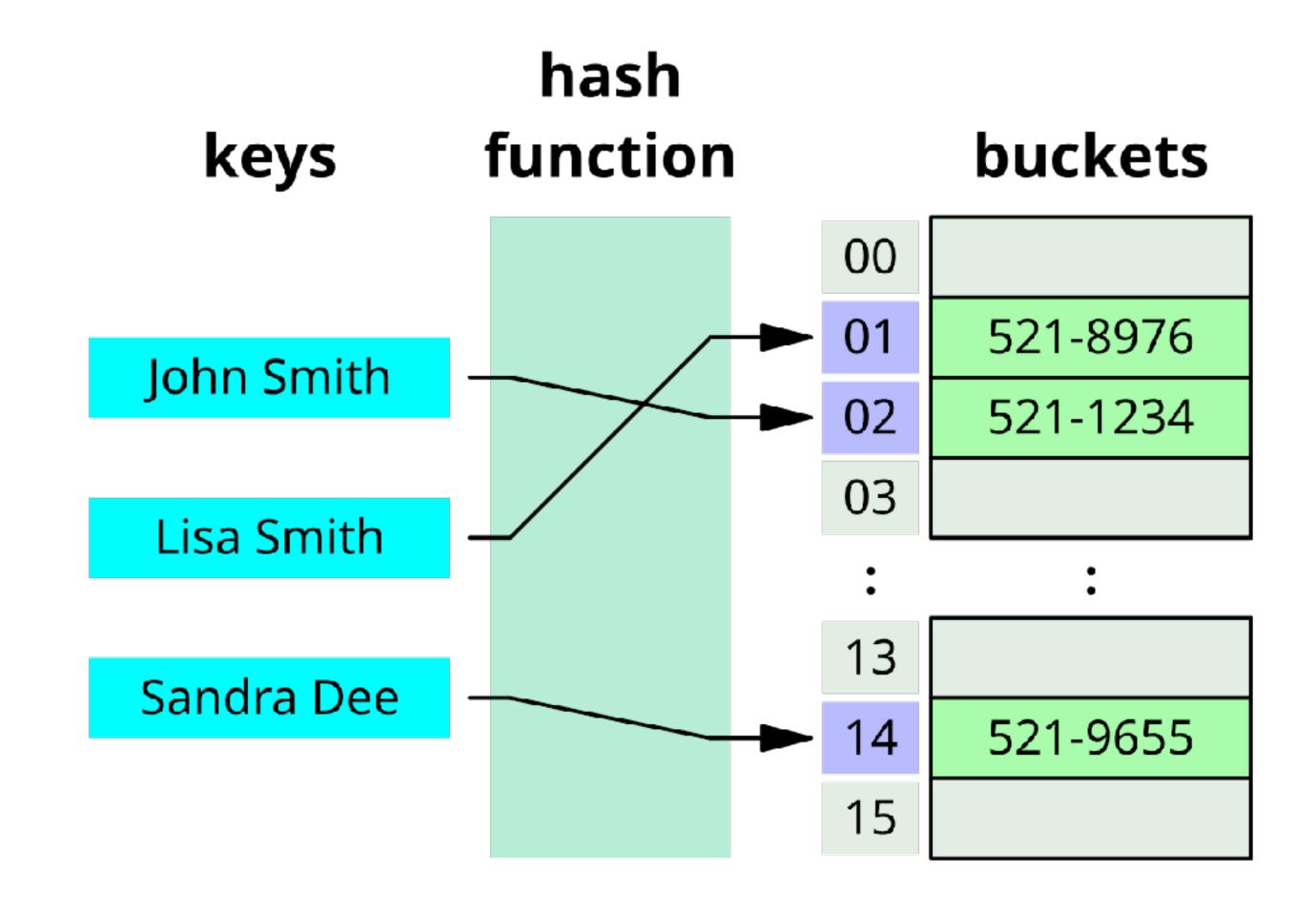
CS62 Class 18: Hash Tables (pt I)

Searching





Agenda

- Motivation
- Deriving hashtables: naive approach
- Separate chaining & handling collisions
- Hash table resizing

Hashtable motivation

We've now seen several implementations of maps...

		Worst case		_	Average cas	e	
	Search	Insert	Delete	Search	Insert	Delete	Notes
BST	n	n	n	log n	log n	\sqrt{n}	Random trees are log n
2-3 Tree	log n	log n	log n	log n	log n	log n	Beautiful idea, very hard to implement
LLRB	log n	log n	log n	log n	logn	log n	Bijection with 2-3 tree, hard to implement

Limits of Search Tree Based Maps (and Sets)

Our search tree based sets require items to be comparable.

- Need to be able to ask "is X < Y?" Not true of all types (ex. How do you compare $\bar{\tau}$ and 橙?).
- Could we somehow avoid the need for objects to be comparable?

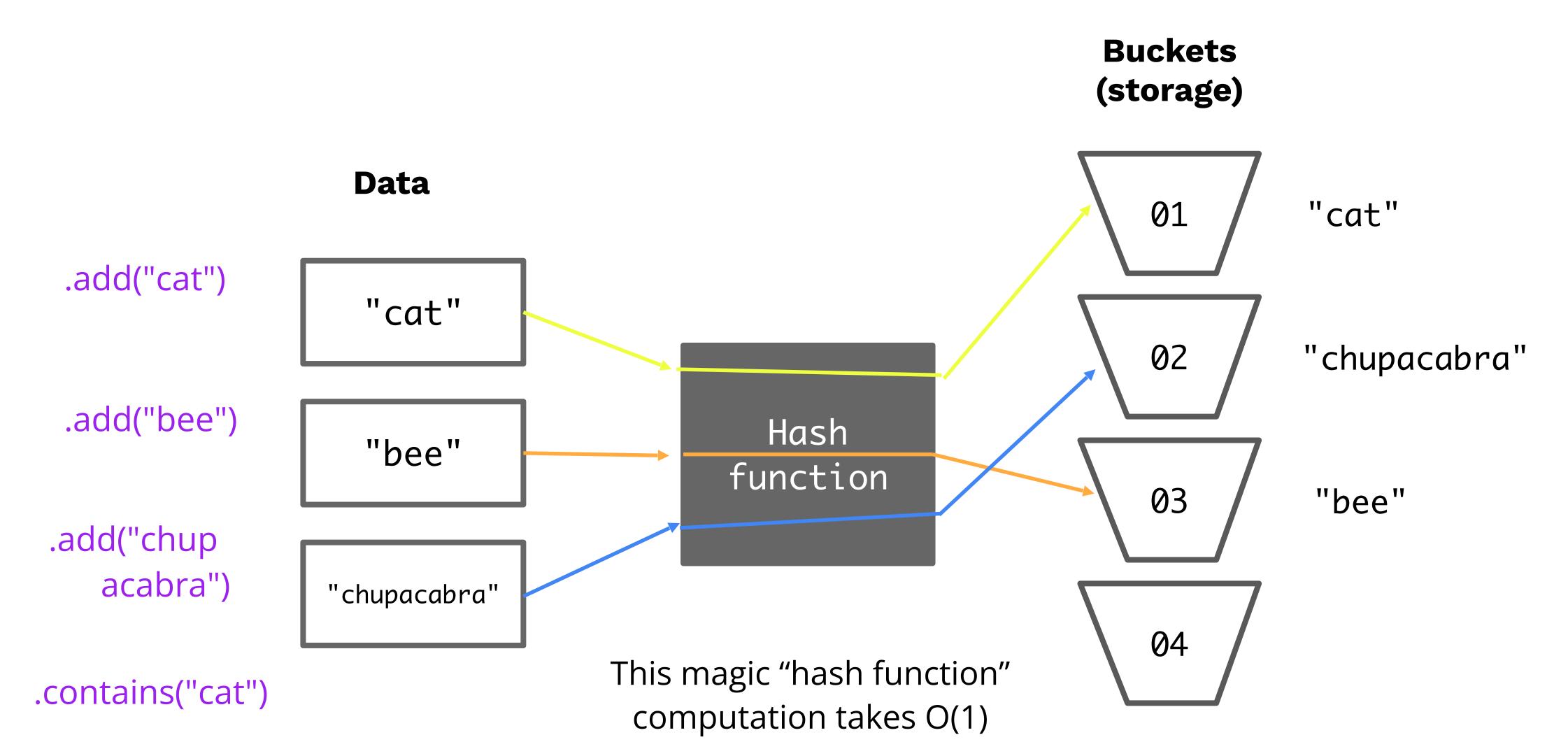
Our search tree sets have excellent performance, but could maybe be better?

- Θ(log N) is amazing. 1 billion items is still only height ~30.
- Could we somehow do better than Θ(log N)?

Today we'll see the answer to both of the questions above is yes.

Hash functions associate data with a storage bucket

Hash tables take data, transform them using a hash function into integer indices
for storage buckets. We'll be focusing on .add() and .contains() today.



Naive approach: Data indexed arrays

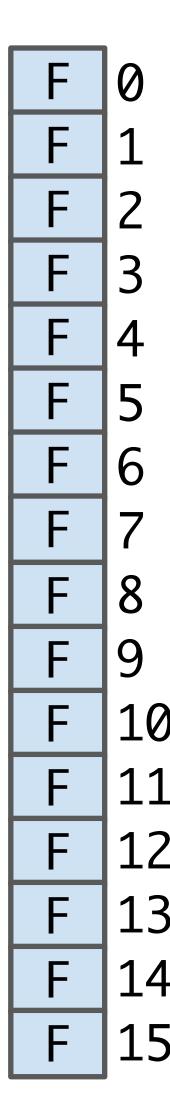
How would we implement extremely fast add/lookups?

One extreme approach: Create an array of booleans indexed by data. Data = integers, which is also the index of the array.

- Initially all values are false
- When an item is added, set appropriate index to true

```
DataIndexedIntegerSet diis;
diis = new DataIndexedIntegerSet();
```

Set containing nothing



Using Data as an Index

One extreme approach: Create an array of booleans indexed by data. Data = integers, which is also the index of the array.

- Initially all values are false
- When an item is added, set appropriate index to true

```
DataIndexedIntegerSet diis;
diis = new DataIndexedIntegerSet();
diis.add(0);
```

```
F
F
F
F
F
Set containing 0
F
```

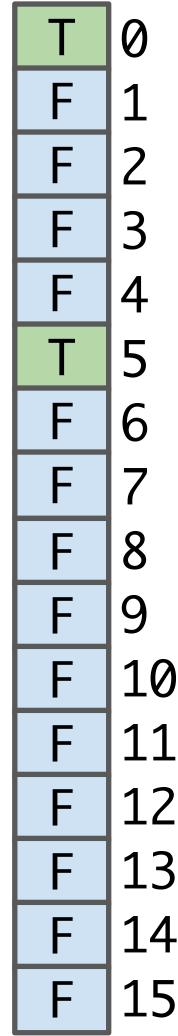
Using Data as an Index

One extreme approach: Create an array of booleans indexed by data. Data = integers, which is also the index of the array.

- Initially all values are false
- When an item is added, set appropriate index to true

```
DataIndexedIntegerSet diis;
diis = new DataIndexedIntegerSet();
diis.add(0);
diis.add(5);
```

```
Set containing 0, 5
```



Using Data as an Index

One extreme approach: Create an array of booleans indexed by data. Data = integers, which is also the index of the array.

- Initially all values are false
- When an item is added, set appropriate index to true

```
DataIndexedIntegerSet diis;
diis = new DataIndexedIntegerSet();
diis.add(0);
diis.add(5);
diis.add(11);
```

Set containing 0, 5, 11

DataIndexedIntegerSet implementation

```
public class DataIndexedIntegerSet {
   private boolean[] present;
   public DataIndexedIntegerSet() {
       present = new boolean [2000000000];
   public void add(int i) {
       present[i] = true;
   public boolean contains(int i) {
           return present[i];
```

Q: What are some downsides to this approach?

add() is a *constant* time operation: just set a flag to true

contains() is a constant time operation: just look up the value in the array

Set containing 0, 5, 11

10

. . .

DataIndexedIntegerSet downsides

```
public class DataIndexedIntegerSet {
   private boolean[] present;
   public DataIndexedIntegerSet() {
       present = new boolean [2000000000];
   public void add(int i) {
       present[i] = true;
   public boolean contains(int i) {
           return present[i];
```

- Extremely wasteful of memory
- Need some way to generalize beyond integers

Set containing 0, 5, 11

. . .

Generalizing to (English) Strings

aka Hashing

Generalizing the DataIndexedIntegerSet Idea

Suppose we want to add("cat")

The key question:

- What is the catth element of a list?
- One idea: Use the first letter of the word as an index.
 - a = 1, b = 2, c = 3, ..., z = 26



What's wrong with this approach?

- Other words start with c.
 - contains("chupacabra"): true
- Can't store ":3"

1 F a b T c d

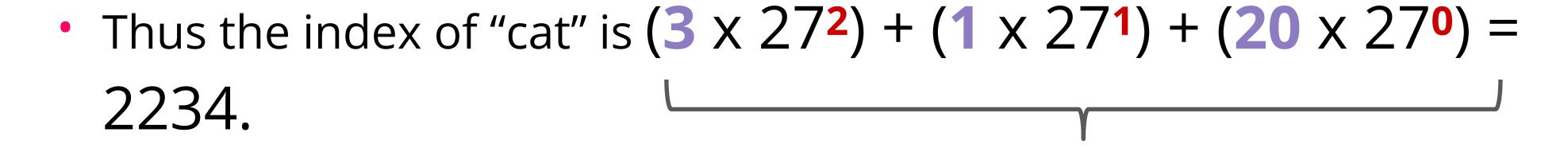
25 F 3

we say that "chupacabra" **collides**with "cat"

Avoiding Collisions

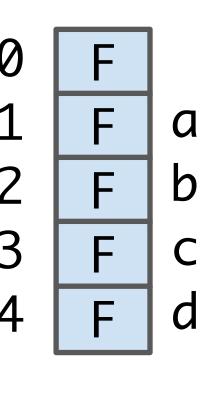
Here's an idea: To get the new index, we will use all the letters by multiplying each by a power of 27.

•
$$a = 1$$
, $b = 2$, $c = 3$, ..., $z = 26$



Why this specific pattern?

Let's review how numbers are represented in decimal.



- - -

. . .

- - -

The Decimal Number System vs. Our System for Strings

In the decimal number system, we have 10 digits: 0, 1, 2, 3, 4, 5, 6, 7, 8, 9

Want numbers larger than 9? Use a sequence of digits.

Example: 7091 in base 10

•
$$7091_{10} = (7 \times 10^3) + (0 \times 10^2) + (9 \times 10^1) + (1 \times 10^0)$$

Our system for strings is almost the same, but with letters (base 27, we don't use 0).

•
$$cat_{27} = (3 \times 27^2) + (1 \times 27^1) + (20 \times 27^0) = 2234_{10}$$

This is the beginnings of a hash function.

Worksheet time!

Convert the word "bee" into a number by using our "powers of 27" strategy. That is, hash "bee".

