

GEMS OF TCS

DATA STRUCTURES

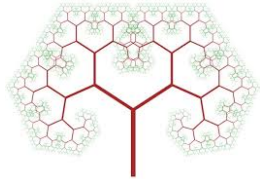
Sasha Golovnev

February 9, 2021

DATA STRUCTURES



Stack, Queue, List, Heap



Search Trees

```
hash(unsigned x) {  
    x ^= x >> (w-m);  
    return (a*x) >> (w-m);  
}
```

Hash Tables

COPING WITH HARD PROBLEMS

- Some problems are too hard to solve exactly

COPING WITH HARD PROBLEMS

- Some problems are too hard to solve exactly
- Approximation

COPING WITH HARD PROBLEMS

- Some problems are too hard to solve exactly
- Approximation
- Randomness

COPING WITH HARD PROBLEMS

- Some problems are too hard to solve exactly
- Approximation
- Randomness
- Today: Preprocessing

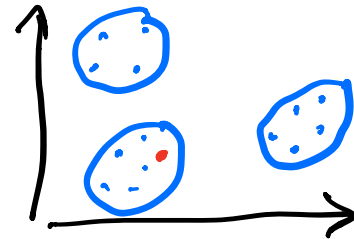
EXAMPLES

- **Graph Distances:** Preprocess a road network in order to efficiently compute distance queries between cities
(Google Maps)

Preprocessing takes forever
Query: Fast

EXAMPLES

- **Graph Distances:** Preprocess a road network in order to efficiently compute distance queries between cities
(Google Maps)
- **Clustering:** Preprocess a set of movies in order to efficiently find closest movie to a query movie
(Netflix recommendations)



DATA STRUCTURES

II Queries

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

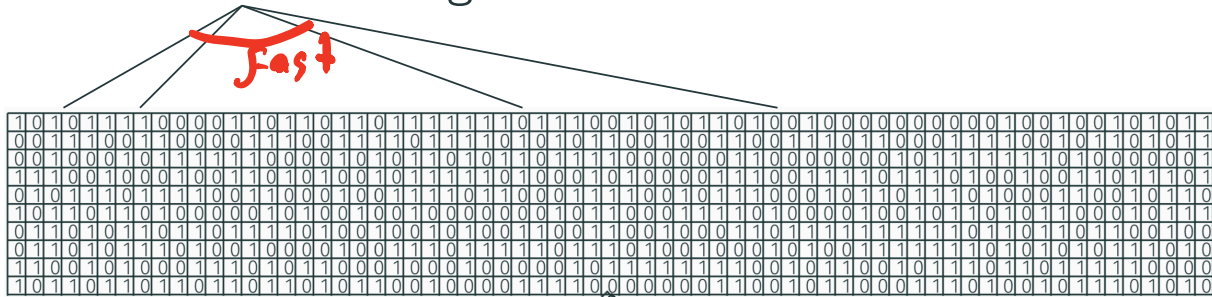
Preprocessing



DATA STRUCTURES

Queries

New York — Washington

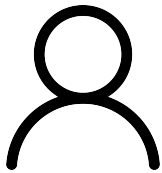


Preprocessing



Stealing Passwords

PASSWORD HASHING



login/pwd



PASSWORD HASHING

haveibeenpwned.com:
Your account has
been compromised



login/pwd



PASSWORD HASHING



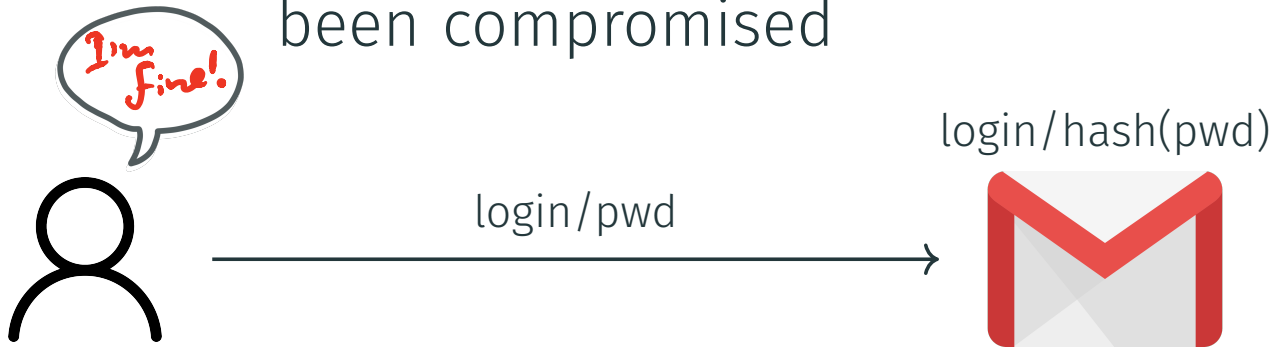
PASSWORD HASHING



hash(qwerty)=1xe4ht
hash(111111)=nh83l0

PASSWORD HASHING

haveibeenpwned.com:
Your account has
been compromised



hash(qwerty)=1xe4ht
hash(111111)=nh83l0

HASHING

- (Cryptographic) hash function maps strings to hash strings such that it's hard to invert ^{pwd}

Ideally, in order to find pwd \rightarrow fixed hash,
you have to brute force all pwds

HASHING

- (Cryptographic) hash function maps strings to strings such that it's hard to invert
- Ideally, to find a password that leads to a fixed hash value, one needs to brute force all possible passwords

HASHING

- (Cryptographic) hash function maps strings to strings such that it's hard to invert
- Ideally, to find a password that leads to a fixed hash value, one needs to brute force all possible passwords
- Hash functions are publicly known (SHA-3)
wiki page

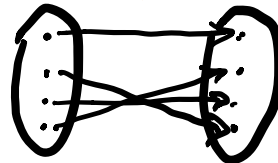
HASHING

- (Cryptographic) hash function maps strings to strings such that it's hard to invert
- Ideally, to find a password that leads to a fixed hash value, one needs to brute force all possible passwords
- Hash functions are publicly known (SHA-3)
- For now, consider hash functions

$f: \{1, \dots, N\} \rightarrow \{1, \dots, N\}$ that are **bijections**

$$f: \{0, 1\}^n \rightarrow \{0, 1\}^n$$

$$N = 2^n$$



INVERTING A BIJECTION

- Let $f: \{1, \dots, N\} \rightarrow \{1, \dots, N\}$ be a **bijection**

Preprocess: years

Query: hash value \rightarrow pwd

given $\gamma \in \{1, \dots, N\}$

find x s.t. $f(x) = \gamma$

2 Naive solution

$$N = 2^{256} \approx 10^{77}$$

I. No preprocessing

$$\text{Space} = O$$

$$\text{Time} = N \approx 10^{77}$$

$< 10^{20}$ operations per second
the age of Universe 10^{15} seconds

II. Preprocess: store

hash \rightarrow pwd

0...0 \rightarrow pwd₁

0...01 \rightarrow pwd₂

⋮

1...11 \rightarrow pwd_N

$$\text{Space} = N \approx 10^{77}$$

$$\text{Time} = \log N$$

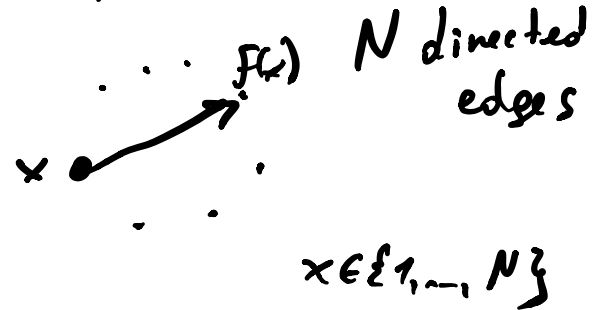
of el. particles in observable

$$\text{Universe} \approx 10^{86}$$

INVERTING A BIJECTION

- Let $f: \{1, \dots, N\} \rightarrow \{1, \dots, N\}$ be a **bijection**
- Invert it in time $T = \sqrt{N}$ and space $S = \sqrt{N}$

