# Data Structures and Algorithms
# Hashing

*the best data structure!*

CS 225
Brad Solomon

November 12, 2025

*when in doubt, use a hash! :)*

UNIVERSITY OF
ILLINOIS
URBANA-CHAMPAIGN

Department of Computer Science

# Randomization in Algorithms → Brad's take

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

  ↳ B ST is more likely  (weak)
      log n height

2. Use randomness inside algorithm to estimate expected running time

  ↳ Randomized quicksort    b/100% correct
      ↳ selection of pivot ⟵ ⟵ randomness is on runtime

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

  ↳ Fermat's primality test    O(1)
      ↳ selection of some constant
         ↑ randomness is correctness

# Learning Objectives

Motivate and formally define a hash table

Discuss what a 'good' hash function looks like
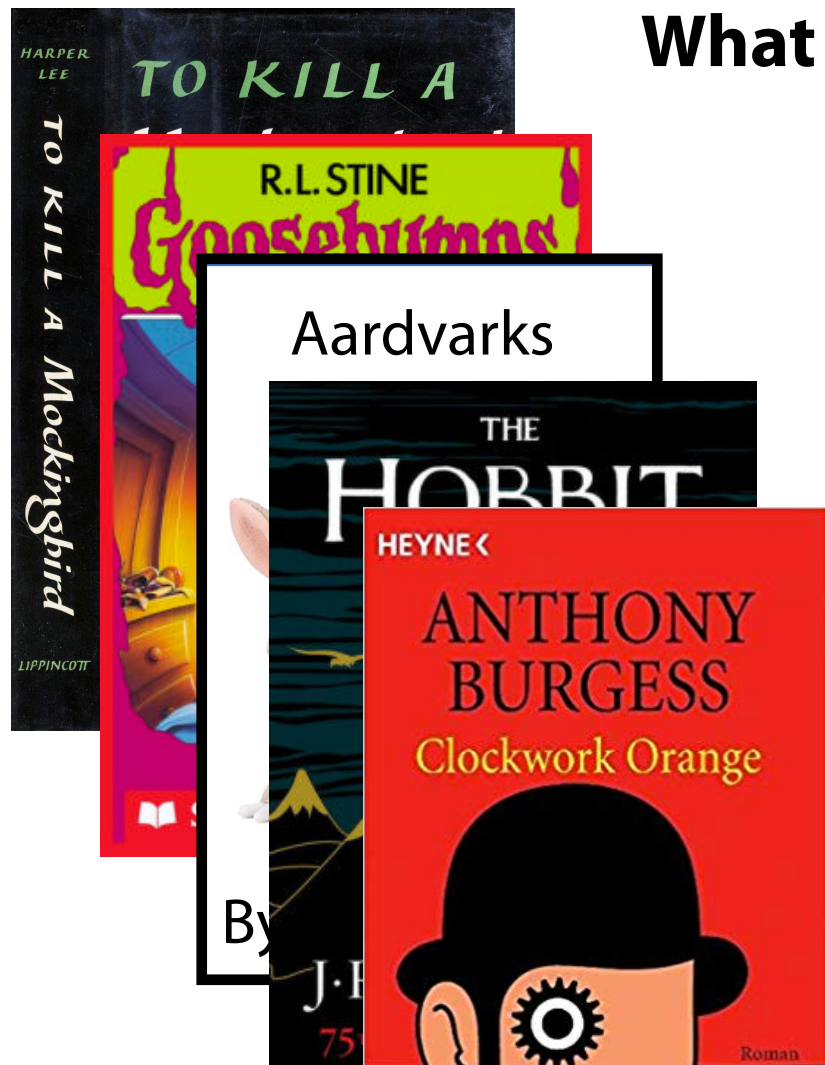
Identify the key weakness of a hash table

