

# COMMONWEALTH OF AUSTRALIA

## Copyright Regulations 1969

### WARNING

This material has been reproduced and communicated to you by or on behalf of the University of Sydney pursuant to Part VB of the Copyright Act 1968 (**the Act**). The material in this communication may be subject to copyright under the Act. Any further copying or communication of this material by you may be the subject of copyright protection under the Act.

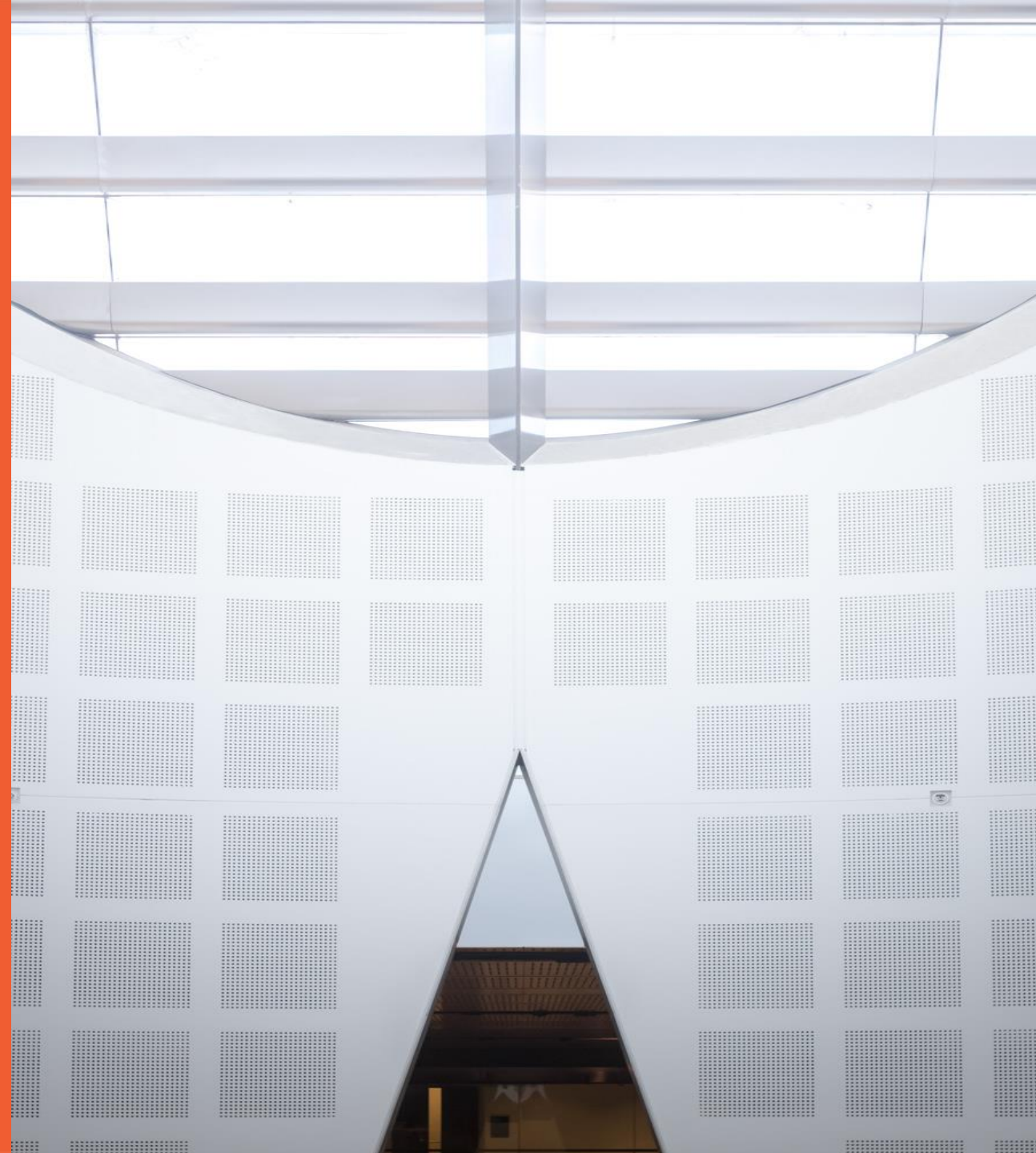
**Do not remove this notice.**

COMPx270: Randomised and  
Advanced Algorithms  