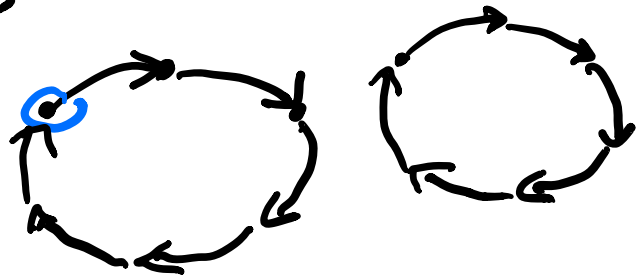
Directed Graph with N vertices



INVERTING A BIJECTION

- Let $f: \{1, \dots, N\} \rightarrow \{1, \dots, N\}$ be a bijection
- Invert it in time $T = \sqrt{N}$ and space $S = \sqrt{N}$
- Let's define a directed graph on N vertices with edges $x \rightarrow f(x)$

$$\begin{aligned} \text{out-degree}(x) &= 1 \\ \text{in-degree}(x) &= 1 \end{aligned}$$

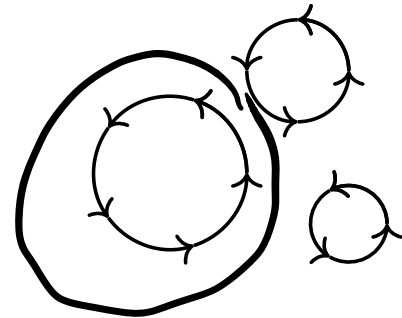


INVERTING A BIJECTION

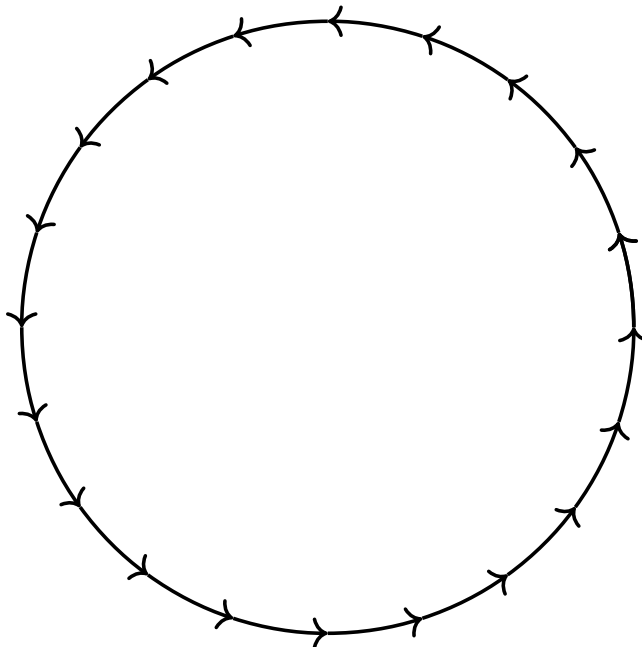
- Let $f: \{1, \dots, N\} \rightarrow \{1, \dots, N\}$ be a **bijection**
- Invert it in time $T = \sqrt{N}$ and space $S = \sqrt{N}$
- Let's define a directed graph on N vertices with edges $x \rightarrow f(x)$
- In- and out-degrees of all vertices are 1

INVERTING A BIJECTION

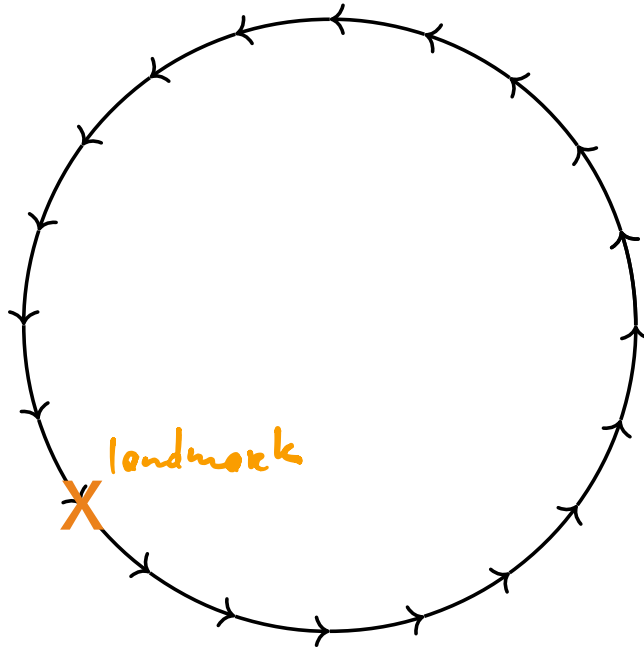
- Let $f: \{1, \dots, N\} \rightarrow \{1, \dots, N\}$ be a **bijection**
- Invert it in time $T = \sqrt{N}$ and space $S = \sqrt{N}$
- Let's define a directed graph on N vertices with edges $x \rightarrow f(x)$
- In- and out-degrees of all vertices are 1
- Thus, this graph is a union of cycles



