

# Late class start: Enjoy the solar eclipse!

Class begins at 2:15 PM.



# Algorithms and Data Structures for Data Science

## Hashing

CS 277

April 8, 2024

Brad Solomon

Graphs

Algorithms



UNIVERSITY OF  
**ILLINOIS**  
URBANA - CHAMPAIGN

Department of Computer Science

Dept

Hashing



# Learning Objectives

Motivate and define a hash table

(Dictionary in Python)

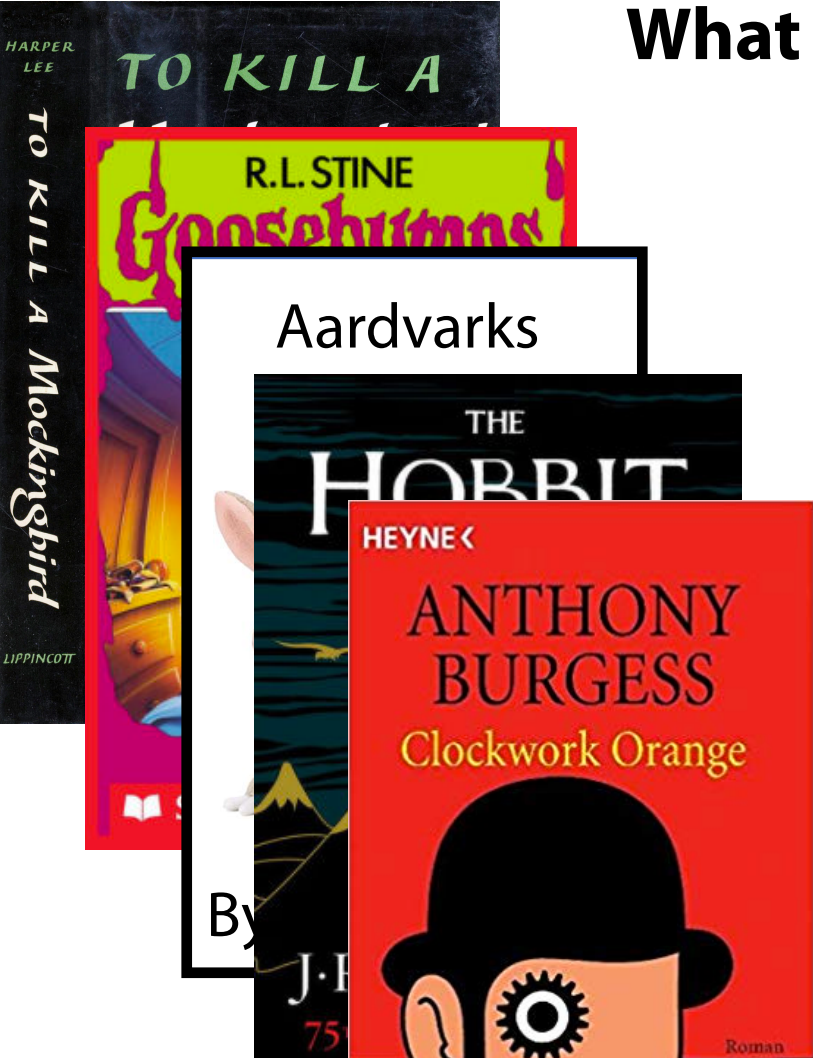
Discuss what a 'good' hash function looks like

Identify a key weakness of the 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?

Key → Value

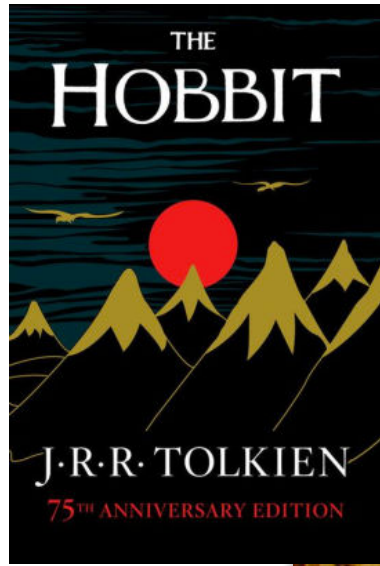
Book

Contents

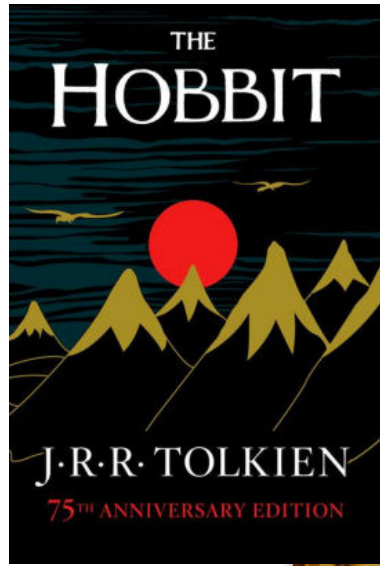
2 lists      keys      values      match their indexing

Tree (Binary Search Tree)

If we recognize that libraries are ordered:  $O(\log n)$

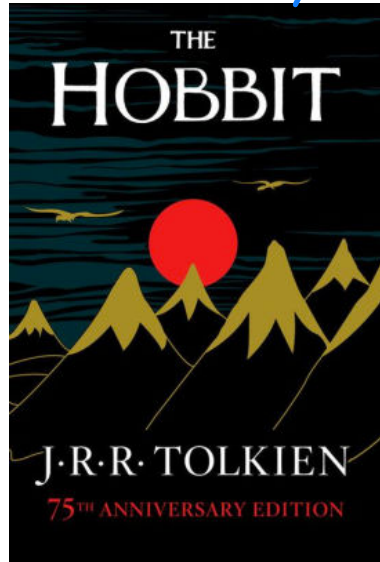


# What if $O(\log n)$ isn't good enough?



# A Hash Table based Dictionary

Key



hash #  
address

ISBN: 9781526602381  
Call #: PR  
6068.093  
H35 1998

ISBN: 9781526602381

Call #: PR

6068.093

H35 1998

Array  
list



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

$O(1)$

Hash  
function

(address)  
# integer

$O(1)$

# A Hash Table based Dictionary



```
1 d = {}  
2 d[k] = v  
3 print(d[k])
```

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

1. A hash function (key  $\rightarrow$  int)
2. A list (stores our data @ int)
3. ???





# Hash Function

A hash function **must** be:

- **Deterministic:** Given same key twice, return same value
- **Efficient:**  $O(1)$
- **Defined for a certain size table:** Universe  $\rightarrow 0, \dots, M-1$   
 $M$  unique values



# Hash Function

I to 1 mapping of these specific names

to this string array

- (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)

*Alkuni*

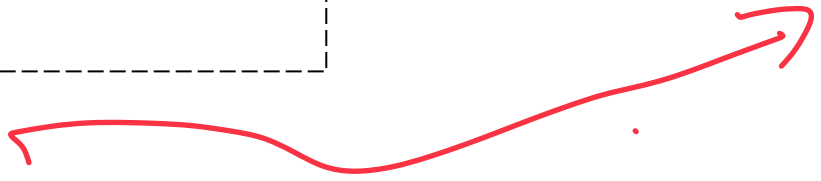
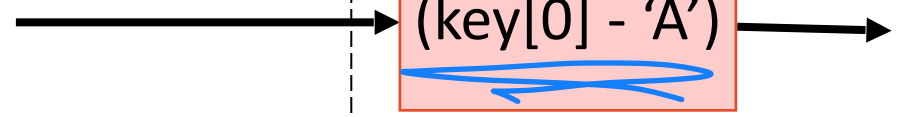
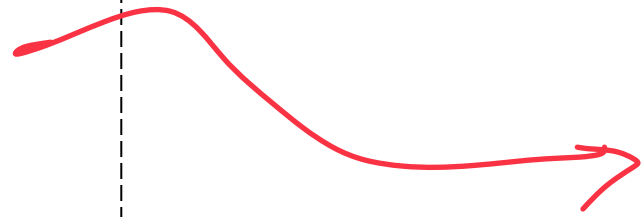
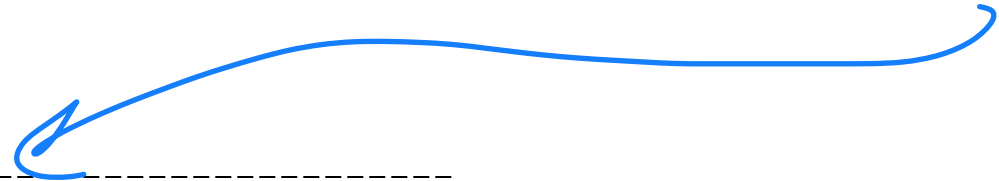
Hash function

