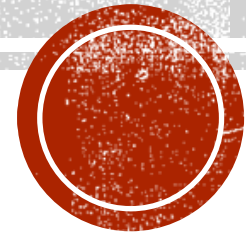




Indian Institute of Information Technology Allahabad

Data Structures and Algorithms

Hashing



Dr. Shiv Ram Dubey

Associate Professor

Department of Information Technology

Indian Institute of Information Technology, Allahabad

Email: srdubey@iiita.ac.in

Web: <https://profile.iiita.ac.in/srdubey/>

DISCLAIMER

The content (text, image, and graphics) used in this slide are adopted from many sources for academic purposes. Broadly, the sources have been given due credit appropriately. However, there is a chance of missing out some original primary sources. The authors of this material do not claim any copyright of such material.

Today

- Hashing!
 - What operations are we trying to support?
 - Hash Functions
 - Dealing with collisions
 - What makes a good hash function?
 - *Universal* hash families are what we're looking for!

Hash Tables Overview

What operations does it support?

The Task

Again, we want to keep track of objects that have keys 5 (aka, **nodes** with **keys**)

The Task

Again, we want to keep track of objects that have keys 5 (aka, **nodes** with **keys**)

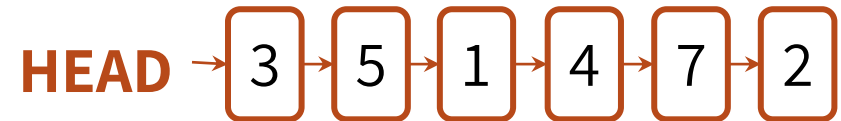
Sorted Arrays



$O(n)$ INSERT/DELETE: first, find the relevant element (via SEARCH) and move a bunch of elements in the array

$O(\log n)$ SEARCH: use binary search to see if an element is in A

Linked Lists



$O(1)$ INSERT: just insert the element at the head of the linked list

$O(n)$ SEARCH/DELETE: since the list is not necessarily sorted, you need to scan the list (delete by manipulating pointers)

Hash Table: Motivation

OPERATION	SORTED ARRAY	UNSORTED LINKED LIST	HASH TABLES (HOPEFULLY)
SEARCH	$O(\log(n))$	$O(n)$	$O(1)$
DELETE	$O(n)$	$O(n)$	$O(1)$
INSERT	$O(n)$	$O(1)$	$O(1)$

Hash Table: Motivation

OPERATION	SORTED ARRAY	UNSORTED LINKED LIST	HASH TABLES (HOPEFULLY)
SEARCH	$O(\log(n))$	$O(n)$	$O(1)$
DELETE	$O(n)$	$O(n)$	$O(1)$
INSERT	$O(n)$	$O(1)$	$O(1)$

What is a *naive* way to achieve these runtimes?

Attempt 1: Direct Addressing

Suppose you're storing numbers from 1 - 1000:



Attempt 1: Direct Addressing

Suppose you're storing
numbers from 1 - 1000:



Reasonable Attempt: *Direct Addressing!*
(each address/bucket stores one type of item)

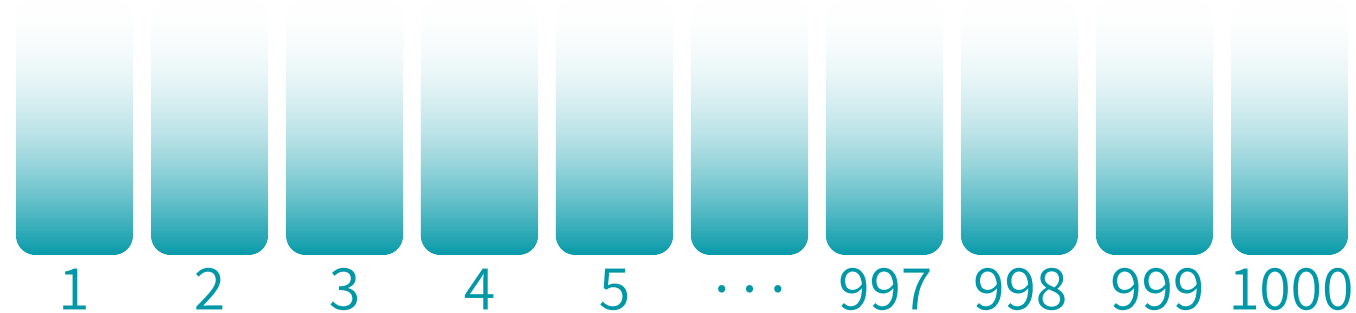
Attempt 1: Direct Addressing

Suppose you're storing numbers from 1 - 1000:



Reasonable Attempt: *Direct Addressing!*

(each address/bucket stores one type of item)



Attempt 1: Direct Addressing

Suppose you're storing numbers from 1 - 1000:



Reasonable Attempt: *Direct Addressing!*
(each address/bucket stores one type of item)



$O(1)$ INSERT/DELETE/SEARCH: Just index into the bucket!

Attempt 1: Direct Addressing

Suppose you're storing numbers from 1 - 1000:

2 4 5 998 999

Not bad!

But what's the issue with this approach?

$O(1)$ INSERT/DELETE/SEARCH: Just index into the bucket!

Attempt 1: Direct Addressing

Suppose you're storing numbers from 1 - 10^{10} :



Attempt 1: Direct Addressing

Suppose you're storing numbers from 1 - 10^{10} :

2 3 1000 1002 10^{10}

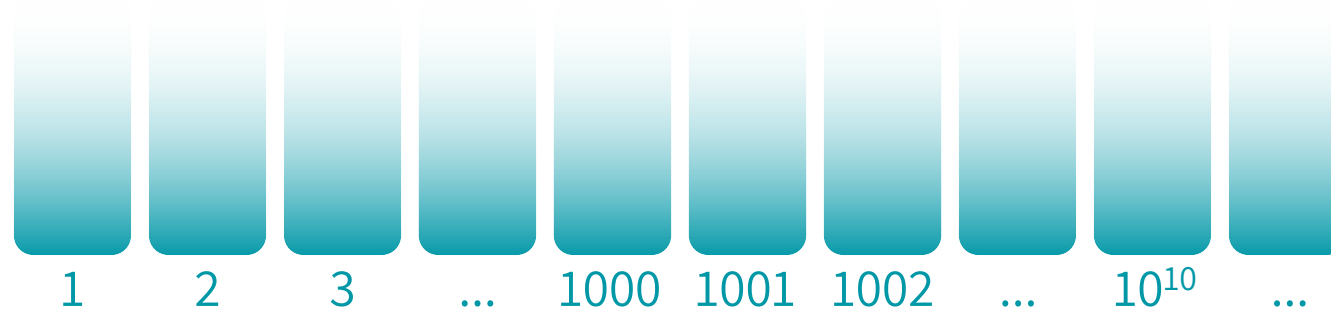
Reasonable Attempt (???): *Direct Addressing!*
(each address/bucket stores one type of item)

Attempt 1: Direct Addressing

Suppose you're storing numbers from 1 - 10^{10} :



Reasonable Attempt (???): *Direct Addressing!*
(each address/bucket stores one type of item)

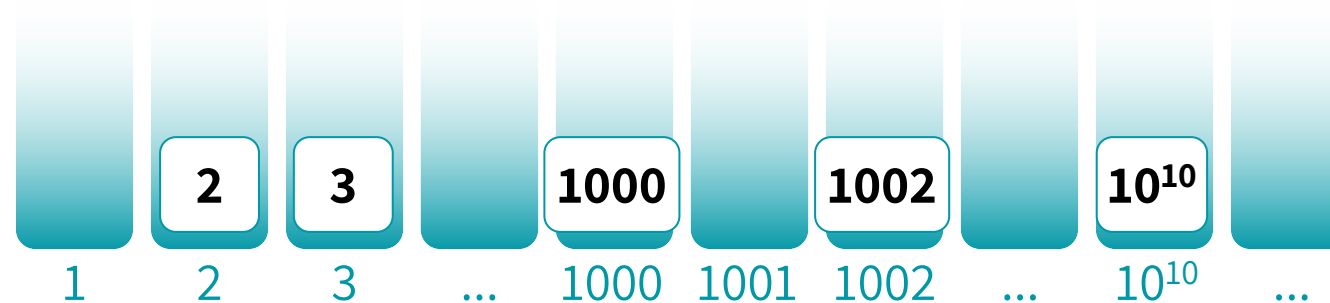


Attempt 1: Direct Addressing

Suppose you're storing numbers from 1 - 10^{10} :



Reasonable Attempt (???): *Direct Addressing!*
(each address/bucket stores one type of item)



$O(1)$ INSERT/DELETE/SEARCH: Just index into the bucket!

Attempt 1: Direct Addressing

Suppose you're storing numbers from 1 - 10^{10} :

2 3 1000 1002 10^{10}

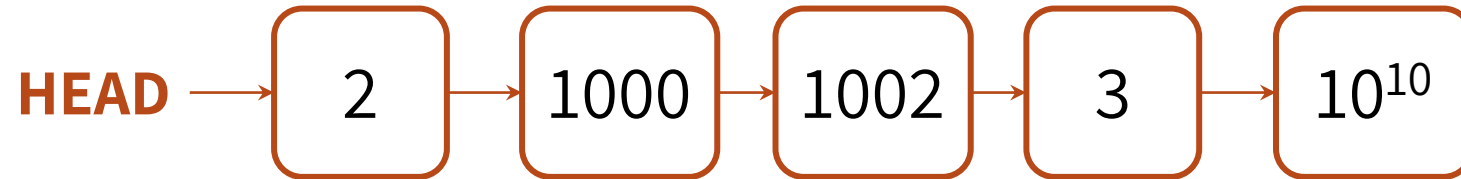
Reasonable Attempt (???): *Direct Addressing!*

**But the space requirement is
HUGE...**

$O(1)$ INSERT/DELETE/SEARCH: Just index into the bucket!

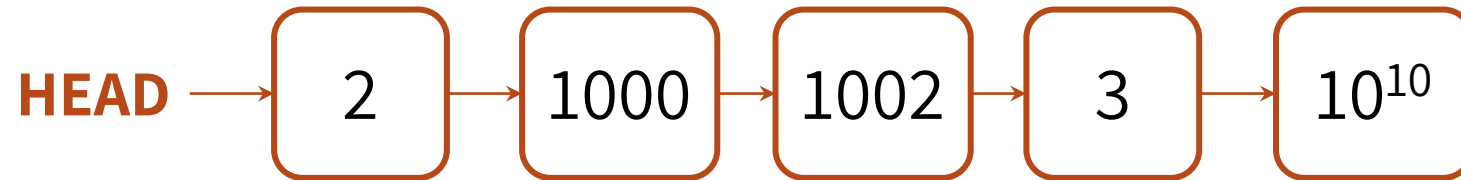
Attempt 2: Back to Linked List

On the other extreme, we could save a lot of space by using linked lists!



Attempt 2: Back to Linked List

On the other extreme, we could save a lot of space by using linked lists!



- **Good news:** Space is now proportional to the number of objects you deal with
- **Bad news:** Searching for an object is now going to scale with the number of inputs you deal with... not close to our desired $O(1)$!
- The direct-addressing approach still has merit because of its fast object search/access

How to improve this?

We like the **functionality of a direct-addressable** array for constant time access
(super fast INSERT/DELETE/SEARCH)

But reserving an bucket/array slot for each possible key leads to unreasonable space requirements

How to improve this?

