



Hash Tables I

Data Structures and
Parallelism

Announcements

Midterm Study Resources

Project 2

- *Slightly* shorter than normal.
- Writeup is extensive. It really is almost half of the assignment.

Exercise 3 due Wed at noon

- You can do either the version with the typo or the corrected version
- Just email me if you are doing the typo version

Wrap Up

AVL Trees:

$\Theta(\log n)$ worst case `find`, `insert`, and `delete`.

Pros:

Much more reliable running times than regular BSTs.

Cons:

Tricky to implement

A little more space to store subtree heights

Other Balanced Tree Dictionaries

There are lots of flavors of self-balancing search trees

“Red-black trees” work on a similar principle to AVL trees.

“Splay trees”

- Get $O(\log n)$ amortized bounds for all operations.

“Scapegoat trees”

“Treaps” – a BST and heap in one (!)

And more!

Similar tradeoffs to AVL trees.

Another Dictionary

Our guiding principle for designing AVL trees was optimizing for the **worst case**.

What if we want to optimize for the **average case**?

That goal will lead us to a totally different data structure: **hash tables**

A Simple Case

Suppose you were promised your keys would be distinct numbers in the range 0 to k .

How would you implement a dictionary. What are the running times for `insert`, `find`, and `delete`?

Just store the values in an array of size $k + 1$.

Store the value associated with i at index i of the array.

$O(1)$ operations for everything!

Generalization (Step 1)

What if the keys are guaranteed to be integers,
But the upper limit is huge.

Why not just use the array from last time?

How could we still use the array of size k ?

Generalization (Step 1)

What if the keys are guaranteed to be integers,
But the upper limit is huge.

Why not just use the array from last time?
WAY too much space

How could we still use the array of size k ?
Map the keys into the range $\{0, \dots, k - 1\}$.

% table size

Map to index $\text{key} \% \text{TableSize}$


indices	0	1	2	3	4	5	6	7	8	9
array	" : (' ' "	"biz"				"bar"			"bop"	

```
put (0, "foo"); 0 % 10 = 0
```

```
put (5, "bar"); 5 % 10 = 5
```

```
put(11, "biz") 11 % 10 = 1
```

```
put (18, "bop"); 18 % 10 = 8
```

```
put (20, ":"); 20 % 10 = 0  Collision!
```

Problem 1: What do we do when the keys collide?

Collision Resolution

Multiple Possible Strategies.

We'll talk about "open addressing" strategies later.

First, the Strategy for P2 is "Separate Chaining"

Idea: If more than one thing goes to the same spot

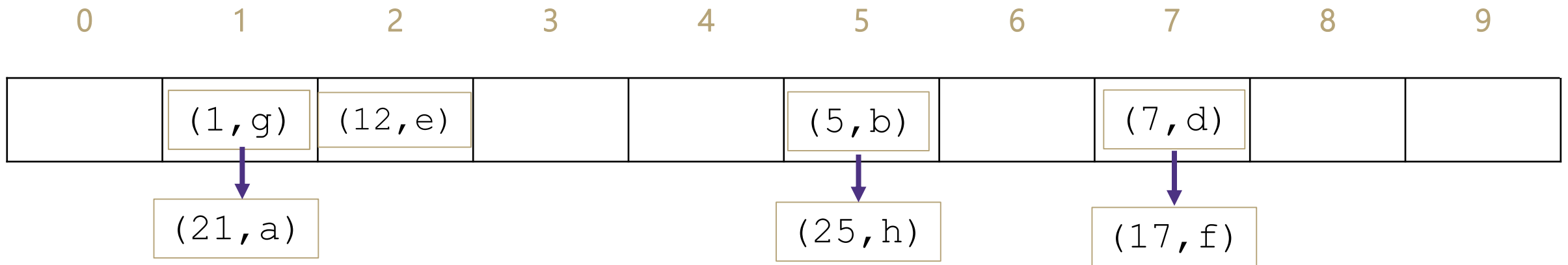
Just stuff them all in that one spot!