$(\text{key}[0] - \text{'A'})$

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

*Solomon 277*

*H?*





# General Hash Function

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

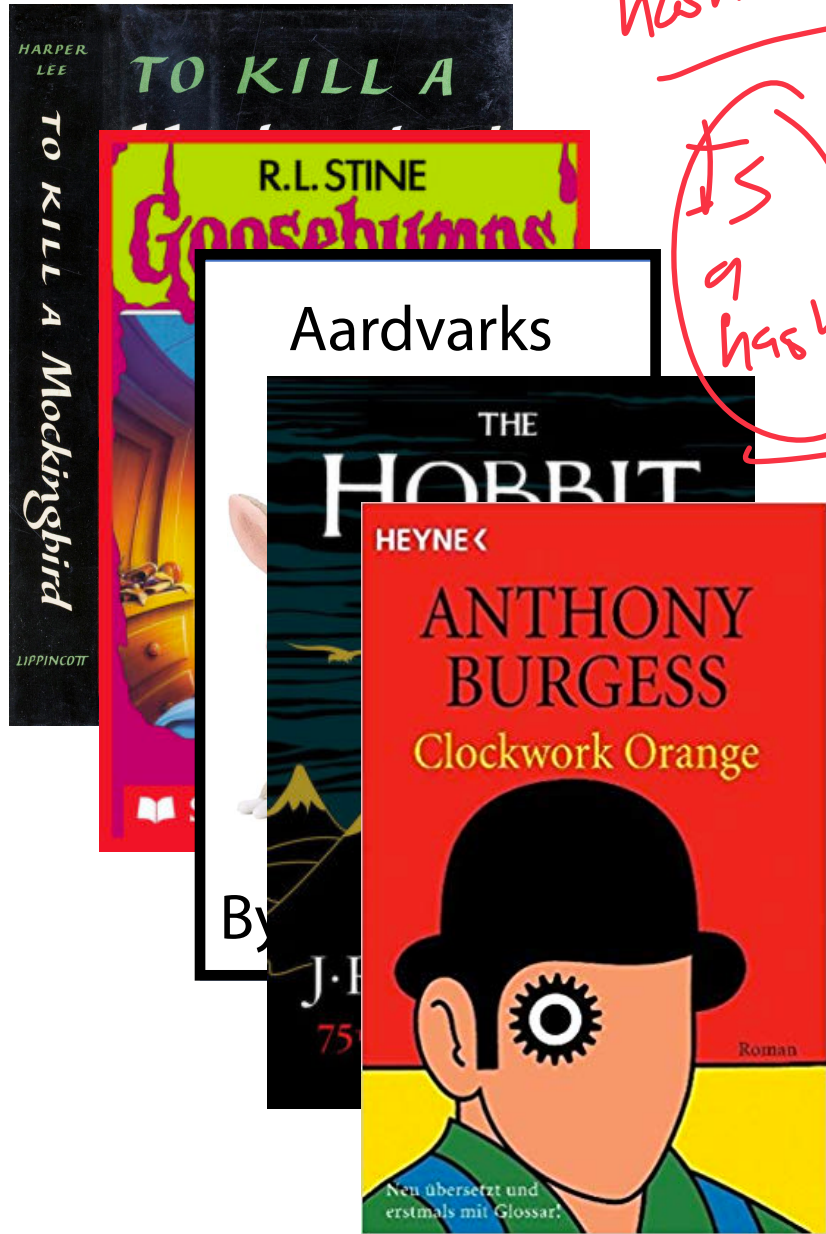
- **A hash:** Function that converts any  $k$  in universe to a  $\#$   
( $H$  could be any number)      hash  $\# \ \% \ \frac{m}{\quad}$
- **A compression:** Takes our  $\#$  and converts to  $[0, \dots, m-1]$   
 $[0, 1] \leftarrow X \% 2$  (Modulo or Remainder)

Choosing a good hash function is tricky...

- Don't create your own!

# Hash Function

Is it  
hash function?



Is  
a  
hash

Authors 2

Authors →

$$h(k) = (k.firstName[0] + k.lastName[0]) \% m$$

↳ Deterministic ✓

↳ Efficient ✓

↳ universe of all books?

↳ (0, ..., m-1)

need work ???

