

CS2040S

Data Structures and Algorithms

Hashing III

Hashing overview

- What is a hash function?
- Collision resolution: chaining
- Java hashing
- Collision resolution: open addressing
- Table (re)sizing

Hashing overview

- What is a hash function?
- Collision resolution: chaining
- Java hashing
- Collision resolution: open addressing
- Table (re)sizing

Abstract Data Types

Symbol Table

public interface SymbolTable

void	insert(Key k, Value v)	<i>insert (k,v) into table</i>
Value	search(Key k)	<i>get value paired with k</i>
void	delete(Key k)	<i>remove key k (and value)</i>
boolean	contains(Key k)	<i>is there a value for k?</i>
int	size()	<i>number of (k,v) pairs</i>

Note: no successor / predecessor queries.

Direct Access Tables

Attempt #1: Use a table, indexed by keys.

0	null
1	null
2	item1
3	null
4	null
5	item3
6	null
7	null
8	item2
9	null

Universe $U = \{0..9\}$ of size $m = 10$.

(key, value)

(2, item1)

(8, item2)

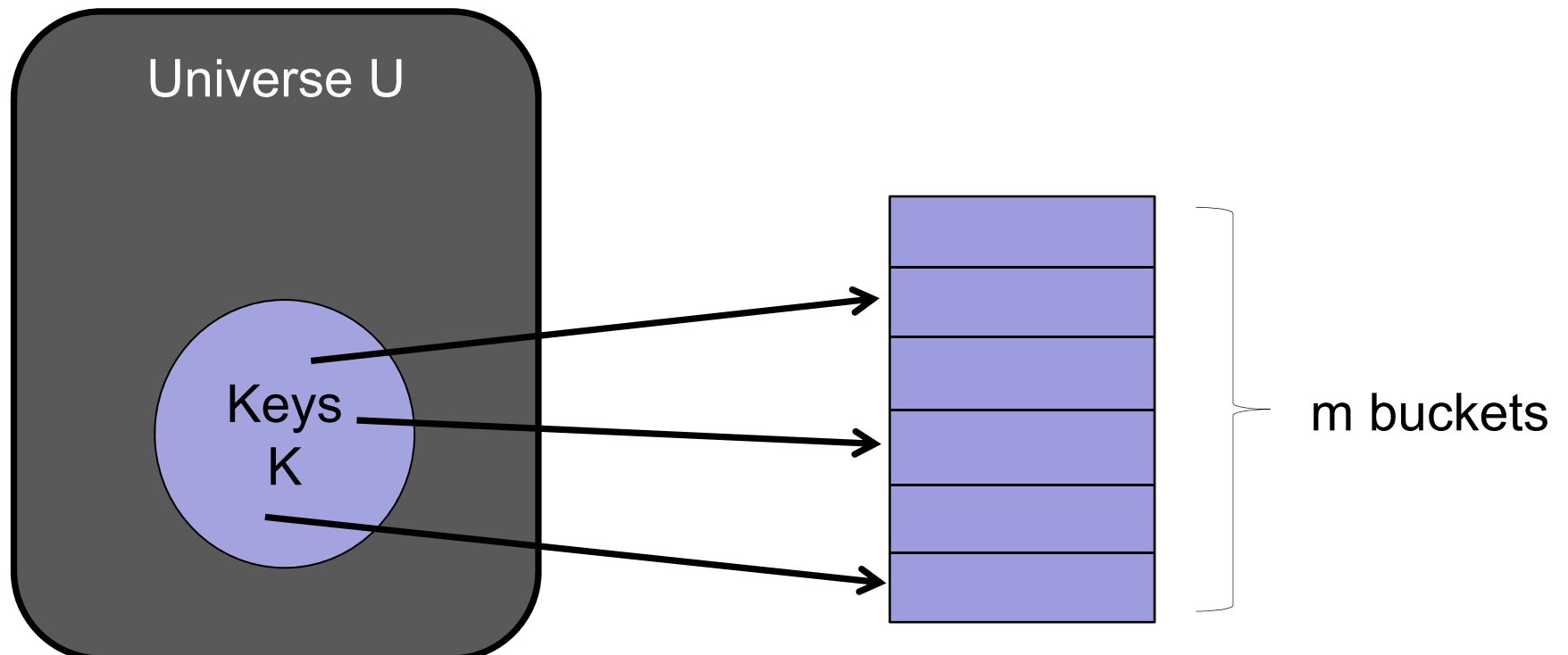
(5, item3)

Assume keys are distinct.

Hash Functions

Problem:

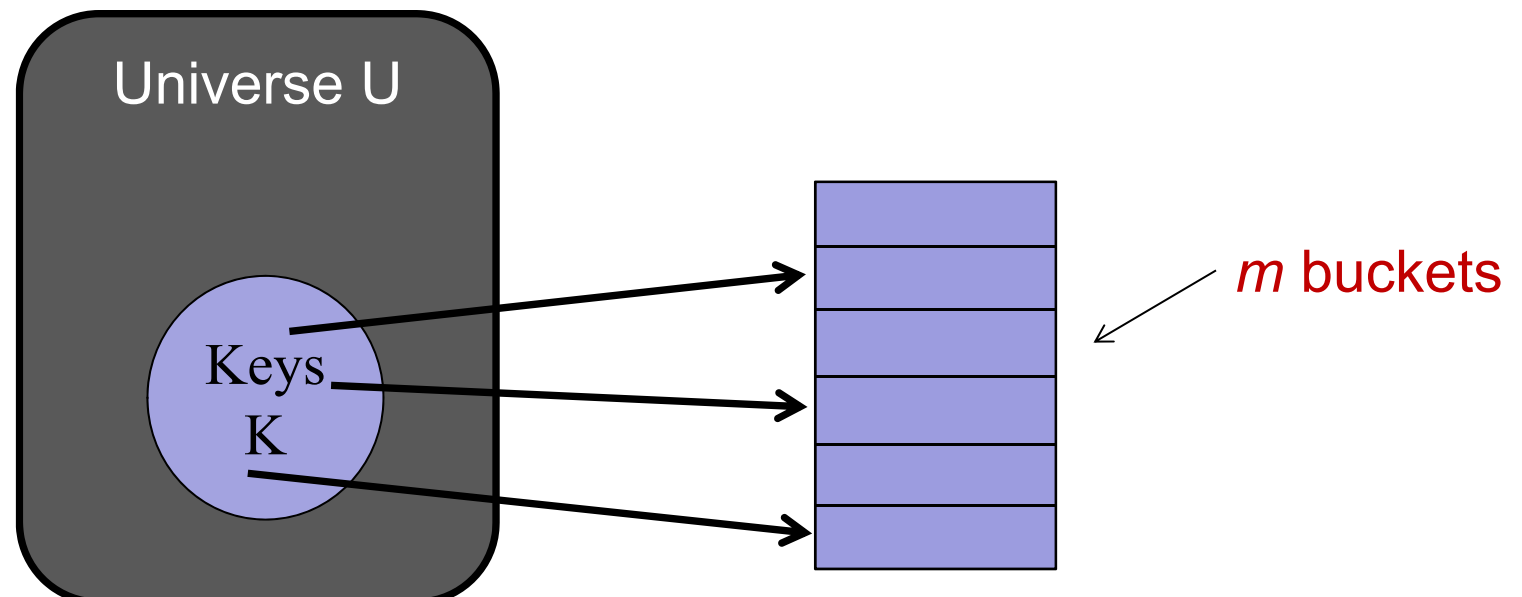
- Huge universe U of possible keys.
- Smaller number n of actual keys.
- How to map n keys to $m \approx n$ buckets?



Hash Functions

Define hash function $h : U \rightarrow \{1..m\}$

- Store key k in bucket $h(k)$.



Hash Functions

Collisions:

- We say that two distinct keys k_1 and k_2 **collide** if:

$$h(k_1) = h(k_2)$$

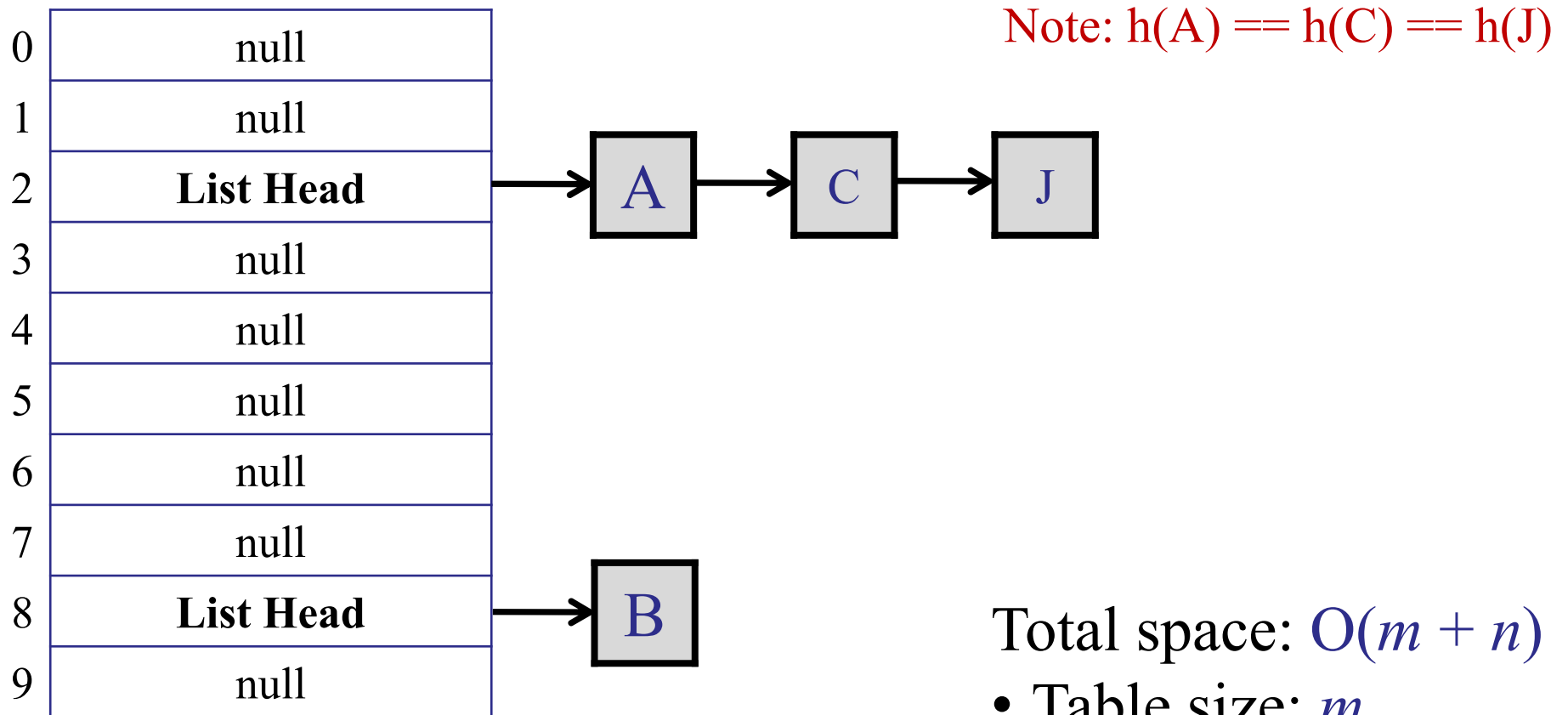
- Unavoidable!
 - The table size is smaller than the universe size.
 - The pigeonhole principle says:
 - There must exist two keys that map to the same bucket.
 - Some keys must collide!

