

Welcome to Algorithms and Data Structures! - CS2100

Tablas hash

Qué son las tablas hash?

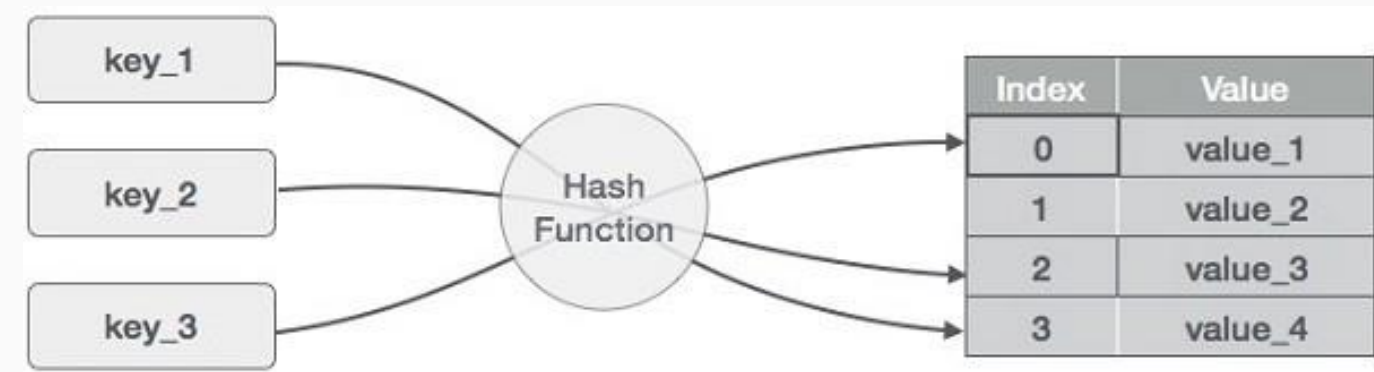
Son estructuras de datos similares a los arrays, con la diferencia que el índice puede ser cualquier tipo de dato comparable y no solo enteros

Restricciones:

- Claves únicas
- No sabemos en qué posición se almacenan los valores

La clave es mapeada a un índice:

```
int index = getIndex(key)  
array[index] = value
```



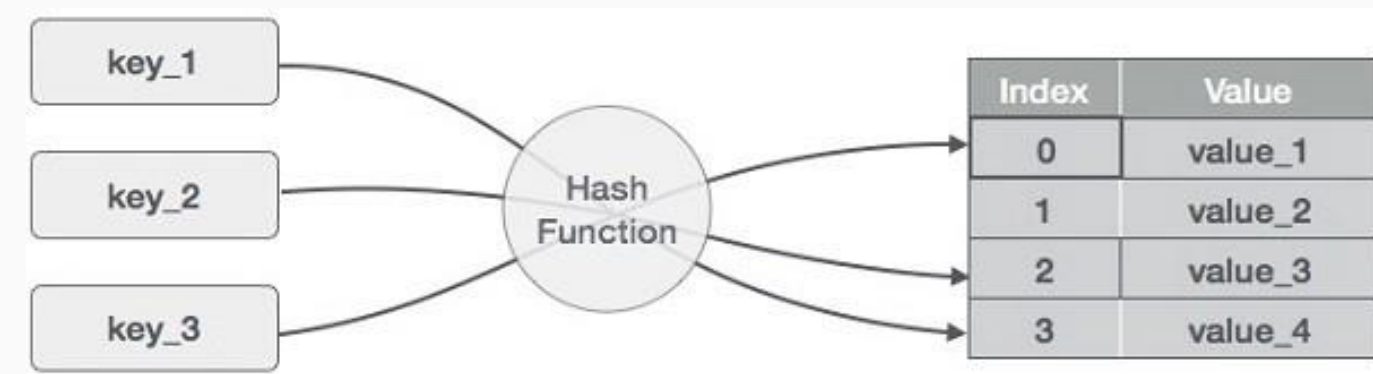
Hashing

Saben lo que es una función hash?

Es un método que nos permite obtener un índice en un array desde una key

Qué problemas podríamos tener con esa función?

- Procesar la key puede ser difícil para ciertos tipos de datos
- Manejar colisiones cuando dos keys tienen el mismo índice
- Espacio-tiempo tradeoff