$$h(k) = (\text{rand()} * k.numPages) \% m$$

↳ Not a hash function

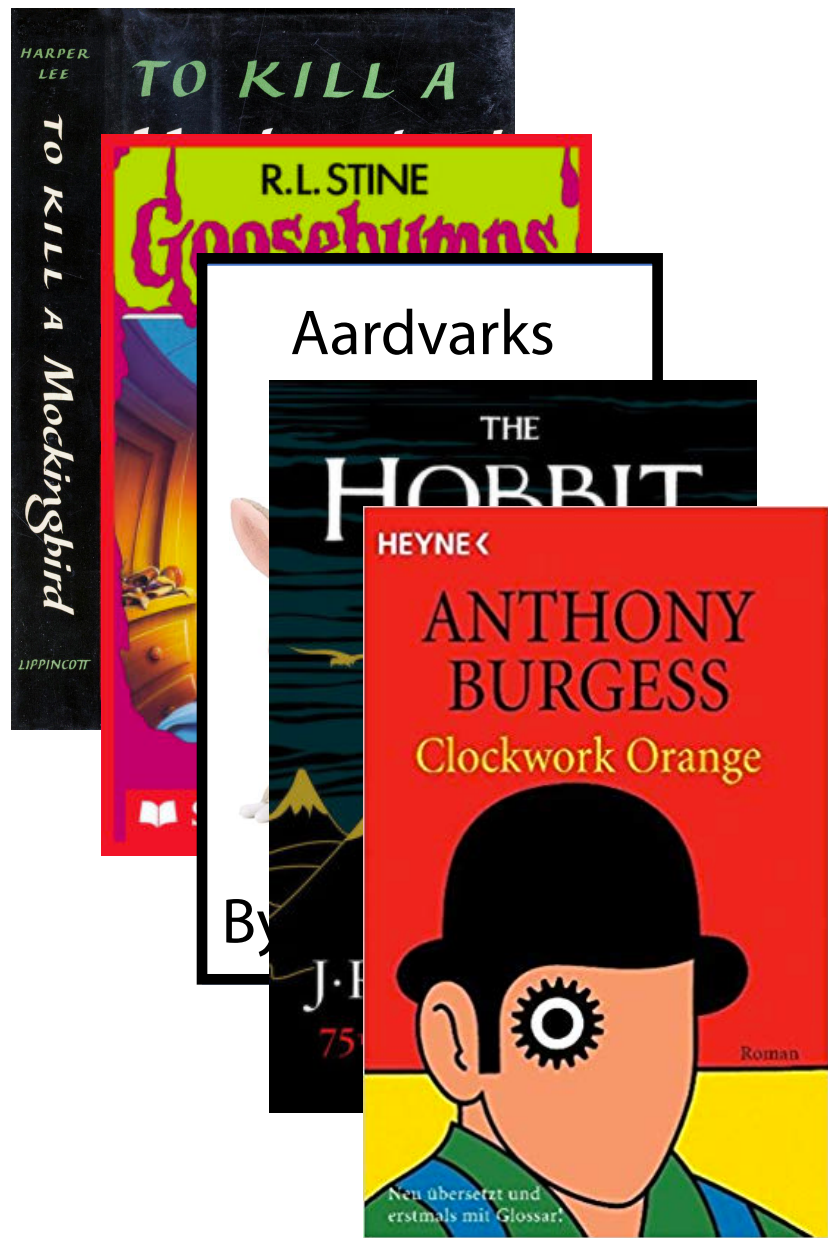
↳ Not deterministic

$$h(k) = (\text{Order I insert} [\text{Order seen}]) \% m$$

↳ Not a hash function

↳ Not deterministic

# Hash Function

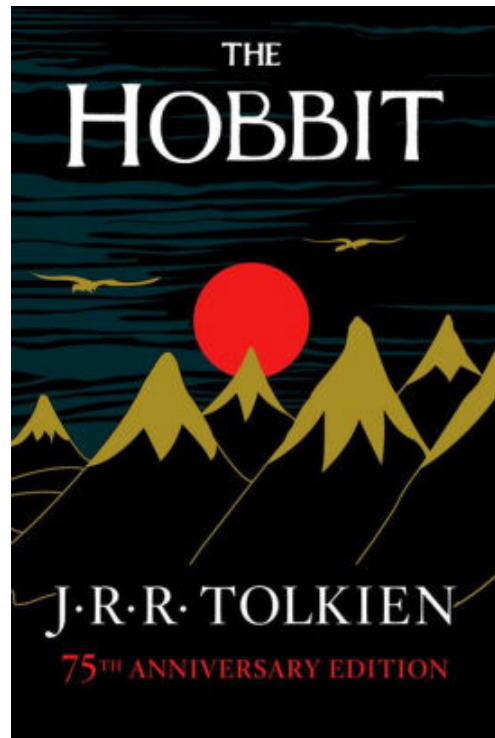


Exercise for viewer

↳ Some hash function

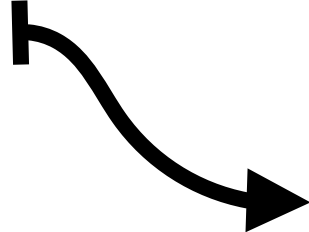
for books

# A Hash Table based Dictionary



Author Name  
Hash Function

'J' + 'T' = 30

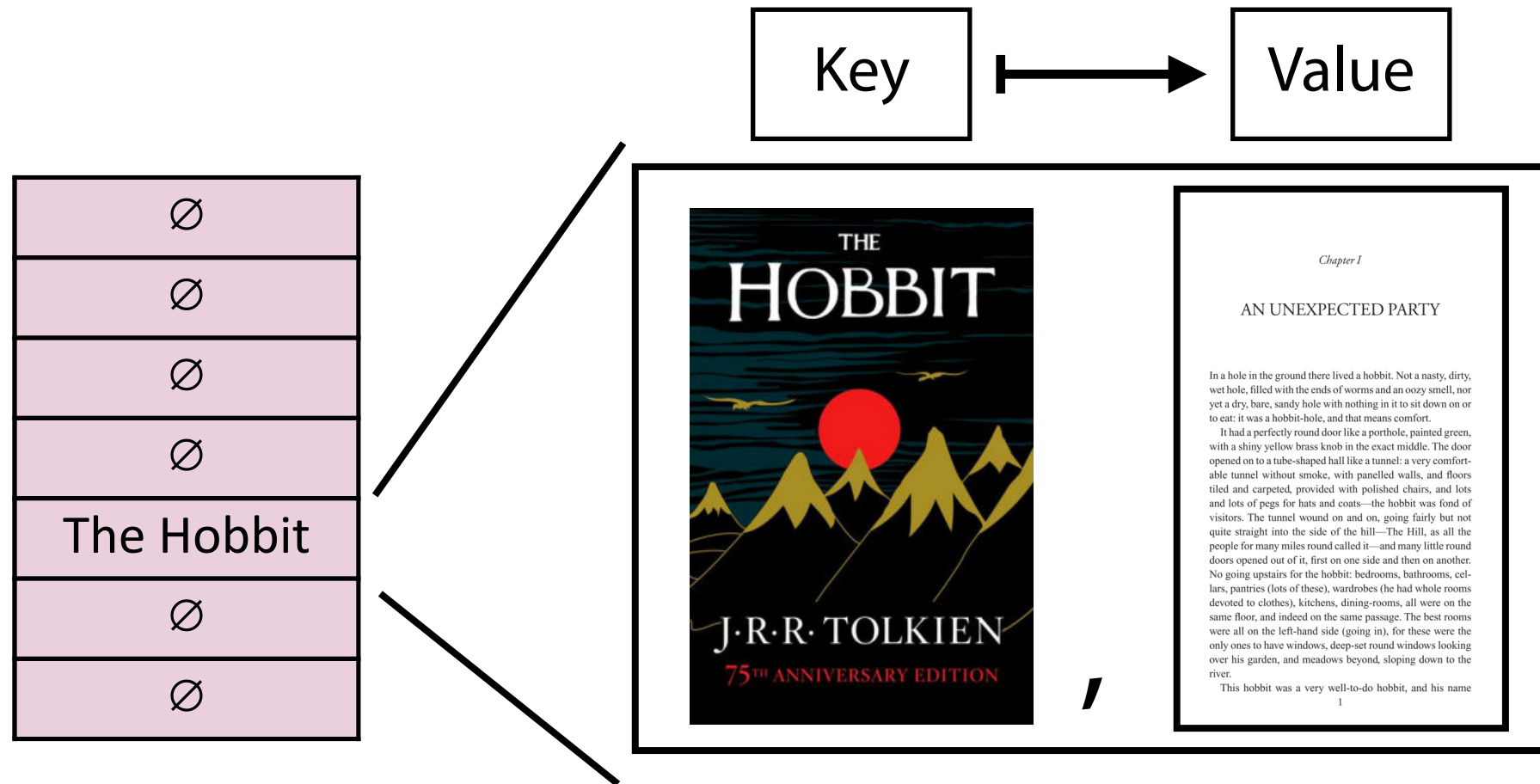


27  
28  
29  
30  
31  
...

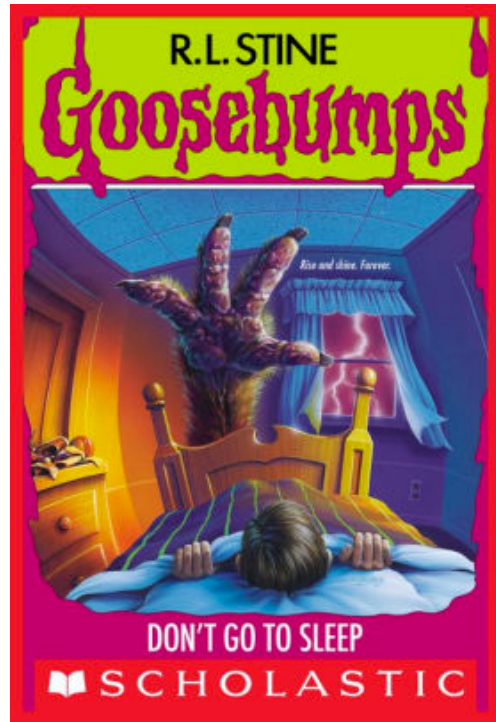
...
∅
∅
∅
The Hobbit
∅
...



# A Hash Table based Dictionary

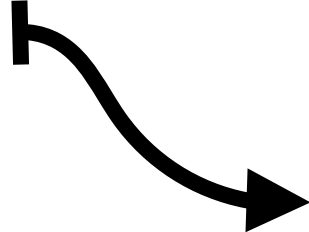


# A Hash Table based Dictionary



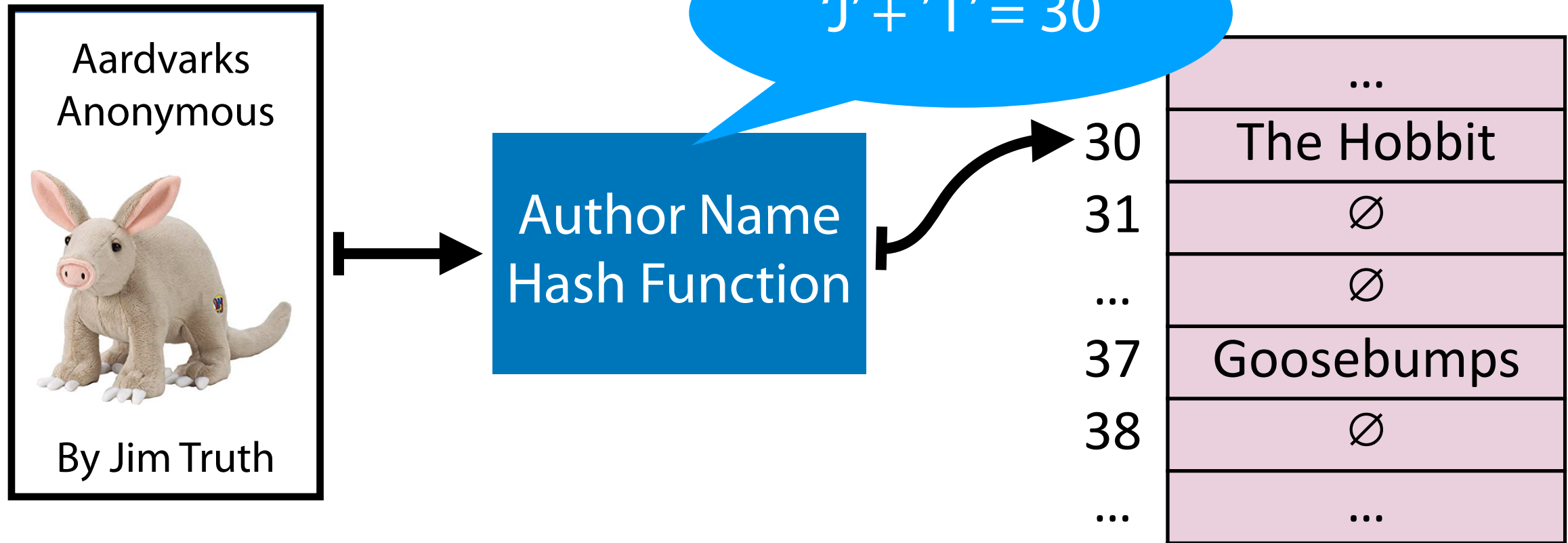
Author Name  
Hash Function

'R' + 'S' = 37



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

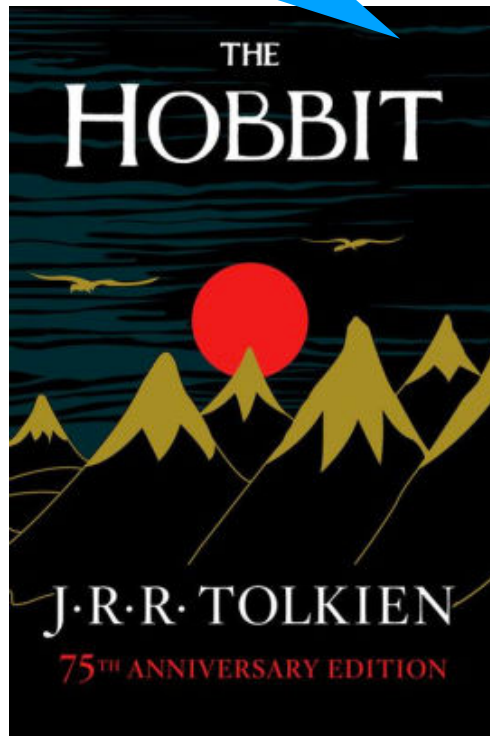
# A Hash Table based Dictionary



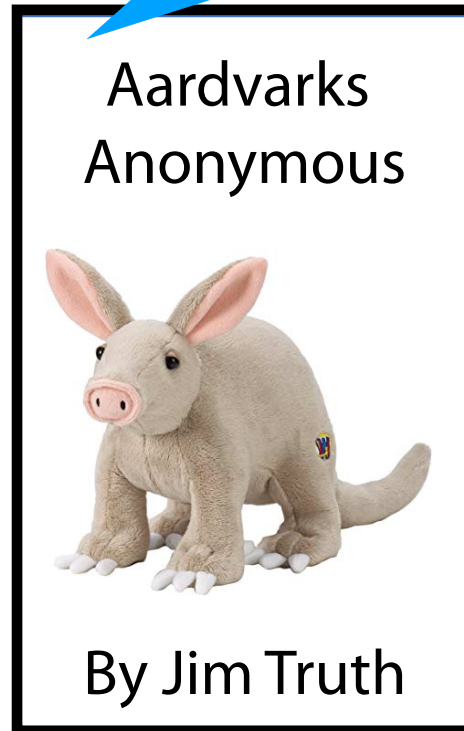
# Hash Collision

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	∅
...	...