Resolving Collisions

- Basic problem:
 - What to do when two items hash to the same bucket?
- Solution 1: Chaining
 - Insert item into a linked list.
- Solution 2: Open Addressing
 - Find another free bucket.

Chaining

Each bucket contains a linked list of items.



Let's be optimistic today.

The Simple Uniform Hashing Assumption

- Every key is equally likely to map to every bucket.
- Keys are mapped independently.

Intuition:

- Each key is put in a random bucket.
- Then, as long as there are enough buckets, we won't get too many keys in any one bucket.

Let's be optimistic today.

The Simple Uniform Hashing Assumption

- Assume:
 - n items
 - m buckets
- Define: $\text{load}(\text{hash table}) = n/m$
= average # items / buckets.
- Expected search time = $1 + n/m$
 - hash function + array access
 - linked list traversal

Let's be optimistic today.

The Simple Uniform Hashing Assumption

– Assume:

- n items
- $m = \Omega(n)$ buckets, e.g., $m = 2n$

– Expected search time = $1 + n/m$
= $O(1)$

Resolving Collisions

- Basic problem:
 - What to do when two items hash to the same bucket?
- Solution 1: Chaining
 - Insert item into a linked list.
- Solution 2: Open Addressing
 - Find another free bucket.

Open Addressing

Advantages:

- No linked lists!
- All data directly stored in the table.
- One item per slot.

0	null
1	null
2	A
3	null
4	null
5	null
6	null
7	null
8	B
9	null

Open Addressing

On collision:

Probe a sequence of buckets until you find an empty one.

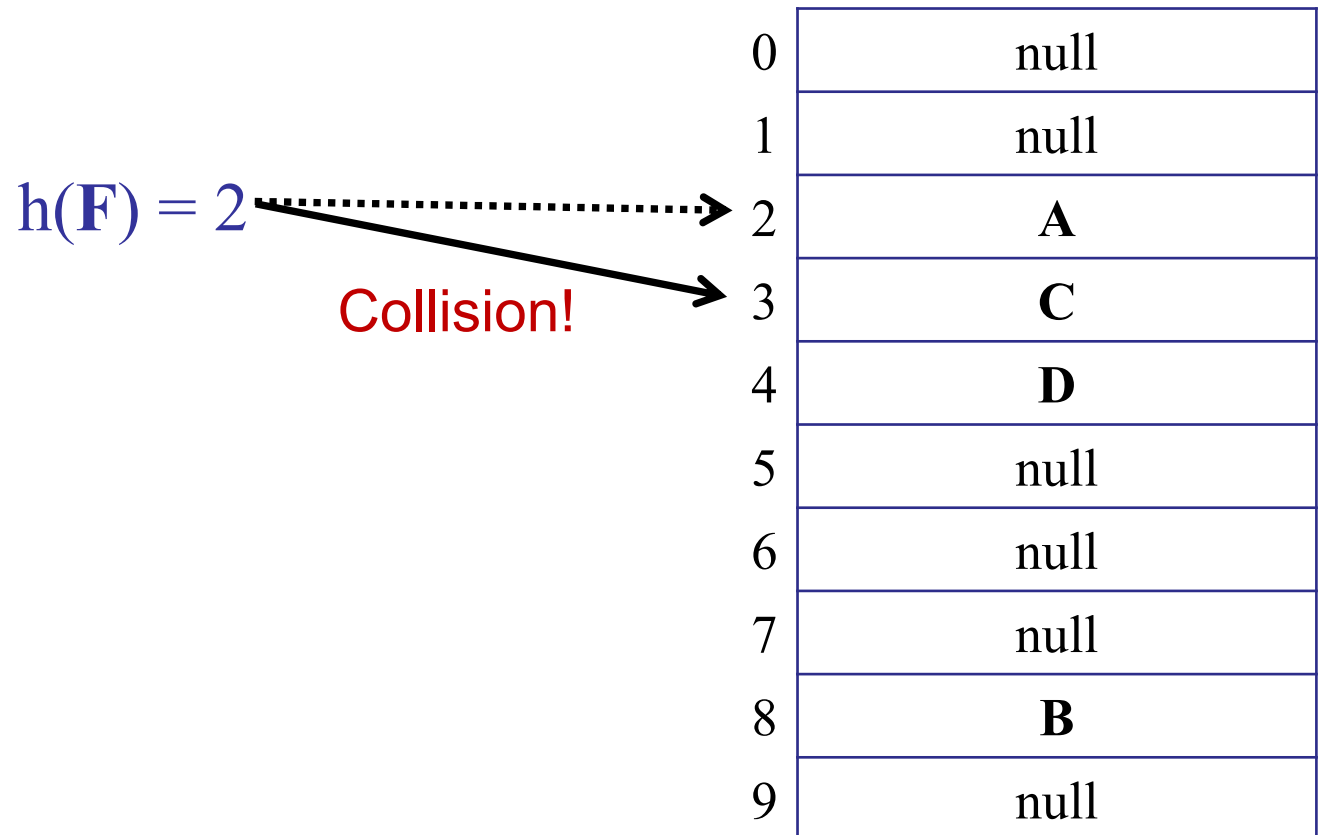
$h(F) = 2$ Collision! 

0	null
1	null
2	A
3	C
4	D
5	null
6	null
7	null
8	B
9	null

Open Addressing

On collision:

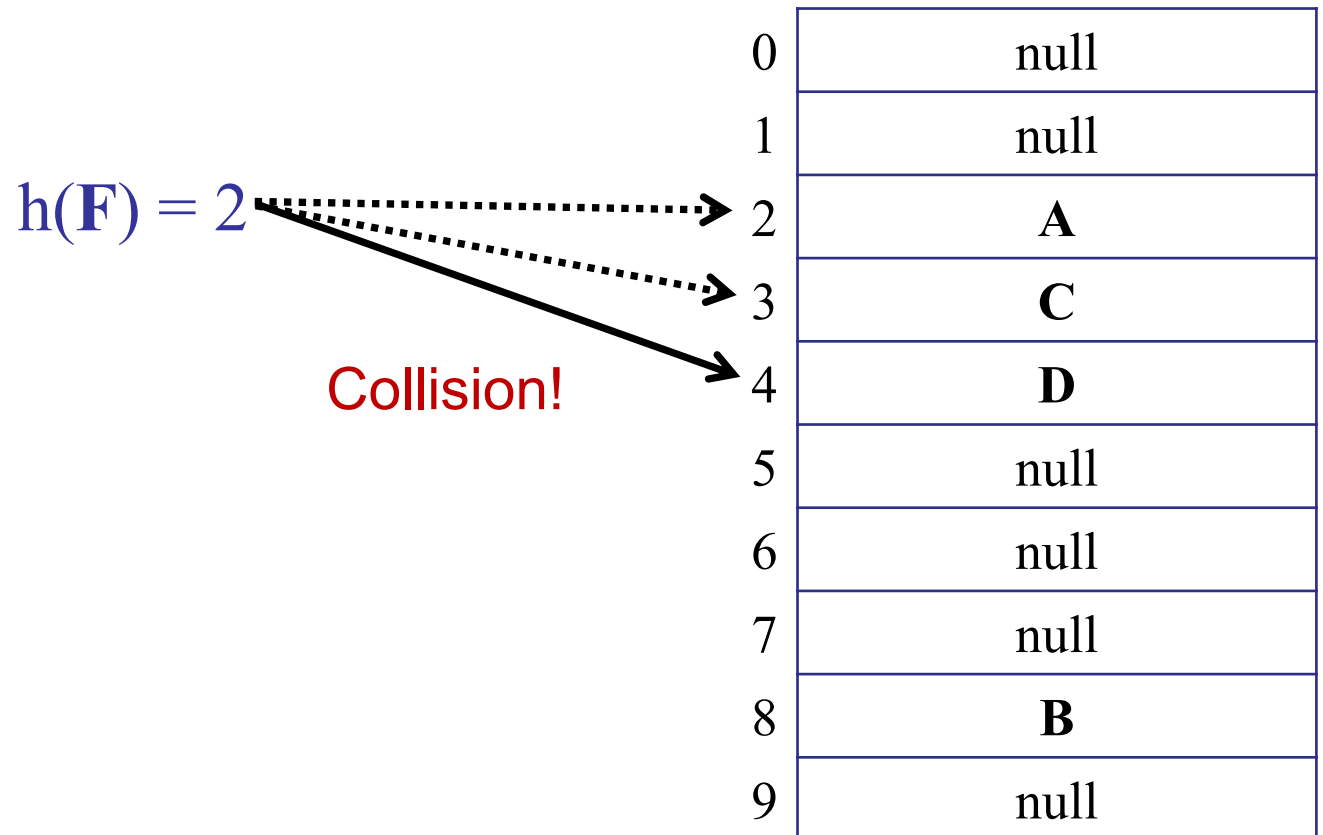
Probe a sequence of buckets until you find an empty one.



Open Addressing

On collision:

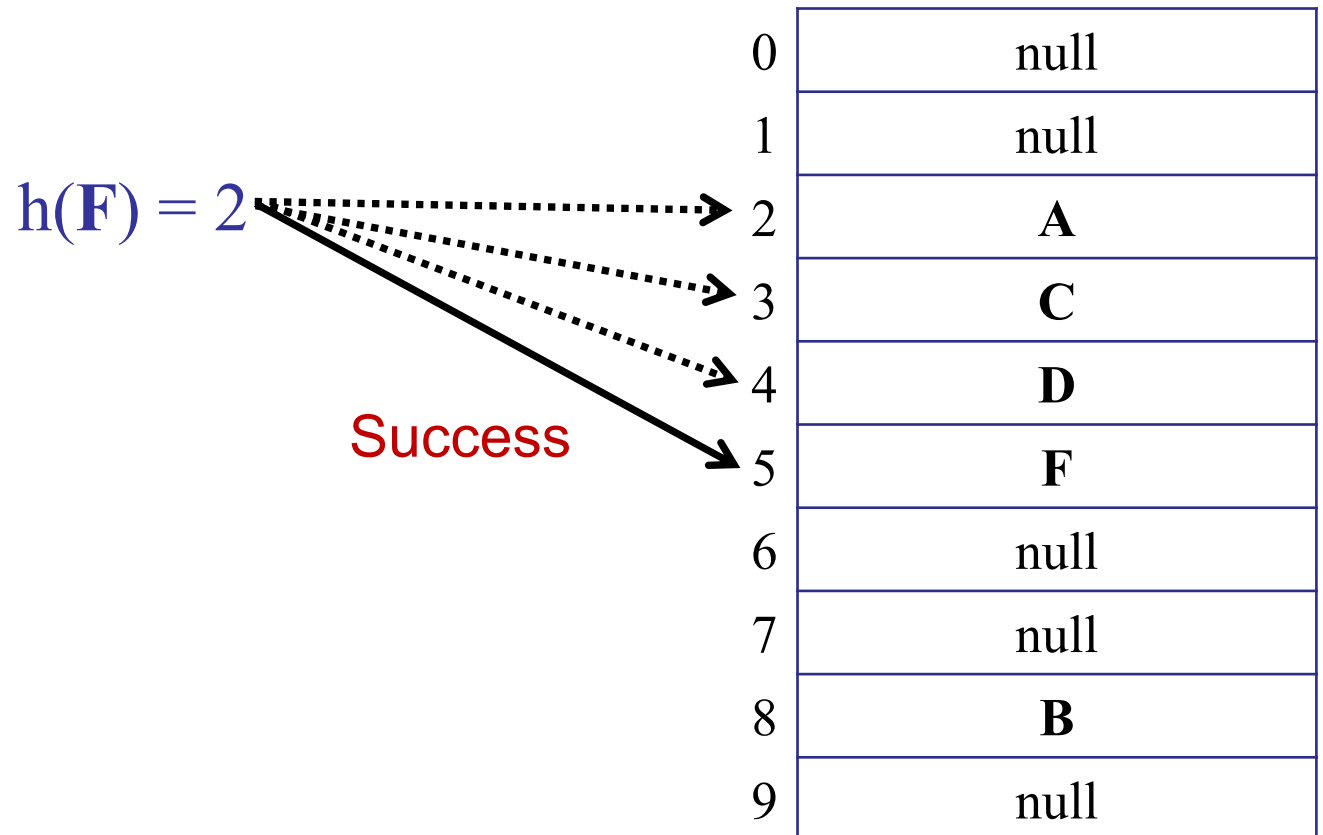
Probe a sequence of buckets until you find an empty one.



