

# Fingerprinting & String Matching

In today's lecture, we will talk about randomization and hashing in a slightly different way. In particular, we use arithmetic modulo prime numbers to probabilistically check if two strings are equal to each other. Building on that, we will get a randomized algorithm (called the *Karp-Rabin fingerprinting scheme*) for checking if a long text  $T$  contains a certain pattern string  $P$  as a substring. This technique is surprisingly elegant and extremely extensible!

## Objectives of this lecture

In this lecture, we will:

- Cover some facts about prime numbers that are useful for randomized hashing schemes
- See a new application of hashing to string equality checking
- Look at the Karp-Rabin pattern matching algorithm

## Recommended study resources

- CLRS, *Introduction to Algorithms*, Section 32.2, The Rabin-Karp algorithm
- DPV, *Algorithms*, Section 1.3.1, Generating random primes

## 1 How to Pick a Random Prime

In this lecture, we will often be picking random primes, so let's talk about that. This is used in many widely used algorithms, for example, you do this when generating RSA public/private key pairs. So, how do we actually pick a random prime in some range  $\{1, \dots, M\}$ ? Here's a straightforward approach:

### Algorithm: Random prime generation

- Pick a random integer  $x$  in the range  $\{1, \dots, M\}$ .
- Check if  $x$  is a prime. If so, output it. Else go back to the first step.

Okay this is not quite complete, we have to fill in some details. How would you pick a random number in the prescribed range? Pick a uniformly random bit string of length  $\lceil \log_2 M \rceil + 1$ . (We assume we have access to a source of random bits.) If it represents a number  $\leq M$ , output it,

else repeat. The chance that you will get a number  $\leq M$  is at least half, so in expectation you have to repeat this process at most twice.

How do you check if  $x$  is prime? You can use the Miller-Rabin randomized primality test<sup>1</sup> (which may produce false positives, but it will only output “prime” when the number is composite with very low probability). There are other randomized primality tests as well, see the Wikipedia page. Or you can use the Agrawal-Kayal-Saxena<sup>2</sup> primality test, which has a worse runtime, but is deterministic and hence guaranteed to be correct. We won’t cover those algorithms in this course, so for now, just know that they exist, we can use them, and know that they run in  $O(\text{polylog } M)$  time.

## 2 How Many Primes?

You have probably seen a proof that there are infinitely many primes. Here’s a different question that we’ll need for this lecture.

*For positive integer  $n$ , how many primes are there in the set  $\{1, 2, \dots, n\}$ ?*

Let there be  $\pi(n)$  primes between 0 and  $n$ . One of the great theorems of the 20<sup>th</sup> century was the Prime Number theorem:

### *Theorem 1: The prime number theorem*

The prime counting function  $\pi(n)$  satisfies

$$\lim_{n \rightarrow \infty} \frac{\pi(n)}{n/(\ln n)} = 1.$$

And while this is just a limiting statement, an older result of Chebyshev (from 1848) says that

### *Theorem 2: Chebyshev*

For  $n \geq 2$ , the prime counting function  $\pi(n)$  satisfies

$$\pi(n) \geq \frac{7}{8} \frac{n}{\ln n} = (1.262\dots) \frac{n}{\log_2 n} > \frac{n}{\log_2 n}$$

Here are two consequences of this theorem. The first is that a random integer between 1 and  $n$  is a prime number with probability at least  $1/\log_2 n$ . We also get the following fact:

### *Corollary 1: Density of primes*

If we want at least  $k \geq 4$  primes between 1 and  $n$ , it suffices to have  $n \geq 2k \log_2 k$ .

<sup>1</sup>[http://en.wikipedia.org/wiki/Miller-Rabin\\_primality\\_test](http://en.wikipedia.org/wiki/Miller-Rabin_primality_test)

<sup>2</sup>[http://en.wikipedia.org/wiki/AKS\\_primality\\_test](http://en.wikipedia.org/wiki/AKS_primality_test)

*Proof.* Just plugging in to Theorem 2, we get

$$\pi(2k \log_2 k) \geq \frac{2k \log_2 k}{\log_2(2k \log_2 k)} \geq \frac{2k \log_2 k}{\log_2 2 + \log_2 k + \log_2 \log_2 k} \geq k.$$

□

## 2.1 Tighter Bounds

The following even tighter set of bounds were proved by Pierre Dusart in 2010.

