

# Data Structures and Algorithms

## Hashing

CS 225

Brad Solomon

April 17, 2026



UNIVERSITY OF  
**ILLINOIS**  
URBANA - CHAMPAIGN

Department of Computer Science

$\wedge$   $\wedge$   
—  
The best  
data structure!

# Randomization in Algorithms

weakest!

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

↳ BST      A random input dataset is probably  $O(\log n)$  height bounded  
data is random ↑

2. Use randomness inside algorithm to estimate expected running time

↳ Randomized Quicksort  
↳ selection of pivot is random      for any dataset  
100% accuracy ↑

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

↳ primality test      ↳ some loss of accuracy here  
↳ All prime #s = 1  
↳ Some not prime #s = 1      # of this is very small

# Learning Objectives

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

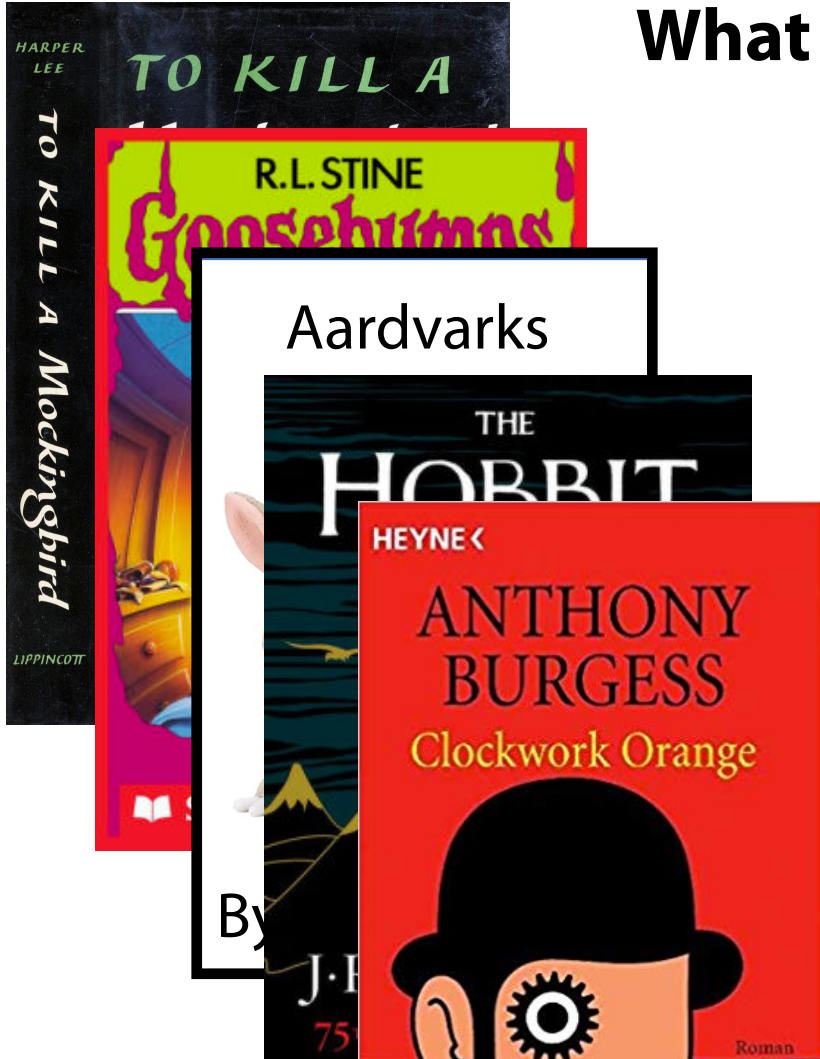
we want like  
hash tables yet

# Data Structure Review

Key: Book title  
Value: Book contents

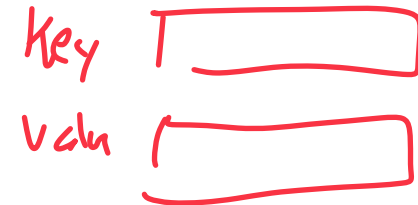
I have a collection of books and I want to store them in a dictionary!

## What data structures can I use here?



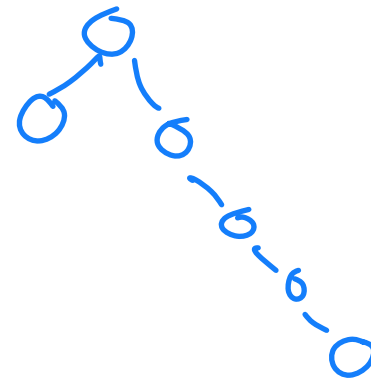
Parallel Arrays

↳ Insert  $O(1)^*$   
↳ Find  $O(n)$



BST Each node K,V pair

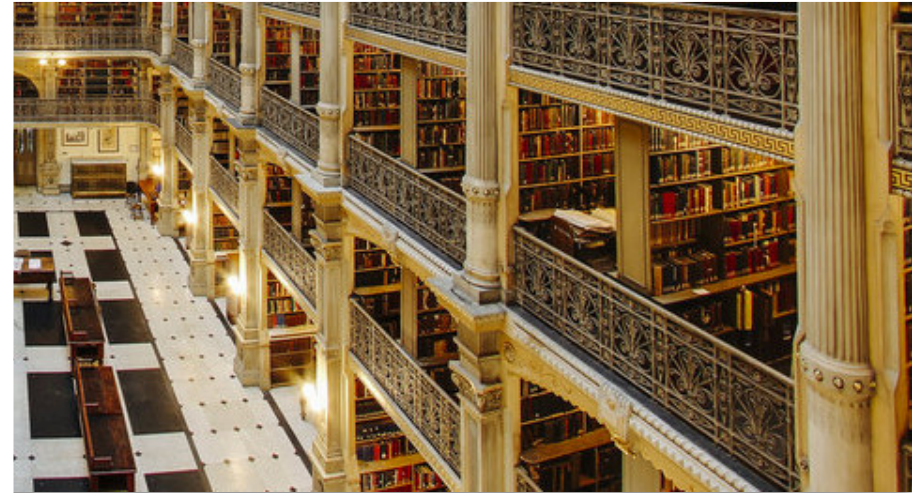
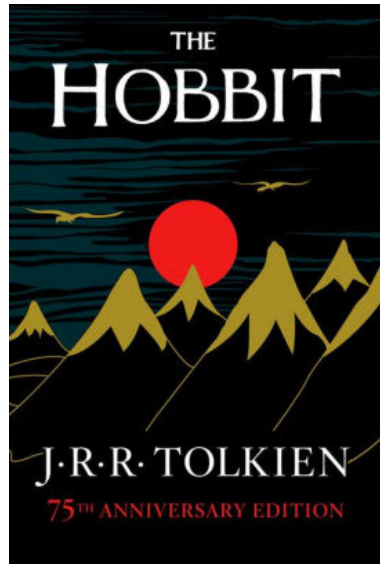
↳ Insert  $O(n) ??$   
↳ Find  $O(n)$



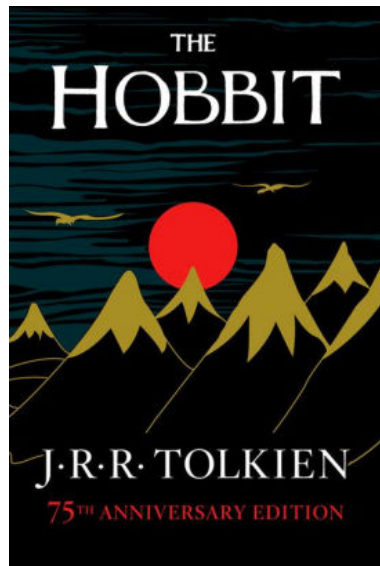
Balanced BST ← \* Balanced  $\log - \log$

↳ Insert  $O(\log n)$   
↳ Find  $O(\log n)$

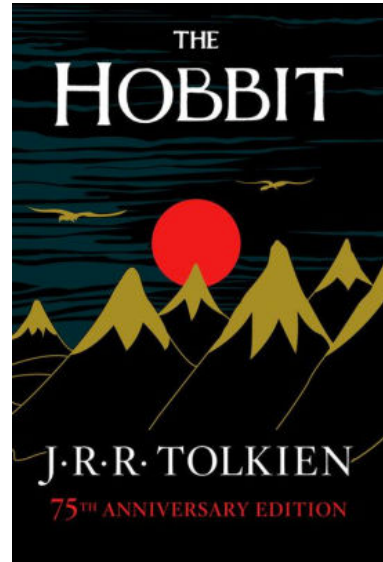
# What if $O(\log n)$ isn't good enough?