Lecture 6: Hashing and Friends

Clément Canonne  
School of Computer Science



THE UNIVERSITY OF  
SYDNEY



## A question ?

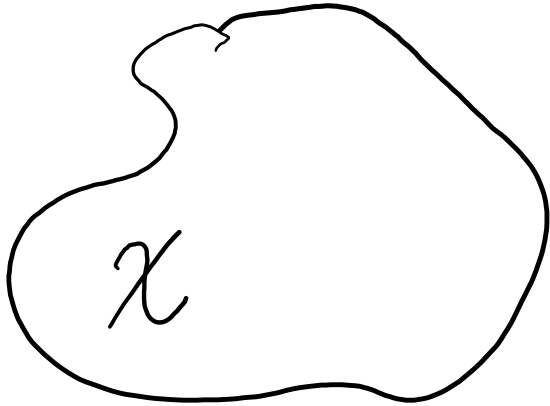
A SID is of the form **450687816** (9 digits, each between 0 and 9). How many distinct SIDs are there?

$$10^9$$

How many distinct students have there even been at U Syd?

$$\leq 250,000 = 10^7$$

# Dictionaries (Maps, Associative arrays...)



↑  
 $m = 10^9$

$S \subseteq X$

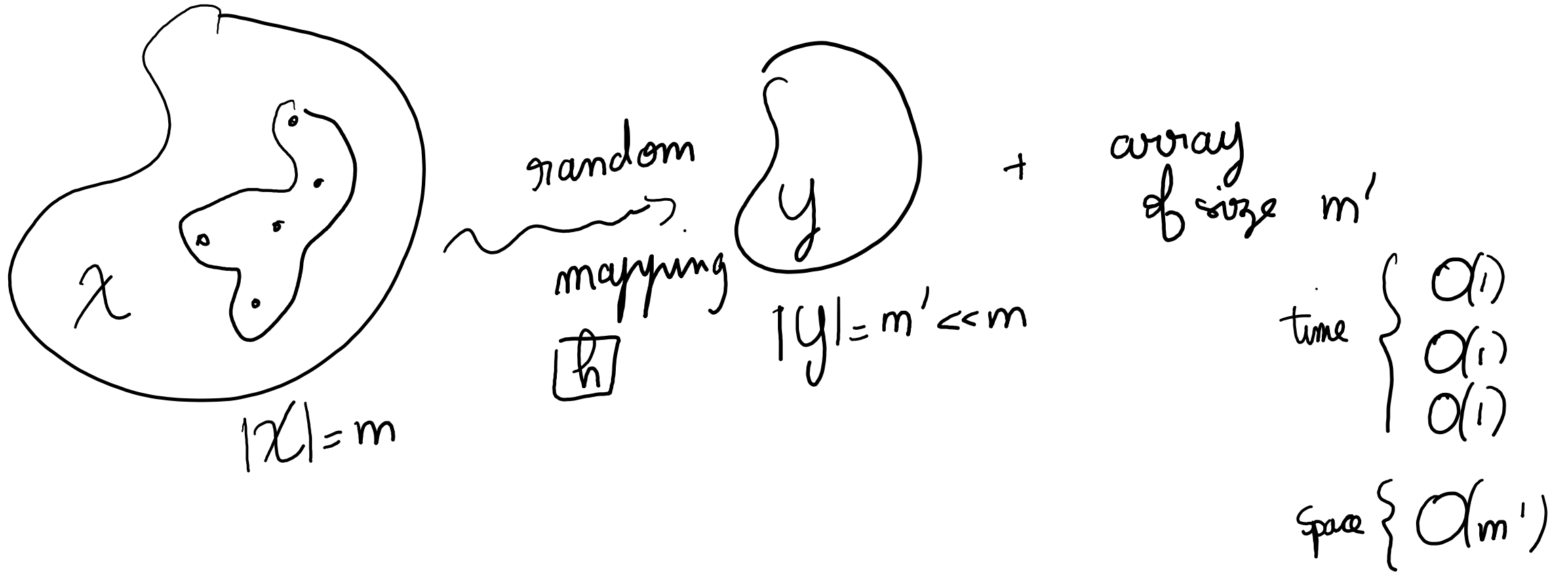
$|S| \ll |X|$   
 ↑                    ↑  
 n                    m