### Theorem 3: Dusart

For all  $n \geq 60184$  we have:

$$\frac{n}{\ln n - 1.1} > \pi(n) > \frac{n}{\ln n - 1}$$

Because this is a two-sided bound, it allows us to deduce a lower bound on the number of primes in a range. For example, the number of 9-digit prime numbers (i.e. primes in the range  $[10^8, 10^9 - 1]$ ) is

$$\pi(10^9 - 1) - \pi(10^8 - 1) > \frac{10^9 - 1}{\ln(10^9 - 1) - 1} - \frac{10^8 - 1}{\ln(10^8 - 1) - 1.1} = 44928097.3 \dots$$

From this we can infer that a randomly generated 9 digit number is prime with probability at least 0.049920.... Thus, the random sampling method would take at most 21 iterations in expectation to find a 9-digit prime.

## 3 The String Equality Problem

Here's a simple problem: we're sending a Mars lander. Alice, the captain of the Mars lander, receives an  $N$ -bit string  $x$ . Bob, back at mission control, receives an  $N$ -bit string  $y$ . Alice knows nothing about  $y$ , and Bob knows nothing about  $x$ . They want to check if the two strings are the same, i.e., if  $x = y$ .<sup>3</sup> One way is for Alice to send the entire  $N$ -bit string to Bob. But  $N$  is very large. And communication is super-expensive between the two of them. So sending  $N$  bits will cost a lot. *Can Alice and Bob share less communication and check equality?*

If they want to be 100% sure that  $x = y$ , then one can show that fewer than  $N$  bits of communication between them will not suffice. But suppose we are OK with being correct with probability 0.9999. Formally, we want a way for Alice and Bob to send a message to Bob so that, at the end of the communication:

- If  $x = y$ , then  $\Pr[\text{Bob says equal}] = 1$ .
- If  $x \neq y$ , then  $\Pr[\text{Bob says equal}] \leq \delta$ .

<sup>3</sup>E.g., this could be an update to the lander firmware, and we want to make sure the file did not get corrupted

Here's a protocol that does almost that. We will *hash* the strings using the hash function  $h_p(x) = (x \bmod p)$  for a random prime  $p$ , then check whether the hashes are equal.

**Algorithm: Randomized string equality test**

1. Alice picks a random prime  $p$  from the set  $\{1, 2, \dots, M\}$  for  $M = \lceil (200N) \cdot \log_2(100N) \rceil$ .
2. She sends Bob the prime  $p$ , and also the value  $h_p(x) := (x \bmod p)$ .
3. Bob checks if  $h_p(x) \equiv y \bmod p$ . If so, he says **equal** else he says **not equal**.

For now, let's not worry about where the particular value of  $M$  came from: it will arise naturally. Let's see how this protocol performs.

**Lemma 1**

If  $x = y$ , then Bob always says **equal**.

*Proof.* Indeed, if  $x = y$ , then  $x \bmod p = y \bmod p$ . So Bob's test will always succeed.  $\square$

**Lemma 2**

If  $x \neq y$ , then  $\Pr[\text{Bob says equal}] \leq \frac{1}{100}$ .

*Proof.* Consider  $x$  and  $y$  and  $N$ -bit binary numbers. So  $x, y < 2^N$ . Let  $D = |x - y|$  be their difference. Bob says **equal** only when  $x \bmod p = y \bmod p$ , or equivalently  $(x - y) = 0 \bmod p$ . This means  $p$  divides  $D = |x - y|$ . In words, the random prime  $p$  we picked happened to be a divisor of  $D$ . What are the chances of that? Let's do the math.

