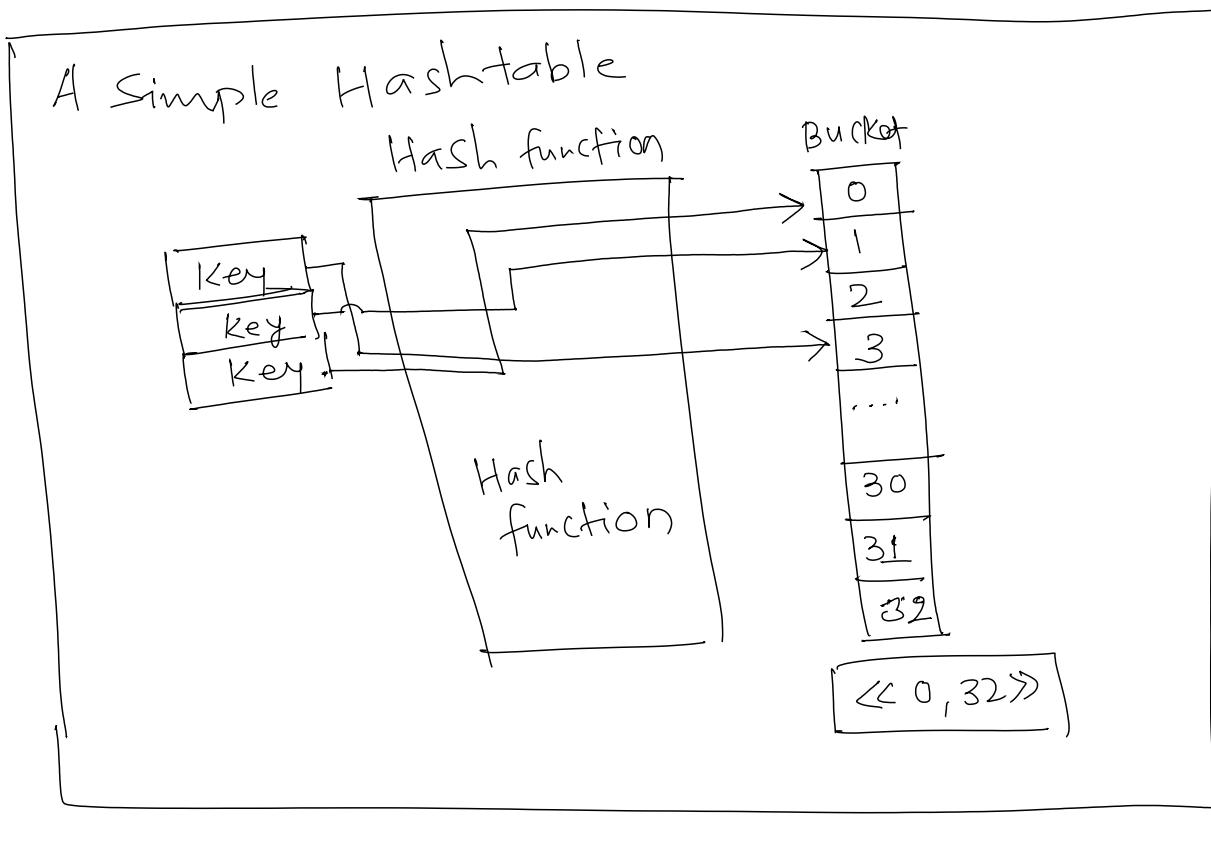


# Hashing and Tries

15/07/2022  
①

what is Hashing ?

Hashing = Hash function + Hash table  
①   ②



let's define Hash function?

"Mathematical function used to map some data (key) to a representative 'index'." This index is also referred to as a hash value

So, for ex:

A Simple Hash function

```
int hash( int Key ) {  
    int index = Key % SIZE;  
    return index;  
}
```

What is a hash table?

"A hash table/map is a data structure that maps keys to values."

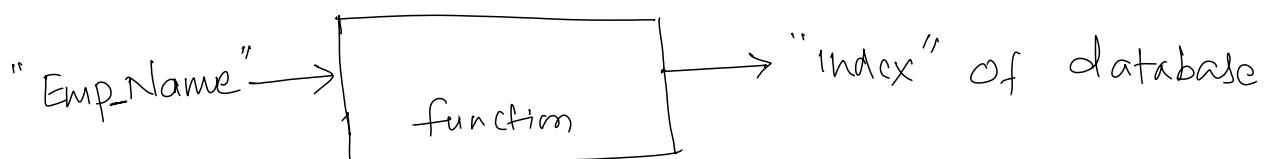
Why do we use Hashing?

In general hashes are used to speedup searches.

As an example, let's say you want to find the EmployeeID associated with a particular Employee Name.

A naive approach is to search an entire database for the particular employee. but this is not the optimal way to do it for large database.

so, to handle better of searching employee from database, hashing could be the alternative.



So, we have to use this function for generating the index of database, and you will store the Employee info at this location, and same way you can search/retrieve also.

---

Hashing is the key concept in performant code for many search specific problems.

Hashing could be a great fit if you can think of ~~con~~ forming conventional & unconventional key to be used.

