

Data Structures and Algorithms

Hashing

CS 225

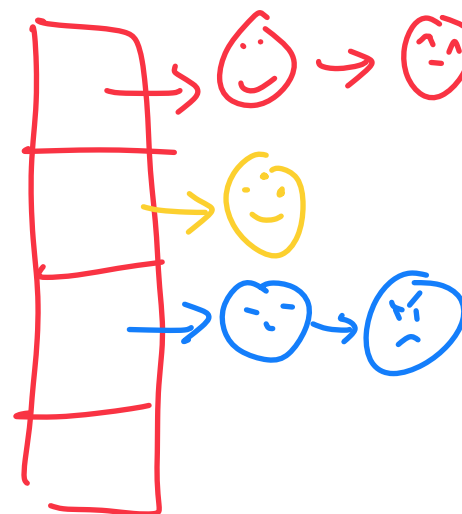
October 27, 2023

Brad Solomon



UNIVERSITY OF
ILLINOIS
URBANA - CHAMPAIGN

Department of Computer Science



Extra Credit Projects — Submit by 10/31

"Approved by Nov 1"

↳ Submitted by 10/31 → Ideally
10/29

Its fine!

Randomization in Algorithms

1. Assume input data is random to estimate average-case performance

2. Use randomness inside algorithm to estimate expected running time

↳ In the slides!

3. Use randomness inside algorithm to approximate solution in fixed time

Learning Objectives

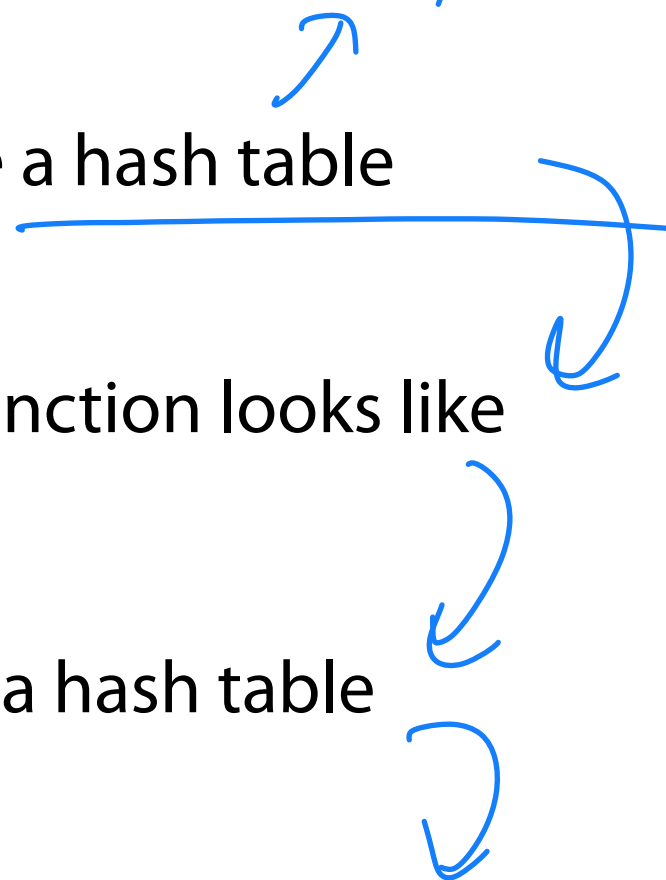
My fav core D.S.

Motivate and formally define a hash table

Discuss what a 'good' hash function looks like

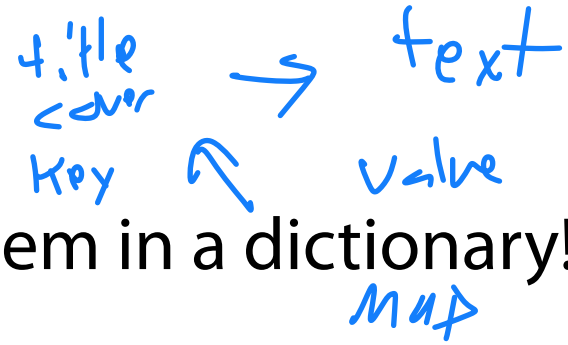
Identify the key weakness of a hash table

Introduce strategies to "correct" this weakness



Data Structure Review

I have a collection of books and I want to store them in a dictionary!



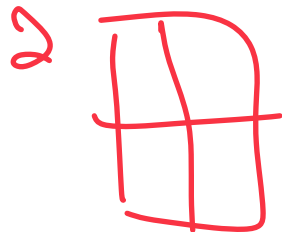
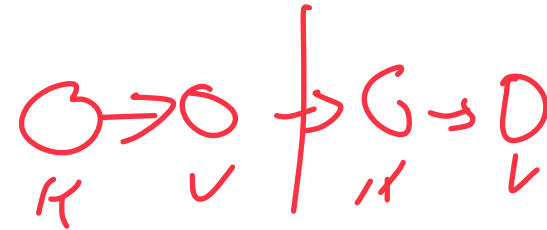
What data structures can I use here?

$O(\log n)$ → tree (AVL tree)

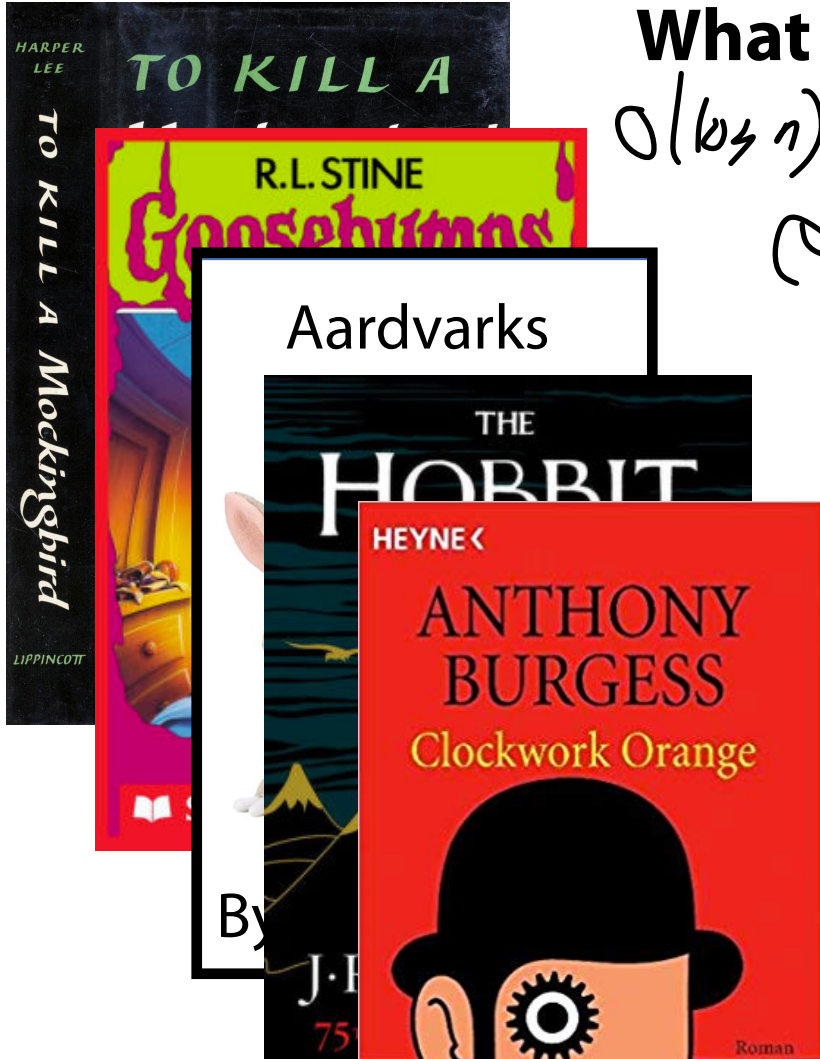
min heap
→ priority on min book

→ 2 ~~linked lists~~ Arrays
key →
value →

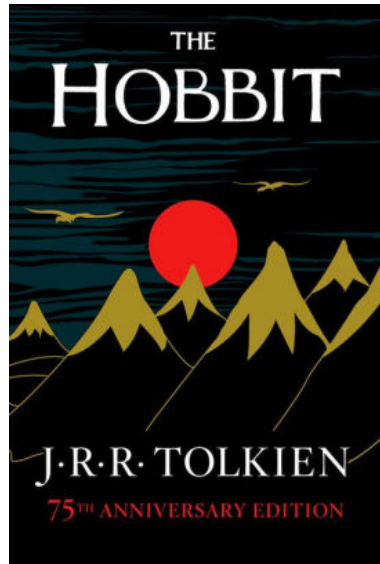
→ I linked list



→ similar?
→ similar trees?

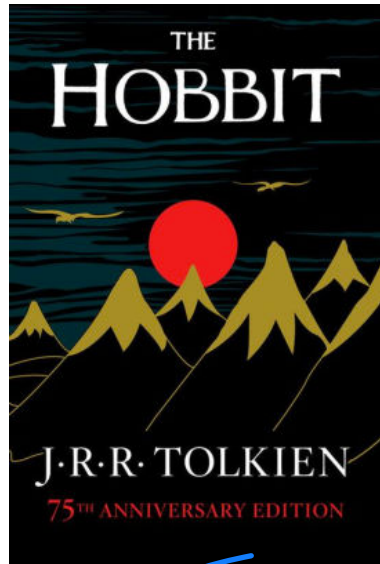


What if $O(\log n)$ isn't good enough?