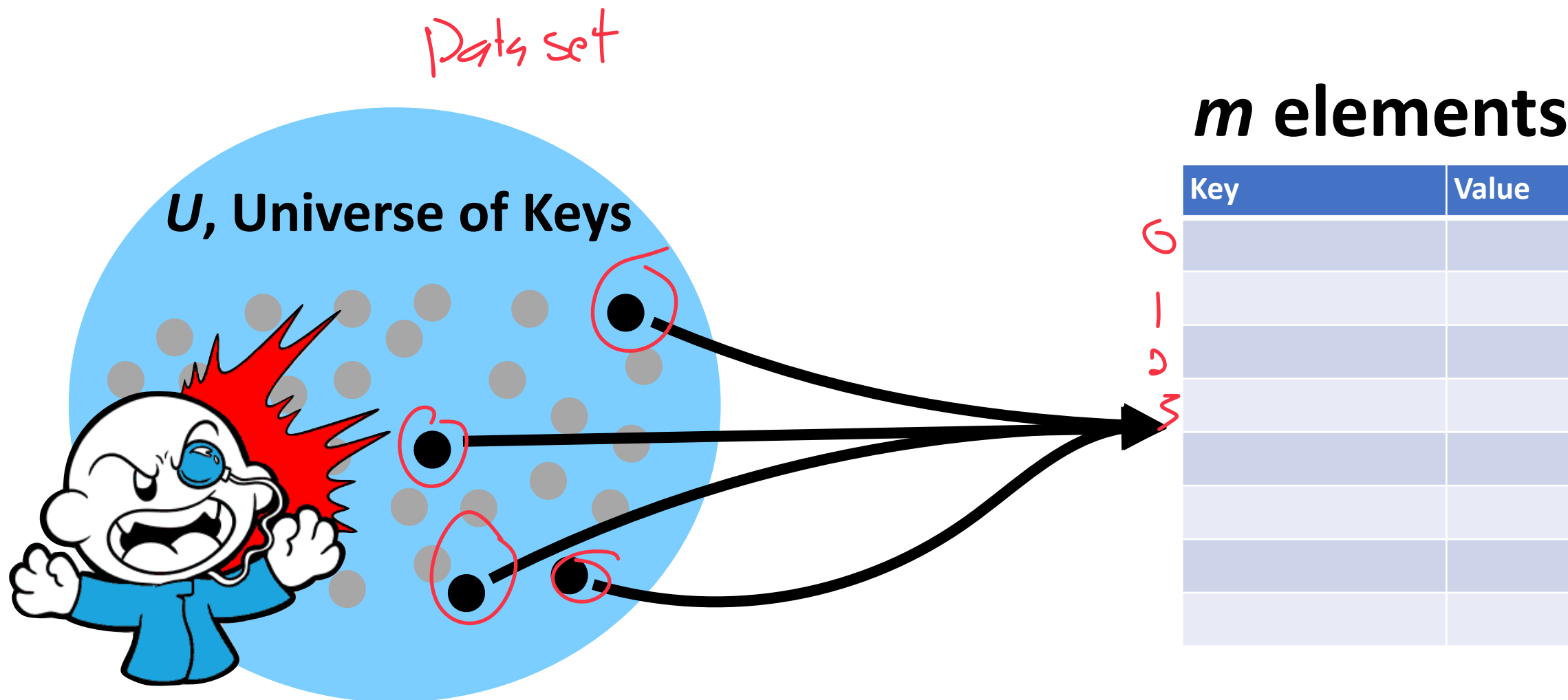






# General Purpose Hashing

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:

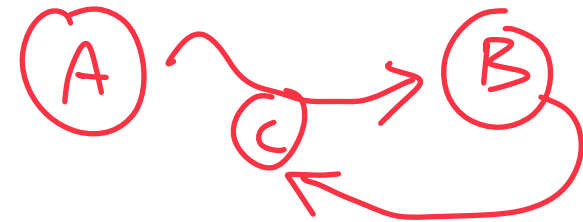
1. A hash function
2. A data storage structure
- 3. A method of addressing *hash collisions***

# Open vs Closed Hashing

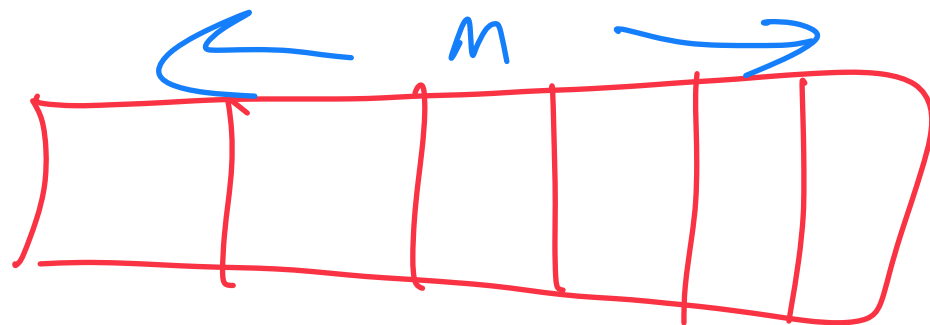
Addressing hash collisions depends on your storage structure.

- **Open Hashing:** Stores key, value externally

↳ Linked list (can store individual nodes anywhere)



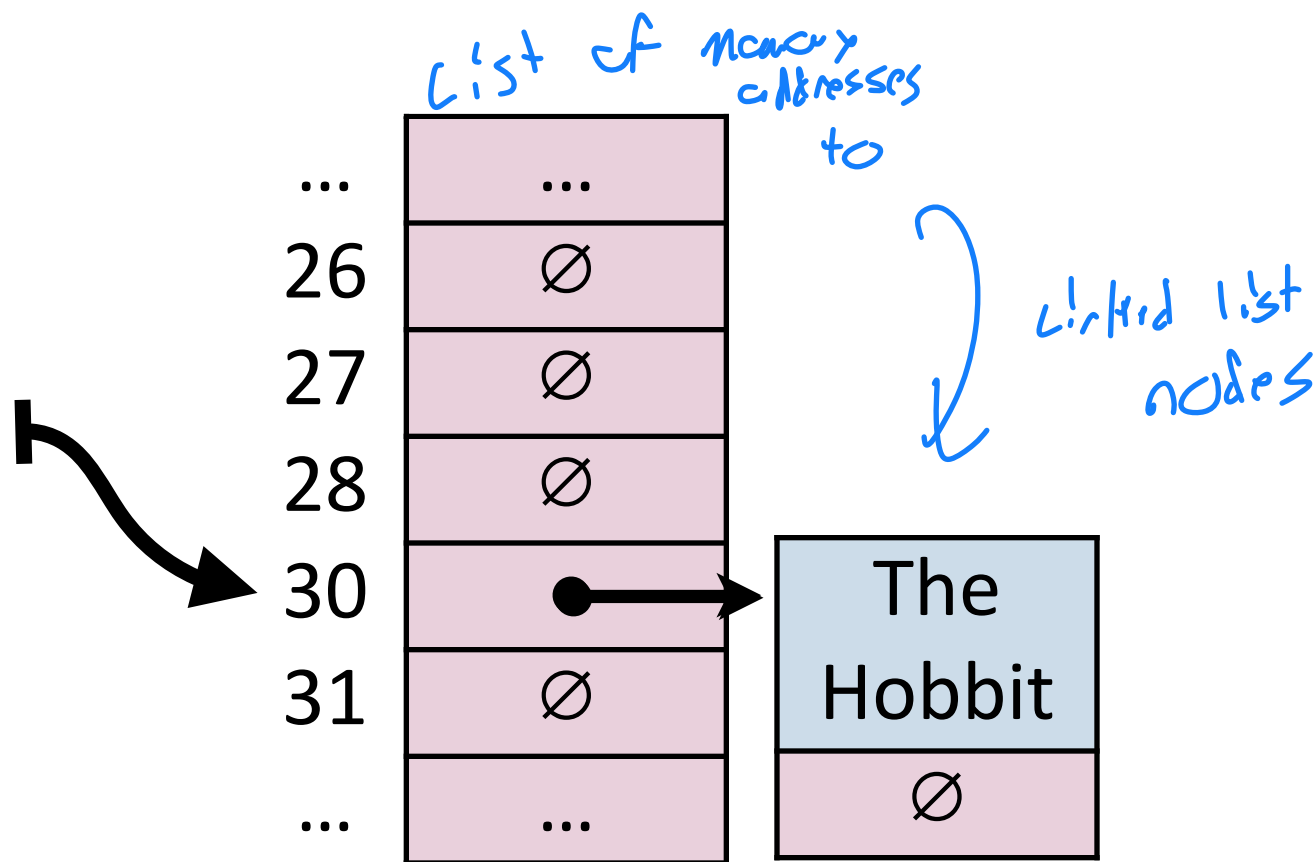
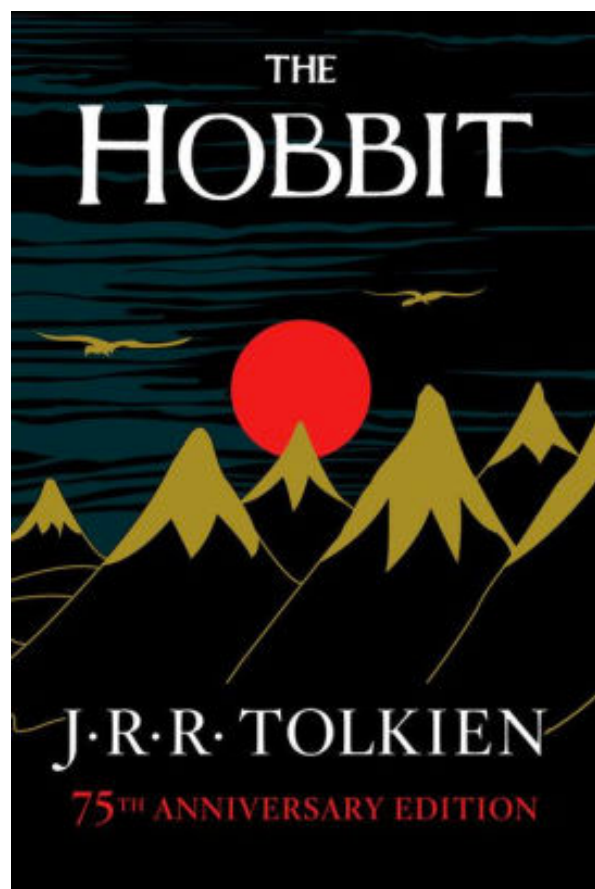
- **Closed Hashing:** Stores  $k, v$  internally