INSERT+  
 LOOKUP  
 REMOVE

SPACE

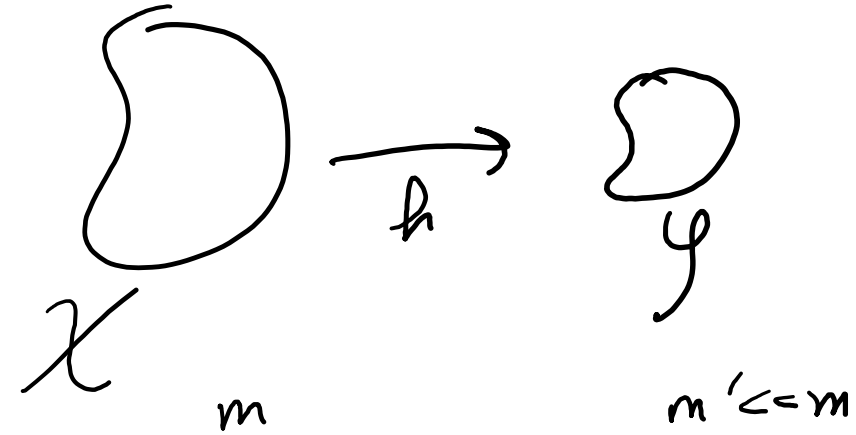
	LIST	ARRAY	BST
INSERT+ LOOKUP REMOVE	$\begin{cases} O(1)/O(n) \\ O(n) \\ O(n) \end{cases}$	$\begin{cases} O(1) \\ O(1) \\ O(1) \end{cases}$	$\begin{cases} O(\log n) \\ O(\log n) \\ O(\log n) \end{cases}$
SPACE	$O(n)$ ✓	$O(m)$ x	$O(n)$ ✓

# Hash tables

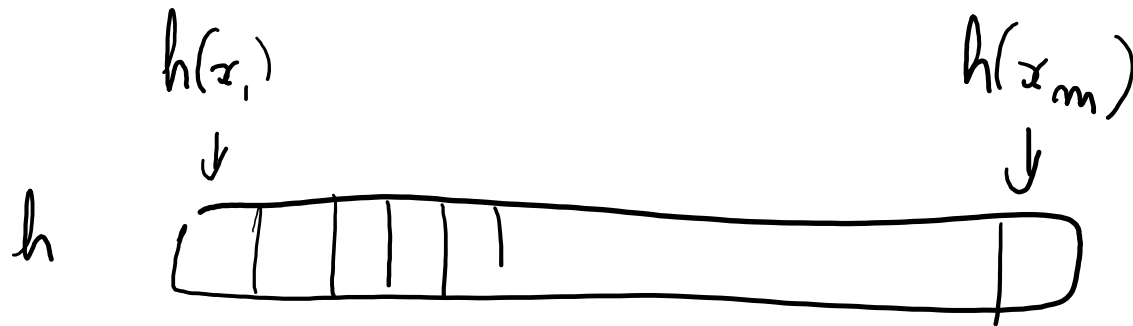


# Hash tables: what is random?

- Dataset?
- $h$  itself  $\rightarrow \Theta(m \log m')$  space  $\times$
- $\mathcal{H}$ : hash family  $\rightarrow \mathcal{O}(\log |\mathcal{H}|)$