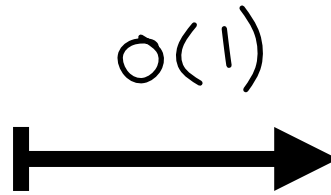
# What if $O(\log n)$ isn't good enough?



# A Hash Table based Dictionary



Key



Black Box



A number  
Address

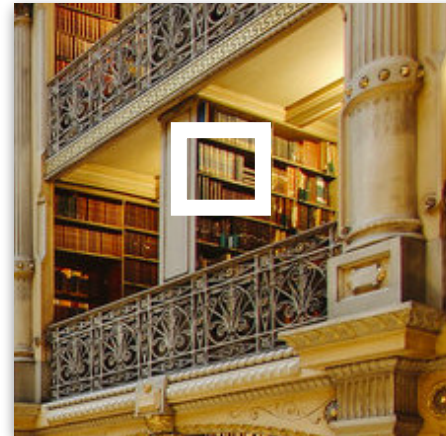
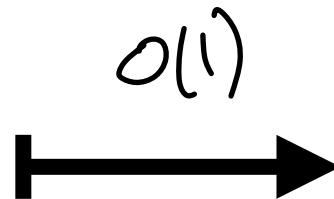
ISBN: 9780062265722

Call #: PR  
6068.093  
H35 1937

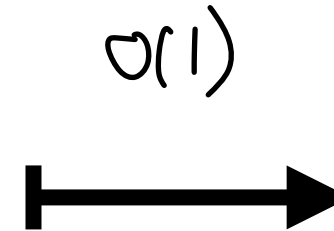
ISBN: 9780062265722

Call #: PR  
6068.093  
H35 1937

Address



Find Value



f return 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**

B BST can find in  $O(\log n)$

↳ As tradeoff it costs  $O(\log n)$  to insert

Randomized data structures rely on **expected** performance

Randomized data structures 'cheat' tradeoffs!

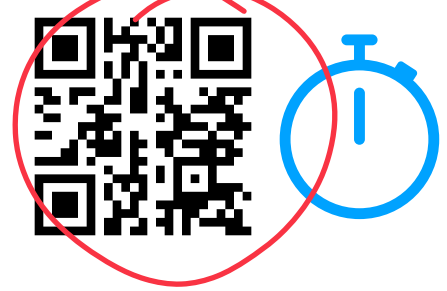
↳ Big  $O$  / worst case may be bad  
but  $\checkmark$  randomness

I can proc in  
expectation  $O(1)^*$

Star #1 of hash table!



# 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  $h(k): K \rightarrow \text{int}$  (address / hash value)
2. A data storage structure — An array <sup>70%</sup>
3.  $\gg$  This is some way to handle "chaos"

# Hash Function

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

A universe of keys

$U = \{\text{Brad's Books}\}$  (A finite set)

$U = \{0 \leq i \leq 255\}$  (continuous range)

$U = \text{All strings possible}$  (infinite range)

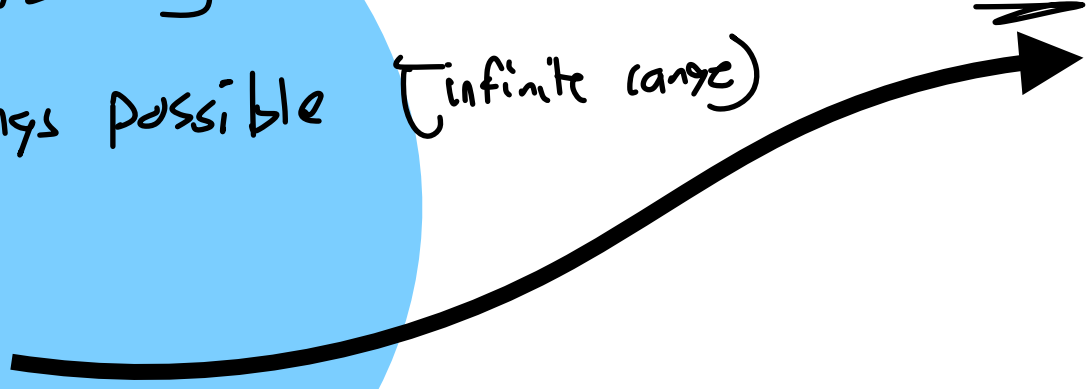
an address

**$m$  elements**

Key	Value

0

$m-1$



# Hash Function

A hash function **must** be:

\* • **Deterministic:**  $\forall k_1, k_2$  if  $(k_1 == k_2)$  then  $h(k_1) = h(k_2)$   
↳ I.E. you give the same key twice, get same address/hash value out!

• **Efficient:**  $O(1)$

• **Defined for a certain size table:** Size  $n$  table  $\rightarrow [0, n-1]$   
↑  
I choose  $n$   
entries



# Hash Function

Problem #1: collisions

1) What is a collision?

Both Angrave  
Alanini want to be @ 0!

(Angrave, CS 241) Alanini: 411

(Beckman, CS 421)

(Challon, CS 125)

(Davis, CS 101)

(Evans, CS 225)

(Fagen-Ulmschneider, CS 107)

(Gunter, CS 422)

(Herman, CS 233)

Hash function

(key[0] - 'A')

Key	Value
Angrave	241
Beckman	421
Challon	125
Davis	101
Evans	225
Fagen-U	107
Gunter	422
Herman	233

Soloman 225

Doesn't fit in array!

# General Hash Function

efficient

to lookup

general purpose



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:  $h(k) \rightarrow \text{int}$

$$h(\text{Brad}) = 9999 \\ = 4$$

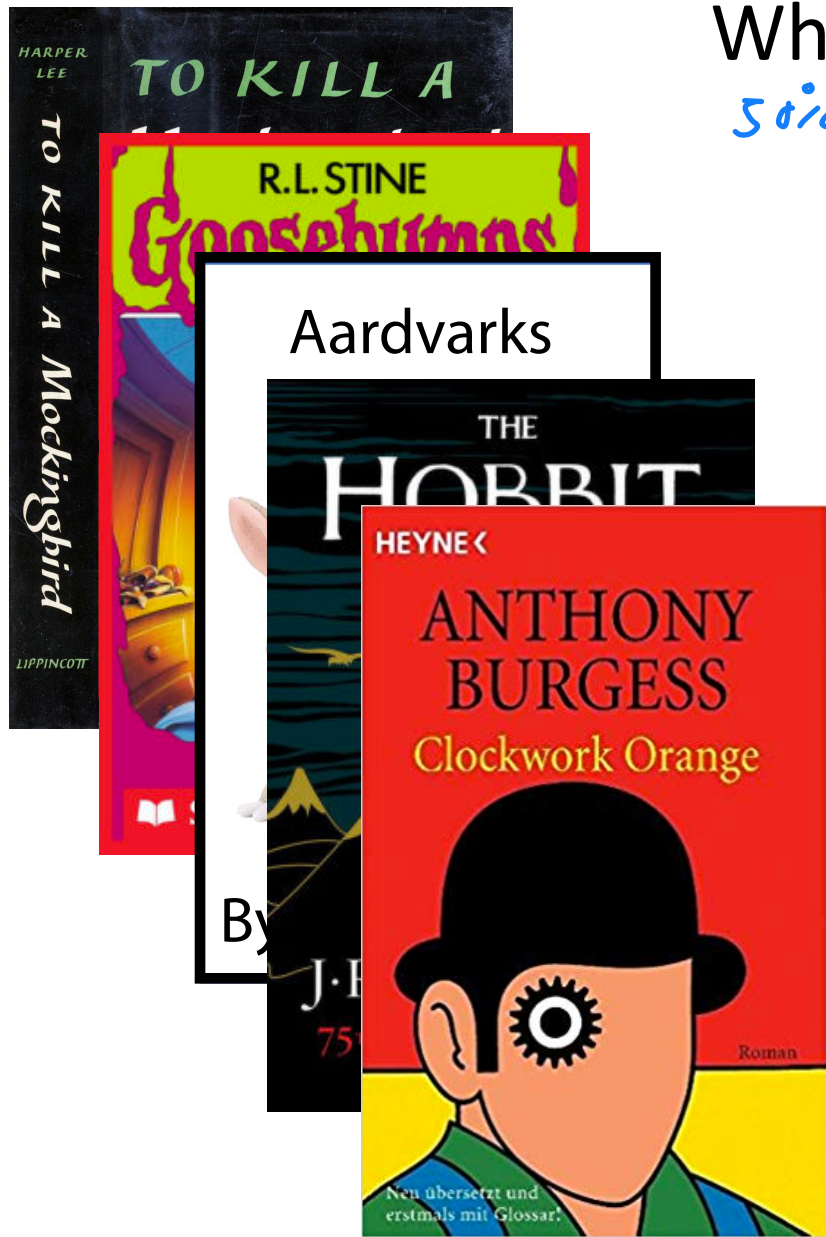
- A compression:  $h(k) \% m$

$\uparrow$   
 $m$   
 $\downarrow$

Choosing a good hash function is tricky...\*

- Don't create your own (yet\*)

# Hash Function



Which of the following are valid hashes?