allocate fixed memory

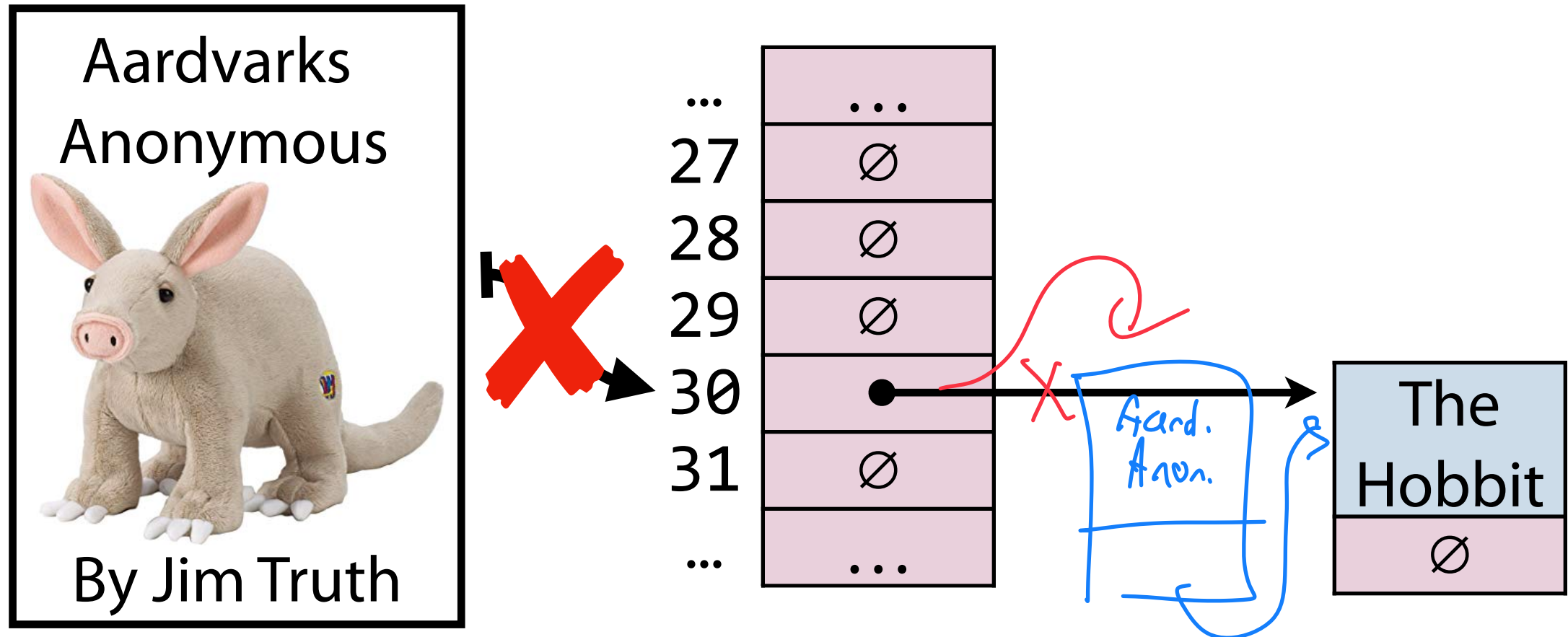
# Open Hashing

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 adding to the linked list. This is called **separate chaining**.

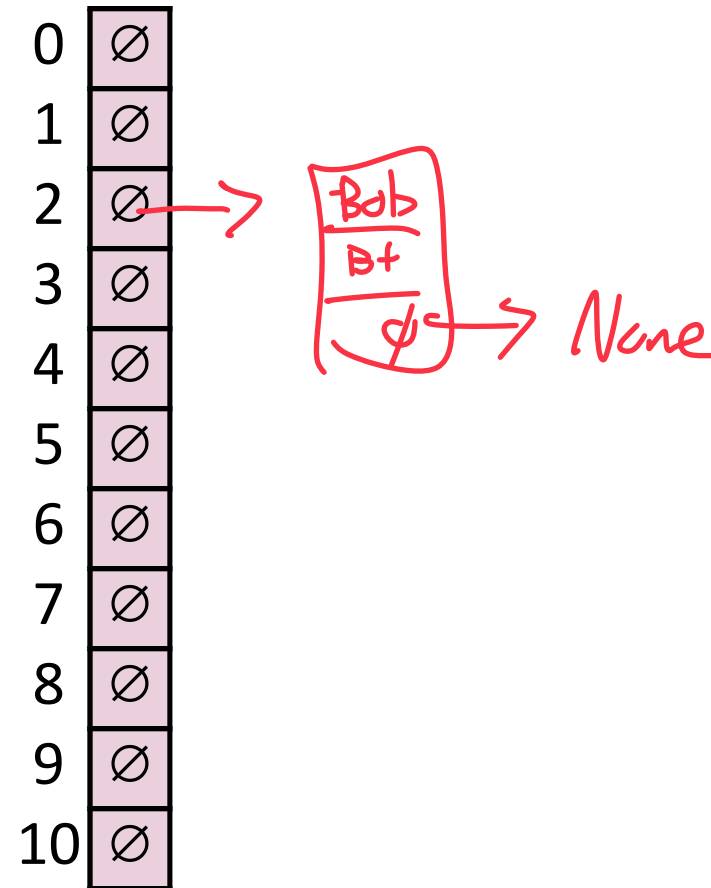


# Insertion (Separate Chaining)

`_insert("Bob")`

`_insert("Anna")`

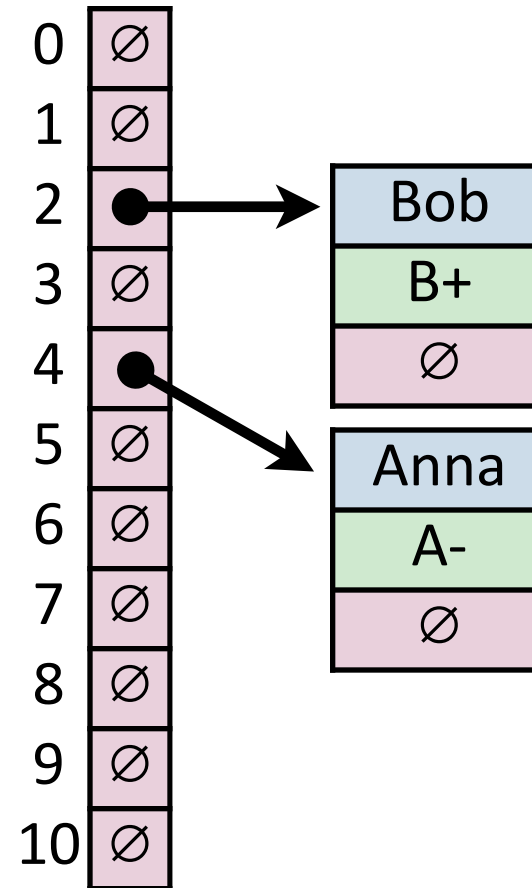
Key	Value	Hash
<b>Bob</b>	<b>B+</b>	<b>2</b>
<b>Anna</b>	<b>A-</b>	<b>4</b>
Alice	A+	4
Betty	B	2
Brett	A-	2
Greg	A	0
Sue	B	7
Ali	B+	4
Laura	A	7
Lily	B+	7