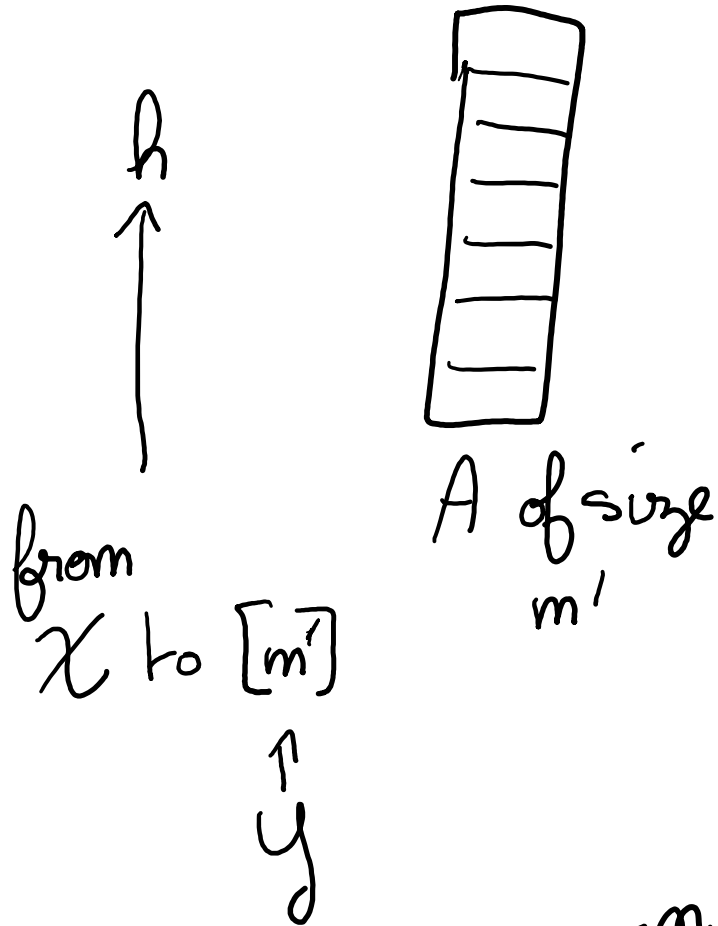


$S$ :  $n$  points



universal hash family  $\xrightarrow{h \sim \mathcal{H}}$   $\forall x \neq x' \Pr [h(x) = h(x')] \leq \frac{1}{m}$

# Hash tables: the data structure



$$m, (m'), (n)$$

INSERT( $x$ ):  $A[h(x)] \leftarrow 1$   
LOOKUP  $A[h(x)] == 1?$   
REMOVE  $A[h(x)] \leftarrow 0$

SPACE:  $O(\log |S| + m')$

Can we set  $m' = O(n)$ ?

Collision  $x \neq x'$   
 $h(x) = h(x')$

INSERT( $x'$ )  
LOOKUP( $x$ )  
REMOVE( $x$ )

LOOKUP( $x'$ )

## Hash tables: no collisions?

- $m' = \Omega(n^2)$  (bad?)

- If  $m' = O(n)$

max # collisions  
in some bucket  $= \Theta\left(\frac{\log n}{\log \log n}\right)$