Introduce strategies to "correct" this weakness

# Data Structure Review

Keys: Title

Value: Contents

I have a collection of books and I want to store them in a dictionary!
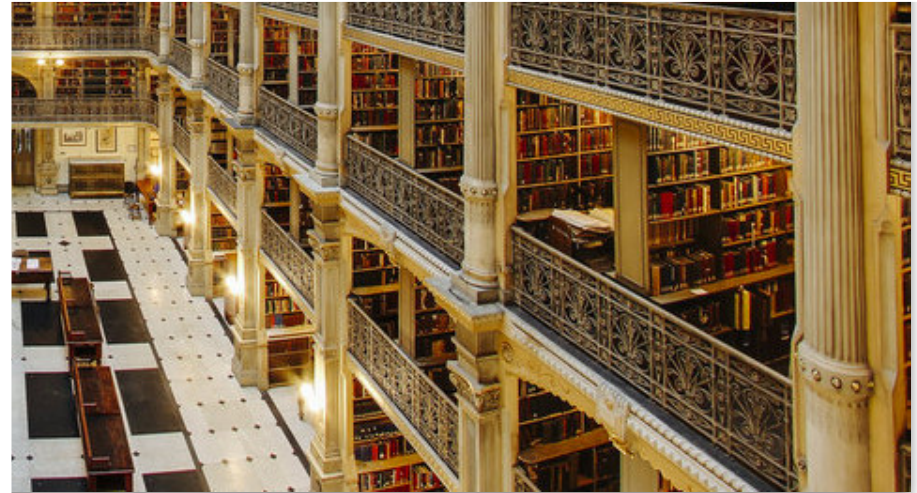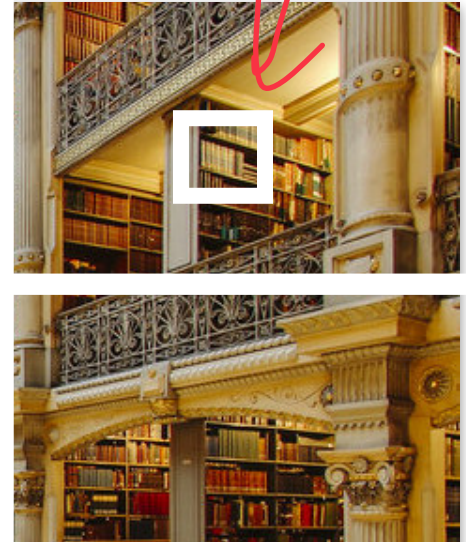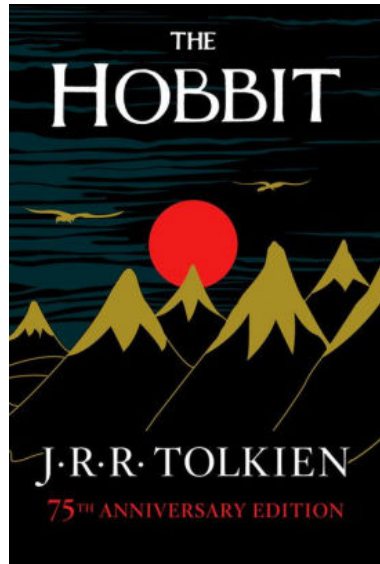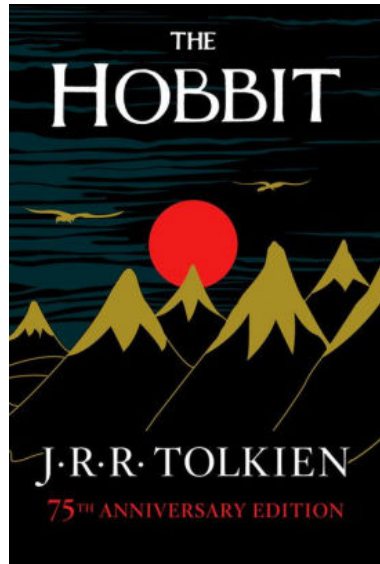
**What data structures can I use here?**

Aardvarks

By

|  | BST | AVL Tree | Min-heap |
|---|---|---|---|
| Insert | $O(h) = O(n)$ | $O(h) = O(\log n)$ | $O(1) / O(n)$ |
| Remove | $O(h) = O(n)$ | $O(\log n)$ |  |
| Find | $O(h) = O(n)$ | $O(\log n)$ |  |