We like the **functionality of a direct-addressable** array for constant time access
(super fast INSERT/DELETE/SEARCH)

But reserving an bucket/array slot for each possible key leads to unreasonable space requirements

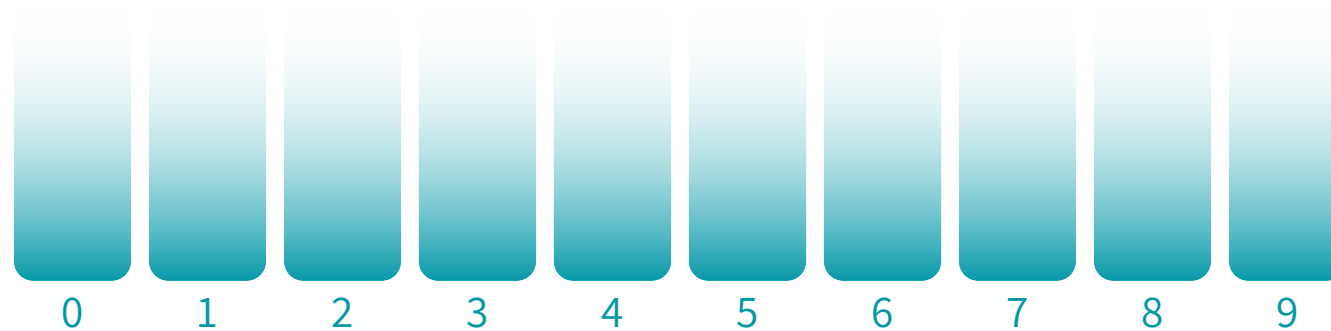
Let's try bucketing **by the least-significant digit...**

Bucketing Attempt 1

Suppose you're storing numbers from 1 - 10^{10} :



Bucket by last digit?

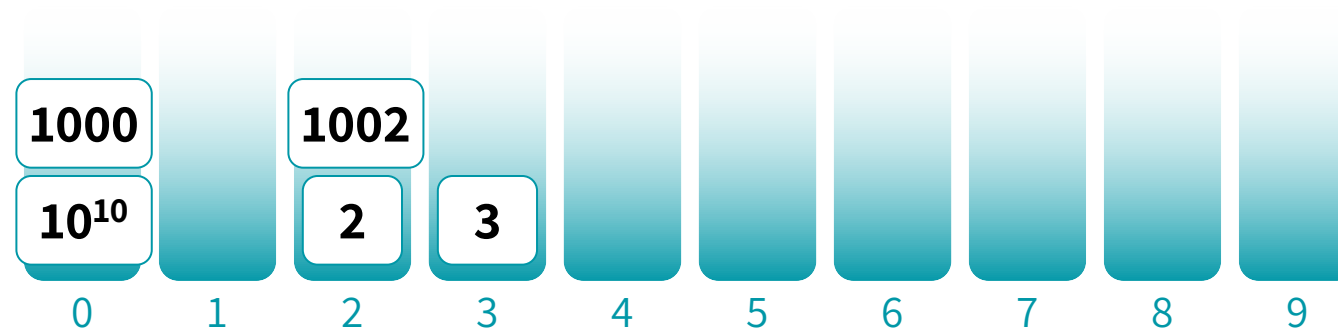


Bucketing Attempt 1

Suppose you're storing numbers from 1 - 10^{10} :



Bucket by last digit?

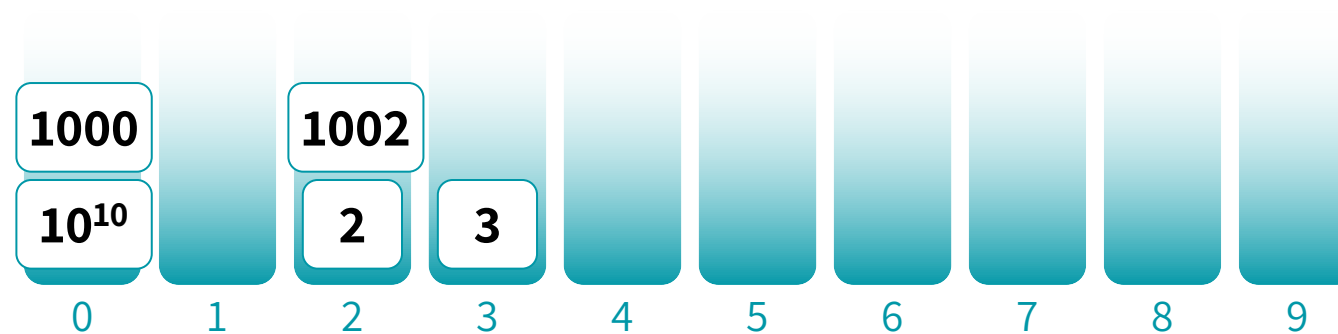


Bucketing Attempt 1

Suppose you're storing numbers from 1 - 10^{10} :



Bucket by last digit?



$O(1)$ INSERT:

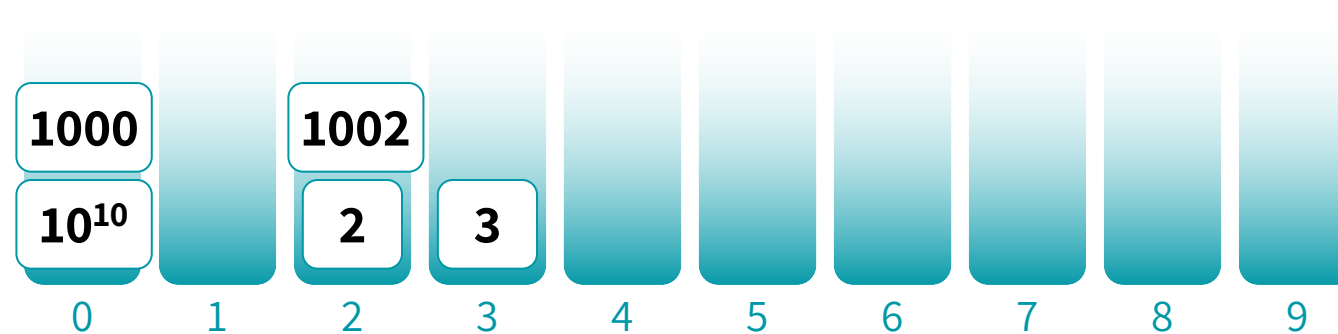
Just index into the bucket (& insert at front of a linked list)!

Bucketing Attempt 1

Suppose you're storing numbers from 1 - 10^{10} :

2 3 1000 1002 10^{10}

Bucket by last digit?

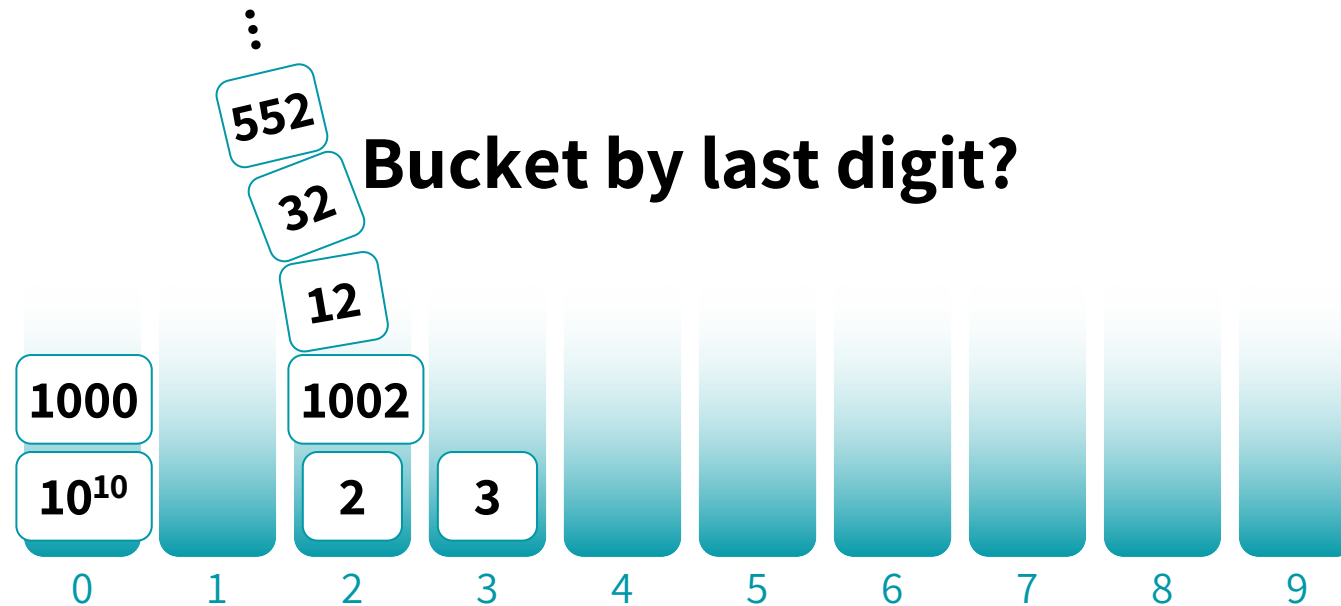


$O(?????)$ SEARCH/DELETE:

Go visit bucket & search through until you find it...

Bucketing Attempt 1

Under this scheme, a bad guy could give us inputs that yields quite ugly worst-case runtimes...

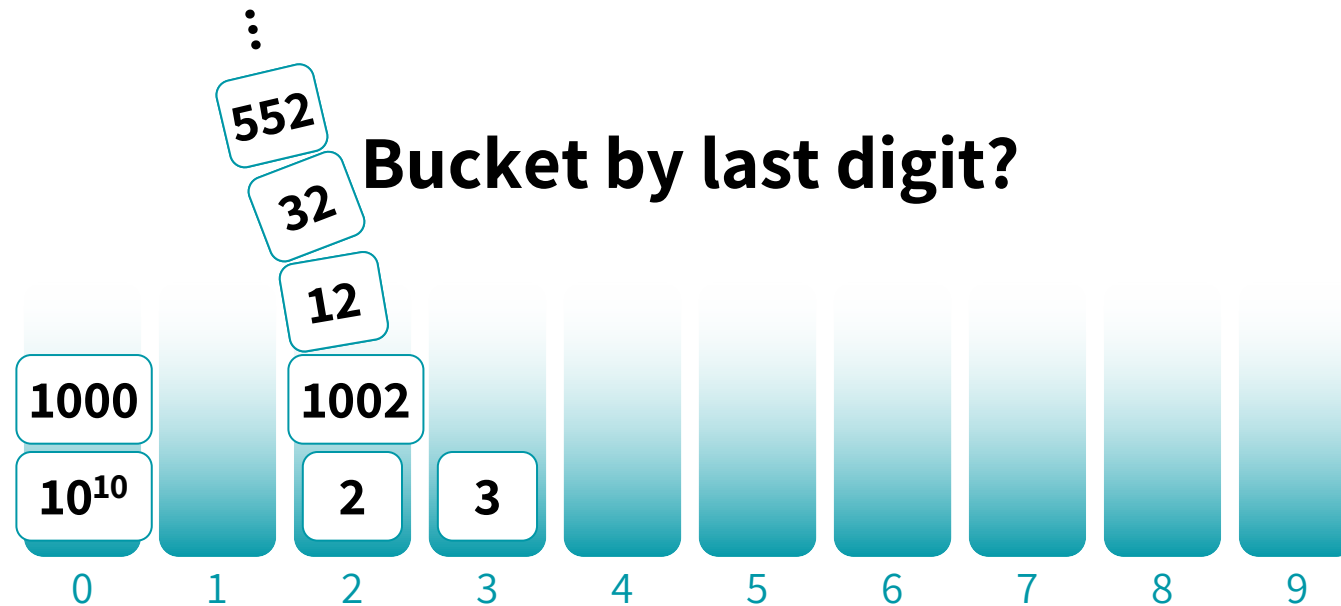


$O(n)$ SEARCH/DELETE:

Go visit bucket & search through until you find it...

Bucketing Attempt 1

Under this scheme, a bad guy could give us inputs that yields quite ugly worst-case runtimes...



Maybe another bucketing scheme?

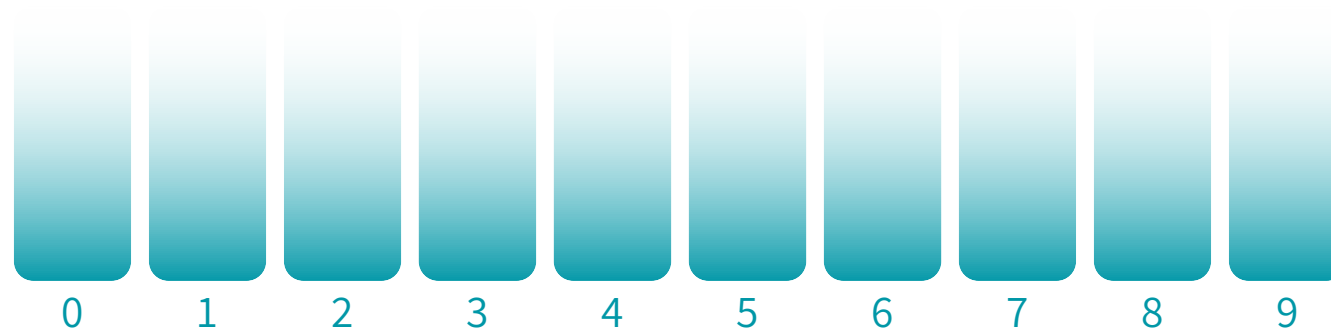
THUD!!...

