	Topics
CSCI 136: Data Structures and Advanced Programming Lecture 29 Hash collisions	Hash collisions Graphs
Williams	
Your to-dos	
	Note about lab 9:

1. Read **before Wed**: *Bailey*, Ch. 16.4.

2. Lab 9 (solo lab), due Tuesday 5/3 by 10pm.

note about lab 9.

You may use the structure5 **Hashtable** implementation.



Hash codes

Hashing so important that every Object in Java has a built-in hash function.



Hash collisions

A hash collision is when two or more distinct keys have the same hash value.



Perfect hash function

A perfect hash function is a hash function that ensures that distinct keys map to distinct indices. I.e., there are no collisions.



Perfect hash function

Problem: It's **pretty darn hard** to come up with a perfect hash function.

- 1. You need to know all possible keys in advance.
- 2. If the number of possible keys is large, it is expensive to compute (O(n²) time) and expensive to store (O(n) space).

With a good hash table implementation, "imperfect" hash functions are usually **good enough**.

Dealing with collisions

There are two approaches to dealing with collisions:

- 1. Change your hash function.
- 2. Change your hash table design.

The easier of the two approaches turns out to be #2.

Open addressing

Open addressing is a method for resolving collisions in a hash table. Collisions are resolved by **probing**, which is a predetermined method for searching the hash table (aka a **probe sequence**). On **insertion**, probing finds the **first available bucket**. On **lookup**, probing searches until either the **key is found** or **an empty space** is found.







Deletion is trivial.



Method	Successful	Unsuccessful
Linear probes	$\frac{1}{2}\left(1+\frac{1}{(1-\alpha)}\right)$	$\frac{1}{2}\left(1+\frac{1}{(1-\alpha)^2}\right)$
Double hashing	$\frac{1}{\alpha} \ln \frac{1}{(1-\alpha)}$	$\frac{1}{1-\alpha}$
External chaining	$1 + \frac{1}{2}\alpha$	$\alpha + e^{-\alpha}$

Figure 15.11 Expected theoretical performance of hashing methods, as a function of α , the current load factor. Formulas are for the number of association compares needed to locate the correct value or to demonstrate that the value cannot be found.

Hash Table Expansion

When a hash table **fills up**, we should **expand**, just as with a Vector. But there are some problems...

Hash Table Expansion

Virtually every hash function relies on the size of the underlying array to do the hashing. Recall:

```
int index(K key) {
  return abs(h(key) % A.length);
}
```

When a hash table expands, we usually address this by rehashing all elements during a copy. Why is this OK?

Hash Table Expansion

Another issue: hash table performance **degrades severely** as it **fills up**.

Recall that we can have an **effectively full** hash table even when there is actually space.

 $h(\text{key}) + c \times i$

where c = 2



Hash Table Expansion

Therefore, we resize before the table is likely to be full.

Let **n** be the **number of elements** stored in a hash table.

Let m be the number of buckets.

Load factor = n / m

When the load factor is reached, the hash table is **resized**.

Hash Table Expansion

There are two ways to find a good load factor.

- 1. Careful analysis of the probability of attempting to insert more than one element into the same bucket, combined with a preference for acceptable average slowdown.
- 2. Empirical measurement, combined with a preference for acceptable average slowdown.

A load factor of 0.7-0.8 is generally accepted to be a good threshold.

Recap & Next Class

Today:

Hash collisions

Graphs

Next class:

Graph algorithms