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

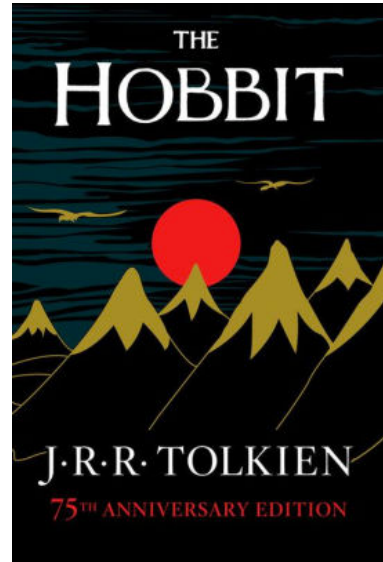# What if $O(log\ n)$ isn't good enough?



$O(1)$ address lookup

$O(1)$ lookup address

# A Hash Table based Dictionary

Some arb input

(Key)

A number / an address

ISBN: 9780062265722

Call #: PR
6068.O93
H35 1937

look up

(Value)

Address

ISBN: 9780062265722

Call #: PR
6068.O93
H35 1937

# Randomized Data Structures

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

↳ Tradeoff arg!   Benefit of sorting data is $O(\log n)$ search

Downside $O(\log n)$ insert

Randomized data structures rely on **expected** performance

Randomized data structures 'cheat' tradeoffs!

# 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 ( key input → A number/address )

2. A data storage structure (A List!)

3. ?? Something to handle 'chaos'

# Hash Function

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

Some universe of values

Brad's Books [finite set]

$0 \leq K \leq 255$ [continuous)

All Possible strings in universe [infinite set)

List

*m* **elements**

| Key | Value |
|-----|-------|
|     |       |
|     |       |
|     |       |
|     |       |
|     |       |
|     |       |
|     |       |
|     |       |
|     |       |

# Hash Function

A hash function *must* be:

- **Deterministic**: If you input same key twice, get same output

- **Efficient**: $O(1)$ is goal

- **Defined for a certain size table**: $h(k) \rightarrow [0 \ldots m-1]$

We have list of size $m$

# Hash Function

(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**

$$h(K) = K[0] - 'A'$$

4 get 0th letter
A = 0
B = 1
C = 2
⋮

| Key | Value |
|---|---|
| Angrave | 241 |
| Beckman | 421 |
| | |
| | |
| | |
| | |
| | |
| | |

# Hash Function

If this is my only dataset, this is a perfect hash

(Angrave, CS 241)     Alumni 411

(Beckman, CS 421)                    Collision!

(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

$\leq$ ??
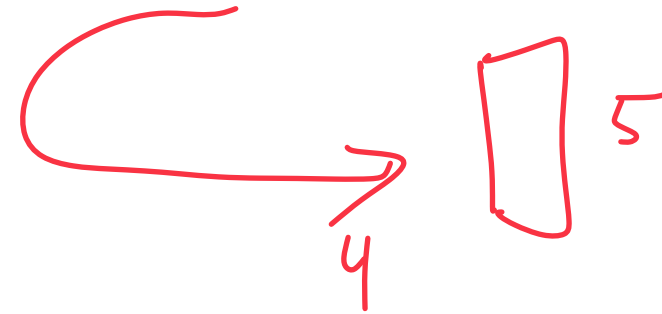
Problem 2: Not mapped to $[0, m-1]$

