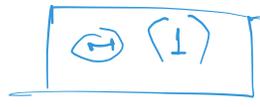


Hashing → some complexity



Dictionaries <key, value>

Look up for a key in the dictionary with $\Theta(1)$. values

Word → key that addresses its definitions.

Eg Two-letter words only! 26 alphabets.

$26 \times 26 = 676$ possible different 2-letter words.



Most of the possible combinations are NOT actual words.

class WORD

char → int
ASCII

```

{
  Encoding
  int hashCode() { returns
  }
}

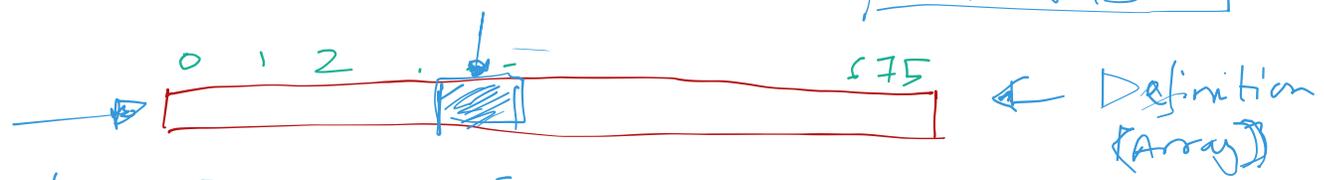
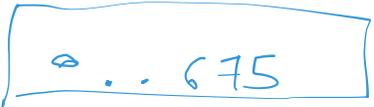
```

$$26 \times ((\text{WORDS}[0] - 'a') \rightarrow 0..25 + (\text{WORDS}[1] - 'a') \rightarrow 0..25 ;$$

```

private: char WORDS[2];
};

```



```

class Dictionary {
public:
  void insert(WORD& w, Value& v) {
    def[w.hashCode()] = v;
  }
  Definition& find(WORD& w)
private:
  Definition def; { return def[w.hashCode()];
};
}

```

Everything is $\Theta(1)$ operation
⇒ Looking up is extremely fast.



Problem: Longest word: 45 letters.
 26^{45} entries to store all possible words.

~ 700,000 words in the dictionary.
 Need to operate on arbitrarily large words!

Hash Tables

- * Insert
- * Find
- * Remove.

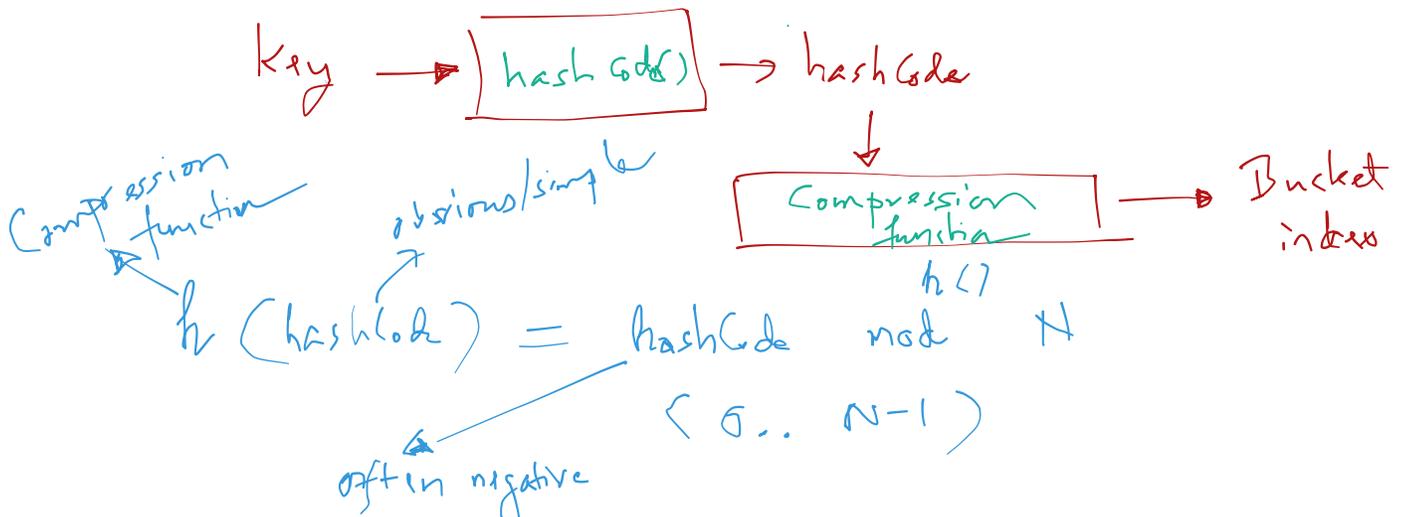
Actual

n : Number of keys stored

Table of N buckets, where N is a bit larger than n (~25% larger)

Way smaller than the total number of possible keys.

A Hash table maps a huge set of possible keys into N buckets by applying a compression function to every hash code.



\Rightarrow Likely to have collision

Collision: Several keys hash to the

same bucket

Iff $h(\text{hashCode1}) == h(\text{hashCode2})$

How to deal with the collision?

