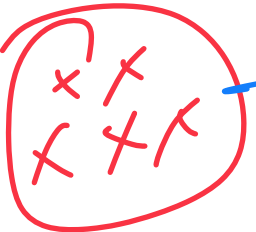


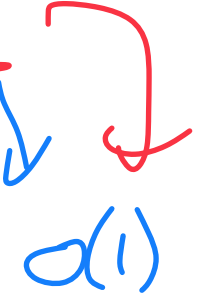
classroom A



$O(1)$



Find is  $O(\log n)$



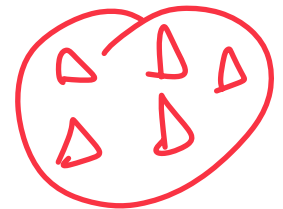
# Data Structures and Algorithms

## Probability in Computer Science

CS 225  
Brad Solomon

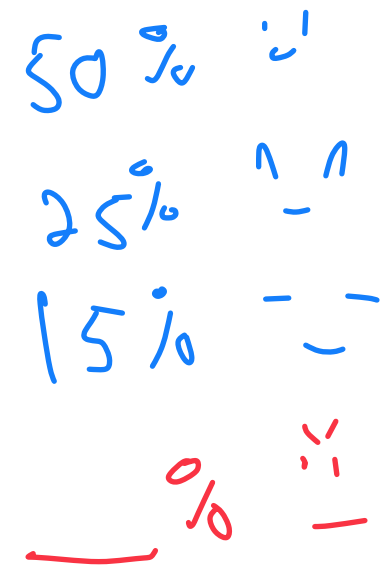
October 23, 2022

classroom B



UNIVERSITY OF  
**ILLINOIS**  
URBANA - CHAMPAIGN

Department of Computer Science



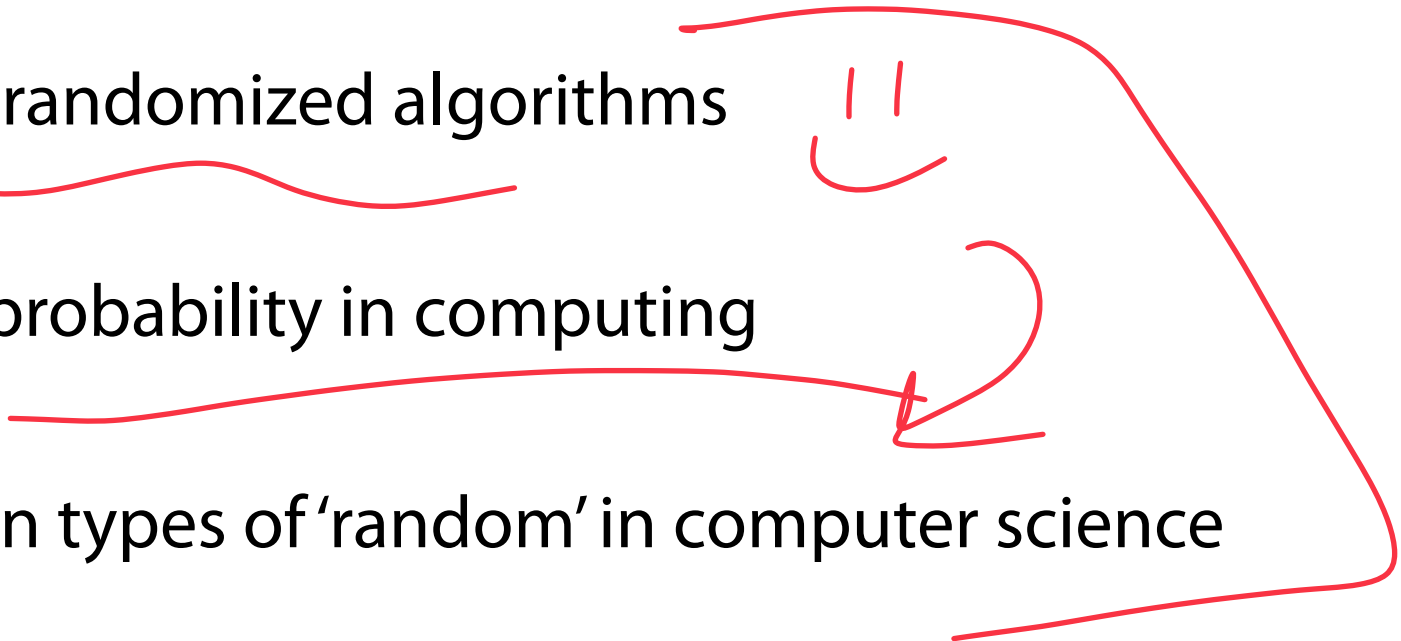
# Learning Objectives

Finish disjoint set analysis (one final proof)

Formalize the concept of randomized algorithms