# General Hash Function

An $O(1)$ deterministic operation that maps all keys in a universe $U$ to a defined range of integers $[0,\ldots,m-1]$
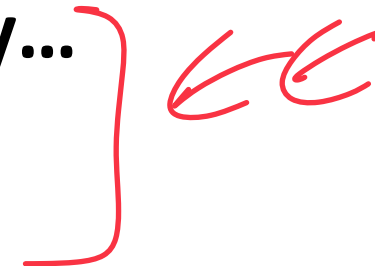
- A **hash**: $Key \rightarrow int$

$$h(Brad) = 999999$$

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

compression

$$\rightarrow \boxed{\phantom{5}}\, 5$$

$4$

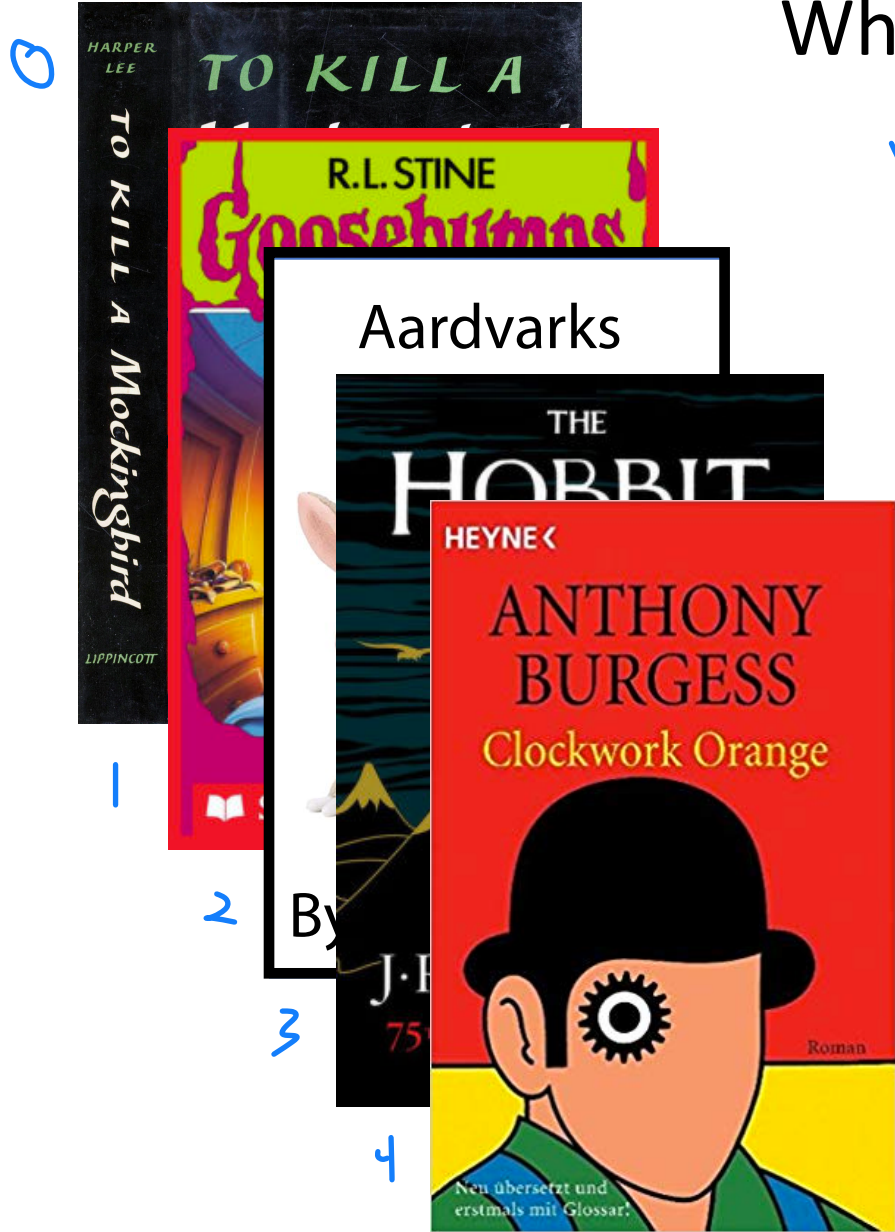**Choosing a good hash function is tricky...**
- Don't create your own (yet*)