Open Addressing

On collision:

Probe a sequence of buckets until you find an empty one.



Open Addressing

On collision:

Probe a sequence of buckets until you find an empty one.

$h(F) = 2$

Success

0	null
1	null
2	A
3	C
4	D
5	F
6	null
7	null
8	B
9	null

Linear Probing:

- $h(k)+1, h(k)+2, h(k)+3 \dots$

Open Addressing

Hash Function re-defined:

$$h(\text{key}, i) : U \rightarrow \{1..m\}$$

Two parameters:

- key : the thing to map
- i : number of collisions

Open Addressing

Hash Function re-defined:

$$h(\text{key}, i) : U \rightarrow \{1..m\}$$

Example: Linear Probing

- $h(k, 1) = \text{hash of key } k$
- $h(k, 2) = h(k, 1) + 1$
- $h(k, 3) = h(k, 1) + 2$
- $h(k, 4) = h(k, 1) + 3$
- ...
- $h(k, i) = h(k, 1) + i \bmod m$

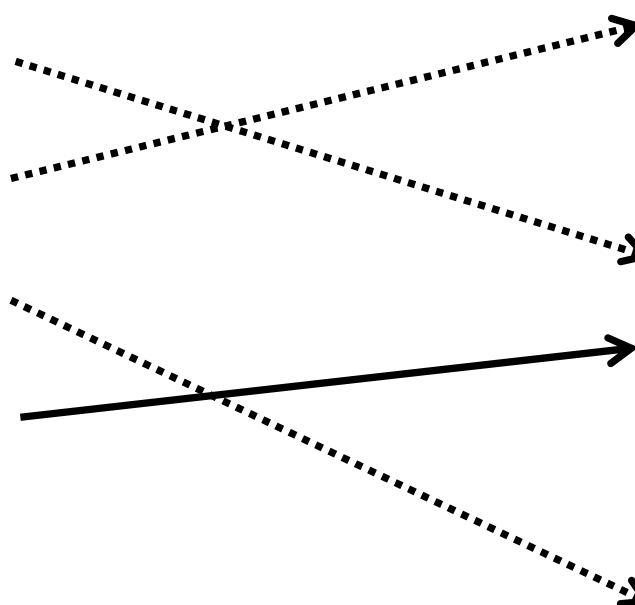
0	null
1	null
2	A
3	C
4	D
5	F
6	null
7	null
8	B
9	null

Open Addressing

Hash Function re-defined:

$$h(\text{key}, i) : U \rightarrow \{1..m\}$$

Example: Weird Probing

- $h(k, 1) = 4$
 - $h(k, 2) = 1$
 - $h(k, 3) = 8$
 - $h(k, 4) = 5$
- 
- The diagram illustrates the mapping of keys to slots in a hash table. On the left, four hash values are listed: $h(k, 1) = 4$, $h(k, 2) = 1$, $h(k, 3) = 8$, and $h(k, 4) = 5$. On the right, a hash table with 10 slots (0-9) is shown. Dotted arrows indicate the initial hash values: $h(k, 1)$ points to slot 4, $h(k, 2)$ points to slot 1, $h(k, 3)$ points to slot 8, and $h(k, 4)$ points to slot 5. Solid arrows indicate the final placement after probing: $h(k, 1)$ points to slot 1, $h(k, 2)$ points to slot 4, $h(k, 3)$ points to slot 5, and $h(k, 4)$ points to slot 8.

0	null
1	G
2	A
3	C
4	D
5	null
6	null
7	null
8	B
9	null

Open Addressing

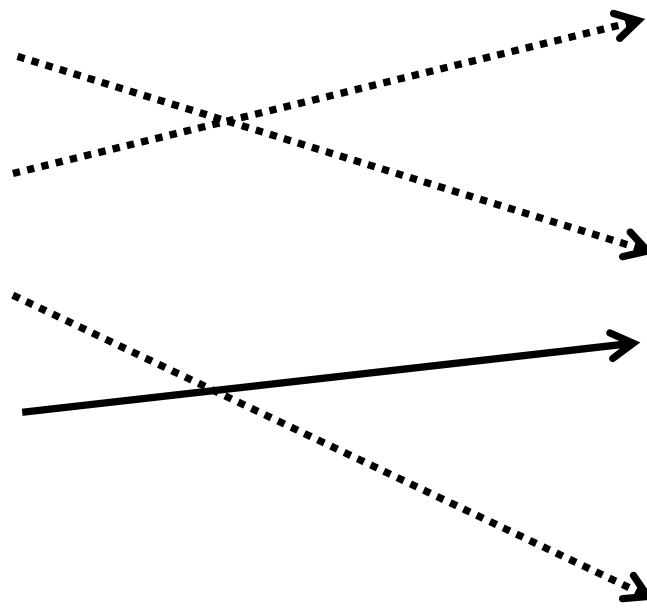
Hash Function re-defined:

$$h(\text{key}, i) : U \rightarrow \{1..m\}$$

search(key)

- $h(\text{key}, 1) = 4$
- $h(\text{key}, 2) = 1$
- $h(\text{key}, 3) = 8$
- $h(\text{key}, 4) = 5$

0	null
1	G
2	A
3	C
4	D
5	key
6	null
7	null
8	B
9	null



Open Addressing

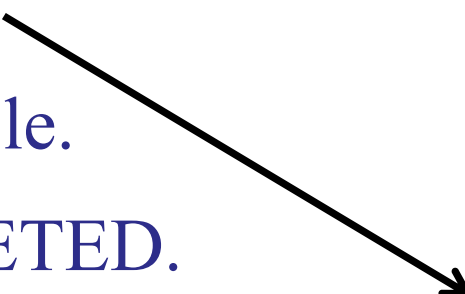
Hash Function re-defined:

$$h(\text{key}, i) : U \rightarrow \{1..m\}$$

delete(key)

- Find key to delete
- Remove it from table.
- Set bucket to DELETED.

(Tombstone value.)




0	null
1	G
2	A
3	C
4	D
5	DELETED
6	null
7	null
8	B
9	null

Hash Functions

Two properties of a good hash function:

1. $h(key, i)$ enumerates all possible buckets.
 - For every bucket j , there is some i such that:
$$h(key, i) = j$$
 - The hash function is permutation of $\{1..m\}$.
 - For linear probing: true!

What goes wrong if the sequence is not a permutation?

1. Search incorrectly returns key-not-found.
2. Delete fails.
3. Insert puts a key in the wrong place
-  4. Returns table-full even when there is still space left.

Hash Functions

Two properties of a good hash function:

2. Simple Uniform Hashing Assumption

Every key is equally likely to be mapped to every bucket, independently of every other key.

For $h(\text{key}, 1)$?

For every $h(\text{key}, i)$?

Hash Functions

Two properties of a good hash function:

2. Uniform Hashing Assumption

Every key is equally likely to be mapped to every *permutation*, independent of every other key.

n! permutations for probe sequence: e.g.,

- 1 2 3 4
- 1 2 4 3
- 1 4 2 3
- 1 4 3 2
- ...

Hash Functions

Two properties of a good hash function:

2. Uniform Hashing Assumption

Every key is equally likely to be mapped to every *permutation*, independent of every other key.

n! permutations for probe sequence: e.g.,

- 1 2 3 4 $\text{Pr}(1/m)$
- 1 2 4 3 $\text{Pr}(0)$
- 1 4 2 3 $\text{Pr}(0)$
- 1 4 3 2 $\text{Pr}(0)$
- ...

NOT Linear Probing

Linear probing

In practice, linear probing is very fast!

- Why? Caching!
- It is *cheap* to access nearby array cells.
 - Example: access $T[17]$
 - Cache loads: $T[10..50]$
 - Almost 0 cost to access $T[18]$, $T[19]$, $T[20]$, ...
- If the table is 1/4 full, then there will be clusters of size: $\theta(\log n)$
 - Cache may hold entire cluster!
 - No worse than wacky probe sequence.

Open Addressing

Properties of a good hash function:

2. Uniform Hashing Assumption

Every key is equally likely to be mapped to every *permutation*, independent of every other key.

n! permutations for probe sequence: e.g.,

- 1 2 3 4
- 1 2 4 3
- 1 4 2 3
- 1 4 3 2
- ...

Double Hashing

- Start with two ordinary hash functions:

$$f(k), g(k)$$

- Define new hash function:

$$h(k, i) = f(k) + i \cdot g(k) \mod m$$

- Note:
 - Since $f(k)$ is good, $f(k, 1)$ is “almost” random.
 - Since $g(k)$ is good, the probe sequence is “almost” random.

Double Hashing

Hash function

$$h(k, i) = f(k) + i \cdot g(k) \mod m$$

Claim: if $g(k)$ is relatively prime to m , then $h(k, i)$ hits all buckets.

- Assume not: then for some distinct $i, j < m$:

$$f(k) + i \cdot g(k) = f(k) + j \cdot g(k) \mod m$$

$$\rightarrow i \cdot g(k) = j \cdot g(k) \mod m$$

$$\rightarrow (i - j) \cdot g(k) = 0 \mod m$$

$$\rightarrow g(k) \text{ not relatively prime to } m, \text{ since } (i, j < m)$$

Double Hashing

Hash function

$$h(k, i) = f(k) + i \cdot g(k) \mod m$$

Claim: if $g(k)$ is relatively prime to m , then $h(k, i)$ hits all buckets.

Example: if $(m = 2^r)$, then choose $g(k)$ odd.

Performance of Open Addressing

If ($m=n$), what is the expected insert time, under uniform hashing assumption?


1. $O(1)$
2. $O(\log n)$
3. $O(n)$
4. $O(n^2)$
- ✓ 5. None of the above.

Performance of Open Addressing

- Chaining:
 - When $(m==n)$, we can still add new items to the hash table.
 - We can still search efficiently.
- Open addressing:
 - When $(m==n)$, the table is full.
 - We cannot insert any more items.
 - We cannot search efficiently.


Performance of Open Addressing

Define:

- Load $\alpha = n / m$  Average # items / bucket
- Assume $\alpha < 1$.

Performance of Open Addressing

Define:

- Load $\alpha = n / m$  Average # items / bucket
- Assume $\alpha < 1$.