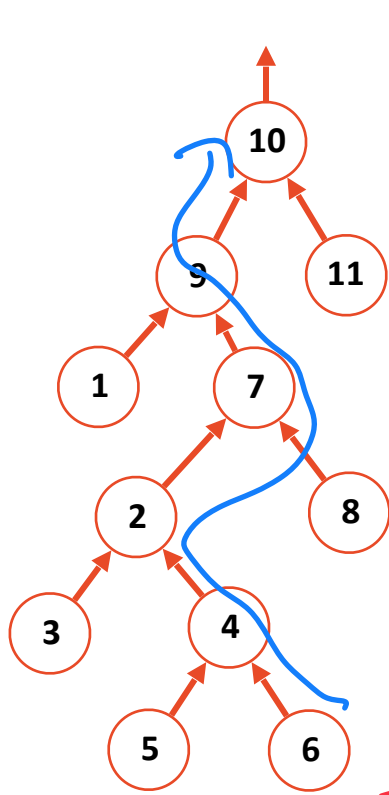
Review fundamentals of probability in computing

Distinguish the three main types of 'random' in computer science



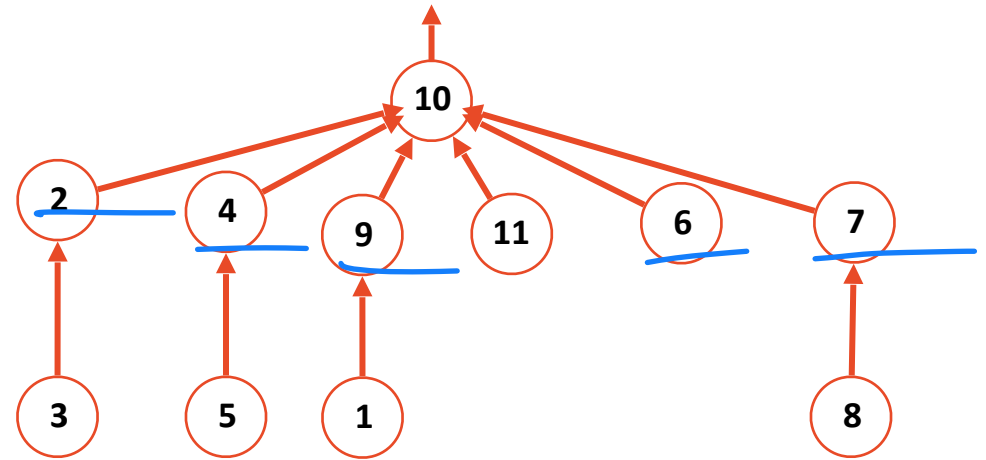
# Disjoint Sets w/ Path Compression (and rank)

How do we observe how the efficiency of a set changes due to PC?



height  
 $\log n$

Single op can be slow!



Amortized time!

Eventually speed up

# Amortized Time Review

We have **n items**. We make **n insert()** calls.

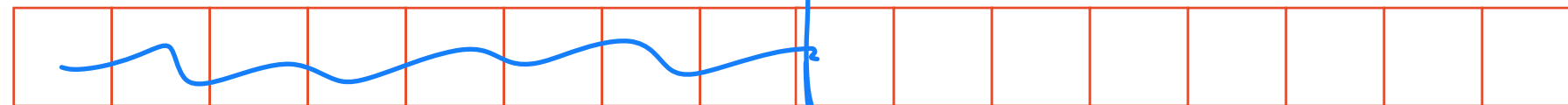
We are interested in the worst case work possible over n calls.