Reminder: $Cat_{27} = (3 \times 27^2) + (1 \times 27^1) + (20 \times 27^0) = 2234_{10}$

Hint: 'b' is letter 2, and 'e' is letter 5. And $27^2 = 729$ (but feel free to use a calculator)

Worksheet answers

Convert the word "bee" into a number by using our "powers of 27" strategy.

Reminder:
$$Cat_{27} = (3 \times 27^2) + (1 \times 27^1) + (20 \times 27^0) = 2234_{10}$$

Hint: 'b' is letter 2, and 'e' is letter 5. And $27^2 = 729$

• bee₂₇ =
$$(2 \times 27^2) + (5 \times 27^1) + (5 \times 27^0) =$$

 $(1458) + (135) + (5) = 1598_{10}$

Avoiding collisions with uniqueness

- $cat_{27} = (3 \times 27^2) + (1 \times 27^1) + (20 \times 27^0) = 2234_{10}$
- bee₂₇ = $(2 \times 27^2) + (5 \times 27^1) + (5 \times 27^0) = 1598_{10}$

As long as we pick a base ≥ 26, this algorithm is guaranteed to give each lowercase English word a unique number

Using base 27, no other words will get the number 1598.

This is an example of a *hash function* for hashing lowercase English words where it's guaranteed that we will *never have a collision*.

(Practice problem: write an englishToInt hash function that will automatically compute the hash for any input english word.)

DataIndexedEnglishWordSet Implementation

```
public class DataIndexedEnglishWordSet {
   private boolean[] present;
   public DataIndexedEnglishWordSet() {
       present = new boolean [20000000000];
   public void add(String s) {
       present[englishToInt(s)] = true;
   public boolean contains(String s) {
           return present[englishToInt(s)];
```

```
We've solved the
first problem of
generalizing from
integers through
introducing a hash
function
englishToInt(s)
public void add(int i) {
   present[i] = true;
                                 cas
                     2233
 ^ old way
                                 cat
                                 cau
```

Set containing "cat"

Generalizing to any String

DataIndexedStringSet

Using only lowercase English characters is too restrictive.

- What if we want to store strings like "2pac" or "eGg!"?
- To understand what value we need to use for our base, let's briefly discuss the ASCII standard.

ASCII Characters

The most basic character set used by most computers is ASCII format.

- Each possible character (128 total) is assigned a value between 0 and 127.
- Characters 33 126 are "printable", and are shown below.
- For example, char c = 'D' is equivalent to char c = 68.

33		49	1	65	Α	81	Q	97	а	113	q
34	"	50	2	66	В	82	R	98	b	114	r
35	#	51	3	67	С	83	S	99	С	115	S
36	\$	52	4	68	D	84	Т	100	d	116	t
37	%	53	5	69	Е	85	U	101	. е	117	u
38	&	54	6	70	F	86	V	102	. f	118	V
39	'	55	7	71	G	87	W	103	g	119	W
40	(56	8	72	Н	88	X	104	- h	120	X
41)	57	9	73	Ι	89	Υ	105	i i	121	У
42	*	58	:	74	J	90	Z	106	j	122	Z
43	+	59	;	75	K	91	[107	' k	123	{
44	,	60	<	76	L	92	\	108	3	124	
45	-	61	=	77	М	93]	109) m	125	}
46		62	>	78	Ν	94	^	110	n	126	~
47	/	63	?	79	0	95	_	111	. 0		
48	0	64	@	80	Р	96	`	112	. p		

DataIndexedStringSet

Using only lowercase English characters is too restrictive.

- What if we want to store strings like "2pac" or "eGg!"?
- Let's use the ASCII standard as an encoding, and take 126 as our base.

Examples:

- bee₁₂₆= $(98 \times 126^2) + (101 \times 126^1) + (101 \times 126^0) = 1,568,675$
- $2pac_{126} = (50 \times 126^3) + (112 \times 126^2) + (97 \times 126^1) + (99 \times 126^0) = 101,809,233$
- $eGg!_{126} = (98 \times 126^3) + (71 \times 126^2) + (98 \times 126^1) + (33 \times 126^0) = 203,178,213$

Implementing the hash function asciiToInt

```
public static int asciiToInt(String s) {
    int intRep = 0;
    for (int i = 0; i < s.length(); i += 1) {
        intRep = intRep * 126;
        intRep = intRep + s.charAt(i);
        }
        return intRep;
}</pre>
Strings are composed of chars, so they
automatically take their
ASCII value in math operations
```

Finally, if you want to use characters beyond ASCII, you can use the **Unicode** encoding. This supports, for instance, characters like 员, **个**, and الطبيعة.

Another problem: Integer Overflow

In Java, the largest possible integer is 2,147,483,647.

- If you go over this limit, you overflow, starting back over at the smallest integer, which is -2,147,483,648.
- In other words, the next number after 2,147,483,647 is -2,147,483,648.

```
int x = 2147483647;
System.out.println(x);
System.out.println(x + 1);
```

```
$ javac BiggestPlusOne.java
$ java BiggestPlusOne
2147483647
-2147483648
```

Consequence of Overflow: Collisions for long words

Because Java has a maximum integer, we won't get the numbers we expect!

- With base 126, we will run into overflow even for short strings.
 - Example: omens₁₂₆= 28,196,917,171, which is much greater than the maximum integer (28 billion versus 2 billion)!
 - asciiToInt('omens') will give us -1,867,853,901 instead due to overflow.
- Overflow can lead to collisions, resulting in wrong answers.

```
public static void main(String[] args) {
  DataIndexedStringSet disi = new DataIndexedStringSet();
  disi.add("melt banana");
  disi.contains("subterrestrial anticosmetic");
  //asciiToInt for these strings is both 839099497
}
returns true!
```

Hash Codes and the Pigeonhole Principle

The official term for the number we're computing is a "hash code", which is the result of a hash function.

- Via Wolfram Alpha: a hash code "projects a value from a set with many (or even an infinite number of) members to a value from a set with a fixed number of (fewer) members."
- Here, our target set is the set of Java integers, which is of size 4,294,967,296 (both negative and positive integers).

Pigeonhole principle tells us that if there are more than 4294967296 possible items, multiple items will share the same hash code.

- There are more than 4294967296 strings.
 - "one", "two", ... "nineteen quadrillion", ...

Bottom line: Collisions are inevitable.



Hash Tables

Two Fundamental Challenges of Hash Tables

How do we resolve hash code collisions?

- How do we compute a hash code for arbitrary objects?
 - Example: Our hash code for "cat" was 2234.
 - For Strings, this was relatively straightforward (treat as a base 126 number).
 - What about for Class PomonaStudent objects? What about for lists?

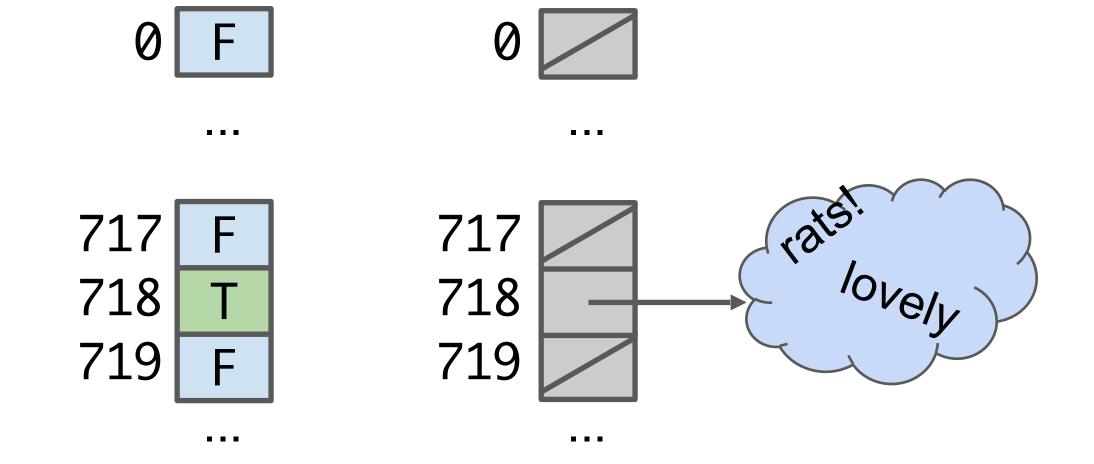
Solution for collisions: buckets

Suppose *N* items have the same numerical representation *h*:

- Example: hash code for "rats!" and "lovely" might both be 718.
- Instead of storing true in position h, store a "bucket" of these N items at position h.

How to implement a "bucket"?

- Conceptually simplest way: Singly Linked List.
- Could also use ArrayLists.
- Will see it doesn't really matter what you do.

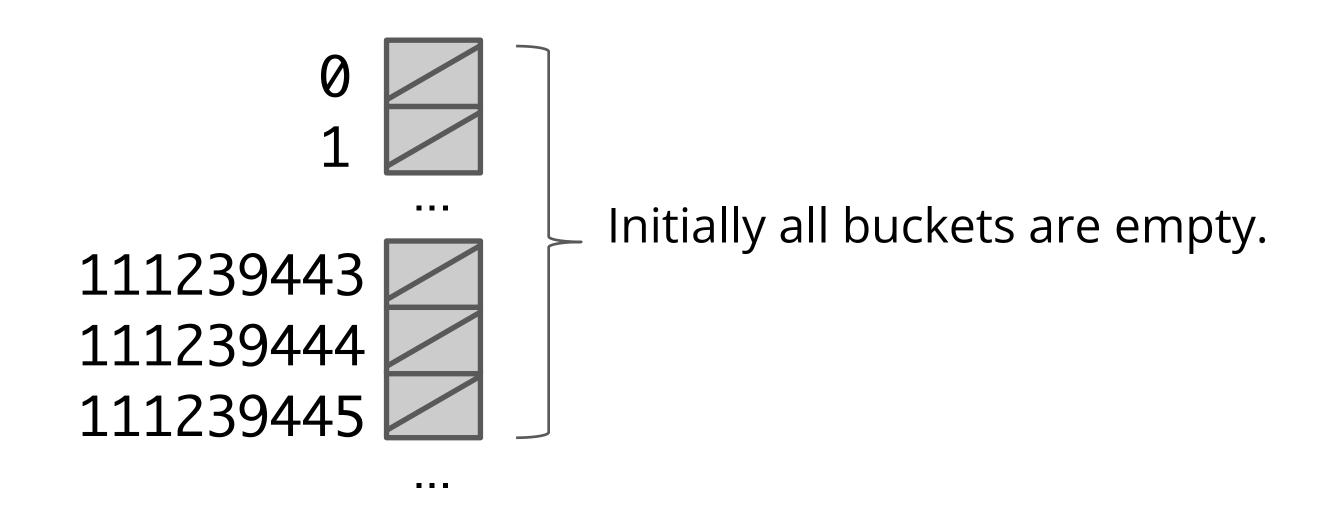


Each bucket in our array is initially empty. When an item x gets added at index h:

- If bucket h is empty, we create a new list containing x and store it at index h.
- If bucket *h* is already a list, we add *x* to this list if it is not already present.

We might call this a "separate chaining data indexed array".

• Bucket #h is a "separate chain" of all items that have hash code h.



Each bucket in our array is initially empty. When an item x gets added at index h:

- If bucket h is empty, we create a new list containing x and store it at index h.
- If bucket *h* is already a list, we add *x* to this list if it is not already present.

We might call this a "separate chaining data indexed array".

• Bucket #*h* is a "separate chain" of all items that have hash code *h*.

add("a")

Bucket 1 now has a length 1 list

111239443
111239444
111239445

Each bucket in our array is initially empty. When an item x gets added at index h:

- If bucket h is empty, we create a new list containing x and store it at index h.
- If bucket *h* is already a list, we add *x* to this list if it is not already present.

We might call this a "separate chaining data indexed array".

• Bucket #*h* is a "separate chain" of all items that have hash code *h*.

```
add("a")
add("chupacabra")

111239443
111239444
111239445
```

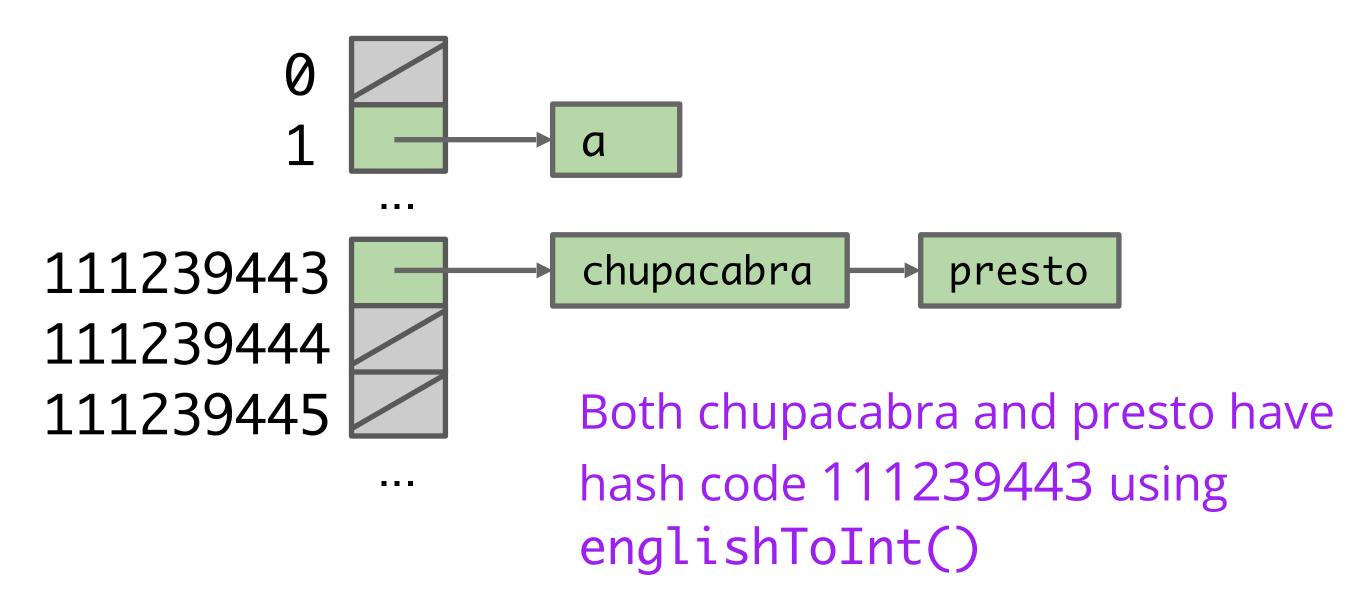
Each bucket in our array is initially empty. When an item x gets added at index h:

- If bucket h is empty, we create a new list containing x and store it at index h.
- If bucket *h* is already a list, we add *x* to this list if it is not already present.

We might call this a "separate chaining data indexed array".

