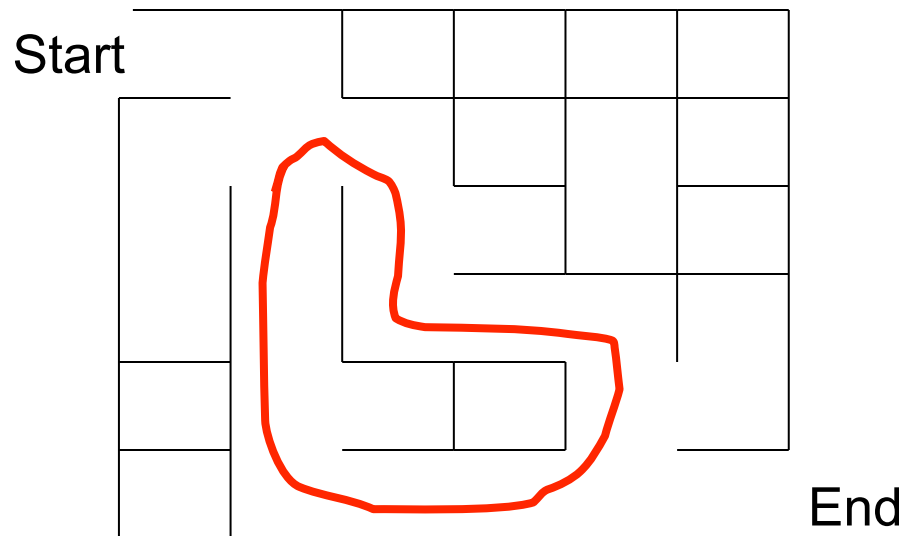
Repeatedly pick random edges to delete

One approach: just keep deleting random edges until you can get from start to finish



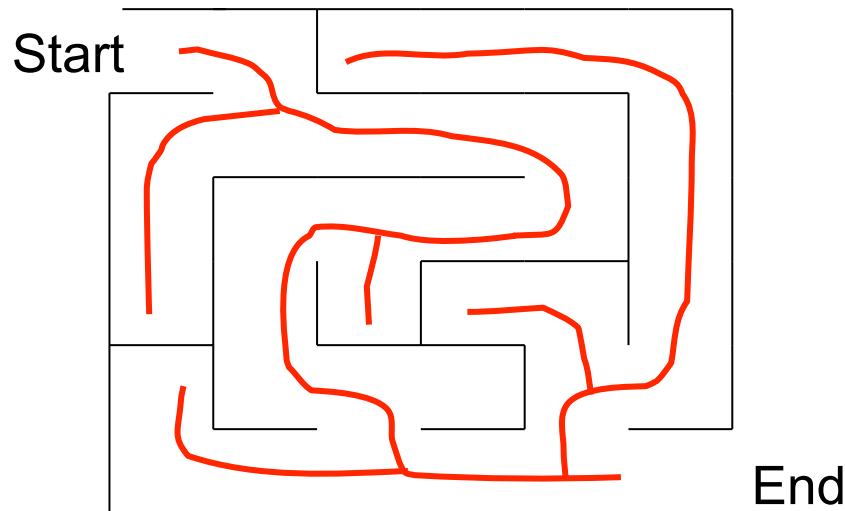
Problems with this approach

1. How can you tell when there is a path from start to finish?
 - We do not really have an algorithm yet
2. We have *cycles*, which a “good” maze avoids
 - Want one solution and no cycles



Revised approach

- Consider edges in random order
- But only delete them if they introduce no cycles (how? TBD)
- When done, will have one way to get from any place to any other place (assuming no backtracking)



- Notice the funny-looking *tree* in red

Cells and edges

- Let's number each cell
 - 36 total for 6 x 6
- An (internal) edge (x,y) is the line between cells x and y
 - 60 total for 6x6: (1,2), (2,3), ..., (1,7), (2,8), ...

