# Estructura de Datos II

David Concha Gómez

Asignatura obligatoria

Segundo cuatrimestre

Créditos: 6

Moodle de la asignatura

Guía docente

- Conjuntos y mapas
- Tablas Hash
  - Introducción
  - Funciones Hash
    - Códigos Hash
    - Función de compresión
  - Colisiones
- Diccionarios
  - unordered\_multiset y unordered\_multimap

- Conjuntos y mapas
- Tablas Hash
  - Introducción
  - Funciones Hash
    - Códigos Hash
    - Función de compresión
  - Colisiones
- Diccionarios
  - unordered\_multiset y unordered\_multimap

• Un conjunto almacena elementos que pueden ser localizados eficientemente O(1).

• Los elementos no se repiten y no tienen orden, no hay un primero, no hay un último.

Ya los conocéis como arrays de booleanos.

• Un mapa es similar a un conjunto, pero guarda parejas clave:valor.

 Solo hay una entrada por clave, aunque varias claves pueden tener el mismo valor.

La búsqueda por clave es eficiente.

• Son muy usados cuando se necesita indexar por índices no enteros.

• También se usan con índices enteros cuando estos no son consecutivos.

Se implementan mediante tablas hash.

- Operaciones generales:
  - o isEmpty y size?
  - T& operator[](k): Devuelve una referencia al valor de clave "k". Si no existe devuelve un elemento vacío.
  - size\_t erase(k): Elimina el elemento con clave "k" si existe. Devuelve cuantos elementos se han borrado (0 o 1).
  - bool contains(k): Devuelve si existe una entrada con clave "k".
  - begin, end: Iteradores a los valores de la estructura.

```
std::unordered map<int, double> m = {
{3, 1.0}, {42, 4.2}, {37, 3.5}};
m[9] = 0.0:
m.erase(42);
for (const auto& p: m) {
 std::cout << "Key: " << p.first;
 std::cout << " Value: " << p.second;
 std:cout << "\n";
```

- Conjuntos y mapas
- Tablas Hash
  - Introducción
  - Funciones Hash
    - Códigos Hash
    - Función de compresión
  - Colisiones
- Diccionarios
  - unordered\_multiset y unordered\_multimap

- Las tablas hash son una forma eficiente de implementar conjuntos y mapas.
  - Se dispone de una función hash que permite convertir la clave en un entero utilizado para indexar.
  - Las operaciones de inserción, borrado y consulta son O(1) amortizado.

- Una tabla hash se componen de:
  - Un almacén para elementos, típicamente un array llamado bucket.
  - Una función hash i = h(k)

Internamente, se almacenan los pares (k, v) en la posición i=h(k).

• Si las claves son enteros en el rango de la tabla, la implementación es directa.



Clave

Valor

- Problemas:
  - Si la tabla es menor que la clave máxima (DNI, número de expediente, ...).
  - Si las claves no son enteros (nombres, código alfanumérico, ...)

 La función hash "mapea" las claves de un determinado tipo a enteros en un determinado rango.

m[09562] = Point{};

Para variables enteras, por ejemplo: h(k) = k %N

El valor se conoce como clave dispersa (hash value).

Clave

Valor

m[23490] = Point{};

m[12385] = Point{};