(3)  
15/07/2022

I mean, conventional key Example : string, Integer ~~etc~~ etc.

Unconventional key example: matrix,  
list of integers( array, linked list etc.)

Please note Unconventional key finding is the bread/buffer of this approach.

#  
#

Hash can sometime also be used to maintain states in certain questions. For example, to maintain visited states in a complicated BFS/DFS

How do we implement Hashing ?

- Hashing would be achieved in traditional C++ by implementing a hash function that maps the key value to an integer.
- The hash table would then simply be an array of class/struct objects (the struct itself can be a list or other DS).

As an example

15/6/2022  
④

Hash function for String : (DJB2)

①

```
unsigned long hash(unsigned char *str) {
    unsigned long h = 5381;
    int c;
    while (*c = *str++)
    {
        h = ((h << 5) + h) + c;
    }
    return h % SIZE;
}
```

Hash fn for integers

②

```
unsigned int hash(
    unsigned int x)
{
    x = ((x >> 16) ^ x) * 0x45d9f3b;
    x = ((x >> 16) ^ x) * 0x45d9f3b;
    x = (x >> 16) ^ x;
    return x % SIZE;
}

* Check the hash value here.
```

How do we implement hashing?

- Unordered maps are internally implemented as a hash table with a simple modulo hash function
- ordered maps are internally maintained as balanced binary tree (RB Tree) — therefore no hashing or collision occurs

For getting better decisions, Please consider the below points

- ① How many buckets should I reserve.
- ② Should I use unordered\_map or ordered\_map
- ③ Is it possible to create a collision-free hash function which I can use?

# The Technical dept of Hashing

15/07/2022  
5

Please note time complexity is not always  $O(1)$  as we understood from beginning.

- A major pitfall is the issue of collisions and collision handling.
- Users also need to be cognizant of the complexity of the hash function being used. (Similarly in ordered\_map - the complexity of Compare function).
- Another pitfall is the blind reliance on a general purpose hash function, when a better/collision-free hash function exists.

---

Let's understand these disadvantages in more detail.

~~What is Hash Collision?~~

---

How to work around Collision?

Chaining

- Chaining is the simple idea of storing a list of values at the index of a hash table - rather than a single value.
- The downside is here you might need key also for verification of exact key match (Space Complexity increased)

## Open Addressing

- The idea is to find the next free space in the hash table to complete an insertion
- The free space can be searched for either in linear fashion or quadratic fashion.

"Drawback": Here depending on the search pattern - our hash table would need to grow dynamically. (we will reach its load factor quickly and require a rehash.)

## Understanding collisions in Unordered-maps

- Unordered\_maps work on the concept of "buckets". A "bucket" refers to the storage unit of one key-value pair.
- The Unordered\_map stores key-value pair by taking the modulo of input value by a prime number (10789 or 126271) and then stores it in hash table.
- When the input data is big and input values are multiples of this prime number a lot of collisions take place and may cause its complexity to  $O(n^2)$ .
- Keeping in mind that it is prudent to change the hash function used by the Unordered-map to a hash function with a better hash value distribution.

## The Bigger Enemy : Rehashing

When using a hash table we need to be aware of a concept known as the load\_factor. This is specially important when using STL Unordered\_maps.

# # Note  
What is load\_factor?

The load\_factor is the percentage of used space in the container.

# - Once the load\_factor reaches a threshold known as max\_load\_factor  
→ a rehash/resize operation of the container is triggered.

- During a rehash → All the elements in the container are rearranged according to their hash value into a new set of buckets. This may alter the order of iteration of elements within the container.

# --- Complexity of rehash =  $O(n \times (\text{Complexity of hash function} + \text{"New memory allocation")})$

## How to prevent unnecessary rehash?

15/07/2022  
⑧

- A simple fix to the rehashing problem is to reserve your unordered\_map with the Unordered\_map::reserve() operation.
- Such reserve() operation is not available in ordered\_map because ordered map is tree based and already optimised for memory uses. The concern of re-hashing is most prominent in case of unordered\_map.
- bucket size of reserve method will be the maximum size of key-value pair in unordered\_map.
- In general use this reserve method for reserving sufficient (more than) to avoid un-necessary re-hashing.

Let's below compare bfr ordered\_map vs unordered\_map.



Concept	<code>std::map</code>	<code>std::unordered_map</code>
Template Prototype	<code>template &lt;class key, class value, class Compare = less&lt;key&gt;&gt;</code>	<code>template &lt;class key, class value, class hash = hash&lt;key&gt;&gt;</code>
Internal Impl.	RBL Tree (self balancing binary tree)	Hash Table based on buckets (collision handled via open addressing)
Ordering	Sorted by Key	unordered
Size	Effective for small and Known Key Space	Effective for larger and Unknown key space
Average Case (ins, del, search)	$O(\log n)$	$O(1)$
Worst Case	$O(n)$	$O(n^{1/l})$ where $l$ is the length of key.
Memory Efficient	High Memory Efficient	Not Memory efficient

## Please note, sometime checking its performance b/w using `unordered_map` vs `map` is very useful. So please consider.