50%

✓  $h_1(k) = (k.firstName[0] + k.lastName[0]) \% m$   
First letter first & last name  
↳ Deterministic? Yes!    ↳ O(1)? Yes!    ↳ Any name? Yes!

A bad function → many possible collisions

generate random # \* num pages

35% ✗  $h_2(k) = (rand() * k.numPages) \% m$

↳ Deterministic? No

If we generate only one # then it is fine!  
↳ c=2

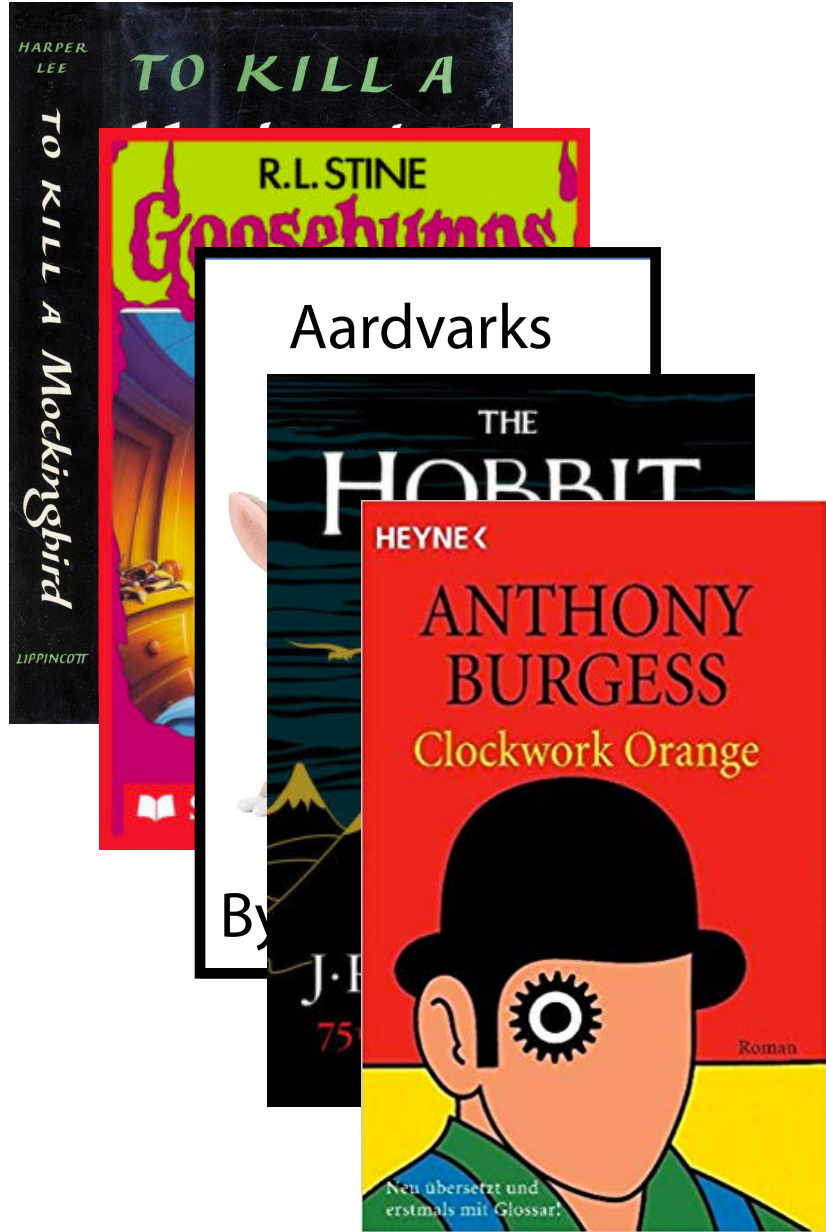
Index position is order seen items

✗  $h_3(k) = (return Count++;) \% m$

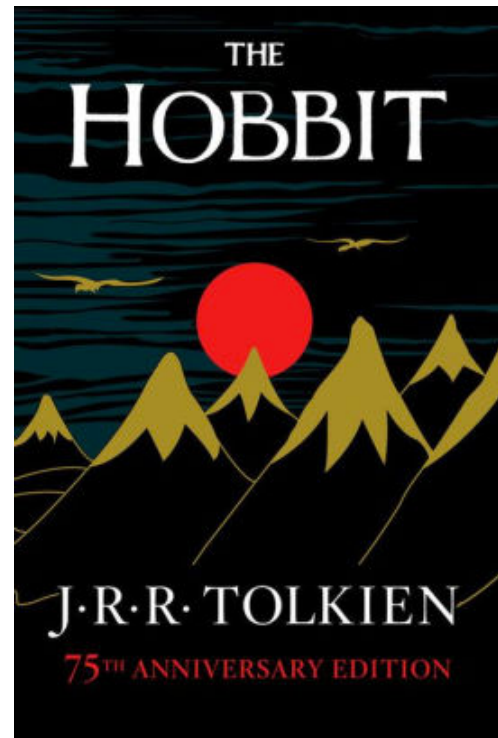
↳ Some books diff count value each time

Neu übersetzt und erstmals mit Glossar!