 $m[36893] = Point{};$ 

23490

\*

09562

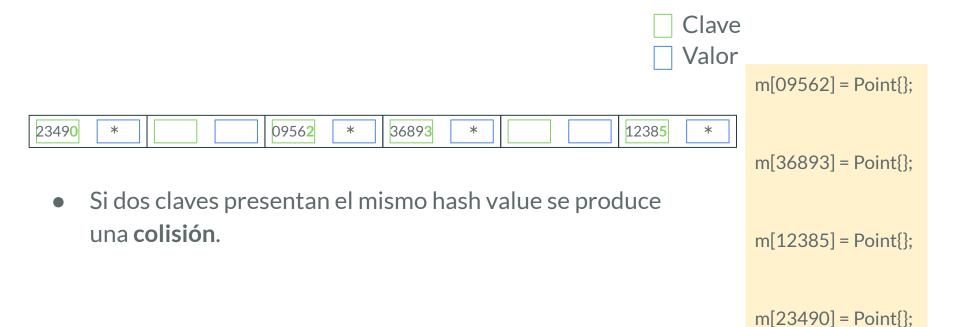
\*

3689<mark>3</mark>

\*

1238<mark>5</mark>

\*



 $m[10233] = Point{};$ 

Es importante seleccionar una buena función de hash:

• Que sea fácil de calcular

Que minimice las colisiones

- Conjuntos y mapas
- Tablas Hash
  - Introducción
  - Funciones Hash
    - Códigos Hash
    - Función de compresión
  - Colisiones
- Diccionarios
  - unordered\_multiset y unordered\_multimap

## **Funciones Hash**

- Las funciones hash se pueden ver como una composición de dos funciones:
  - Código hash:

$$h_H: Key \rightarrow int$$

Función de compresión:

$$h_c: int \rightarrow [0, N-1]$$

- El código hash transforma la clave en un entero.
  - Es dependiente de la clave y el problema.

- La función de compresión limita el entero a un rango aceptable para la tabla.
  - Es dependiente de la tabla hash.

# Código hash

- Deben ser consistentes:
  - O Dos claves que se consideren iguales deben presentar el mismo código hash.

- Se mapea un objeto a un valor entero.
  - En general definir un buen código hash es complicado.

• Si se quiere que un objeto pueda actuar como clave será necesario definir una función hash adecuada.

# Código hash

#### Ejemplos:

- Dirección de memoria del objeto.
- Interpretación de los bits de la clave como un entero.
- o Dividir los bits en bloques, sumarlos e interpretarlos como un entero
- Concatenar los atributos de un objeto y proceder con algún otro método.
- Acumulación polinómica:
  - Dividir los bits de la clave en componentes de la misma longitud.

$$k = \{k_0, k_1, k_2, ..., k_{n-1}\}$$

Dada una constante a calcular el polinomio:

$$p_0(a) = k_{n-1}$$
  
 $k_0 + k_1 a + k_2 a^2 + ... + k_{n-1} a^{n-1}$   
 $p_i(a) = k_{n-i-1} + a p_{i-1}(a)$ 

■ Para a = 33 se generan unas 6 colisiones para 50.000 palabras.

# Funciones de compresión

- Ejemplos:
  - División:

$$h_c(y) = y \mod N$$

- El tamaño de la tabla N suele ser un número primo.
- Ejemplo: Para {200, 205, 210, ... 600}
  - Si N = 100 cada clave colisiona con otras 4
  - Si N = 101 no hay colisión
- o MAD Multiplicar, sumar (Add) y dividir

$$hc(y) = ((ay + b) \mod p) \mod N$$

- p es un primo mayor que N.
- a y b son enteros no negativos aleatorios en el rango [0, p-1]

- Conjuntos y mapas
- Tablas Hash
  - Introducción
  - Funciones Hash
    - Códigos Hash
    - Función de compresión
  - Colisiones
    - Encadenamiento separado
    - Direccionamiento abierto
- Diccionarios
  - unordered\_multiset y unordered\_multimap

## **Colisiones**

- Dos claves pueden coincidir en la misma dirección.
  - Principio del palomar

A este fenómeno se le denomina colisión.

- Existen técnicas efectivas de resolver los conflictos:
  - Encadenamiento separado.
  - Direccionamiento abierto.

# **Encadenamiento Separado**

- Cada celda de la tabla tiene guarda todas las entradas con la misma clave hash.
  - La clave del par clave:valor es distinta.

