# CSE 250
## Data Structures

Dr. Eric Mikida
epmikida@buffalo.edu
208 Capen Hall

# Lec 32: Hash Table Variants

# Announcements

- PA3 testing Autolab now open (testing due Sunday)

# Recap of HashTables (so far...)

**Current Design:** HashTable with Chaining
- Array of buckets
- Each bucket is the head of a linked list (a "chain" of elements)

# Runtime for `contains(x)`

**Expected Runtime:**

# Runtime for `contains(x)`

**Expected Runtime:**

1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$

# Runtime for `contains(x)`

**Expected Runtime:**

1.  Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2.  Find the record in the bucket: $O(\alpha \cdot c_{equality}) = O(1)$

# Runtime for `contains(x)`

**Expected Runtime:**

1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Find the record in the bucket: $O(\alpha \cdot c_{equality}) = O(1)$

**Remember**: we don't let $\alpha$ exceed a constant value

# Runtime for `contains(x)`

**Expected Runtime:**

1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Find the record in the bucket: $O(\alpha \cdot c_{equality}) = O(1)$
3. **Total: $O(c_{hash} + \alpha \cdot c_{equality}) = O(1)$**

# Runtime for `contains(x)`

**Expected Runtime:**

1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Find the record in the bucket: $O(\alpha \cdot c_{equality}) = O(1)$
3. **Total: $O(c_{hash} + \alpha \cdot c_{equality}) = O(1)$**

**Unqualified Worst-Case:**

# Runtime for `contains(x)`

**Expected Runtime:**

1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Find the record in the bucket: $O(\alpha \cdot c_{equality}) = O(1)$
3. **Total:** $O(c_{hash} + \alpha \cdot c_{equality}) = O(1)$

**Unqualified Worst-Case:**

1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$

# Runtime for `contains(x)`

**Expected Runtime:**
1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Find the record in the bucket: $O(\alpha \cdot c_{equality}) = O(1)$
3. **Total:** $O(c_{hash} + \alpha \cdot c_{equality}) = O(1)$

**Unqualified Worst-Case:**
1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Find the record in the bucket: $O(n \cdot c_{equality}) = O(n)$

# Runtime for `contains(x)`

**<u>Expected Runtime:</u>**
1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Find the record in the bucket: $O(\alpha \cdot c_{equality}) = O(1)$
3. **Total:** $O(c_{hash} + \alpha \cdot c_{equality}) = O(1)$

**<u>Unqualified Worst-Case:</u>**
1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Find the record in the bucket: $O(n \cdot c_{equality}) = O(n)$

**Note:** The expected number of equality checks and the worst-case number of equality checks are where these costs differ

# Runtime for `contains(x)`

**Expected Runtime:**
1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Find the record in the bucket: $O(\alpha \cdot c_{equality}) = O(1)$
3. **Total: $O(c_{hash} + \alpha \cdot c_{equality}) = O(1)$**

**Unqualified Worst-Case:**
1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Find the record in the bucket: $O(n \cdot c_{equality}) = O(n)$
3. **Total: $O(c_{hash} + n \cdot c_{equality}) = O(n)$**

# Runtime for `remove(x)`

**Expected Runtime:**
1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Find the record in the bucket: $O(\alpha \cdot c_{equality}) = O(1)$
3. Remove (by reference): $O(1)$
4. **Total: $O(c_{hash} + \alpha \cdot c_{equality} + 1) = O(1)$**

**Unqualified Worst-Case:**
1. Find the record in the bucket: $O(n \cdot c_{equality}) = O(n)$
2. **Total: $O(c_{hash} + n \cdot c_{equality} + 1) = O(n)$**

# Runtime for `remove(x)`

**Expected Runtime:**
1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Find the record in the bucket: $O(\alpha \cdot c_{equality}) = O(1)$
3. Remove (by reference): $O(1)$
4. **Total: $O(c_{hash} + \alpha \cdot c_{equality} + 1) = O(1)$**    Only one extra constant-time step to remove

**Unqualified Worst-Case:**
1. Find the record in the bucket: $O(n \cdot c_{equality}) = O(n)$
2. **Total: $O(c_{hash} + n \cdot c_{equality} + 1) = O(n)$**