Bucketing Attempt 2

Suppose you're storing numbers from 1 - 10^{10} :



Bucket by last digit of (number * 7) mod 3

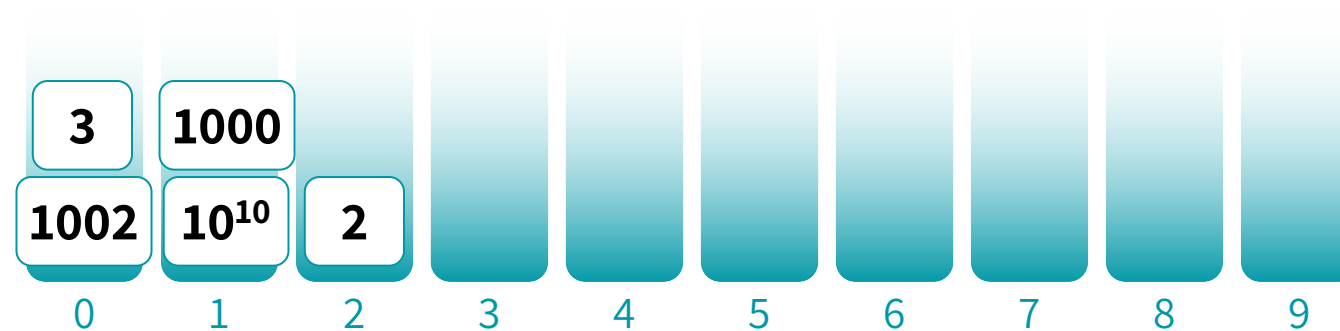


Bucketing Attempt 2

Suppose you're storing numbers from 1 - 10^{10} :



Bucket by last digit of (number * 7) mod 3

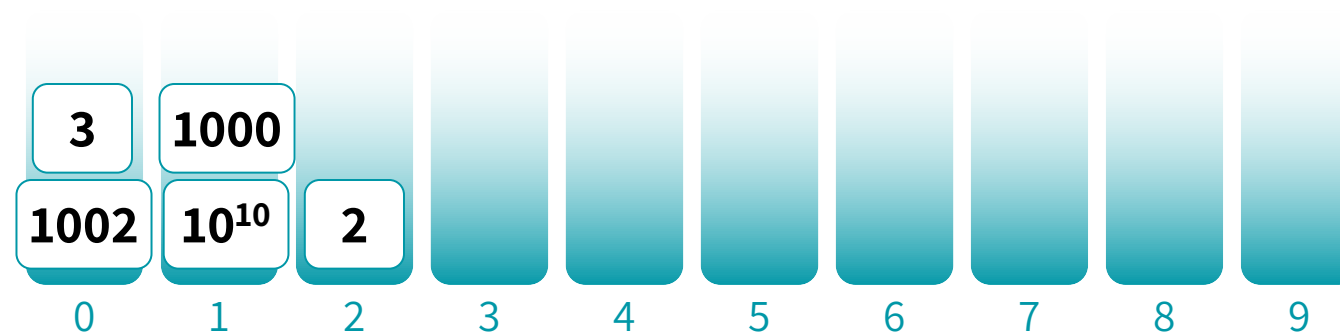


Bucketing Attempt 2

Suppose you're storing numbers from 1 - 10^{10} :



Bucket by last digit of (number * 7) mod 3



$O(1)$ INSERT:

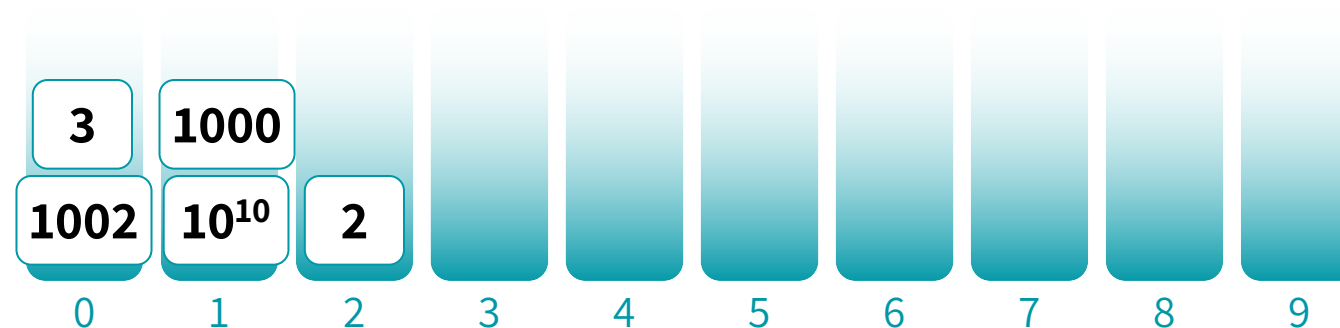
Just index into the bucket (& insert at front of a linked list)!

Bucketing Attempt 2

Suppose you're storing numbers from 1 - 10^{10} :



Bucket by last digit of (number * 7) mod 3

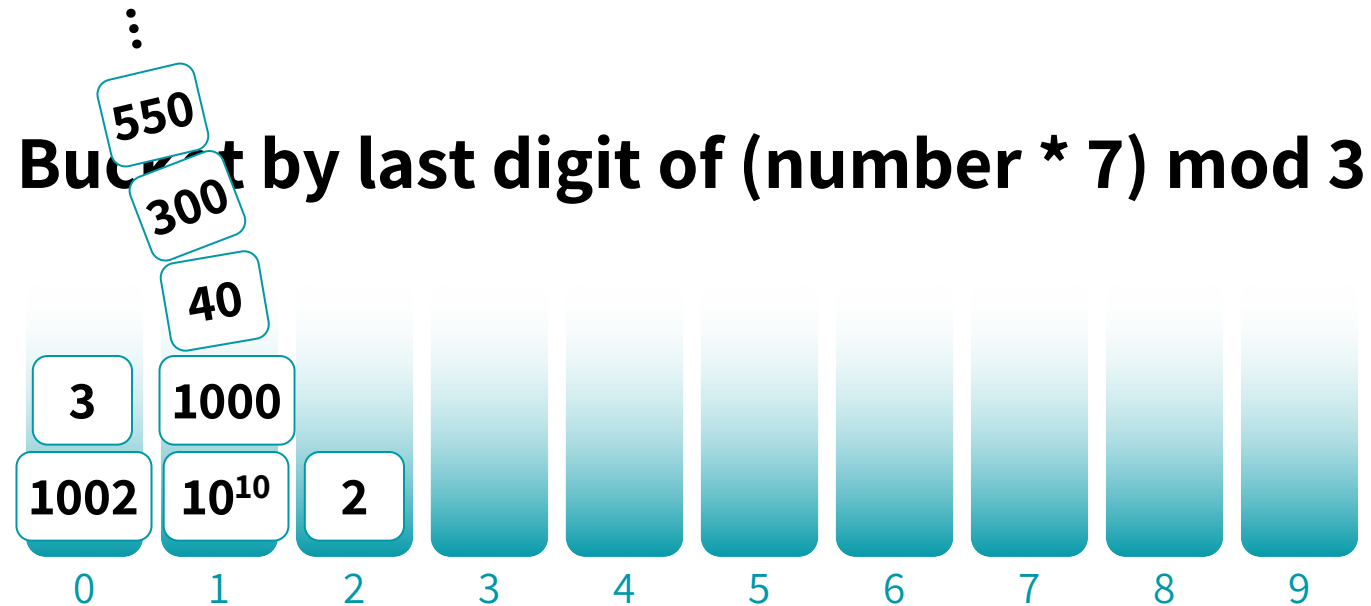


$O(\text{??????})$ SEARCH/DELETE:

Go visit bucket & search through until you find it...

Bucketing Attempt 2

Under this scheme, a bad guy could give us inputs that yields quite ugly worst-case runtimes...



$O(n)$ SEARCH/DELETE:

Go visit bucket & search through until you find it...

Bucketing Attempt 2

Under this scheme, a bad guy could give us inputs that yields quite ugly worst-case runtimes...

⋮

Seems like a bad guy could still thwart us.
There are other bucketing schemes we could use,
so to reason about them more formally,
let's talk about **HASH FUNCTIONS**.

Go visit bucket & search through until you find it...

Hash Functions

What are “good” hash functions?

Some Terminology

There exists a universe **U** of keys, with size M .

Generally, M is *really big*. Examples:

- U = the set of all ASCII strings of length 20. $M = 26^{20}$
- U = the set of all IPv4 addresses. $M = 2^{32}$
- U = the set of all possible YouTube view stats. $M = 6.8$ billion

Some Terminology

There exists a universe **U** of keys, with size M .

Generally, M is *really big*. Examples:

- U = the set of all ASCII strings of length 20. $M = 26^{20}$
- U = the set of all IPv4 addresses. $M = 2^{32}$
- U = the set of all possible YouTube view stats. $M = 6.8$ billion

Our job is to store **n** keys, and we assume $M \gg n$

- Only a few (at most n) elements of U are ever going to show up.
- We don't know which ones will show up in advance.

Some Terminology

There exists a universe **U** of keys, with size M .

Generally, M is *really big*. Examples:

- U = the set of all ASCII strings of length 20. $M = 26^{20}$
- U = the set of all IPv4 addresses. $M = 2^{32}$
- U = the set of all possible YouTube view stats. $M = 6.8$ billion

Our job is to store **n** keys, and we assume $M \gg n$

- Only a few (at most n) elements of U are ever going to show up.
- We don't know which ones will show up in advance.

A hash function **$h: U \rightarrow \{1, \dots, n\}$**
maps elements of U to buckets $1, \dots, n$

Some Terminology

There exists a universe **U** of keys, with size M .

NOTE:

For this lecture, I'm assuming that the
of elements I receive = # of buckets (both are n).