# Hash Function

**Which of the following are valid hashes?**

first and last name first character

✓

$$h_1(k) = \left( k.firstName[0] + k.lastName[0] \right) \% m$$

↳ Deterministic    ↳ O(1)    ↳ fixed range

↳ Can have collisions

0

Generate random # multiply by page count

✗

$$h_2(k) = \left( rand() * k.numPages \right) \% m$$

10%

↳ Not deterministic

1

✗

$$h_3(k) = (\text{Order I insert [Order seen]}) \% m$$

30%

↳ Not deterministic

No way of connecting query to an previous address

Insert BooK

2   By

3

4

Stl

Aardvarks

# Hash Function

HT

30

Insert

find
H(goosebumps)

Aardvarks

By

0
1
2
3
4

Goosebumps → 1

0 →
1 →
2 →
3 →
4 →

prev address
Not
deterministic

re-
address

find
H(Goosebumps)
4 5

# A Hash Table based Dictionary

# A Hash Table based Dictionary

Key $\longrightarrow$ Value

# A Hash Table based Dictionary

# A Hash Table based Dictionary

# Hash Collision

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

# Perfect Hashing

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

**S, a finite Keyspace**

## m elements

| Key | Value |
| --- | --- |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |

# General Purpose Hashing

In CS 225, we want our hash functions to work *in general.*

**U, Universe of Keys**

**m elements**

| Key | Value |
|-----|-------|
|     |       |
|     |       |
|     |       |
|     |       |
|     |       |
|     |       |
|     |       |
|     |       |

# General Purpose Hashing

If $m < U$, there must be at least one hash collision.

# 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, U in Some external data structure

- **Closed Hashing:** Store K,u in One fixed Size array

# 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 ___extent the chain_____. This is called **separate chaining.**

# Insertion (Separate Chaining)

`_insert("Bob")`

`_insert("Anna")`

$h(k)$

| 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 ∅ → Bob
3 ∅
4 ∅ → Anna
5 ∅
6 ∅
7 ∅
8 ∅
9 ∅
10 ∅

# Insertion (Separate Chaining) `_insert("Alice")`

**Where does Alice end up in the hash table?**

| 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 |

# 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)

$O(1)$

| Key | Hash |
|-----|------|
| Sue | 7 |

$O(1)$

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

$O(n)$ :(

| Lily | | Laura | | Sue |
|------|--|-------|--|-----|
| B+ | | A | | B |
| ● | | ● | | ∅ |

# Remove (Separate Chaining)  _remove("Betty")

Finding Betty is O(n)

O(1)

| Key | Hash |
|-----|------|
| Betty | 2 |

O(1)

```
0  ∅
1  ∅
2  ●———→  Brett ——→ Betty ——→  Bob
3  ∅         A-        B        B+
4  ∅                             ∅
5  ∅
6  ∅
7  ∅
8  ∅
9  ∅
10 ∅
```

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: $O(n)$ _____

# 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

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

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

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

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

# Separate Chaining Under SUHA

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

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

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

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

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

# Separate Chaining Under SUHA

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

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

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

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

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

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

# Separate Chaining Under SUHA

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

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

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

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

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

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

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

# Separate Chaining Under SUHA

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

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

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

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

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

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

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

# Separate Chaining Under SUHA

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

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

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

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

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

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

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

$$\boxed{\mathbf{E}[\alpha_\mathbf{j}] = \frac{\mathbf{n}}{\mathbf{m}}}$$
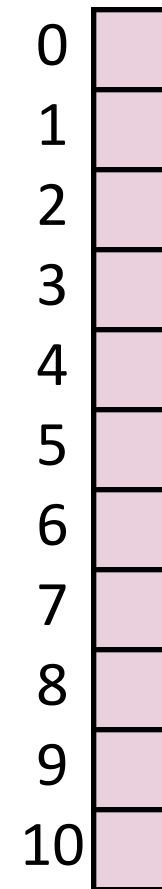
# Separate Chaining Under SUHA

**Under SUHA, a hash table of size _m_ and _n_ elements:**

Find runs in: _____.

Insert runs in: _____.

Remove runs in: _____.

0
1
2
3
4
5
6
7
8
9
10

# Separate Chaining Under SUHA

**Pros:**

**Cons:**

# Next time: Closed Hashing

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

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

$h(k) = k \% 7$

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | |
| 6 | |