# Data Structures and Algorithms
# Probability in Computer Science

CS 225
Carl Evans

**UNIVERSITY OF ILLINOIS**
**URBANA-CHAMPAIGN**

Department of Computer Science

# Learning Objectives

Formalize the concept of randomized algorithms

Review fundamentals of probability in computing

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

# Randomized Algorithms

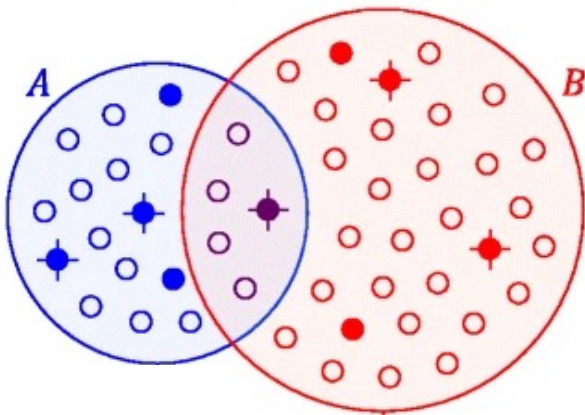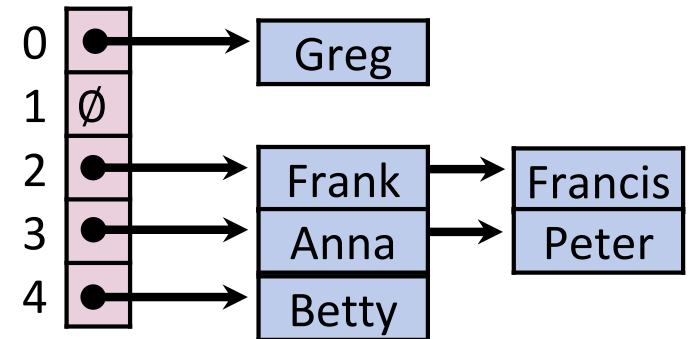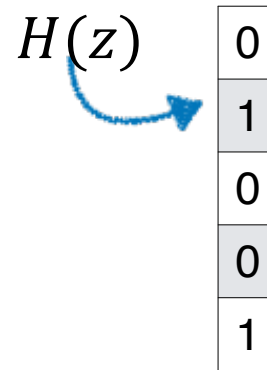A **randomized algorithm** is one which uses a source of randomness somewhere in its implementation.



Figure from Ondov et al 2016

$H(z)$

| | |
|---|---|
| 0 | • → Greg |
| 1 | ∅ |
| 2 | • → Frank → Francis |
| 3 | • → Anna → Peter |
| 4 | • → Betty |

$H(z)$

| |
|---|
| 0 |
| 1 |
| 0 |
| 0 |
| 1 |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| $H(x)$ | 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 |

# Quick Primes with Fermat's Primality Test

If $p$ is prime and $a$ is not divisible by $p$, then $a^{p-1} \equiv 1 \pmod{p}$

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

# 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.

# 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.

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

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

# Fundamentals of Probability

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

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

# 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 $Y$,

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

# 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 $Y$,

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

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

$$] = \sum_x x \sum_y Pr\{X = x, Y = y\} + \sum_y y \sum_x Pr\{X = x, Y = y\}$$

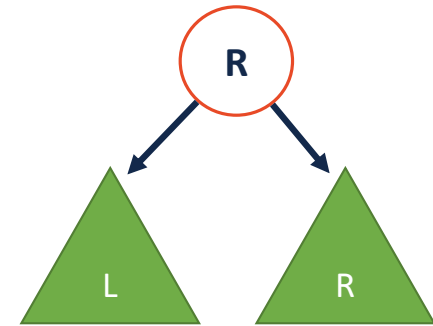$$] = \sum_x x \cdot Pr\{X = x\} + \sum_y y \cdot Pr\{Y = y\}$$

# 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

# Average-Case Analysis: BST



Smallest ——————————————— Largest

# Average-Case Analysis: BST

Let $S(n)$ be the average **total internal path length** over all BSTs that can be constructed by uniform random insertion of $n$ objects

**Claim:** $S(n)$ is $O(n \log n)$
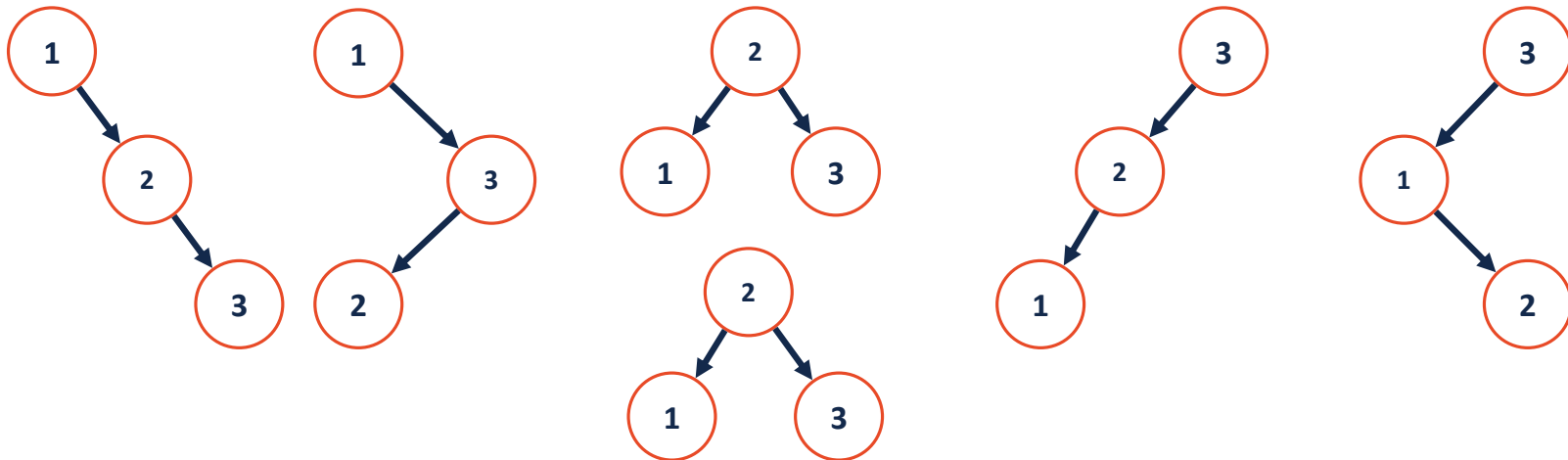
**N=0:**                                  **N=1:**

