

Lecture #14

- Hash Tables

- The Modulus Operator
- Closed hash tables
- Open hash tables
- Hash table efficiency and "load factor"
- Hashing non-numeric values
- `unordered_map`: A hash-based STL map class
- (Database) Tables

How can we write a `hashFunc` that converts our large **ID#** into a **bucket #** that falls within our 100,000 element array?

```
int hashFunc(int idNum)
{
    const int ARRAY_SIZE = 100000;
    int bucket = idNum % ARRAY_SIZE;
    return bucket;
}
```

This line takes an input value `idNum` and returns an output value between **0** and **ARRAY_SIZE - 1**. (0 to 99,999)

And this corresponding value can be used to pick a **bucket** in our 100,000 element array!

RIGHT! The C++ **% operator** (aka the **modulus division operator**) does exactly what we want!!!

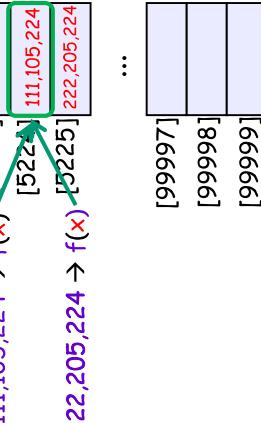
So now for each input **ID#** we can compute a corresponding value between **0-99,999**!

19 Closed Hash Table with Linear Probing: Insertion

Linear Probing Insertion:

As before, we use our hash function to locate the right bucket in our array.

If the target bucket is empty, we can store our value there.



However, instead of storing **true** in the bucket, we store our **full original value** - this prevents ambiguity!

If the bucket is occupied, scan down from that bucket until we hit the first open bucket. Put the new value there.

20 Closed Hash Table with Linear Probing: Insertion

Linear Probing Insertion:

Sometimes, you'll need to insert an item near the end of the table...

For instance, let's say we want to insert a new value of **640,099,998** into our hash table.

If you run into a collision on the last bucket, and go past the end... You simply wrap back around the top!



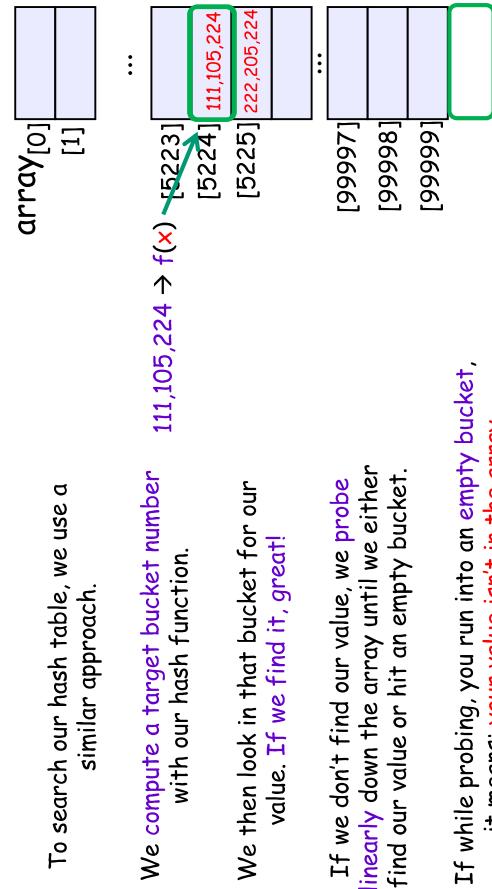
Closed Hash Table with Linear Probing: Searching

21

Closed Hash Table with Linear Probing: Linear Probing

22

Linear Probing Searching:



To search our hash table, we use a similar approach.

We compute a target bucket number with our hash function.

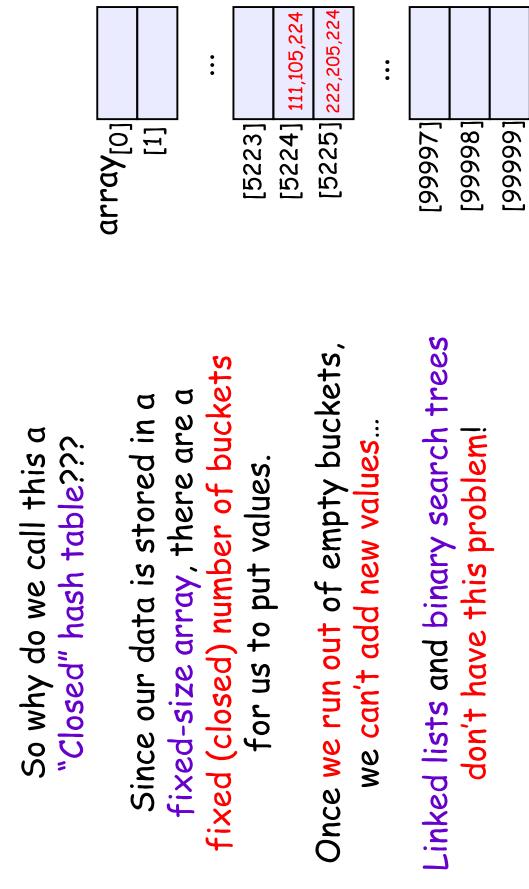
We then look in that bucket for our value. If we find it, great!

