# COMPSCI 614: Randomized Algorithms with Applications to Data Science

Prof. Cameron Musco

University of Massachusetts Amherst. Spring 2024.
Lecture 6

- Problem Set 1 is due tomorrow at midnight.
- I am holding office hours directly after class today.
- No class or office hours on Thursday.
- Problem Set 2 will be posted later this week.

## Summary

**Last Time:**

- Stronger concentration bounds for sums of independent random variables. I.e., exponential concentration bounds.

- Chernoff and Bernstein bound.

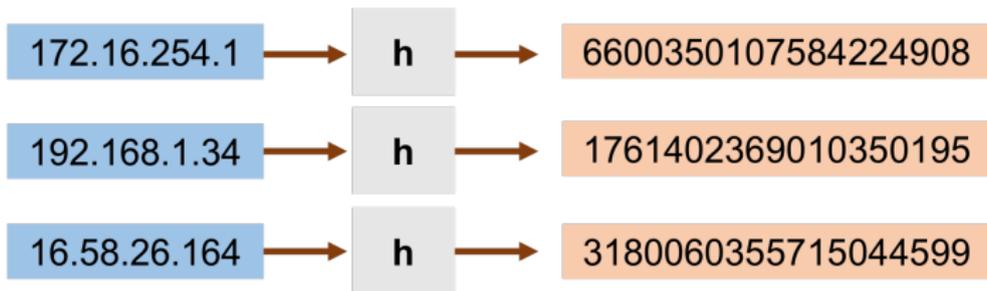- Application to balls-into-bins and linear probing analysis.

**Today:**

- Random hash functions and fingerprinting.

- Applications to pattern matching and communication complexity.

# Random Hashing and Fingerprinting

A random hash function maps inputs to random outputs.

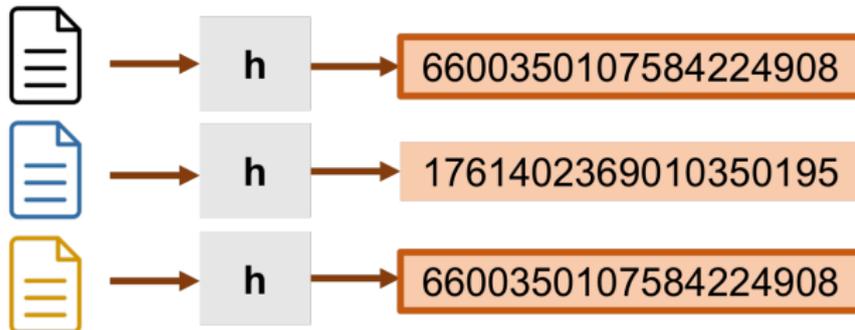| | | |
|---|---|---|
| 172.16.254.1 | **h** | 6600350107584224908 |
| 192.168.1.34 | **h** | 1761402369010350195 |
| 16.58.26.164 | **h** | 3180060355715044599 |

h is picked randomly, but after it is picked it is fixed – so a single input is always mapped to the same output.

```
import random
a = random.randint(1,100)
b = random.randint(1,100)
def myHash(x):
  return (a*x+b) % 100
```

```
import random
def myHash(x):
  a = random.randint(1,100)
  b = random.randint(1,100)
  return (a*x+b) % 100
```

Random hash functions are often used to reduce large files down to hash 'fingerprints', which can be used to check equality of files (deduplication), detect updates/corruptions, etc.



- Key requirement is that two distinct files are unlikely to have the same hash – low collision probability.
- In practice $h$ is often a deterministic 'cryptographic' hash function like SHA or MD5 – hard to analyze formally.

# Rabin Fingerprint

**Rabin Fingerprint:** Interpret a bit string $x_1, x_2, \ldots, x_n$ as the binary representation of the integer $x = \sum_{i=1}^{n} x_i \cdot 2^{i-1}$. Let

$$h(x) = x \mod p,$$

where $p$ is a randomly chosen prime in $[1, tn \log tn]$.

**Prime Number Theorem:** There are $\approx \frac{tn \log tn}{\log(tn \log tn)} = \Theta(tn)$ primes in $[1, tn \log tn]$. So $p$ is chosen randomly from $\Theta(tn)$ possible values.