# Hash Function

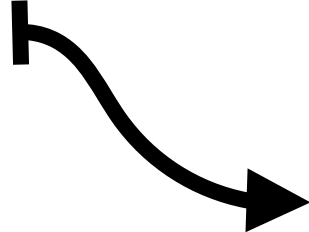


# A Hash Table based Dictionary



Author Name  
Hash Function

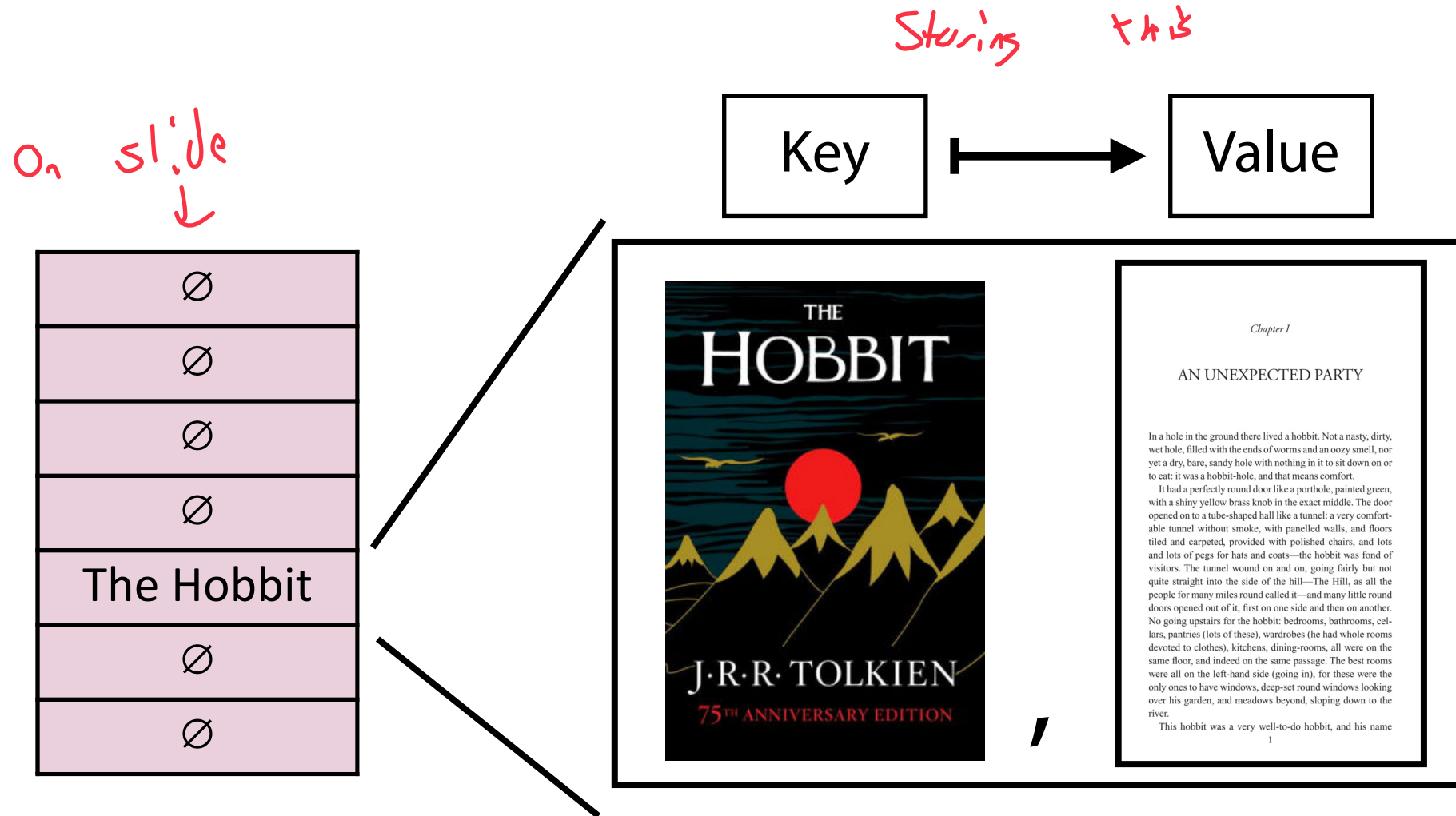
'J' + 'T' = 30



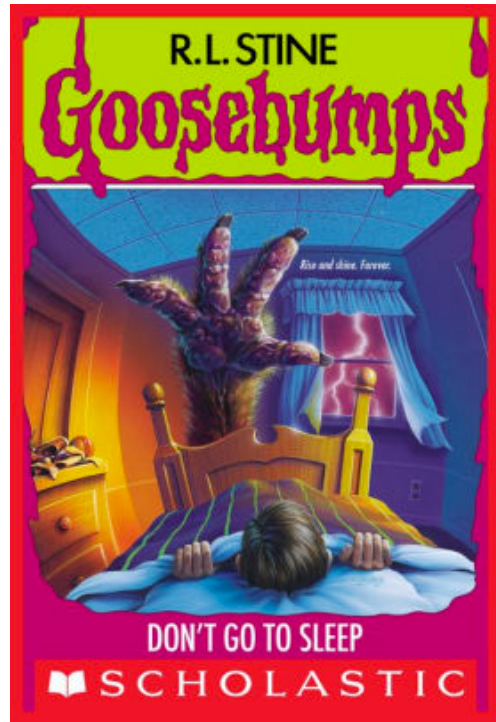
27  
28  
29  
30  
31  
...

...
∅
∅
∅
The Hobbit
∅
...

# A Hash Table based Dictionary

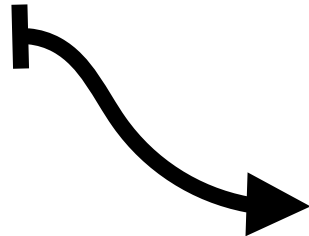


# A Hash Table based Dictionary



Author Name  
Hash Function


'R' + 'S' = 37



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

# A Hash Table based Dictionary

Aardvarks  
Anonymous

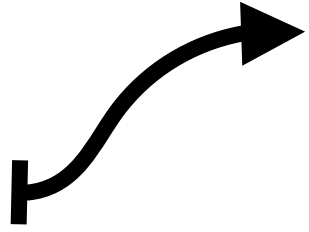


By Jim Truth



Author Name  
Hash Function

'J' + 'T' = 30



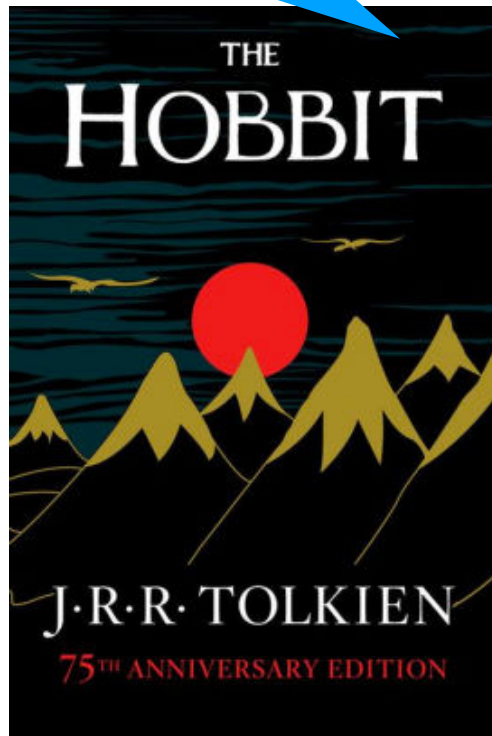
...	...
30	The Hobbit
31	∅
...	∅
37	Goosebumps
38	∅
...	...

# Hash Collision

How to handle this?

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

J.R.R Tolkien = 30!



Jim Truth = 30!