If we don't find our value, we probe linearly down the array until we either find our value or hit an empty bucket.

If while probing, you run into an empty bucket, it means: your value isn't in the array.

Closed Hash Table with Linear Probing

23



So why do we call this a "closed" hash table???

Since our data is stored in a fixed-size array, there are a fixed (closed) number of buckets for us to put values.

Once we run out of empty buckets, we can't add new values...

Linked lists and binary search trees don't have this problem!

Ok, let's see the C++ code now!

Linear Probing Hash Table: The Details

24

In a Linear Probing Hash Table, each bucket in the array is just a C++ struct.

Each bucket holds two items:

1. A variable to hold your value (e.g., an int for an ID#)
2. A "used" field that indicates if this bucket in the hash table has been filled or not.

```
struct BUCKET {  
    // a bucket stores a value (e.g. an ID#)  
    int idNum;  
    bool used; // is bucket in-use?  
};
```

If this field is false, it means that this Bucket in the array is empty. If the field is true, then it means this Bucket is already filled with valid data.

What Can you **Store** in your Hash Table?

```
#define NUM_BUCK 10
class HashTable {
public:
    void insert(int idNum)
    {
        int bucket = hashFunc(idNum);
        for (int tries=0;tries<NUM_BUCK;tries++)
        {
            if (m_buckets[bucket].used == false)
            {
                m_buckets[bucket].idNum = idNum;
                m_buckets[bucket].used = true;
                return;
            }
            bucket = (bucket + 1) % NUM_BUCK;
        }
        // no room left in hash table!!!
    }
private:
    int hashFunc(int idNum) const
    {
        return idNum % NUM_BUCK;
    }
    BUCKET m_buckets[NUM_BUCK];
};
```

Since our array has 10 slots, we will **loop up to 10 times looking for an empty space**. If we don't find an empty space after 10 tries, our table is full!

We'll store our new item in the **first unused bucket** that we find, starting with the bucket selected by our hash function.

If the current **bucket** is already occupied by an item, **advance to the next bucket** (wrapping around from slot 9 back to slot 0 when we hit the end).

Here's our hash function. As before, we compute our bucket number by **dividing the ID number by the total # of buckets** and then **taking the remainder (%)**.

Our hash table has 10 slots, aka "buckets."

```
struct Bucket
{
    int idNum;
    string name;
    float GPA;
    bool used;
};

bool search(int id, string &name, float &GPA)
{
    int bucket = hashFunc(idNum);
    for (int tries=0;tries<NUM_BUCK;tries++)
    {
        if (m_buckets[bucket].used == false)
            return false;
        if (m_buckets[bucket].idNum == idNum)
        {
            name = m_buckets[bucket].name;
            GPA = m_buckets[bucket].GPA;
            return true;
        }
        bucket = (bucket + 1) % NUM_BUCK;
    }
    return false; // not in the hash table
}
```

Oh, and if you like, you can include additional associated values (e.g., a **name**, **GPA**) in each bucket!

For instance, what if I want to also store the student's **name** and **GPA** in each bucket along with their ID#?

You can do that!

Now when you look up a student by their ID#, you can **ALSO** get their **name** and **GPA**!

How about searching our Open hash table?

The "Open" Hash Table

Idea: Instead of storing our values directly in the array, each array bucket points to a linked list of values.

To search for an item:

1. As before, compute a bucket # with your **hash function**:
`bucket = hashFunc(idNum);`
 2. Search the linked list at `array[bucket]` for your item
 3. If we reach the end of the list without finding our item, it's not in the table!
- Cool! Since the linked list in **each bucket** can hold an **unlimited** numbers of values... Our open hash table is **not size-limited** like our closed one!
- Insert the following values: 1, 3, 11, 25, 101

Linear Probing: Deleting?