What if $O(\log n)$ isn't good enough?

$O(1)$ →

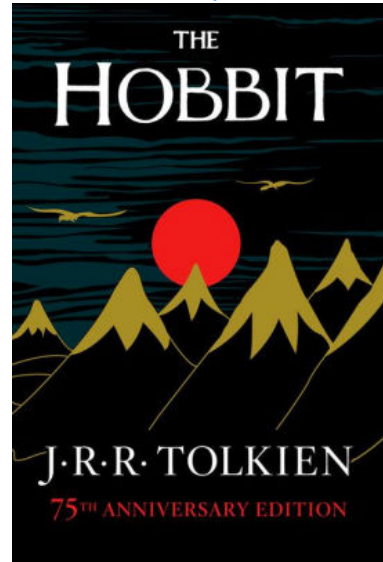


→ $O(1)$



A Hash Table based Dictionary

Key



Black box



O(1)



(Memory)

Address

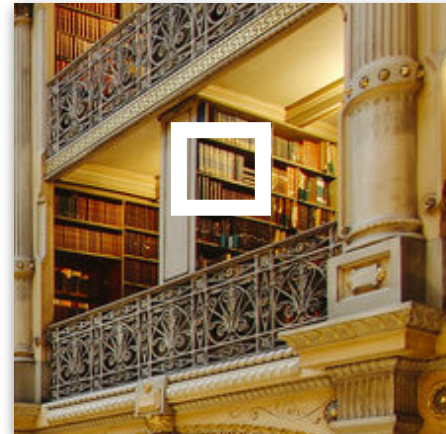
ISBN: 9780062265722

Call #: PR
6068.093
H35 1937

address

ISBN: 9780062265722

Call #: PR
6068.093
H35 1937



O(1)



Value

Chapter I

AN UNEXPECTED PARTY

In a hole in the ground there lived a hobbit. Not a nasty, dirty, wet hole, filled with the ends of worms and an oozy smell, nor yet a dry, bare, sandy hole with nothing in it to sit down on or to eat: it was a hobbit-hole, and that means comfort.

It had a perfectly round door like a porthole, painted green, with a shiny yellow brass knob in the exact middle. The door opened on to a tube-shaped hall like a tunnel: a very comfortable tunnel without smoke, with panelled walls, and floors tiled and carpeted, provided with polished chairs, and lots and lots of pegs for hats and coats—the hobbit was fond of visitors. The tunnel wound on and on, going fairly but not quite straight into the side of the hill—The Hill, as all the people for many miles round called it—and many little round doors opened out of it, first on one side and then on another. No going upstairs for the hobbit: bedrooms, bathrooms, cellars, pantries (lots of these), wardrobes (he had whole rooms devoted to clothes), kitchens, dining-rooms, all were on the same floor, and indeed on the same passage. The best rooms were all on the left-hand side (going in), for these were the only ones to have windows, deep-set round windows looking over his garden, and meadows beyond, sloping down to the river.

This hobbit was a very well-to-do hobbit, and his name

1

Randomized Data Structures

Sometimes a data structure can be **too ordered / too structured**

↳ AVL tree $(\log n)$ ←

Randomized data structures rely on **expected** performance

Expected performance good!

Randomized data structures 'cheat' tradeoffs!

Worst Case is Bad

A Hash Table based Dictionary



User Code (is a map):

```
1 Dictionary<KeyType, ValueType> d;  
2 d[k] = v;
```



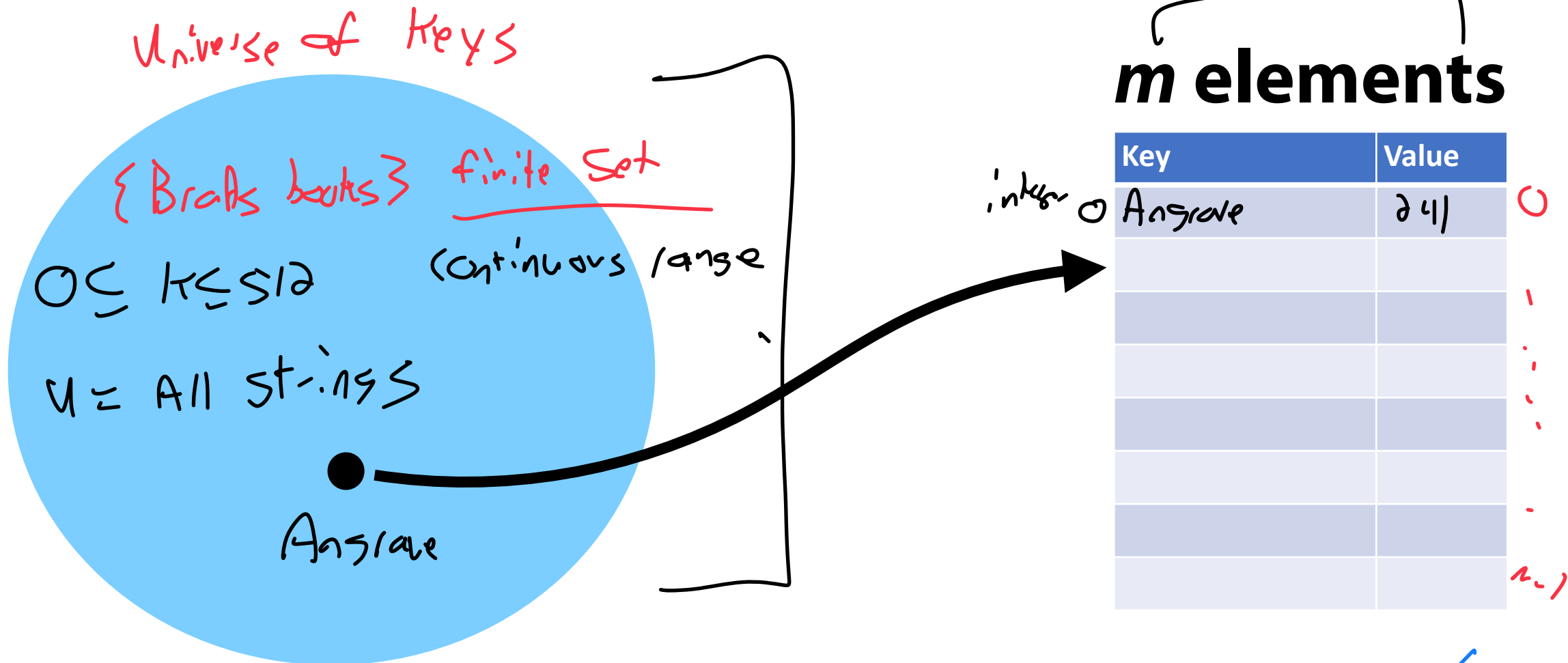
(KeyType) \rightarrow int
uint64_t
32-bit

A Hash Table consists of three things:

1. A hash function Key as input \rightarrow returns an address
2. A data storage structure (An array is good!) given index \rightarrow O(1) lookup
3. ??? A way of handling chaos

Hash Function

Maps a **keyspace**, a (mathematical) description of the keys for a set of data, to a set of integers.

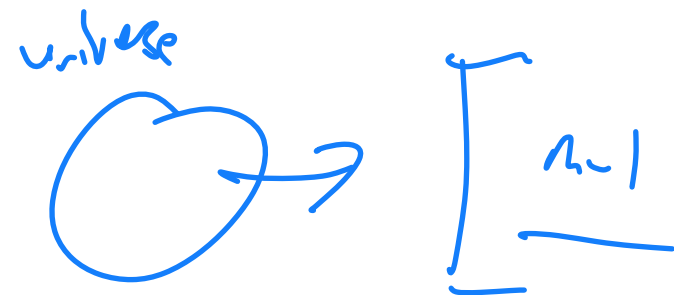


Hash Function

A hash function **must** be:

- **Deterministic:** Given the same key, return the same value
↳ If math major obvious!
- **Efficient:** Our goal is $O(1)$. The hash function is run every time [Be aware of our key size but indep of key size]
- **Defined for a certain size table:**

↳ % m (size of table)



Hash Function

Key

Value

(Angrave, CS 241)
(Beckman, CS 421)
(Challon, CS 125)
(Davis, CS 101)
(Evans, CS 225)
(Fagen-Ulmschneider, CS 107)
(Gunter, CS 422)
(Herman, CS 233)

Hash function

$K[K] = 'A'$

Key	Value
Angrave	241
Beckman	421

~ ~ ~

Hash Function

Good
↳ A perfect bijection!

Are all letters equally likely?
↳ Bad hash function!

