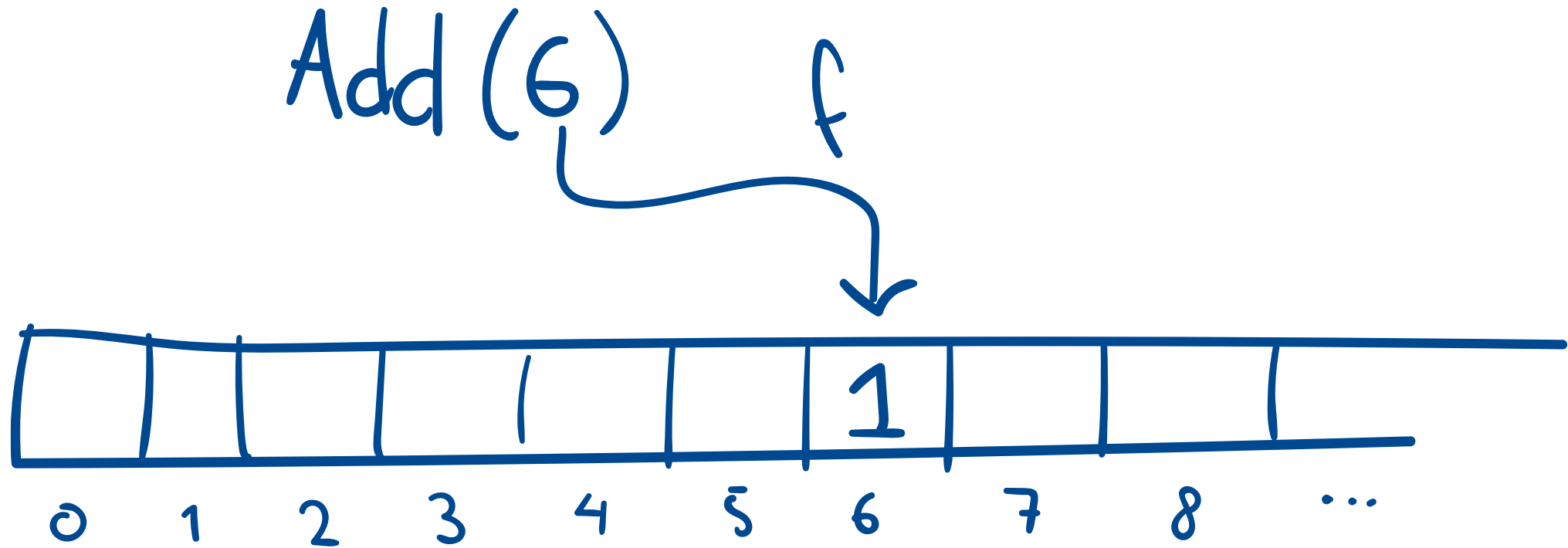# Overview

– Set membership data structures

– Why are false positives acceptable

– A Bloom filter in a few steps

– Bloom filter tricks

– GloBiMaps

# SetMembership

|  | List | Tree | Hash Map |
|---|---|---|---|
| $Add(x:U)$ | $O(1)$ | $O(\log n)$ | $O(1)$ $O(n)$ |
| $Remove(x:U)$ | $O(n)$ | $O(\log n)$ | $O(1)$ $O(n)$ |
| $Test(x:U) : bool$ | $O(n)$ | $O(\log n)$ | $O(1)$ $O(n)$ |
| "Search" | | | |
| Needs: | $=$ | $\leq$ | hash $=$ |

# Bit Set

– Bijection between elements and array of bits

# SetMembership

| | List | Tree | Hash Map | | Bit Set |
|---|---|---|---|---|---|
| $Add(x:U)$ | $O(1)$ | $O(\log n)$ | $O(1)$ | $O(n)$ | $O(1)$ |
| ~~$Remove(x:U)$~~ | ~~$O(n)$~~ | ~~$O(\log n)$~~ | ~~$O(1)$~~ | ~~$O(n)$~~ | ~~$O(1)$~~ |
| $Test(x:U):bool$ | $O(n)$ | $O(\log n)$ | $O(1)$ | $O(n)$ | $O(1)$ |
| allow false positives | Needs: | $=$ | $\leq$ | hash | index |
| | Space | $\vdash$———— | $O(n)$ | ———$\dashv$ | $|U|$ bits |

Also store the elements

# Space/Time Trade-offs in Hash Coding with Allowable Errors

BURTON H. BLOOM
*Computer Usage Company, Newton Upper Falls, Mass.*

In this paper trade-offs among certain computational factors in hash coding are analyzed. The paradigm problem considered is that of testing a series of messages one-by-one for membership in a given set of messages. Two new hash-coding methods are examined and compared with a particular conventional hash-coding method. The computational factors considered are the size of the hash area (space), the time required to identify a message as a nonmember of the given set (reject time), and an allowable error frequency.