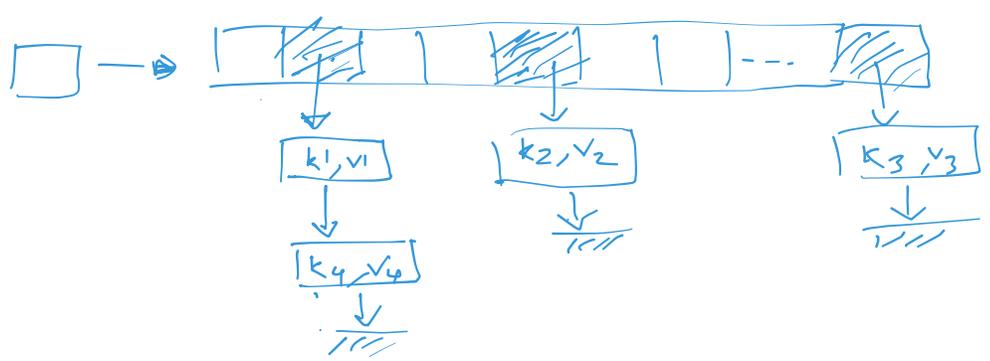
Chaining: Each bucket points to a linked list of entries called a chain.

=> store both words + definitions in the chain

Entry = <key, value>



definition Table



Operations:

- 1) Entry* insert (key, value)
 - Compute the key's hash code
 - Compare to determine the bucket
 - Insert the entry into the bucket's chain.

- 2) Entry* find (key)
 - Compute the key's hash code & the bucket
 - Search the chain for entry with the given key.
 - If found, return the entry. Otherwise, return NULL.

- 3) Entry* remove (key)
 - Hash key

std::unordered_map

- Search the chain.
- = Remove the entry from the chain.
- Return the entry or NULL.

Multiple copies of the same key?

- Insert both; find can return arbitrarily one of them.
- Replace the old value with the new. → ONLY ONE ENTRY.

Performance of a hash table

$$\text{Load factor of a hash table} = \frac{n \rightarrow \# \text{ keys}}{N \rightarrow \# \text{ buckets}}$$

Good load factor ≈ 0.8 or $1/2$

If the load factor stays low (close to a small constant) and hash code and compression function are good and no duplicate keys, then the chains are short and each operation takes $\Theta(1)$ time.

keys keeps growing

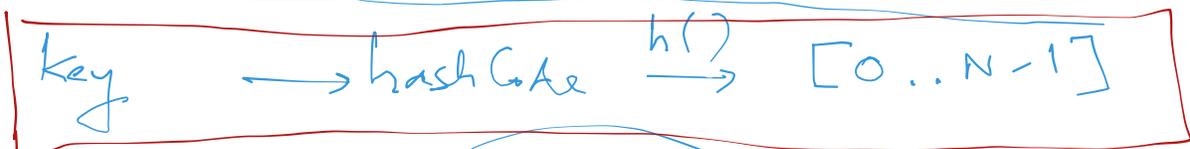
If the load factor gets BIG ($n \gg N$), $\Theta(n)$ time.

N is BIG → cost to set up uses a lot of memory.

② How to minimize collision?

mod

Good Hash Code and Compression function



$17 \bmod 3 \rightarrow 2$
 $[\bmod 3]$

0, 1, 2

5 mod 3 = 2

$\frac{N}{3} \rightarrow$ Take Remainder
 $N \bmod 3 \rightarrow$ Remainder when $N \div 3$.

hashCode mod N \rightarrow $[0, N-1]$ \rightarrow for fewer collisions

Ideal: Map each key to a random bucket with each bucket being equally likely.

Bad Compression function

Suppose keys are all ints.

hashCode (i) = i

$h(\text{hashCode}) = \text{hashCode mod } N$

Suppose $N = 10000$

& keys are divisible by 4,
 $h()$ is divisible by 4, too.

$\Rightarrow \left(\frac{3}{4}\right)^{\text{th}}$ of the buckets are NEVER used.

Same compression function can be better if

N is PRIME

Better: $h(\text{hashCode}) = (a \times \text{hashCode} + b) \bmod p \bmod N$

$a, b, p \rightarrow$ Positive int.

p is a large prime.

$p \gg N$

scrambling of the bits.

need not be prime

$\Rightarrow N$ does not need to be a prime.

Good Hash Code for strings

```

int hashCode (string key) {
    int hashVal = 0;
    for (int i = 0; i < key.length(); ++i) {
        hashVal = (127 * hashVal + key[i])
                % 16908799;
    }
    return hashVal;
}
}

```

↪ 127?

⇒ Works well in practice.

Bad hashCode : Creating biases
for string keys in data

①. Sum ASCII values of characters.

→ Rarely crosses 500.

Most entries mapped into these small set of buckets.

eg - Anagrams like "pat", "apt", "tap" collide.

②. First 3 letters of a word with 2^3 buckets.

Biases in English : lot of "pre", "pro", etc.

No word begins with "exp"...

⇒ Lots of collision.

③. Suppose change prime modulus to 127

$(127 * \text{hashVal}) \% 127 = 0$

⇒ Only the last char in the

Resizing Hash Tables:

If the load factor $\left(\frac{n}{N}\right)$ gets too big,
we lose $\Theta(1)$ time.

- Enlarge the hash table when $\frac{n}{N} > c$ 0.5
↑
- Allocate new array (at least twice as large) ↑
Hash Table
- Walk through the old array, rehash entries into the new hash table.

Cannot simply copy from old chains into the new table.

Compression function changes \rightarrow Entries go into diff. buckets.

If $\frac{n}{N} < 0.25$, shrink the hash table to free up some memory