• Bucket #*h* is a "separate chain" of all items that have hash code *h*.

```
add("a")
add("chupacabra")
add("presto")
```



"Separate chaining data indexed array"

Each bucket in our array is initially empty. When an item x gets added at index h:

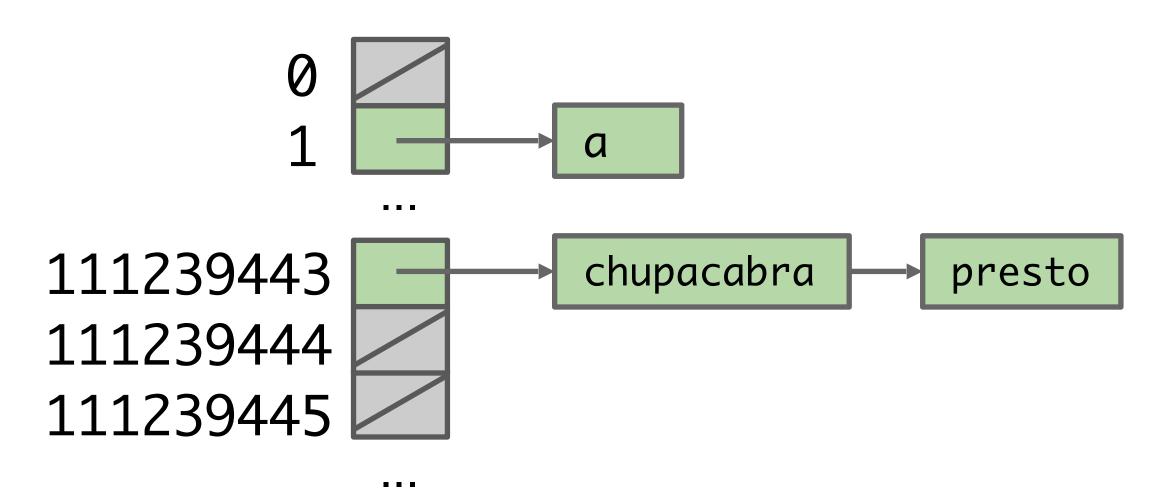
- If bucket h is empty, we create a new list containing x and store it at index h.
- If bucket *h* is already a list, we add *x* to this list if **it is not already present.**Note, if we were storing key/value pairs instead, we would update the value instead.

We might call this a "separate chaining data indexed array".

• Bucket #h is a "separate chain" of all items that have hash code h.

```
add("a")
add("chupacabra")
add("presto")
add("chupacabra")

Doesn't do anything, because
chupacabra is already present
```



"Separate chaining data indexed array"

Each bucket in our array is initially empty. When an item x gets added at index h:

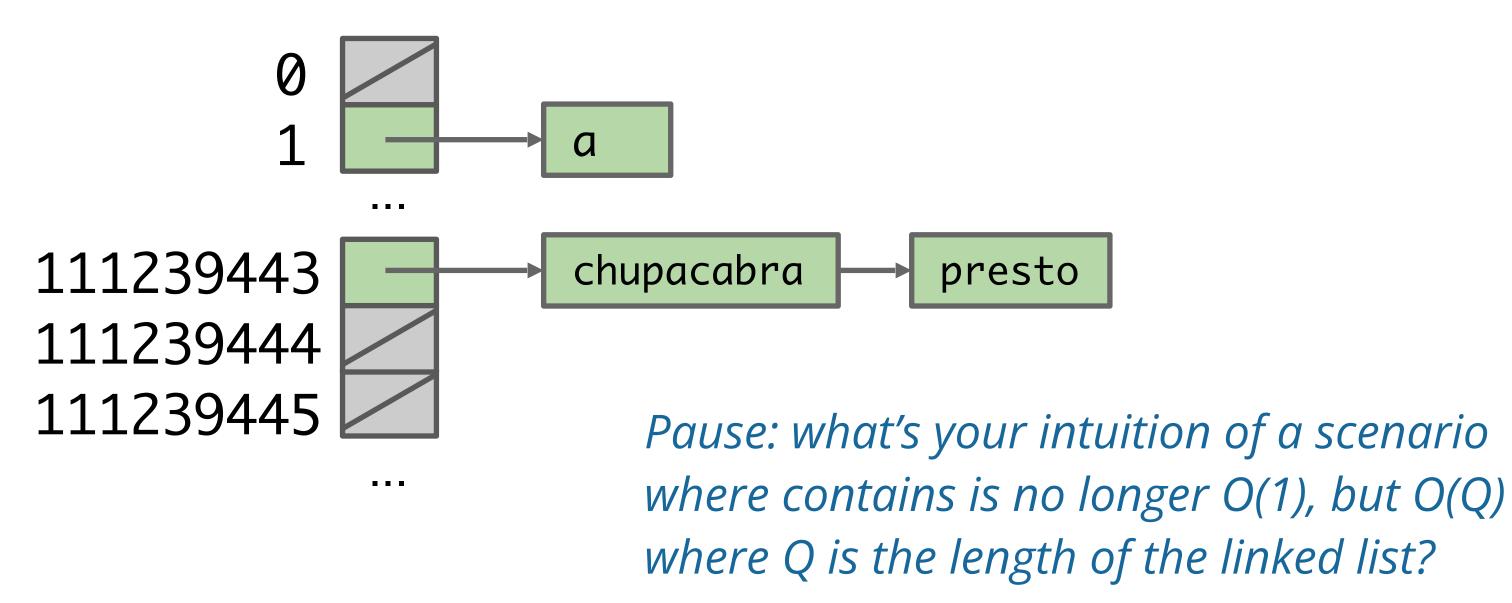
- If bucket *h* is empty, we create a new list containing *x* and store it at index *h*.
- If bucket h is already a list, we add x to this list if it is not already present.

We might call this a "separate chaining data indexed array".

• Bucket #*h* is a "separate chain" of all items that have hash code *h*.

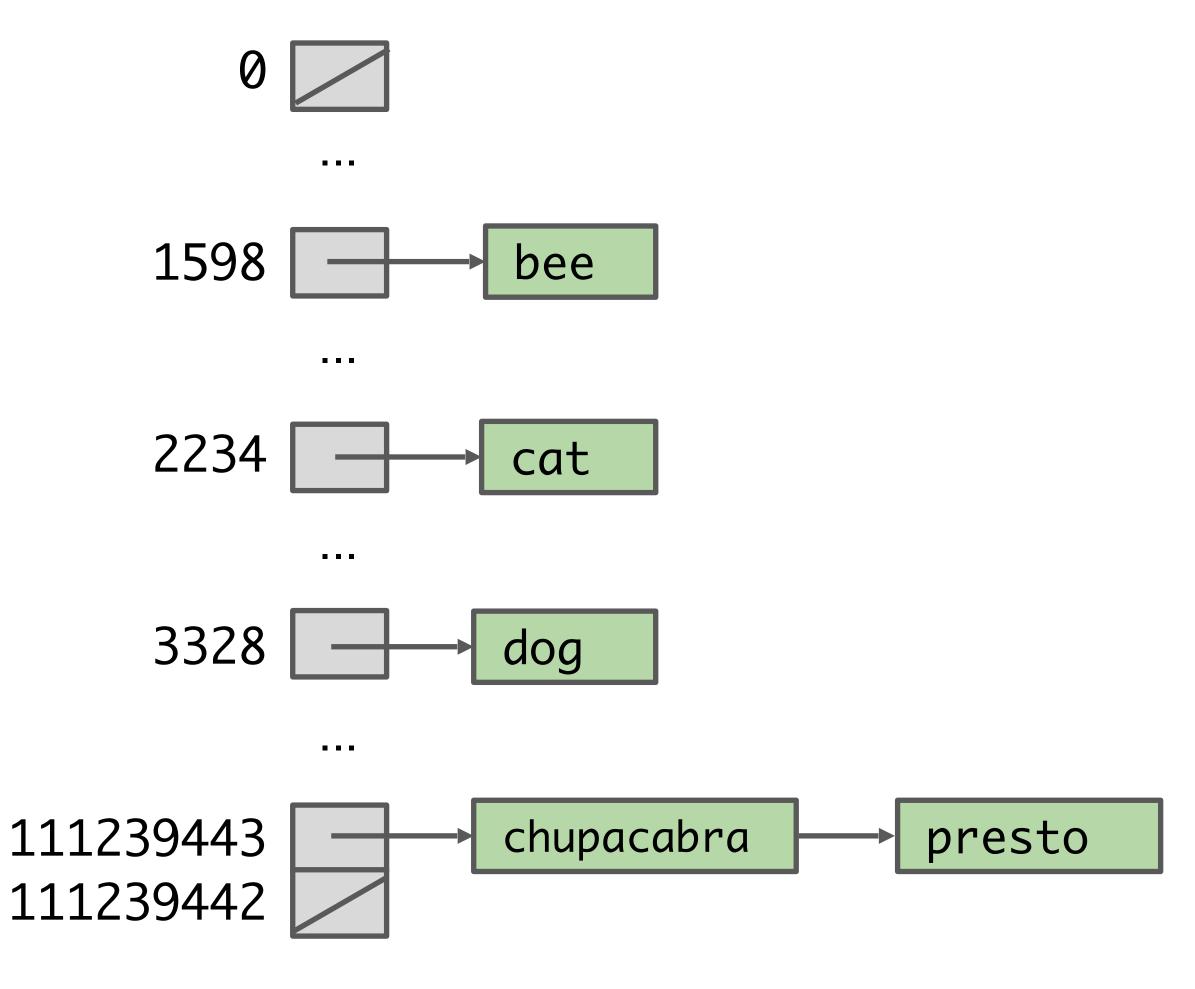
```
add("a")
add("chupacabra")
add("presto")
add("chupacabra")
contains("presto")
```

Traverse the list in bucket 111239443 to see if presto exists

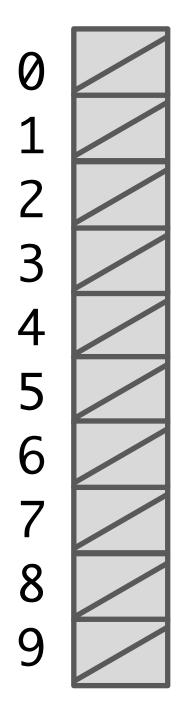


Saving Memory Using Separate Chaining

Observation: We don't really need billions of buckets.

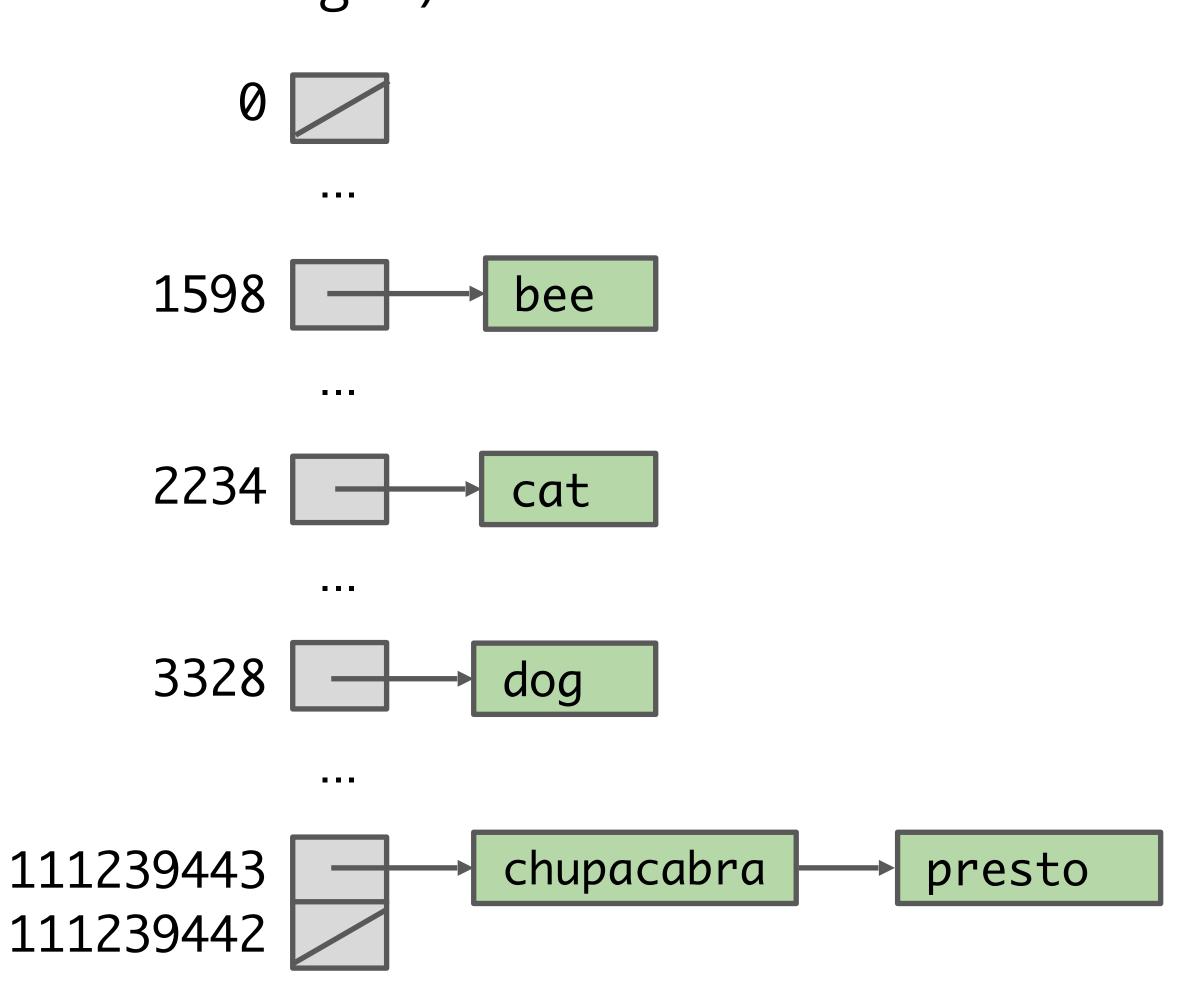


Q: If we use the 10 buckets on the right, where should our five items go?