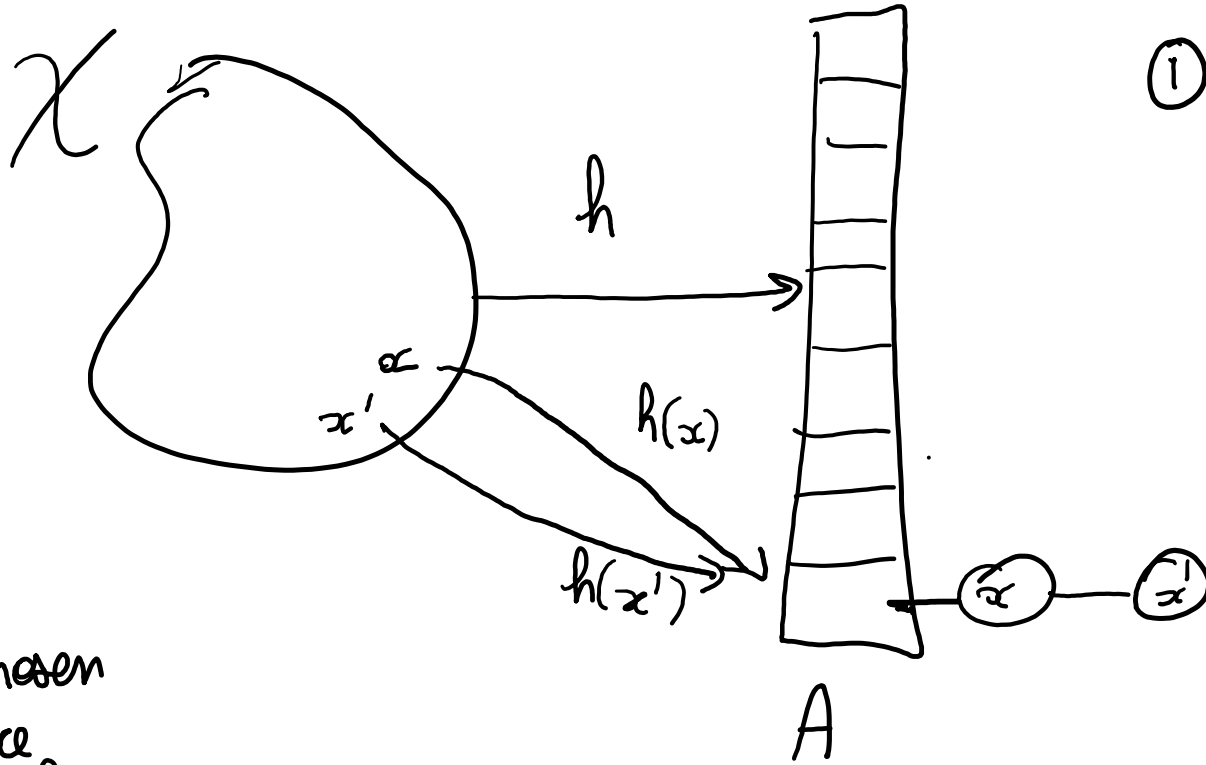


# Hash tables: ~~no~~ collisions

Load factor

$$\alpha = \frac{n}{m'}$$

① Chaining



(  $h$  chosen once and for all from "good"  $H$  )