# Insertion (Separate Chaining)

`_insert("Alice")`

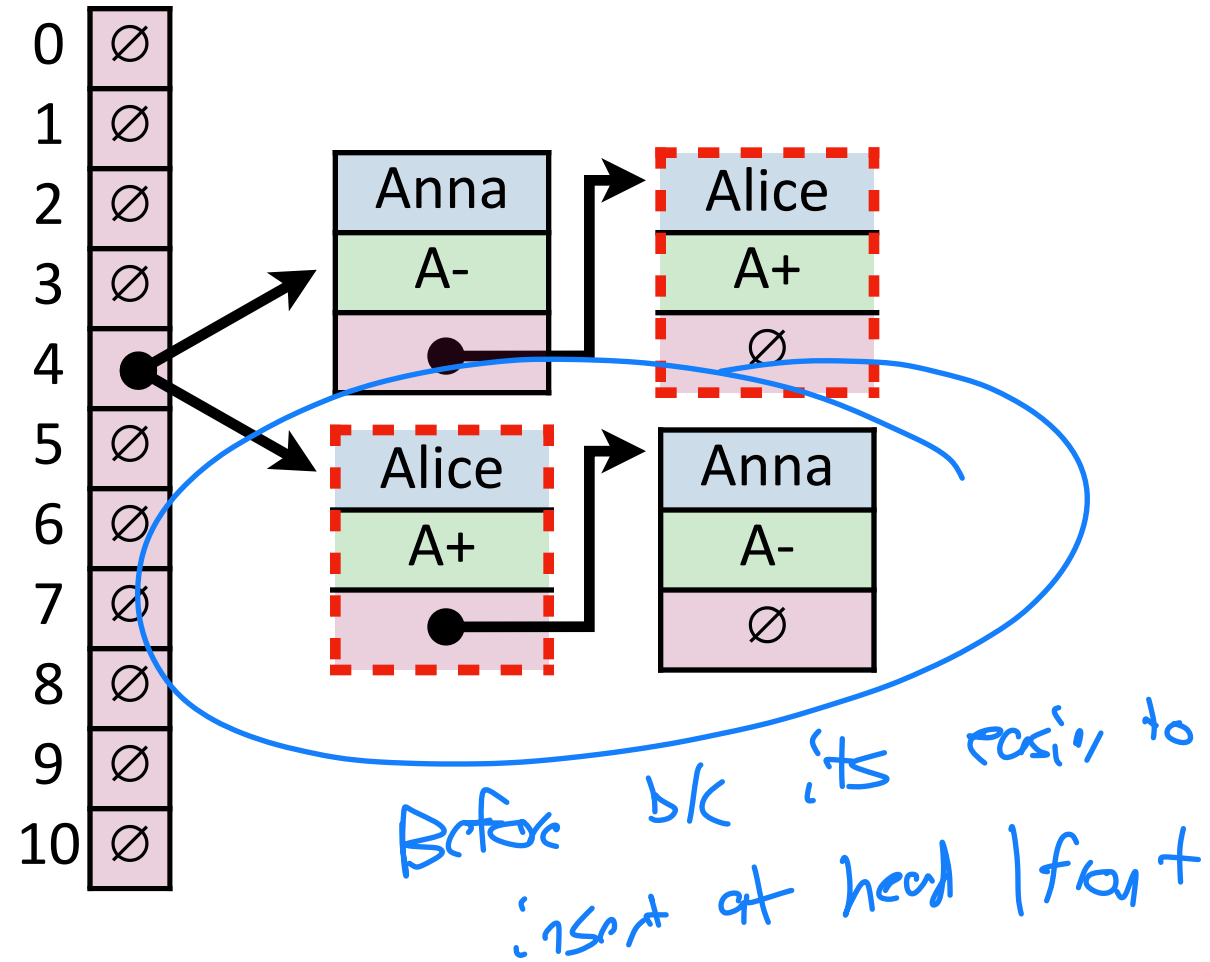
Key	Value	Hash
Bob	B+	2
Anna	A-	4
<b>Alice</b>	<b>A+</b>	<b>4</b>
Betty	B	2
Brett	A-	2
Greg	A	0
Sue	B	7
Ali	B+	4
Laura	A	7
Lily	B+	7



# Insertion (Separate Chaining)

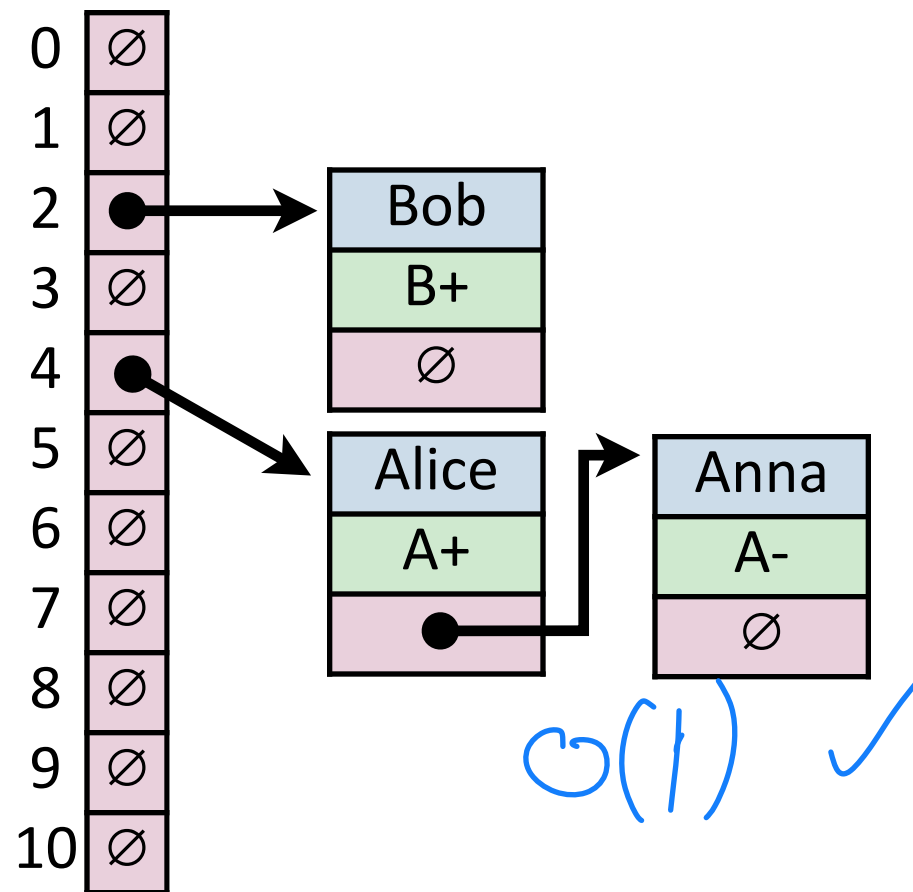
Where does Alice end up relative to Anna in the chain?

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)

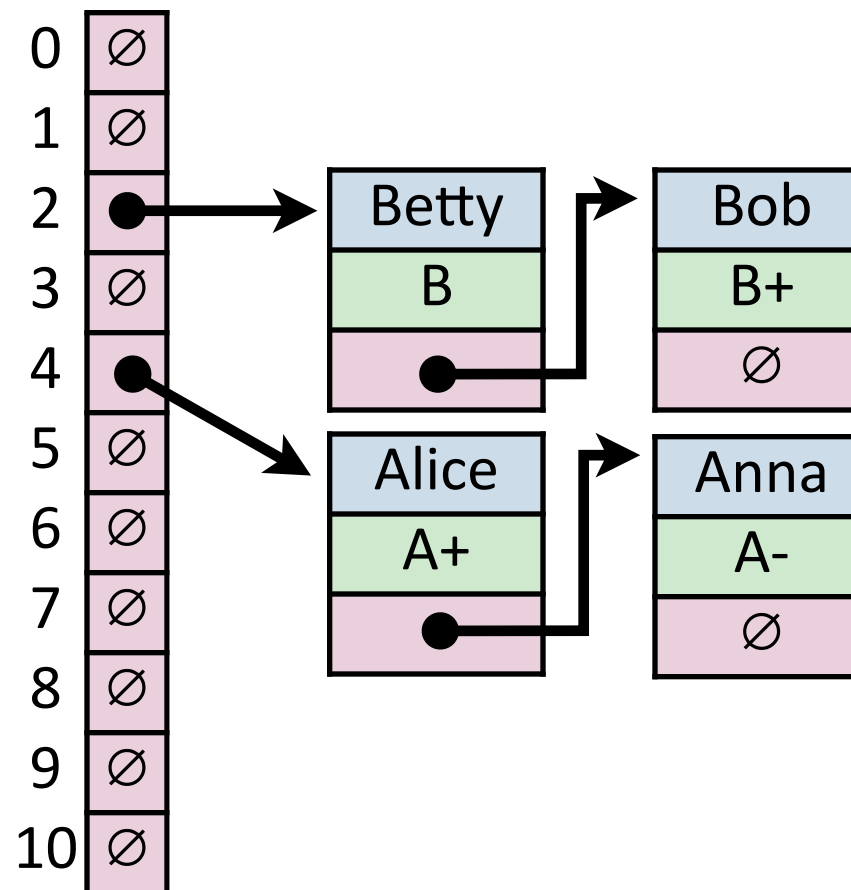
Key	Value	Hash
Bob	B+	2
Anna	A-	4
Alice	A+	4
<b>Betty</b>	<b>B</b>	<b>2</b>
Brett	A-	2
Greg	A	0
Sue	B	7
Ali	B+	4
Laura	A	7
Lily	B+	7