(Angrave, CS 241)

Alquini

(Beckman, CS 421)

(Challon, CS 125)

(Davis, CS 101)

Daves

(Evans, CS 225)

(Fagen-Ulmschneider, CS 107)

(Gunter, CS 422)

(Herman, CS 233)

Only works for 26 letters

Hash function

(key[0] - 'A')

Key	Value	
Angrave	241	0
Beckman	421	1
Challon	125	2
Davis	101	3
Evans	225	4
Fagen-U	107	5
Gunter	422	6
Herman	233	7

Solomon

223 8

Bad! Cant have two people w/ same first letter.



General Hash Function $H: \text{Keytype} \rightarrow \text{int}$

An $O(1)$ deterministic operation that maps all keys in a universe U to a defined range of integers $[0, \dots, m - 1]$

- A hash: $\text{Key} \rightarrow \text{int}$

- A compression:

$\text{Mod. } H(x) \% m$

Lab-Hash
 $\hookrightarrow \text{int} \rightarrow \text{int}$
 $\hookrightarrow \text{return } x$

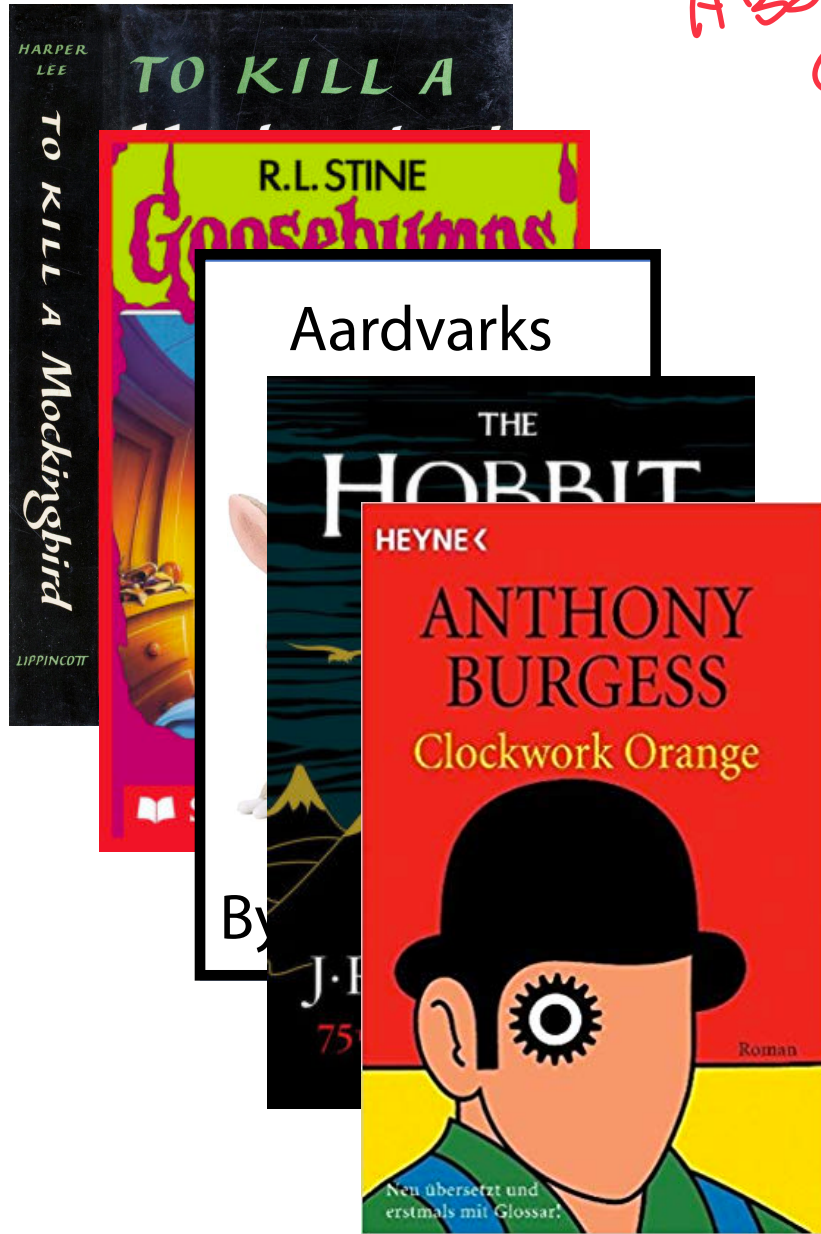
Choosing a good hash function is tricky...

- Don't create your own (yet*)

Hash Function

Are those hash functions?

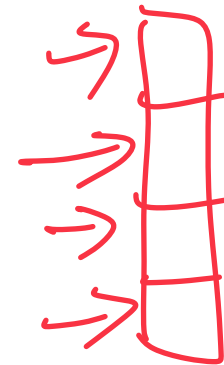
Also 1 book per author



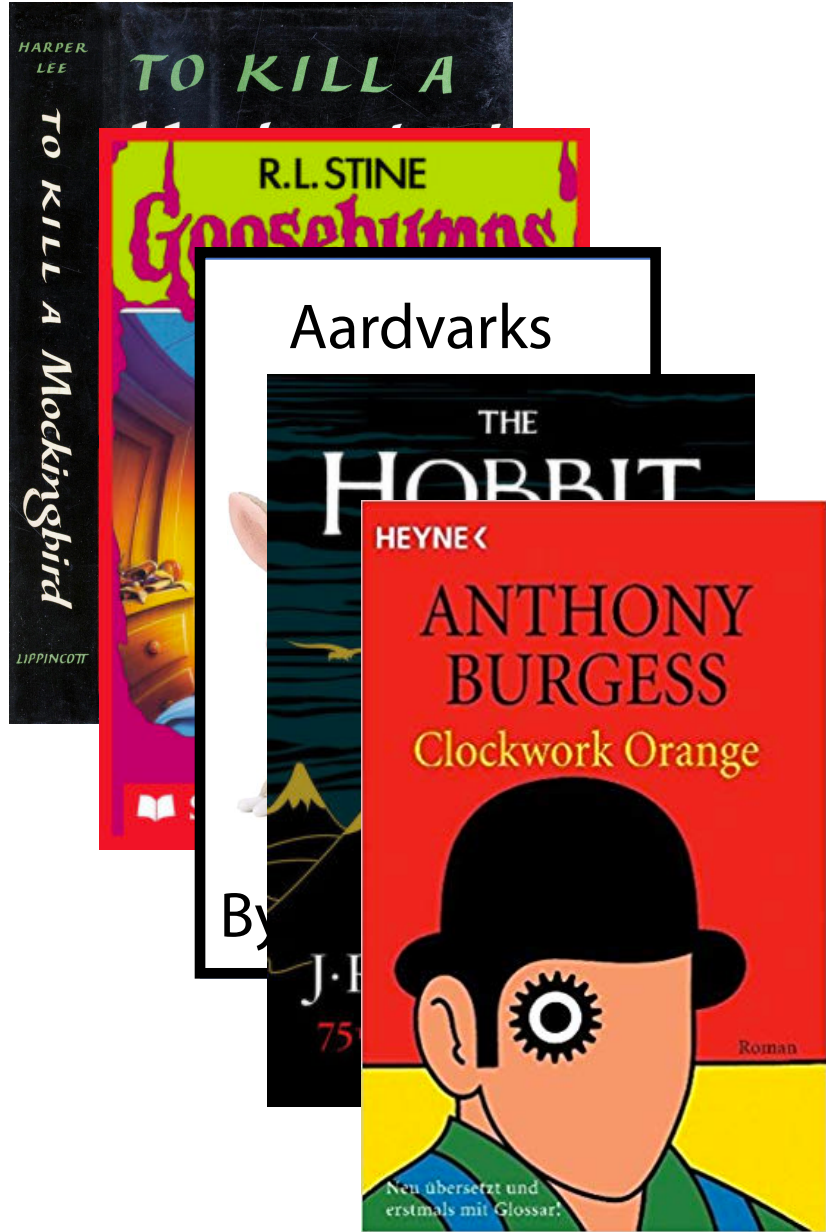
$h(k) = (k.firstName[0] + k.lastName[0]) \% m$
Is $O(1)$ ✓
deterministic ✓
First & last names are not uniformly distributed

$h(k) = (\text{rand}() * k.numPages) \% m$
Not a hash function!
Can't have randomness inside function!
Started my program by generating one #

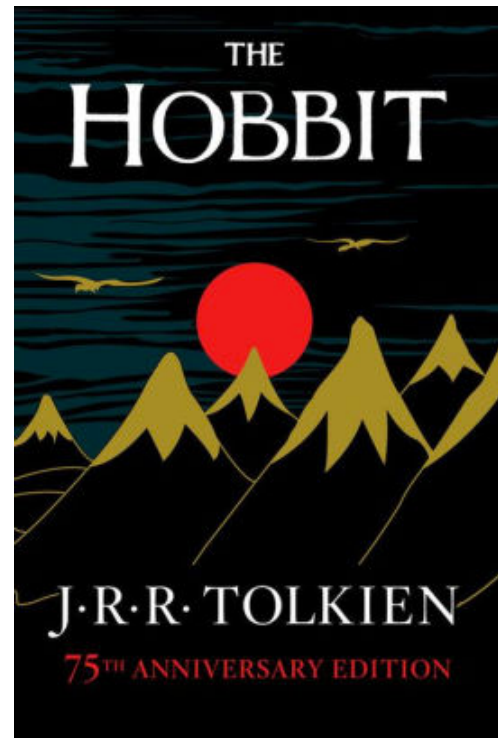
$h(k) = (\text{Order I insert} [\text{Order seen}]) \% m$
Not deterministic!