Expected time complexity of

INSERT  
LOOKUP  
REMOVE

$$O(1 + \alpha)$$

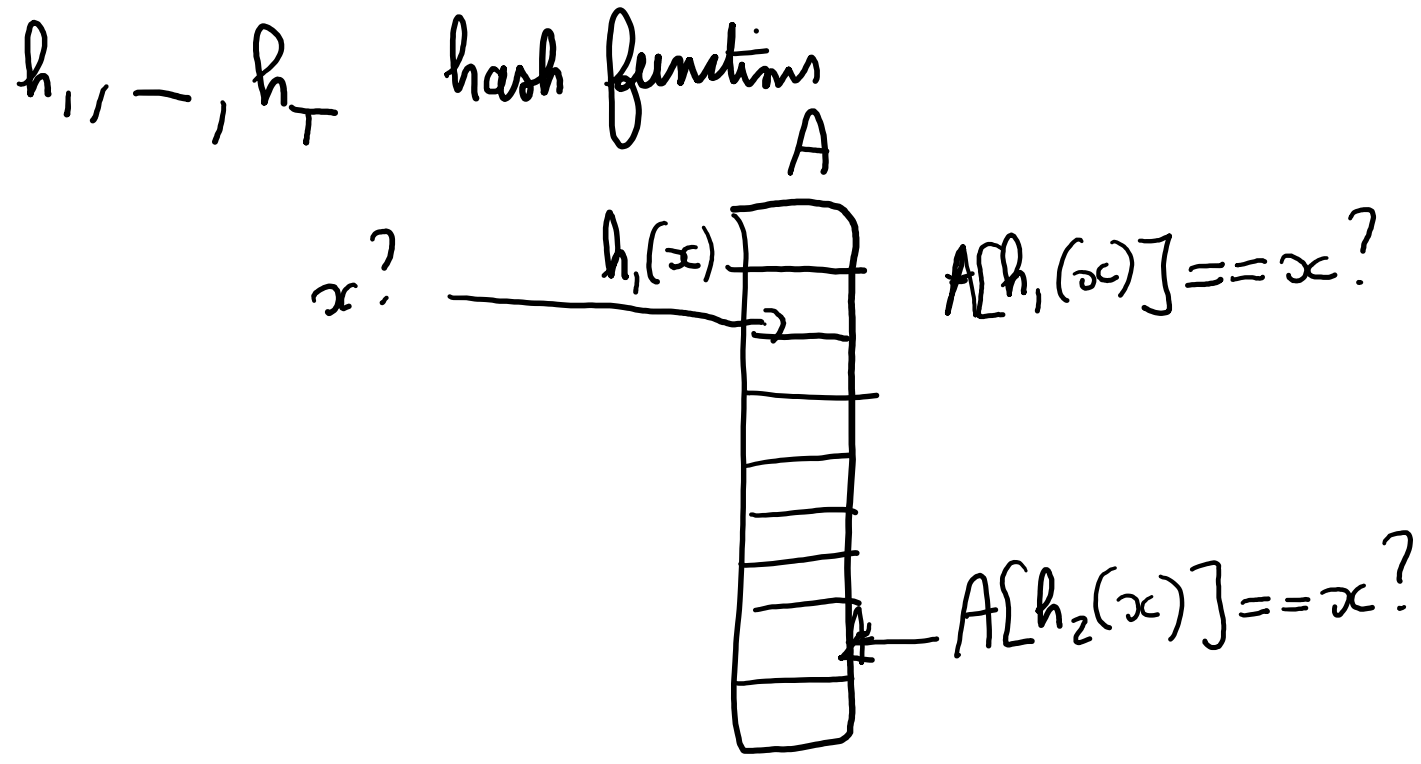
Worst-case  $O(n)$

Will happen  $O\left(\frac{\log n}{\log \log n}\right)$

# Hash tables: collisions

$$\alpha = \frac{n}{m'}$$

Open addressing  
( $\alpha \leq 1$ )



$\forall x, (h_1(x), \dots, h_T(x))$  is a permutation of  $[m']$

$\uparrow$   
 $T=m'$

- Query complexity can be bad:  $O(\alpha n)$
- Space:  $O(\underline{m'} \log m + m')$

# Handling collisions: separate chaining

# Handling collisions: open addressing

$h_1, \dots, h_m$

INSERT( $x$ )

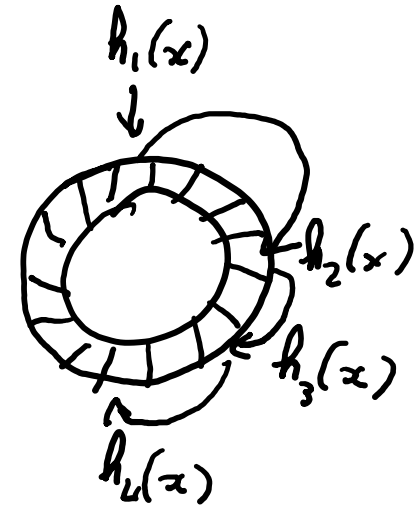
For  $1 \leq t \leq m$

if  $A[h_t(x)] == x$  return

if  $A[h_t(x)] == \emptyset$  or  $A[h_t(x)] == \perp$

$A[h_t(x)] \leftarrow x$

return



REMOVE( $x$ )

For  $1 \leq t \leq m$

if  $A[h_t(x)] == x$

$A[h_t(x)] \leftarrow \perp$ ; return

if  $A[h_t(x)] == \emptyset$  return

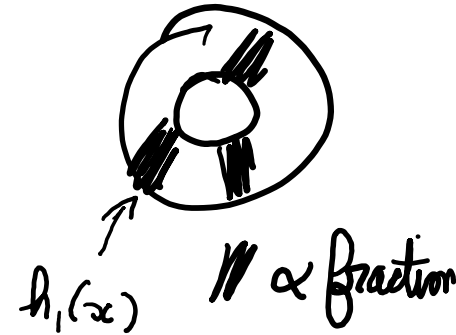
# Handling collisions: open addressing

Assuming [stuff], expected  
time complexity of LOOKUP is  $O\left(\frac{1}{1-\alpha}\right)$

[stuff]  $\forall x, (h_1(x), \dots, h_m(x))$  is a u.a.r permutation of  $[m]$

$$\begin{aligned} E[T_{n,m'}] &= O(1) + \alpha \cdot E[T_{n-1,m'-1}] \\ &\leq \underbrace{O(1)}_{\text{"hardware"}} + \alpha E[T_{n,m'}] \end{aligned}$$