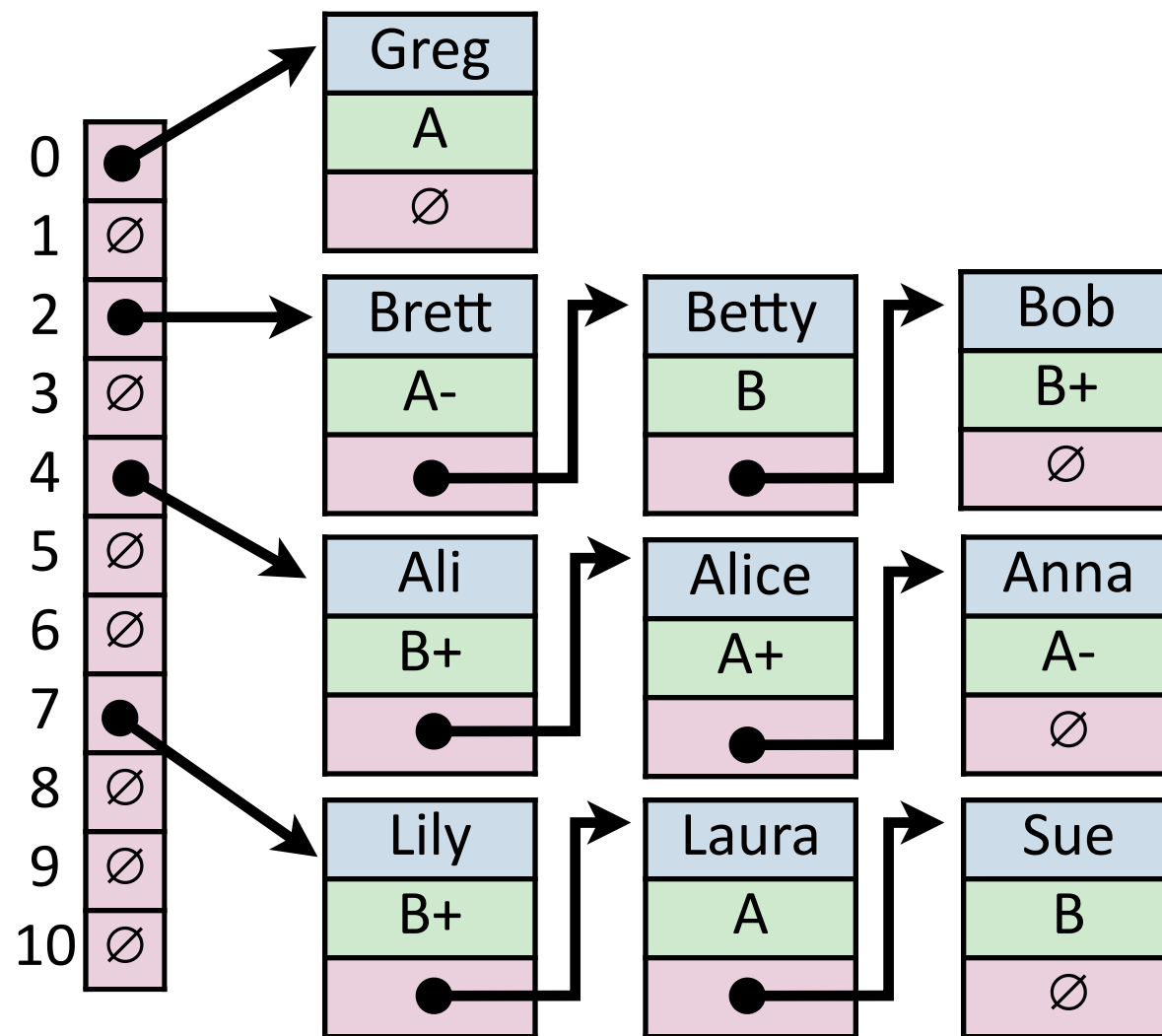
# Insertion (Separate Chaining)

Key	Value	Hash
Bob	B+	2
Anna	A-	4
Alice	A+	4
<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")`

$\hookrightarrow O(n)$

$\frac{1}{n}$

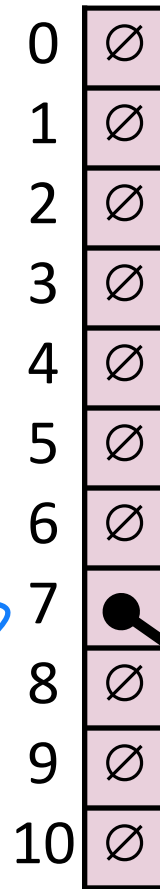
Key	Hash
Sue	7

lookup

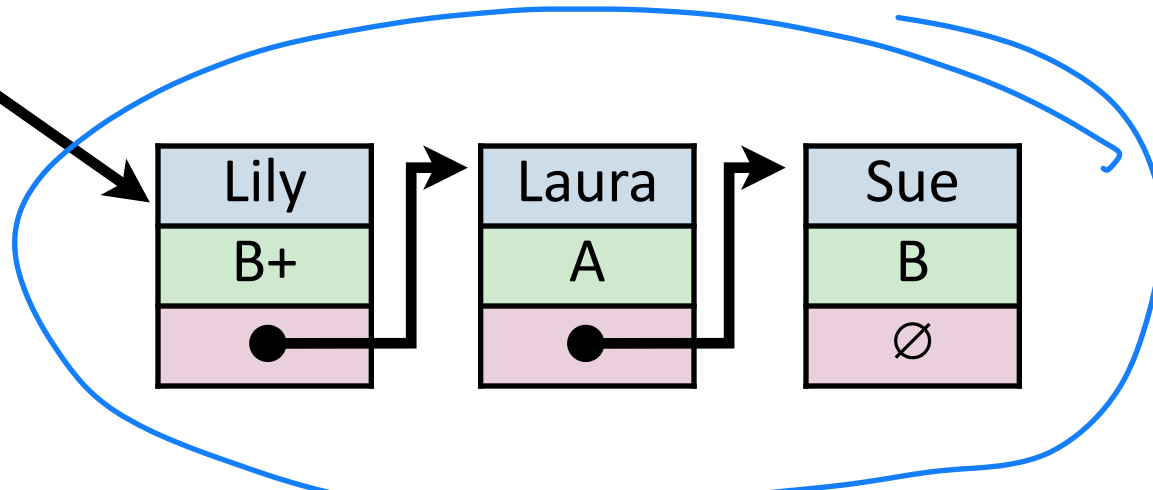
$O(1)$

$O(1)$

Hash



$O(n)$   
 $n$  is # of items



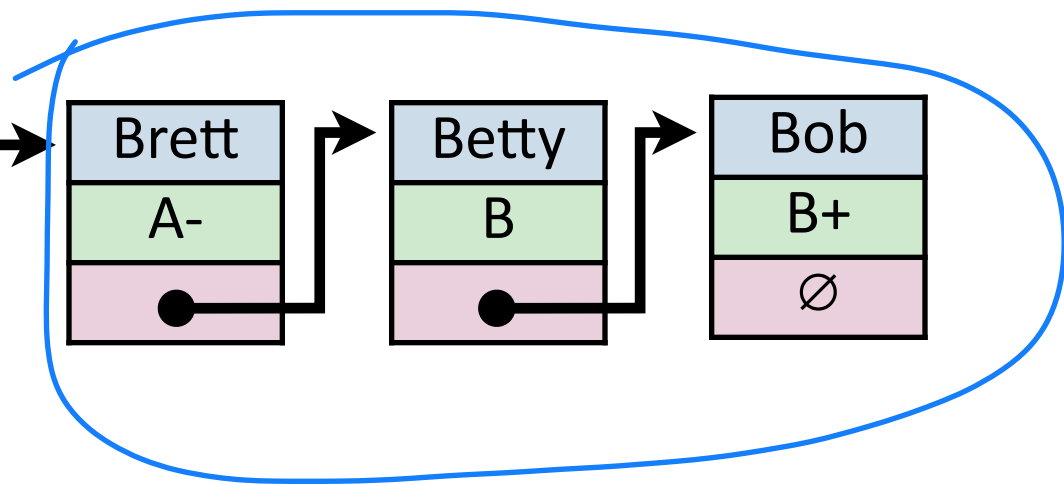
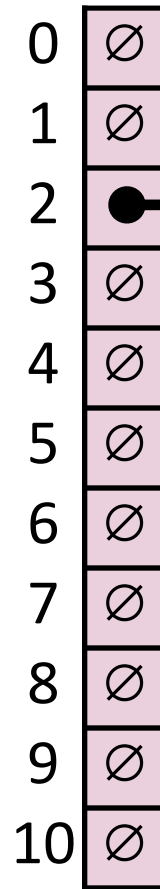
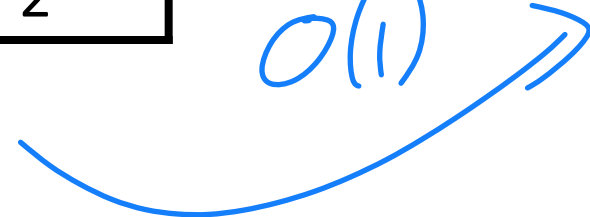
# Remove (Separate Chaining)

`_remove("Betty")`

$\hookrightarrow O(n)$

Key	Hash
Betty	2

$O(1)$



1) Find Betty  $O(n)$

2) Remove Betty  $O(1)$

# Hash Table (Separate Chaining)

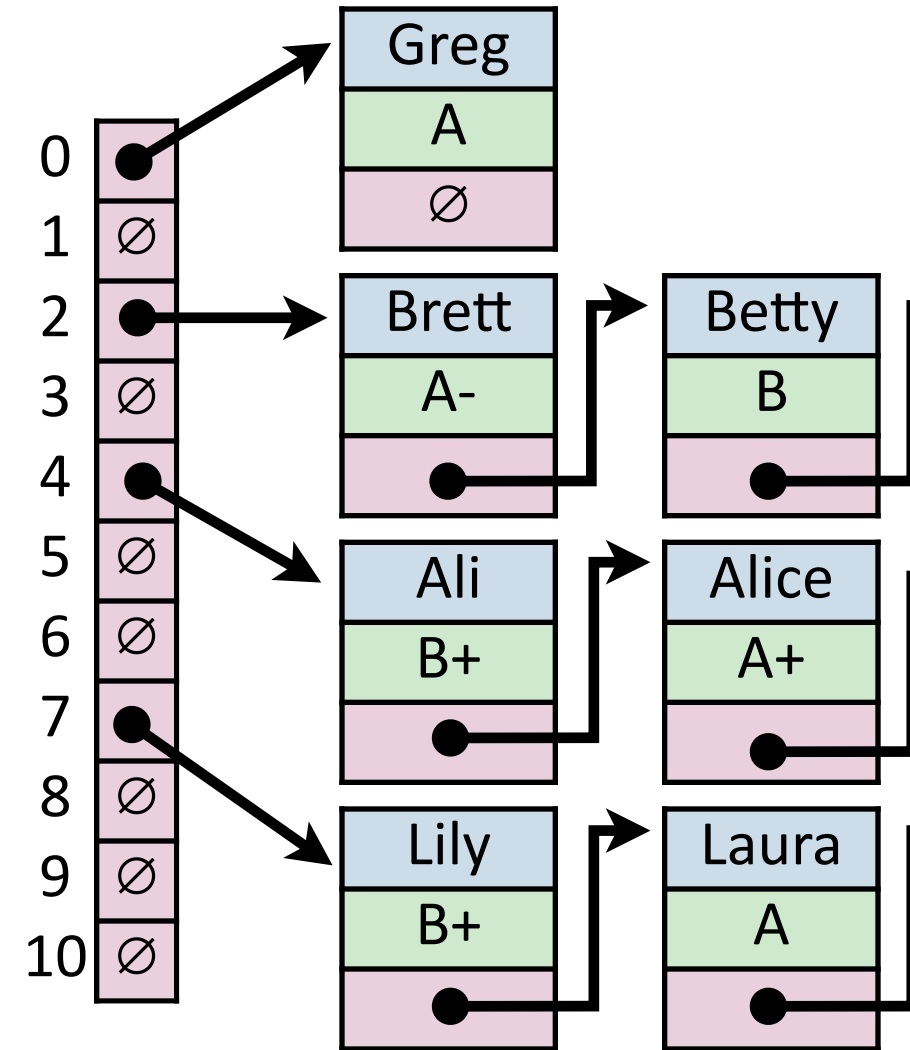


For hash table of size  $m$  and  $n$  elements:

Find runs in:  $O(n)$

Insert runs in:  $O(1)$  ☺

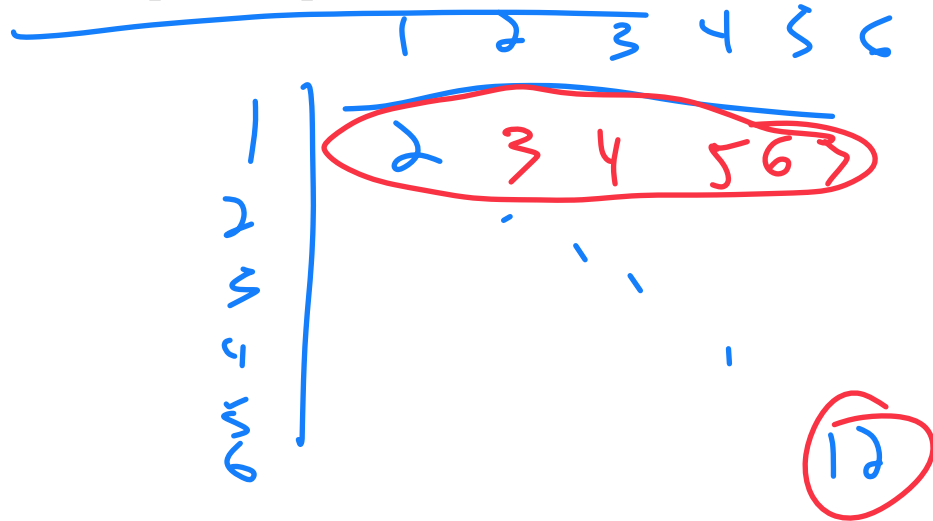
Remove runs in:  $O(n)$



# Fundamentals of Probability

Imagine you roll a pair of six-sided dice.

The **sample space**  $\Omega$  is the set of all possible outcomes.



An **event**  $E \subseteq \Omega$  is any subset.

↳ D1 rolls 1

↳ D1 & D2 roll 12

# Fundamentals of Probability

Imagine you roll a pair of six-sided dice. What is the expected value?

The **expectation** of a (discrete) random variable is:

$$E[X] = \sum_{x \in \Omega} \underbrace{Pr\{X = x\}}_{\text{probability}} \cdot \underbrace{x}_{\text{outcome}}$$

Average expected outcome

$$\frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 2 + \frac{1}{6} \cdot 3 + \dots$$

$$= 3.5$$

# Fundamentals of Probability

Imagine you roll a pair of six-sided dice. What is the expected value?

**Linearity of Expectation:** For any two random variables  $X$  and  $Y$ ,

$$E[X + Y] = E[X] + E[Y]$$



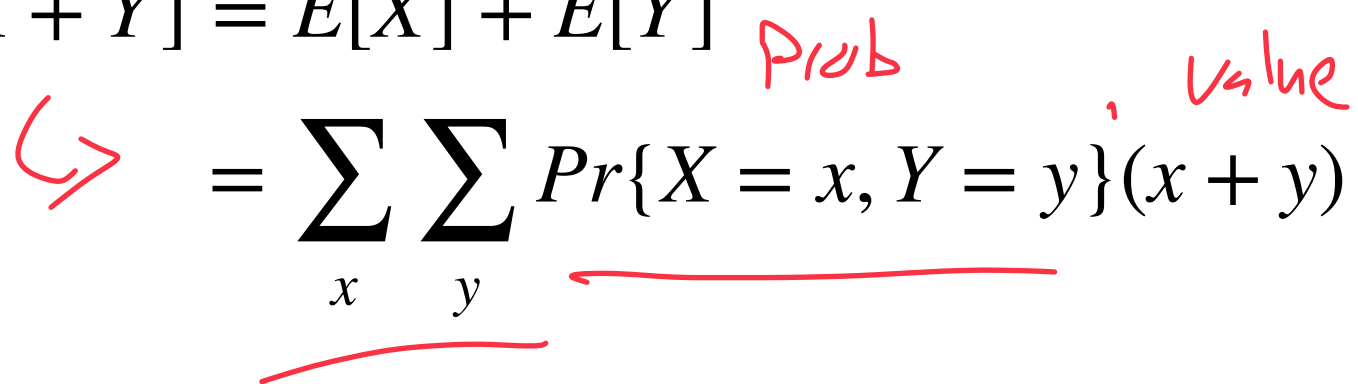
# Fundamentals of Probability

Imagine you roll a pair of six-sided dice. What is the expected value?

**Linearity of Expectation:** For any two random variables  $X$  and  $Y$ ,

$$E[X + Y] = E[X] + E[Y]$$

*Prob* , *Value*

$$\hookrightarrow = \sum_x \sum_y \Pr\{X = x, Y = y\} (x + y)$$


# Fundamentals of Probability

Imagine you roll a pair of six-sided dice. What is the expected value?

**Linearity of Expectation:** For any two random variables  $X$  and  $Y$ ,

$$E[X + Y] = E[X] + E[Y]$$

$$= \sum_x \sum_y \Pr\{X = x, Y = y\}(x + y)$$

$$= \sum_x x \sum_y \Pr\{X = x, Y = y\} + \sum_y y \sum_x \Pr\{X = x, Y = y\}$$

sums up to 1

# Fundamentals of Probability

Imagine you roll a pair of six-sided dice. What is the expected value?

**Linearity of Expectation:** For any two random variables  $X$  and  $Y$ ,

$$E[X + Y] = E[X] + E[Y]$$

$$= \sum_x \sum_y \Pr\{X = x, Y = y\}(x + y)$$

$$= \sum_x x \sum_y \Pr\{X = x, Y = y\} + \sum_y y \sum_x \Pr\{X = x, Y = y\}$$

$$= \sum_x x \cdot \Pr\{X = x\} + \sum_y y \cdot \Pr\{Y = y\}$$

# Fundamentals of Probability



Imagine you roll a pair of six-sided dice. What is the expected value?

**Linearity of Expectation:** For any two random variables  $X$  and  $Y$ ,

$$E[X + Y] = E[X] + E[Y]$$

$$\begin{aligned} & 3.5 + 3.5 \\ & = 7 \end{aligned}$$

# Hash Table

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

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

↳ Simple uniform hash assumption

↳ every  $o$

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

# Simple Uniform Hashing Assumption (SUHA)

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

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

Same hash value

$m$  is table size

**Uniform:** Is a flat line

↳ everything is equally likely to hash to every position

---

**Independent:**

↳ Every item hashes independently of every other item

---

# 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  $j$

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

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

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

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

# 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  $j$

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

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

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

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

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

# 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  $j$

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

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

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

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

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

# 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:  $O(1+2)$ .

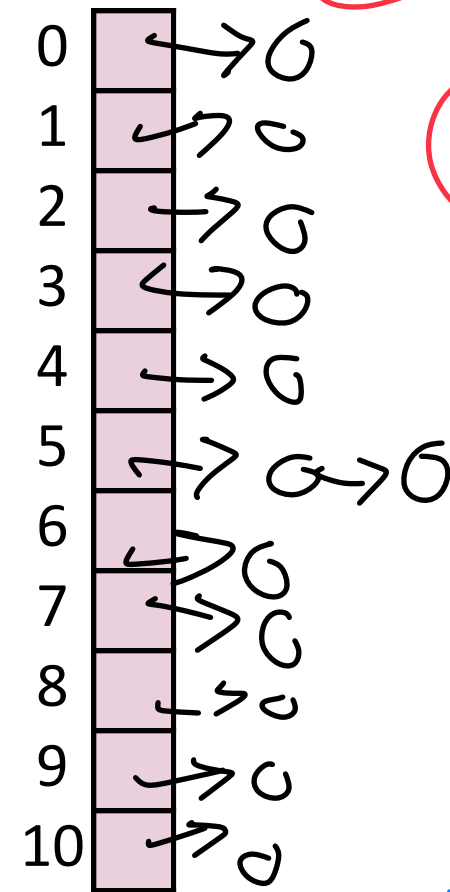
$$\alpha = 1/m$$

Insert runs in:  $O(1)$ .

$$O(1)^*$$

Remove runs in:  $O(1 + \alpha)$ .

\* - we  
(control  
constant  
factors



Expected  $H$   
of items  
is

$$1/m$$

is a  
constant

$$x = 1/m$$

$$m = 1/x$$

You give me  $n$  & I set  $m$  to be good

# Separate Chaining Under SUHA



**Pros:**

**Cons:**

# Next time: Closed Hashing

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

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

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