In such applications, it is envisaged that overall performance could be improved by using a smaller core resident hash area in conjunction with the new methods and, when necessary, by using some secondary and perhaps time-consuming test to "catch" the small fraction of errors associated with the new methods. An example is discussed which illustrates possible areas of application for the new methods.
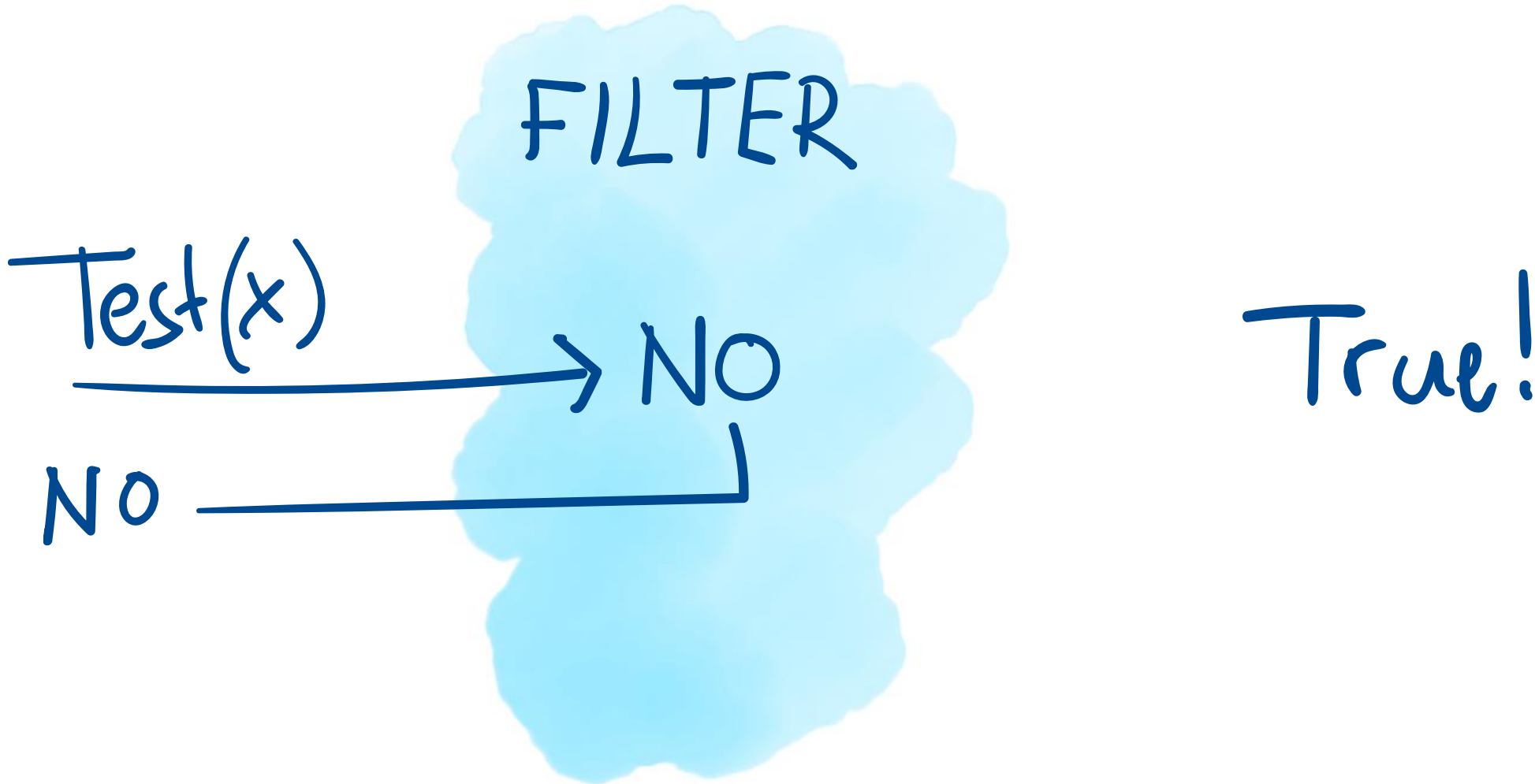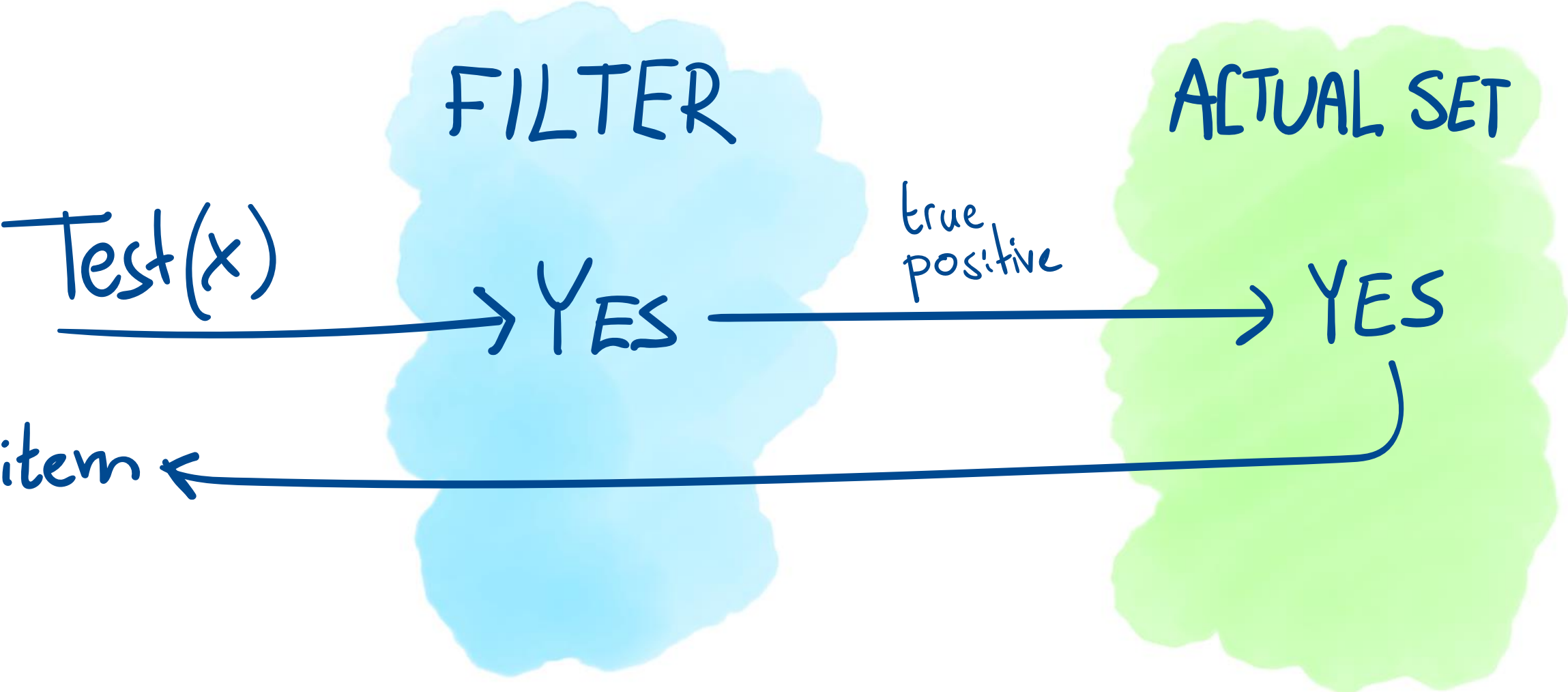
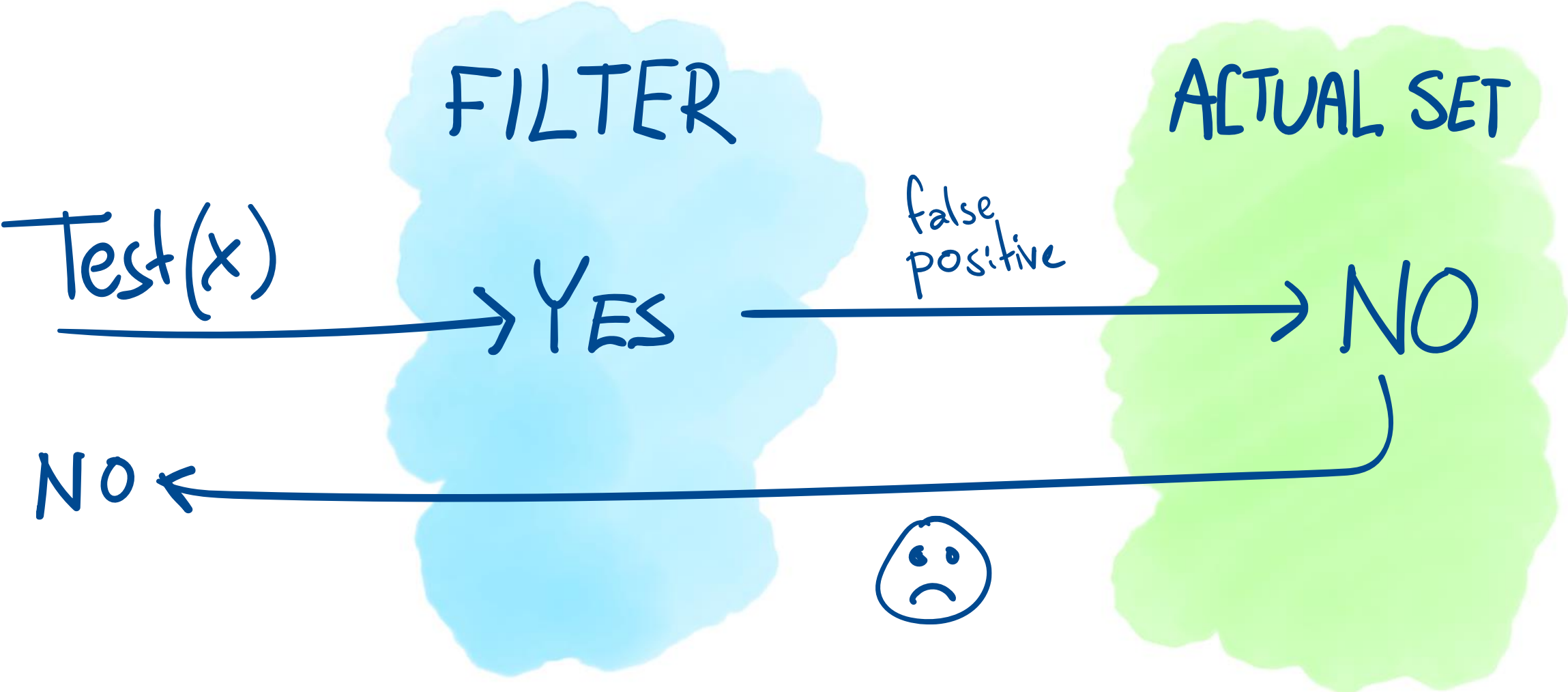# Allow false positives