Hash Function



A Hash Table based Dictionary



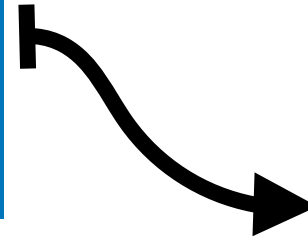
key



Author Name
Hash Function

$O(1)$

'J' + 'T' = 30



27
28
29
30
31
...

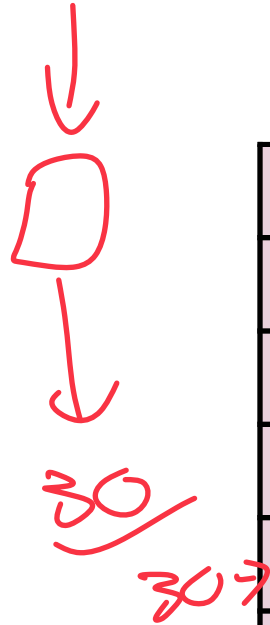
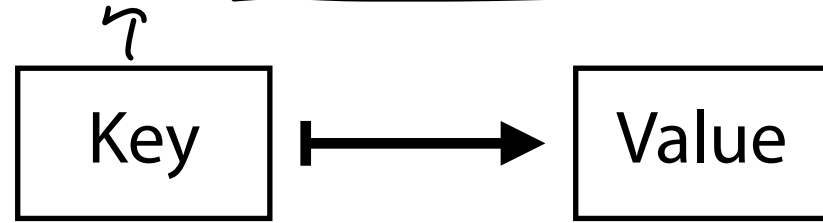
int

...
∅
∅
∅
The Hobbit
∅
...

A Hash Table based Dictionary

The hobbit

informs my address

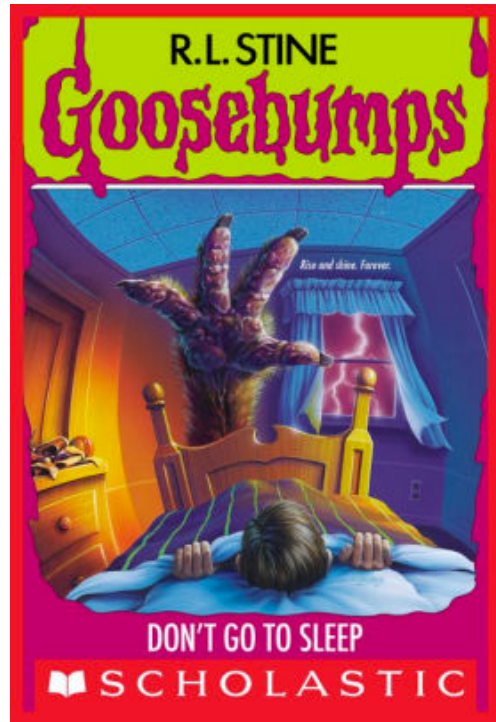


∅
∅
∅
∅
The Hobbit
∅
∅

A large rectangular frame containing two images. On the left is the cover of the book "The Hobbit" by J.R.R. Tolkien, 75th Anniversary Edition, featuring a red sun and green mountains. On the right is a page of text from the book, titled "Chapter 1 AN UNEXPECTED PARTY". The text describes a hobbit's home and the beginning of the story.

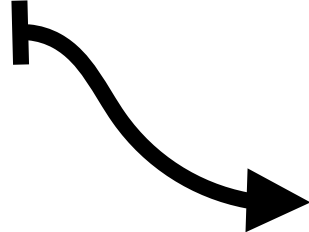
B/c collisions

A Hash Table based Dictionary



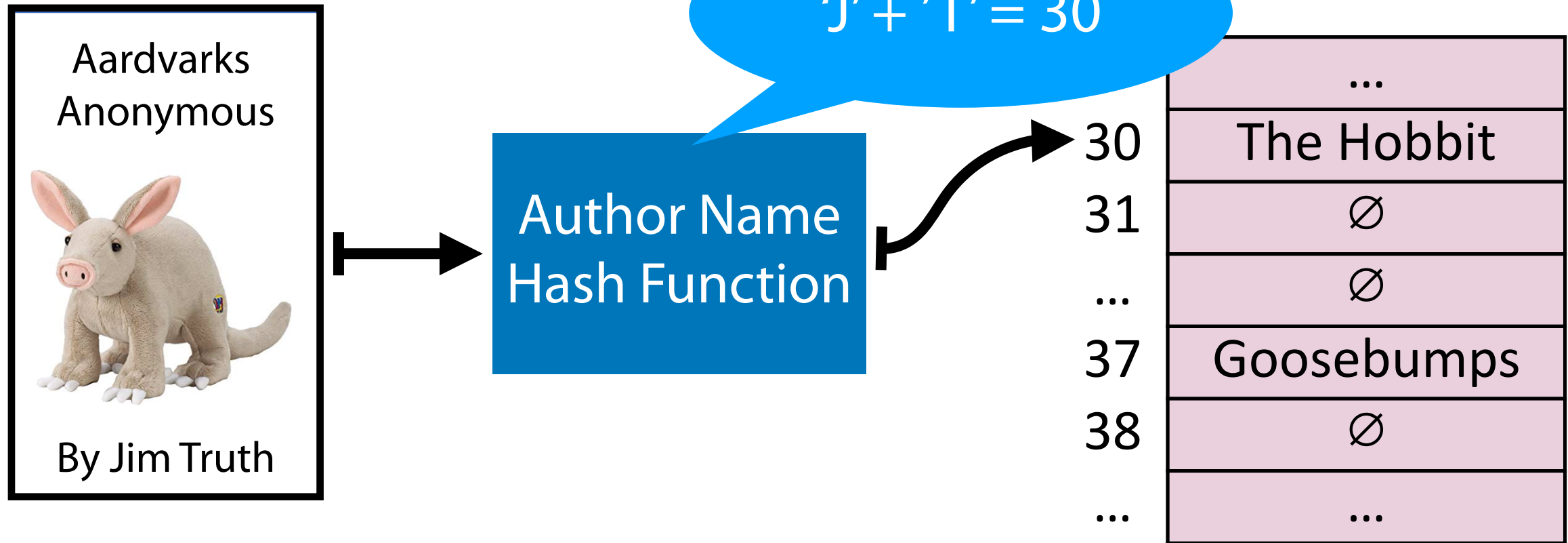
Author Name
Hash Function

'R' + 'S' = 37



...	...
30	The Hobbit
31	∅
...	∅
37	Goosebumps
38	∅
...	...

A Hash Table based Dictionary

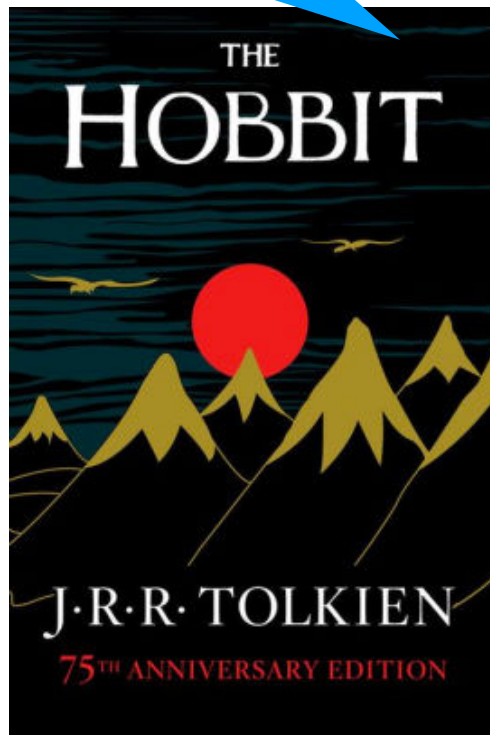


Hash Collision

A *hash collision* occurs when multiple unique keys hash to the same value

hash

J.R.R Tolkien = 30!