This doesn't have to be the case, but we usually aim for

#buckets = O (# elements that show up)
(otherwise, we're using "too much" space)

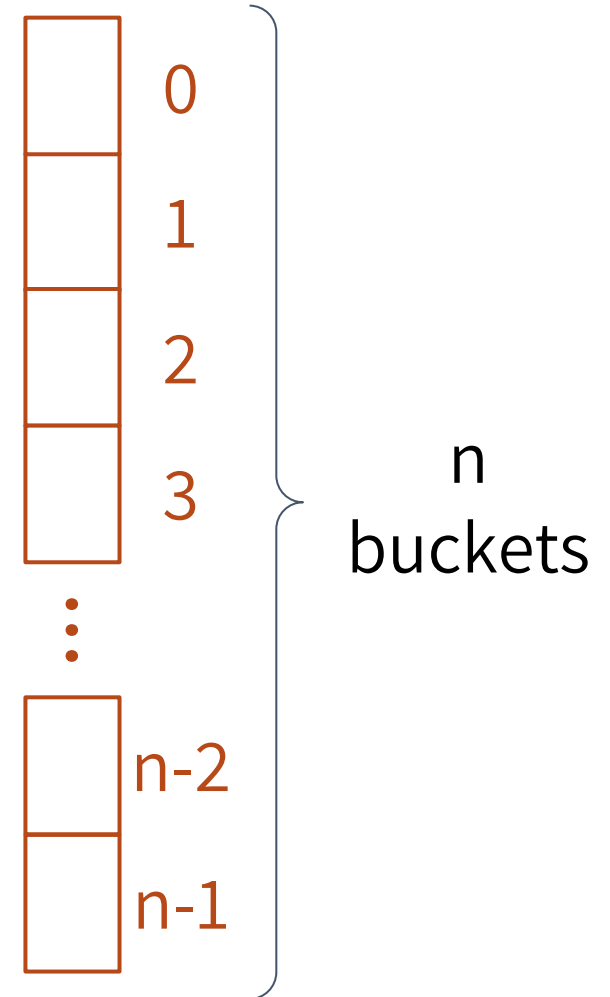
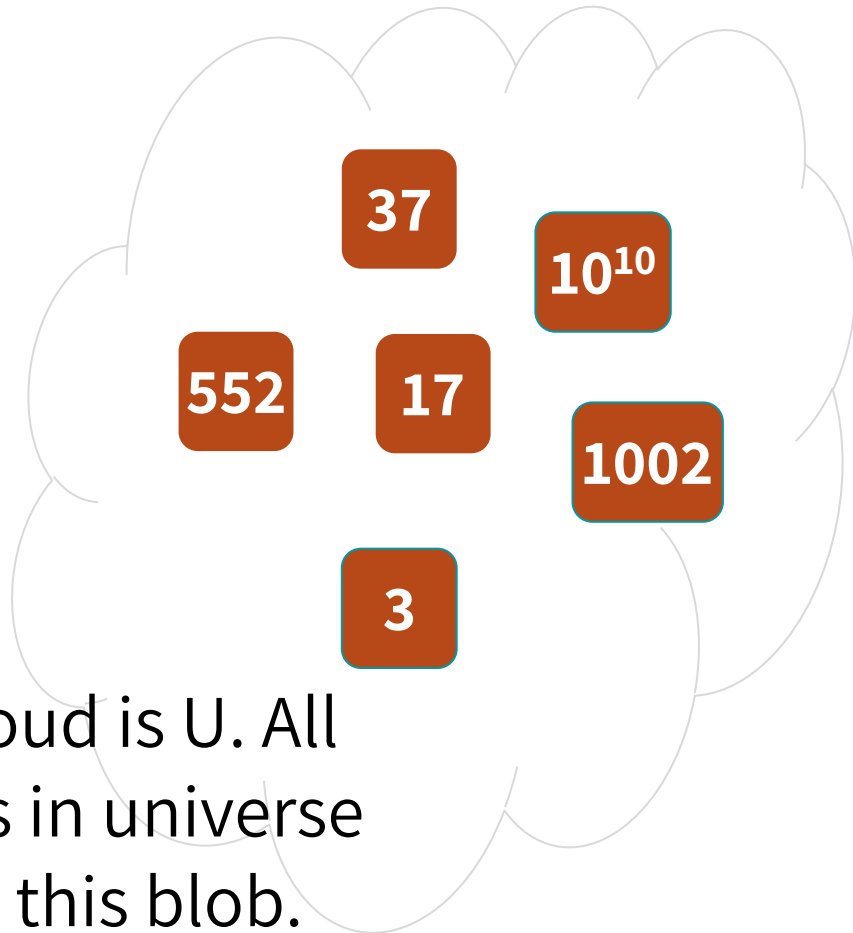
n

Some Terminology

A hash function $h: U \rightarrow \{1, \dots, n\}$
maps elements of U to buckets $1, \dots, n$

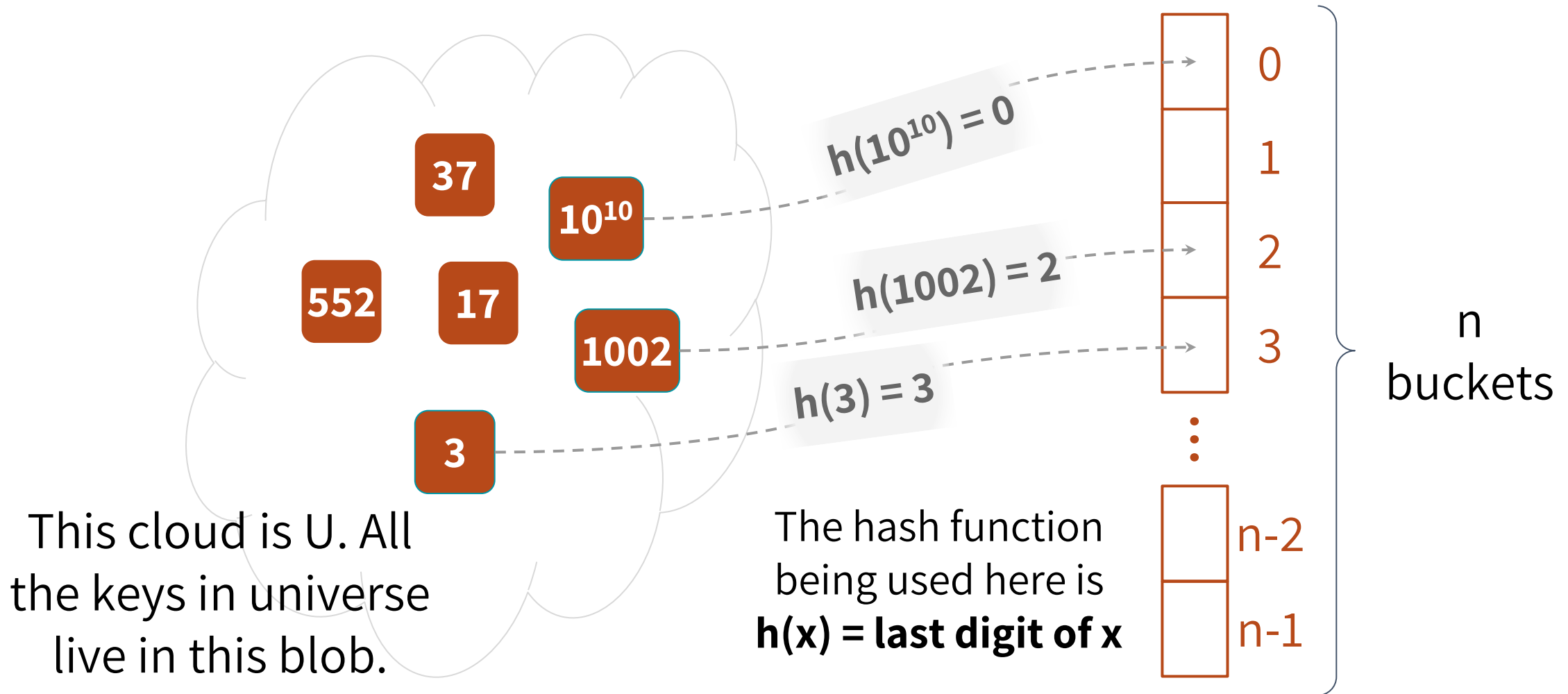
Some Terminology

A hash function $h: U \rightarrow \{1, \dots, n\}$
maps elements of U to buckets $1, \dots, n$



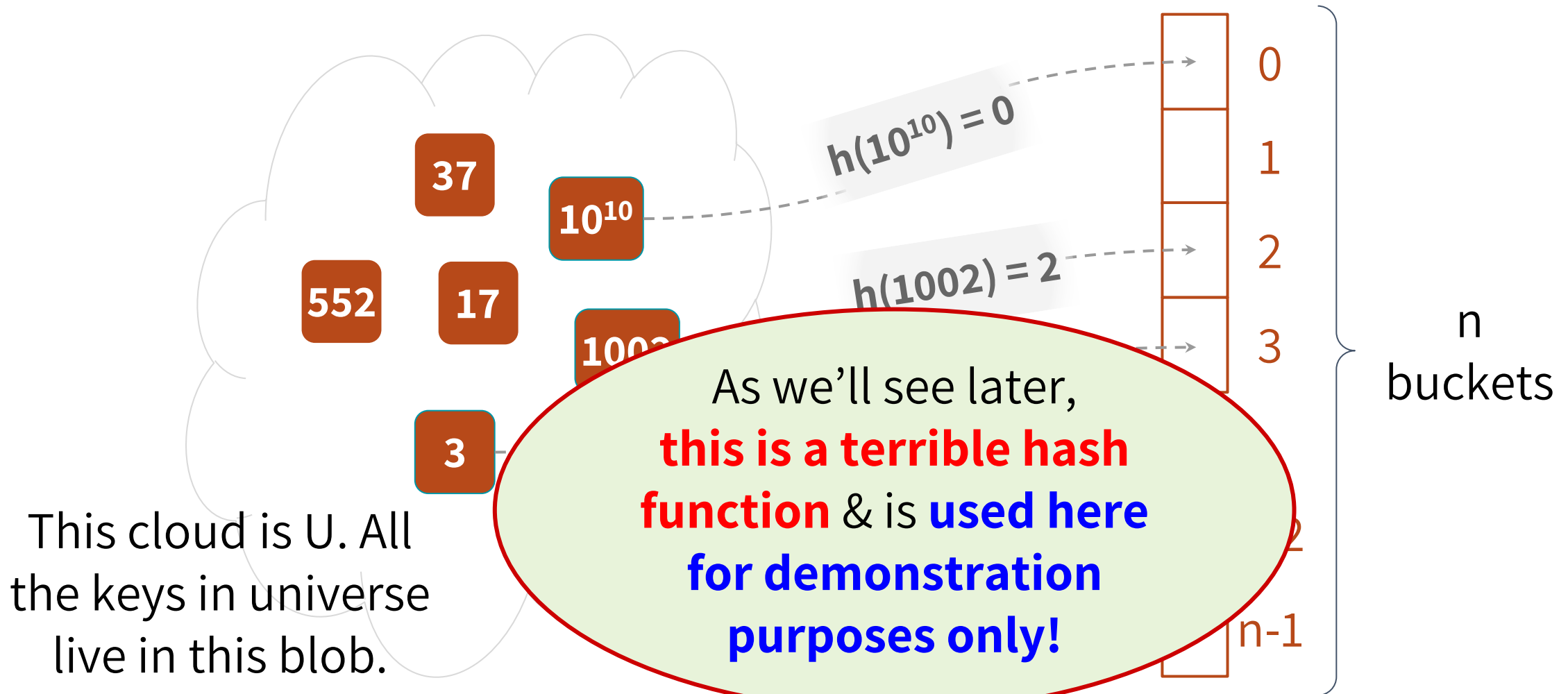
Some Terminology

A hash function $h: U \rightarrow \{1, \dots, n\}$
maps elements of U to buckets $1, \dots, n$



Some Terminology

A hash function $h: U \rightarrow \{1, \dots, n\}$
maps elements of U to buckets $1, \dots, n$



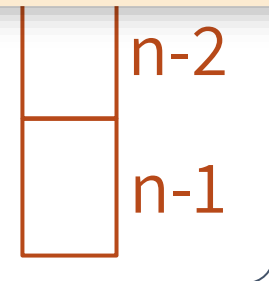
Some Terminology

A hash function $h: U \rightarrow \{1, \dots, n\}$
maps elements of U to buckets $1, \dots, n$

- A hash function tells you where to start looking for an object.
- For example, if a particular hash function h has $h(1002) = 2$, then we say “1002 *hashes to* 2”, and we go to bucket 2 to search for 1002, or insert 1002, or delete 1002.

This cloud is U . All the keys in universe live in this blob.