The difference  $D$  is a  $N$ -bit integer, so  $D \leq 2^N$ . So  $D$  can be written (uniquely) as  $D = p_1 p_2 \cdots p_k$ , each  $p_i$  being a prime, where some of the primes might repeat<sup>4</sup>. Each prime  $p_i \geq 2$ , so  $D = p_1 p_2 \cdots p_k \geq 2^k$ . Hence  $k \leq N$ : the difference  $D$  has at most  $N$  prime divisors. The probability that the randomly chosen prime  $p$  is one of them is

$$\frac{N}{\text{number of primes in } \{1, 2, \dots, M\}}.$$

We want this to be at most  $1/100$ , i.e., we would like that the number of primes in  $\{1, 2, \dots, M\}$  is at least  $100N$ . But Corollary 1 says that choosing  $M \geq 200N \log_2 100N$  will give us at least  $100N$  primes. Hence

$$\Pr[\text{Bob falsely says equal}] \leq \frac{N}{\text{number of primes in } \{1, 2, \dots, M\}} \leq \frac{N}{100N} \leq \frac{1}{100}.$$

$\square$

<sup>4</sup>This unique prime-factorization theorem is known as the fundamental theorem of arithmetic.

### 3.1 Communication cost

Naïvely, Alice could have sent  $x$  over to Bob. That would take  $N$  bits. Now she sends the prime  $p$ , and  $x \bmod p$ . That's two numbers at most  $M = 200N \log_2 100N$ . The number of bits required for that:

$$2 \log_2 M = 2 \log_2(200N \log_2 100N) = O(\log N).$$

To put this in perspective, suppose  $x$  and  $y$  were two copies of all of Wikipedia. Say that's about 25 billion characters (25 GB of data!). Say 8 bits per character, so  $N \approx 2 \cdot 10^{11}$  bits. With the new approach, Alice sends over  $2 \log_2(200N \log_2 100N) \approx 100$  bits, or 13 bytes of data. That's a lot less communication!

### 3.2 Reducing the Error Probability

If you don't like the probability of error being 1%, here are two ways to reduce it.

**Approach #1** Have Alice repeat this process multiple times independently with different random primes, with Bob saying **equal** if and only if in all repetitions, the test passes. For example, for 5 repetitions, the chance that he will make an error (i.e., say **equal** when  $x \neq y$ ) is only

$$(1/100)^5 = 10^{-10}$$

That's a 99.999999999% chance of success! In general, if we repeat  $R$  times, we get the probability of error is at most  $(1/100)^R$ , so if we desire an error probability of  $\delta$ , we should do  $R = \log_{100}(1/\delta)$  repetitions. Since each round requires communicating  $O(\log N)$  bits, the total number of bits that Alice must communicate is

$$O(\log(1/\delta) \log N).$$

Can we do better than this?

**Approach #2** Have Alice choose a random prime from a larger set. For some integer  $s \geq 1$ , if we choose  $M = 2 \cdot sN \log_2(sN)$ , then the arguments above show that the number of primes in  $\{1, \dots, M\}$  is at least  $sN$ . And hence the probability of error is  $1/s$ . If we desire an error probability of  $\delta$ , then we must choose  $s = 1/\delta$ . For example, to obtain 99.999% chance of success, we would pick  $s = 1/10^{-6} = 10^6$ . Now Alice is communicating two integers at most  $2 \cdot sN \log_2(sN)$ , so the number of bits is

$$\begin{aligned} 2 \log_2(2 \cdot sN \log_2(sN)) &= 2 \log_2 s + 2 \log_2 N + 2 \log_2(\log_2(sN)) + 2, \\ &= O(\log s + \log N), \\ &= O(\log(1/\delta) + \log N). \end{aligned}$$

This is much better than Approach #1!

## 4 The Karp-Rabin Algorithm (the “Fingerprint” method)

Let’s use this idea to solve a different problem.

### *Problem: Pattern matching*

In the *pattern matching* problem, we are given, over some alphabet  $\Sigma$ ,

- A text  $T$ , of length  $n$ .
- A pattern  $P$ , of length  $m$ .

The goal is to output the locations of all the occurrences of the pattern  $P$  inside the text  $T$ . E.g., if  $T = \text{abracadabra}$  and  $P = \text{ab}$  then the output should be  $\{0, 7\}$ .

There are many ways to solve this problem, but today we will use randomization to solve this problem. This solution is due to Karp and Rabin. The idea is smart but simple, elegant and effective—like in many great algorithms. To simplify the presentation, we will start by assuming that  $\Sigma = \{0, 1\}$ , i.e., all of our strings are written in binary, but the ideas generalize to larger alphabets.