# Runtime for `insert(x)`

**Expected Runtime:**

1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Remove $x$ from bucket if present: $O(\alpha \cdot c_{equality} + 1)$
3. Prepend to bucket: $O(1)$
4. Rehash if needed: $O(n \cdot c_{hash} + N)$ **(amortized $O(1)$)**
5. **Total:** $O(c_{hash} + \alpha \cdot c_{equality} + 3) = O(1)$

**Unqualified Worst-Case:**

1. Remove $x$ from bucket if present: $O(n \cdot c_{equality} + 1) = O(n)$
2. **Total:** $O(c_{hash} + n \cdot c_{equality} + 3) = O(n)$

# Runtime for `insert(x)`

**Expected Runtime:**

1. Find the bucket (call our hash function): $O(c_{hash}) = O(1)$
2. Remove *x* from bucket if present: $O(\alpha \cdot c_{equality} + 1)$
3. Prepend to bucket: $O(1)$
4. Rehash if needed: $O(n \cdot c_{hash} + N)$ (amortized $O(1)$)
5. **Total:** $O(c_{hash} + \alpha \cdot c_{equality} + 3) = O(1)$

One additional constant-time step to prepend, and then potentially the need to rehash, but that is amortized $O(1)$

**Unqualified Worst-Case:**

1. Remove *x* from bucket if present: $O(n \cdot c_{equality} + 1) = O(n)$
2. **Total:** $O(c_{hash} + n \cdot c_{equality} + 3) = O(n)$

# HashTable Drawbacks?

…So the expected runtime of all operations is $O(1)$

*Why would you ever use any other data structure?*

# HashTable Drawbacks?

…So the expected runtime of all operations is $O(1)$

*Why would you ever use any other data structure?*

- HashTables do not preserve ordering
- HashTables may waste a lot of memory
- Rehashing can be expensive
- Only **guarantee** on lookup time is that it is $O(n)$

# HashTable Drawbacks?

…So the expected runtime of all operations is $O(1)$

*Why would you ever use any other data structure?*

- HashTables do not preserve ordering
- HashTables may waste a lot of memory
- Rehashing can be expensive
- Only **guarantee** on lookup time is that it is $O(n)$

These can be partially addressed by some HashTable variations

# Collision Resolution

- When two records are assigned to the same bucket, this is called a **collision**
  - With chaining, collisions are resolved by treating each bucket as a list
  - May result in even more empty buckets (more wasted space)
- Two more collision resolution techniques try to help with this issue
  - Open Addressing
  - Cuckoo Hashing

# HashTables with Chaining

hash(A) = 4

hash(B) = 5

hash(C) = 5

hash(D) = 2

hash(E) = 6

hash(F) = 2

# HashTables with Chaining

hash(A) = 4

hash(B) = 5

**hash(C) = 5**

hash(D) = 2

hash(E) = 6

**hash(F) = 2**

| 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
|   |   | D F | | A | B C | E |

Collisions are resolved by adding the element to the buckets linked list

# HashTables with Open Addressing

**hash(A) = 4 ← no collision**

hash(B) = 5

hash(C) = 5

hash(D) = 2

hash(E) = 6

hash(F) = 4

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

With Open Addressing collisions are resolved by "cascading" to the next available bucket

# HashTables with Open Addressing

hash(A) = 4

**hash(B) = 5 ← no collision**

hash(C) = 5

hash(D) = 2

hash(E) = 6

hash(F) = 4

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

With Open Addressing collisions are resolved by "cascading" to the next available bucket

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

**hash(C) = 5 ← collision! Search for next free bucket**

hash(D) = 2

hash(E) = 6

hash(F) = 4

| 0 | 1 | 2 | 3 | A<br>4 | B<br>5 | 6 |
|---|---|---|---|---|---|---|

C ✖

With Open Addressing collisions are resolved by "cascading" to the next available bucket

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

**hash(C) = 5 ← collision! Search for next free bucket**

hash(D) = 2

hash(E) = 6

hash(F) = 4

| 0 | 1 | 2 | 3 | 4 A | 5 B | 6 C |
|---|---|---|---|---|---|---|

With Open Addressing collisions are resolved by "cascading" to the next available bucket