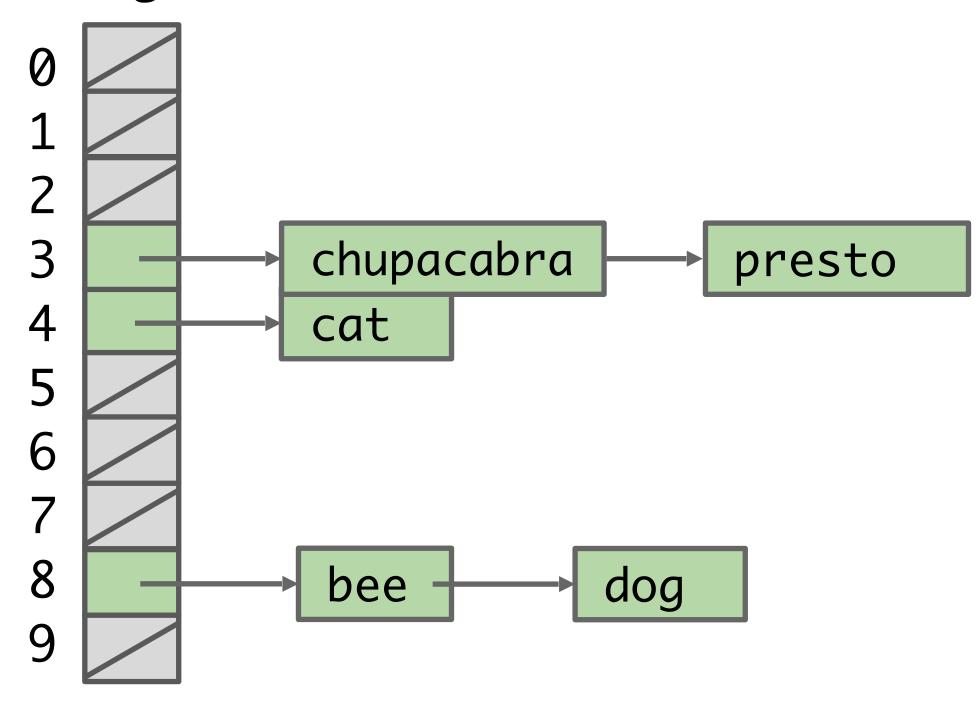
Saving Memory Using Separate Chaining and modulus

We can use the modulus operator to reduce bucket count. (Downside is that the lists will be longer.)



Q: If we use the 10 buckets on the right, where should our five items go?

A: Try bucket = hash code % 10 (i.e., look at the last digit of the hash code)



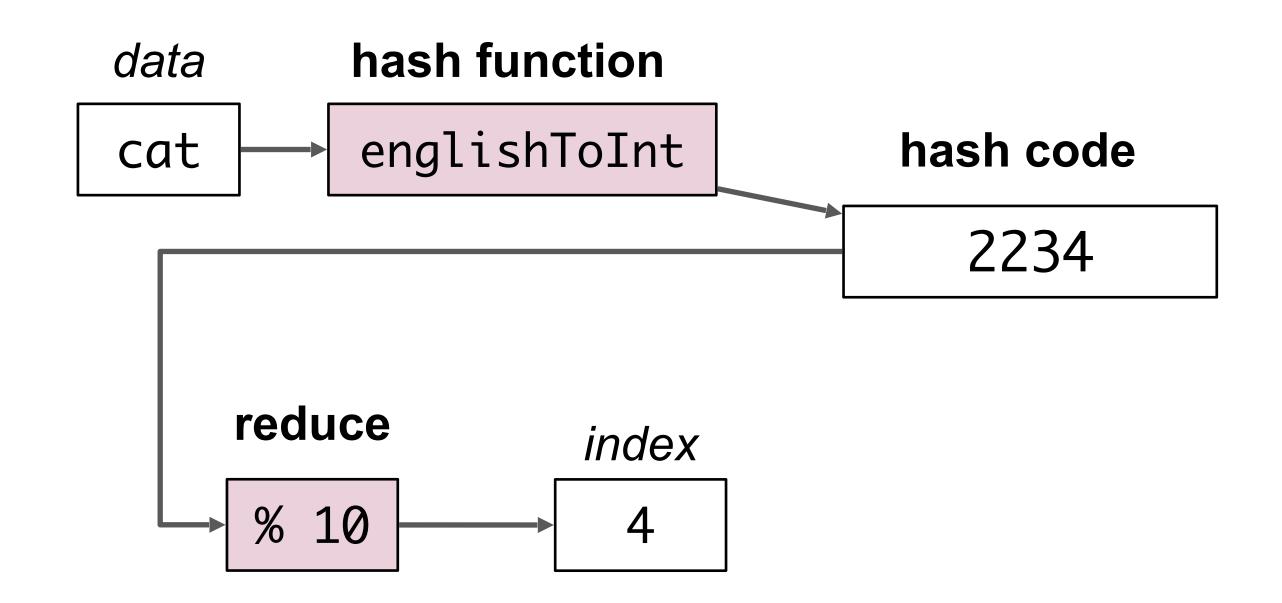
Finally: The Hash Table

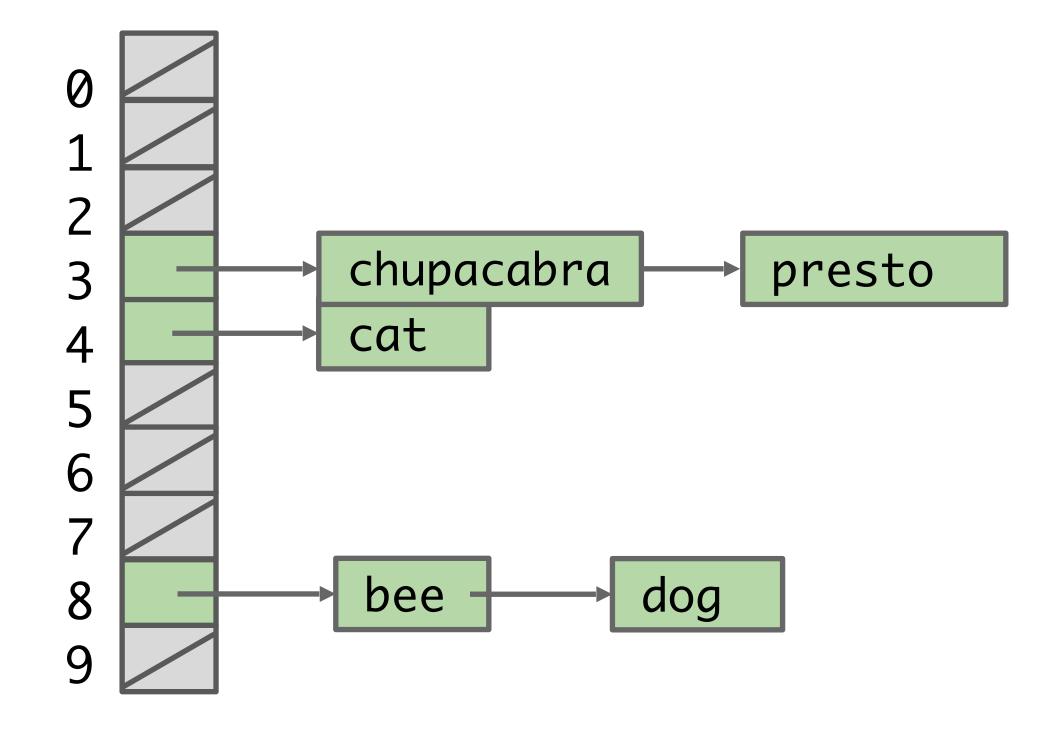
What we've just created here is called a hash table.

• **Data** is converted by a hash function into an integer representation called a hash code.

• The hash code is then reduced to a valid *index*, usually using the modulus

operator, e.g. 2348762878 % 10 = 8.



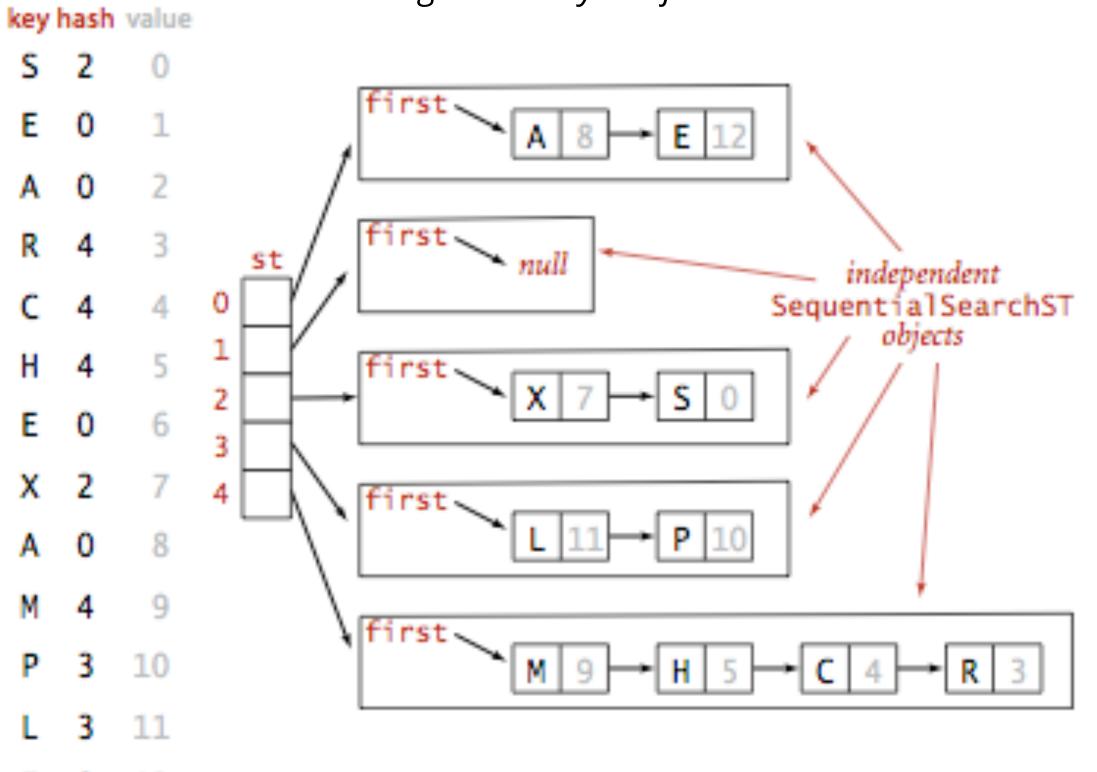


A hash table!

Separate/External Chaining (Closed Addressing)

- The formal name for what we've learned is called "separate chaining", "external chaining", or "closed addressing". (Why can't computer scientists just give one concept one name?)
- Use an array of m < n distinct linked lists (chains) m = # of buckets/chains in a hashtable [H.P. Luhn, IBM 1953].
 - Hash: Map key to integer i between 0 and m-1.
 - Insert: Put key-value pair **at front** of i-th chain (if not already there in which case we only update the associated value).
 - Search: Need to only search the i-th chain.

Note: In our textbook example, we store key-value pairs, while we've just been talking about keys so far in lecture



Hashing with separate chaining for standard indexing client

Note: in lecture, we saw putting it at the back of the chain, but the textbook example puts it at the front. (Why?) (More recent data quicker to access.)

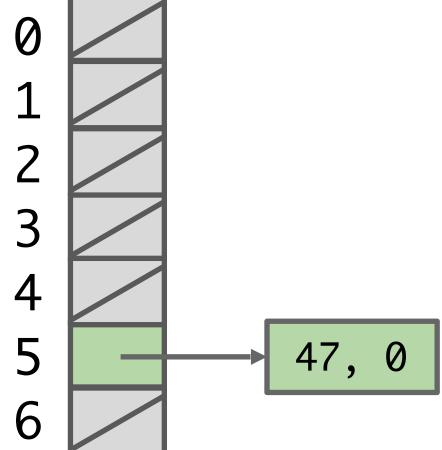
Worksheet time!

- Assume a dictionary implemented using hashing and separate chaining for handling collisions.
- Let m = 7 be the hash table size.
- For simplicity, we will assume that keys are integers and that the hash value for each key k is calculated as h(k) = k % m.

Insert the key-value pairs (47, 0), (3, 1), (28, 2), (14, 3), (9,4), (47,5) and show the resulting hash table.

As an example, for (47, 0)

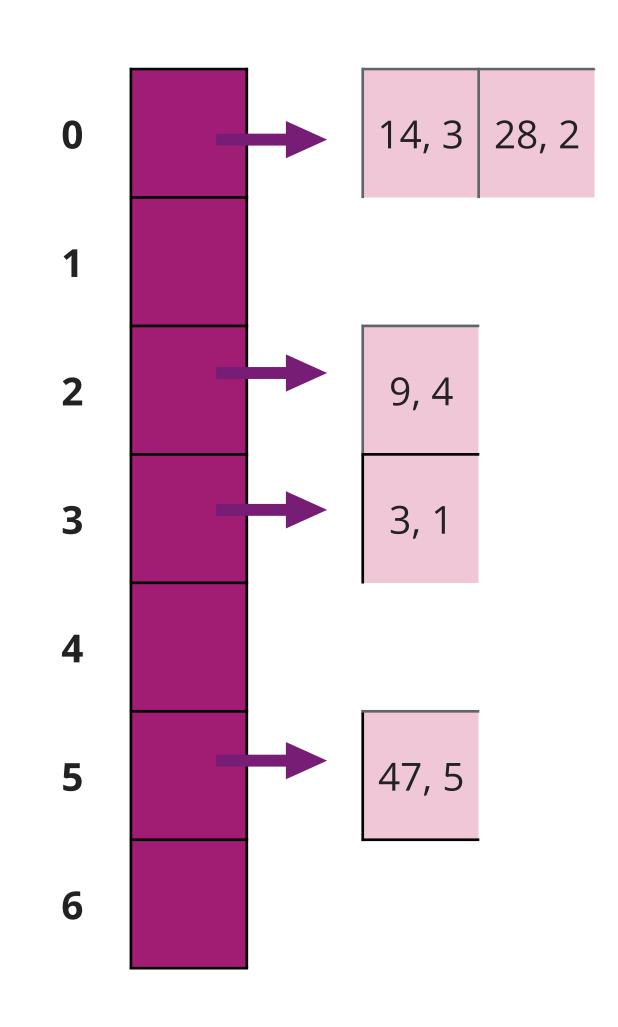
Hash value of 47 = 47 % 7 = 5



store key/value pairs!