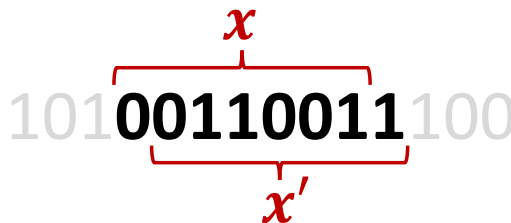
### 4.1 The Karp-Rabin Idea: “Rolling the hash”

As in the last section, interpret the string written in binary *as an integer* and use the hash

$$h_p(x) = (x \bmod p)$$

for some randomly chosen prime  $p$ .

Now look at the string  $x'$  obtained by dropping the leftmost bit of  $x$ , and adding a bit to the right end. E.g., if  $x = 0011001$  then  $x'$  might be  $0110010$  or  $0110011$ . If I told you  $h_p(x) = z$ , can you compute  $h_p(x')$  fast?



Let  $x'_l$  be the lowest-order (rightmost) bit of  $x'$ , and  $x_h$  be the highest order (leftmost) bit of  $x$ . Now observe that

- removing the high-order bit ( $x_h$ ) is just equivalent to subtracting  $x_h \cdot 2^{m-1}$ ,
- shifting all of the remaining bits to one higher position is equivalent to multiplying by 2,
- appending the low-order bit  $x'_l$  is equivalent to just adding  $x'_l$ .

Therefore, we can write

$$x' = 2(x - x_h \cdot 2^{m-1}) + x'_l$$

Since  $h_p(a + b) = (h_p(a) + h_p(b)) \bmod p$ , and  $h_p(2a) = 2h_p(a) \bmod p$ , we then have

$$h_p(x') = (2h_p(x) - x_h \cdot h_p(2^m) + x'_l) \bmod p.$$

Take a moment to understand the significance of this fact. Given the hash  $h_p(x)$  for the substring  $x$  and the value  $h_p(2^m)$ , we can compute the hash of the next adjacent substring  $h_p(x')$  in just a constant number of arithmetic operations modulo  $p$ . This is an enormous speedup compared to computing  $h_p(x')$  from scratch which would take  $O(m)$  arithmetic operations.

## 4.2 The pattern matching algorithm

To keep things short, let  $T_{i\dots j}$  denote the string from the  $i^{\text{th}}$  to the  $j^{\text{th}}$  positions of  $T$ , inclusive. So the string matching problem is: output all the locations  $i \in \{0, 1, \dots, n - m\}$  such that

$$T_{i\dots i+(m-1)} = P.$$

Here's the algorithm.

### Algorithm: Karp-Rabin pattern matching

1. Pick a random prime  $p$  in  $\{1, \dots, M\}$  for  $M = \lceil 2sm \log_2(sm) \rceil$  (we'll choose  $s$  later.)
2. Compute  $h_p(P)$  and  $h_p(2^m)$ , and store these results.
3. Compute  $h_p(T_{0\dots m-1})$ , and check if it equals  $h_p(P)$ . If so, output **match at position 0**.
4. For each  $i \in \{0, \dots, n - m - 1\}$ 
  - (i) compute  $h_p(T_{i+1\dots i+m})$  using  $h_p(T_{i\dots i+m-1})$
  - (ii) If  $h_p(T_{i+1\dots i+m}) = h_p(P)$ , output **match at position  $i + 1$** .

Notice that we'll never have a false negative (i.e., miss a match) but we may erroneously output location that are not matches (have a false positive) if we get a hash collision! Let's analyze the error probability, and the runtime.

**Probability of Error** Since the Karp-Rabin algorithm can encounter false positives, we should analyze the probability of them occurring. This will also show us how large of a prime number of we need to pick in order to achieve a desired false positive probability.

### Theorem: Error probability of Karp-Rabin

We can achieve an error probability of  $\delta$  using the Karp-Rabin algorithm with a prime of  $O(\log(\frac{1}{\delta}) + \log m + \log n)$  bits.