FILTER

Test(x)

→ NO

NO

True!

# Allow false positives

# Is this a good idea?

- Depends on "Filter" versus "Actual Set" cost
- Depends on hit rates

- If filter can be amazingly small
- Maybe you don't actually have the set!
- Fast rejection most of the time

|  | True | False |
|---|---|---|
| Negative | Filter | Filter |
| Positive | Filter + Set Hit | Filter + Set Miss |

# Applications 1/3

- **Bloom '70**: hyphenation
  - Most words covered by a few rules
  - Make a set containing the exceptions
  - Hypenation algo: check set, else use rules
  - Let's add a filter! False positives?
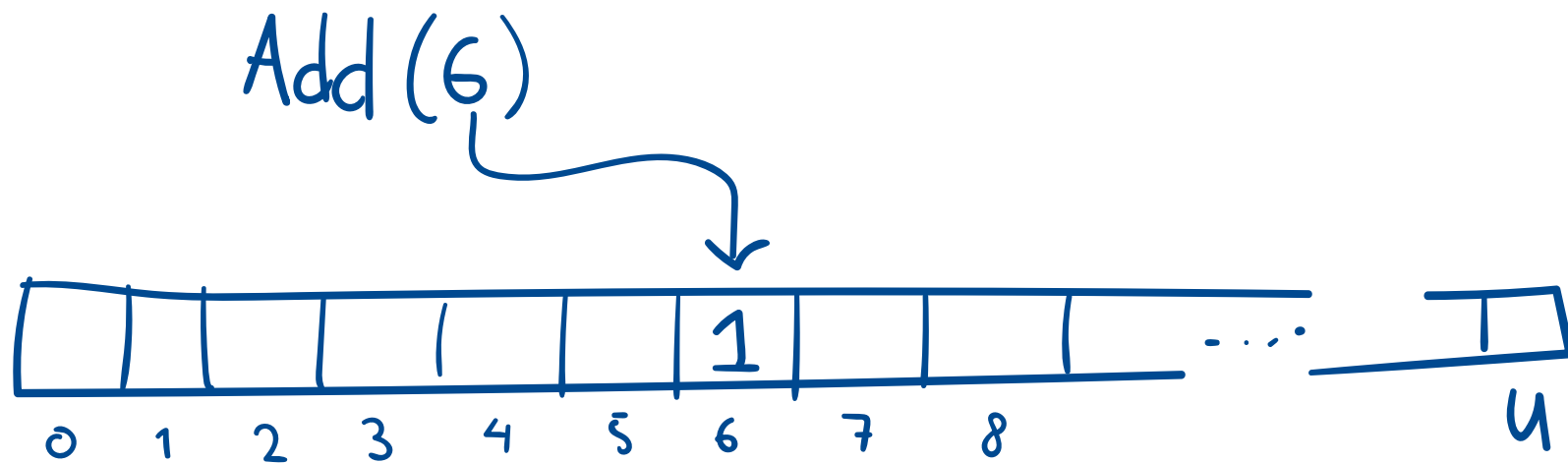    - Unnecessary lookup; still correct