0	idNum: 79	used: T
1	idNum: 1	used: F
2	idNum: 3	used: F
3	idNum: 25	used: F

So far, we've seen how to **insert** items into our **Linear Probe** hash table.

What if we want to **delete** a value from our hash table?

Let's take a naive approach and see what happens...

To delete the value, let's just zero out our value and set the **used** field to **False**...

If we delete a value where a collision happened... When we try to search again, we may prematurely abort our search, failing to find the sought-for value.

So, as you can see, if we simply delete an item from our hash table, we have problems!

There are ways to solve this problem with a **Linear Probing** hash table, but they're **not recommended**!

So, in summary, **only** use Closed/Linear Probing hash tables when you **don't intend to delete** items from your hash table.

Like if you're building a hash table that holds words for a **dictionary**... You'll just add words, never delete any, right?

Hash Table Efficiency

38

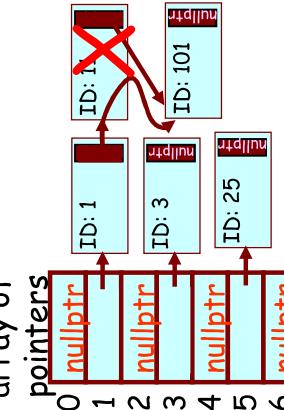
Oh - and there's no reason why we have to use a linked-list to deal with collisions...

Question:
How do you delete an item from an open hash table?

Answer:
You just remove the value from the linked list.

Let's delete the student with ID=11 and see what happens...

Cool! Unlike a closed hash table, you can easily delete items from an open hash table!



If you plan to repeatedly insert and delete values into the hash table, then the **Open table** is your best bet!
Also, you **can insert more than N items** into your table and still have great performance!

Hash Table Efficiency

Let's assume we have a completely (or nearly) empty hash table...

What's the maximum number of steps required to insert a new value?

Right! There's zero chance of collision, so we can add our new value in one step!

And finding an item in a nearly-empty hash table is just as fast!

We have no collisions so either we find an item right away or we know it's not in the hash table...

40

0	idNum: -1	GPA: etc...
1	idNum: -1	GPA: etc...
2	idNum: -1	GPA: etc...
3	idNum: -1	GPA: etc...
4	idNum: -1	GPA: etc...
5	idNum: -1	GPA: etc...
6	idNum: -1	GPA: etc...
7	idNum: -1	GPA: etc...
8	idNum: -1	GPA: etc...
9	idNum: -1	GPA: etc...

Hash Table Efficiency

Ok, but what if our hash table is nearly full?

What's the **maximum number of steps** required to insert a new value?

Right! It could take up to **N steps!**

And searching can take just as long in the worst case...

So technically, a hash table can be up to **O(N)** when it's nearly full!

So how big must we make our hash table so it runs quickly? To figure this out, we first need to learn about the "**load**" **concept**...

0	idNum: 89	GPA: 3.87
1	idNum: 21	GPA: 4.0
2	idNum: 12	GPA: 3.2
3	idNum: 42	GPA: 3.9
4	idNum: 34	GPA: 1.10
5	idNum: 06	GPA: 3.89
6	idNum: 67	GPA: 3.4
7	idNum: 78	GPA: 1.7
8	idNum: 29	GPA: 2.1
9	idNum: 11	GPA: etc...

41

Sizing your Hash Table

(Challenge)

If you want to store up to **1000 items** in an Open Hash Table and be able to find any item in roughly **1.25 searches**, **how many buckets** must your hash table have?

Remember: Expected # of Checks = $1 + \frac{L}{2}$

If our hash table has **2000 buckets** and we're inserting a maximum of **1000 values**, we are guaranteed to have an average of **1.25 steps per insert/search!**

This result means:

"If you want to be able to find/insert items into your open hash table in an **average of 1.25 steps**, you need a **load of .5**, or roughly **2x more buckets** than the **maximum number of values** you'll put into your table."

Answer:

Part 1: Set the equation above equal to 1.25 and solve for L:

$$1.25 = \frac{.5}{\text{Required hash table size}} \rightarrow .25 = \frac{.5}{\text{Required hash table size}} \rightarrow .5 = \frac{1}{\text{Required hash table size}}$$

Part 2: Use the load formula to solve for "Required size":

$$\frac{\# \text{ of items to insert}}{\text{Required hash table size}} \rightarrow \frac{1000}{\text{Required hash table size}} \rightarrow \frac{1000}{.5} = 2000 \text{ buckets}$$

What Happens If ...