Claim:

For n items, in a table of size m , assuming *uniform hashing*, the expected cost of an operation is:

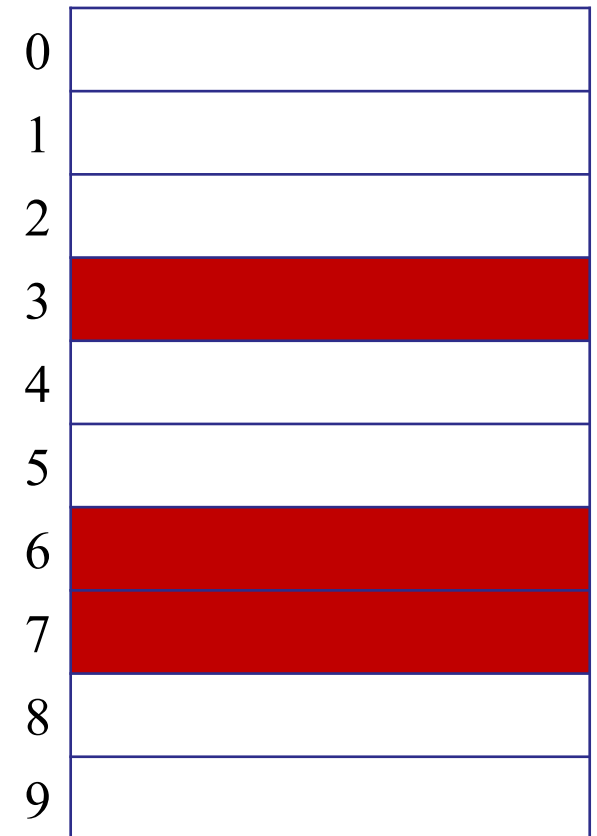
$$\leq \frac{1}{1 - \alpha}$$

Example: if ($\alpha=90\%$), then $E[\# \text{ probes}] = 10$

Performance of Open Addressing

Proof of Claim:

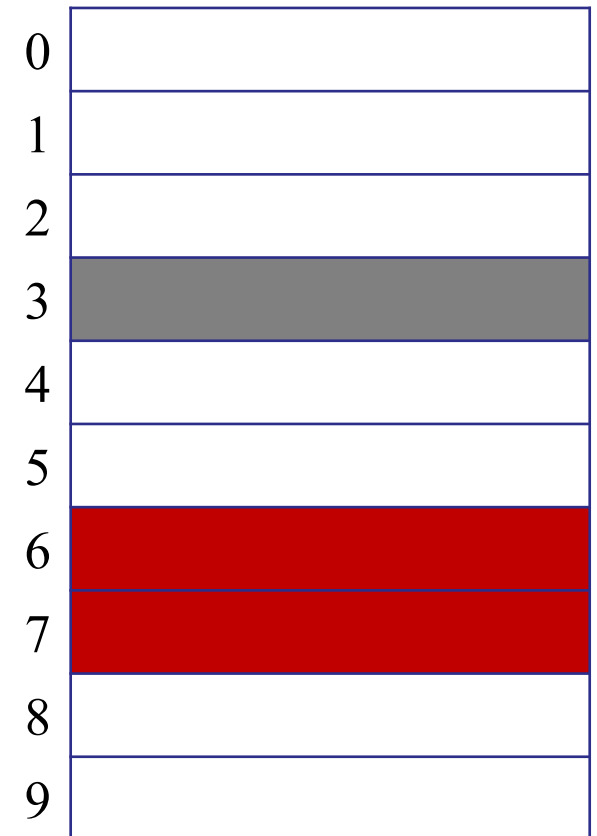
- First probe: probability that first bucket is full is: n/m



Performance of Open Addressing

Proof of Claim:

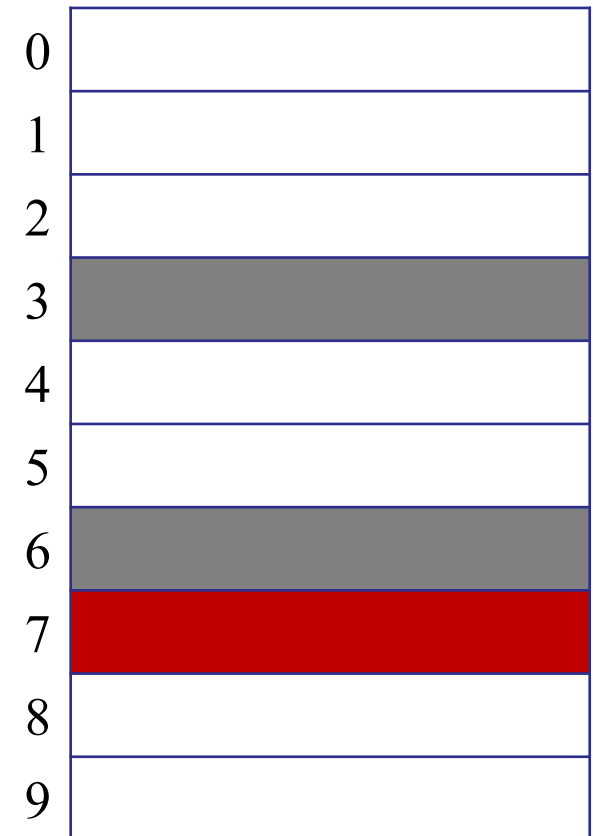
- First probe: probability that first bucket is full is: n/m
- Second probe: if first bucket is full, then the probability that the second bucket is also full: $(n - 1) / (m - 1)$



Performance of Open Addressing

Proof of Claim:

- First probe: probability that first bucket is full is: n/m
- Second probe: if first bucket is full, then the probability that the second bucket is also full: $(n - 1) / (m - 1)$
- Third probe: probability is full: $(n - 2) / (m - 2)$



Performance of Open Addressing

Proof of Claim:

– Expected cost:

$$1 + \frac{n}{m} \left(\text{Expected cost of remaining probes} \right)$$

The diagram illustrates the components of the expected cost formula. A black rounded rectangle contains the text "Expected cost of remaining probes". An arrow points from the text "First probe" to the constant "1" in the formula. Another arrow points from the text "Probability of collision on first probe" to the fraction $\frac{n}{m}$.

First probe

Probability of collision on first probe

Performance of Open Addressing

Proof of Claim:

– Expected cost:

$$1 + \frac{n}{m} \left(1 + \frac{n-1}{m-1} \left(\text{Expected cost of remaining probes} \right) \right)$$

The diagram illustrates the components of the recursive formula for the expected cost of open addressing. Three red text labels at the bottom are connected by arrows to parts of the formula above:

- First probe**: An arrow points from this label to the initial '1' in the formula.
- Probability of collision on first probe**: An arrow points from this label to the fraction $\frac{n}{m}$.
- Probability of collision on second probe**: An arrow points from this label to the inner recursive term $1 + \frac{n-1}{m-1} (\dots)$.

The inner recursive term is enclosed in a dark rounded rectangle and labeled **Expected cost of remaining probes**.

Performance of Open Addressing

Proof of Claim:

- Expected cost:

$$1 + \frac{n}{m} \left(1 + \frac{n-1}{m-1} \left(1 + \frac{n-2}{m-2} \left(\square \square \square \right) \right) \right)$$

First probe

Second probe

Third probe

Performance of Open Addressing

Proof of Claim:

– Expected cost:

$$1 + \frac{n}{m} \left(1 + \frac{n-1}{m-1} \left(1 + \frac{n-2}{m-2} \left(\square \square \square \right) \right) \right)$$

– Note:

$$\frac{n-i}{m-i} \leq \frac{n}{m} \leq \alpha$$

Performance of Open Addressing

Proof of Claim:

– Expected cost:

$$1 + \frac{n}{m} \left(1 + \frac{n-1}{m-1} \left(1 + \frac{n-2}{m-2} \left(\begin{array}{ccc} \square & \square & \square \end{array} \right) \right) \right)$$

$$\leq 1 + \alpha (1 + \alpha (1 + \alpha (\dots)))$$

Performance of Open Addressing

Proof of Claim:

– Expected cost:

$$1 + \frac{n}{m} \left(1 + \frac{n-1}{m-1} \left(1 + \frac{n-2}{m-2} \left(\begin{array}{ccc} \square & \square & \square \end{array} \right) \right) \right)$$

$$\leq 1 + \alpha \left(1 + \alpha \left(1 + \alpha (\dots) \right) \right)$$

$$\leq 1 + \alpha + \alpha^2 + \alpha^3 + \dots$$

Performance of Open Addressing

Proof of Claim:

– Expected cost:

$$1 + \frac{n}{m} \left(1 + \frac{n-1}{m-1} \left(1 + \frac{n-2}{m-2} \left(\begin{array}{ccc} \square & \square & \square \end{array} \right) \right) \right)$$


$$\leq 1 + \alpha (1 + \alpha (1 + \alpha (\dots)))$$

$$\leq 1 + \alpha + \alpha^2 + \alpha^3 + \dots$$

$$\leq \frac{1}{1 - \alpha}$$

Performance of Open Addressing

Define:

- Load $\alpha = n / m$  Average # items / bucket
- Assume $\alpha < 1$.

Claim:

For n items, in a table of size m , assuming *uniform hashing*, the expected cost of an operation is:

$$\leq \frac{1}{1 - \alpha}$$

Example: if ($\alpha=90\%$), then $E[\# \text{ probes}] = 10$

Advantages...

Open addressing:

- Saves space
 - Empty slots vs. linked lists.
- Rarely allocate memory
 - No new list-node allocations.
- Better cache performance
 - Table all in one place in memory
 - Fewer accesses to bring table into cache.
 - Linked lists can wander all over the memory.

Disadvantages...

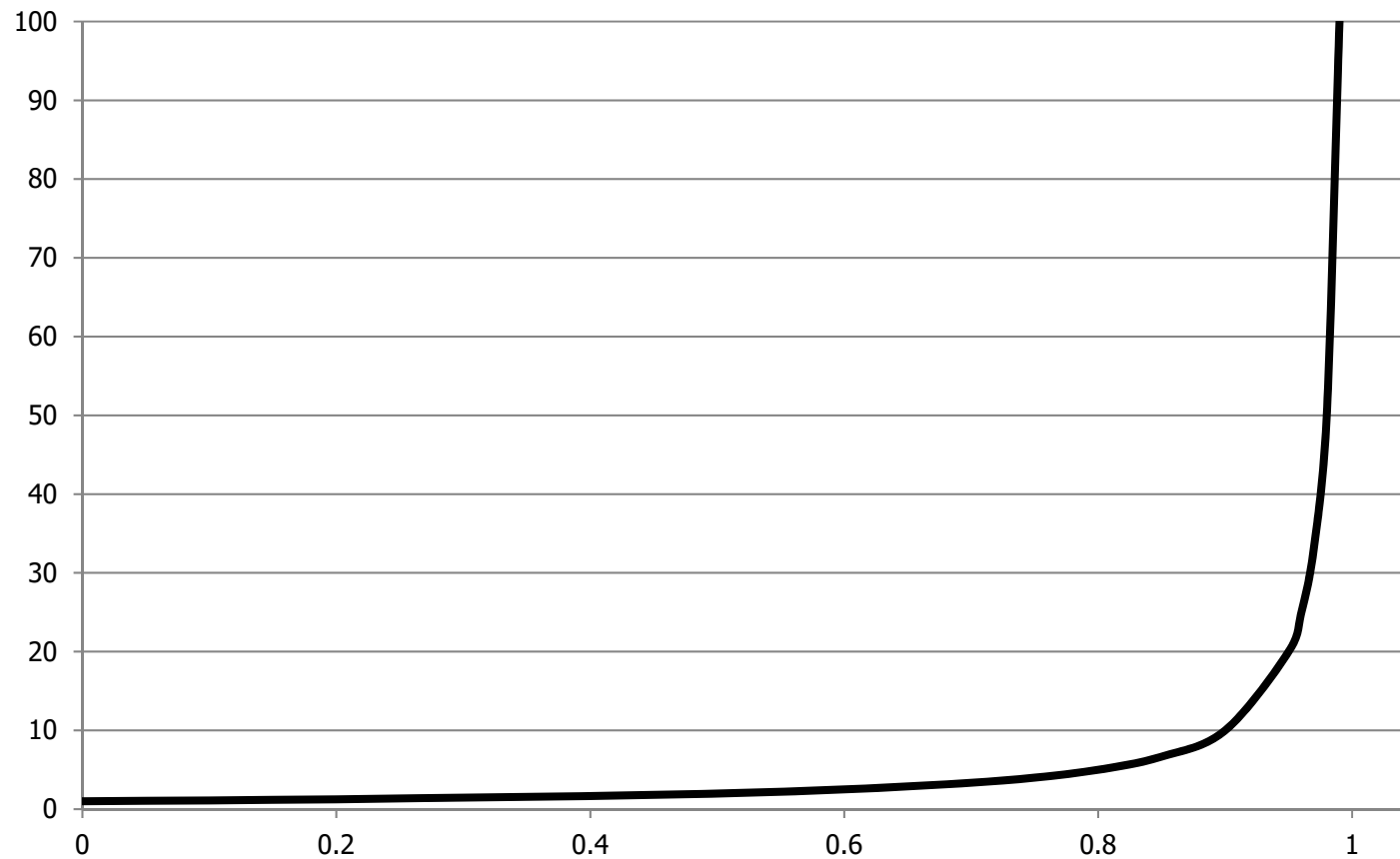
Open addressing:

- More sensitive to choice of hash functions.
 - Clustering is a common problem.
 - See issues with linear probing.
- More sensitive to load.
 - Performance degrades badly as $\alpha \rightarrow 1$.