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

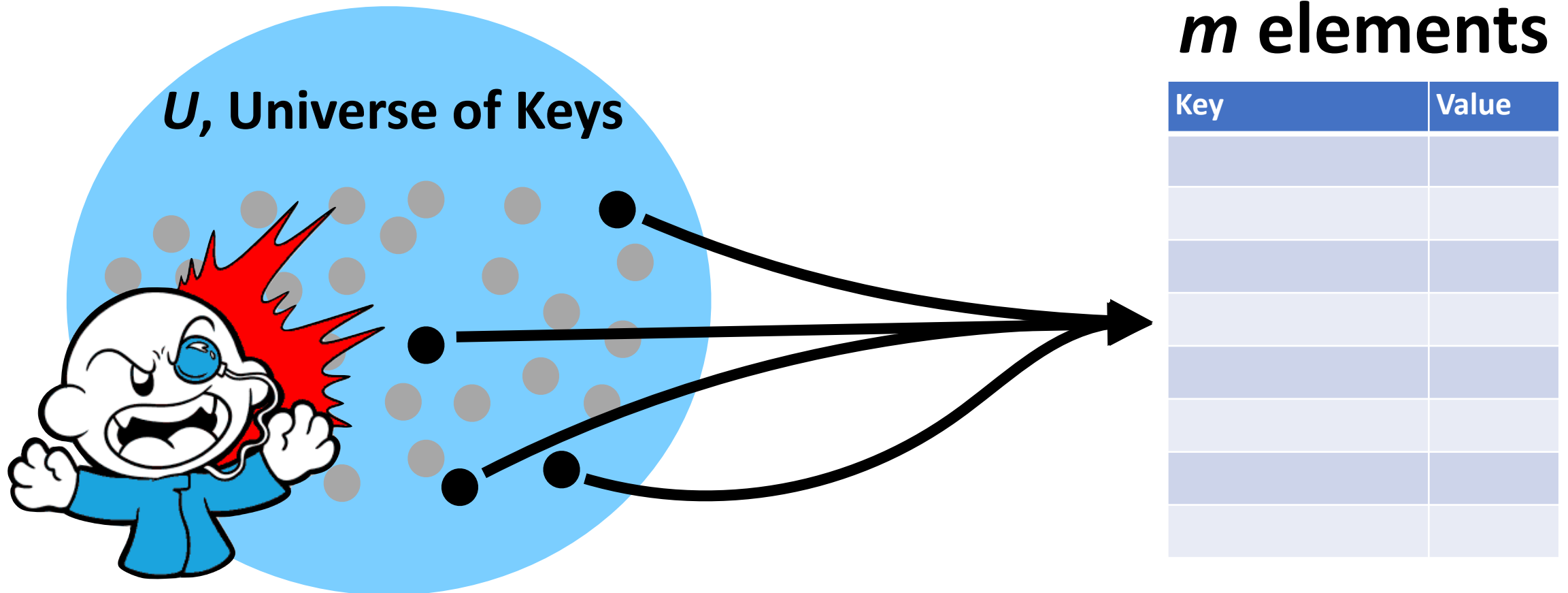






# General Purpose Hashing

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;
```

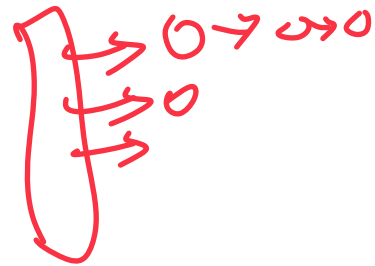
A **Hash Table** consists of three things:

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 to hash table



- **Closed Hashing:**

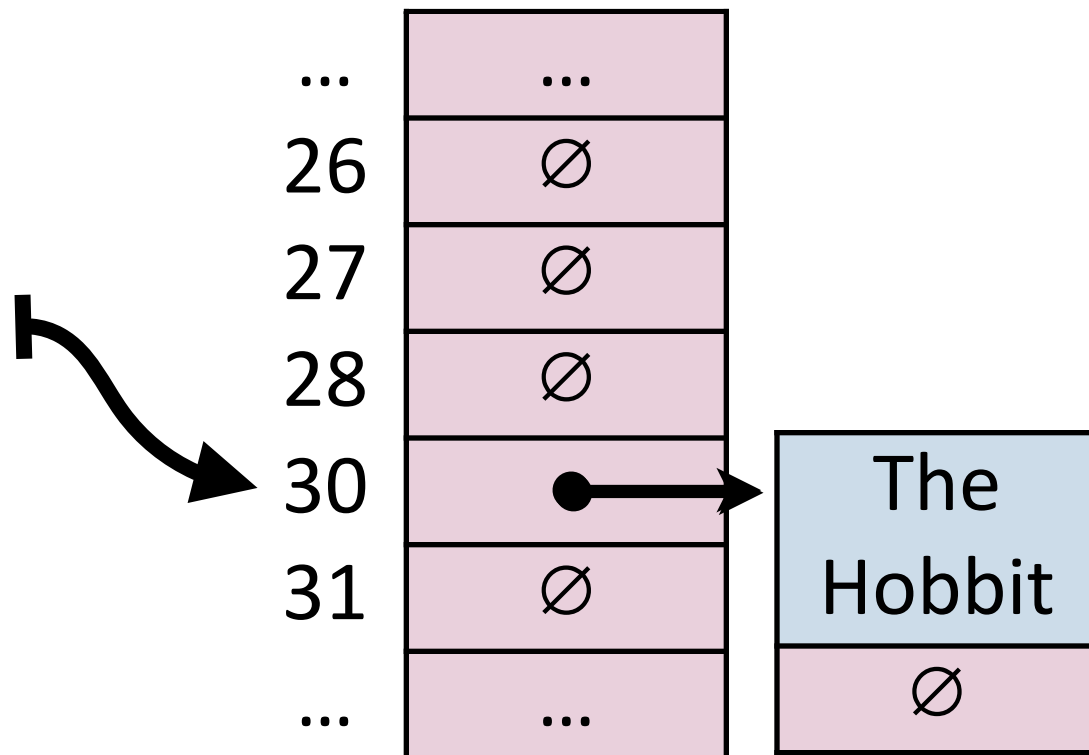
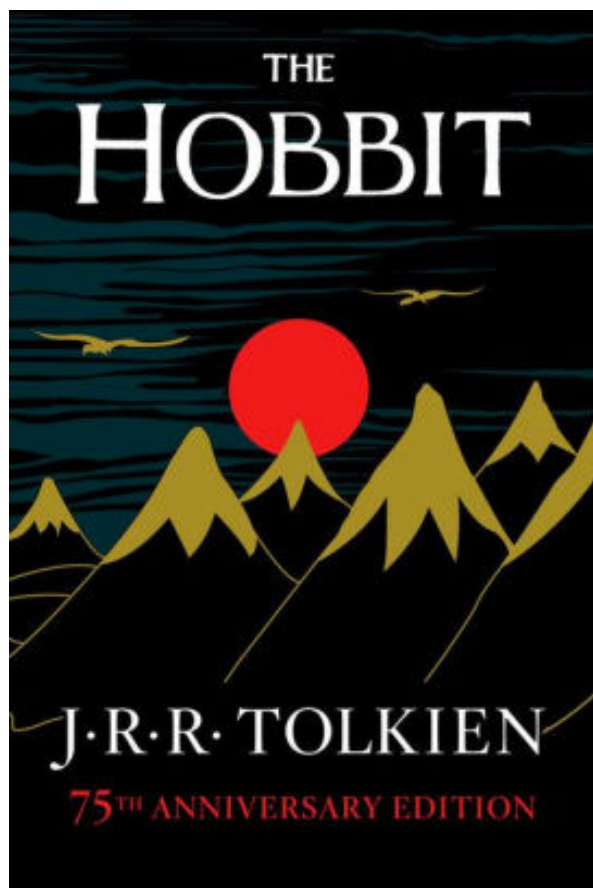
internally



# Open Hashing

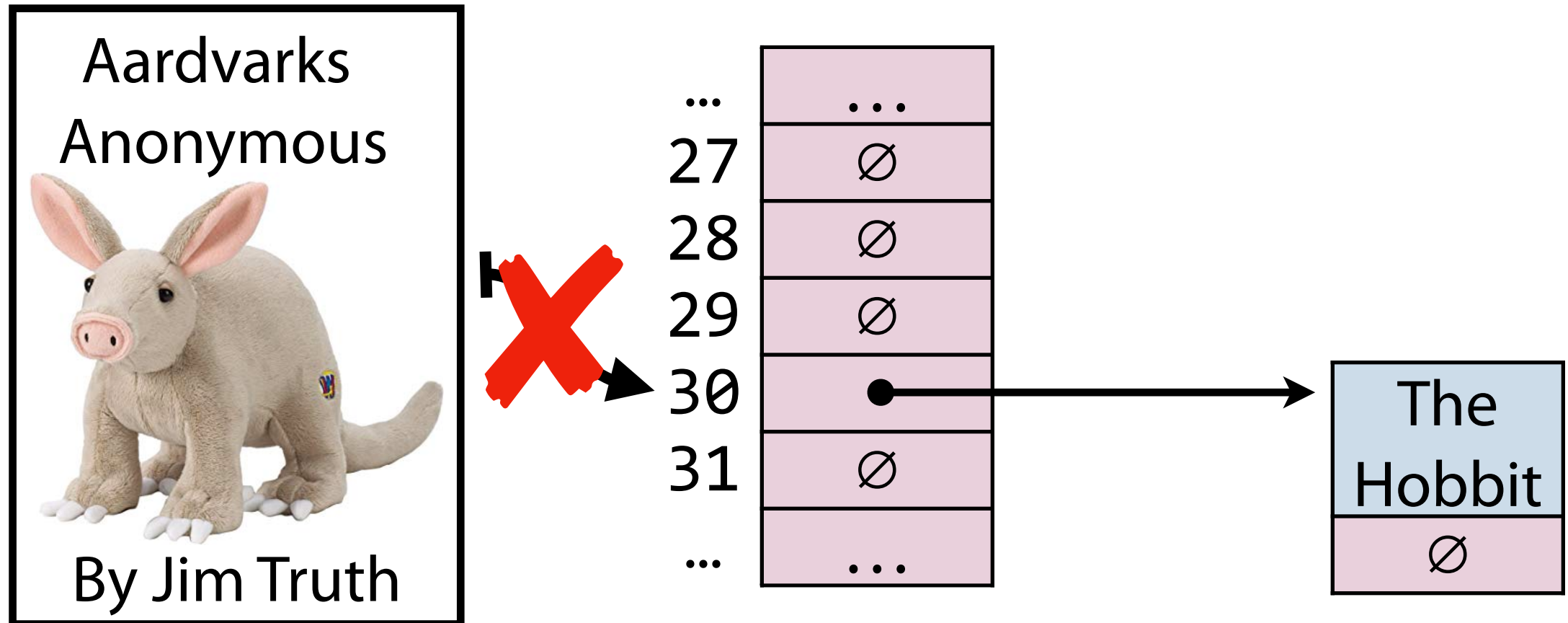
In an *open hashing* scheme, key-value pairs are stored externally (for example as a linked list).

↳ a list of k,v pairs!



# Hash Collisions (Open Hashing)

A **hash collision** in an open hashing scheme can be resolved by adding to linked list. This is called **separate chaining**.



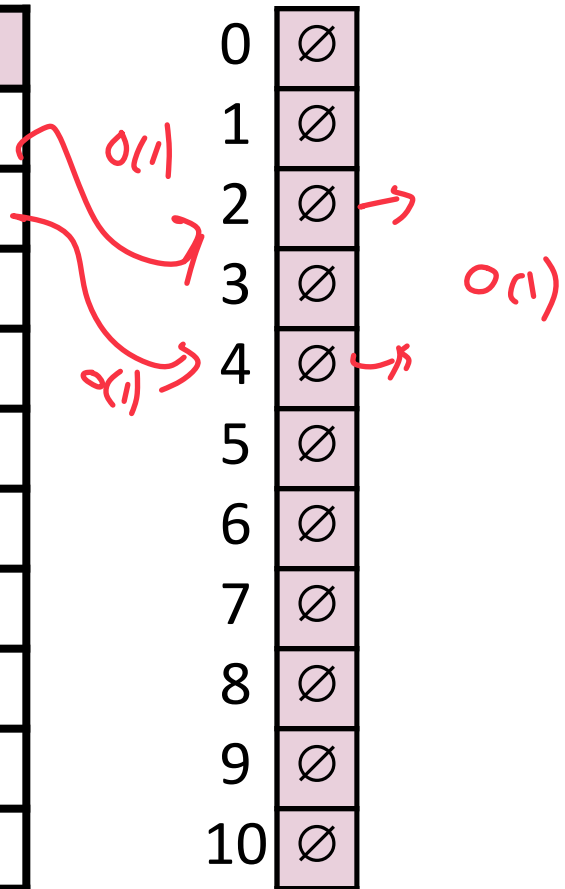
# Insertion (Separate Chaining)

`_insert("Bob")`

`_insert("Anna")`

$O(1)$

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

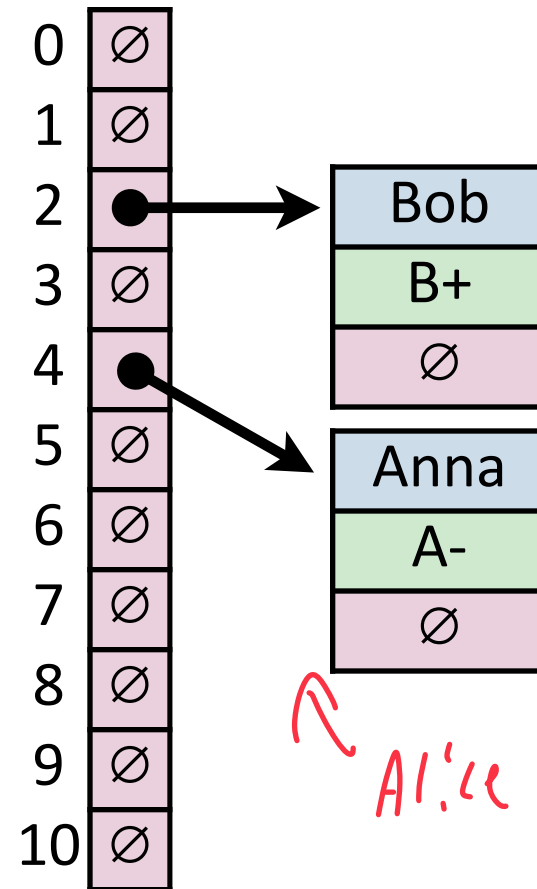


# Insertion (Separate Chaining)

`_insert("Alice")`

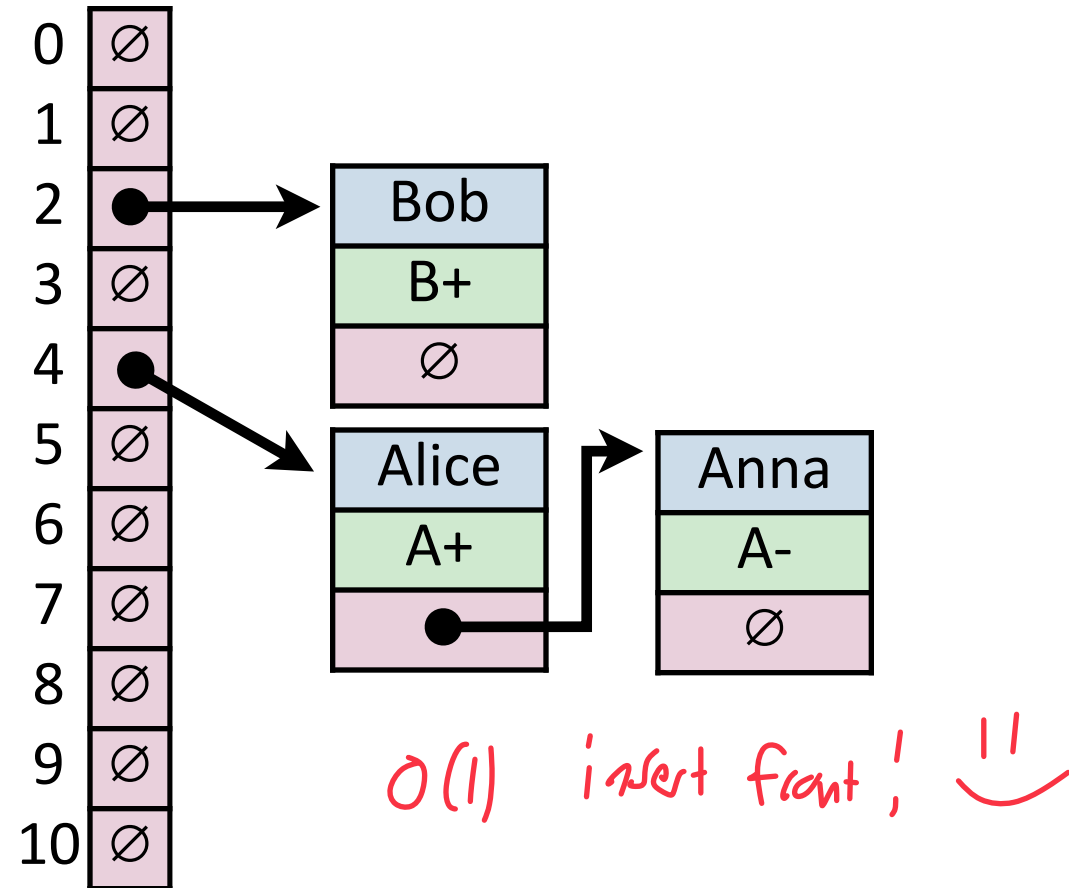
Where does Alice end up in the hash table?

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



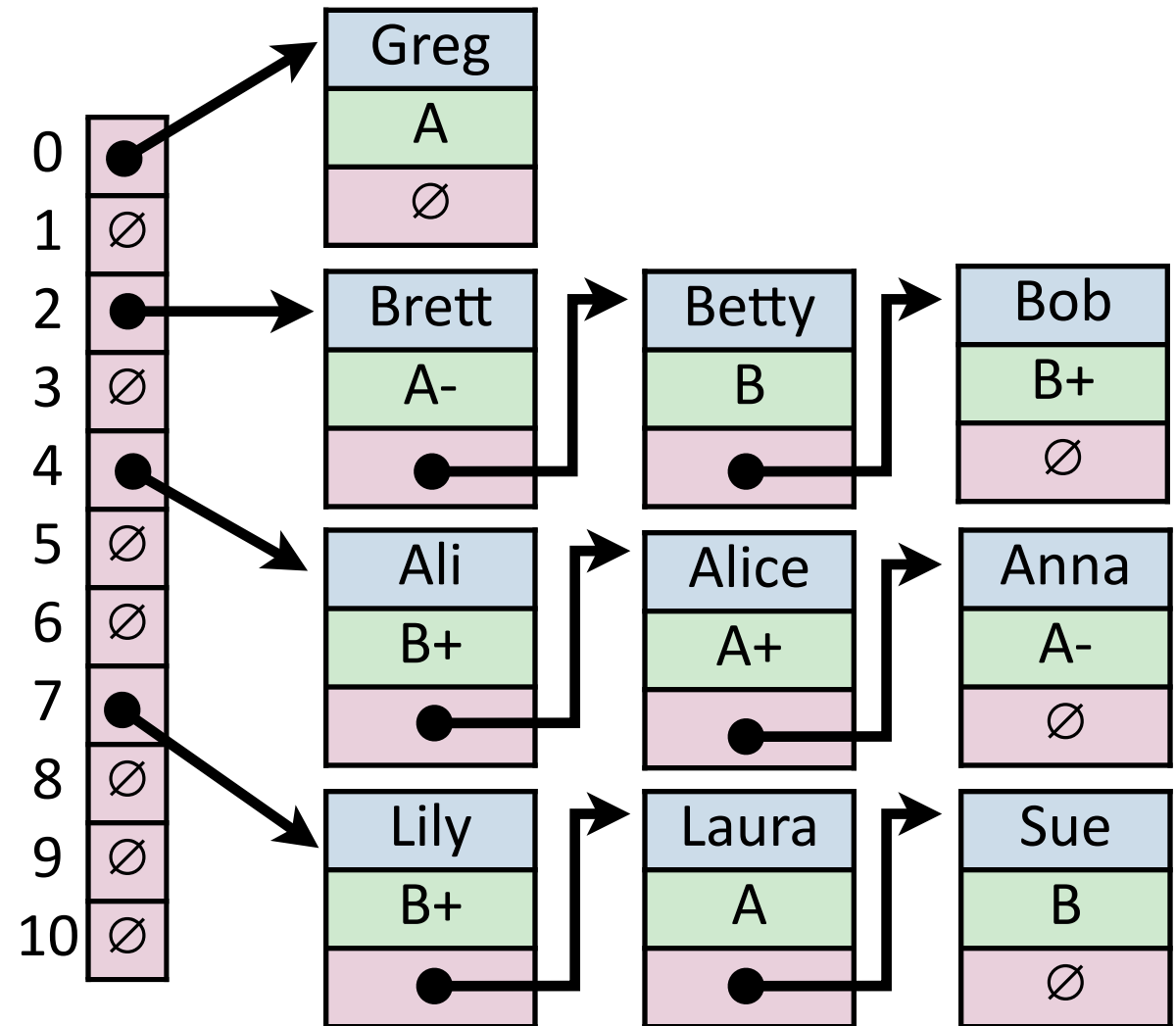