A insert(B)  $O(n) + O(1) + O(1) + O(1) + O(1) + \dots$

$\downarrow$   
A B  $\rightarrow$  only every n steps



$\swarrow$  refill up to n planets



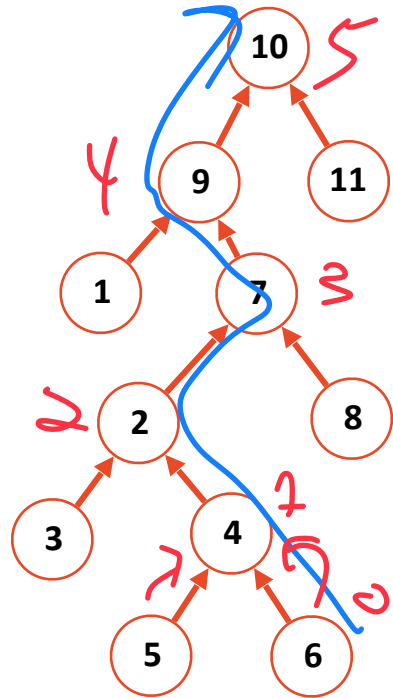
# Amortized Time (Rank w/ Path Compression)

We have **n items** in an Uptree. We make **m find()** calls. → as I do this I can save work!

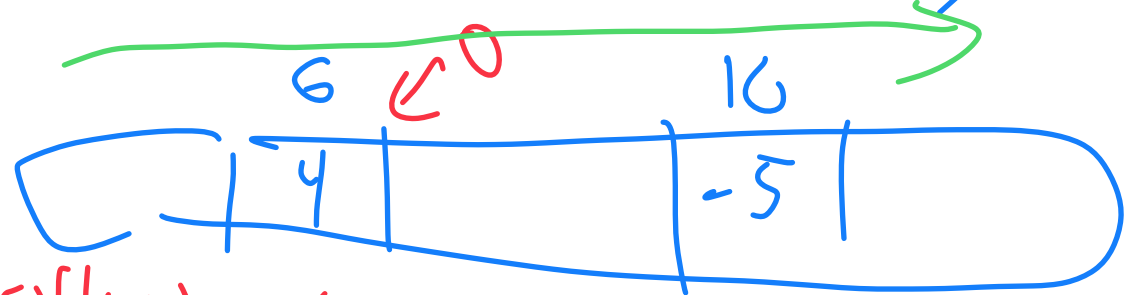
We are interested in the **worst case work** possible **over m calls**.

Once we find once, save work in future!

↳ Max work we can do → ∞



$m \rightarrow \infty$  walk  $O(n)$  and  $O(1)$



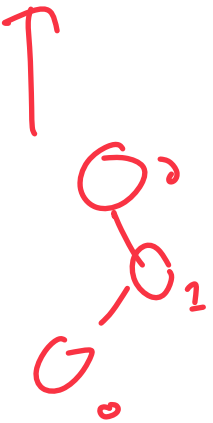
B/C this is optimal storage of set 10  $O(1)$

1 find -  $O(\log n)$

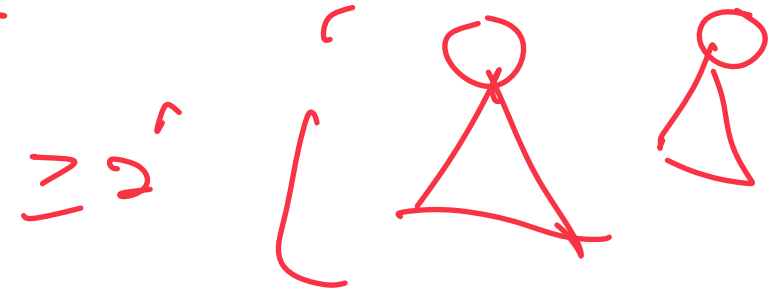
# Key Properties of UpTree by rank w/ PC

The parent of a node is always higher rank than the node.