*Proof.* We do  $n - m$  different comparisons, each has a probability  $1/s$  of failure. So, by a union bound, the probability of having at least one false positive is at most  $n/s$ . Hence, setting  $s = 100n$  will make sure we have at most a  $\frac{1}{100}$  chance of even a single mistake.

This means we set  $M = (200 \cdot mn) \log_2(100 \cdot mn)$ , which requires  $\log_2 M + 1 = O(\log m + \log n)$  bits to store. Hence our prime  $p$  is also at most  $O(\log m + \log n)$  bits.<sup>5</sup> Generalizing this, picking a prime from the range  $M = (2/\delta \cdot mn) \log_2(1/\delta \cdot mn)$  achieves a failure probability of  $\delta$  with a prime of  $O(\log(\frac{1}{\delta}) + \log m + \log n)$  bits.  $\square$

**Running Time** Continuing to work in the word-RAM model, note that since the inputs consists of an  $n$ -length string and an  $m$ -length string, our word-RAM must have  $w \geq \log(n)$  and  $w \geq \log(m)$  to be able to index into the input. It is common in randomized algorithms to seek *polynomial* failure probability, i.e., the probability of failure should be proportional to

$$\delta = \frac{1}{O(\text{poly}(n, m))},$$

where the notation  $\text{poly}(n, m)$  means any polynomial expression in  $n$  and  $m$ . Since  $p < M$ , we have that  $\log M = O(\log \text{poly}(n, m)) = O(\log n + \log m)$  and hence all arithmetic on integers mod  $p$  can be done in constant time!

#### *Theorem: Running time of Karp-Rabin*

The Karp-Rabin pattern matching algorithm runs in  $O(m + n)$  time.

*Proof.* Since  $p < M$  and all of our calculations are done mod  $p$ , each individual arithmetic operation takes constant time.

- Computing  $h_p(x)$  for  $m$ -bit  $x$  can be done in  $O(m)$  time. So each of the hash function computations in Steps 2 and 3 take  $O(m)$  time.
- Now, using the idea in Section 4.1, we can compute each subsequent hash value in  $O(1)$  time! So iterating over all the values of  $i$  takes  $O(n)$  time.

That's a total of  $O(m + n)$  time! You can't do much faster, since the input is  $m + n$  bits long. We did not talk about the complexity of picking the prime since we did not cover the complexity of the best primality testing algorithms, but this step can also be shown to run in  $O(\text{polylog}(n, m))$  time, which is dominated by  $O(n + m)$ .  $\square$

<sup>5</sup>If we do the math, and say  $m, n \geq 10$ , then  $\log_2 M \leq 4(\log_2 m + \log_2 n)$ . Now, just for perspective, if we were looking for a  $n = 1024$ -bit phrase in Wikipedia, this means the prime  $p$  is only  $4(\log_2 2^{38} + \log_2 2^{10}) \leq 192$  bits long.



## Exercises: Fingerprinting

**General alphabets** For simplicity, we looked at the case where  $\Sigma = \{0, 1\}$ . The Karp-Rabin algorithm generalizes naturally to larger alphabets. Instead of treating the input as a number in binary, treat it as base- $|\Sigma|$ . For example, if the text contains only lower-case English words, we would use base-26. The formula for rolling the hash still works, except we replace 2 with  $|\Sigma|$ , and the range  $\{1, \dots, M\}$  from which we should select our prime becomes slightly larger.

**Problem 1.** Suppose we use an alphabet  $\Sigma$  which has size  $|\Sigma|$  for Karp-Rabin. What should we use as our new value of  $M$ , and how does this affect the number of bits required to store the prime  $p$ ?

**Other pattern matching algorithms and problems** Though there are many algorithms for pattern matching, the advantage of the Karp-Rabin approach is not only the simplicity, but also the extensibility. You can, for example, solve the following 2-dimensional problem using the same idea.

**Problem 2.** Given a  $(m_1 \times m_2)$ -bit rectangular text  $T$ , and a  $(n_1 \times n_2)$ -bit pattern  $P$  (where  $n_i \leq m_i$ ), find all occurrences of  $P$  inside  $T$ . Show that you can do this in  $O(m_1 m_2)$  time, where we assume that you can do modular arithmetic with integers of value at most  $\text{poly}(m_1 m_2)$  in constant time.