# Insertion (Separate Chaining)

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



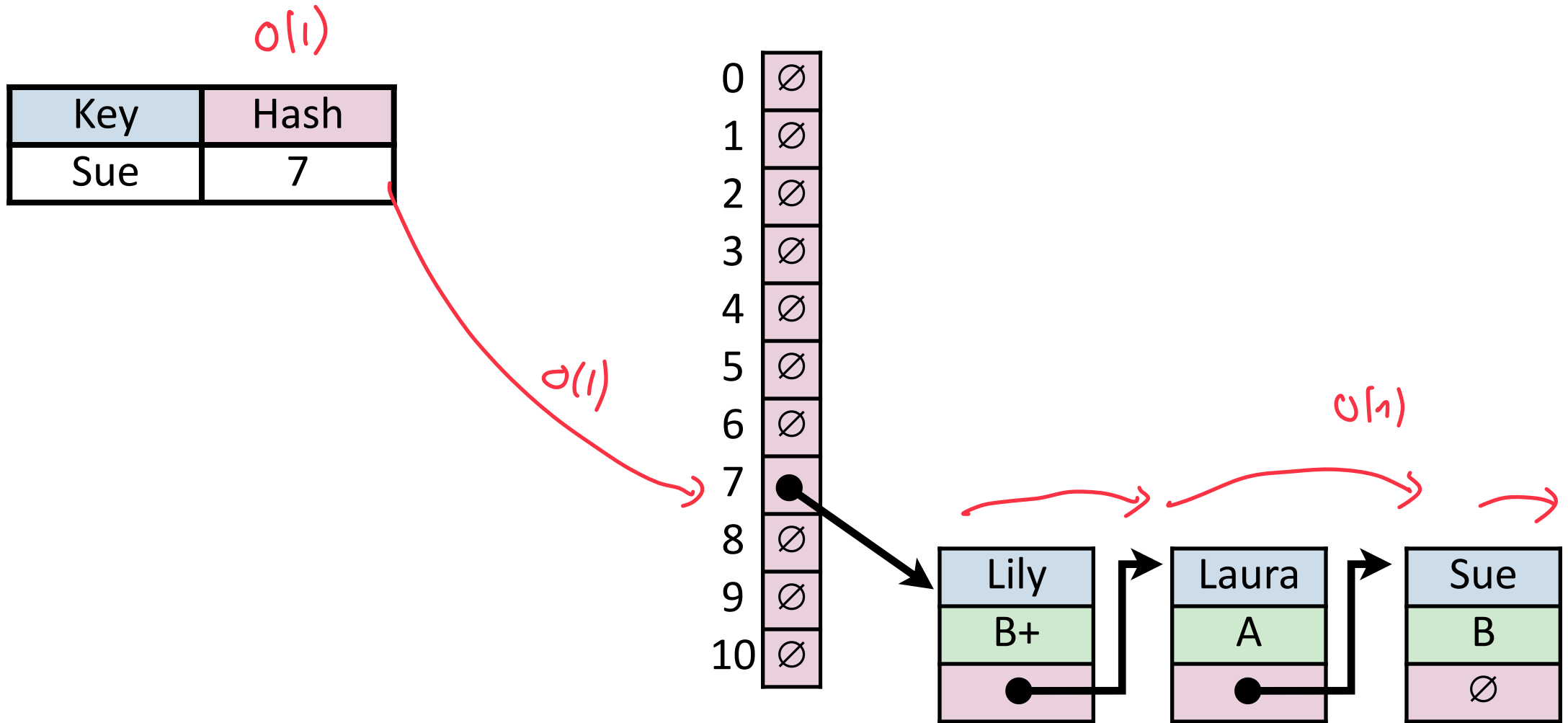
# Insertion (Separate Chaining)

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



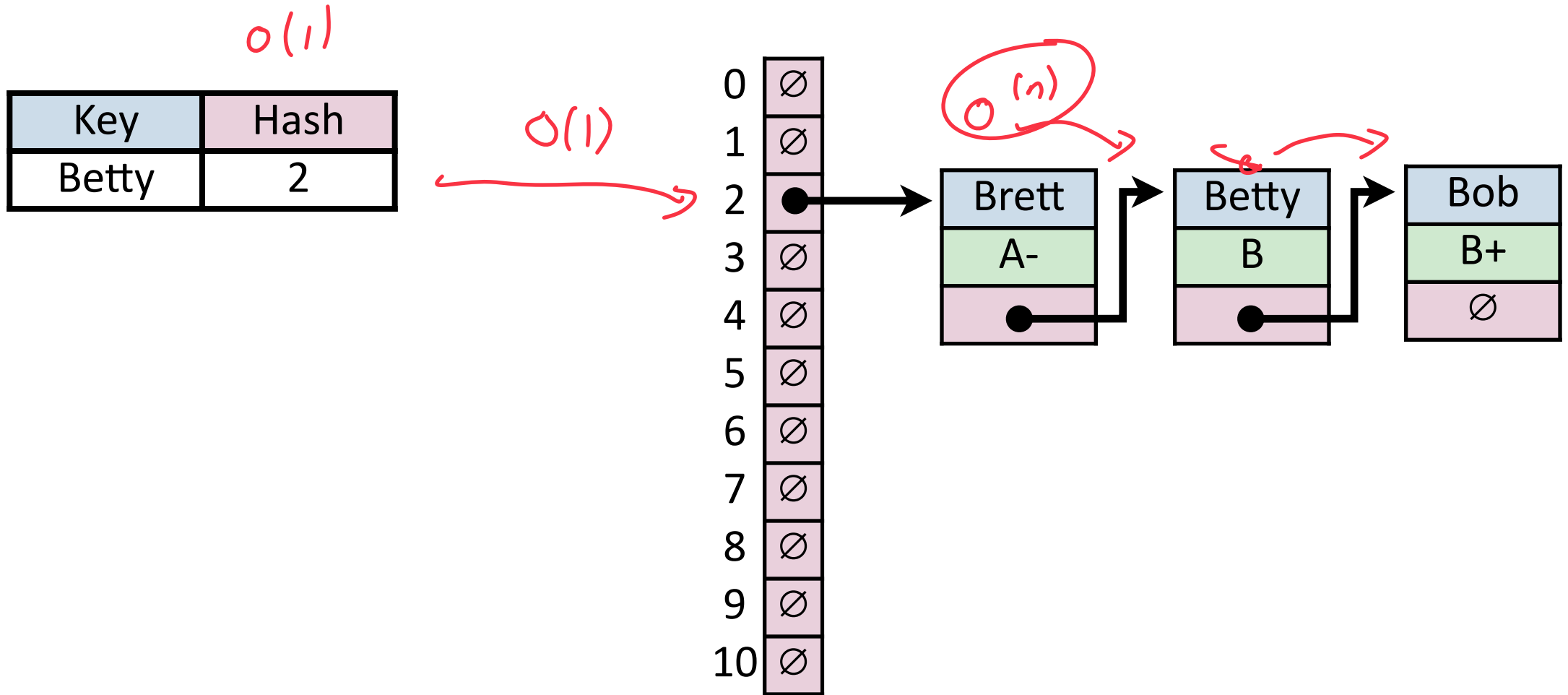
# Find (Separate Chaining)

`_find("Sue")`



# Remove (Separate Chaining)

`_remove("Betty")`





# Hash Table (Separate Chaining)

For hash table of size  $m$  and  $n$  elements:

Find runs in:  $O(n)$

(by key)

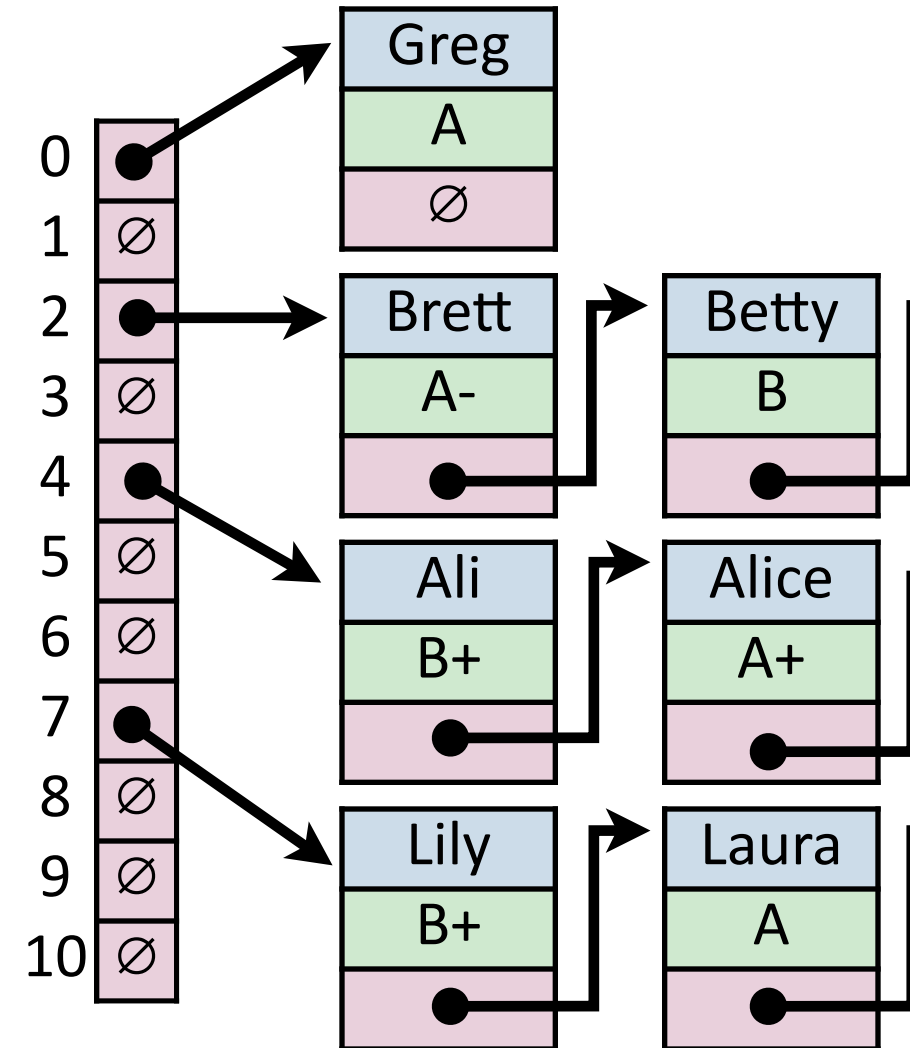
Insert runs in:  $O(1)$

Remove runs in:  $O(n)$

(by key)

∨  
=∨  
∧

How can randomness save us?



# Hash Table

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

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

2) Create a *universal hash function family*:

# 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:**

**Independent:**

# Separate Chaining Under SUHA

Table Size:  $m$

Num objects:  $n$

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

$\alpha_j$  = expected # of items hashing to position  $j$

$$\alpha_j = \sum_i H_{i,j}$$

$$H_{i,j} = \begin{cases} 1 & \text{if item } i \text{ hashes to } j \\ 0 & \text{otherwise} \end{cases}$$

# Separate Chaining Under SUHA

Table Size:  $m$