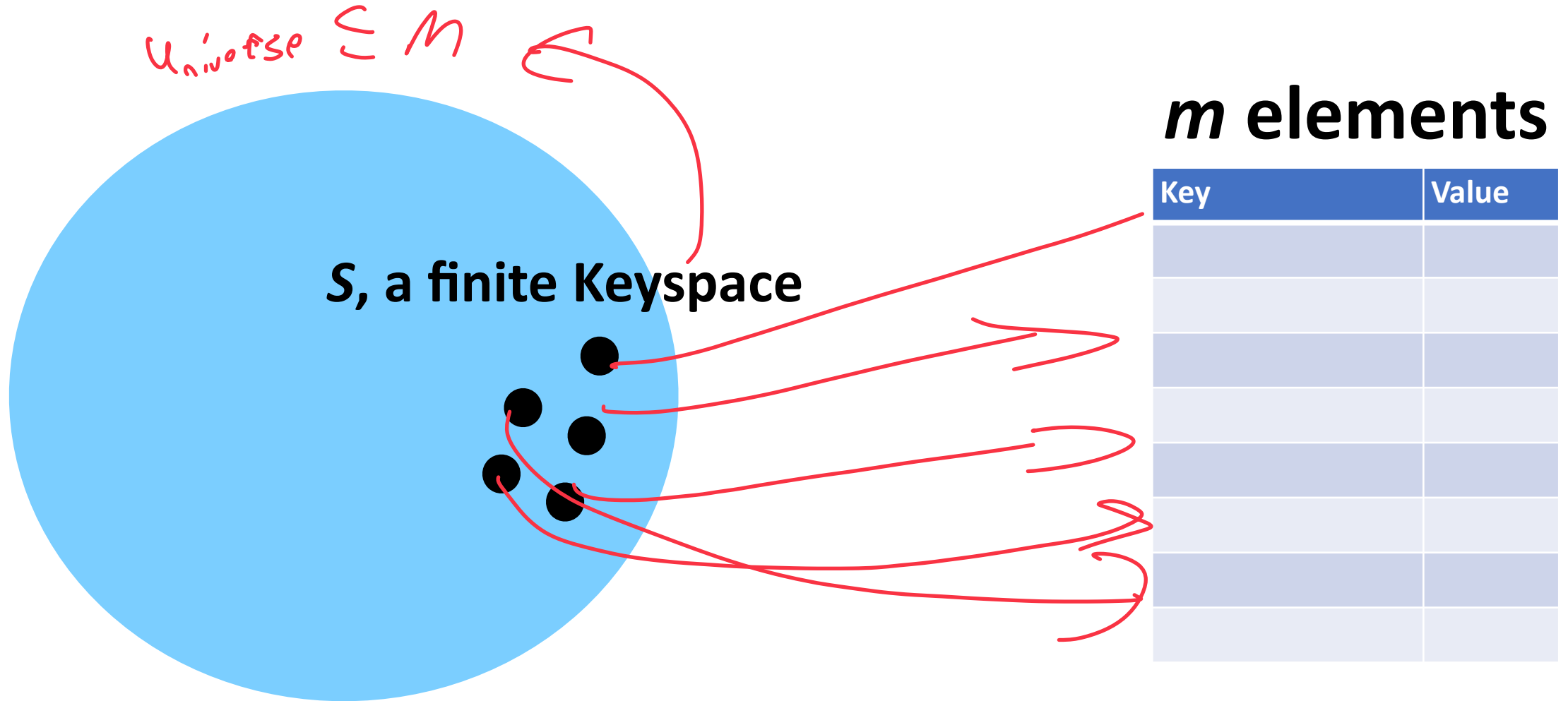
Jim Truth = 30!



...	...
30	???
31	∅
...	∅
37	Goosebumps
38	∅
...	...

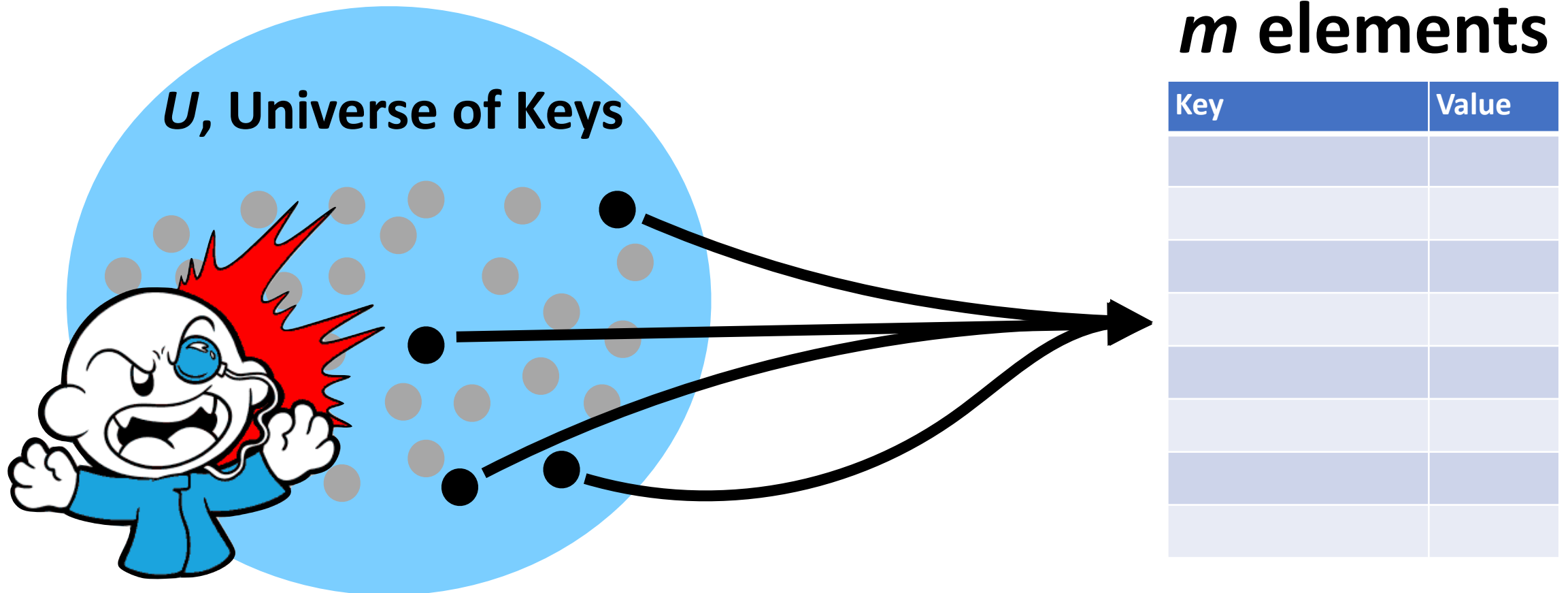
Perfect Hashing

If $m \geq S$, we can write a *perfect* hash with no collisions



General Purpose Hashing

By fixing h , we open ourselves up to adversarial attacks.





A Hash Table based Dictionary

User Code (is a map):

```
1 Dictionary<KeyType, ValueType> d;  
2 d[k] = v;
```

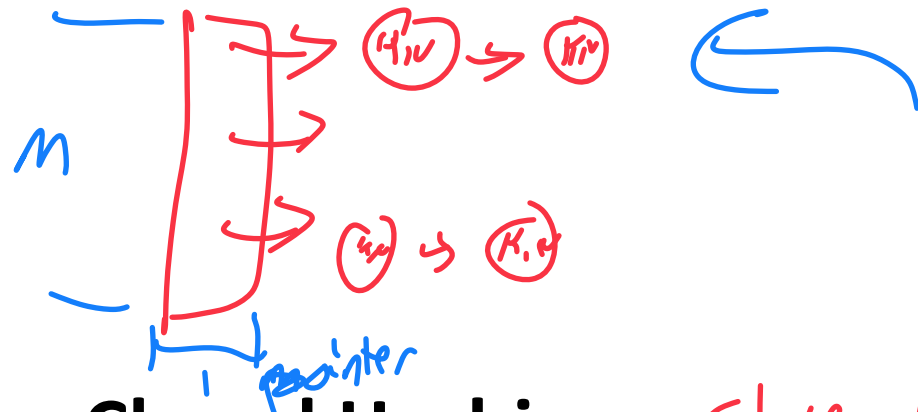
A **Hash Table** consists of three things:

1. A hash function
2. A data storage structure
3. A method of addressing *hash collisions*

Open vs Closed Hashing

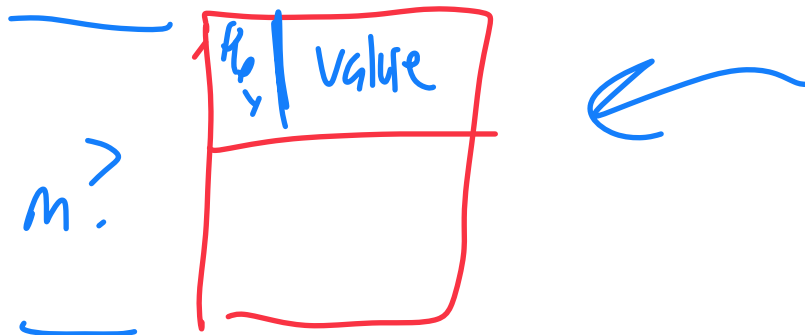
Addressing hash collisions depends on your storage structure.

- **Open Hashing:** Store (k, v) pairs externally (outside my table)



Easy!