Separate Chaining

Instead of an array of values

Have an array of (say) LinkedLists of values.

Insert the following keys: (1, a) (5,b) (21,a) (7,d) (12,e) (17,f) (1,g) (25,h)



Running Times

What are the running times for:

`insert`

Best:

Worst:

`find`

Best:

Worst:

`delete`

Best:

Worst:

Running Times

What are the running times for:

`insert`

Best: $O(1)$

Worst: $O(n)$

`find`

Best: $O(1)$

Worst: $O(n)$

`delete`

Best: $O(1)$

Worst: $O(n)$

Average Case

What about on average?

Let's **assume** that the keys are randomly distributed

What is the average running time if the size of the table *TableSize* and we've inserted n keys?

insert

find

delete

Average Case

What about on average?

Let's **assume** that the keys are randomly distributed

What is the average running time if the size of the table *TableSize* and we've inserted n keys?

insert $O(1)$

find $O\left(1 + \frac{n}{TableSize}\right)$

delete $O\left(1 + \frac{n}{TableSize}\right)$

Average Case

What about on average?

Let's **assume** that the keys are uniformly distributed

What is the average running time if the size of the table *TableSize* and we've inserted n keys?

insert $O(1)$

find $O(1 + \lambda)$

delete $O(1 + \lambda)$

We'll denote $\frac{n}{TableSize}$ by λ .

Often called "load factor"

When λ Grows

If we keep inserting things into the array, λ will keep increasing.

We'll never *really* run out of room.

When should we resize?

When it slows us down, i.e. when λ is a constant.

Heuristic: for separate chaining λ between 1 and 3 is a good time to resize.

Resizing

How long does it take to resize?

Need to:

Remake the table

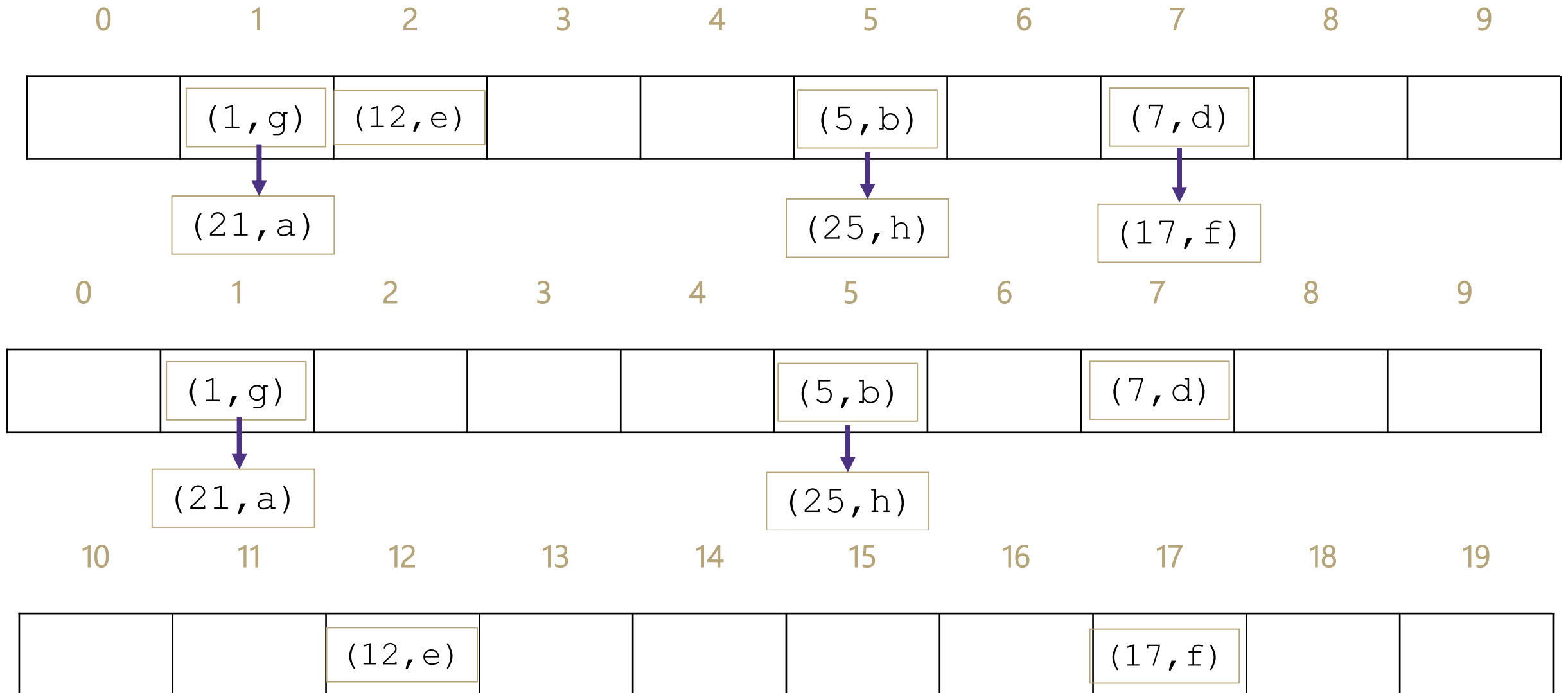
Evaluate the hash function over again.