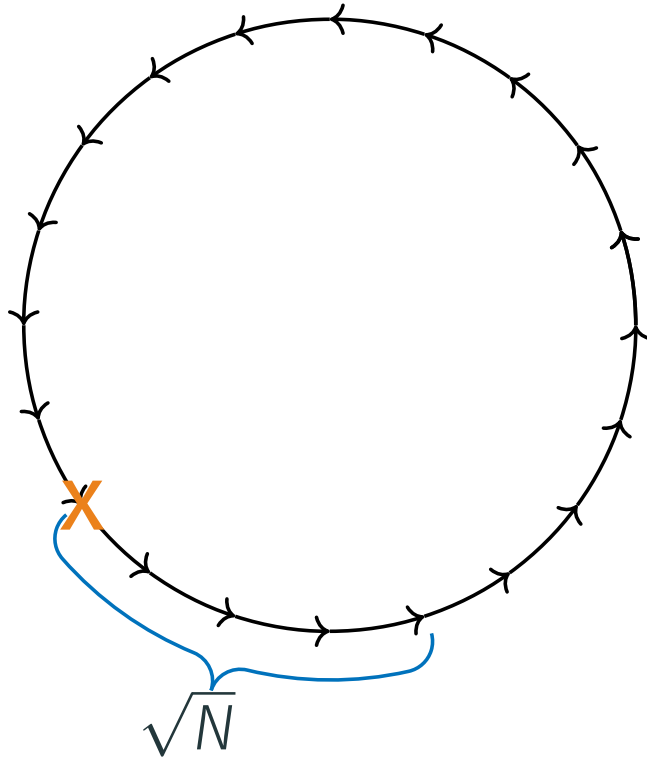
INVERTING A BIJECTION



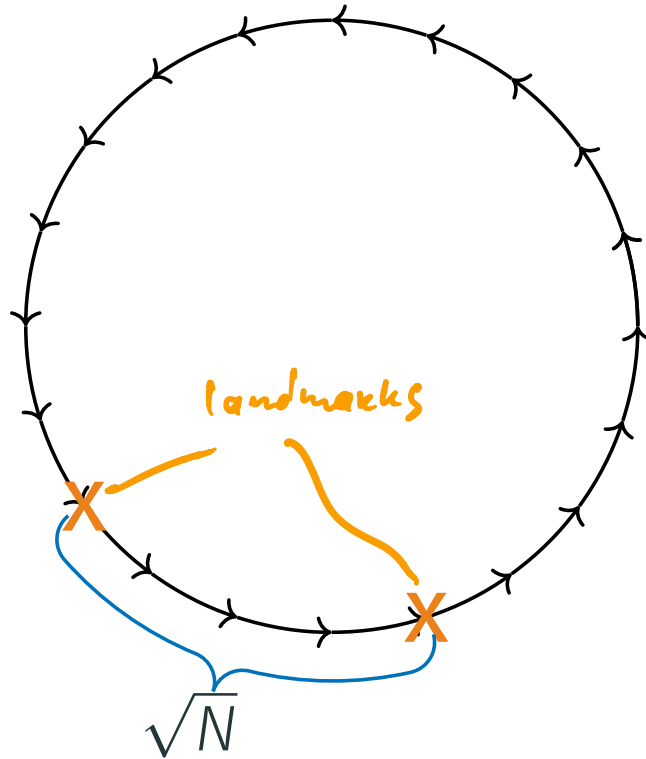
INVERTING A BIJECTION



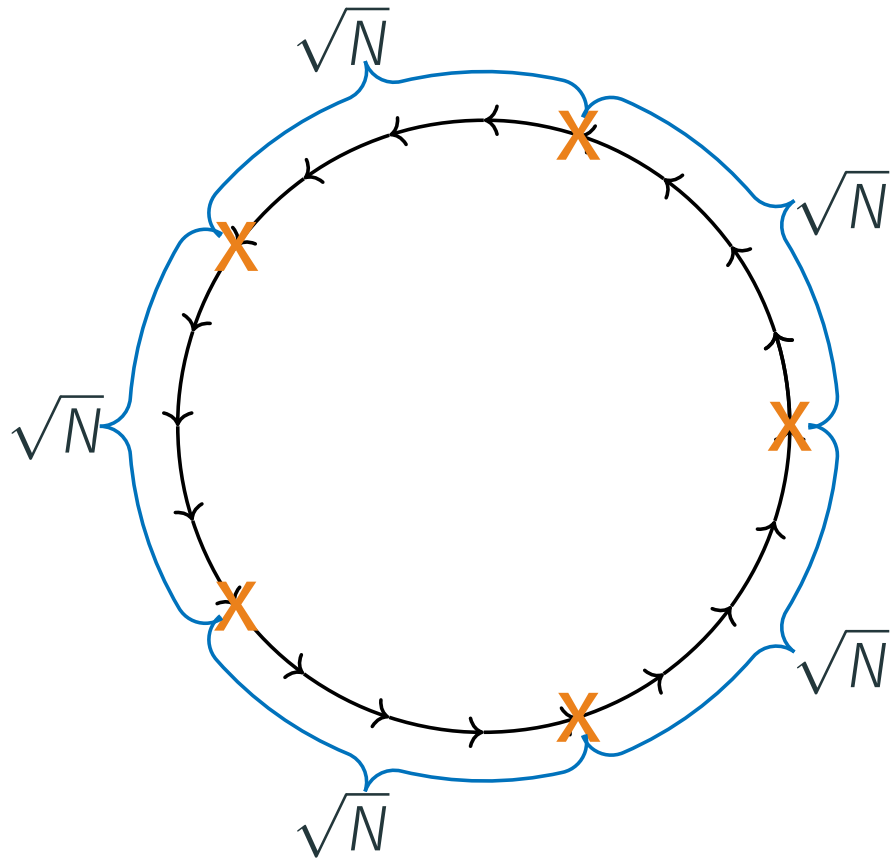
INVERTING A BIJECTION



INVERTING A BIJECTION

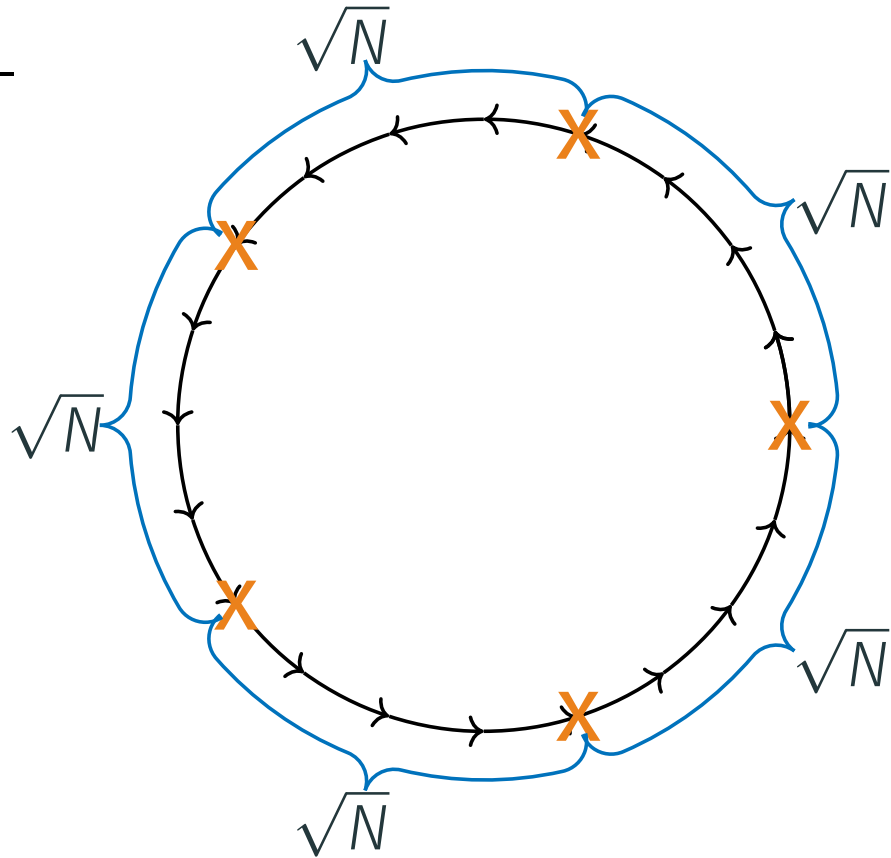


INVERTING A BIJECTION



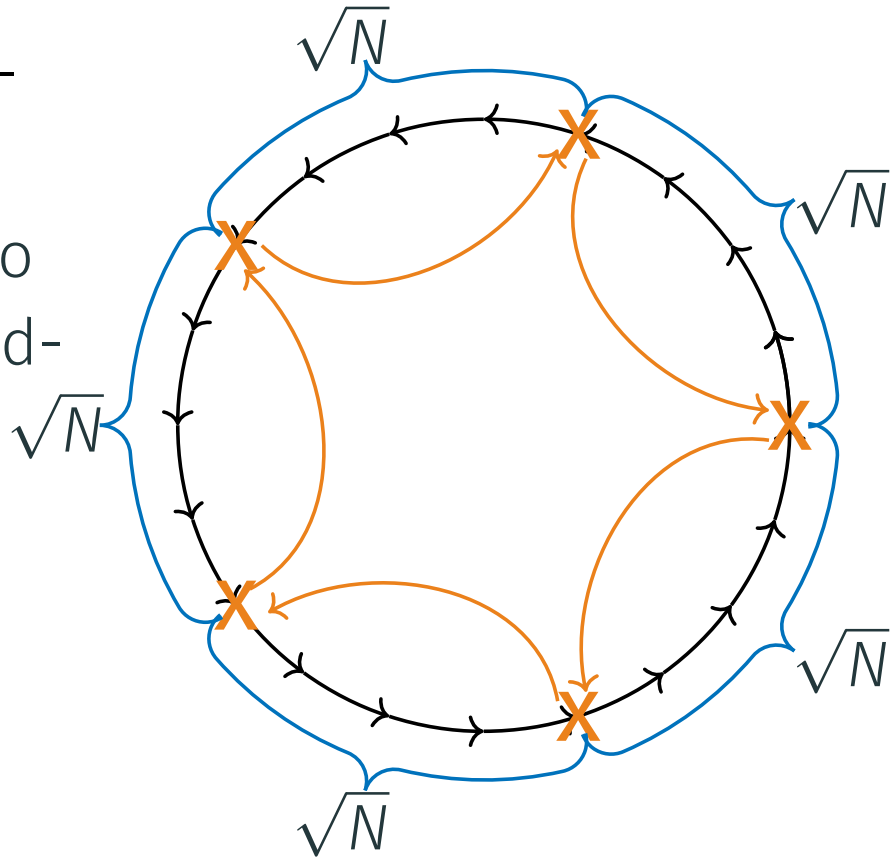
INVERTING A BIJECTION

Store x landmarks,
marks,



INVERTING A BIJECTION