Worksheet answers

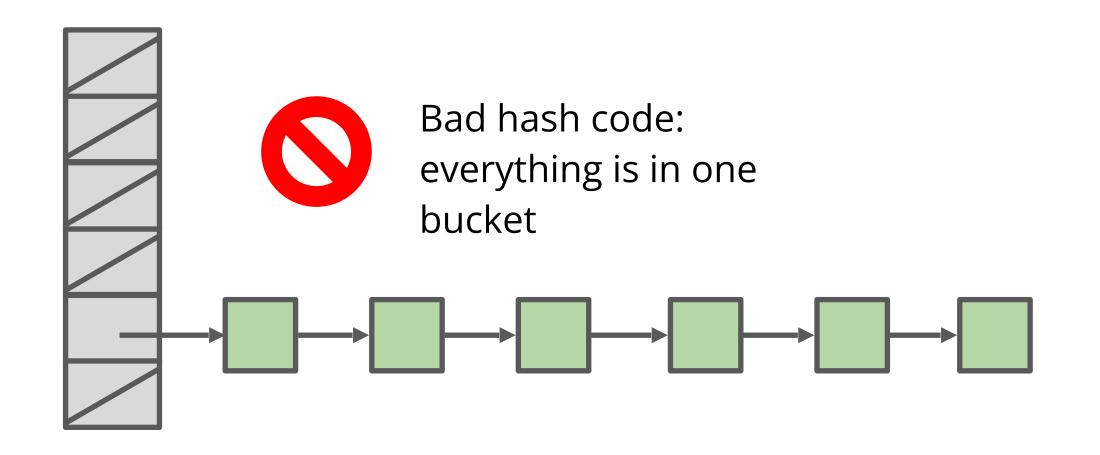
Key	Hash	Value
47	5	0
3	3	1
28	0	2
14	0	3
9	2	4
47	5	5

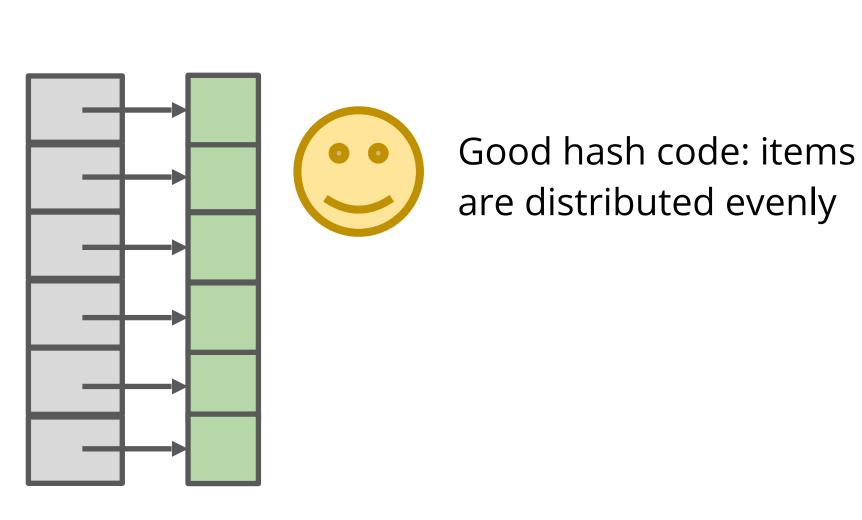


Hash Table Resizing

What makes a good hash function?

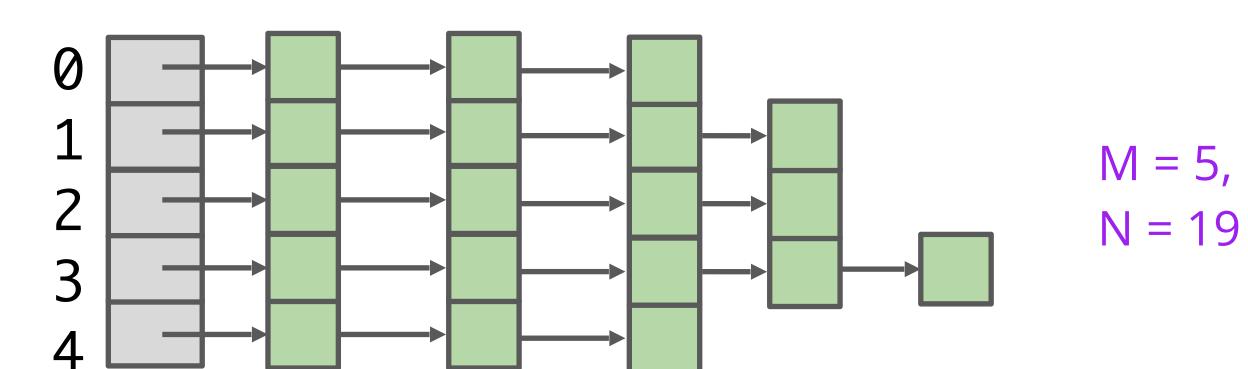
- We want a hash function that spreads things out nicely on real data.
 - Example #1: return 0 is a bad hash code function.
 - Example #2: just returning the first character of a word, e.g. "cat" \rightarrow 3 was also a bad hash function.
 - Example #3: adding chars together is bad. "ab" collides with "ba".
 - Example #4: returning string treated as a base B number can be good!
- A good hash function is hard to write, but it should scramble data seemingly randomly, so they will be evenly distributed over the hash table.







Improving the hash table



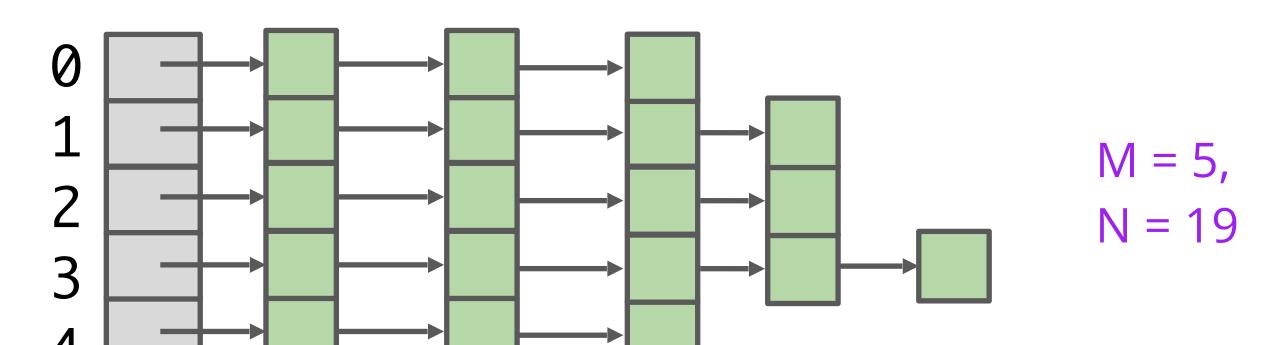
Suppose we have:

- A fixed number of buckets M.
- An increasing number of items N.

Major problem: Even if items are spread out evenly, lists are of length Q = N/M.

- For M = 5, that means the list length Q will scale to the number of items N, which results in linear time operations.
 - The best case is all items are evenly distributed, so Q is N/5. The worst case is all the items are in one bucket, so Q is N.
- Our goal: How can we improve our design to guarantee constant time operations? In other words, how can we make N/M = O(1)?

Improving the hash table via resizing

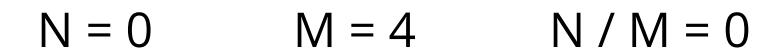


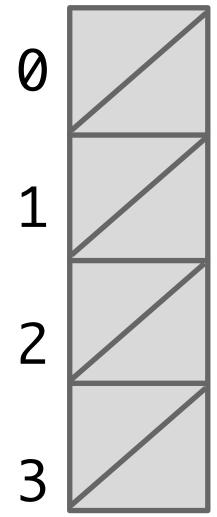
Suppose we have:

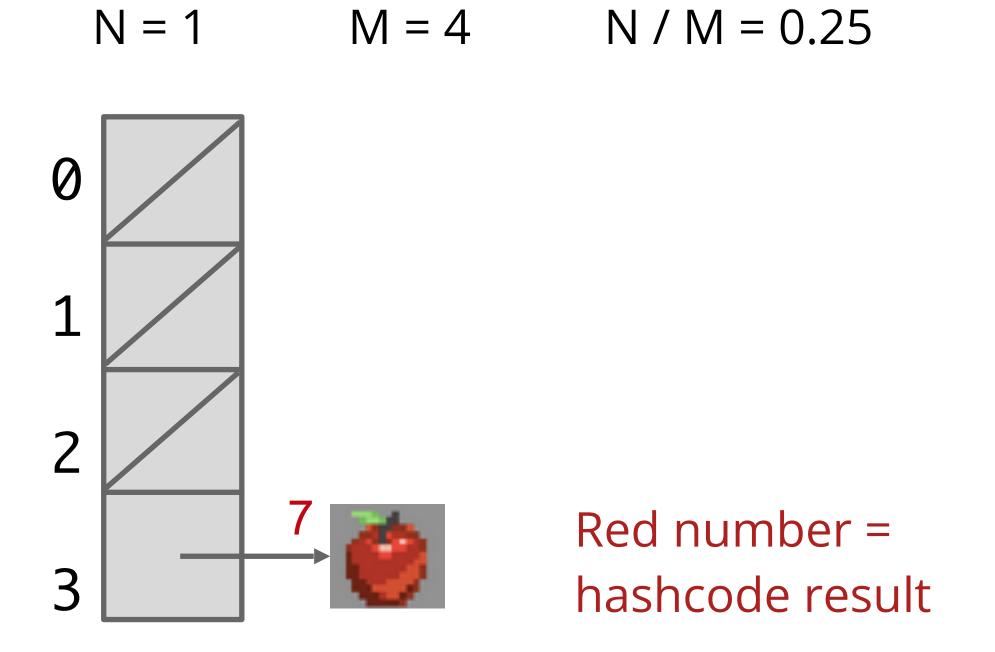
- An increasing number of buckets M.
- An increasing number of items N.

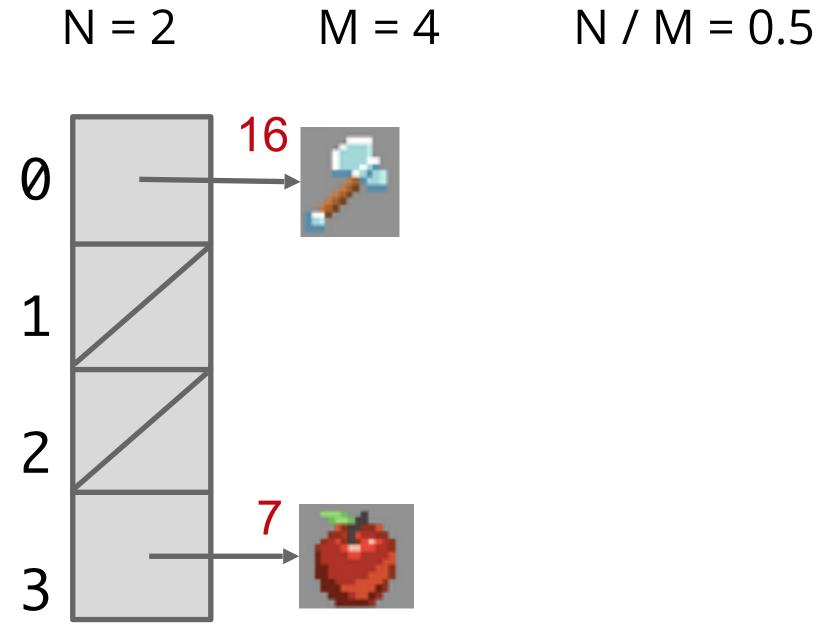
Major problem: Even if items are spread out evenly, lists are of length Q = N/M.

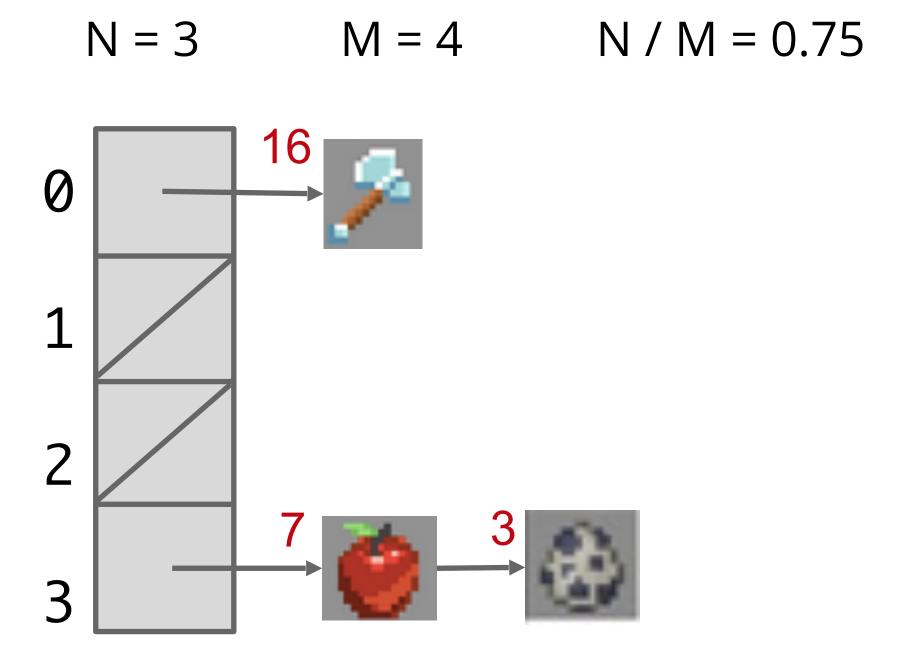
- We can resize the hash table and increase the number of buckets to scale to the number of items.
 - For example, if we want an average of 1.5 (a constant) number of items in our lists, then N/M = 1.5, so here, instead we would have M = 1.5 buckets for our N = 1.5 items.
- N/M is called the "load factor" how "full" the hash table is.