**Claim:** For $x, y \in [0, 2^n]$ with $x \neq y$, $\Pr[h(x) = h(y))] = O(1/t)$.

- If $h(x) = h(y)$, then it must be that $x - y \mod p = 0$. I.e., $p$ divides $x - y$. So we must bound the probability of this occuring.
- **Note:** This is not a cryptographic hash function – it is relatively easy to find $x, y$ with $h(x) = h(y)$ given $p$, or blackbox access to $h$. However, this is fine in many applications.
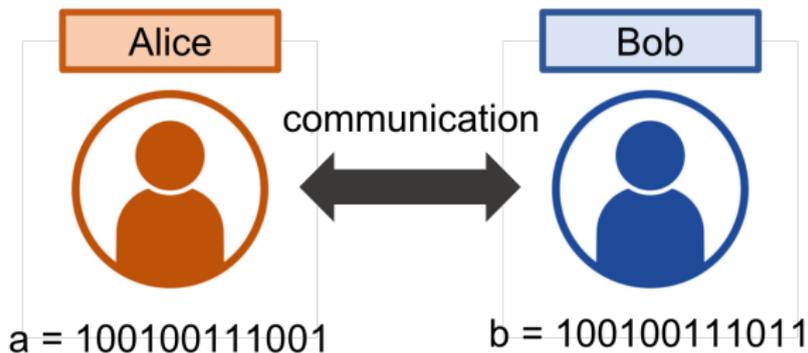
**Think-Pair-Share 1:** How many unique prime factors can an integer in $[-2^n, 2^n]$ have?

**Think-Pair-Share 2:** What is the probability that a random prime $p$ chosen from $[1, tn \log tn]$ divides $x - y \in [-2^n, 2^n]$? I.e., that $h(x) = h(y)$? Recall: There are $\Theta(tn)$ primes in the range $[1, tn \log tn]$.

# Fingerprinting Application 1: Communication Complexity

**Equality Testing Communication Problem:** Alice has some bit string $a \in \{0, 1\}^n$. Bob has some string $b \in \{0, 1\}^n$. How many bits do they need to communicate to determine if $a = b$ with probability at least 2/3?

**Equality Testing Protocol:**

- Alice picks a random prime $p \in [1, tn \log tn]$ for some large constant $t$.
- Alice sends $p$, along with the Rabin fingerprint $\mathsf{h}(a) := a \bmod p$ to Bob. $[O(\log p) = O(\log n)$ bits$]$
- Bob uses $p$ to compute $\mathsf{h}(b) := b \bmod p$.
- If $\mathsf{h}(a) = \mathsf{h}(b)$, Bob sends 'YES' to Alice. Else, he sends 'No'. [1 bit]

**Correctness:** If $a = b$ both Alice and Bob always output 'YES'. If $a \neq b$ they output 'NO' with probability $1 - O(1/t) \geq 2/3$ if $t$ is set large enough.

**Complexity:** Uses just $O(\log p) = O(\log n)$ bits of communication in total.

How many bits must Alice and Bob send if they want to check equality of $a, b \in \{0,1\}^n$ without using randomness?
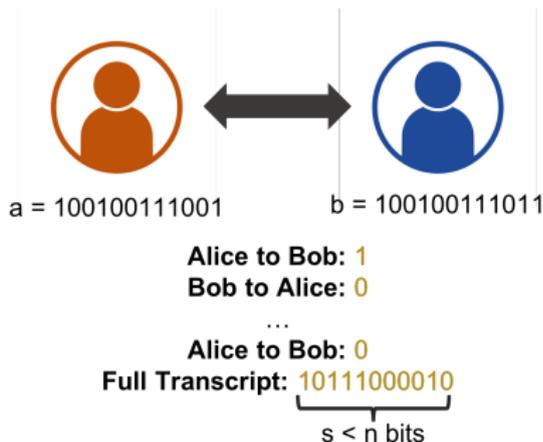
**Claim:** Any deterministic protocol for equality testing requires sending $\Omega(n)$ bits.