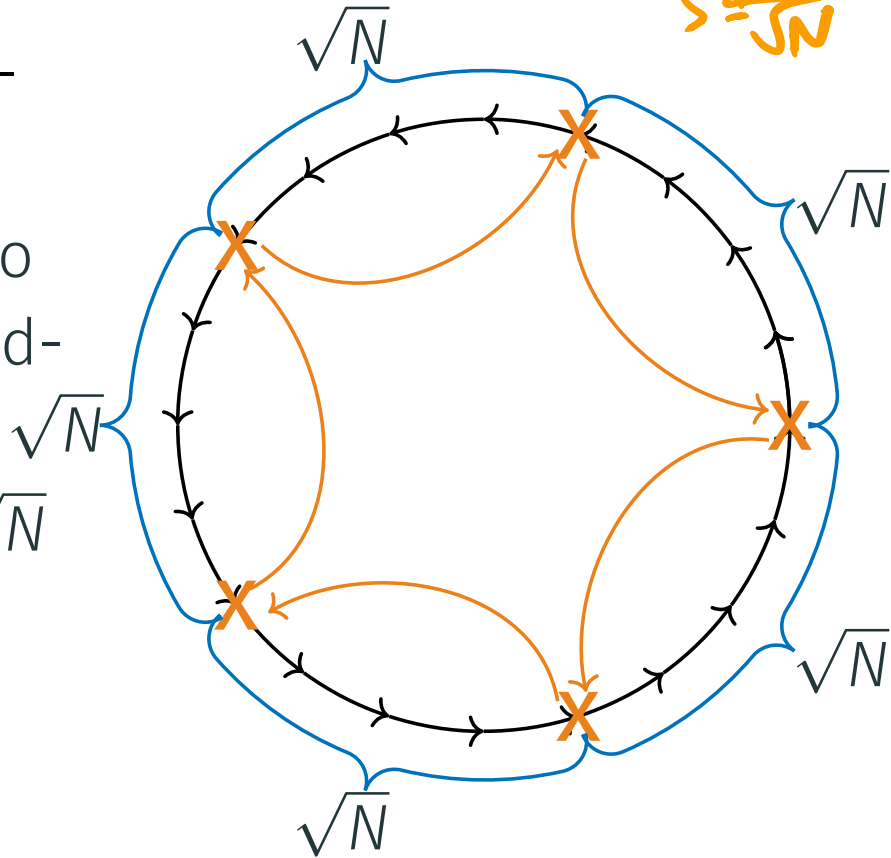
Store x landmarks,
and links \curvearrowright to
previous landmarks



INVERTING A BIJECTION

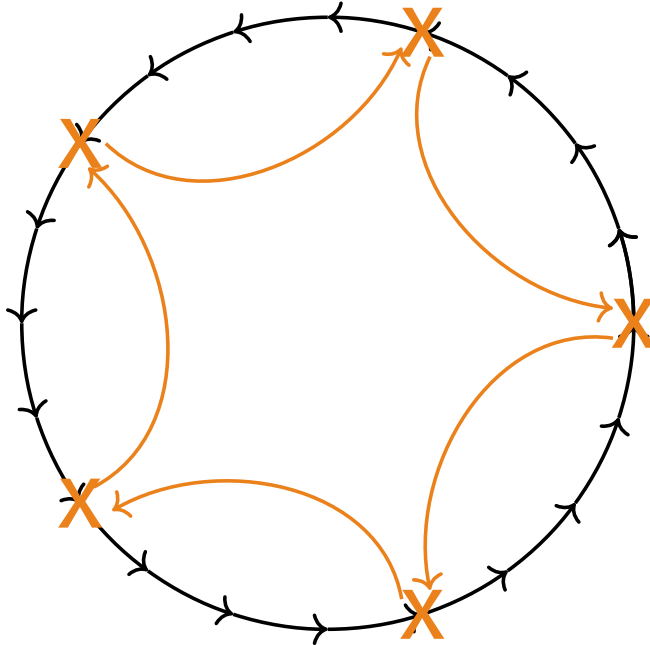
$$S \approx \frac{N}{\sqrt{N}} = \sqrt{N}$$

Store x landmarks,
and links \curvearrowright to
previous landmarks
space $S \approx \sqrt{N}$



INVERTING A BIJECTION

Store x landmarks,
and links \curvearrowright to
previous landmarks
space $S \approx \sqrt{N}$