 La forma de guardar las distintas entradas es utilizar una lista.

```
    La lista añade complejidad a las operaciones
del conjunto/mapa, pero, idealmente, solo
hay un elemento por lista.
```

```
void operator[](K k, T e) {
 auto& I = m[h(k)];
 if (contains(k)) {
    l.remove(k);
 l.insert({k, e});
bool contains(K k) {
 return m[h(k)].size();
```

## Direccionamiento abierto

- Cuando una entrada colisiona, se guarda en otra posición.
  - o Prueba lineal: se inserta en la siguiente posición libre.

$$h(k, p) = h_c(k) + p$$

p es la prueba (probe), cuantas veces se ha llamado a la función hash para resolver una colisión. Inicialmente 0.

## Direccionamiento abierto

- Cuando una entrada colisiona, se guarda en otra posición.
  - Prueba lineal: se inserta en la siguiente posición libre.

$$h(k, p) = h_c(k) + p$$

Clúster primario.

• Prueba cuadrática: Se desplaza dando saltos definidos por dos constantes ( $c_1 y c_2$ )

$$h(k, p) = h_c(k) + c_1 p + c_2 p^2$$

p es la prueba (probe), cuantas veces se ha llamado a la función hash para resolver una colisión. Inicialmente 0.

## Direccionamiento abierto

- Cuando una entrada colisiona, se guarda en otra posición.
  - Prueba lineal: se inserta en la siguiente posición libre.

$$h(k, p) = h_c(k) + p$$

Clúster primario.

Prueba cuadrática: Se desplaza dando saltos definidos por dos constantes ( $c_1 y c_2$ )

$$h(k, p) = h_c(k) + c_1 p + c_2 p^2$$

Clúster secundario.

Hashing doble:

$$h(k, p) = (h_c(k) + d(k)p) \mod N$$

$$d(k) = q - (h_c(k) \mod q)$$

p es la prueba (probe), cuantas veces se ha llamado a la función hash para resolver una colisión. Inicialmente 0.

Factor de carga (λ).

$$\lambda = ceil(n/N)$$

- Indica la ocupación del array interno.
- Es importante que sea menor que 1.
  - Encadenamiento separado < 0.75
  - Direccionamiento abierto < 0.5

- Si no se cumple habrá que redimensionar la tabla:
  - Añadir más buckets
  - Definir una nueva función de hash para adaptarse al nuevo tamaño
  - Rehasing de las entradas ya existentes.

- La complejidad en el peor caso es O(n) pero se considera O(1) amortizada.
  - Solo ocurrirá si todas las claves colisionan.

• El factor de carga afecta al rendimiento.

- Aplicaciones:
  - Pequeñas BBDD
  - Procesamiento de texto
  - Caché para juegos y navegadores

- Conjuntos y mapas
- Tablas Hash
  - Introducción
  - Funciones Hash
    - Códigos Hash
    - Función de compresión
  - Colisiones
    - Encadenamiento separado
    - Direccionamiento abierto
- Diccionarios
  - unordered\_multiset y unordered\_multimap

## **Diccionarios**

• Existe una variante de los conjuntos y mapas que permite claves repetidas.

En c++ se llaman unordered\_multiset y unordered\_multimap.

Cuando una estructura guarda claves:valor y permite repeticiones
 (unordered\_multimap) se le suele denominar diccionario desordenados.

### Resumen

 Los conjuntos, mapas y diccionarios desordenados son eficientes en inserción, borrado y consulta.

- Para c++ la única diferencia entre un conjunto y un mapa es:
  - Si la clave y el valor son lo mismo, es un conjunto.
  - Si la clave es distinta al valor, es un mapa.

 No confundir la versión desordenada (unordered\_map por ejemplo) con la versión ordenada (map por ejemplo). Los set, multiset, map y multimap los veremos en el próximo tema.

# Estructura de Datos II

David Concha Gómez

Asignatura obligatoria

Segundo cuatrimestre

Créditos: 6

Moodle de la asignatura

Guía docente