What happens if we want to allow the user to search by the **student's name** instead of their **ID number**?

Well, our original hash function won't quite work:

```
int hashFunc(int ID)
{
    return(ID % 1000000);
} // what do we do?
```

Now we need a hash function that can convert from a **string of letters** to a number between **0** and **N-1**.

But this hash function isn't so good. Why not?

How can we fix it?

So basically it's a tradeoff!

You could always use a **really big hash table** with **way-too-many buckets** and ensure **really fast searches**...

But then you'll end up **wasting lots of memory**...

On the other hand, if you have a **really small hash table** (with just barely enough room), it'll **be slower**.

Finally, when **choosing the exact size** of your hash table (the number of buckets)...

Always try to choose a **prime number** of buckets...

Instead of **2000 buckets**,
give your hash table **2017 buckets**.

This causes **more even distribution** and **fewer collisions**!

A Hash Function for Strings

Here's one possibility for a hash function that can convert a string into a number between **0** and **N-1**.

```
int hashFunc(string &name)
{
    int i, total=0;
    for (i=0;i<name.length(); i++)
        total = total + name[i];
    total = total % HASH_TABLE_SIZE;
    return(total);
}
```

Hint:

What happens if we hash "**BAT**"?

What happens if we hash "**TAB**"?

A Better Hash Function for Strings

Here's better version of our string hashing function - while not perfect, it disperses items more uniformly in the table.

```
int hashFunc(string &name)
{
    int i, total=0;

    for (i=0;i<name.length(); i++)
        total = total + (i+1) * name[i];

    total = total % HASH_TABLE_SIZE;
    return(total);
}
```

Now "**BAT**" and "**TAB**" hash to different slots in our array since this version takes character position into account.

A GREAT Hash Function for Strings

Rather than write your own hash function from scratch, why not use one written by the pros?

```
#include <functional>
Make sure to
#includes <functional>
to use C++'s hash function!
First you define a
C++ string hashing
object.
#include <functional>
unsigned int yourHashFunction(const std::string &hashMe)
{
    std::hash<std::string> str_hash;
    unsigned int hashValue = str_hash(hashMe); // creates a string hasher!
    Then you use the
    object to hash
    your input string.
    Then just add your own modulo
    unsigned int bucketNum = hashValue % NUM_BUCKETS;
    return bucketNum;
}
Finally, you apply your own modulo
function and return a bucket # that
fits into your hash tables array.
```

Writing Your Own Hash Function

Great! But what if you need to write a hash function for some **non-standard data type**?

```
unsigned int yourHashFunction(const SomeCrazyTypeOfData &hashMe)
{
    Like hashing...
    Geospatial coordinates
    An array of N numbers
    The contents of a data file
    ...
}
```

This is a non-trivial exercise!

You really need to understand the "nature" of the data you're hashing...

Then **design your algorithm**, **analyze it**, and **iterate**.

The unordered_map: A hash-based version of a map

```
#include <unordered_map>
#include <iostream>
#include <string>
using namespace std;
int main()
{
    unordered_map<string,int> hm; // define a new U_M
    unordered_map<string,int>::iterator iter; // define an iterator for a U_M
    hm["Carey"] = 10; // insert a new item into the U_M
    hm["David"] = 20;
    iter = hm.find("Carey");
    // find Carey in the hash map
    if (iter == hm.end())
        cout << "Carey was not found!";
    else
    {
        cout << "When we look up " << iter->first; // "When we look up Carey"
        cout << " we find " << iter->second; // "we find 10"
    }
}
```

Hash Tables vs. Binary Search Trees

"Tables"



Hash Tables

Speed $O(1)$ regardless of # of items

Binary Search Trees
 $O(\log_2 N)$

Let's say you want to want to write a program to keep track of all your BFFs...

Simplicity Easy to implement
More complex to implement

Max Size **Closed:** Limited by array size
Open: Not limited, but high load impacts performance

Space Efficiency Wastes a lot of space if you have a large hash table holding few items

Ordering No ordering (random) Alphabetical ordering

"Tables"

In CS lingo, a group of related data is called a "record."

Each record has a bunch of "fields" like Name, Phone #, Birthday, etc. that can be filled in with values.

If we have a bunch of records, we call this a "table." Simple!

While you may have many records with the same `Name` field value (e.g., John Smith) or the same `Birthday` field value (e.g., Jan 1st)...

Some fields, like `Social Security Number`, will have unique values across all records - this type of field is useful for searching and finding a unique record!