The hash function being used here is $h(x) = \text{last digit of } x$

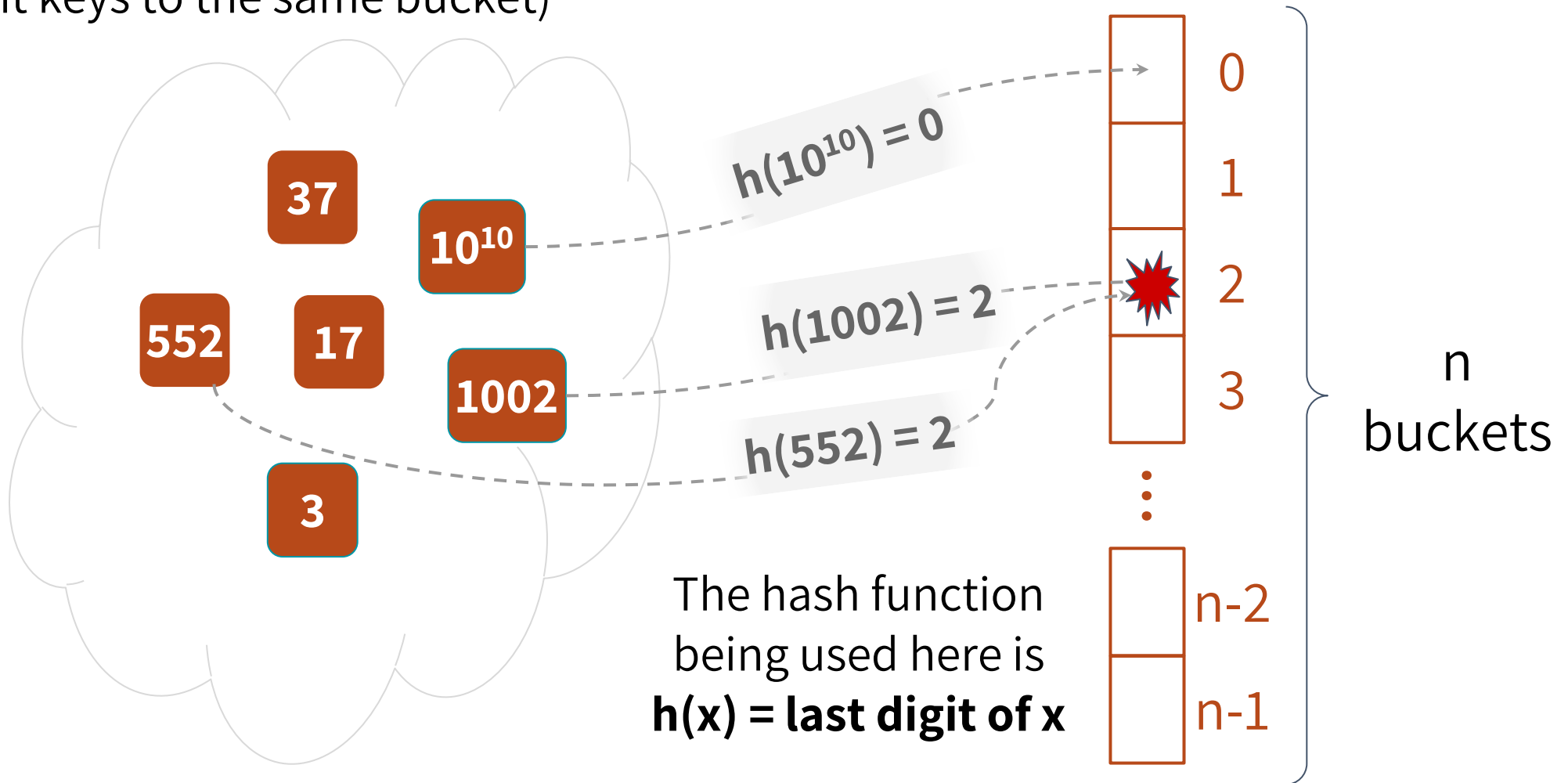


Collisions

Collisions are inevitable!

(when a hash function would map 2 different keys to the same bucket)

This is because of the *Pigeonhole Principle*. Since the size of universe $U > \#$ of buckets, every hash function (no matter how clever), suffers from at least one collision.



Collision Resolving strategies

- Separate chaining
- Open addressing
 - Linear Probing
 - Quadratic Probing
 - Double Probing

Collision Resolution: Linear Probing

- The idea:

- Table remains a simple array of size N
- On **insert(x)**, compute $f(x) \bmod N$, if the cell is full, find another by sequentially searching for the next available slot
 - Go to $f(x)+1$, $f(x)+2$ etc..
- On **find(x)**, compute $f(x) \bmod N$, if the cell doesn't match, look at $f(x)+1$, $f(x)+2$ etc..
- Linear probing function can be given by
 - $h(x, i) = (f(x) + i) \bmod N$ ($i=1,2,\dots$)

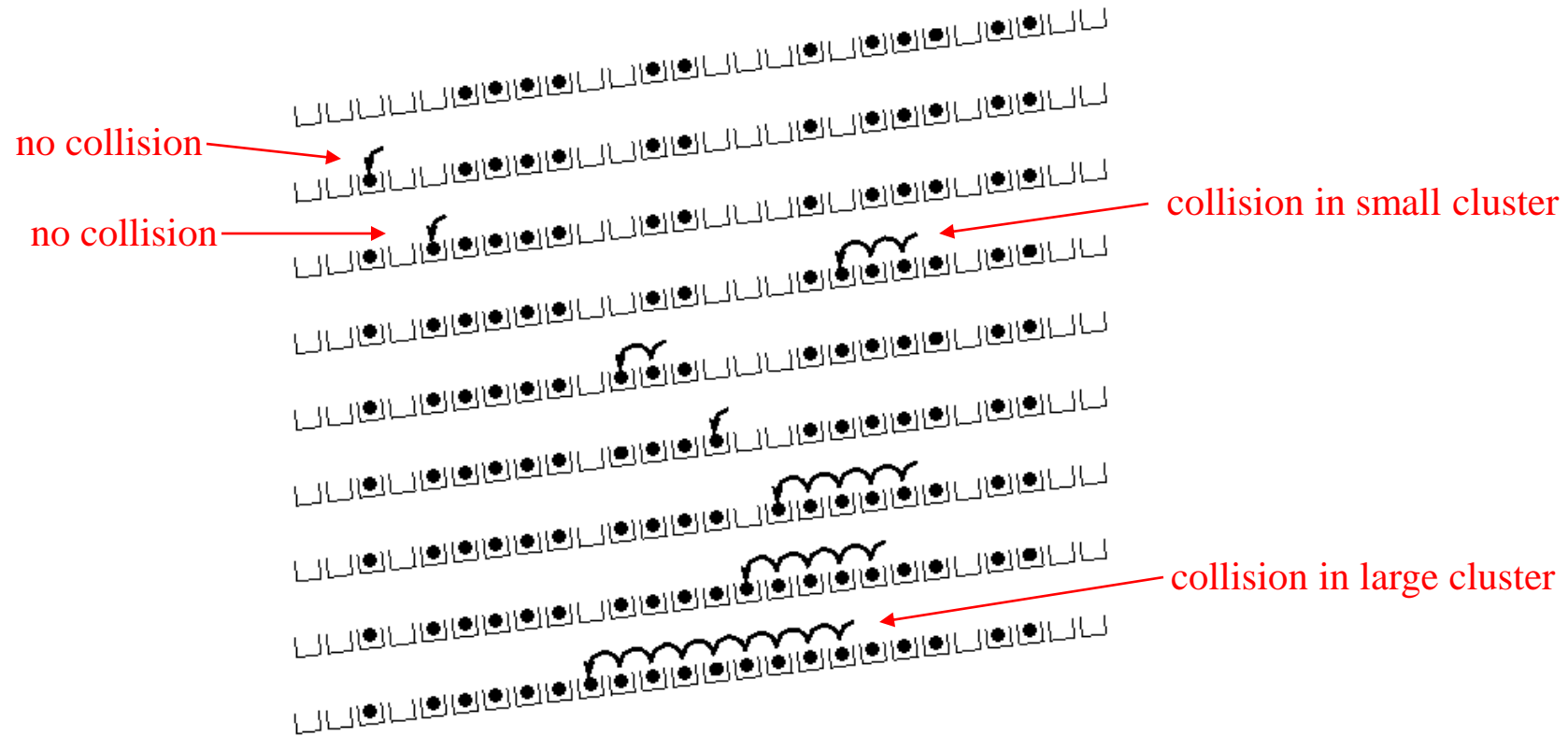
Collision Resolution: Linear Probing

hash (89, 10) = 9
hash (18, 10) = 8
hash (49, 10) = 9
hash (58, 10) = 8
hash (9, 10) = 9

	<i>After insert 89</i>	<i>After insert 18</i>	<i>After insert 49</i>	<i>After insert 58</i>	<i>After insert 9</i>
0			49	49	49
1				58	58
2					9
3					
4					
5					
6					
7					
8		18	18	18	18
9	89	89	89	89	89

Figure 20.4
Linear probing
hash table after
each insertion

Collision Resolution: Linear Probing: Clustering Problem



[R. Sedgwick]

Collision Resolution: Linear Probing: Deletion

- How about deleting items from Hash table?
 - Item in a hash table **connects** to others in the table
 - Deleting items will affect finding the others
 - “**Lazy Delete**” – Just mark the items as inactive (keep Tombstone or marker) rather than removing it.

Collision Resolution: Linear Probing: Load Factor

- The **load factor** λ of a probing hash table is the fraction of the table that is full.
- The load factor ranges from 0 (empty) to 1 (completely full).
- It is better to keep the **load factor under 0.7**
- **Double** the table size and **rehash** if load factor gets high
- Cost of Hash function $f(x)$ must be minimized
- When collisions occur, linear probing can always find an empty cell
 - But clustering can be a problem

Collision Resolution: Quadratic Probing

- Resolve collisions by examining certain cells (1,4,9,...) away from the original probe point
- Collision policy:
 - Define $h_0(k), h_1(k), h_2(k), h_3(k), \dots$ where $h_i(k) = (\text{hash}(k) + i^2) \bmod \text{size}$

Collision Resolution: Quadratic Probing

hash (89, 10) = 9
hash (18, 10) = 8
hash (49, 10) = 9
hash (58, 10) = 8
hash (9, 10) = 9

	<i>After insert 89</i>	<i>After insert 18</i>	<i>After insert 49</i>	<i>After insert 58</i>	<i>After insert 9</i>
0			49	49	49
1					
2				58	58
3					9
4					
5					
6					
7					
8		18	18	18	18
9	89	89	89	89	89

Figure 20.6

A quadratic probing hash table after each insertion (note that the table size was poorly chosen because it is not a prime number).

Collision Resolution: Quadratic Probing

- Caveat:

- May not find a vacant cell!
 - Table must be less than half full ($\lambda < 1/2$)
- (Linear probing always finds a cell.)

- Another issue

- Suppose the table size is 16.
- Probe offsets that will be tried:

1	mod 16	=	1
4	mod 16	=	4
9	mod 16	=	9
16	mod 16	=	0
25	mod 16	=	9
36	mod 16	=	4
49	mod 16	=	1
64	mod 16	=	0
81	mod 16	=	1

only four different values!

Collision Resolution: Double Probing

$$f(i) = i * g(k)$$

where g is a second hash function

- Probe sequence:

$$0^{\text{th}} \text{ probe} = h(k) \text{ mod size}$$

$$1^{\text{th}} \text{ probe} = (h(k) + g(k)) \text{ mod size}$$

$$2^{\text{th}} \text{ probe} = (h(k) + 2 * g(k)) \text{ mod size}$$

$$3^{\text{th}} \text{ probe} = (h(k) + 3 * g(k)) \text{ mod size}$$

...

$$i^{\text{th}} \text{ probe} = (h(k) + i * g(k)) \text{ mod size}$$

Collision Resolution: Double Probing

Double Hashing Example

$$h(k) = k \bmod 7 \text{ and } g(k) = 5 - (k \bmod 5)$$

	76	93	40	47	10	55
0						
1				47	47	47
2		93	93	93	93	93
3					10	10
4						55
5			40	40	40	40
6	76	76	76	76	76	76
Probes	1	1	1	2	1	2

Collision Resolution: Open Addressing: Rehashing

When the table gets too full, create a bigger table (usually 2x as large) and hash all the items from the original table into the new table.

- When to rehash?
 - half full ($\lambda = 0.5$)
 - when an insertion fails
 - some other threshold
- Cost of rehashing?

Collision Resolution: Open Addressing: Conclusion

- Open addressing in hash tables offers an average-case time complexity of $O(1)$ for insertion, deletion, and search operations, assuming a **low load factor**.
- Performance degrades with a worst-case complexity of $O(n)$ when the **table is nearly full**, causing high clustering.