Disadvantages...

Open addressing:

- Performance degrades badly as $\alpha \rightarrow 1$.



Hashing overview

- What is a hash function?
- Collision resolution: chaining
- Java hashing
- Collision resolution: open addressing
- Table (re)sizing

Table Size

How large should the table be?

- Assume: Hashing with Chaining
- Assume: Simple Uniform Hashing
- Expected search time: $O(1 + n/m)$
- Optimal size: $m = \Theta(n)$
 - if $(m < 2n)$: too many collisions.
 - if $(m > 10n)$: too much wasted space.
- Problem: we don't know n in advance.

Table Size

Idea:

- Start with small (constant) table size.
- Grow (and shrink) table as necessary.

Example:

- Initially, $m = 10$.
- After inserting 6 items, table too small! Grow...
- After deleting $n-1$ items, table too big! Shrink...

Table Size

How to grow the table:

1. Choose new table size m .
2. Choose new hash function h .
 - Hash function depends on table size!
 - Remember: $h : U \rightarrow \{1..m\}$
3. For each item in the old hash table:
 - Compute new hash function.
 - Copy item to new bucket.

Not like Java hashCode!



Table Size

Time complexity of growing the table:

– Assume:

- Let m_1 be the size of the old hash table.
- Let m_2 be the size of the new hash table.
- Let n be the number of elements in the hash table.

– Costs:

- Scanning old hash table: $O(m_1)$
- Inserting each element in new hash table: $O(1)$
- Total: $O(m_1 + n)$

Table Size

Time complexity of growing the table:

– Assume:

- Size $m_1 < n$.
- Size $m_2 > 2n$

– Costs:

- Total: $O(m_1 + n)$.
 $= O(n)$

Table Size

Time complexity of growing the table:

Wait! What is the cost of initializing the new table?

- Initializing a table of size x takes x time!
- Costs:

Total: $O(m_1 + m_2 + n)$

Table Size

Time complexity of growing the table:

– Assume:

- Let m_1 be the size of the old hash table.
- Let m_2 be the size of the new hash table.
- Let n be the number of elements in the hash table.

– Costs:

- Scanning old hash table: $O(m_1)$
- Creating new hash table: $O(m_2)$
- Inserting each element in new hash table: $O(1)$
- Total: $O(m_1 + m_2 + n)$

How fast to grow?

Idea 1: Increment table size by 1

- if ($n == m$): $m = m+1$

- Cost of resize:

- Size $m_1 = n$.
- Size $m_2 = n+1$.
- Total: $O(n)$

Initially: $m = 8$

What is the cost of inserting n items?

1. $O(n)$
2. $O(n \log n)$
- ✓ 3. $O(n^2)$
4. $O(n^3)$
5. None of the above.

How fast to grow?

Idea 1: Increment table size by 1

- When ($n == m$): $m = m+1$
- Cost of each resize: $O(n)$

Table size	8	8	9	10	11	12	...	$n+1$
Number of items	0	7	8	9	10	11	...	n
Number of inserts		7	1	1	1	1	...	1
Cost		7	8	9	10	11		n

- Total cost: $(7 + 8 + 9 + 10 + 11 + \dots + n) = O(n^2)$

How fast to grow?

Idea 2: Double table size

- if ($n == m$): $m = 2m$
- Cost of resize:
 - Size $m_1 = n$.
 - Size $m_2 = 2n$.
 - Total: $O(n)$

How fast to grow?

Idea 2: Double table size

- When $(n == m)$: $m = 2m$
- Cost of each resize: $O(n)$

Table size	8	8	16	16	16	16	16	16	16	16	32	32	32	...	2n
# of items	0	7	8	9	10	11	12	13	14	15	16	17	18	...	n
# of inserts		7	1	1	1	1	1	1	1	1	1	1	1	...	1
Cost		7	8	1	1	1	1	1	1	1	16	1	1		n

- Total cost: $(7 + 15 + 31 + \dots + n) = O(n)$

How fast to grow

Idea 2: Double table size

Cost of Resizing:

Table size	Total Resizing Cost
8	8
16	$(8 + 16)$
32	$(8 + 16 + 32)$
64	$(8 + 16 + 32 + 64)$
128	$(8 + 16 + 32 + 64 + 128)$
...	...
m	$<(1+2+4+8+\dots+m) \leq O(m)$

How fast to grow?

Idea 2: Double table size

- if ($n == m$): $m = 2m$
 - Cost of resize: $O(n)$
 - Cost of inserting n items + resizing: $O(n)$
- Most insertions: $O(1)$
- Some insertions: linear cost (expensive)
- Average cost: $O(1)$

How fast to grow?

Idea 3: Square table size

- When $(n == m)$: $m = m^2$

Table size	Total Resizing Cost
8	?
64	?
4,096	?
16,777,216	?
...	...
m	?

Assume: square table size

What is the cost of inserting n items?

1. $O(\log n)$
2. $O(\sqrt{n})$
3. $O(n)$
4. $O(n \log n)$
5. $O(n^2)$
6. $O(2^n)$
7. None of the above.

How fast to grow?

Idea 3: Square table size

- if $(n == m)$: $m = m^2$
- Cost of resize:
 - Size $m_1 = n$.
 - Size $m_2 = n^2$.
 - Total: $O(m_1 + m_2 + n)$
 $= O(n + n^2 + n)$
 $= O(n^2)$

How fast to grow?

Idea 3: Square table size

- When $(n == m)$: $m = m^2$

# Items	Total Resizing Cost
8	64
64	$(64 + 4,096)$
4,096	$(64 + 4,096 + \dots)$
...	...
n	$> n^2$
	$= O(n^2)$

How fast to grow?

Idea 3: Square table size

- When $(n == m)$: $m = m^2$

# Items	Resizing Cost	Insert Cost
8	64	8
64	$(64 + 4,096)$	64
4,096	$(64 + 4,096 + \dots)$	4,096
...
n	$> n^2$	n
	$< O(n^2)$	$O(n)$

How fast to grow?

Idea 3: Square table size

- if ($n == m$): $m = m^2$

- Cost of resize:

- Total: $O(n^2)$

- Cost of inserts:

- Total: $O(n)$

Why else is squaring the table size bad?

1. Resize takes too long to find items to copy.
2. Inefficient space usage.
3. Searching is more expensive in a big table.
4. Inserting is more expensive in big table.
5. Deleting is more expensive in a big table.

Deleting Elements

Basic procedure: (chained hash tables)

Delete(*key*)

1. Calculate hash of *key*.
2. Let *L* be the linked list in the specified bucket.
3. Search for item in linked list *L*.
4. Delete item from linked list *L*.

Cost:

- Total: $O(1 + n/m)$

Deleting Elements

What happens if too many items are deleted?

- Table is too big!
- Shrink the table...
- Try 1:
 - If $(n == m)$, then $m = 2m$.
 - If $(n < m/2)$ then $m = m/2$.

Deleting Elements

Rules for shrinking and growing:

– Try 1:

- If $(n == m)$, then $m = 2m$.
- If $(n < m/2)$ then $m = m/2$.

– Example problem:

- Start: $n=100, m=200$
- Delete: $n=99, m=200 \rightarrow$ shrink to $m=100$
- Insert: $n=100, m=100 \rightarrow$ grow to $m=200$
- Repeat...

Deleting Elements

Example execution:

- Start: $n=100$, $m=200$
- cost=100 • Delete: $n=99$, $m=200 \rightarrow$ shrink to $m=100$
- cost=100 • Insert: $n=100$, $m=100 \rightarrow$ grow to $m=200$
- cost=100 • Delete: $n=99$, $m=200 \rightarrow$ shrink to $m=100$
- cost=100 • Insert: $n=100$, $m=100 \rightarrow$ grow to $m=200$
- cost=100 • Delete: $n=99$, $m=200 \rightarrow$ shrink to $m=100$
- cost=100 • Insert: $n=100$, $m=100 \rightarrow$ grow to $m=200$
- Repeat...

Deleting Elements

Rules for shrinking and growing:

– Try 2:

- If $(n == m)$, then $m = 2m$.
- If $(n < m/4)$, then $m = m/2$.

– Claim:

- Every time you double a table of size m , at least $m/2$ new items were added.
- Every time you shrink a table of size m , at least $m/4$ items were deleted.

Amortized Analysis

Technique for analyzing “average” cost:

- Common in data structure analysis
- Like paying rent:
 - You don’t pay rent every day!
 - Pay \$900/month = \$30/day.