↳ we only store rank of roots but we can label all nodes



The min(nodes) in a set with a root of rank  $r$  has  $\geq 2^r$  nodes.



For any integer  $r$ , there are at most  $\frac{n}{2^r}$  nodes of rank  $r$ .

↳ total elements



# Amortized Time (Rank w/ Path Compression)

Let's  $0, \dots, X$  be a set of buckets

Put every non-root node in a bucket by rank!

Invent buckets

Ranks	Bucket
0	0
1	1
2 - 3	2
4 - 15	3
16 - 65535	4
65536 - $2^{65536}-1$	5

Structure buckets to store ranks  $[r, 2^r - 1]$

Grows very very slowly  $\log^*(n)$

After class question:

This partition makes  $\log^*$  buckets!  
(By design)

$2^{65536}$   $\approx$  larger than atoms in universe.

# Iterated Logarithm Function ( $\log^* n$ )

$\log^* n$  is piecewise defined as

$$0 \text{ if } n \leq 1$$

otherwise

$$1 + \log^*(\log n)$$

the # of times we need to call log  
to get  $n \leq 1$



# Amortized Time (Rank w/ Path Compression)

Let  $|B_r|$  be the size of the bucket with min rank  $r$ .

What is  $\max(|B_r|)$ ?  $\sim 1/2^r$

$r=0$ : All items w rank 0

⊗ ⊗ ⊗ ⊗ ⊗ ...

$r=1$ :  $\sim 1/2$

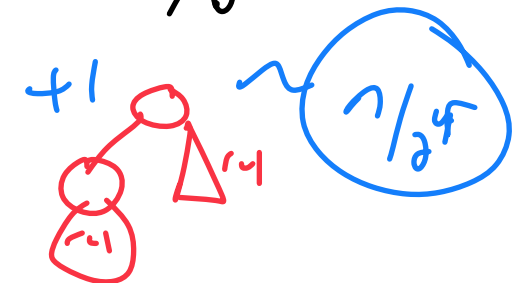
$r=2$ :  $\sim 1/4$

1) We can sum all possible ranks in each bucket

$$1/2^r + 1/2^{r+1} + 1/2^{r+2} + \dots + 1/2^{2^r-1} \approx 2/2^r$$

2) To get rank  $r$ ,  $\pm$  combined two  $(r-1)$  items

Ranks	Bucket
0	0
1	1
2-3	2
4-15	3
16-65535	4
65536- $2^{65536}-1$	5



# Amortized Time (Rank w/ Path Compression)

The work of find(x) are the steps taken on the path from a node x to the root (or immediate child of the root) of the UpTree containing x

We can split this into two cases:

**Case 1:** We take a step from one bucket to another bucket.

↳ Easy!