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

hash(C) = 5

**hash(D) = 2 ← no collision!**

hash(E) = 6

hash(F) = 4

| 0 | 1 | 2 D | 3 | 4 A | 5 B | 6 C |
|---|---|---|---|---|---|---|

With Open Addressing collisions are resolved by "cascading" to the next available bucket

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

hash(C) = 5

hash(D) = 2

**hash(E) = 6  ← collision! cascade to 0**

hash(F) = 4

| E | | D | | A | B | C |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 |

With Open Addressing collisions are resolved by "cascading" to the next available bucket

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

hash(C) = 5

hash(D) = 2

hash(E) = 6

**hash(F) = 4  ← collision! Cascade all the way to 1**

| E | F | D | | A | B | C |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 |

With Open Addressing collisions are resolved by "cascading" to the next available bucket

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

hash(C) = 5

hash(D) = 2

hash(E) = 6

**hash(F) = 4  ← collision! Cascade all the way to 1**



| | | | | | | |
|---|---|---|---|---|---|---|
| E | F | D | | A | B | C |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 |

With Open Addressing collisions are resolved by "cascading" to the next available bucket

**How does lookup work?**

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

hash(C) = 5

hash(D) = 2

hash(E) = 6

hash(F) = 4

`contains(F)`

| E | 1 | D | 3 | A | B | C |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 |

Bucket 4 does not contain F. Are we sure F does not exist?

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

hash(C) = 5

hash(D) = 2

hash(E) = 6

hash(F) = 4

`contains(F)`

| 0 E | 1 | 2 D | 3 | 4 A | 5 B | 6 C |
|---|---|---|---|---|---|---|

Bucket 4 does not contain F. Are we sure F does not exist? **No...it could have cascaded!**

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

hash(C) = 5

hash(D) = 2

hash(E) = 6

hash(F) = 4

`contains(F)`

| 0 E | 1 | 2 D | 3 | 4 A | 5 B | 6 C |
|---|---|---|---|---|---|---|

Bucket 5 does not contain F. Are we sure F
does not exist? **No...it could have cascaded!**

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

hash(C) = 5

hash(D) = 2

hash(E) = 6

hash(F) = 4

`contains(F)`

| E | | D | | A | B | C |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 |

Bucket 6 does not contain F. Are we sure F does not exist? **No...it could have cascaded!**

35

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

hash(C) = 5

hash(D) = 2

hash(E) = 6

hash(F) = 4

`contains(F)`

| 0 E | 1 | 2 D | 3 | 4 A | 5 B | 6 C |
|-----|---|-----|---|-----|-----|-----|

Bucket 0 does not contain F. Are we sure F does not exist? **No…it could have cascaded!**

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

hash(C) = 5

hash(D) = 2

hash(E) = 6

hash(F) = 4

`contains(F)`

| 0 E | 1 | 2 D | 3 | 4 A | 5 B | 6 C |
|-----|---|-----|---|-----|-----|-----|

Bucket 1 does not contain F. Are we sure F does not exist? **Yes! If F existed it would be here, so `apply(F)` returns False.**

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

hash(C) = 5

hash(D) = 2

hash(E) = 6

hash(F) = 4

`contains(F)`

| 0 E | 1 | 2 D | 3 | 4 A | 5 B | 6 C |
|---|---|---|---|---|---|---|

Bucket 1 does not contain F. Are we sure F does not exist? **Yes! If F existed it would be here, so `apply(F)` returns False.**

**What if we insert F then remove E?**

# HashTables with Open Addressing

hash(A) = 4

hash(B) = 5

hash(C) = 5

hash(D) = 2

hash(E) = 6

hash(F) = 4

| 0 | 1 F | 2 D | 3 | 4 A | 5 B | 6 C |
|---|---|---|---|---|---|---|

`contains(F)` would fail in this case because it would check bucket 0 and conclude F doesn't exist!

Remove must also deal with potential cascading!