# Average-Case Analysis: BST

Let $S(n)$ be the average **total internal path length** over all BSTs that can be constructed by uniform random insertion of $n$ objects

**N=3:**

# Average-Case Analysis: BST

Let $S(n)$ be the average **total internal path length** over all BSTs that can be constructed by uniform random insertion of $n$ objects

IH for all $0 \leq k < n$ $S(k)$ is $O(k \log k)$

# Average-Case Analysis: BST

Let $S(n)$ be the average **total internal path length** over all BSTs that can be constructed by uniform random insertion of $n$ objects

Let $0 \leq i \leq n - 1$ be the number of nodes in the left subtree.

Then for a fixed $i$, $S(n) = (n - 1) + S(i) + S(n - i - 1)$

# Average-Case Analysis: BST

Let $S(n)$ be the **average** total internal path length **over all BSTs** that can be constructed by uniform random insertion of $n$ objects

$$S(n) = (n - 1) + \frac{1}{n} \sum_{i=1}^{n-1} S(i) + S(n - i - 1)$$

## Average-Case Analysis: BST

$$S(n) = (n-1) + \frac{2}{n}\sum_{i=1}^{n-1} S(i)$$

$$S(n) = (n-1) + \frac{2}{n}\sum_{i=1}^{n-1} (ci \ln i)$$

$$S(n) \leq (n-1) + \frac{2}{n}\int_{1}^{n} (cx \ln x)\,dx$$

$$S(n) \leq (n-1) + \frac{2}{n}\left(\frac{cn^2}{2}\ln n - \frac{cn^2}{4} + \frac{c}{4}\right) \approx cn \ln n$$

# Average-Case Analysis: BST

Let $S(n)$ be the average **total internal path length** over all BSTs that can be constructed by uniform random insertion of $n$ objects Since $S(n)$ is $O(n \log n)$, if we assume we are randomly choosing a node to insert, find, or delete* then each operation takes:

# Average-Case Analysis: BST

**Summary:** All operations are on average $O(logn)$

**Randomness:**

**Assumptions:**

# Expectation Analysis: Randomized Quicksort

| 6 | 1 | 0 | 3 | 7 | 9 | 2 | 4 |
|---|---|---|---|---|---|---|---|

| 1 | 0 | 3 | 2 | 4 | 9 | 6 | 7 |
|---|---|---|---|---|---|---|---|

| 1 | 0 | 3 | 2 | 4 | 9 | 6 | 7 |
|---|---|---|---|---|---|---|---|

| 1 | 0 | 2 | 3 | 4 | 6 | 7 | 9 |
|---|---|---|---|---|---|---|---|

| 1 | 0 | 2 | 3 | 4 | 6 | 7 | 9 |
|---|---|---|---|---|---|---|---|

| 0 | 1 | 2 | 3 | 4 | 6 | 7 | 9 |
|---|---|---|---|---|---|---|---|

# Expectation Analysis: Randomized Quicksort

# Expectation Analysis: Randomized Quicksort

In **randomized quicksort**, the selection of the pivot is random.

**Claim:** The expected time is $O(n \log n)$ *for any input!*

# Expectation Analysis: Randomized Quicksort

In **randomized quicksort**, the selection of the pivot is random.

**Claim:** The expected time is $O(n \log n)$ *for any input!*

Let $X$ be the total comparisons and $X_{ij}$ be an **indicator variable**:

$$X_{ij} = \left\{ \begin{array}{l} 1 \text{ if } ith \text{ object compared to } jth \\ 0 \text{ if } i \text{ object not compared to } jth \end{array} \right.$$

Then…

# Expectation Analysis: Randomized Quicksort

**Claim:** $E[X_{i,j}] = \frac{2}{j-i+1}$.

**Base Case:** (N=2)

# Expectation Analysis: Randomized Quicksort

**Claim:** $E[X_{i,j}] = \dfrac{2}{j-i+1}$   **Induction:** Assume true for all inputs of $< n$

# Expectation Analysis: Randomized Quicksort

$$E[X] = \sum_{i=0}^{n-1} \sum_{j=i+1}^{n-1} E[X_{ij}] \qquad E[X_{ij}] = \frac{2}{j-i+1}$$

# Expectation Analysis: Randomized Quicksort

$$E[X] = \sum_{i=1}^{n} \sum_{j=i+1}^{n} E[X_{ij}] \quad E[X_{ij}] = \frac{2}{j-i+1}$$

$$E[X] = \sum_{i=1}^{n} 2(\frac{1}{2} + \frac{1}{3} + \ldots + \frac{1}{n-i+1})$$

$$E[X] = \sum_{i=1}^{n} 2(H_{n-1} - 1) \leq 2n \cdot H_n \leq 2n \ln n$$

# Expectation Analysis: Randomized Quicksort

**Summary:** Randomized quick sort is $O(nlogn)$ regardless of input

**Randomness:**

**Assumptions:**