$$\rightarrow E[T_{n,m'}] \leq \frac{O(1)}{1-\alpha}$$



# Handling collisions: open addressing (linear probing)

$h_1, \dots, h_m$  are one hash function in disguise

$$h_1(x) = h(x)$$

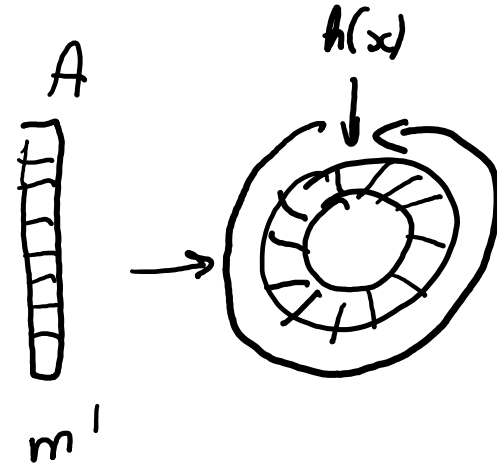
$$h_2(x) = h(x) + 1 \quad [m']$$



$$h_{m'}(x) = h(x) + m' - 1 \quad [m']$$

Theorem

Under some assumptions, expected time complexity of LOOKUP is  $O\left(\frac{1}{(1-\alpha)^2}\right)$



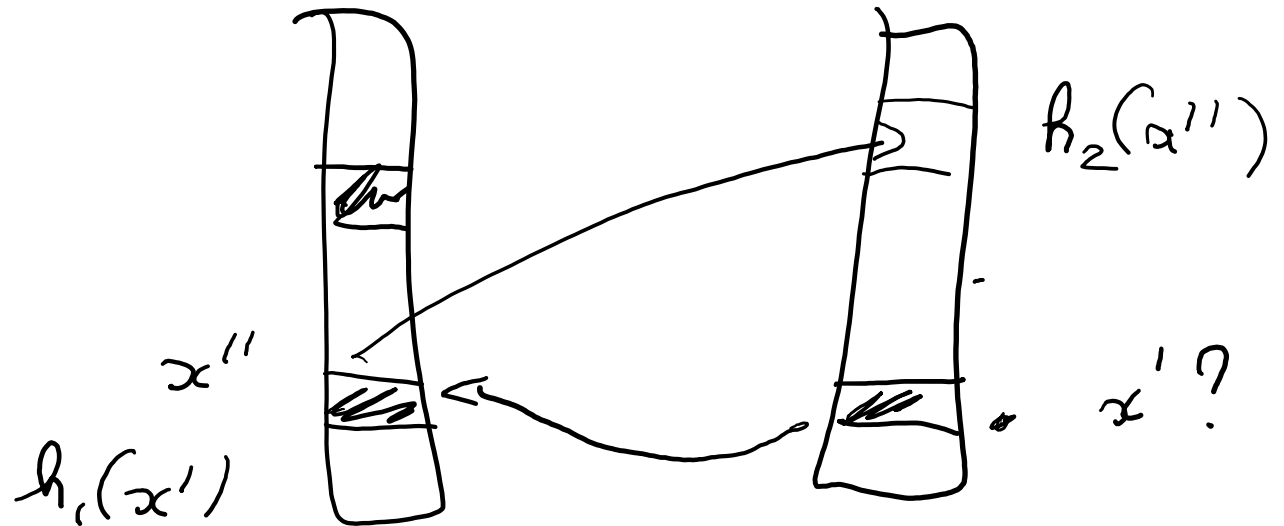
# Handling collisions: open addressing (cuckoo hashing)

$h_1, A_1$

$h_2, A_2$

$$x \rightarrow A_1[h_1(x)] \quad \text{or} \quad A_2[h_2(x)]$$

LOOKUP }  $O(1)$  worst-case  
REMOVE }  
INSERT }  $O(1)$  expected

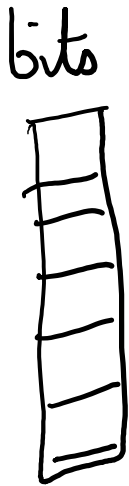


# Hash tables: summary



# Can we do better? Bloom filters!

Bloom filter:  $\left\{ \begin{array}{l} \text{faster } (O(1) \text{ for read}) \\ \text{less space (cst factors)} \end{array} \right.$



But LOOKUP is sometimes wrong  
(false positives)

$$A[h_1(x)] \wedge A[h_2(x)] \wedge \dots \wedge A[h_k(x)]$$