# Applications 2/3

- **MrIlroy '82**: early UNIX spell-checkers
  - Store correct words. False positives?
    - Just accept them 🙁
  - Amazingly small filter!
- **Spafford '92**: unsuitable passwords
  - Store the set. False positives?
    - Not really harmful
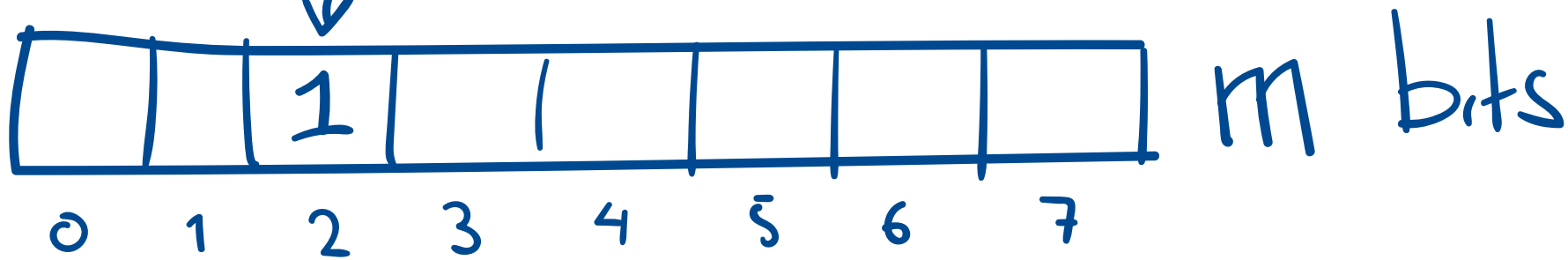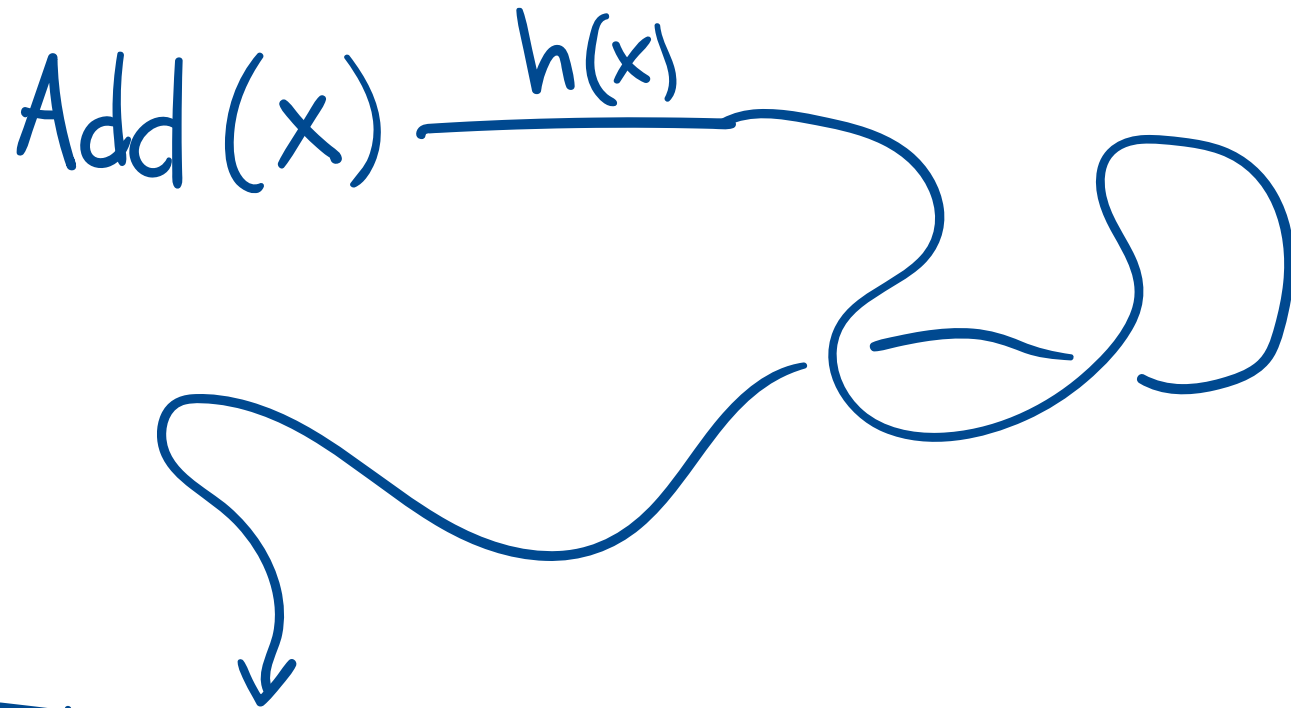
# Applications 3/3

- **Chrome**: local filter for malicious URLs
  - Mostly misses
  - Google doesn't see where you try to go
  - You don't get the list
  - False positives?
    - Unnecessary warning; or ask Google.

# Bit Set

- Good constant factors

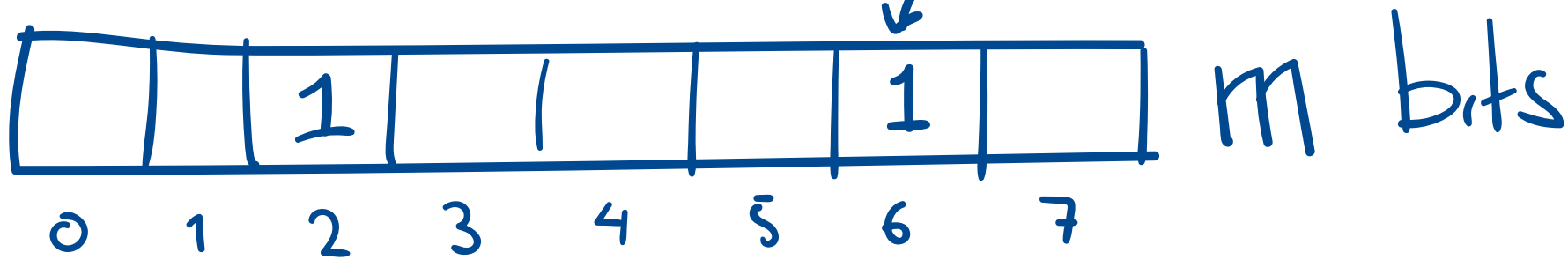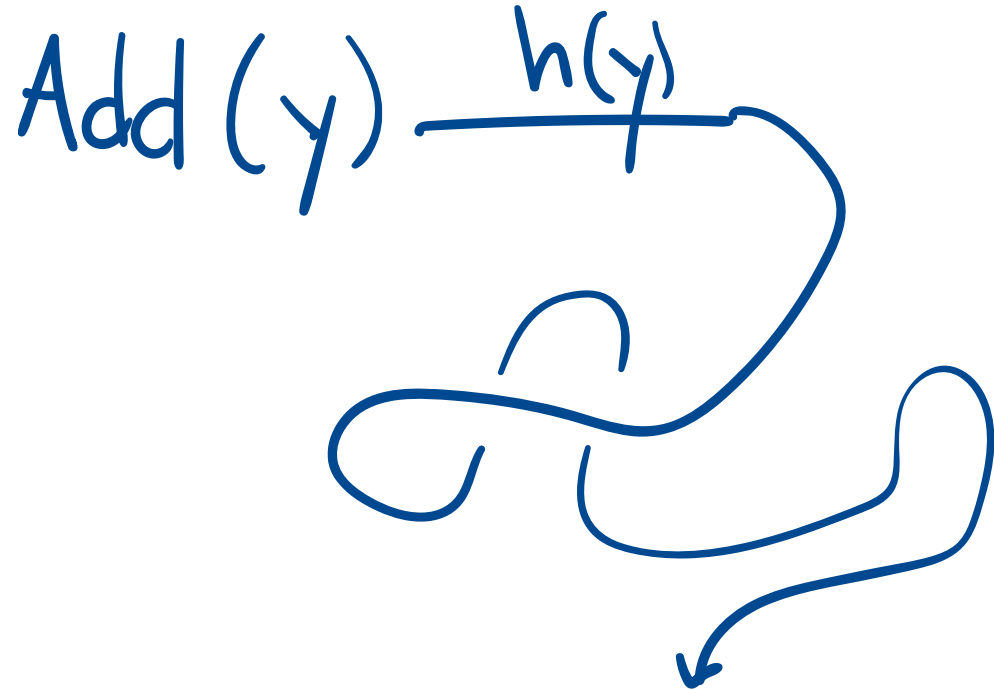- Too large when universe U is large