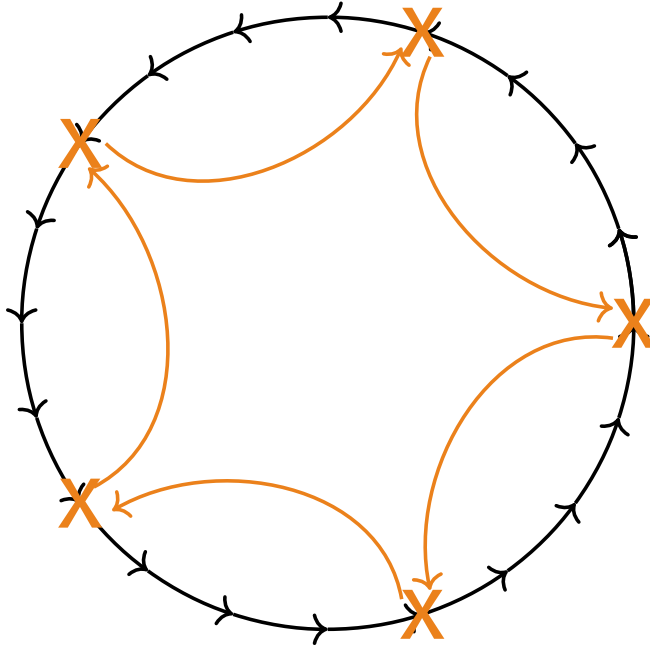


INVERTING A BIJECTION

Store x landmarks,
and links \curvearrowright to
previous landmarks

space $S \approx \sqrt{N}$

time $T \approx \sqrt{N}$:



INVERTING A BIJECTION

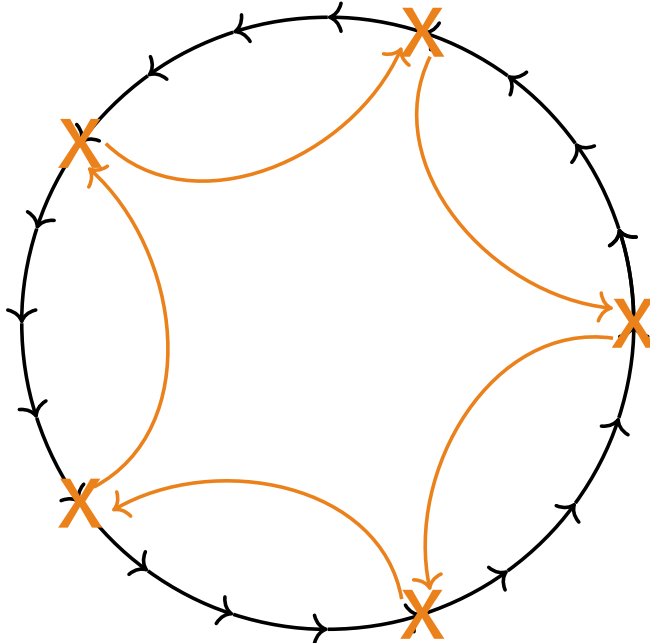
Store x landmarks,
and links \curvearrowright to
previous landmarks

space $S \approx \sqrt{N}$

time $T \approx \sqrt{N}$:

Invert $y = f(x)$

given hash, find pwd



INVERTING A BIJECTION

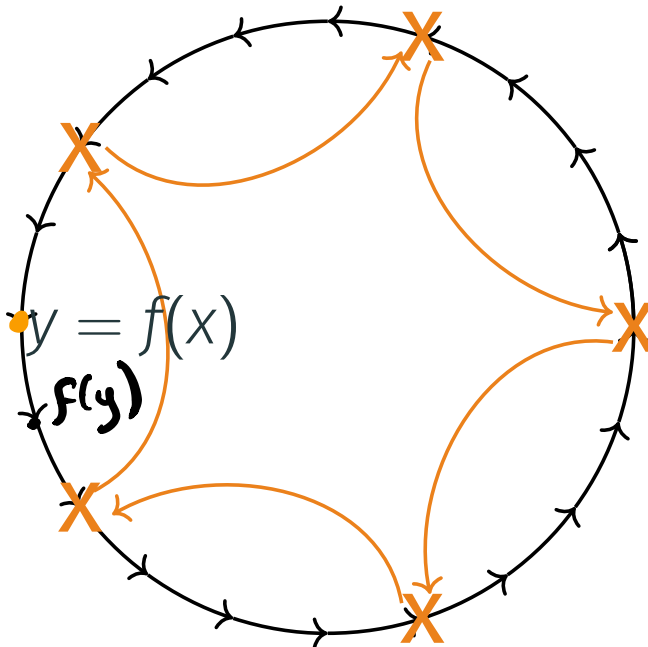
Recall f is publicly known
 $F(y)$ $y \rightarrow f(y)$

Store x landmarks,
and links \curvearrowright to
previous landmarks

space $S \approx \sqrt{N}$

time $T \approx \sqrt{N}$:

Invert $\boxed{y} = f(x)$



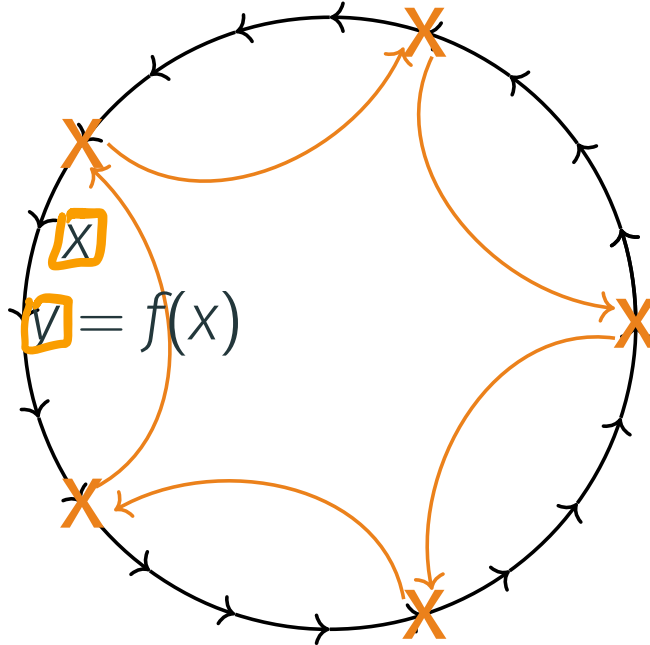
INVERTING A BIJECTION

Store x landmarks,
and links \curvearrowright to
previous landmarks

space $S \approx \sqrt{N}$

time $T \approx \sqrt{N}$:

Invert $y = f(x)$



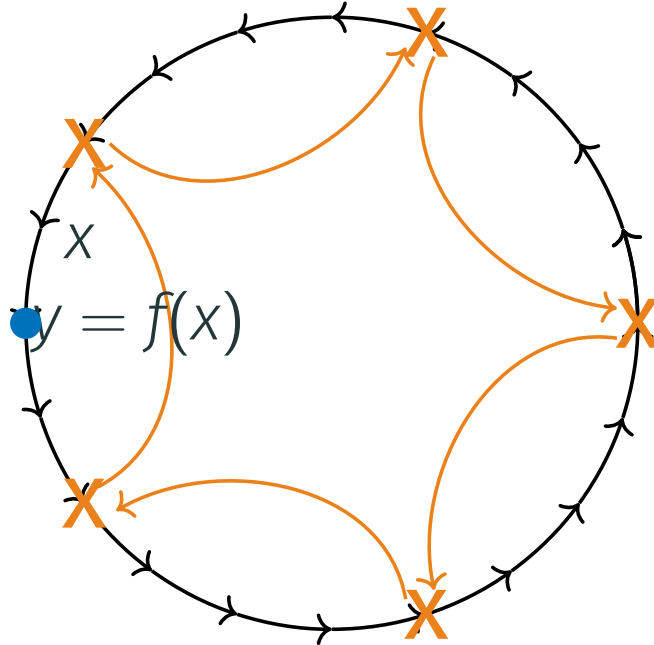
INVERTING A BIJECTION

Store x landmarks,
and links \curvearrowright to
previous landmarks

space $S \approx \sqrt{N}$

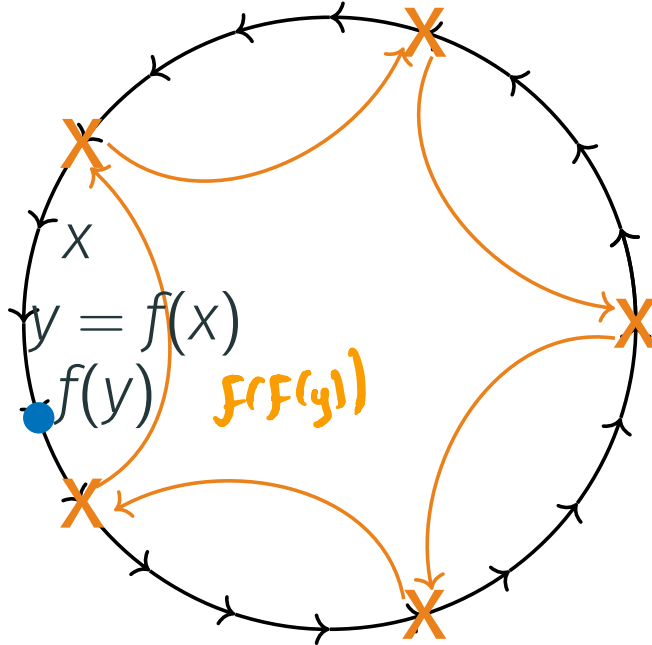
time $T \approx \sqrt{N}$:

Invert $y = f(x)$



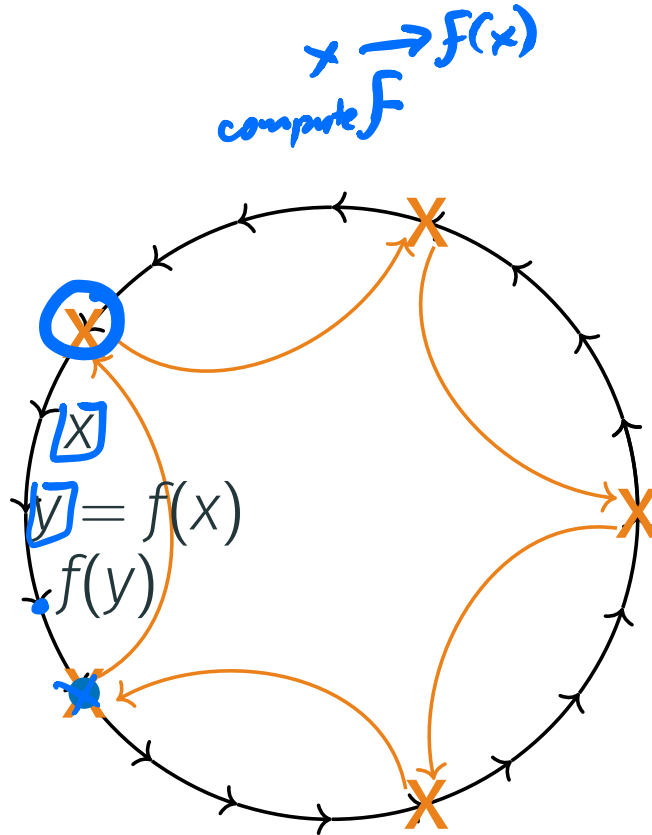
INVERTING A BIJECTION

Store x landmarks,
and links \curvearrowright to
previous landmarks
space $S \approx \sqrt{N}$
time $T \approx \sqrt{N}$:
Invert $y = f(x)$



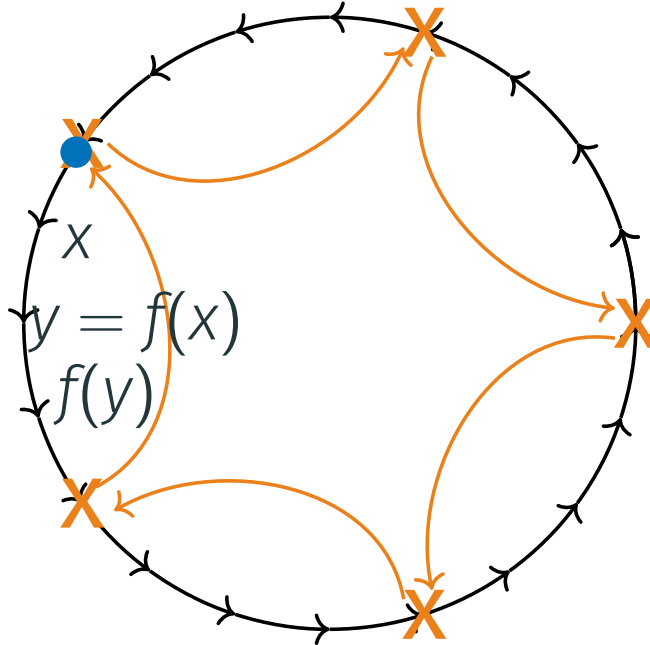
INVERTING A BIJECTION

Store x landmarks,
and links \curvearrowright to
previous landmarks
space $S \approx \sqrt{N}$
time $T \approx \sqrt{N}$:
Invert $y = f(x)$



INVERTING A BIJECTION

Store x landmarks,
and links \curvearrowright to
previous landmarks
space $S \approx \sqrt{N}$
time $T \approx \sqrt{N}$:
Invert $y = f(x)$