Num objects:  $n$

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

$\alpha_j$  = expected # of items hashing to position  $j$

$$\alpha_j = \sum_i H_{i,j}$$

$$H_{i,j} = \begin{cases} 1 & \text{if item } i \text{ hashes to } j \\ 0 & \text{otherwise} \end{cases}$$

$$E[\alpha_j] = E\left[\sum_i H_{i,j}\right]$$

# Separate Chaining Under SUHA

Table Size:  $m$

Num objects:  $n$

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

$\alpha_j$  = expected # of items hashing to position  $j$

$$\alpha_j = \sum_i H_{i,j}$$

$$H_{i,j} = \begin{cases} 1 & \text{if item } i \text{ hashes to } j \\ 0 & \text{otherwise} \end{cases}$$

$$E[\alpha_j] = E\left[\sum_i H_{i,j}\right]$$

$$E[\alpha_j] = \sum_i Pr(H_{i,j} = 1) * 1 + Pr(H_{i,j} = 0) * 0$$

# Separate Chaining Under SUHA

Table Size:  $m$

Num objects:  $n$

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

$\alpha_j$  = expected # of items hashing to position  $j$

$$\alpha_j = \sum_i H_{i,j}$$

$$H_{i,j} = \begin{cases} 1 & \text{if item } i \text{ hashes to } j \\ 0 & \text{otherwise} \end{cases}$$

$$E[\alpha_j] = E\left[\sum_i H_{i,j}\right]$$

$$E[\alpha_j] = \sum_i Pr(H_{i,j} = 1) * 1 + Pr(H_{i,j} = 0) * 0$$

$$E[\alpha_j] = n * Pr(H_{i,j} = 1)$$