Collision Resolution: Chaining

To resolve collisions, one common method is to use **chaining!**

We're just giving a formal name to our bucketing example from earlier:

represent each bucket's contents as a ***linked list!***

(Another method is called "Open Addressing")

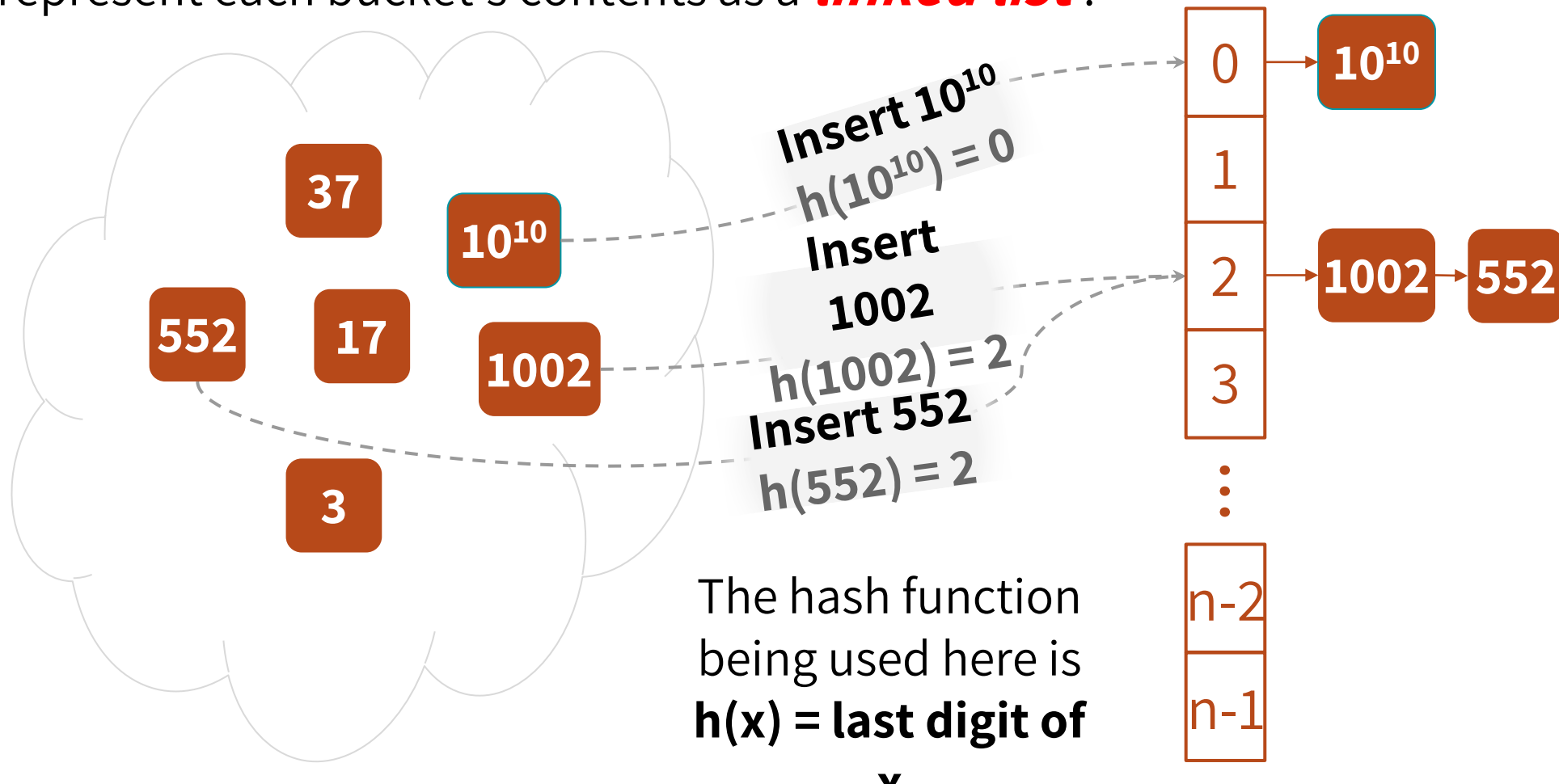
Collision Resolution: Chaining

To resolve collisions, one common method is to use **chaining!**

We're just giving a formal name to our bucketing example from earlier:

represent each bucket's contents as a **linked list!**

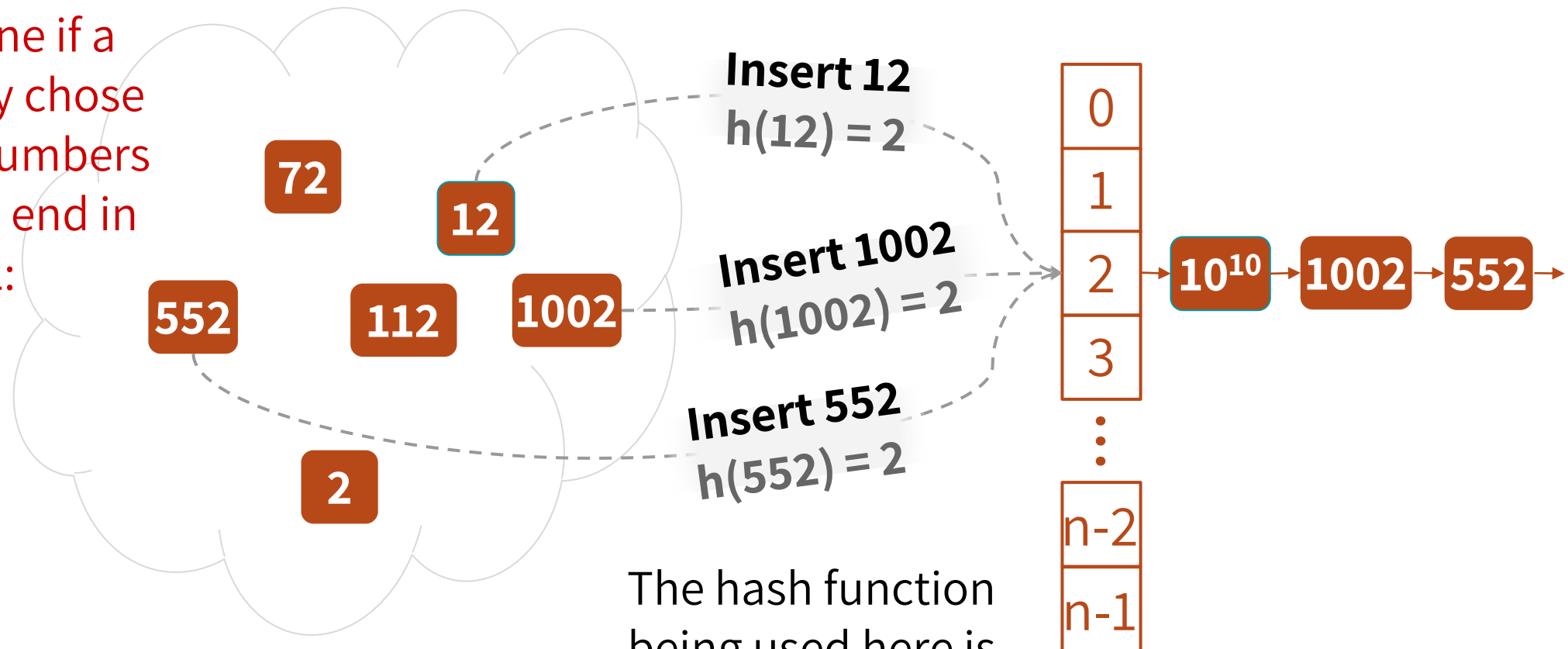
(Another method is called "Open Addressing")



Collision Resolution: Chaining

But if the items are all clumped together in a single bucket, SEARCH/DELETE may be very slow because of the linked list traversal...

Imagine if a bad guy chose these numbers that all end in 2:



The hash function being used here is $h(x) = \text{last digit of } x$

Hash Table Goals

Remember worst-case analysis:

OUR GOAL: Design a function $h: \mathbf{U} \rightarrow \{1, \dots, n\}$ so that no matter what n items of \mathbf{U} a bad guy chooses & the operations they choose to perform, the buckets will be balanced.

(Here, balanced means $O(1)$ entries per bucket)

**Then we'd achieve our dream of $O(1)$
INSERT/DELETE/SEARCH.**

Hash Table Goals

Remember worst-case analysis:

OUR GOAL: Design a function $h: \mathbf{U} \rightarrow \{1, \dots, n\}$ so that no matter what n items of \mathbf{U} a bad guy chooses & the operations they choose to perform, the buckets will be balanced.

(Here, balanced means $O(1)$ entries per bucket)

**Then we'd achieve our dream of $O(1)$
INSERT/DELETE/SEARCH.**

Can you come up with such a function?

Hash Table Goals

OUR GOAL: Design a function $h: \mathbf{U} \rightarrow \{1, \dots, n\}$ so that no matter what n items of \mathbf{U} a bad guy chooses, the buckets will be balanced (have $O(1)$ size)

Can you come up with such a function? **No.**

No *deterministic* hash function can defeat worst-case input!

Hash Table Goals

OUR GOAL: Design a function $h: \mathbf{U} \rightarrow \{1, \dots, n\}$ so that no matter what n items of \mathbf{U} a bad guy chooses, the buckets will be balanced (have $O(1)$ size)

Can you come up with such a function? **No.**

No *deterministic* hash function can defeat worst-case input!

- The universe \mathbf{U} has \mathbf{M} items
- They get hashed into n buckets
- At least 1 bucket has at least \mathbf{M}/n items hashed to it (Pigeonhole)
- \mathbf{M} is wayyyy bigger than n , so \mathbf{M}/n is bigger than n

The n items the bad guy chooses are items that all land in this very full bucket. That bucket has size $\Omega(n)$.

Hash Table Goals

OUR GOAL: Design a function $h: \mathbf{U} \rightarrow \{1, \dots, n\}$ so that no matter what n items of \mathbf{U} a bad guy chooses, the buckets will be balanced (have $O(1)$ size)

Can you come up with such a function? **No.**

No *deterministic* hash function can defeat worst-case input!

- The universe \mathbf{U} has \mathbf{M} items
- They get hashed into n buckets
- At least 1 bucket has at least \mathbf{M}/n items hashed to it (Pigeonhole)
- \mathbf{M} is wayyyy bigger than n , so \mathbf{M}/n is bigger than n

The n items the bad guy chooses are items that all land in this very full bucket. That bucket has size $\Omega(n)$.

The problem is that the bad guy knows our hash function beforehand.

Hash Table Goals

OUR GOAL: Design a function $h: \mathbf{U} \rightarrow \{1, \dots, n\}$ so that no matter what n items of \mathbf{U} a bad guy chooses, the buckets will be balanced (have $O(1)$ size)

Maybe there's a way to weaken the adversary...

LET'S BRING IN SOME

RANDOMNESS!

- The univer
- They get h
- At least 1 b
- **M** is wayyy

The problem is that the bad guy knows our hash function beforehand.

Hash Functions and Randomness

What it means to weaken the adversary & ways to do it

Intuition

So, our strategy is to define a set of hash functions, and then we randomly choose a hash function h from this set to use!

You can think of it like a game:

1. You announce your set of hash functions, H .
2. The adversary chooses n items for your hash function to hash.
3. You then randomly pick a hash function h from H to hash the n items.

Intuition

So, our strategy is to define a set of hash functions, and then we randomly choose a hash function h from this set to use!

You can think of it like a game:

1. You announce your set of hash functions, H .
2. The adversary chooses n items for your hash function to hash.
3. You then randomly pick a hash function h from H to hash the n items.

What would make a “good” set of hash functions H ?

Intuition

What would make a “good” set of hash functions H ?

What we want

Design a set $\mathbf{H} = \{h_1, h_2, h_3, \dots, h_k\}$ where $h_i : U \rightarrow \{1, \dots, n\}$, such that if we chose a random \mathbf{h} in \mathbf{H} and after an adversary chooses \mathbf{n} items $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\}$ to hash,

for any item \mathbf{u}_i ,
the **expected** # of items in \mathbf{u}_i 's bucket is **$O(1)$**

Let's see an example of a set of hash functions H that achieves this goal!