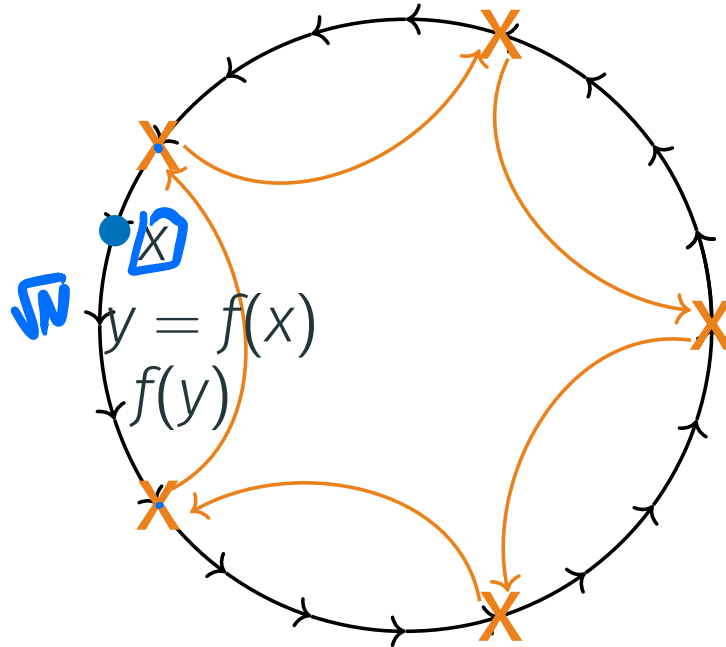
INVERTING A BIJECTION

Store x landmarks,
and links \curvearrowright to
previous landmarks

space $S \approx \sqrt{N}$

time $T \approx \sqrt{N}$:

Invert $y = f(x)$



DATA STRUCTURE

space

query time

- Let $ST = N$

Ex. $S = N^{1/3}$

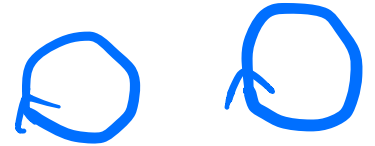
$$T = N^{2/3}$$

DATA STRUCTURE

- Let $ST = N$
- Let's define a directed graph on N vertices with edges $x \rightarrow f(x)$

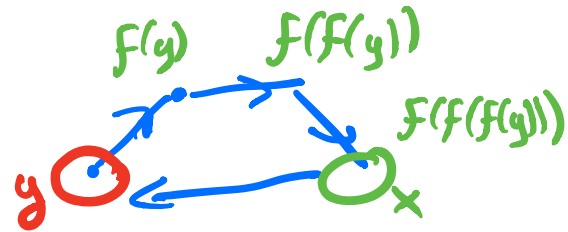
DATA STRUCTURE

- Let $ST = N$
- Let's define a directed graph on N vertices with edges $x \rightarrow f(x)$
- Partition the graph into cycles



DATA STRUCTURE

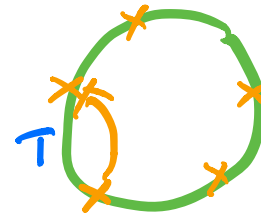
- Let $ST = N$
- Let's define a directed graph on N vertices with edges $x \rightarrow f(x)$
- Partition the graph into cycles
- Ignore cycles of length $\leq T$



DATA STRUCTURE

- Let $ST = N$
- Let's define a directed graph on N vertices with edges $x \rightarrow f(x)$
- Partition the graph into cycles
- Ignore cycles of length $\leq T$
- In all other cycles store every T th vertex as a landmark

$$S \leq N/T$$



DATA STRUCTURE

someone publishes
 \sqrt{N} carefully chosen
hash value

- Let $ST = N$
- Let's define a directed graph on N vertices with edges $x \rightarrow f(x)$
- Partition the graph into cycles
- Ignore cycles of length $\leq T$
- In all other cycles store every T th vertex as a landmark
- Space: S , query time: T

Rainbow tables
attacks

Now use hash functions: big N

Prohibited Passwords


PROHIBITED PASSWORDS

- Check if entered password is in the list of m prohibited passwords

PROHIBITED PASSWORDS

- Check if entered password is in the list of m prohibited passwords
- We can store m strings, check in $\sim \log m$ time

PROHIBITED PASSWORDS

- Check if entered password is in the list of m prohibited passwords
 - We can store m strings, check in $\sim \log m$ time
 - Bloom filters: store $\sim m$ bits, check in $O(1)$ time
- losing information*
- 

PROHIBITED PASSWORDS

- Check if entered password is in the list of m prohibited passwords
- We can store m strings, check in $\sim \log m$ time
- Bloom filters: store $\sim \underline{m \text{ bits}}$, check in $O(1)$ time
- We'll be wrong with small probability

DATA STRUCTURE

- We want a data structure that supports two functions

DATA STRUCTURE

- We want a data structure that supports two functions
 - $\text{Insert}(x)$

prohibited pwd



DATA STRUCTURE

- We want a data structure that supports two

functions

- Insert(x)
- Lookup(x)

*Naive: list of prohibited
pool*

*would use too much
space*

*Lookup would be less
efficient*

DATA STRUCTURE

- We want a data structure that supports two functions
 - Insert(x)
 - Lookup(x)
 - Hashtables: less efficient but don't make mistakes
- we're storing $\approx m$ bits*
- store $\approx m \log m$ bits*

DATA STRUCTURE

- We want a data structure that supports two functions
 - Insert(x) ↙
 - Lookup(x) ↖
- Hashtables: less efficient but don't make mistakes
- Bloom filter will use array of n bits $A[0], \dots, A[n - 1]$, initialized with zeros

DATA STRUCTURE

- We want a data structure that supports two functions
 - $\text{Insert}(x)$
 - $\text{Lookup}(x)$
- Hashtables: less efficient but don't make mistakes
- Bloom filter will use array of n bits $A[0], \dots, A[n - 1]$, initialized with zeros
- We'll use $k = O(1)$ hash functions

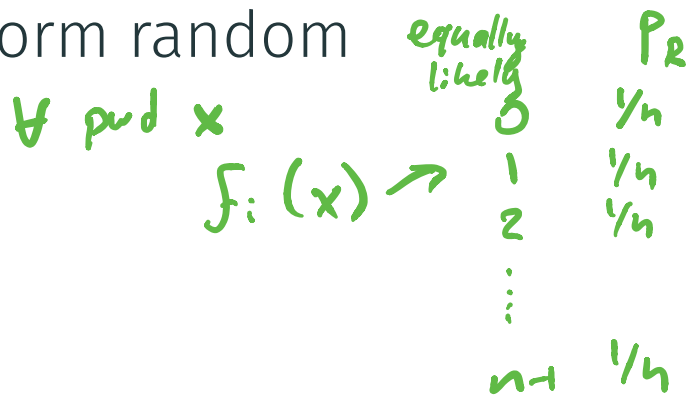
HASH FUNCTIONS

- We have k hash functions f_1, \dots, f_k from strings to $\{0, \dots, n - 1\}$ *-integers*
- prohibited words*
↓

HASH FUNCTIONS

- We have k hash functions f_1, \dots, f_k from strings to $\{0, \dots, n - 1\}$

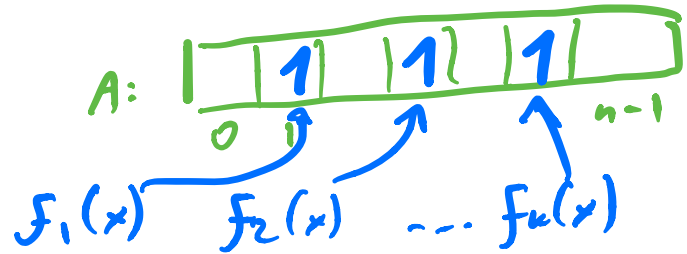
- Assume that functions are independent and uniform random