At most 6

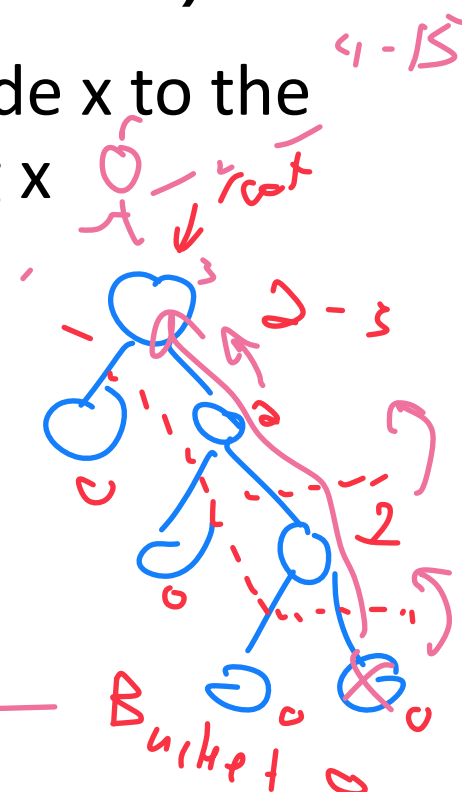
Every find will do this work

2 on fibos

$$O(\log^* n)$$

$$M \log^* n$$

$$O(1)$$



**Case 2:** We take a step from one item to another inside the same bucket.

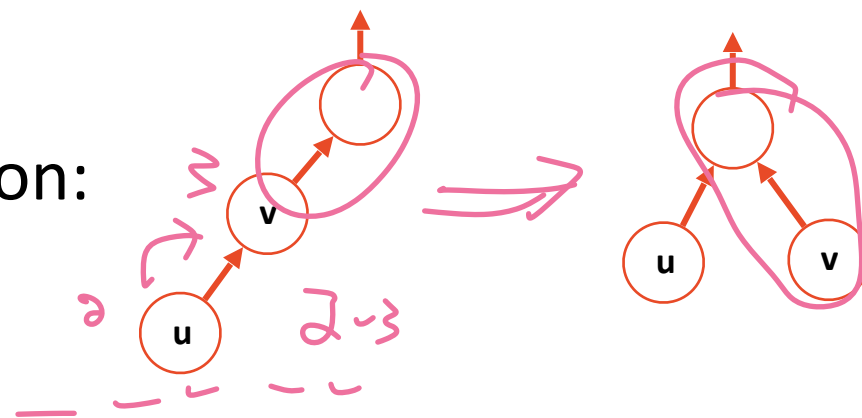
# Amortized Time (Rank w/ Path Compression)

**Case 2:** We take a step from one item to another *inside* the same bucket.

Let's call this the step from  $u$  to  $v$ .

Every time we do this, we do path compression:

We set  $\text{parent}(u)$  a little closer to root



How many total times can I do this for each  $u$  in  $|B_r|$ ? # of ranks in my bucket

B/c parent has higher rank

$1, \dots, 2^r - 1$

$2^r - r - 1 \approx 2^r$

How many nodes are in  $|B_r|$ ?

$\approx 2^r$

Each node ( $\approx 1/8^r$ ) can increment at most  $2^r$   $\downarrow$   
 $\hookrightarrow$   $n$  total work  $\approx \log^* n$



# Final Result

We have **n items** in an Uptree. We make **m find()** calls. Total work is:

cost between buckets

$$m \cdot \log^* n$$

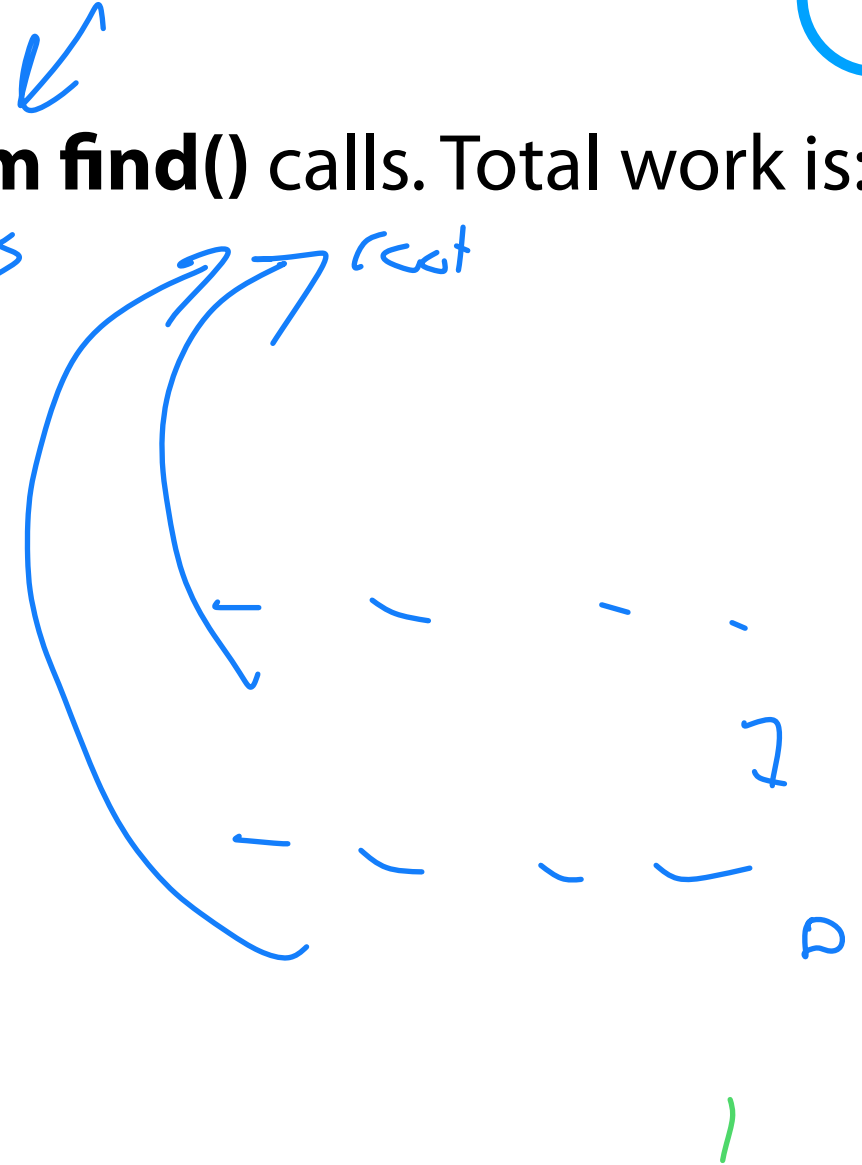
cost inside buckets

$$+ n \cdot \log^* n$$

$$(m+n) \log^* n$$



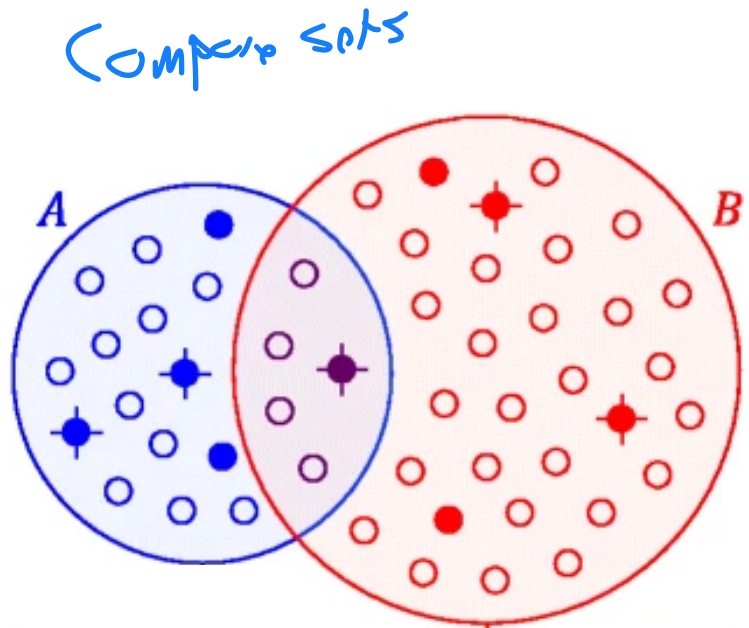
≠ actually my finds we just this



# Randomized Algorithms



A **randomized algorithm** is one which uses a source of randomness somewhere in its implementation.