  - Especially annoying if $n \ll |U|$
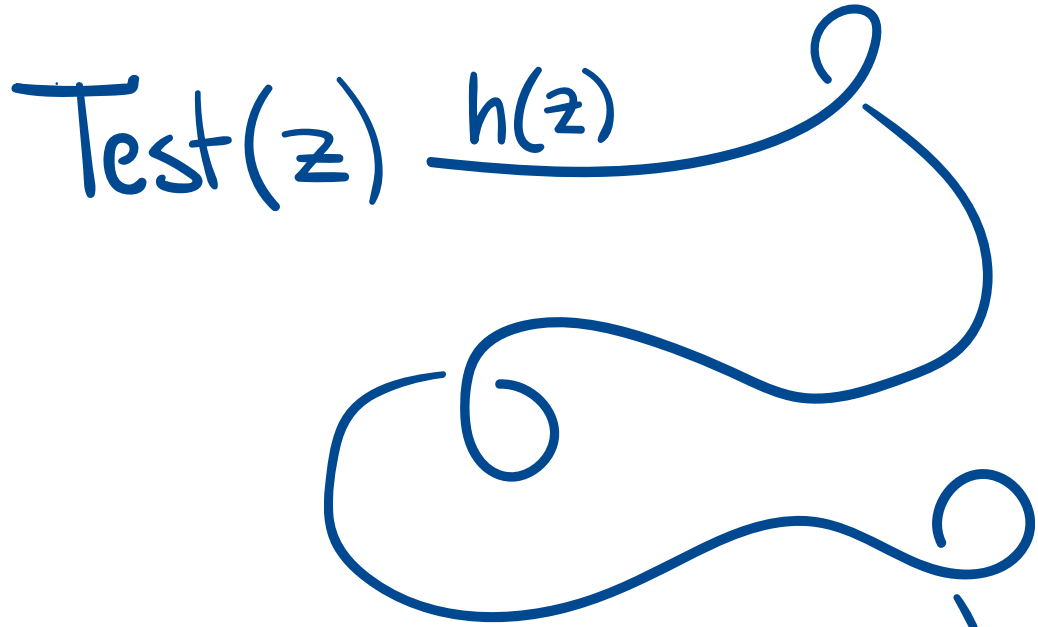
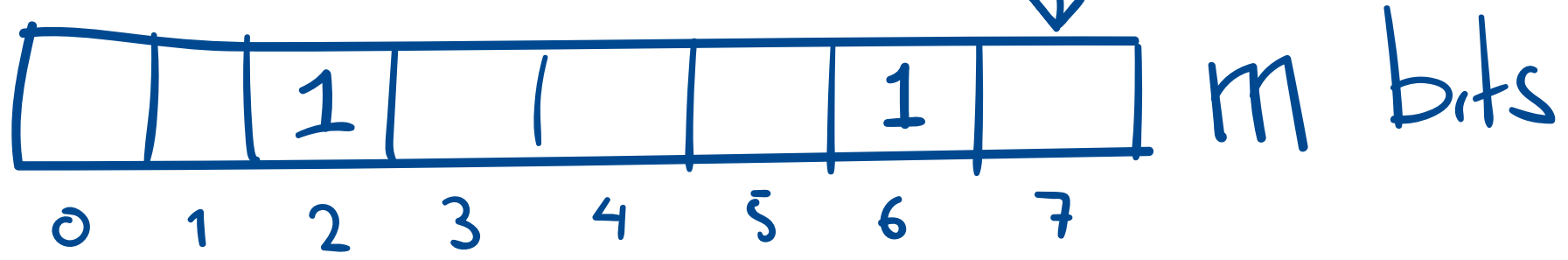Add (6)

$h: U \rightarrow \mathbb{Z}_m$

$Add(x) \xrightarrow{h(x)}$

| | | 1 | 1 | | | | |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

m bits

$$h: U \to \mathbb{Z}_m$$

$\text{Add}(y) \xrightarrow{h(y)}$



| | | 1 | 1 | | 1 | | | m bits |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | |

$$h: U \longrightarrow \mathbb{Z}_m$$

Test(z) $\underline{h(z)}$

$z \overset{?}{\in} S$

NO!

| | | 1 | 1 | | 1 | | | m bits
|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

$$h: U \rightarrow \mathbb{Z}_m$$

Test(x) $\xrightarrow{h(x)}$

$x \overset{?}{\in} S$

| | | 1 | | 1 | | 1 | | m bits
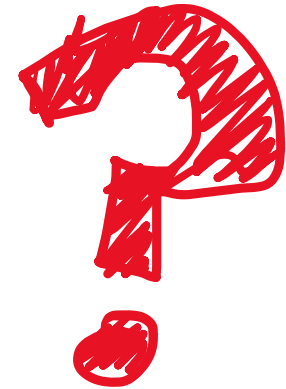|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

Remove(x)?

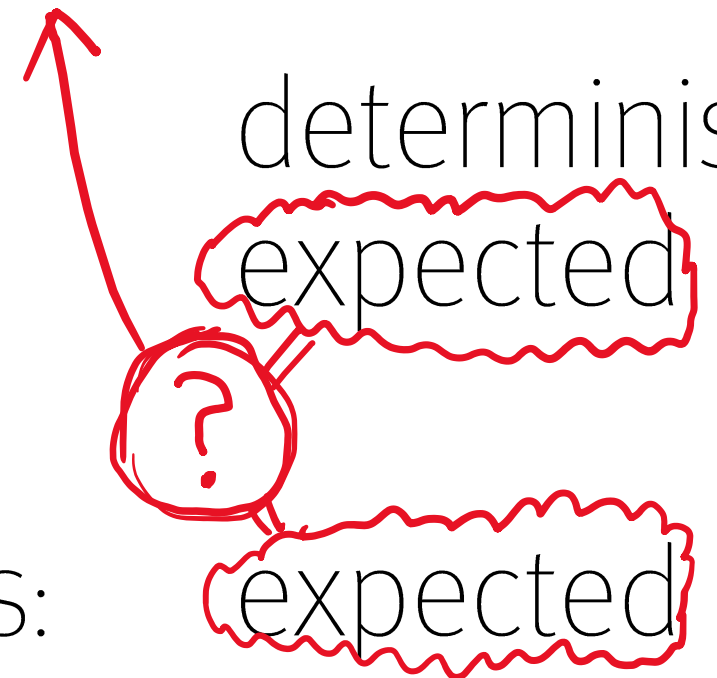# The function h

– Take a "hash function"

– The function is deterministic

– For the analysis, we will assume it gives independent uniformly random indices
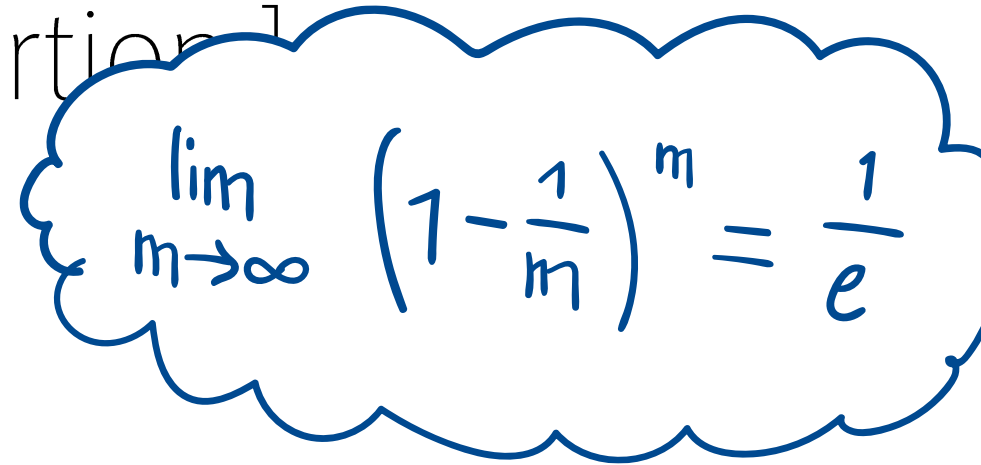
# "Probabilistic data structure"

– Hashmaps: deterministic correctness
expected runtime

– Bloom filters: expected correctness
deterministic runtime

# False positive probability

P[ bit $i$ is **0** after the first insertion ]

$$1 - \frac{1}{m}$$

$$\lim_{m \to \infty} \left(1 - \frac{1}{m}\right)^m = \frac{1}{e}$$

P[ bit $i$ is still **0** after the first $n$ insertions ]

$$\left(1 - \frac{1}{m}\right)^n = \left[\left(1 - \frac{1}{m}\right)^m\right]^{n/m} \approx e^{-\frac{n}{m}}$$

# False positive probability

P[ bit $i$ is **1** after $n$ insertions ]

$$\approx 1 - e^{-\frac{n}{m}}$$

- Consider $Test(x)$ for nonmember $x$
- Bit $h(x)$ is set with probability...