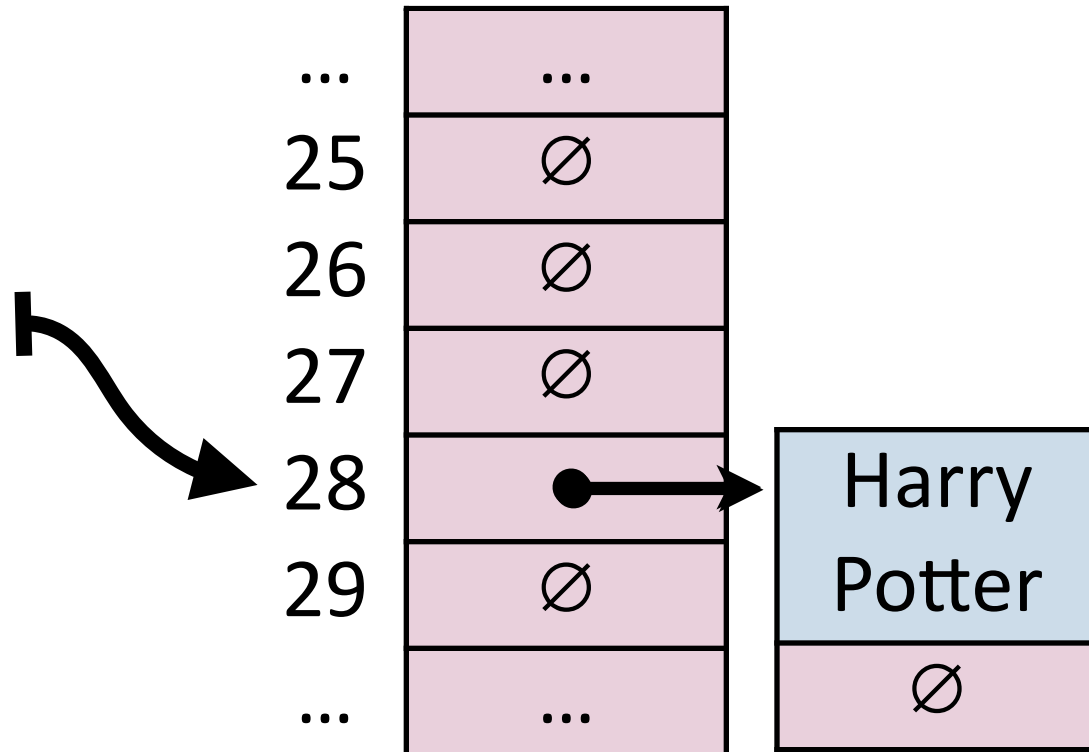
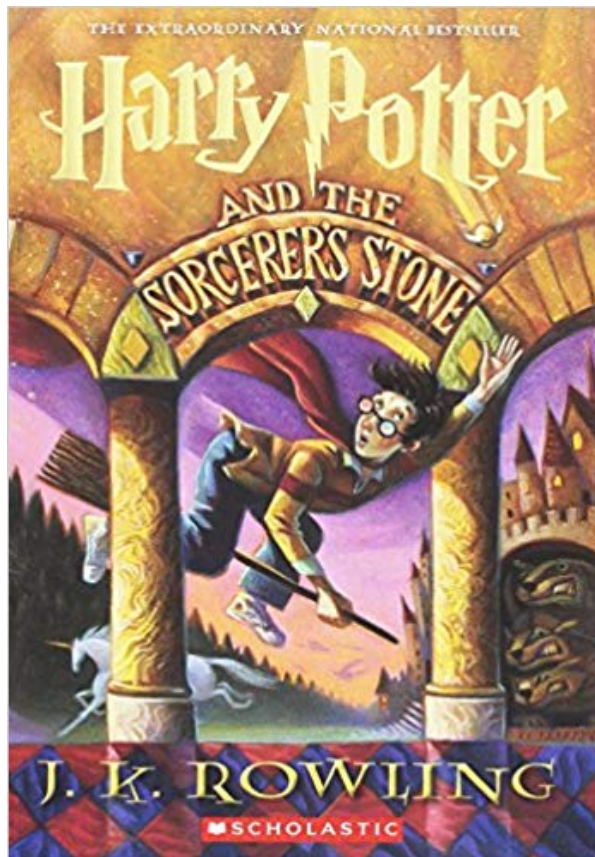
- **Closed Hashing:** Store (k, v) internally (in my table)



Hard!

Open Hashing

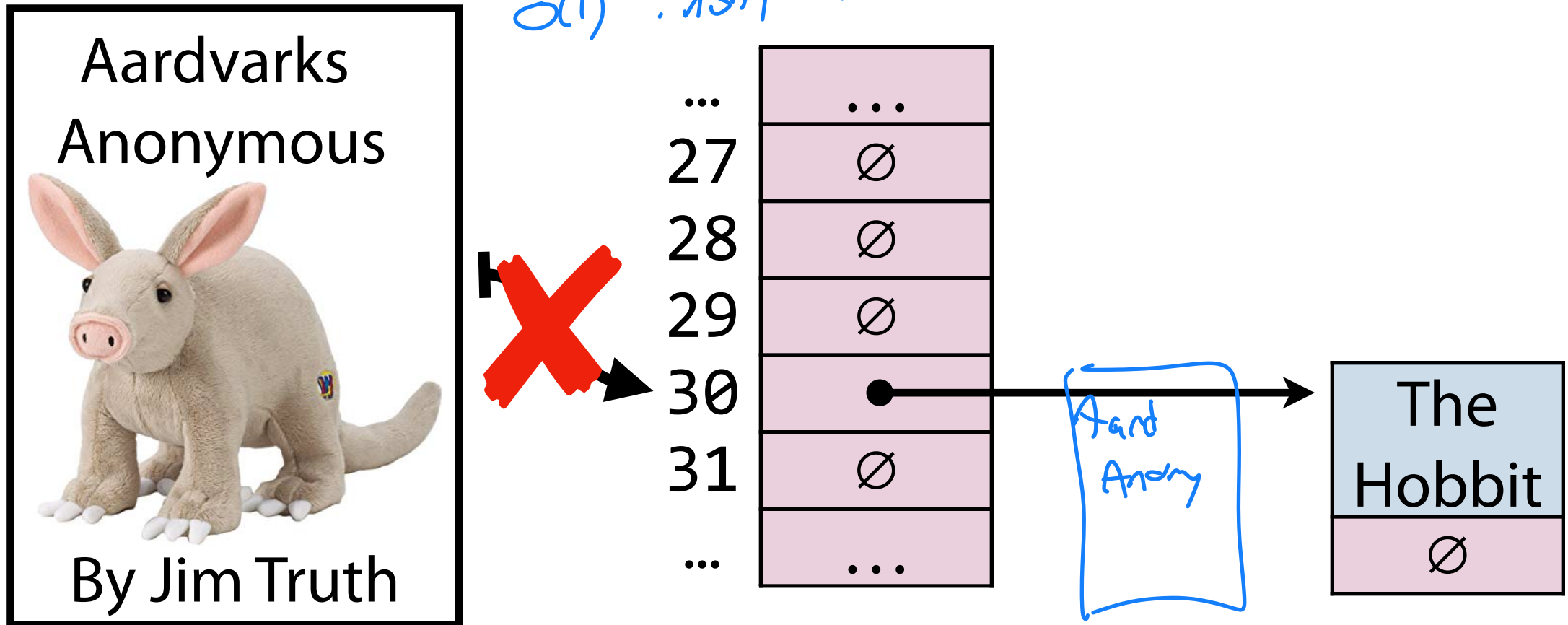
In an *open hashing* scheme, key-value pairs are stored externally (for example as a linked list).



Hash Collisions (Open Hashing)

A **hash collision** in an open hashing scheme can be resolved by inserting at front. This is called **separate chaining**.

$O(1)$ insert at front



Insertion (Separate Chaining)

`_insert("Bob")`

`_insert("Anna")`

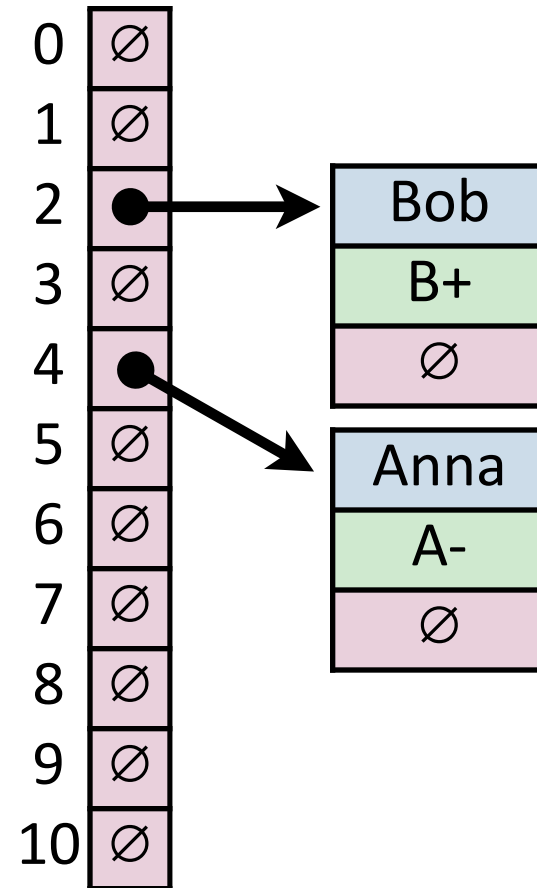
Key	Value	Hash
Bob	B+	2
Anna	A-	4
Alice	A+	4
Betty	B	2
Brett	A-	2
Greg	A	0
Sue	B	7
Ali	B+	4
Laura	A	7
Lily	B+	7