- An exponential separation between randomized and deterministic protocols!
- Unlike for running times, for communication complexity problems there are often large provable separations between randomized and deterministic protocols.
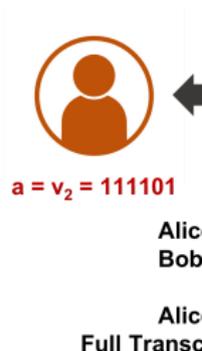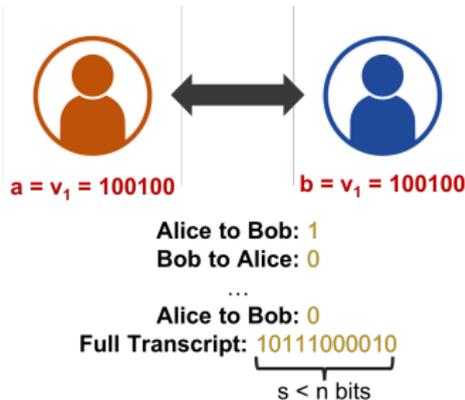
# Deterministic Equality Testing Lower Bound

**Claim:** Any deterministic protocol for equality testing requires sending $\Omega(n)$ bits.

- Assume without loss of generality that Alice and Bob alternate sending 1 bit at a time – at most doubles the number of bits.

- If Alice and Bob send $s < n$ bits, in total, there are $2^s$ possible conversations they may have.



a = 100100111001    b = 100100111011

**Alice to Bob:** 1
**Bob to Alice:** 0
...
**Alice to Bob:** 0
**Full Transcript:** 10111000010
$s < n$ bits

If Alice and Bob send $s < n$ bits, in total, there are $2^s$ possible conversations they may have.



a = 100100111001    b = 100100111011

**Alice to Bob: 1**
**Bob to Alice: 0**
...
**Alice to Bob: 0**
**Full Transcript:** 10111000010
s < n bits

$a = v_1 = 100100$    $b = v_1 = 100100$

**Alice to Bob: 1**
**Bob to Alice: 0**
...
**Alice to Bob: 0**
**Full Transcript:** 10111000010
s < n bits

$a = v_2 = 111101$

**Alic**
**Bob**

**Alic**
**Full Transc**

- Since there are $2^n > 2^s$ possible inputs, there must be two different inputs $v_1 \neq v_2$, such that given $a = b = v_1$ or $a = b = v_2$, the protocol outputs 'YES' and has identical transcripts.

- But then the players will send the same messages and output 'YES' also when Alice is given $a = v_1$ and Bob is given $b = v_2$. This violates correctness!

12

Application 2: Pattern Matching

Given some document $x = x_1 x_2 \ldots x_n$ and a pattern $y = y_1 y_2 \ldots y_m$, find some $j$ such that

$$x_j x_{j+1}, \ldots, x_{j+m-1} = y_1 y_2 \ldots y_m.$$

x = The quick brown **fox** jumped across the pond…

y = fox

Can assume without loss of generality that the strings are binary strings.

What is the 'naive' running time required to solve this problem?

## Rolling Hash

We will use the fact that the Rabin fingerprint is a rolling hash.

- Letting $X_j = \sum_{i=0}^{m-1} x_{j+i} \cdot 2^{m-1-i}$ be the integer value represented by the binary string $x_j x_{j+1}, \ldots, x_{j+m-1}$, we have

$$X_{j+1} = 2 \cdot X_j - 2^m x_j + x_{j+m}.$$

- Thus, since for any $X$, $h(X) = X \mod p$,

$$h(X_{j+1}) = 2 \cdot h(X_j) - 2^m x_j + x_{j+m} \mod p.$$

- Given $h(X_j)$, this hash value can be computed using just $O(1)$ arithmetic operations.

## Rabin-Karp Algorithm

The Rabin-Karp pattern matching algorithm is then:

- Pick a random prime $p \in [1, tm \log mt]$, for $t = cn$.
- Let $Y = \mathsf{h}(y)$ be the Rabin fingerprint of the pattern.
- Let $H = \mathsf{h}(X_1)$ be the Rabin fingerprint of the first block of text.
- For $j = 1, \ldots, x_{n-m+1}$
  - If $Y == H$, return $j$.
  - Else, $H = 2 \cdot H - 2^m x_j + x_{j+m} \mod p$.

**Runtime:** Takes $O(m + n)$ time in total. $O(m)$ for the initial hash computations, and $O(1)$ for each iteration of the for loop.

**Correctness:** The probability of a false positive at any step is upper bounded by $\frac{1}{t} = \frac{1}{cn}$. Thus, via a union bound, the probably of a false positive overall is at most $\frac{n}{cn} = \frac{1}{c}$.

# Questions on Random Hashing?

Interesting topics I am not covering:

- Constructions of universal hash functions.
- Constructions of $k$-wise independent hash functions.
- Concentration bounds and hash table analysis using $k$-wise independent hash functions. See Lectures 3-4 of Jelani Nelson's course notes for some material on this (link on schedule page).
- Connections to pseudorandom number generators (PRGs).