# Parameters

n     Number of items

m    Number of bits
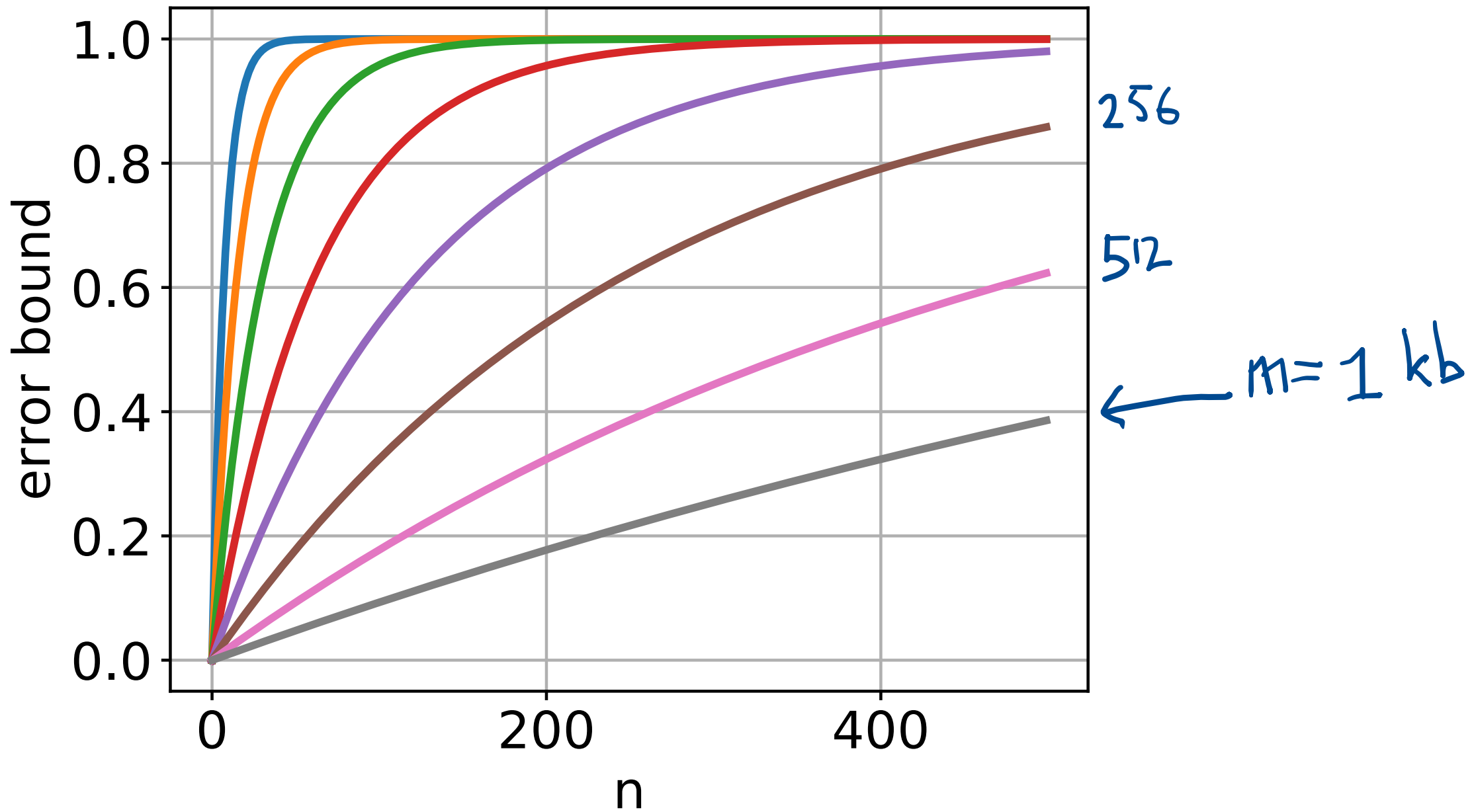
ε     Error bound

Probably fixed

{
- Maybe fixed
- Costs space
}

# Dictionary example

| | n | m | ε |
|---|---|---|---|
| **English dictionary** | 500.000 | 100 kB | ≈ 46% |
| (≈3 MB ASCII) | | 1 MB | ≈ 5% |
| | | 3 MB | ≈ 1.9% |

$$h_1, h_2 : U \longrightarrow \mathbb{Z}_m$$

Test(z)   $h_1(z)$

$h_2(z)$

$z \overset{?}{\in} S$

NO!



| | | 1 | | 1 | | | | 1 | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | | |

m bits

# Bloom filter

- Fix **k** hash functions $h_i$
- Storage: array of **m** bits, all start unset

*Add(x)*: set all bits $h_i(x)$

*Test(x)*: are all bits $h_i(x)$ set?

# What is the effect of k?

- Increases runtime
- It does not affect the space!

- Error probability?
  - Check more bits: accidents *less* likely
  - Set more bits: accidents are *more* likely

# False positive probability

– A particular bit is **0** after the first insertion:

$$\left(1 - \frac{1}{m}\right)^{k} = \left[\left(1 - \frac{1}{m}\right)^{m}\right]^{k/m} \approx e^{-\frac{k}{m}}$$

– A particular bit is still **0** after $n$ insertions:

$$\left(1 - \frac{1}{m}\right)^{kn} \approx e^{-\frac{kn}{m}}$$

# False positive probability

- A particular bit is **1** after n insertions:

$$\approx 1 - e^{-\frac{kn}{m}}$$

- False positive $Test(x)$

$$\left[ 1 - \left(1 - \frac{1}{m}\right)^{kn} \right]^k \overset{\text{(handwave!)}}{\approx} \left(1 - e^{-\frac{kn}{m}}\right)^k$$

# Parameters

$n$  Number of items

*Probably fixed*

$m$  Number of bits

- Maybe fixed
- Costs space

$k$  Number of hash functions  ← • Costs time

$\varepsilon$  Error bound

# Picking k

– Idea: given $n$ and $m$, pick $k$ to minimize $\varepsilon$.

$$\varepsilon \approx \left(1 - e^{-\frac{kn}{m}}\right)^k = \exp\left(k \ln\left(1 - e^{-\frac{kn}{m}}\right)\right)$$