H = Exhaustive Set of All Hash Functions

WHAT WE WANT:

Design a set $\mathbf{H} = \{h_1, h_2, h_3, \dots, h_k\}$ where $h_i : U \rightarrow \{1, \dots, n\}$, such that if we chose a uniformly random \mathbf{h} in \mathbf{H} and after an adversary chooses \mathbf{n} items $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\}$ to hash,

for any item \mathbf{u}_i ,

the **expected** # of items in \mathbf{u}_i 's bucket is **$O(1)$**

\mathbf{H} = the exhaustive set of all hash functions that map elements in the universe U to buckets 1 to n .

\mathbf{H} contains a total of n^M hash functions.

H = Exhaustive Set of All Hash Functions

WHAT WE WANT:

Design a set $\mathbf{H} = \{h_1, h_2, h_3, \dots, h_k\}$ where $h_i : U \rightarrow \{1, \dots, n\}$, such that if we chose a uniformly random \mathbf{h} in \mathbf{H} and after an adversary chooses \mathbf{n} items $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\}$ to hash,

for any item \mathbf{u}_i ,

the **expected** # of items in \mathbf{u}_i 's bucket is **$O(1)$**

\mathbf{H} = the exhaustive set of all hash functions that map elements in the universe U to buckets 1 to n .

\mathbf{H} contains a total of n^M hash functions.

Here is an example

where

$\mathbf{U} = \{\text{"a"}, \text{"b"}, \text{"c"}\}$

so $\mathbf{M} = 3$. Also, we have $\mathbf{n} = 2$.

	h_1	h	h	h	h	h	h
"a"	0	0	0	0	1	1	1
"b"	0	0	1	1	0	0	1
"c"	0	1	0	1	0	1	0

The 0's and 1's represent the binary buckets i.e. h_8 will hash "b" to bucket 1.

H = Exhaustive Set of All Hash Functions

H = the exhaustive set of all hash functions that map elements in the universe U to buckets 1 to n .

H contains a total of n^M hash functions.

$\mathbb{E}[\# \text{ of items in } u_i \text{ 's bucket}] =$

H = Exhaustive Set of All Hash Functions

H = the exhaustive set of all hash functions that map elements in the universe U to buckets 1 to n .

H contains a total of n^M hash functions.

$$\mathbb{E}[\# \text{ of items in } u_i \text{ 's bucket}] = \sum_{j=1}^n P[h(u_i) = h(u_j)]$$

This probability is taken over the random choice of hash function!

H = Exhaustive Set of All Hash Functions

H = the exhaustive set of all hash functions that map elements in the universe U to buckets 1 to n .

H contains a total of n^M hash functions.

$$\begin{aligned}\mathbb{E}[\# \text{ of items in } u_i \text{ 's bucket}] &= \sum_{j=1}^n P[h(u_i) = h(u_j)] \\ &= P[h(u_i) = h(u_i)] + \sum_{j \neq i} P[h(u_i) = h(u_j)]\end{aligned}$$

This probability is taken over the random choice of hash function!

H = Exhaustive Set of All Hash Functions

H = the exhaustive set of all hash functions that map elements in the universe U to buckets 1 to n .

H contains a total of n^M hash functions.

$$\begin{aligned}\mathbb{E}[\# \text{ of items in } u_i \text{ 's bucket}] &= \sum_{j=1}^n P[h(u_i) = h(u_j)] \leftarrow \text{This probability is taken over the random choice of hash function!} \\ &= P[h(u_i) = h(u_i)] + \sum_{j \neq i} P[h(u_i) = h(u_j)] \\ &= 1 + \sum_{j \neq i} P[h(u_i) = h(u_j)]\end{aligned}$$

H = Exhaustive Set of All Hash Functions

H = the exhaustive set of all hash functions that map elements in the universe U to buckets 1 to n .

H contains a total of n^M hash functions.

$$\begin{aligned}\mathbb{E}[\# \text{ of items in } u_i \text{ 's bucket}] &= \sum_{j=1}^n P[h(u_i) = h(u_j)] \\ &= P[h(u_i) = h(u_i)] + \sum_{j \neq i} P[h(u_i) = h(u_j)] \\ &= 1 + \sum_{j \neq i} P[h(u_i) = h(u_j)] \\ &= 1 + \sum_{j \neq i} \frac{1}{n}\end{aligned}$$

This probability is taken over the random choice of hash function!

How do we know that $P[h(u_i) = h(u_j)] = 1/n$?

H = Exhaustive Set of All Hash Functions

H = the exhaustive set of all hash functions that map elements in the universe U to buckets 1 to n .

H contains a total of n^M hash functions.

$$\begin{aligned}\mathbb{E}[\# \text{ of items in } u_i \text{ 's bucket}] &= \sum_{j=1}^n P[h(u_i) = h(u_j)] \\ &= P[h(u_i) = h(u_i)] + \sum_{j \neq i} P[h(u_i) = h(u_j)] \\ &= 1 + \sum_{j \neq i} P[h(u_i) = h(u_j)] \\ &= 1 + \sum_{j \neq i} \frac{1}{n} \\ &= 1 + \frac{n-1}{n}\end{aligned}$$

This probability is taken over the random choice of hash function!

How do we know that $P[h(u_i) = h(u_j)] = 1/n$?

H = Exhaustive Set of All Hash Functions

H = the exhaustive set of all hash functions that map elements in the universe U to buckets 1 to n .

H contains a total of n^M hash functions.

$$\begin{aligned}\mathbb{E}[\# \text{ of items in } u_i \text{ 's bucket}] &= \sum_{j=1}^n P[h(u_i) = h(u_j)] \\ &= P[h(u_i) = h(u_i)] + \sum_{j \neq i} P[h(u_i) = h(u_j)] \\ &= 1 + \sum_{j \neq i} P[h(u_i) = h(u_j)] \\ &= 1 + \sum_{j \neq i} \frac{1}{n} \\ &= 1 + \frac{n-1}{n} \leq 2\end{aligned}$$

This probability is taken over the random choice of hash function!

How do we know that $P[h(u_i) = h(u_j)] = 1/n$?

$O(1)$
This is what we wanted!

H = Exhaustive Set of All Hash Functions

H = the exhaustive set of all hash functions that map elements in the universe U to buckets 1 to n

Good News:

H achieves our goal!

If we choose a *uniformly random hash function* from Exhaustive Set of All Hash Functions, then INSERT/DELETE/SEARCH on any n elements will have **expected runtime of $O(1)$** .

H = Exhaustive Set of All Hash Functions

H = the exhaustive set of all hash functions that map elements in the universe U to buckets 1 to n

Bad News:

How many bits does it take to store a uniformly random hash function?

A lot!

How many bits does it take to store a uniformly random hash function?

We'd use a lookup table: one entry per element of U , each storing which bucket to hash that element to.

$(M \text{ elements}) * (\log(n) \text{ bits to write down a bucket \#}) = M \log n \text{ bits}$
This is HUGE... (& enough to do direct addressing!)

How many bits does it take to store a uniformly random hash function?

We'd use a lookup table: one entry per element of U , each storing which bucket to hash that element to.

(M elements) * ($\log(n)$ bits to write down a bucket #) = $M \log n$ bits
This is HUGE... (& enough to do direct addressing!)

How do we fix this size issue?

Universal Hash Families

“Good” sets of hash functions that aren’t as large!

What we wanted

H = the exhaustive set of all hash functions that map elements in the universe U to buckets 1 to n .

H contains a total of n^M hash functions.

$$\begin{aligned}\mathbb{E}[\# \text{ of items in } u_i \text{ 's bucket}] &= \sum_{j=1}^n P[h(u_i) = h(u_j)] \\ &= P[h(u_i) = h(u_i)] + \sum_{j \neq i} P[h(u_i) = h(u_j)]\end{aligned}$$

The fact that
 $P[h(u_i)=h(u_j)] = 1/n$
did all the work
here

$$\begin{aligned}&= 1 + \sum_{j \neq i} P[h(u_i) = h(u_j)] \\ &= 1 + \sum_{j \neq i} \frac{1}{n} \\ &= 1 + \frac{n-1}{n} \leq 2\end{aligned}$$

$O(1)$
This is what
we wanted!

What we wanted

H = the exhaustive set of all hash functions that map elements in the universe U to buckets 1 to n .

H contains a total of n^M hash functions.

The exhaustive set of all hash functions achieved our goal but was way too big, so let's pick **h** from

a *smaller* hash family where

$$P[h(u_i) = h(u_j)] \leq 1/n$$

T
P[h(
dic

$$= 1 + \frac{1}{n} \leq 2$$

Universal Hash Family

A **hash family** is a fancy name for a set of hash functions.

A hash family \mathbf{H} is a **universal hash family** if, when h is chosen uniformly at random from \mathbf{H} ,

for all $u_i, u_j \in U$ with $u_i \neq u_j$,

$$P_{h \in H} [h(u_i) = h(u_j)] \leq \frac{1}{n}$$

Universal Hash Family

A **hash family** is a fancy name for a set of hash functions.

A hash family \mathbf{H} is a **universal hash family** if, when \mathbf{h} is chosen uniformly at random from \mathbf{H} ,

for all $u_i, u_j \in U$ with $u_i \neq u_j$,

$$P_{h \in H} [h(u_i) = h(u_j)] \leq \frac{1}{n}$$

Then if we randomly choose \mathbf{h} from a universal hash family \mathbf{H} , we'll be guaranteed that:

$$E[\# \text{ of items in } u_i\text{'s bucket}] \leq 2 = O(1)$$

Flashback of the Math

A hash family \mathbf{H} is a **universal hash family** if, when h is chosen uniformly at random from \mathbf{H} ,

for all $u_i, u_j \in U$ with $u_i \neq u_j$,

$$P_{h \in H} [h(u_i) = h(u_j)] \leq \frac{1}{n}$$

$$\begin{aligned} \mathbb{E}[\# \text{ of items in } u_i \text{ 's bucket}] &= \sum_{j=1}^n P[h(u_i) = h(u_j)] \\ &= P[h(u_i) = h(u_i)] + \sum_{j \neq i} P[h(u_i) = h(u_j)] \end{aligned}$$

This inequality is now what a universal hash family guarantees!

$$= 1 + \sum_{j \neq i} P[h(u_i) = h(u_j)]$$

$$\leq 1 + \sum_{j \neq i} \frac{1}{n}$$

$$= 1 + \frac{n-1}{n} \leq 2$$

O(1)

This is what we wanted!

A Small Universal Hash Family?

Are there smaller ones universal hash families?

A Non-Example

$H = \{h_0, h_1\}$ where

$h_0 = \text{MOST_SIGNIFICANT_DIGIT}$

$h_1 = \text{LEAST_SIGNIFICANT_DIGIT}$

Why is this not a universal hash family?

A Non-Example

$\mathbf{H} = \{\mathbf{h}_0, \mathbf{h}_1\}$ where

$\mathbf{h}_0 = \text{MOST_SIGNIFICANT_DIGIT}$

$\mathbf{h}_1 = \text{LEAST_SIGNIFICANT_DIGIT}$

Why is this not a universal hash family?

$$P_{h \in H} [h(153) = h(173)] = 1 > \frac{1}{n}$$

There's a $\frac{1}{2}$ probability of choosing \mathbf{h}_0 , and $\mathbf{h}_0(153) = \mathbf{h}_0(173) = \mathbf{bucket}$

1

There's a $\frac{1}{2}$ probability of choosing \mathbf{h}_1 , and $\mathbf{h}_1(153) = \mathbf{h}_1(173) = \mathbf{bucket}$

3

Probability that a randomly chosen \mathbf{h} from \mathbf{H} collides 153 & 173 is 1!