Definition:

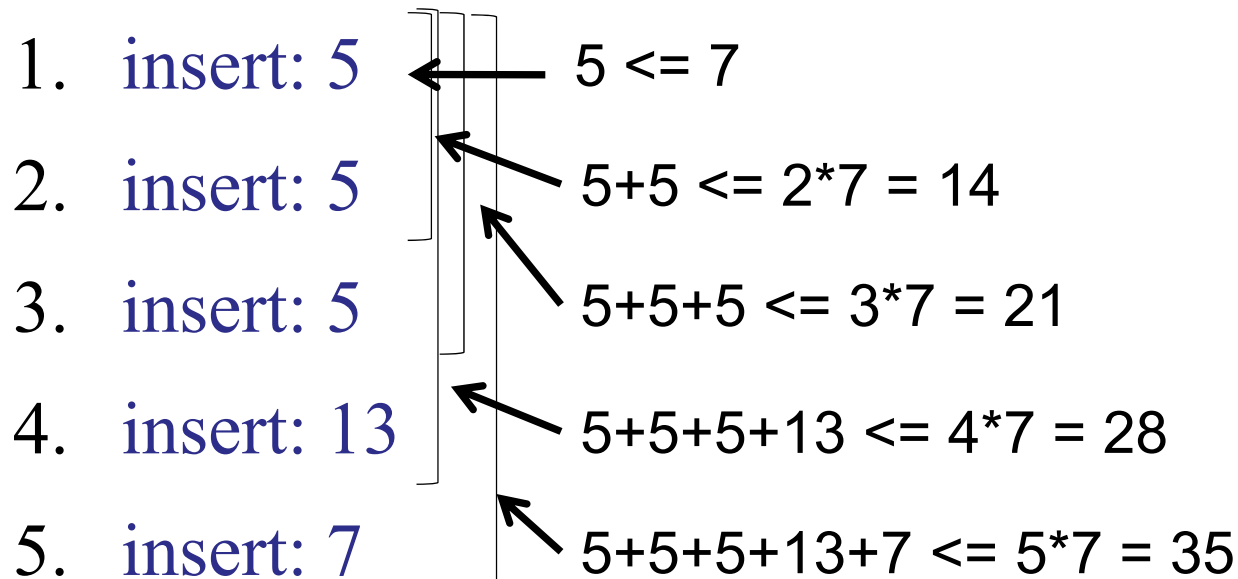
- Operation has amortized cost $T(n)$ if for every integer k , the cost of k operations is $\leq k T(n)$

Amortized Analysis

Definition:

- Operation has amortized cost $T(n)$ if for every integer k , the cost of k operations is $\leq k T(n)$

Example: amortized cost = 7



Amortized Analysis

“amortized” is NOT “average”

Definition:

- Operation has amortized cost $T(n)$ if for every integer k , the cost of k operations is $\leq k T(n)$

Example: amortized cost **NOT** 7

-
1. insert: 13 $13 > 7$
2. insert: 5 $13+5 > 2*7 = 14$
3. insert: 5 $13+5+5 > 3*7 = 21$
4. insert: 5 $13+5+5+5 \leq 4*7 = 28$
5. insert: 7 $5+5+5+13+7 \leq 5*7 = 35$
- The diagram illustrates a sequence of operations on a data structure. A vertical line represents the state of the structure. Brackets on the left group the operations: the first operation (insert: 13) is alone; the next three (insert: 5, insert: 5, insert: 5) are grouped together; and the final operation (insert: 7) is alone. Arrows point from the cumulative cost calculations on the right to the corresponding points on the vertical line.

Amortized Analysis

Definition:

- Operation has amortized cost $T(n)$ if for every integer k , the cost of k operations is $\leq k T(n)$

Example: (Hash Tables)

- Inserting k elements into a hash table takes time $O(k)$.
- Conclusion:

The insert operation has amortized cost $O(1)$.

Amortized Analysis

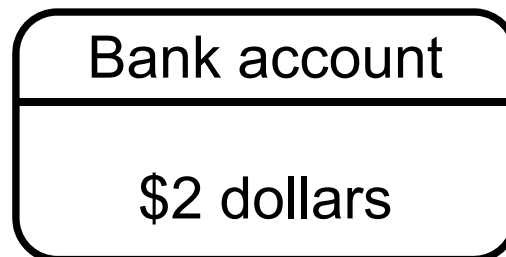
Accounting Method (paying rent)

- Imagine a bank account **B**.
- Each operation adds money to the bank account.
- Every step of the algorithm spends money:
 - Immediate money: to perform the operation.
 - Deferred money: from the bank account.
- Total cost execution = total money
 - Average time / operation = money / num. ops

Amortized Analysis

Accounting Method Example (Hash Table)

- Each table has a bank account.
- Each time an element is added to the table, it adds $O(1)$ dollars to the bank account, uses $O(1)$ dollars to insert element.
- A table with k new elements since last resize has k dollars in bank.



0	null
1	null
2	(k_1, A)
3	null
4	null
5	null
6	null
7	null
8	(k_2, B)
9	null

Amortized Analysis

Accounting Method Example (Hash Table)

- Each table has a bank account.
- Each time an element is added to the table, it adds $O(1)$ dollars to the bank account.
- Claim:
 - Resizing a table of size m takes $O(m)$ time.
 - If you resize a table of size m , then:
 - at least $m/2$ new elements since last resize
 - bank account has $\Theta(m)$ dollars.

Amortized Analysis

Accounting Method Example (Hash Table)

- Each table has a bank account.
- Each time an element is added to the table, it adds $O(1)$ dollars to the bank account.
- Pay for resizing from the bank account!
- Strategy:
 - Analyze inserts ignoring cost of resizing.
 - Ensure that bank account always is big enough to pay for resizing.

Amortized Analysis

Total cost: Inserting k elements costs:

- Deferred dollars: $O(k)$ (to pay for resizing)
- Immediate dollars: $O(k)$ for inserting elements in table
- Total (Deferred + Immediate): $O(k)$

Amortized Analysis

Total cost: Inserting k elements costs:

- Deferred dollars: $O(k)$ (to pay for resizing)
- Immediate dollars: $O(k)$ for inserting elements in table
- Total (Deferred + Immediate): $O(k)$

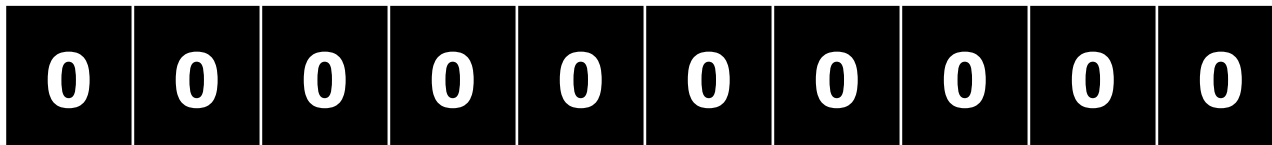
Cost per operation:

- Deferred dollars: $O(1)$
- Immediate dollars: $O(1)$
- Total: $O(1)$ / per operation

Example: Binary Counter

Counter ADT:

- `increment()`
- `read()`



Example: Binary Counter

Counter ADT:

- increment()
- read()

increment()

[illegible]

Example: Binary Counter

Counter ADT:

- increment()
- read()

increment(), increment()

[illegible]

Example: Binary Counter

Counter ADT:

- increment()
- read()

increment(), increment(), increment()

[illegible]

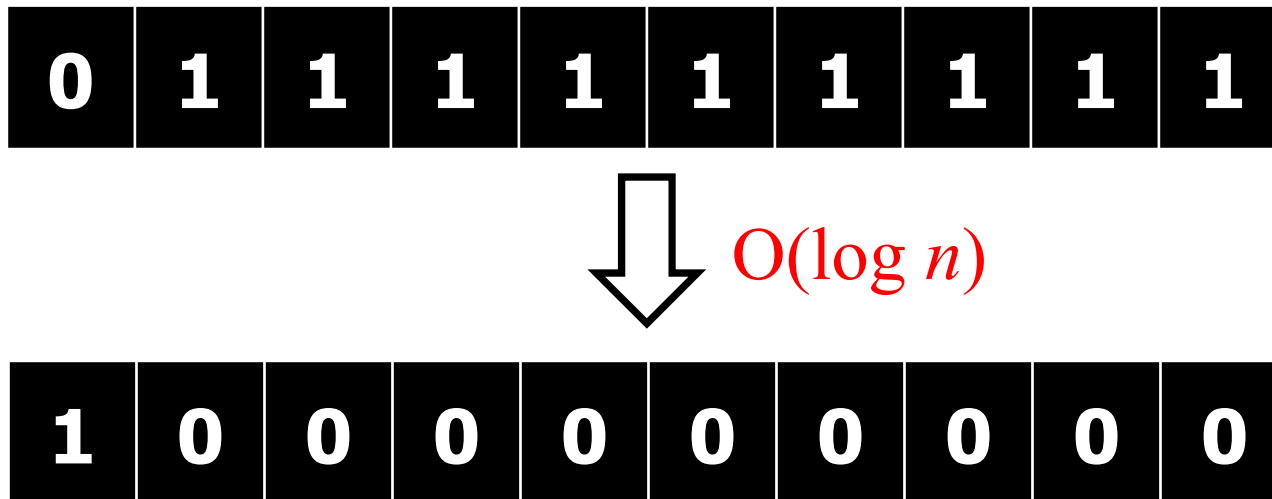
What is the worst-case cost of incrementing a counter with max-value n ?

1. $O(1)$
- ✓ 2. $O(\log n)$
3. $O(n)$
4. $O(n^2)$
5. I have no idea.

Example: Binary Counter

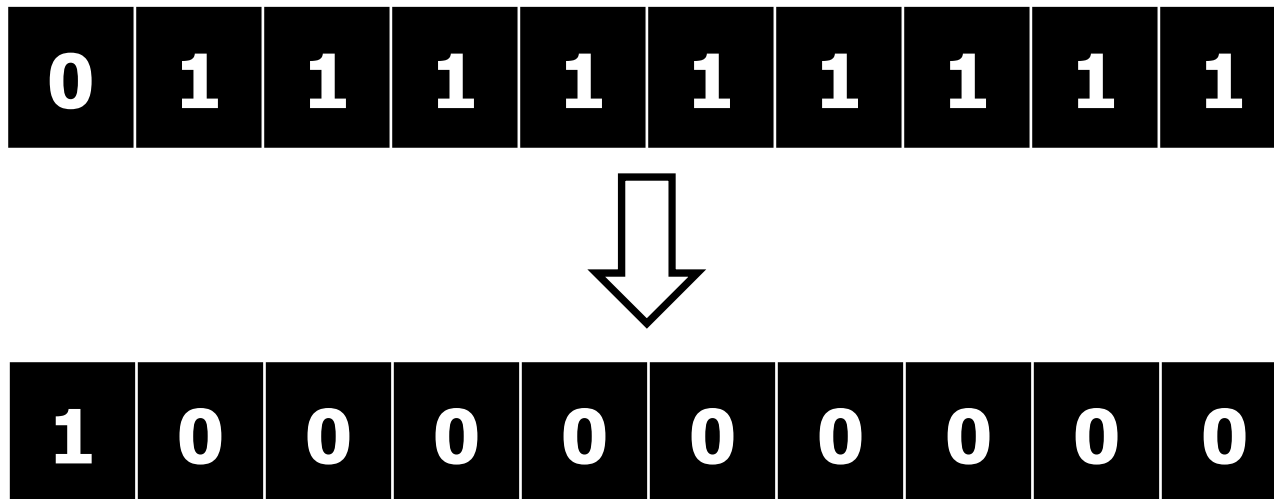
Question: If we increment the counter to n , what is the amortized cost per operation?

- Easy answer: $O(\log n)$
- More careful analysis....



Observation:

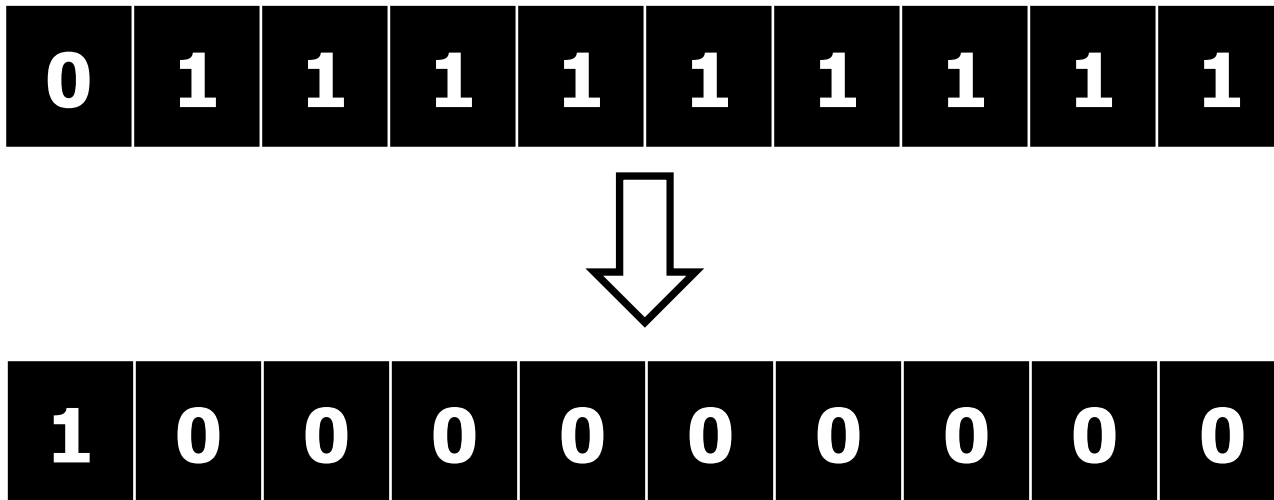
During each increment, only one bit is changed
from: $0 \rightarrow 1$



Observation:

Accounting method: each bit has a bank account.