approximate membership

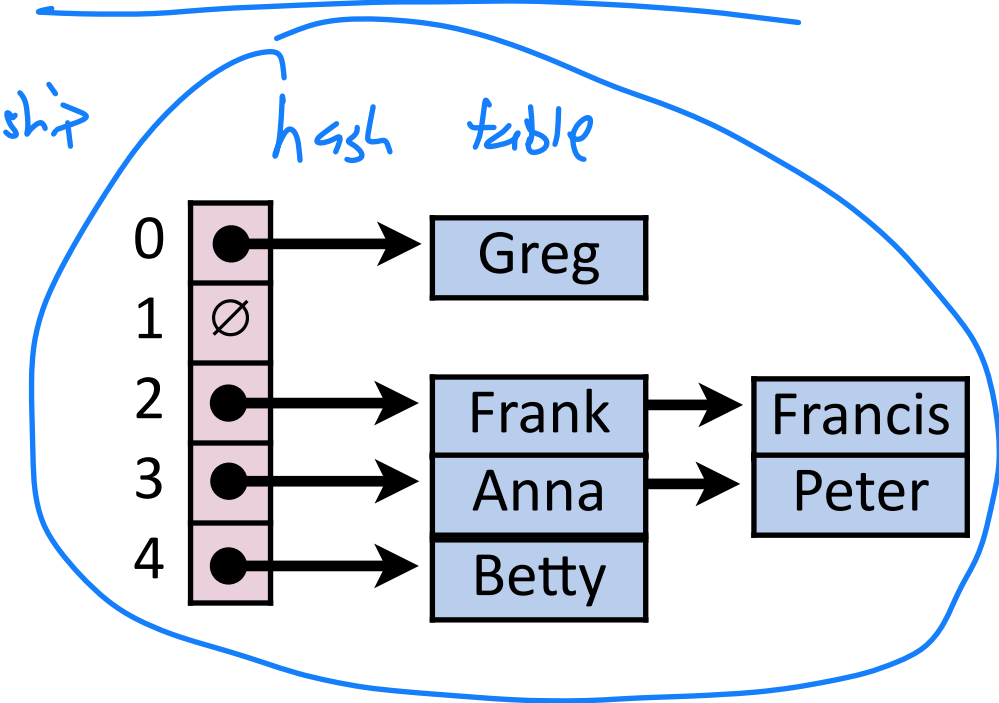
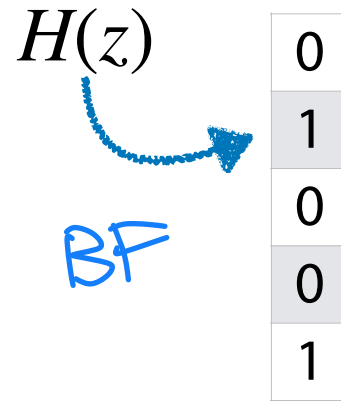


Figure from Ondov et al 2016

Min hash

$H(x)$	0	2	1	0	0	4	0	2	0	6
$H(y)$	1	0	2	3	1	0	3	4	0	1
$H(z)$	2	1	0	2	0	1	0	0	7	2

Approx  
< min

# A faulty list

Imagine you have a list ADT implementation **except**...

all on a coin  
FLIP

Every time you called **insert**, it would fail 50% of the time.

Website caching

Scrambling data

Buffering, throwing out some data

Encryption  $\rightarrow$  lossy!  $\rightarrow$  we don't want random

A lot!

A PQ approximate counting!  $\sim 1/2$  real count

# Quick Primes with Fermat's Primality Test

*is this prime?*  
If  $p$  is prime and  $a$  is not divisible by  $p$ , then  $a^{p-1} \equiv 1 \pmod{p}$  *the always* ✓  
*↳ a random #*

But... **sometimes** if  $n$  is composite and  $a^{n-1} \equiv 1 \pmod{n}$  ✓

↳  $a=2 \rightarrow$  21,853 Pseudo Prime  $\rightarrow$  not prime satisfy prime

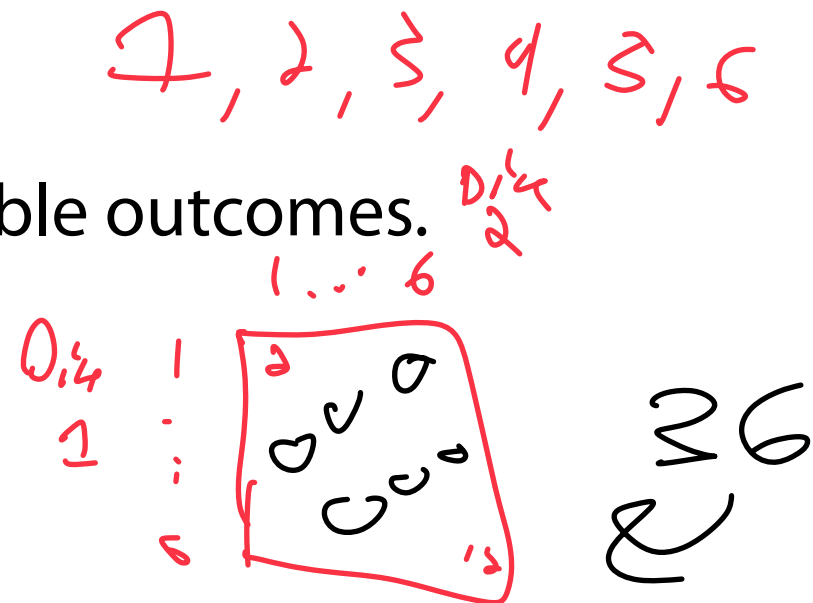
$25 \cdot 10^9$

↳ very very low chance of failure!