Start	1	2	3	4	5	6	
	7	8	9	10	11	12	
	13	14	15	16	17	18	
	19	20	21	22	23	24	
	25	26	27	28	29	30	
	31	32	33	34	35	36	End

The trick

- Partition the cells into **disjoint sets**: “are they connected”
 - Initially every cell is in its own subset
- If an edge would connect two different subsets:
 - then remove the edge and **union** the subsets
 - else leave the edge because removing it makes a cycle

Start

1	2	3	4	5	6
7	8	9	10	11	12
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30
31	32	33	34	35	36

Start

1	2	3	4	5	6
7	8	9	10	11	12
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30
31	32	33	34	35	36

End

The algorithm

- **P** = **disjoint sets** of connected cells, initially each cell in its own 1-element set
- **E** = **set** of edges not yet processed, initially all (internal) edges
- **M** = **set** of edges kept in maze (initially empty)

while P has more than one set {

– Pick a random edge (x,y) to remove from E

– **u** = **find**(x)

– **v** = **find**(y)

– if **u**==**v**

 then add (x,y) to M // same subset, do not create cycle

 else **union**(u,v) // do not put edge in M, connect subsets

}

Add remaining members of E to M, then output M as the maze

Example step

Pick (8,14)

Start	1	2	3	4	5	6	
	7	8	9	10	11	12	
	13	14	15	16	17	18	
	19	20	21	22	23	24	
	25	26	27	28	29	30	
	31	32	33	34	35	36	End

P
{1,2,7,8,9,13,19}
{3}
{4}
{5}
{6}
{10}
{11,17}
{12}
{14,20,26,27}
{15,16,21}
{18}
{25}
{28}
{31}
{22,23,24,29,30,32
33,34,35,36}

Example step

P

{1,2,7,8,9,13,19}

{3}

{4}

{5}

{6}

{10}

{11,17}

{12}

{14,20,26,27}

{15,16,21}

{18}

{25}

{28}

{31}

{22,23,24,29,30,32

33,34,35,36}

Find(8) = 7

Find(14) = 20

Union(7,20)



P

{1,2,7,8,9,13,19,14,20,26,27}

{3}

{4}

{5}

{6}

{10}

{11,17}

{12}

{15,16,21}

{18}

{25}

{28}

{31}

{22,23,24,29,30,32

33,34,35,36}

Add edge to M step

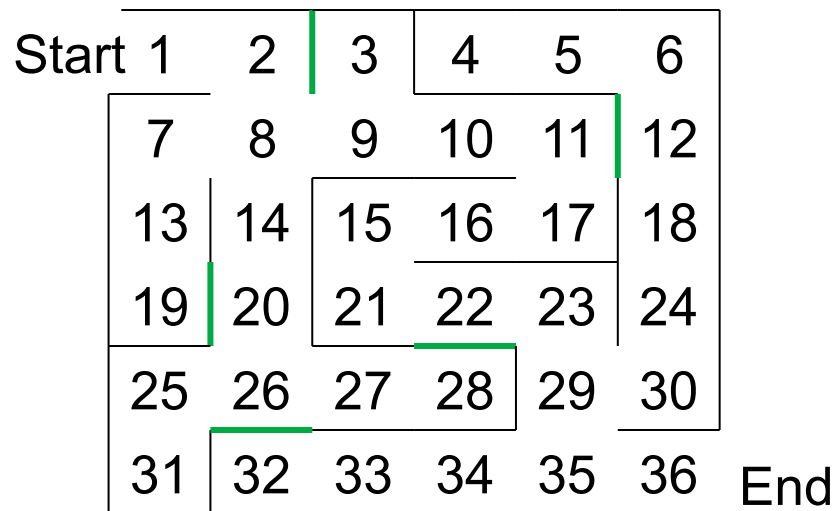
Pick (19,20)

Start	1	2	3	4	5	6	
	7	8	9	10	11	12	
	13	14	15	16	17	18	
	19	20	21	22	23	24	
	25	26	27	28	29	30	
	31	32	33	34	35	36	End

P
 {1,2,7,8,9,13,19,14,20,26,27}
 {3}
 {4}
 {5}
 {6}
 {10}
 {11,17}
 {12}
 {15,16,21}
 {18}
 {25}
 {28}
 {31}
 {22,23,24,29,30,32
 33,34,35,36}

At the end

- Stop when P has one set
- Suppose green edges are already in M and black edges were not yet picked
 - Add all black edges to M



P
{1,2,3,4,5,6,7,... 36}

Other applications

- Maze-building is:
 - Cute
 - Homework 4 😊
 - A surprising use of the union-find ADT
- Many other uses (which is why an ADT taught in CSE373):
 - Road/network/graph connectivity (will see this again)
 - “connected components” e.g., in social network
 - Partition an image by connected-pixels-of-similar-color
 - Type inference in programming languages
- Not as common as dictionaries, queues, and stacks, but valuable because implementations are very fast, so when applicable can provide big improvements