Whenever you change it from $0 \rightarrow 1$, add one dollar.



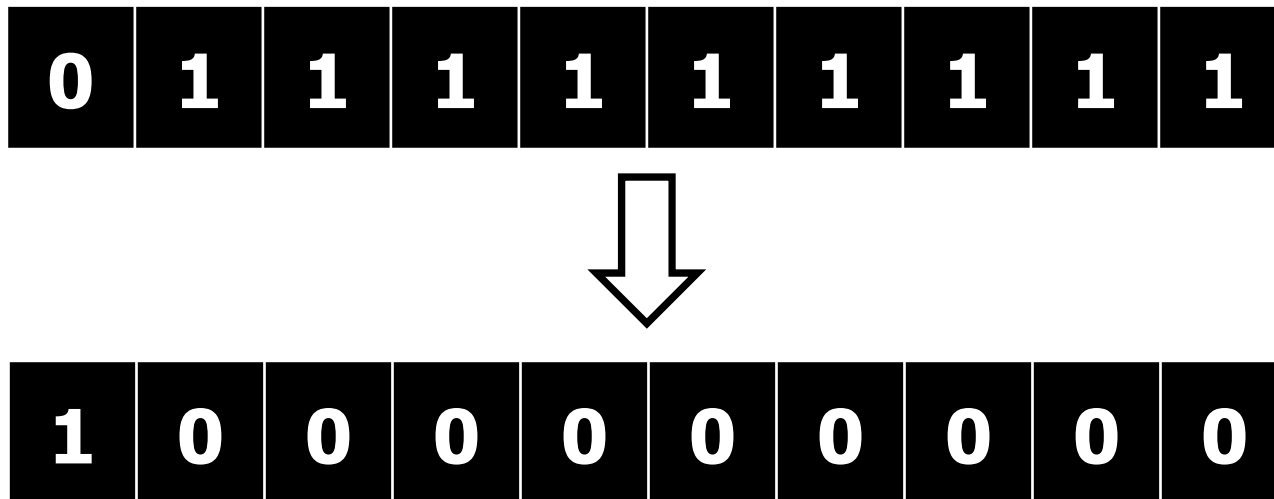
Example: Binary Counter

Observation:

Accounting method: each bit has a bank account.

Whenever you change it from $0 \rightarrow 1$, add one dollar.

Whenever you change it from $1 \rightarrow 0$, pay one dollar.



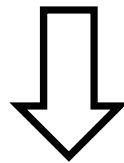
Example: Binary Counter

Counter ADT

increment()

0	0	0	0	0	0	0	0	0	0
----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

0 0 0 0 0 0 0 0 0 0



0	0	0	0	0	0	0	0	0	1
----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

0 0 0 0 0 0 0 0 0 1

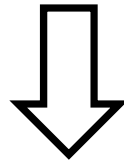
Example: Binary Counter

Counter ADT

increment(), increment()

0	0	0	0	0	0	0	0	0	1
----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

0 0 0 0 0 0 0 0 0 1



0	0	0	0	0	0	0	0	1	0
----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

0 0 0 0 0 0 0 0 1 0

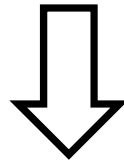
Example: Binary Counter

Counter ADT

increment(), increment(), increment()

0	0	0	0	0	0	0	0	1	0
----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

0 0 0 0 0 0 0 0 1 0



0	0	0	0	0	0	0	0	1	1
----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

0 0 0 0 0 0 0 0 1 1

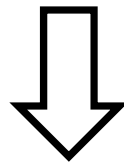
Example: Binary Counter

Counter ADT

increment()

0	1	1	1	1	1	1	1	1	1
----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

0 1 1 1 1 1 1 1 1 1



1	0	0	0	0	0	0	0	0	0
----------	----------	----------	----------	----------	----------	----------	----------	----------	----------

1 0 0 0 0 0 0 0 0 0

Example: Binary Counter

Observation:

Amortized cost of increment: 2

- One operation to switch one $0 \rightarrow 1$
- One dollar (for bank account of switched bit).

(All switches from 1 \rightarrow 0 paid for by bank account.)

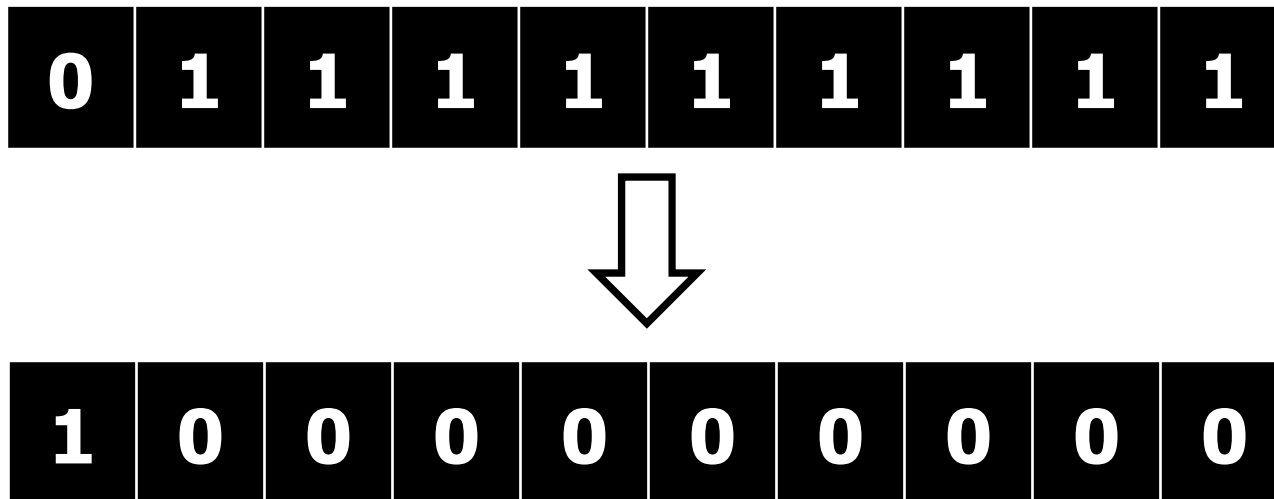


Table Size Rules

Rules for shrinking and growing:

- If $(n == m)$, then $m = 2m$.
- If $(n < m/4)$, then $m = m/2$.

– Claim:

- Every time you double a table of size m , at least $m/2$ new items were added.
- Every time you shrink a table of size m , at least $m/4$ items were deleted.

Amortized Analysis

Accounting Method

- Each table has a bank account.
- Each time an element is added to the table, it adds $O(1)$ dollars to the bank account.
- Claim:
 - Resizing a table of size m takes $O(m)$ time.
 - If you resize a table of size m , then:
 - at least $m/2$ new elements since last resize
 - bank account has $\Theta(m)$ dollars.

Amortized Analysis

Total cost: Inserting k elements costs:

- Deferred dollars: $O(k)$ (to pay for resizing)
- Immediate dollars: $O(k)$ for inserting elements in table
- Total (Deferred + Immediate): $O(k)$

Cost per operation:

- Deferred dollars: $O(1)$
- Immediate dollars: $O(1)$
- Total: $O(1)$ / per operation

Hash Table Resizing

Conclusion: Hashing with Chaining

- with Simple Uniform Hashing Assumption (SUHA)

Cost per operation:

- Insert operation: **amortized $O(1)$**
- Search operation: **expected $O(1)$**

Notes:

- Inserts are amortized because of table resizing.
- Inserts are not randomized (because no searching for duplicates).
- Searches are expected (but not amortized) since no resizing on a search.

Hashing overview

- What is a hash function?
- Collision resolution: chaining
- Java hashing
- Collision resolution: open addressing
- Table (re)sizing

Reality Fights Back

Simple Uniform Hashing doesn't exist.

- Keys are not random.
 - Lots of regularity.
 - Mysterious patterns.
- Patterns in keys can induce patterns in hash functions unless you are very careful.

Problem Hash Functions

Example:

- One bucket for each letter [a..z]
- Hash function: $h(\text{string}) = \text{first letter}$.
 - E.g., $h(\text{“hippopotamus”}) = \text{h}$.
- Bad hash function: why??

Problem Hash Functions

Example:

- One bucket for each letter [a..z]
- Hash function: $h(\text{string}) = \text{first letter}$.
 - E.g., $h(\text{“hippopotamus”}) = \text{h}$.
- Bad hash function: many fewer words start with the letter x than start with the letter s .

Problem Hash Functions

Example:

- One bucket for each number from $[1..26*28]$
- Hash function: $h(\text{string}) = \text{sum of the letters}$.
 - E.g., $h(\text{"hat"}) = 8 + 1 + 20 = 29$.
- Bad hash function: why??

Problem Hash Functions

Example:

- One bucket for each number from $[1..26*28]$
- Hash function: $h(\text{string}) = \text{sum of the letters}$.
 - E.g., $h(\text{"hat"}) = 8 + 1 + 20 = 29$.
- Bad hash function: lots of words collide, and you don't get a uniform distribution (since most words are short).

Problem Hash Functions

But pretty good hash functions do exist...

- Optimism pays off!

Moral of the story:

- Don't design your own hash functions.
- Ever.
- Unless you really need to.

Designing Hash Functions

Goal: find a hash function whose values *look* random.

- Similar to pseudorandom generators:
 - When you use Java random, there is no real randomness.
 - Instead, it generates a sequence of numbers that looks random.
- For every hash function, some set of keys is bad!
- If you know the keys in advance, you can choose a hash function that is always good!
 - But if you change the keys, then it might be bad again.

Designing Hash Functions

Two common hashing techniques...

- Division Method
- Multiplication Method

Designing Hash Functions

Division Method

- $h(k) = k \bmod m$
 - For example: if $m=7$, then $h(17) = 3$
 - For example: if $m=20$, then $h(100) = 0$
 - For example: if $m=20$, then $h(97) = 17$
- Two keys k_1 and k_2 collide when:
$$k_1 = k_2 \bmod m$$
- Collision unlikely if keys are random.