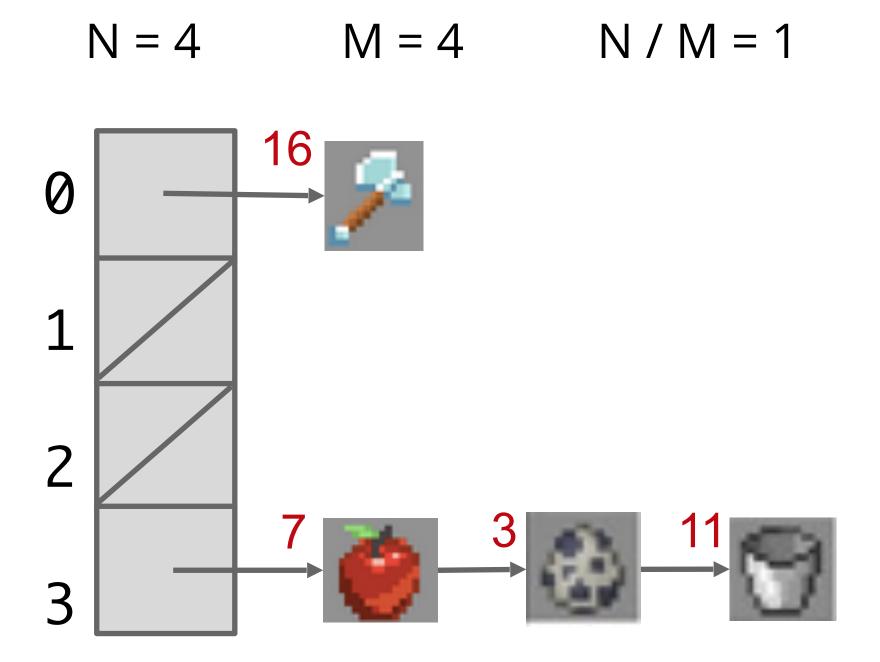


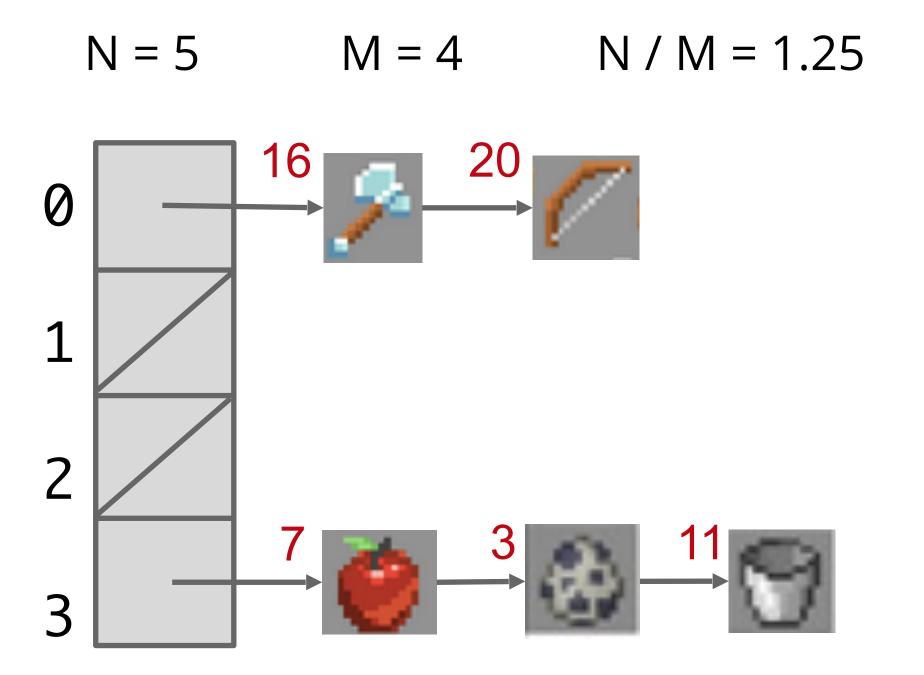




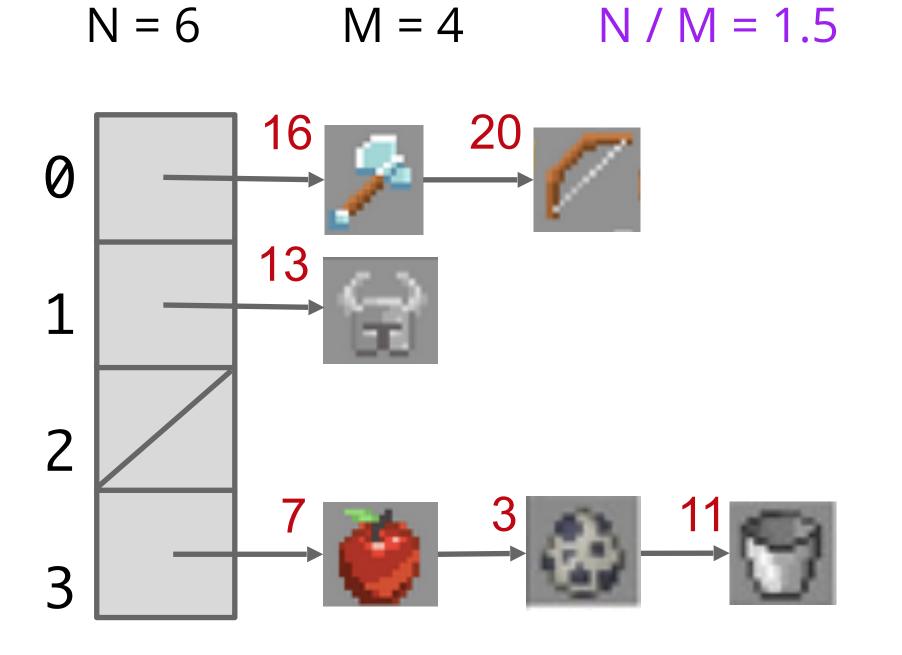








When N/M is \geq 1.5, then double M.



N/M is too large.
Time to double the number of buckets!

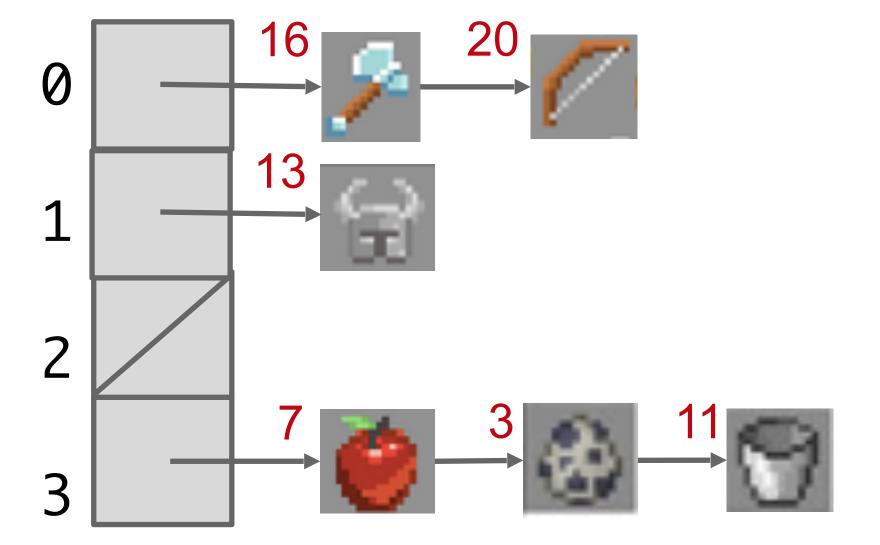
Worksheet time!

Where do all the existing elements go on the new hashtable? Assume we are placing items at the end, instead of beginning, of the list.

(Hint: how do we reduce hashes to bucket indices? How does that change with resize?)

When N/M is \geq 1.5, then double M.

N = 6 M = 4 N / M = 1.5



?

. | ?

2 ?

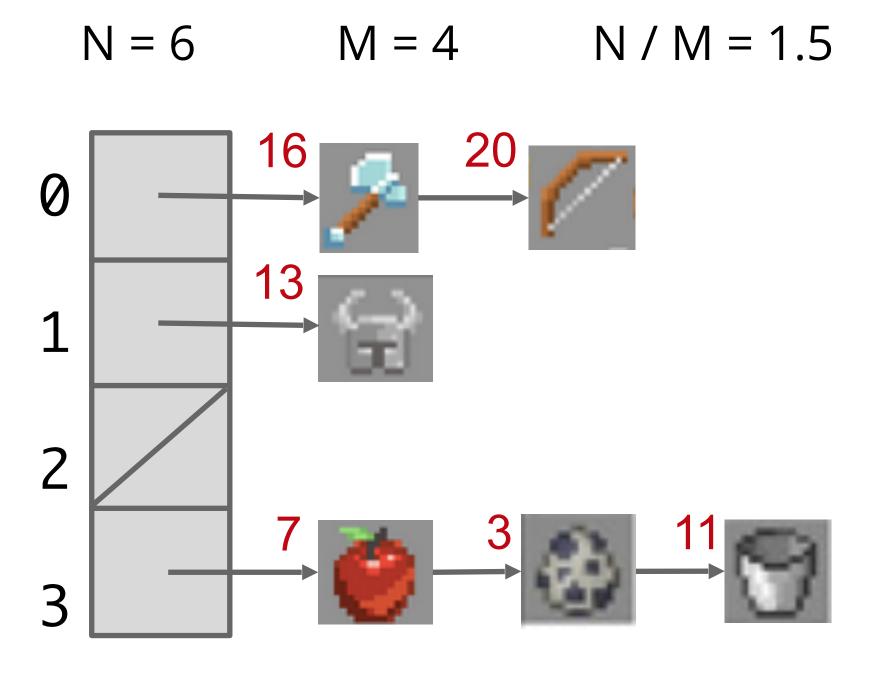
3 ?

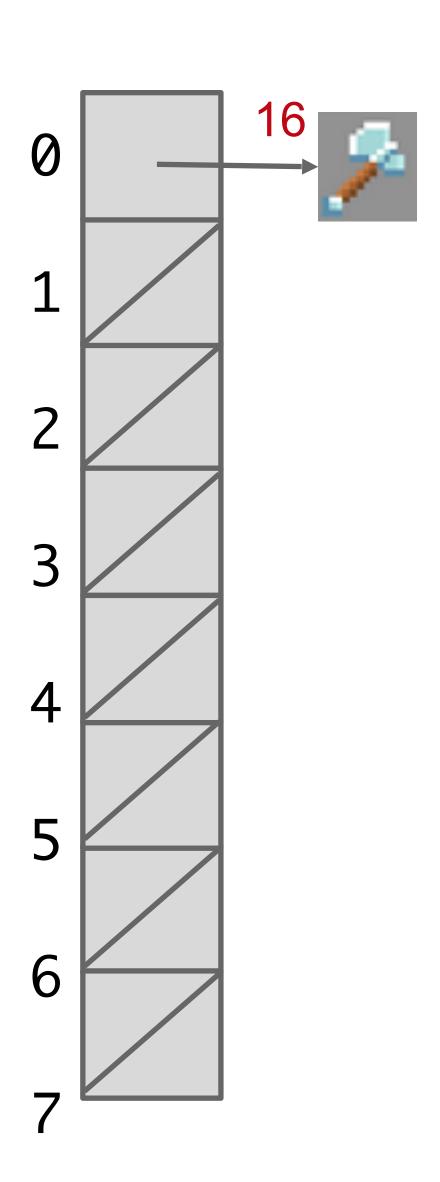
1 | ?

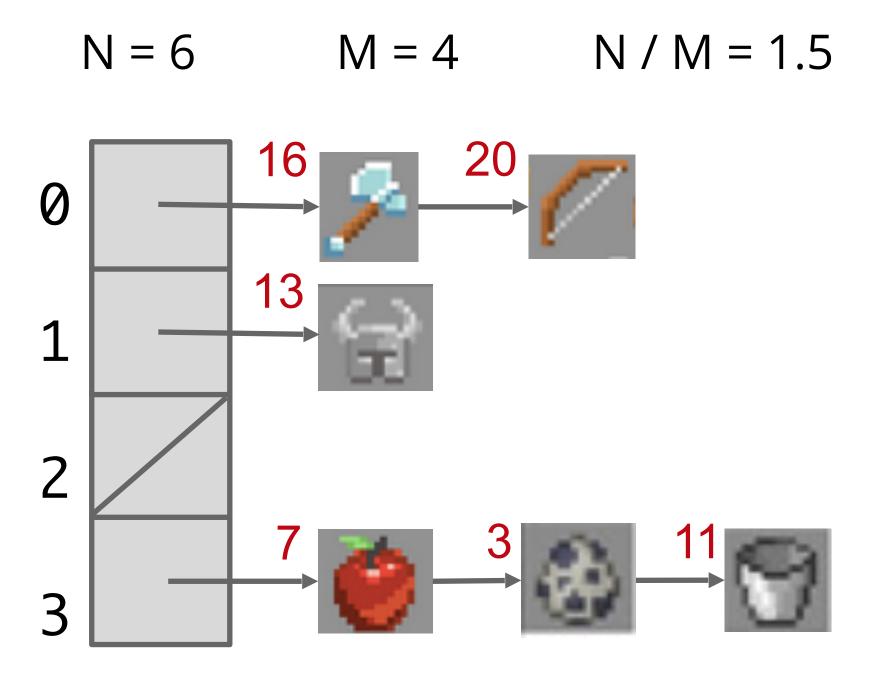
?

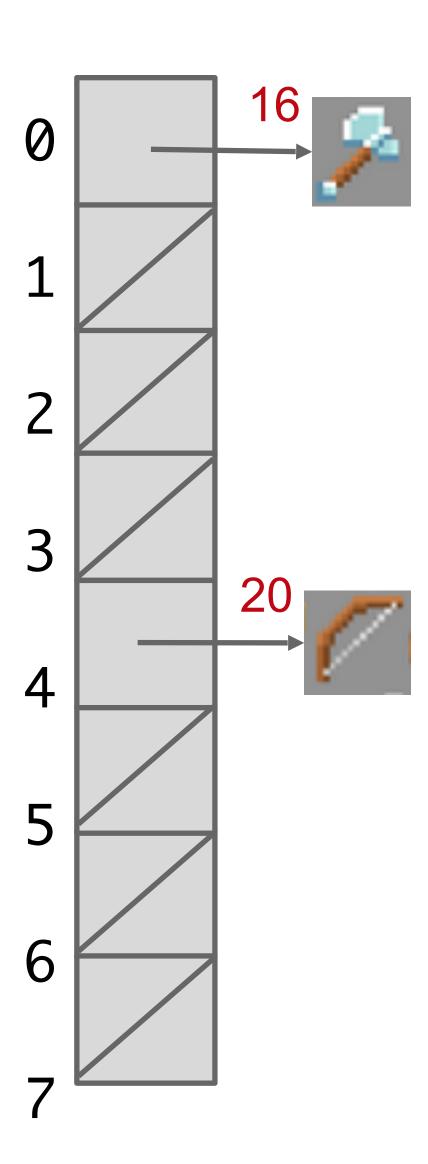
6 | ?