Re-insert.

Total time: $O(n + TableSize) = O(n)$ if λ is a constant.

Resizing Redux

Let's resize by doubling the size of the array.



Resizing Redux

That didn't work very well!

It turned out that most of the keys that were equal mod 10 were also equal mod 20.

This is likely with real data.

Don't just double the table size

Instead make the table size some new prime number.

Collisions can still happen, but patterns with multiple prime numbers are rarer in real data than patterns with powers of 2.

Reaching the Average Case

In general our keys might not be integers.

Given an arbitrary object type E, how do we get an array index?



Hash
Function

17,423

% TableSize

Usually Object writer's
responsibility

Usually HashTable
writer's responsibility

How do we make our assumption (keys are uniformly distributed) true?
Or at least true-ish?

Designing a Hash Function

For simplicity, let's start with Strings.

Question: How many Strings are there compared to ints?

WAY more strings.

Can we always avoid collisions

-NO!

We can try to minimize them though.

Some Possible Hash Functions

For each of these hash functions, think about

- what Strings will cause collisions
- how long it will take to evaluate

Keys: strings of form $s_0s_1 \dots s_{k-1}$ (s_i are chars in range $[0,256]$)

$$h(K) = s_0$$

$$h(K) = \sum_{i=0}^{k-1} s_i$$

A Better Hash Function

$$h(K) = \sum_{i=0}^{k-1} s_i \cdot 31^i$$

Can we do this fast? Avoid calculating 31^i directly

```
for (i=k-1; i>=0; i--) {  
    h = 31*h + s[i];  
}
```


Other Classes

Should we use that same hash function if the strings are all URLs?

Other Classes

Should we use that same hash function if the strings are all URLs?

No! "https://www." is worthless, use the rest of the string

Other Classes

Should we use that same hash function if the strings are all URLs?

No! "https://www." is worthless, use the rest of the string

Person Class

String firstname; String middlename; String lastname; Date birthdate;

Tradeoff between speed and collision avoidance.

What to hash is often just an unprincipled guess.

General Principles

You have 32 bits, use them.

If you have multiple pieces, have the hashes stretch across bits

Bitwise xor if you have to combine

DON'T DO THIS IF YOU DON'T HAVE TO

Rely on others to get this right if you can.

Java Specific Notes

Every object in Java implements the `hashCode` method.

If you define a new Object, and want to use a hash table, you might want to override `hashCode`.

But if you do, you also need to override `equals`

Such that

If `a.equals(b)` then `a.hashCode() == b.hashCode()`

This is part of the contract. Other code makes this assumption!

What about the converse?

Can't require it, but you should try to make it true as often as possible.