$\xrightarrow{\ln}$  $g \overset{\text{def}}{=} k \ln(1 - p)$

$$\frac{dg}{dk} \overset{!}{=} 0$$

# Picking k

– Idea: given $n$ and $m$, pick $k$ to minimize $\varepsilon$.

$$\varepsilon \approx \left(1 - e^{-\frac{kn}{m}}\right)^k = \exp\left(k \ln\left(1 - e^{-\frac{kn}{m}}\right)\right)$$

$\xrightarrow{\ln}$ $g \stackrel{\text{def}}{=} k \ln(1-p) = -\frac{m}{n} \ln(p) \ln(1-p)$

$$= m \frac{n}{k} \cdot \ln\left(\frac{mk}{n}\right) \, \ln(1-p)$$

# Picking k

$n, m, k \in \mathbb{N}$

– Idea: given $n$ and $m$, pick $k$ to minimize $\varepsilon$.

$$\varepsilon \approx \left(1 - e^{-\frac{kn}{m}}\right)^k = \exp\left(k \ln\left(1 - e^{-\frac{kn}{m}}\right)\right)$$

$$\xrightarrow{\ln} \quad g \overset{\text{def}}{=} k \ln(1-p) = -\frac{m}{n} \ln(p) \ln(1-p)$$

$$\underbrace{\qquad\qquad\qquad}_{\text{min} \Rightarrow p = \frac{1}{2}}$$

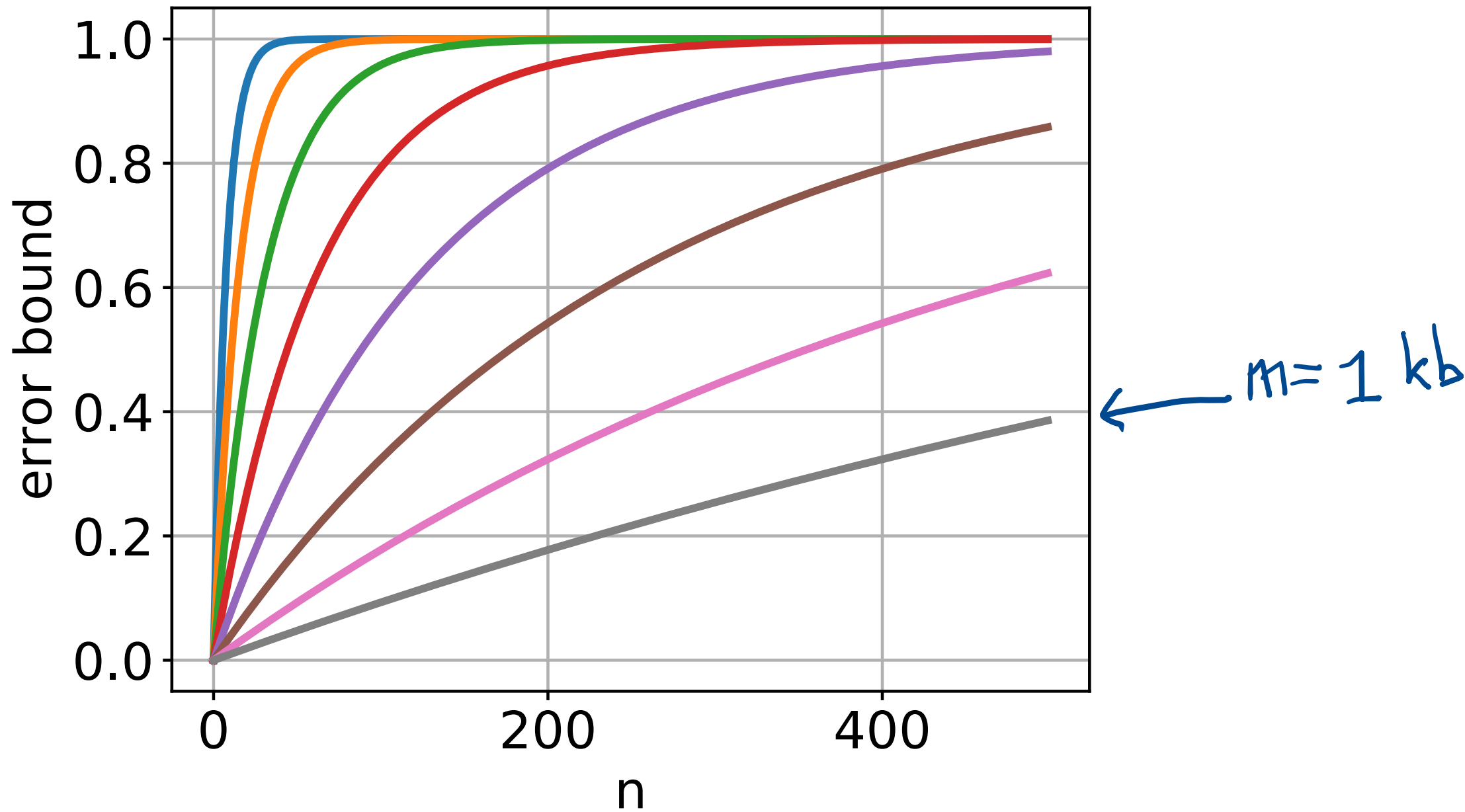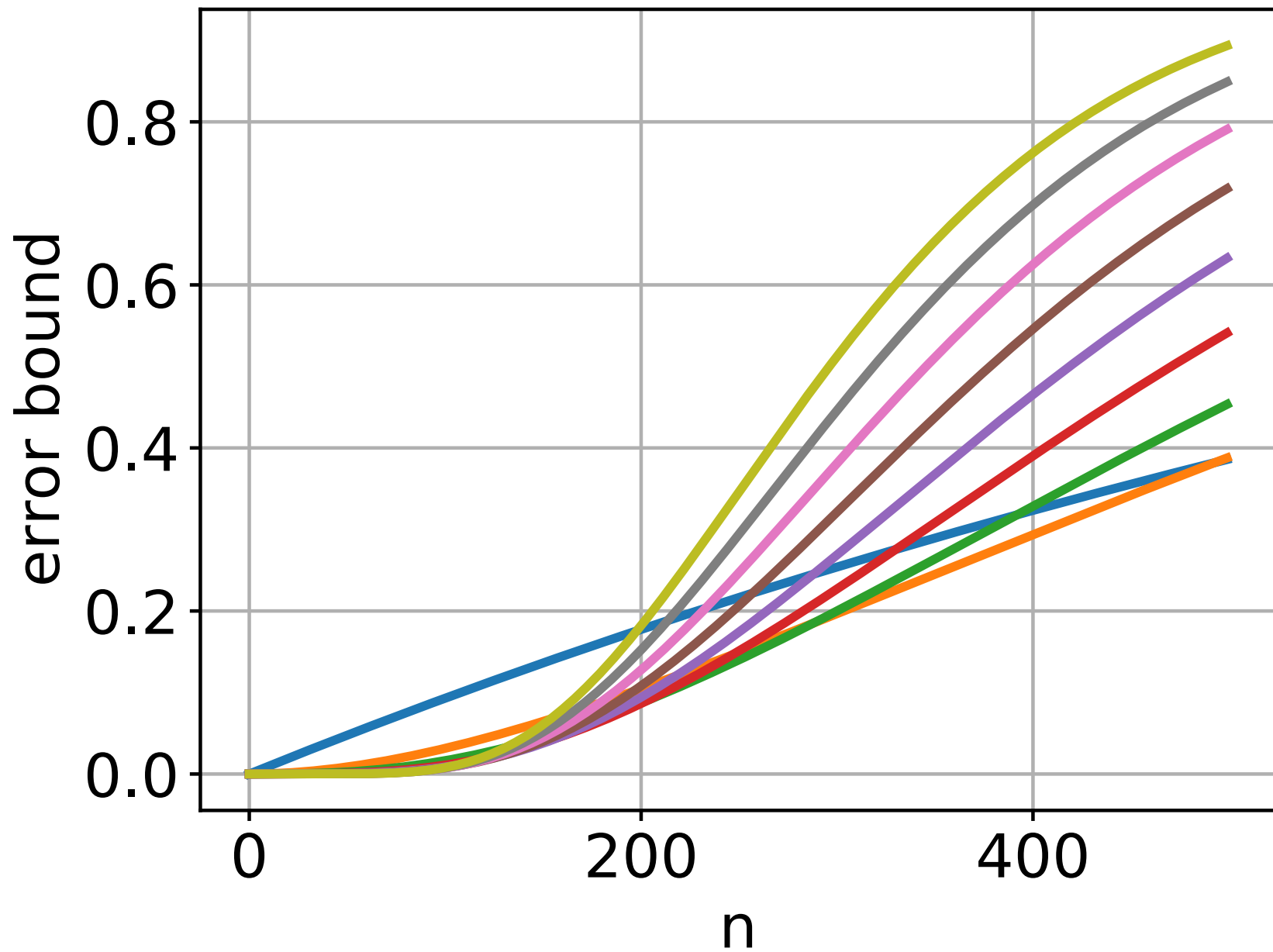$$e^{-\frac{kn}{m}} = \frac{1}{2} \Rightarrow k = \ln 2 \cdot \frac{m}{n}$$