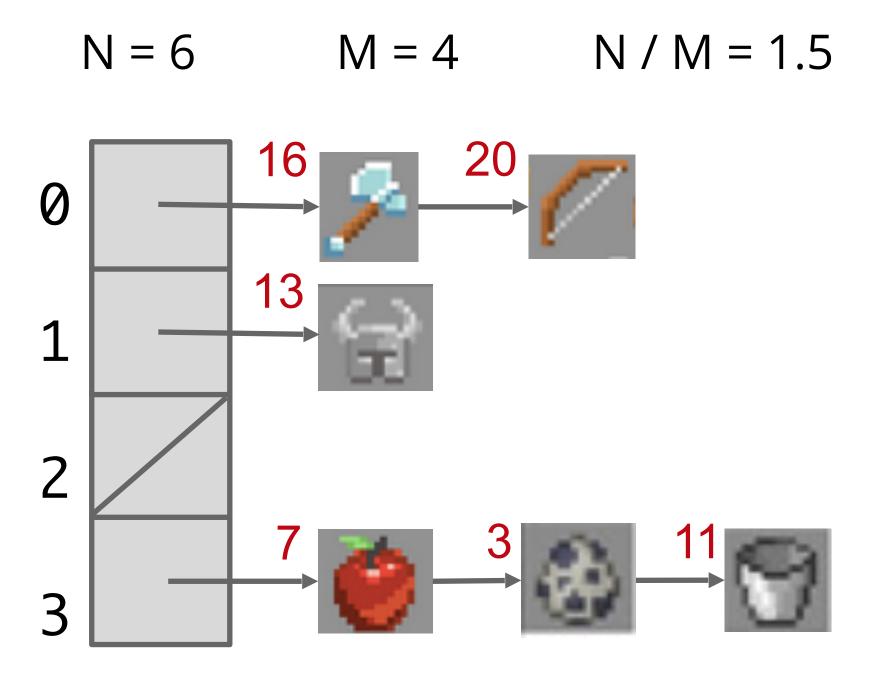
7 |

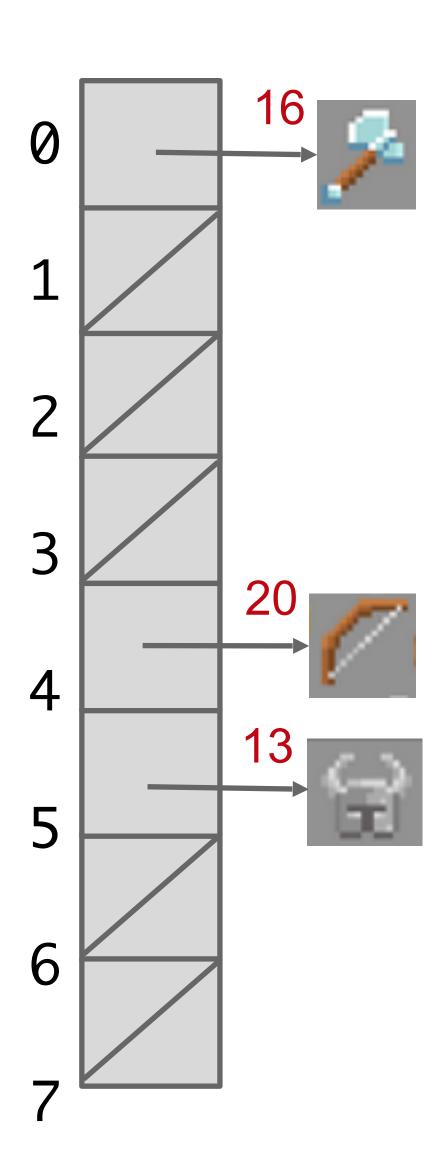


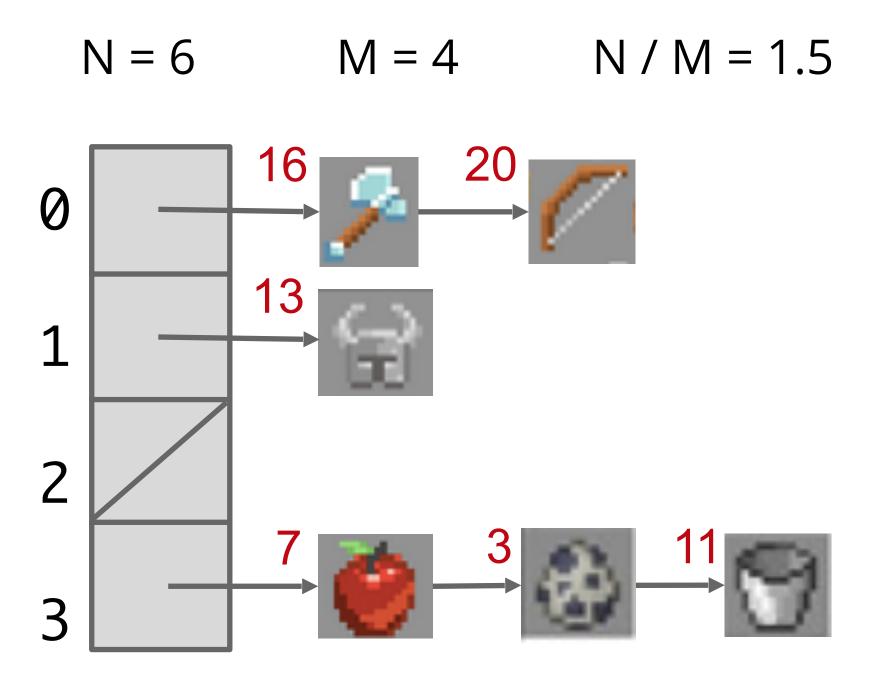


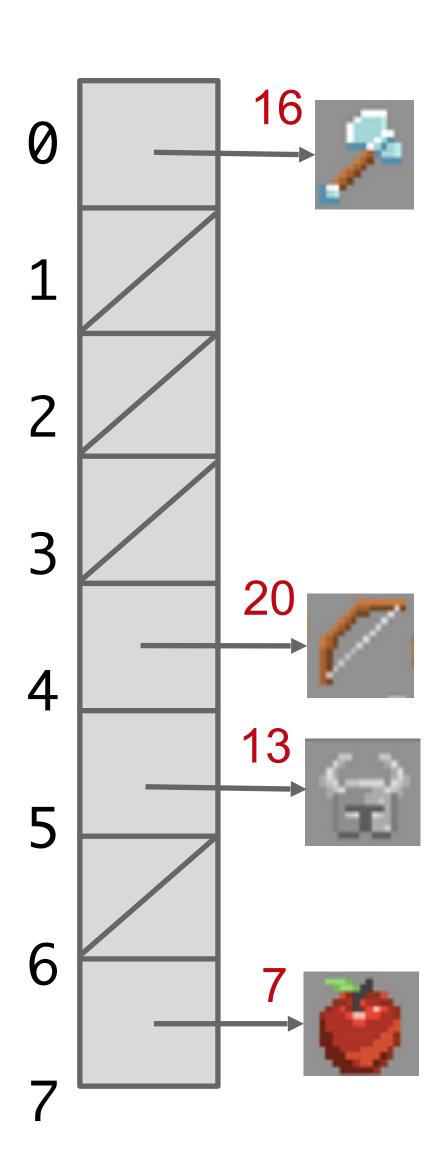


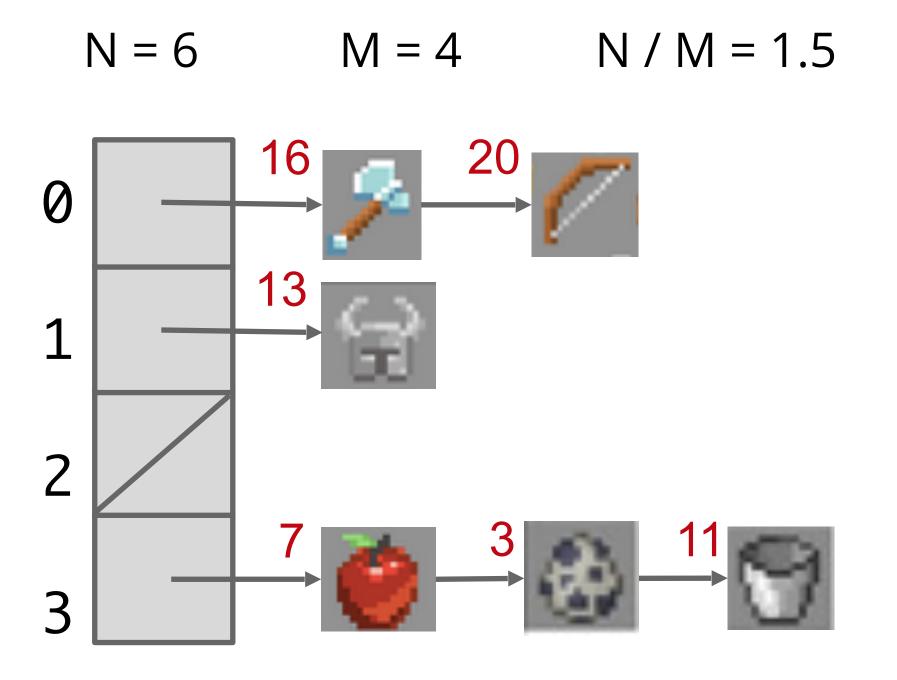


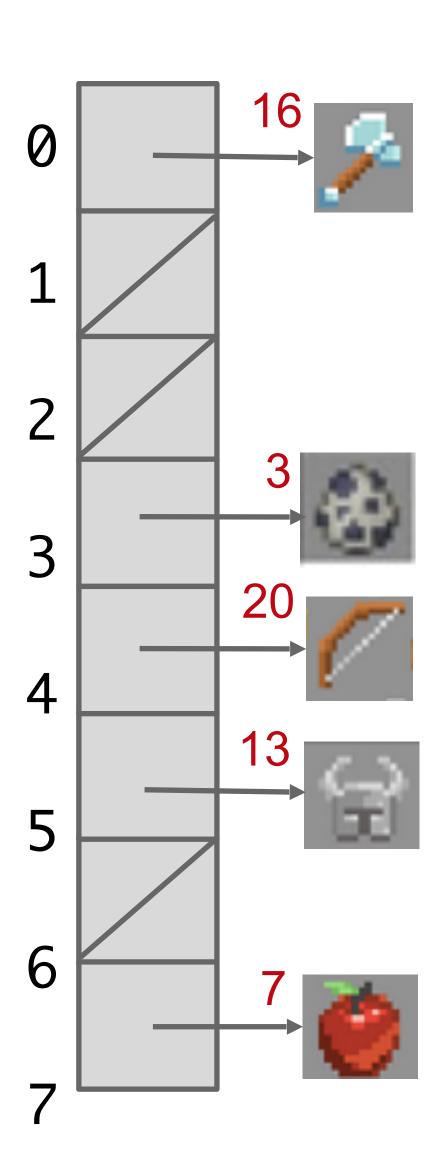


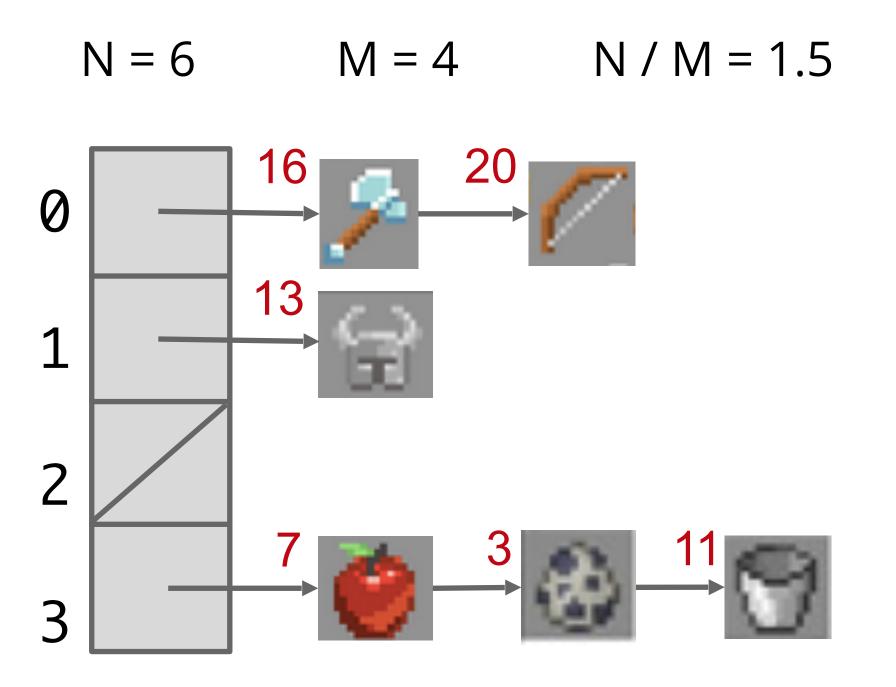


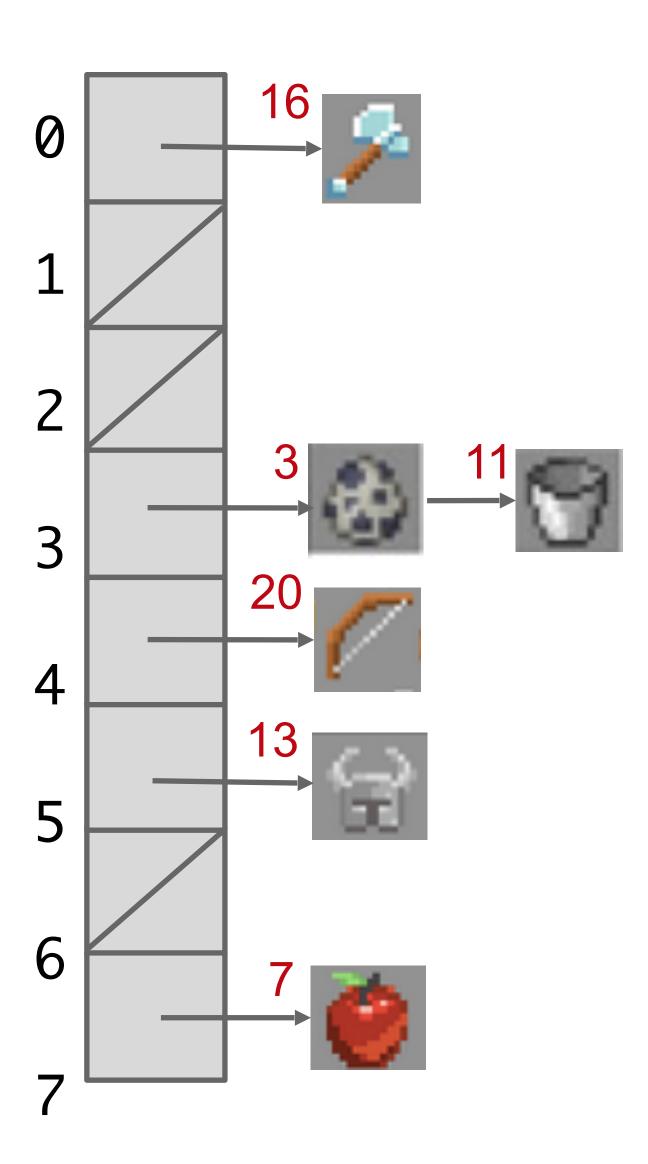


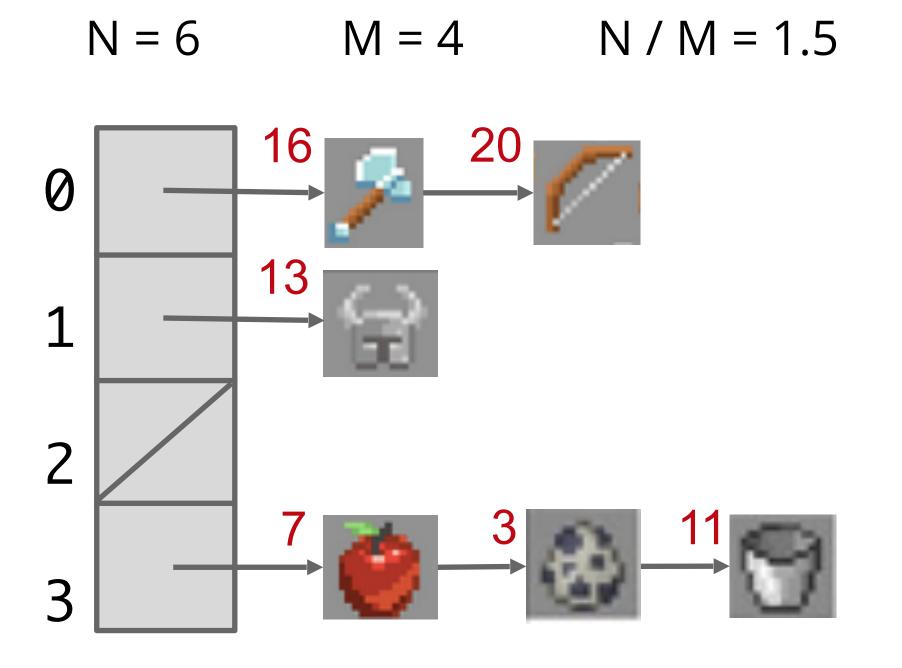


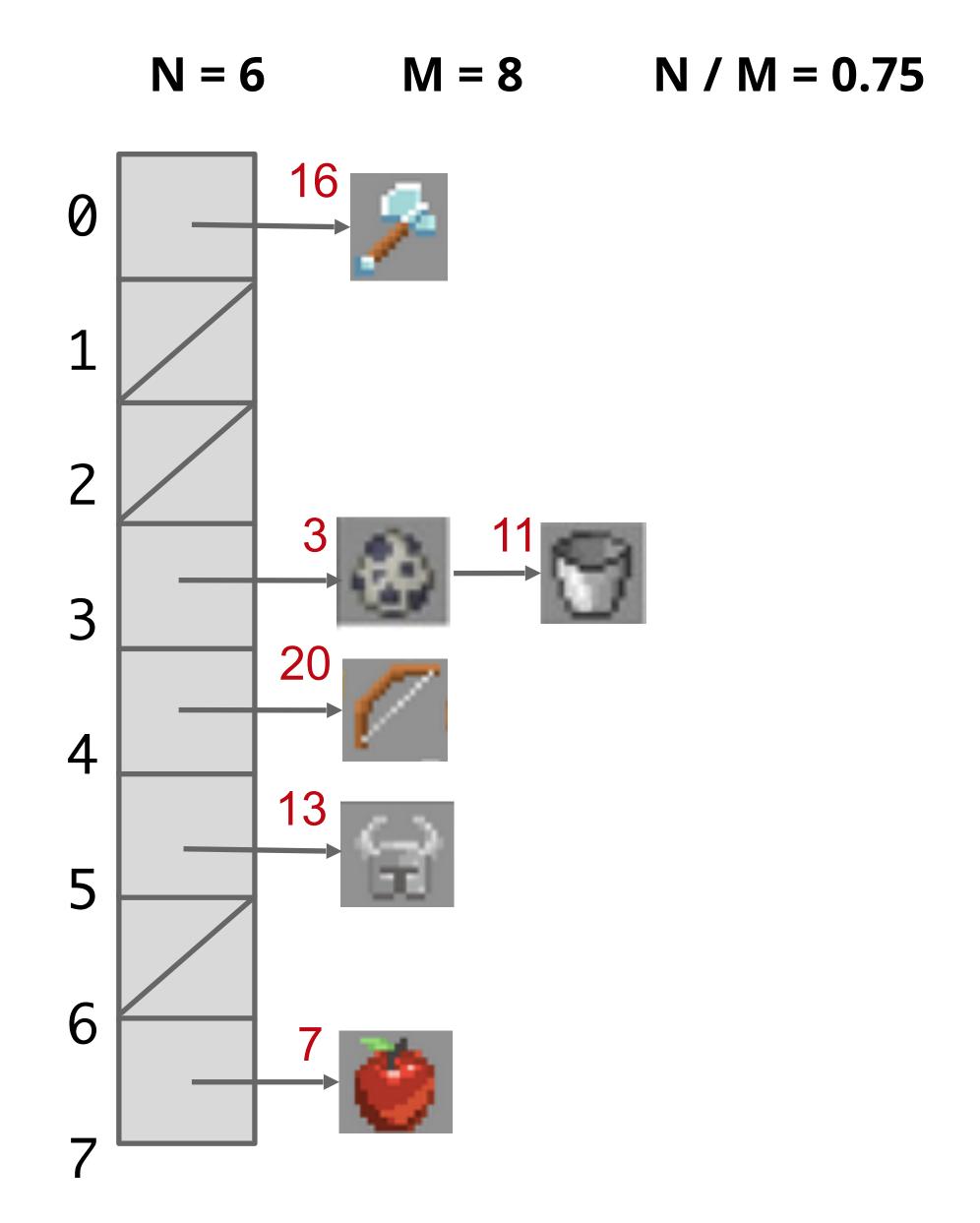












Final notes

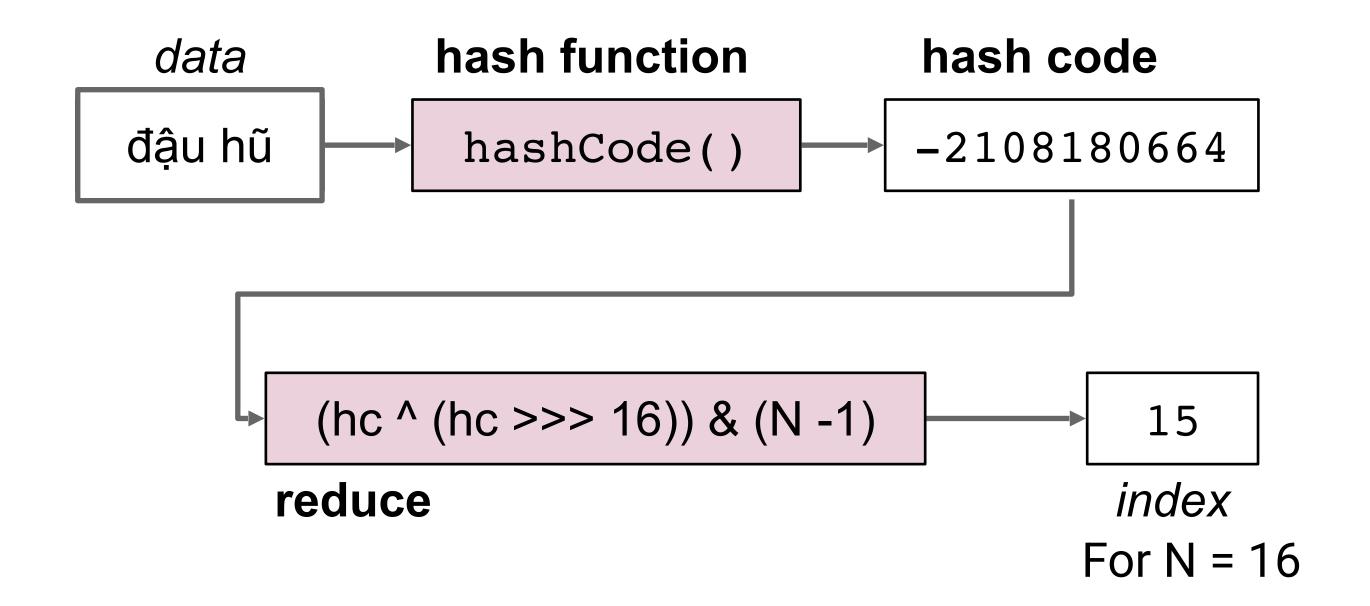
HashMap source code

We can then look at the code that implements the HashMap in Java:

 https://github.com/openjdk/jdk/blob/master/src/java.base/share/classes/java/util/ HashMap.java

Reading the code, we can see that:

- Hash table starts at size 16, then doubles every time N exceeds load factor which defaults to 0.75.
- The reduce function is a bit complicated using bitwise operations you'll learn in CS105.



Another Interesting Optimization

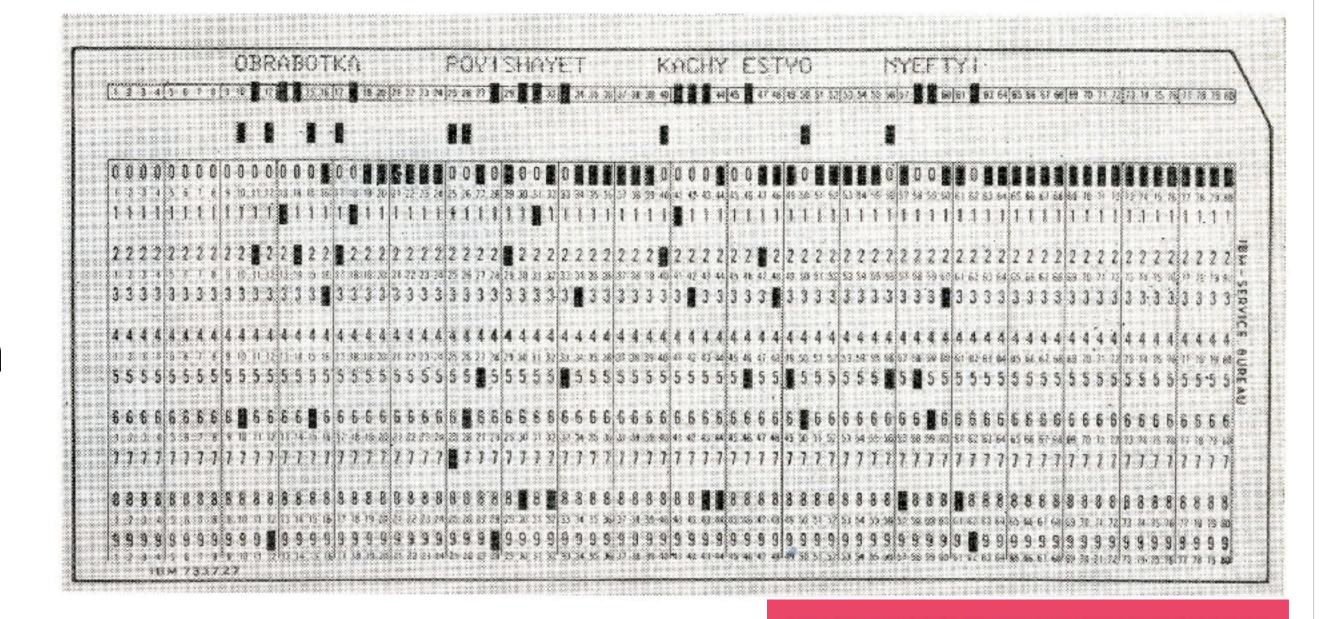
If we ctrl-F for "red-black" we find that that if a bin gets too full, it is converted into a red-black tree!

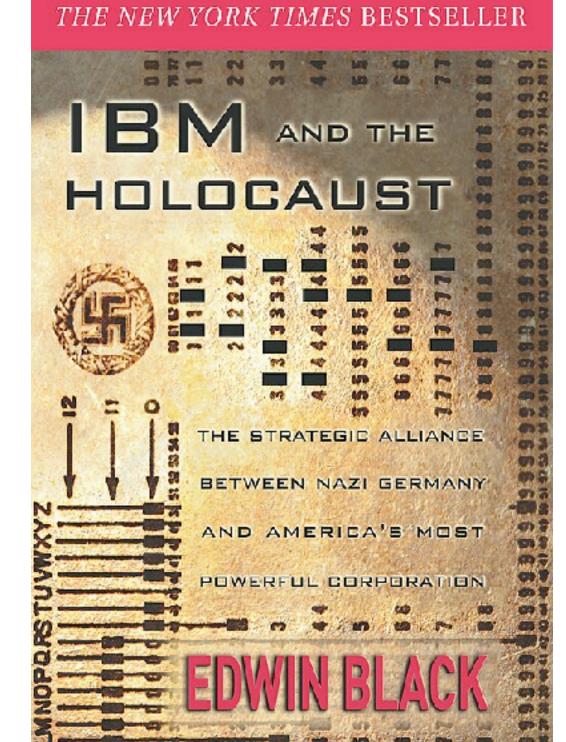
- "This map usually acts as a binned (bucketed) hash table, but when bins get too large, they are transformed into bins of TreeNodes, each structured similarly to those in java.util.TreeMap. Most methods try to use normal bins, but relay to TreeNode methods when applicable (simply by checking instanceof a node). Bins of TreeNodes may be traversed and used like any others, but additionally support faster lookup when overpopulated. However, since the vast majority of bins in normal use are not overpopulated, checking for existence of tree bins may be delayed in the course of table methods."
- This is well beyond the scope of our course.

"The most useful algorithms are, unfortunately, not always the most beautiful." - Josh Hug

Hashing history

- Herman Hollerith invented using punch cards to organize data for the 1890 US census (can physically sort cards based on hole location)
- This inspired IBM to make punch cards in the 1928s to store all kinds of data - the precursor to hash tables!
- During WWII, IBM's German branch (Dehomag) supplied the Nazi regime with punch cards to track things like a person's gender, age, race, and religion...so they could track who was Jewish to deport them and put them in concentration camps
- At the same time, in the US, punch cards were used to track interned Japanese Americans





History repeats itself

• Just like how IBM merged with state power and state interests during WWII, we see similarities today (e.g. Palantir and ICE)

ICE to Use ImmigrationOS by Palantir, a New Al System, to Track Immigrants' Movements