# Separate Chaining Under SUHA

Table Size:  $m$

Num objects:  $n$

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

$\alpha_j$  = expected # of items hashing to position  $j$

$$\alpha_j = \sum_i H_{i,j}$$

$$H_{i,j} = \begin{cases} 1 & \text{if item } i \text{ hashes to } j \\ 0 & \text{otherwise} \end{cases}$$

$$E[\alpha_j] = E\left[\sum_i H_{i,j}\right]$$

$$Pr[H_{i,j} = 1] = \frac{1}{m}$$

$$E[\alpha_j] = n * Pr(H_{i,j} = 1)$$

# Separate Chaining Under SUHA



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

$\alpha_j$  = expected # of items hashing to position  $j$

**Num objects:  $n$**

$$\alpha_j = \sum_i H_{i,j}$$

$$H_{i,j} = \begin{cases} 1 & \text{if item } i \text{ hashes to } j \\ 0 & \text{otherwise} \end{cases}$$

$$E[\alpha_j] = E\left[\sum_i H_{i,j}\right]$$

$$Pr[H_{i,j} = 1] = \frac{1}{m}$$

$$E[\alpha_j] = n * Pr(H_{i,j} = 1)$$

$$\mathbf{E}[\alpha_j] = \frac{\mathbf{n}}{\mathbf{m}}$$

# Separate Chaining Under SUHA

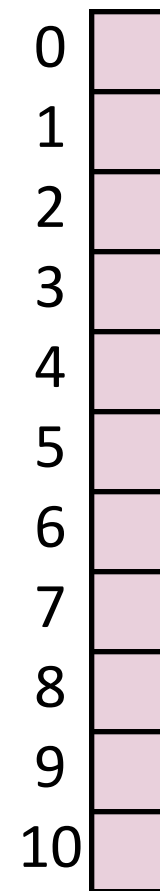


**Under SUHA, a hash table of size  $m$  and  $n$  elements:**

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

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

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