**What if we insert F then remove E?**

# Removals with Open Addressing

**To remove elements with Open Addressing:**

1. First find the element (if it exists)
2. Remove the element
   a. Check all following elements in a contiguous block and move them up
   b. Don't move any element Y to a position that comes before hash(Y)

# Open Addressing Runtime

**Cascading to the next bucket(s) is called probing**
- **Linear Probing:** If collision, cascade to hash(X) + $ci$
- **Quadratic Probing:** If collision, cascade to hash(X) + $ci^2$

**Runtime Costs:**
- Chaining is dominated by searching the chain
- Open Addressing is dominated by probing
  - In both cases, with low $\alpha$ we expect operations to be ***O(1)***
  - Open addressing will occupy more buckets (waste less space)

# Cuckoo Hashing

**Open Addressing can have arbitrarily long chains**

*Can we reduce the chance of cascading for some operations?*

# Cuckoo Hashing

**Idea:** Use two hash functions, hash$_1$ and hash$_2$

To insert a record **X**:

1. If hash$_1$(**X**) and hash$_2$(**X**) are both available, pick one at arbitrarily
2. If only one of those buckets is available, pick the available bucket
3. If neither is available, pick one arbitrarily and evict the record there
   a. Insert **X** in this bucket
   b. Insert the evicted record following the same procedure

# HashTables with Cuckoo Hashing

**hash$_1$(A) = 1**    hash$_2$(A) = 3

hash$_1$(B) = 2    hash$_2$(B) = 4

hash$_1$(C) = 2    hash$_2$(C) = 1

hash$_1$(D) = 4    hash$_2$(D) = 6

hash$_1$(E) = 3    hash$_2$(E) = 4

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

# HashTables with Cuckoo Hashing

$hash_1(A) = 1$      $hash_2(A) = 3$

$\mathbf{hash_1(B) = 2}$      $hash_2(B) = 4$

$hash_1(C) = 2$      $hash_2(C) = 1$

$hash_1(D) = 4$      $hash_2(D) = 6$

$hash_1(E) = 3$      $hash_2(E) = 4$

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

# HashTables with Cuckoo Hashing

$hash_1(A) = 1$     $hash_2(A) = 3$

$hash_1(B) = 2$     $hash_2(B) = 4$

**$hash_1(C) = 2$**     **$hash_2(C) = 1$**

$hash_1(D) = 4$     $hash_2(D) = 6$

$hash_1(E) = 3$     $hash_2(E) = 4$

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

C

*C* can't go in either bucket, so evict one at random (let's say *B*) and reinsert the evicted element

# HashTables with Cuckoo Hashing

$hash_1(A) = 1$  $hash_2(A) = 3$

**$hash_1(B) = 2$**  **$hash_2(B) = 4$**

$hash_1(C) = 2$  $hash_2(C) = 1$

$hash_1(D) = 4$  $hash_2(D) = 6$

$hash_1(E) = 3$  $hash_2(E) = 4$

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

B

*B* can only go in 4 now, but 4 is free

# HashTables with Cuckoo Hashing

$hash_1(A) = 1$     $hash_2(A) = 3$

$hash_1(B) = 2$     $hash_2(B) = 4$

$hash_1(C) = 2$     $hash_2(C) = 1$

$hash_1(D) = 4$     $hash_2(D) = 6$

$hash_1(E) = 3$     $hash_2(E) = 4$

| 0 | 1 A | 2 C | 3 | 4 B | 5 | 6 |
|---|---|---|---|---|---|---|

*B* can only go in 4 now, but 4 is free

# HashTables with Cuckoo Hashing

$hash_1(A) = 1$    $hash_2(A) = 3$

$hash_1(B) = 2$    $hash_2(B) = 4$

$hash_1(C) = 2$    $hash_2(C) = 1$

**$hash_1(D) = 4$**    **$hash_2(D) = 6$**

$hash_1(E) = 3$    $hash_2(E) = 4$

| 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
|   | A | C |   | B |   | D |