# Fundamentals of Probability

Imagine you roll a pair of six-sided dice.

The **sample space**  $\Omega$  is the set of all possible outcomes.



An **event**  $E \subseteq \Omega$  is any subset.

$\hookrightarrow D_1 = 2, D_2 = 4 \rightarrow$  Prob of occurrence

$\hookrightarrow D_1 + D_2 = >$

$D_1$  is odd



We can't use most CSP

# Fundamentals of Probability

Imagine you roll a pair of six-sided dice. What is the expected value?

A **random variable** is a function from events to numeric values.

Let  $D_1$  be the value of the first dice  
 $D_2$  Second dice

↓ simplification?

The **expectation** of a (discrete) random variable is:

$$E[X] = \sum_{x \in \Omega} \text{Pr}\{X = x\} \cdot x$$

Sum of states not R.V. (under  $\sum$ )

Sum over sample space (under  $\sum$ )

prob of one state (over  $\text{Pr}\{X = x\}$ )

value of the state (over  $x$ )

$$E[D_1] = \frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 2 + \frac{1}{6} \cdot 3 + \frac{4}{6} + \frac{5}{6} + \frac{6}{6}$$

$$= 3.5$$

$$E[D_1 + D_2] = \frac{1}{36} (1+1) + \frac{1}{36} (1+2) \dots$$

# Fundamentals of Probability

Imagine you roll a pair of six-sided dice. What is the expected value?

**Linearity of Expectation:** For any two random variables  $X$  and  $Y$ ,

$$E[X + Y] = E[X] + E[Y] \text{ (Claim)}$$

$$= \sum_x \sum_y \underbrace{Pr(X=x, Y=y)}_{\text{Prob event}} \cdot \underbrace{(x+y)}_{\text{value}}$$

Split

$$= \sum_x x \sum_y Pr(X=x, Y=y) + \sum_y y \sum_x Pr(X=x, Y=y)$$

All prob sum to 1

$$Pr(X=x) \cdot x + \sum_y Pr(Y=y) \cdot y$$

# Fundamentals of Probability

Imagine you roll a pair of six-sided dice. What is the expected value?

**Linearity of Expectation:** For any two random variables  $X$  and  $Y$ ,

$$E[X + Y] = E[X] + E[Y]$$

$$\begin{aligned} E[X + Y] &= \sum_x \sum_y \Pr\{X = x, Y = y\} (x + y) \\ &= \sum_x x \sum_y \Pr\{X = x, Y = y\} + \sum_y y \sum_x \Pr\{X = x, Y = y\} \\ &= \sum_x x \cdot \Pr\{X = x\} + \sum_y y \cdot \Pr\{Y = y\} \end{aligned}$$

*Handwritten notes:*

- Blue arrow pointing to the double sum: "split in half"
- Red 'X' marks over  $Y = y$  in the second term of the second line and  $X = x$  in the second term of the third line.
- Blue circles around the inner sums in the second line.
- Red underline under the final expression.
- Red labels  $E[X]$  and  $E[Y]$  under the first and second terms of the final expression, respectively.

# Fundamentals of Probability



Imagine you roll a pair of six-sided dice. What is the expected value?

**Linearity of Expectation:** For any two random variables  $X$  and  $Y$ ,

$$E[X + Y] = E[X] + E[Y]$$

# Randomization in Algorithms

1. Assume input data is random to estimate average-case performance
2. Use randomness inside algorithm to estimate expected running time
3. Use randomness inside algorithm to approximate solution in fixed time