Designing Hash Functions

Division Method

- (Bad) idea: choose $m = 2^x$

Very fast to calculate $k \bmod m$ via shifts

Recall: $001001 \gg 1 = 00100$

$001001 \gg 2 = 0010$

$001001 \gg 3 = 001$

Designing Hash Functions

Division Method

- (Bad) Idea: choose $m = 2^x$

Very fast to calculate $k \bmod m$ via shifts:

$$k \bmod 2^x = k - ((k \gg x) \ll x)$$

Designing Hash Functions

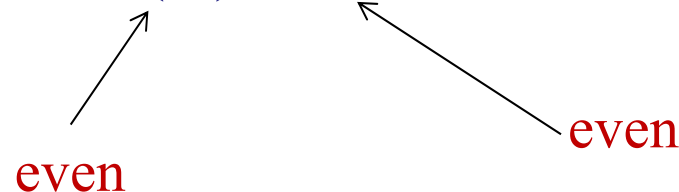
Division Method

- (Bad) Idea: choose $m = 2^x$

Very fast to calculate $k \bmod m$ via shifts

- Problem: Regularity
 - Input keys are often regular
 - Assume input keys are even.
 - Then $h(k) = k \bmod m$ is even!

$$k \bmod m + i(m) = k$$



Designing Hash Functions

Division Method

- Assume k and m have common divisor d .

$$k \bmod m + i * m = k$$

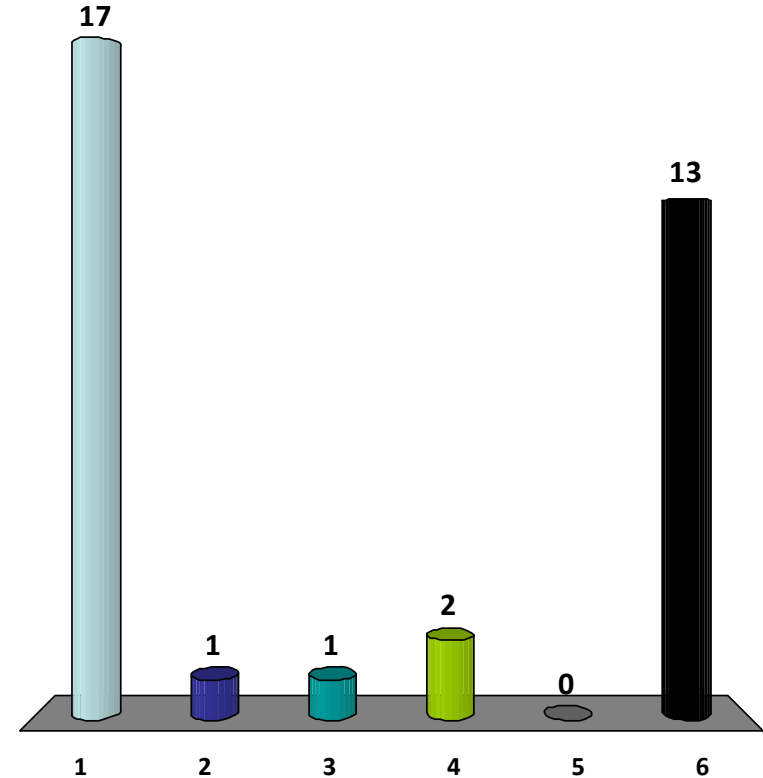


divisible by d

- Implies that $h(k) = k \bmod m$ is divisible by d .

If d is a divisor of m and every key k , then what percentage of the table is used?

- ✓ 1. $1/d$
- 2. $1/k$
- 3. $1/m$
- 4. d/n
- 5. m/n
- 6. d/m



Designing Hash Functions

Division Method

- Assume k and m have common divisor d .

$$k \bmod m + i * m = k$$

divisible by d

- Implies that $h(k)$ is divisible by d .
- If all keys are divisible by d , then you only use 1 out of every d slots

0	A
1	null
2	null
d = 3	B
4	null
5	null
2d = 6	C
7	null
8	null
3d = 9	D

Designing Hash Functions

Division Method

- $h(k) = k \bmod m$
- Choose $m =$ prime number
 - Not too close to a power of 2.
 - Not too close to a power of 10.
- Division method is popular (and easy), but not always the most effective.
- Division is slow.

Designing Hash Functions

Two common hashing techniques...

- Division Method
- Multiplication Method

Designing Hash Functions

Multiplication Method

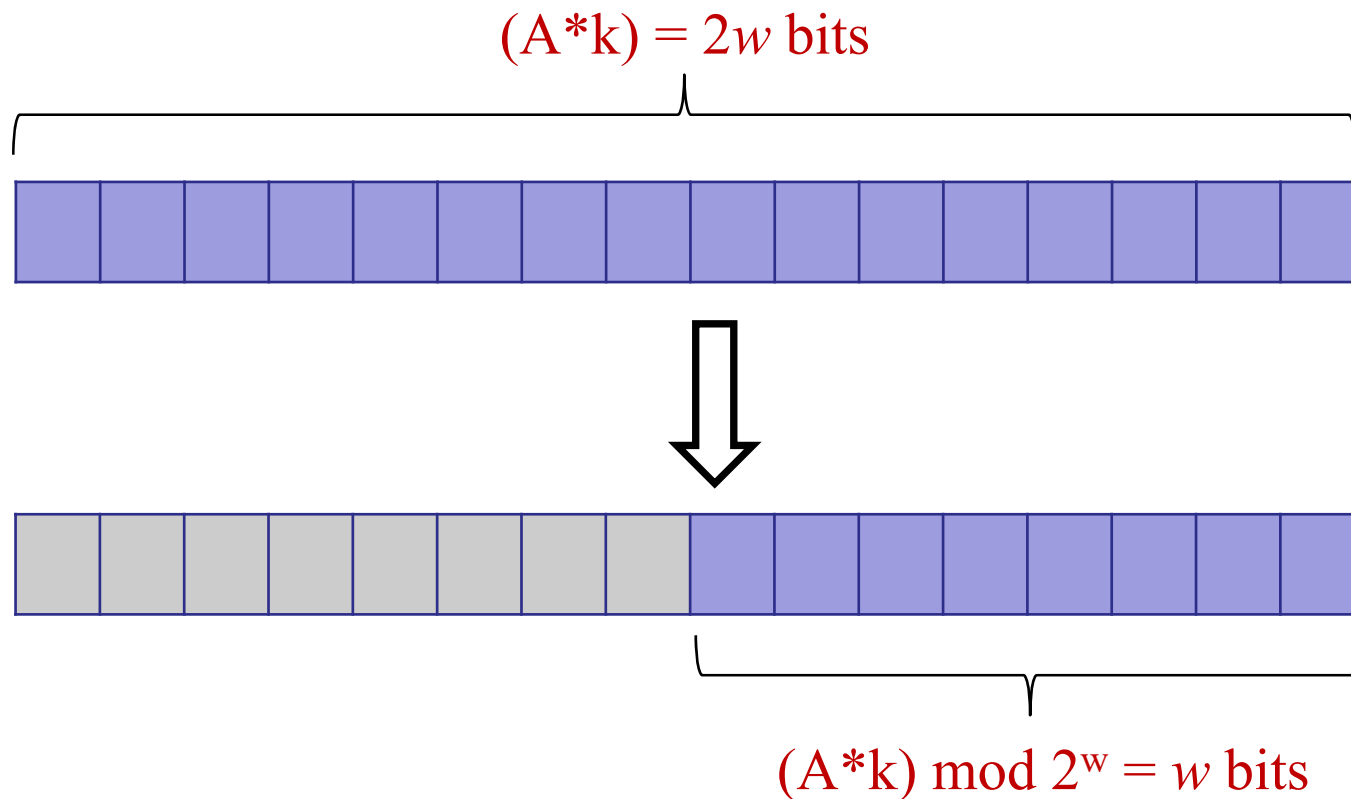
- Fix table size: $m = 2^r$, for some constant r .
- Fix word size: w , size of a key in bits.
- Fix (odd) constant A .

$$h(k) = (Ak) \bmod 2^w \gg (w - r)$$

Designing Hash Functions

Multiplication Method

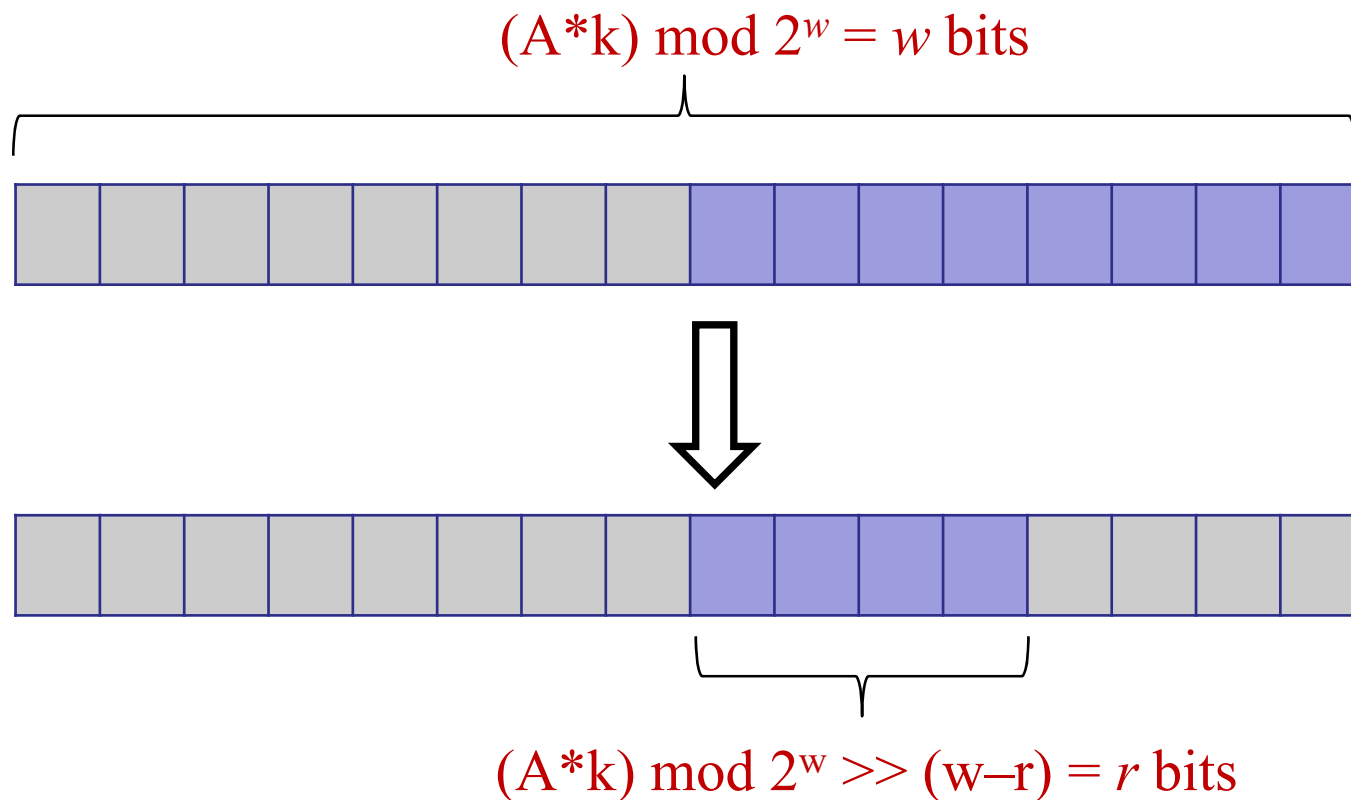
- Given m, w, r, A : $h(k) = (Ak) \bmod 2^w \gg (w - r)$



Designing Hash Functions

Multiplication Method

- Given m, w, r, A : $h(k) = (Ak) \bmod 2^w \gg (w - r)$



Designing Hash Functions

Multiplication Method

- Faster than Division Method
 - Multiplication, shifting faster than division
- Works reasonably well when A is an odd integer $> 2^{w-1}$
 - Odd: if it is even, then lose at least one bit's worth
 - Big enough: use all the bits in A .

Hashing overview

- What is a hash function?
- Collision resolution: chaining
- Java hashing
- Collision resolution: open addressing
- Table (re)sizing