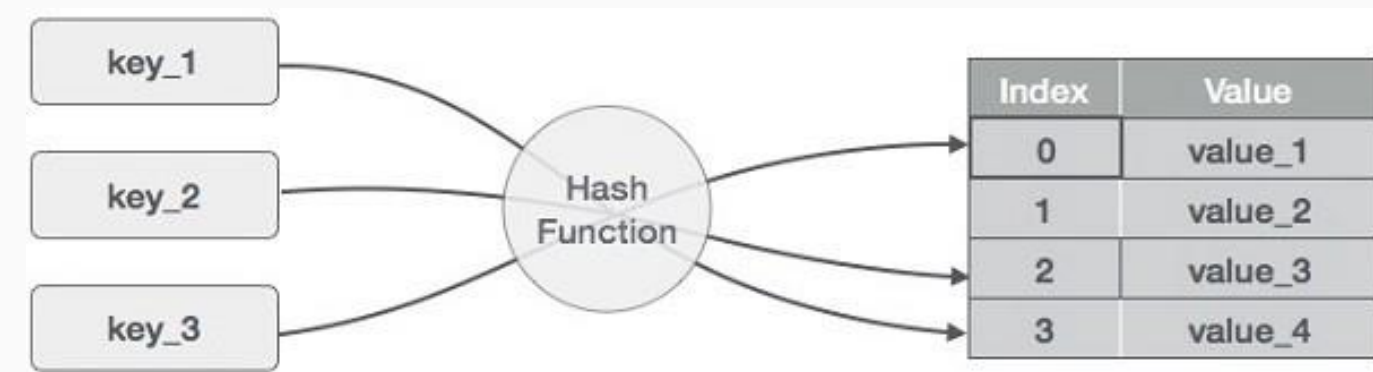


Hashing

Existen varios algoritmos para hashing, algunos retornan un output más grande que otros

Propiedades:

- **Estable:** El output siempre debería ser el mismo para el mismo input (invariante)
- **Uniforme:** Los valores hash deben ser distribuidos de manera uniforme (reducir colisiones)
- **Eficiencia:** Debe ser balanceado de acuerdo a las necesidades en espacio y tiempo
- **Seguridad:** Dado un valor hash, obtener un valor que me pueda generar ese valor hash no debería ser posible

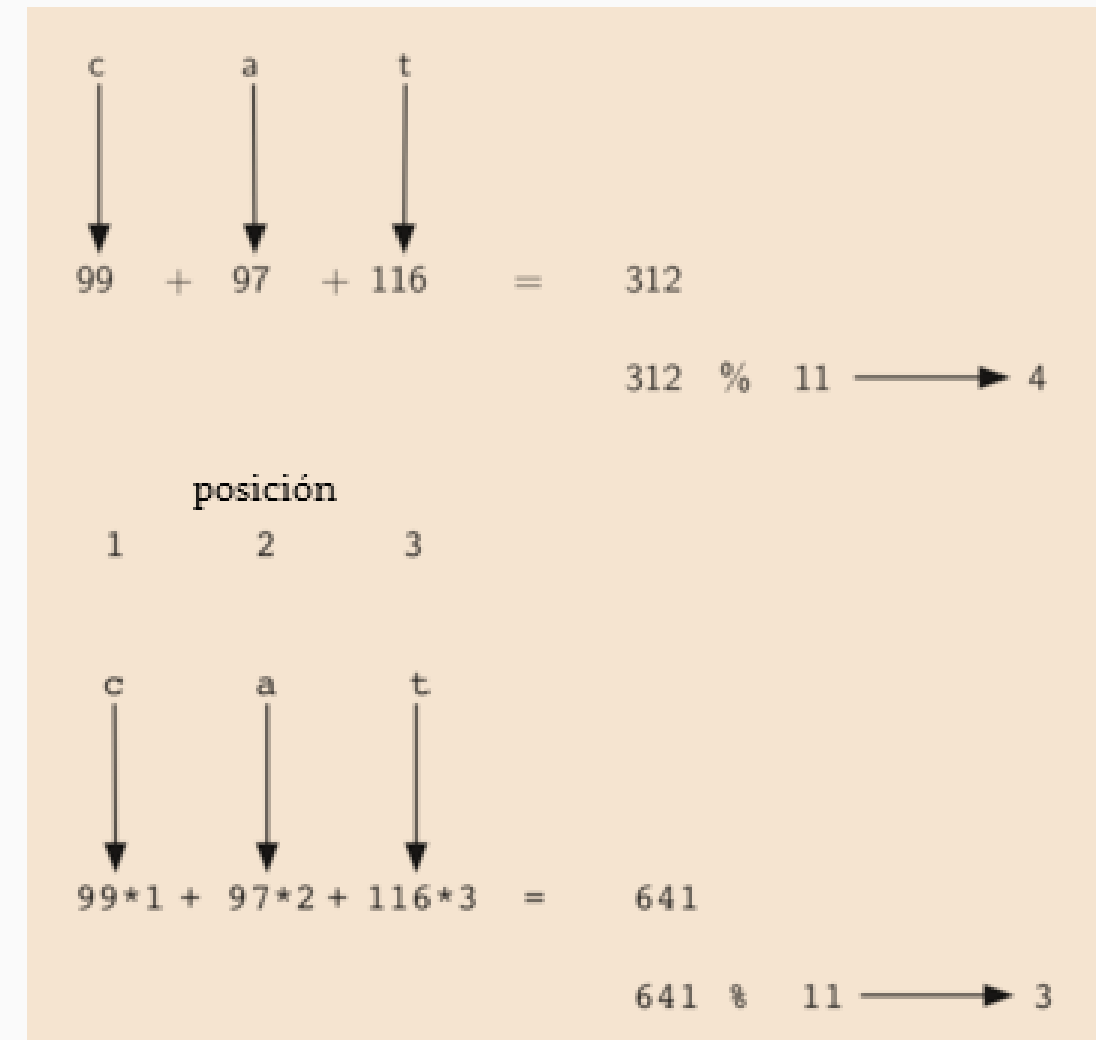


Hashing (strings)

- Implementación ingenua (suma):

```
int additiveHash(string key) {  
    int total = 0;  
    for (char& character : key) {  
        total += (int) character;  
    }  
    return total;  
}
```

Cuál es el problema de esta función hash?