Published: August 21, 2025

Author: Steven Hubbard

Topics: Interior Enforcement

Questions:

- When does abstraction become dangerous? How might simplifying human identity into keys and categories—whether via punch cards or hash functions—erase complexity in ethically troubling ways?
- In what ways do modern technologies, like hash-based recommendation systems or surveillance tools, reflect the values of the corporations that build them? How should we understand the role of tech companies in shaping the moral direction of our digital infrastructure?

Lecture 18 wrap-up

- Checkpoint 2 regrades due 11:59pm next Thursday
- No HW this week, but final group project details released in lab tonight
- Lab tonight is project specs/details + coding interview practice
 - Practice writing code on the board and thinking through your process orally in groups
 - don't worry, you'll get full marks for just participating!

Resources

- Hashtable history (it's really dark. More next lecture): https://cs.pomona.edu/classes/cs62/history/
 hashtables/
- Reading from textbook: Chapter 3.4 (Pages 458-477); https://algs4.cs.princeton.edu/34hash/
- Hashtable visualization: https://visualgo.net/en/hashtable
- Practice problems behind this slide
- Most of these slides were from my hashtable teaching demo from when I was applying to teaching jobs ...:)

Problem 1

• Insert the keys E, A, S, Y, Q, U, E, S T, I, O, N in that order into an initially empty table of **m=5** lists, using separate chaining. Use the hash function **11*k%m** to transform the k-th letter of the English alphabet into a table index.

Problem 2

Try to write a function englishToInt that can convert English strings to integers by adding characters scaled by powers of 27.

Examples:

• a: 1

• z: 26

• aa: 28

• bee: 1598

• cat: 2234

dog: 3328

potato: 237,949,071

Answer 1

- Insert the keys E, A, S, Y, Q, U, E, S, T, I, O, N in that order into an initially empty table of m=5 lists, using separate chaining. Use the hash function 11*k%m to transform the k-th letter of the English alphabet into a table index.
- 0 -> O->T->Y->E
- 1 -> U -> a
- 2 -> Q
- 3 -> null
- 4 -> N -> I -> S

Answer 2

```
/** Converts ith character of String to a letter number.
    * e.g. 'a' -> 1, 'b' -> 2, 'z' -> 26 */
public static int letterNum(String s, int i) {
        int ithChar = s.charAt(i);
        if ((ithChar < 'a') || (ithChar > 'z'))
        { throw new IllegalArgumentException(); }
        return ithChar - 'a' + 1;
public static int englishToInt(String s) {
        int intRep = 0;
        for (int i = 0; i < s.length(); i += 1) {
        intRep = intRep * 27;
        intRep = intRep + letterNum(s, i);
         return intRep;
```