Implementing Tables

```
struct Student
{
    string name;
    int IDNum;
    float GPA;
    string phone;
    ...
};
```

How could you create a **record** in C++?

Answer: Just use a `struct` or `class` to represent a record of data!

How can you create a **table** in C++?

Answer: You can simply create an array or vector of your struct!

```
// algorithm to search by the name field
int SearchByName(vector<Student> &table, string &findName)
{
    for (int s = 0; s < table.size(); s++)
        if (findName == table[ s ].name)
            return( s ); // the student you're looking for is in slot s
    return( -1 );
}
```



A field (like the SSN) that has unique values across all records is called a **"key field."**

```
// algorithm to search by the phone field
int SearchByPhone(vector<Student> &table, string &findPhone)
{
    for (int s = 0; s < table.size(); s++)
        if (findPhone == table[ s ].phone)
            return( s ); // the student you're looking for is in slot s
    return( -1 );
}
```



Implementing Tables

Heck, why not just create a whole C++ class for our table?

```
class TableOfStudents
{
public:
    TableOfStudents(); // construct a new table
    ~TableOfStudents(); // destructure our table
    void addStudent(Student& s); // add a new Student
    Student getStudent(int s); // retrieve Students from slot s
    int searchByName(string &name); // name is a searchable field
    int searchByPhone(int phone); // phone is a searchable field
    ...
private:
    vector<Student> m_students;
};
```

Tables

In the `TableOfStudents` class, we used a **vector** to hold our table and a **linear search** to find Students by their **name** or **phone**.

This is a perfectly valid table - but it's **slow** to find a student! How can we make it **more efficient**?

Well, we could alphabetically **sort** our vector of records by their **names**...

Then we could use a **binary search** to quickly locate a record based on a person's **name**.

But then every time we add a new record, we have to **re-sort** the whole table. Yuck!

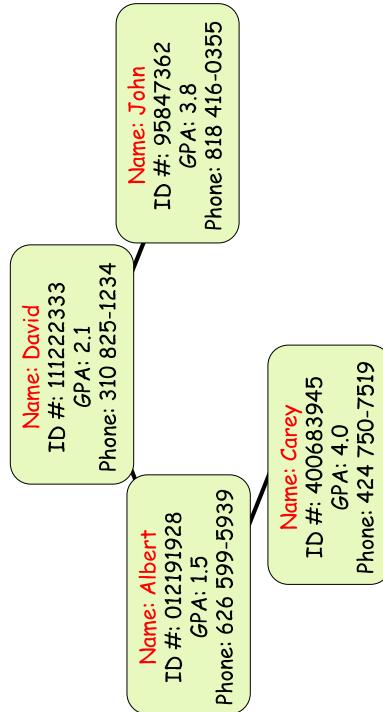
And if we **sort** by **name**, we can't search efficiently by other fields like **phone #** or **ID #**!

In the `TableOfStudents` class, we used a **vector** to hold our table and a **linear search** to find Students by their **name** or **phone**.



Tables

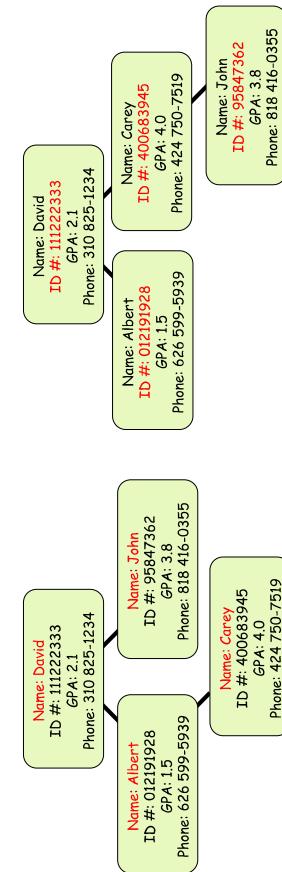
Hmm... What if we stored our records in a **binary search tree** (e.g., a **map**) organized by **name**? Would that fix things?



62

Tables

Hmm... What if we **create two tables**, ordering the first by **name** and the second by **ID#**?



63

Tables

Hmm... What if we **create two tables**, ordering the first by **name** and the second by **ID#**?

Well, now we can search the table efficiently by **name**... But we still can't search efficiently by **ID#** or **Phone #**....

That works... Now I can quickly find people by **name** or **ID#**!
But now we have **two copies of every record**, one in each tree!
If the records are big, that's a waste of space!

So what can we do? Let's see!

61