0	∅
1	∅
2	∅
3	∅
4	∅
5	∅
6	∅
7	∅
8	∅
9	∅
10	∅

Insertion (Separate Chaining)

`_insert("Alice")`

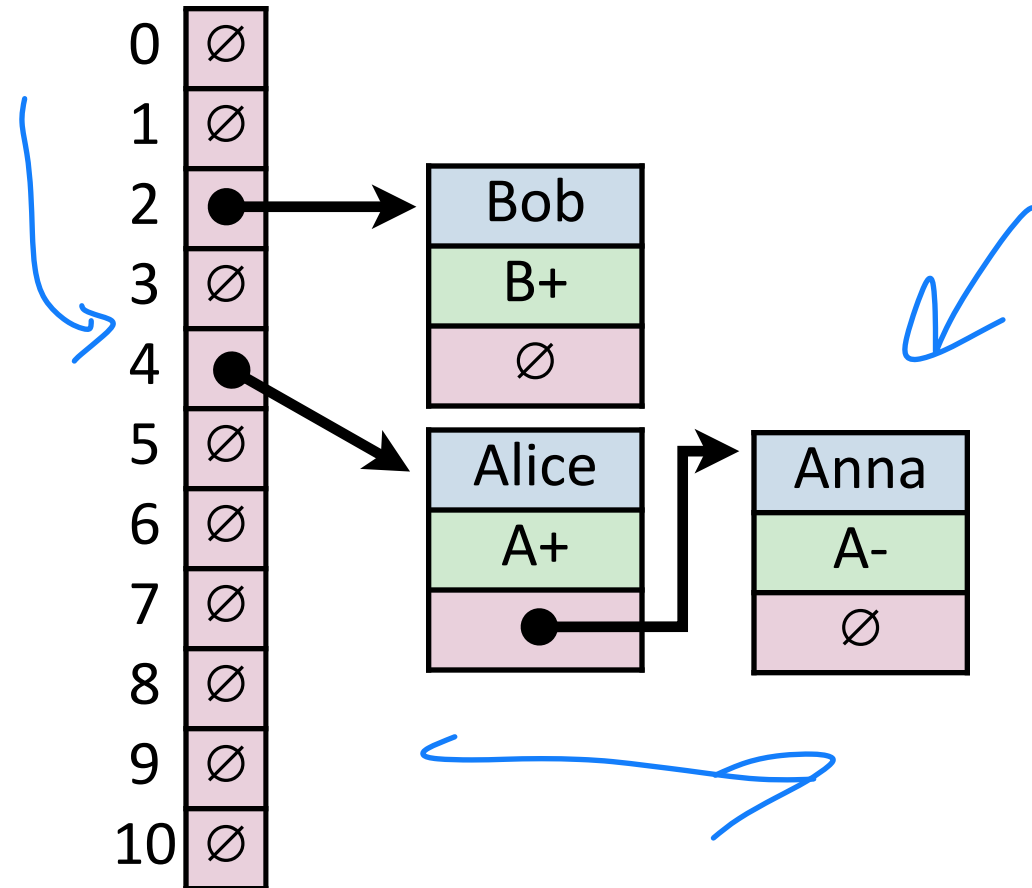
Key	Value	Hash
Bob	B+	2
Anna	A-	4
Alice	A+	4
Betty	B	2
Brett	A-	2
Greg	A	0
Sue	B	7
Ali	B+	4
Laura	A	7
Lily	B+	7



Insertion (Separate Chaining)

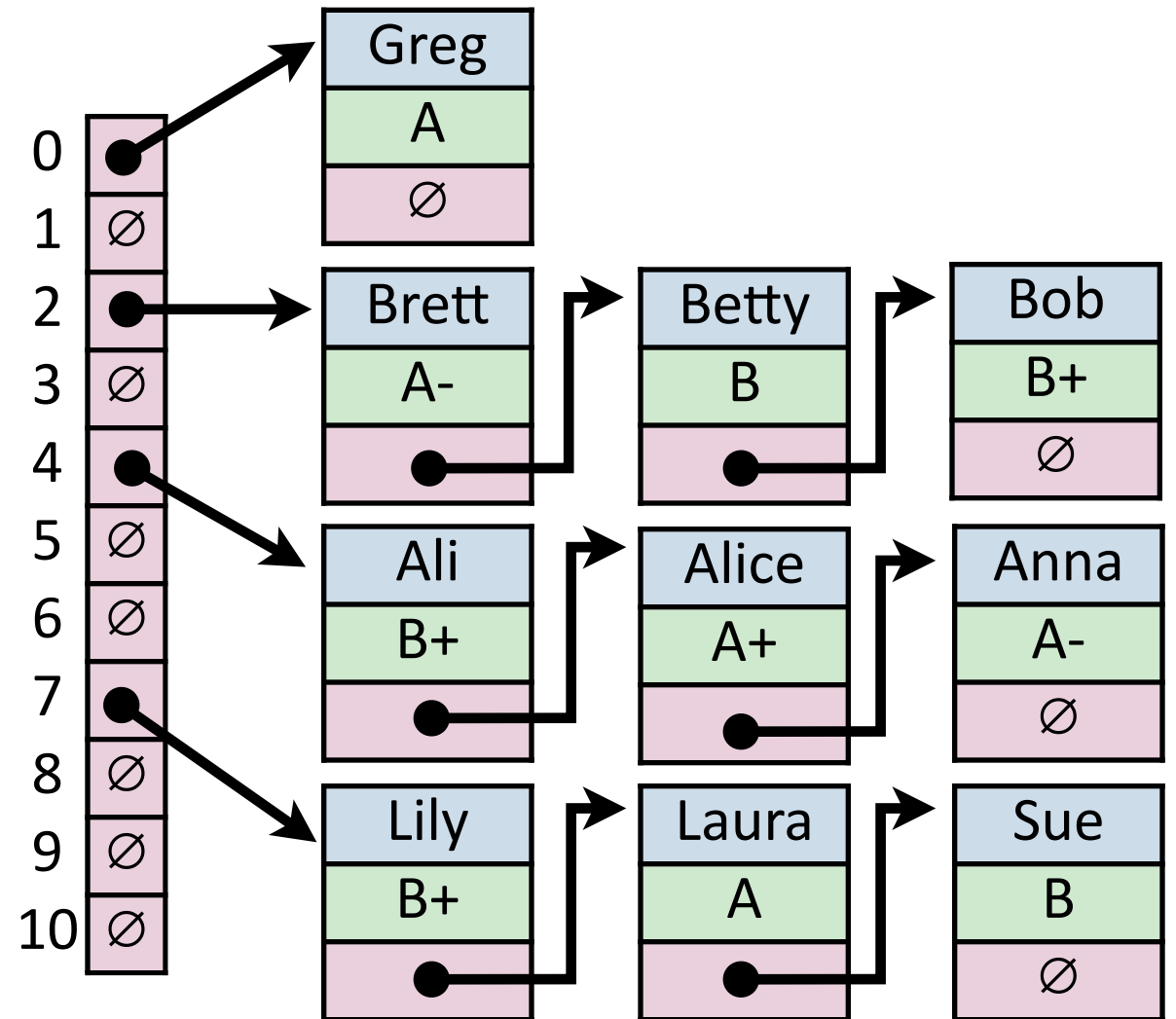
find(Anna)

Key	Value	Hash
Bob	B+	2
Anna	A-	4
Alice	A+	4
Betty	B	2
Brett	A-	2
Greg	A	0
Sue	B	7
Ali	B+	4
Laura	A	7
Lily	B+	7



Insertion (Separate Chaining)

Key	Value	Hash
Bob	B+	2
Anna	A-	4
Alice	A+	4
Betty	B	2
Brett	A-	2
Greg	A	0
Sue	B	7
Ali	B+	4
Laura	A	7
Lily	B+	7



Find (Separate Chaining)