# Picking k

$n, m, k \in \mathbb{N}$

- Idea: given $n$ and $m$, pick $k$ to minimize $\varepsilon$.

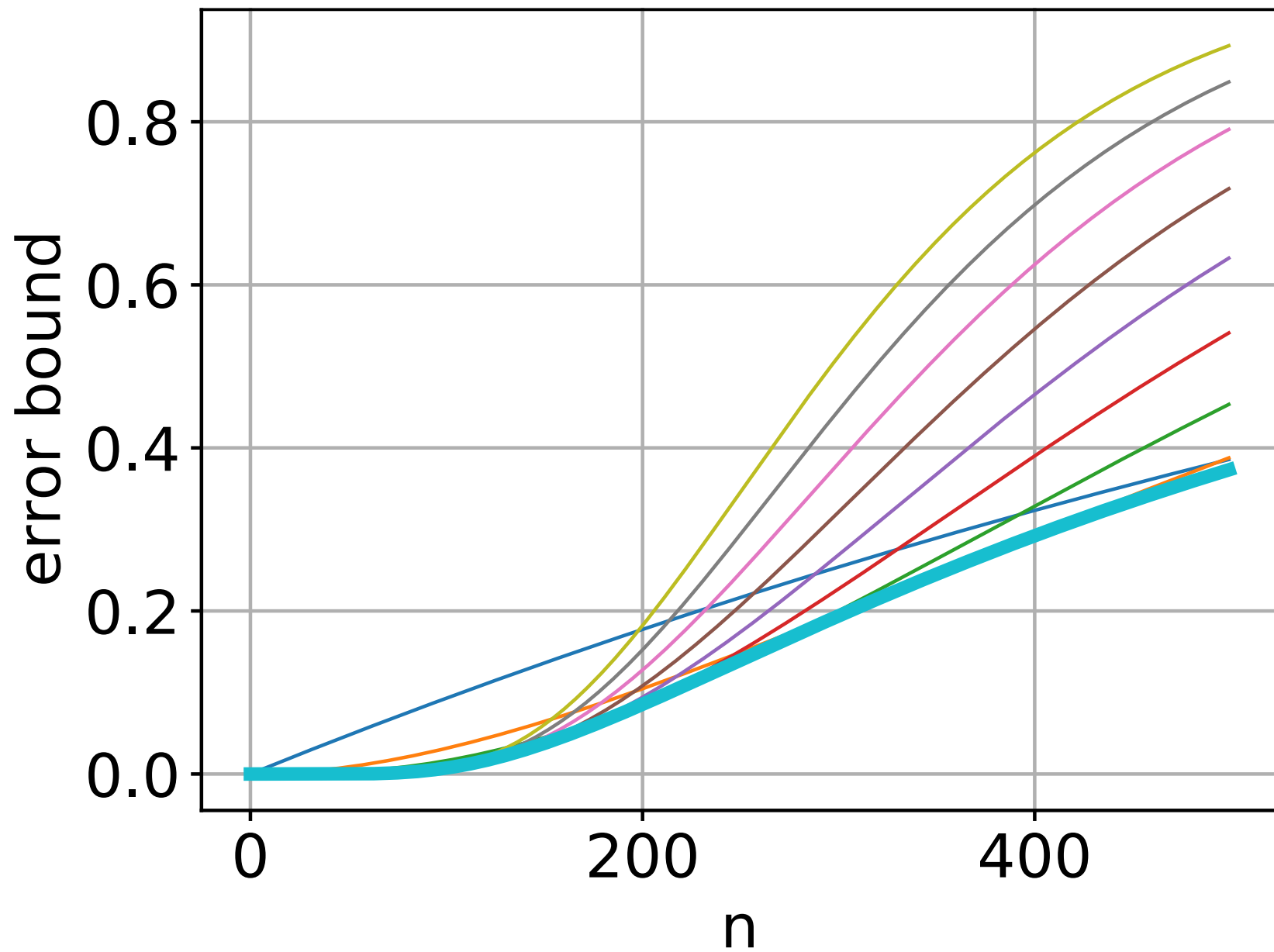- Optimal: $k = \ln 2 \cdot \dfrac{m}{n}$

$$\varepsilon = \left(\tfrac{1}{2}\right)^k \approx 0.6185^{\frac{m}{n}}$$

Values
of **k**

m = 1 kb

Optimal
$k$

$m = 1\ kb$

# Dictionary example

| | n | m | k | ε |
|---|---|---|---|---|
| **English dictionary** | 500.000 | 100 kB | 1 | ≈ 46 % |
| (≈3 MB ASCII) | | 1 MB | 1 | ≈ 5 % |
| | | 3 MB | 1 | ≈ 1.9 % |
| | | 512 kB | 1 | ≈ 11 % |
| | | | 2 | ≈ 4 % |
| | | | 6* | ≈ 1 % |
| | | 1 MB | 12* | ≈ 0.03 % |
| | | 3 MB | 35* | < $10^{-10}$ |