# HashTables with Cuckoo Hashing

$hash_1(A) = 1$       $hash_2(A) = 3$

$hash_1(B) = 2$       $hash_2(B) = 4$

$hash_1(C) = 2$       $hash_2(C) = 1$

$hash_1(D) = 4$       $hash_2(D) = 6$

$\mathbf{hash_1(E) = 3}$       $hash_2(E) = 4$

| 0 | 1 A | 2 C | 3 E | 4 B | 5 | 6 D |
|---|-----|-----|-----|-----|---|-----|

# HashTables with Cuckoo Hashing

$hash_1(A) = 1$     $hash_2(A) = 3$

$hash_1(B) = 2$     $hash_2(B) = 4$

$hash_1(C) = 2$     $hash_2(C) = 1$

$hash_1(D) = 4$     $hash_2(D) = 6$

**$hash_1(E) = 3$**     $hash_2(E) = 4$

| 0 | 1 A | 2 C | 3 E | 4 B | 5 | 6 D |
|---|-----|-----|-----|-----|---|-----|

What if we try to insert *F* which hashes to either 1 or 3?

# HashTables with Cuckoo Hashing

hash$_1$(A) = 1     hash$_2$(A) = 3

hash$_1$(B) = 2     hash$_2$(B) = 4

hash$_1$(C) = 2     hash$_2$(C) = 1

hash$_1$(D) = 4     hash$_2$(D) = 6

**hash$_1$(E) = 3**     hash$_2$(E) = 4

| 0 | A 1 | C 2 | E 3 | B 4 | 5 | D 6 |
|---|---|---|---|---|---|---|

What if we try to insert **F** which hashes to either 1 or 3? **We will loop infinitely trying to evict…so limit the number of eviction attempts then do a full rehash**

# Cuckoo Hashing

So with Cuckoo Hashing, we may have to rehash early, and may follow long chains of evictions inserting, but...

What is the runtime of contains/remove?

# Cuckoo Hashing

**So with Cuckoo Hashing, we may have to rehash early, and may follow long chains of evictions inserting, but...**

**What is the runtime of contains/remove?**

1. Check 2 different buckets: $O(1)$
2. That's it...no chaining, cascading etc...

**Apply and remove are <u>GUARANTEED</u> $O(1)$ with Cuckoo Hashing**

# HashTables as Sets

We've now seen HashTable's as an implementation of Sets

- **HashSet** in Java -> Expected O(1) runtime for **add**, **contains**, **remove**

What about **HashMap**? What is a map??

# HashTables as Sets

We've now seen HashTable's as an implementation of Sets

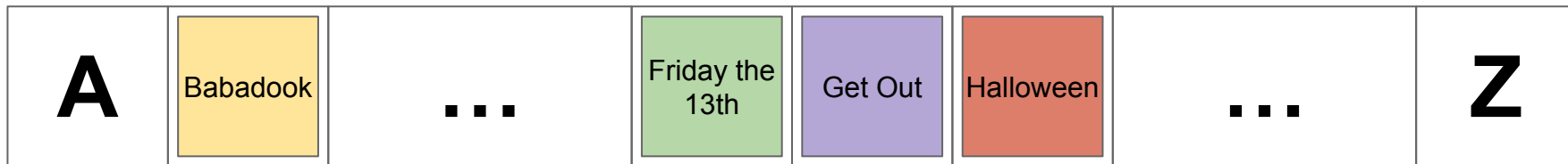- **HashSet** in Java -> Expected O(1) runtime for **add**, **contains**, **remove**

What about **HashMap**? What is a map??

- A map IS as set. It is a set of key-value pairs!

# HashSets vs HashMaps

This was an example of a `HashSet` that stored movie titles (with a bad hash function...but ignore that for now)

| A | Babadook | ... | Friday the 13th | Get Out | Halloween | ... | Z |
|---|---|---|---|---|---|---|---|

# HashSets vs HashMaps

This is an example of a `HashMap` that stores key value pairs where the key is a movie title and the value is the movie object associated with that title



| A | (Babadook,) | ... | (Friday the 13th,) | (Get Out,) | (Halloween,) | ... | Z |

**name:** "Babadook"
**runtime:** 92
**year:** 2014

**name:** "Friday the 13th"
**runtime:** 95
**year:** 1980

**name:** "Get Out"
**runtime:** 103
**year:** 2017

**name:** "Halloween"
**runtime:** 91
**year:** 1978