`_find("Sue")`

Key	Hash
Sue	7

$O(1)$

Black box hash

$O(1)$

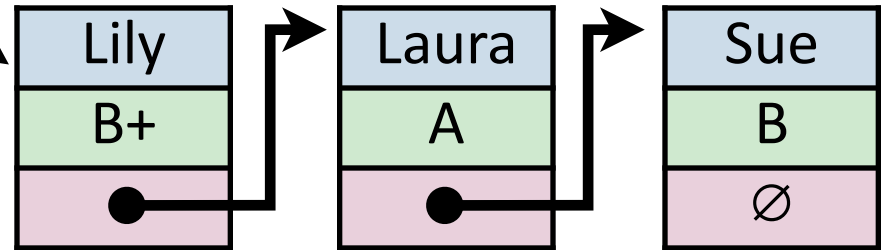
$H[7]$

H

0	∅
1	∅
2	∅
3	∅
4	∅
5	∅
6	∅
7	●
8	∅
9	∅
10	∅

$O(n)$

⋮



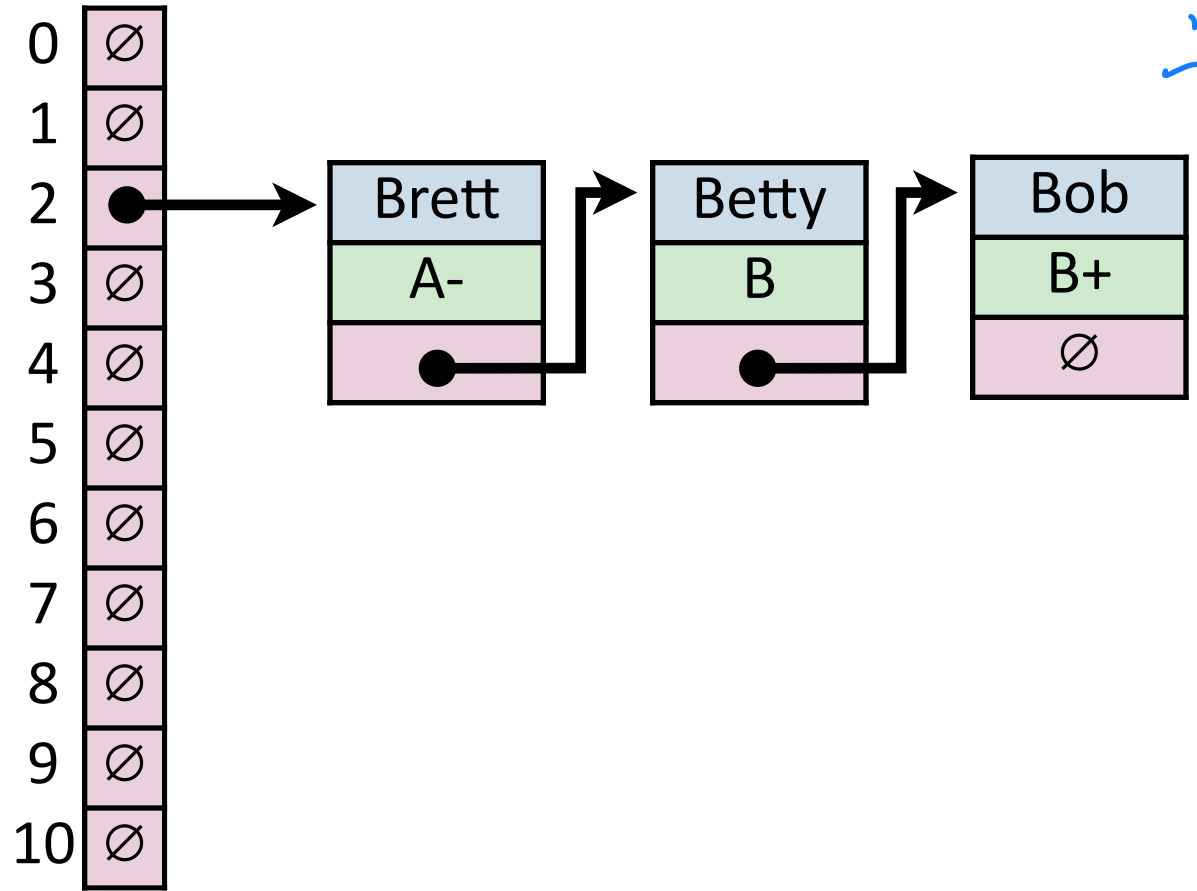
Remove (Separate Chaining)

`_remove("Betty")`

Key	Hash
Betty	2

$O(1)$

$O(1)$



LL remove $O(n)$

Hash Table (Separate Chaining)



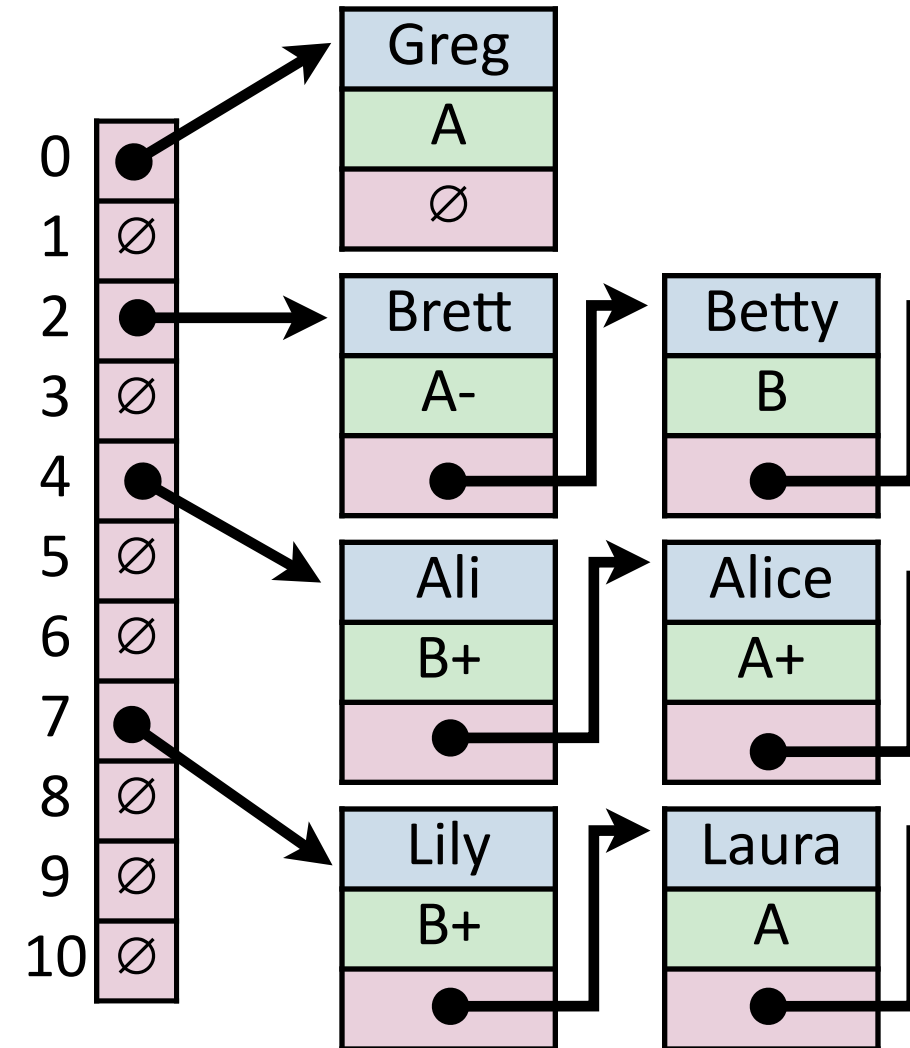
For hash table of size m and n elements:

Find runs in: $O(n)$

Insert runs in: $O(1)$ ☺

Remove runs in: $O(n)$

↳ B/c of collisions!



Hash Table

T_n expectation
↑

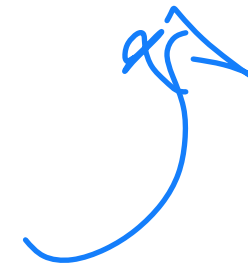
Worst-Case behavior is bad — but what about randomness?

1) Fix h , our hash, and assume it is good for ***all keys***:

Simple uniform hashing assumption

2) Create a ***universal hash function family***:

↳ secret (SHA in real life)



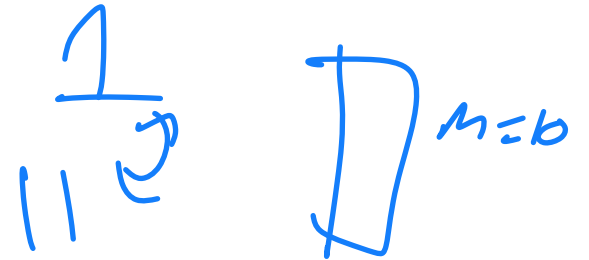
Simple Uniform Hashing Assumption

Given table of size m , a simple uniform hash, h , implies

$$\forall k_1, k_2 \in U \text{ where } k_1 \neq k_2, \Pr(h[k_1] = h[k_2]) = \frac{1}{m}$$

Uniform: My hash is uniform

Independent: All keys hash independently



Separate Chaining Under SUHA

Given table of size m and n inserted objects

Claim: Under SUHA, expected length of chain is $\frac{n}{m}$

Separate Chaining Under SUHA

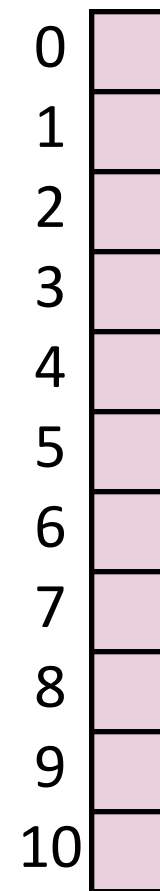


Under SUHA, a hash table of size m and n elements:

Find runs in: _____.

Insert runs in: _____.

Remove runs in: _____.



Separate Chaining Under SUHA



Pros:

Cons:

Next time: Closed Hashing

Closed Hashing: store k, v pairs in the hash table

$S = \{ 1, 8, 15 \}$

$$h(k) = k \% 7$$