An Example

Here is one of the more well-studied universal hash families:

Pick a prime $p \geq M$

Define $h_{a,b}(x) = ((ax + b) \bmod p) \bmod n$

$H = \{ h_{a,b} : a \in \{1, \dots, p - 1\}, b \in \{0, \dots, p - 1\} \}$

An Example

Here is one of the more well-studied universal hash families:

Pick a prime $p \geq M$

Define $h_{a,b}(x) = ((ax + b) \bmod p) \bmod n$

$$H = \{ h_{a,b} : a \in \{1, \dots, p - 1\}, b \in \{0, \dots, p - 1\} \}$$

Example: Suppose $n = 3$, and $p = 5$. Here's $h_{2,4}$:

$$h_{2,4}(1) = ((2 \cdot 1 + 4) \bmod 5) \bmod 3 = (6 \bmod 5) \bmod 3 = 1 \bmod 3 = \mathbf{1}$$

$$h_{2,4}(4) = ((2 \cdot 4 + 4) \bmod 5) \bmod 3 = (12 \bmod 5) \bmod 3 = 2 \bmod 3 = \mathbf{2}$$

$$h_{2,4}(3) = ((2 \cdot 3 + 4) \bmod 5) \bmod 3 = (6 \bmod 5) \bmod 3 = 1 \bmod 3 = \mathbf{1}$$

An Example

Here is one of the more well-studied universal hash families:

Pick a prime $p \geq M$

Define $h_{a,b}(x) = ((ax + b) \bmod p) \bmod n$

$$\mathbf{H} = \{ h_{a,b} : a \in \{1, \dots, p - 1\}, b \in \{0, \dots, p - 1\} \}$$

To draw a hash function h from \mathbf{H} :

Pick a random a
in $\{1, \dots, p - 1\}$.

&

Pick a random b
in $\{0, \dots, p - 1\}$.

An Example

Here is one of the more well-studied universal hash families:

To store $h_{a,b}$, you just need to store two numbers: **a** and **b!**

Since **a** and **b** are at most $p-1$, we need **$\sim 2 \cdot \log(p)$ bits.**

p is a prime that's close-ish to M , so this means the space needed =

$O(\log M)$

This is so much better than $O(M \log n)$!

$\{1, \dots, p-1\}$.

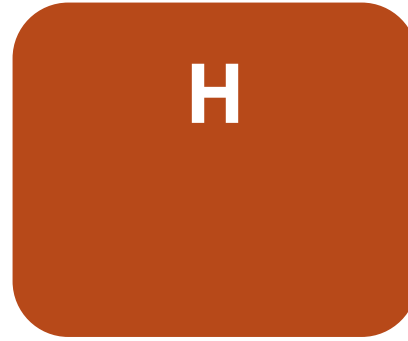
$\{0, \dots, p-1\}$.

Hash Tables

Putting everything together, what's the scheme?

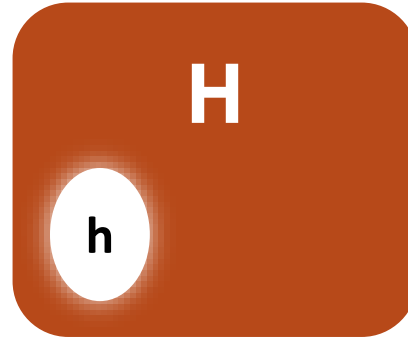
The Whole Scheme

You choose your set of hash functions **H**, a universal hash family like $H = \text{mod } p \text{ mod } n$.



The Whole Scheme

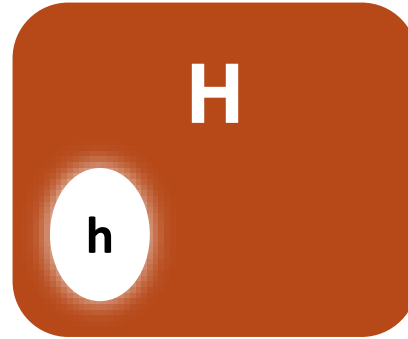
You choose your set of hash functions **H**, a universal hash family like $H = \text{mod } p \text{ mod } n$.



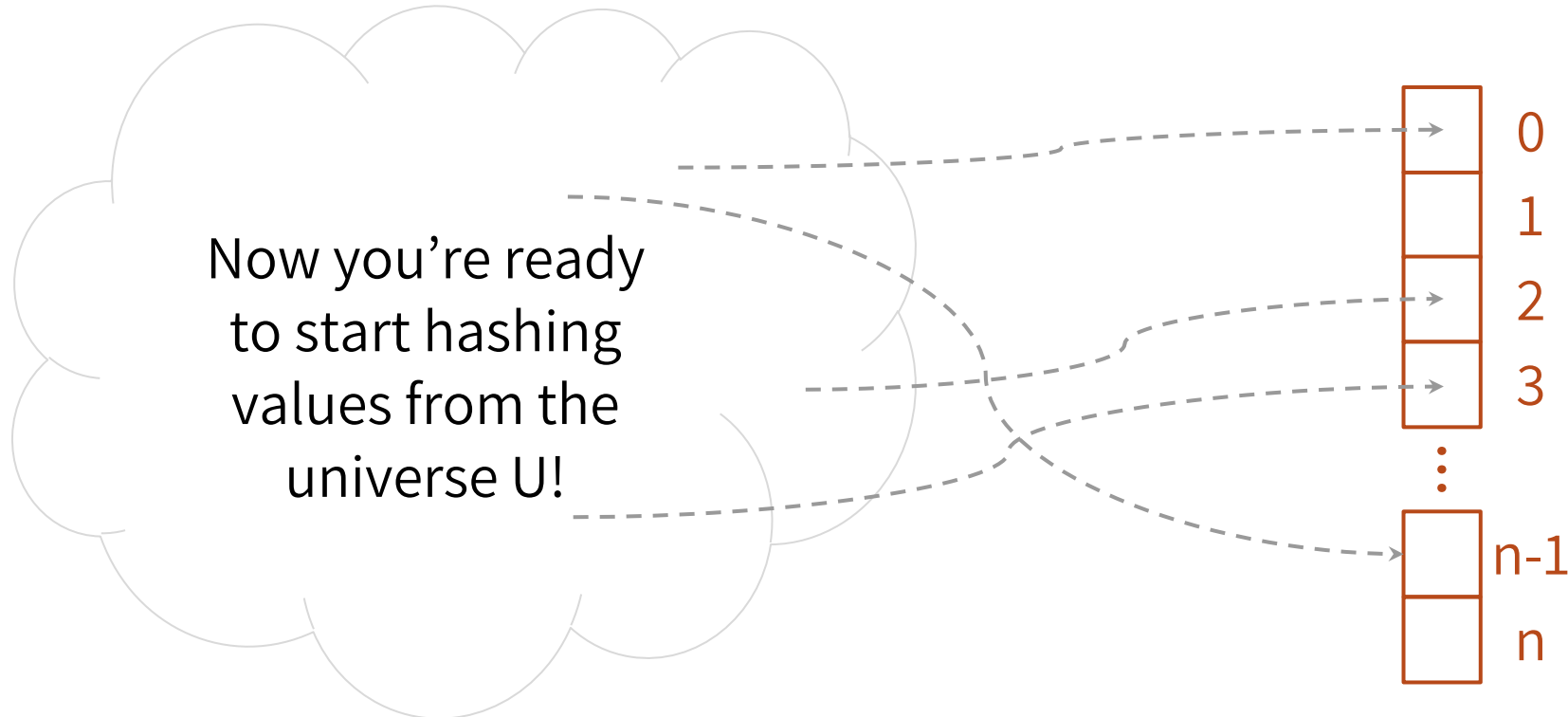
When the client initializes a hash table, randomly pick a hash function **h** from **H** to use in the hash table to hash the items.

The Whole Scheme

You choose your set of hash functions \mathbf{H} , a universal hash family like $H = \text{mod } p \text{ mod } n$.

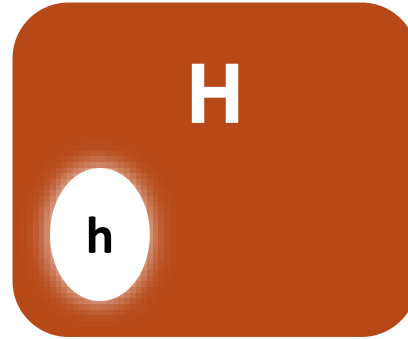


When the client initializes a hash table, randomly pick a hash function h from \mathbf{H} to use in the hash table to hash the items.

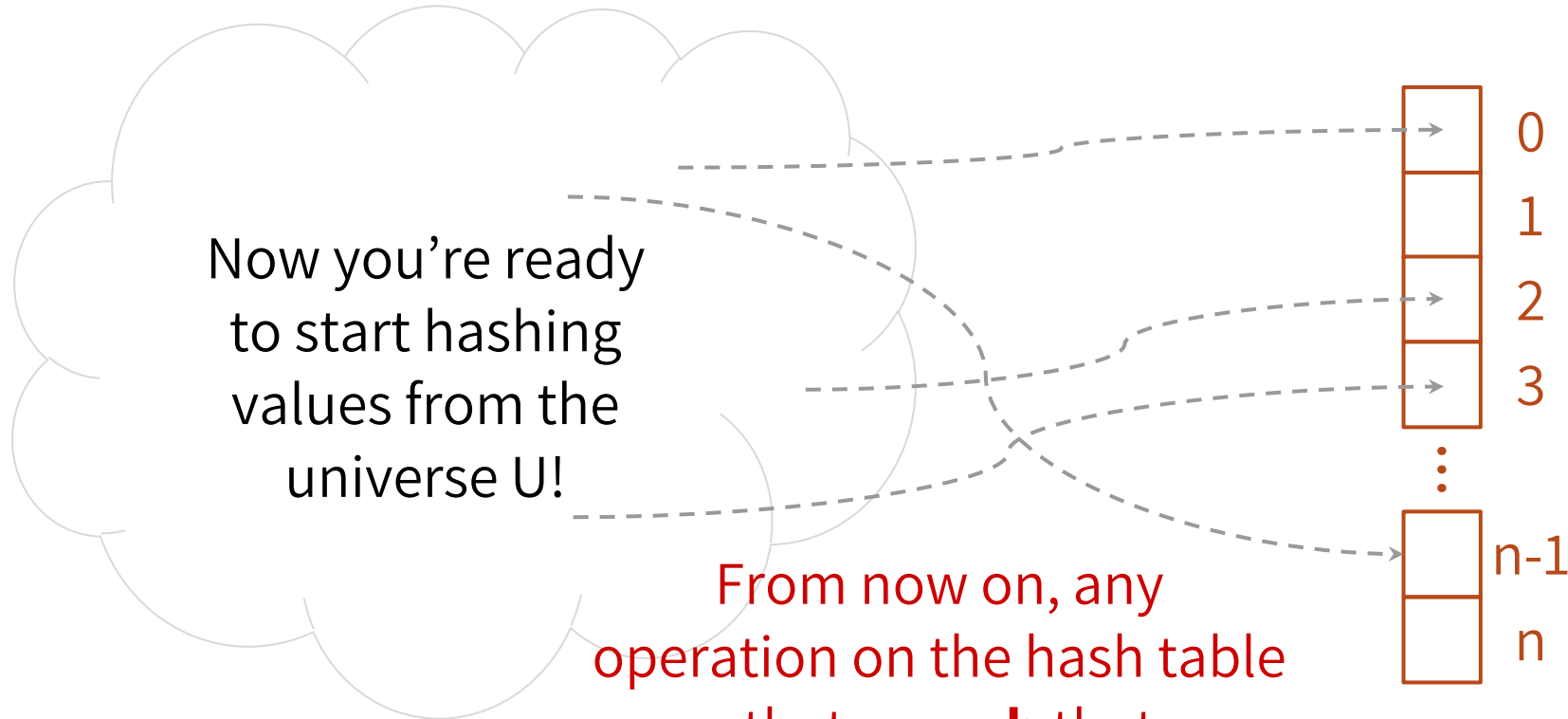


The Whole Scheme

You choose your set of hash functions \mathbf{H} , a universal hash family like $H = \text{mod } p \text{ mod } n$.



When the client initializes a hash table, randomly pick a hash function \mathbf{h} from \mathbf{H} to use in the hash table to hash the items.



From now on, any operation on the hash table uses that same \mathbf{h} that you randomly selected from \mathbf{H}

We can now expect that these buckets will be pretty balanced

Hash Table: Motivation

OPERATION	SORTED ARRAY	UNSORTED LINKED LIST	HASH TABLES (HOPEFULLY)
SEARCH	$O(\log(n))$	$O(n)$	$O(1)$
DELETE	$O(n)$	$O(n)$	$O(1)$
INSERT	$O(n)$	$O(1)$	$O(1)$

*** Assuming we implement it cleverly with a “good” hash function**

Acknowledgement

- Stanford University

Thank You