BLOOM FILTER

- Insert(x):
 - for $i = 1, \dots, k$,
 - $A[f_i(x)] \leftarrow 1$

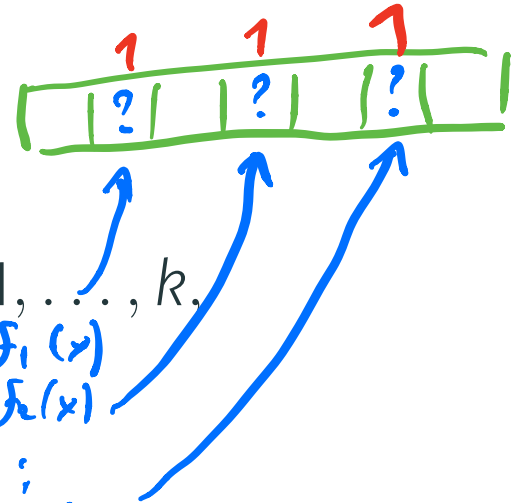
whether or not
bit is already set to 1



BLOOM FILTER

- Insert(x):
 - for $i = 1, \dots, k$,
 - $A[f_i(x)] \leftarrow 1$

- Lookup(x):
 - return 1 iff for every $i = 1, \dots, k$,
 $A[f_i(x)] = 1$



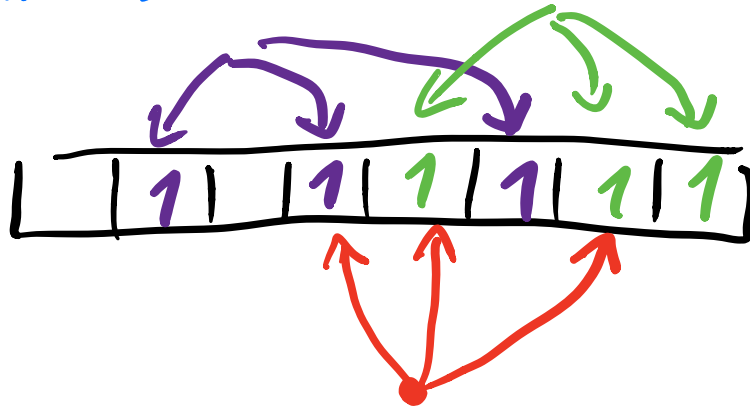
Obs: x was inserted, lookup(x) correctly says x is in data structure

Problem:

ANALYSIS

Problem:

x was not inserted
Lookup(x) may sometimes say that
 x was inserted



Parameters

m - # of prohibited pwds

n - # of bits in array

k - # of hash fns

Remains: analyze Pr of mistake

Analysis to set params n, k .

We've inserted all m strings

$$\Pr[A[0] = 0] = ?$$

Insert 1st pwd:

f_1

$$\Pr[A[0] = 0] = 1 - \frac{1}{n}$$

\vdots

f_k

$$\Pr[A[0] = 0] = 1 - \frac{1}{n}$$

mth

pwd

f_1

\vdots
 f_k

$$\Pr[A[0] = 0] = 1 - \frac{1}{n}$$

Pr $(1 - 1/n)^{m \cdot k}$ $A[0] = 0$ even
 after I inserted
 m els.

$$\text{Pr}[A[1] = 0] = (1 - 1/n)^{m \cdot k}$$

$$e^x = 1 + x + \frac{x^2}{2} + \frac{x^3}{6} + \dots$$

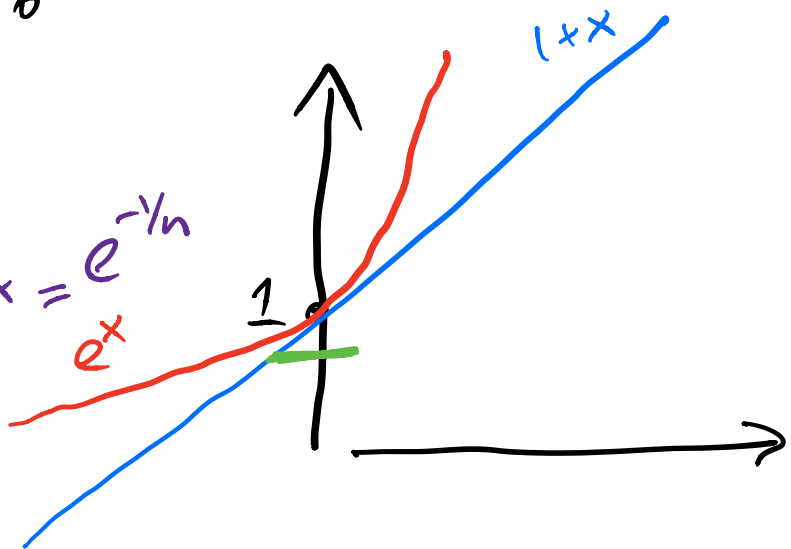
When x is small ($x = \pm 1/n$),

$\frac{x^2}{2} + \frac{x^3}{6} + \dots$ don't matter

$$e^x \approx 1 + x$$

$$x = -1/n$$

$$(1 - 1/n) = (1 + x) \approx e^x = e^{-1/n}$$



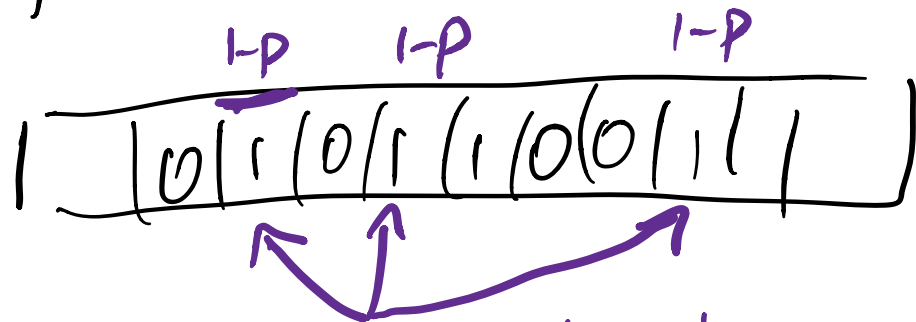
$$Pr[A[0] = 0] = (1 - 1/n)^{m \cdot k}$$

$$(1 - 1/n) \approx e^{-1/n}$$

$$P = (1 - 1/n)^{m \cdot k} \approx (e^{-1/n})^{m \cdot k} =$$

$$= e^{-m \cdot k / n}$$

Pr of error?



all k hash values happen to be ones

$$Pr[\text{error}] = (1 - p)^k$$

Turns out, optimal $p = \frac{1}{2}$

$$\text{Pr}[\text{error}] = (1-p)^k = \boxed{\frac{1}{2^k}}$$

$$\frac{1}{2} = e^{-mk/n}$$

$$\rightarrow \boxed{k} = \frac{n}{m} \ln 2$$

Ex. $n = 8m$ - 8 bits per string

$$k = 8 \cdot \ln 2 \approx 6 \text{ hash func}$$

$$\text{Pr}[\text{error}] \approx 1\%$$

Ex. $n = 32m$ - 32 bits per string

$$k \approx 22 \text{ hash func}$$

$$\text{Pr}[\text{error}] \approx 10^{-7}$$