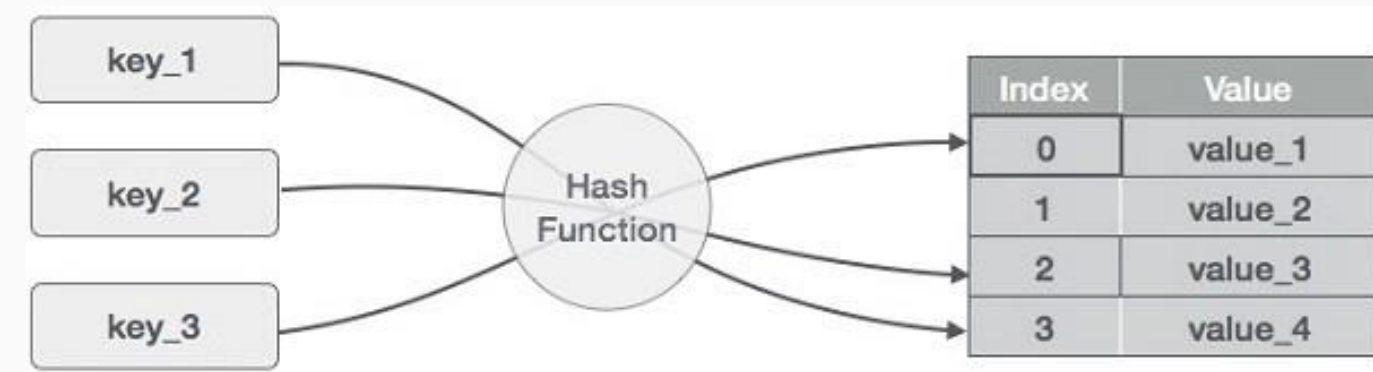


Hashing (strings)

- Algoritmo plegable:

```
const int getNextBytes(int startIndex, string str) {  
    int currentFourBytes = 0;  
    currentFourBytes += getByte(str, startIndex);  
    currentFourBytes += getByte(str, startIndex + 1) << 8;  
    currentFourBytes += getByte(str, startIndex + 2) << 16;  
    currentFourBytes += getByte(str, startIndex + 3) << 24;  
    return currentFourBytes;  
}
```

Cuál es el problema de esta función hash?



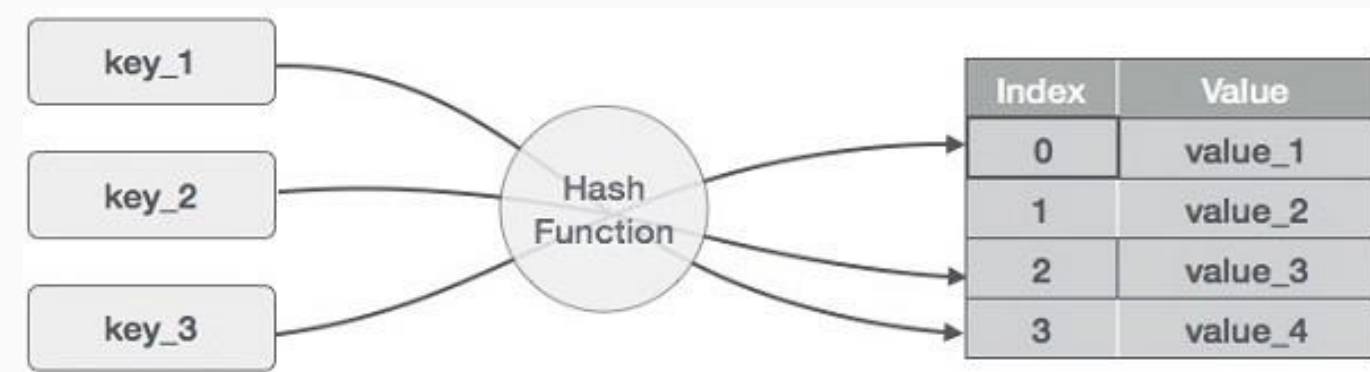
Hashing

No escriban sus propios algoritmos de hash! Utilicen las librerías ya hechas para esto

Escojan el algoritmo hash que encaje mejor con sus necesidades

MD5 solía ser un algoritmo hash bastante usado, tiene un buen balance entre estabilidad y uniformidad pero ya se han detectado colisiones.

SHA-2 es un algoritmo bastante seguro pero no es muy eficiente



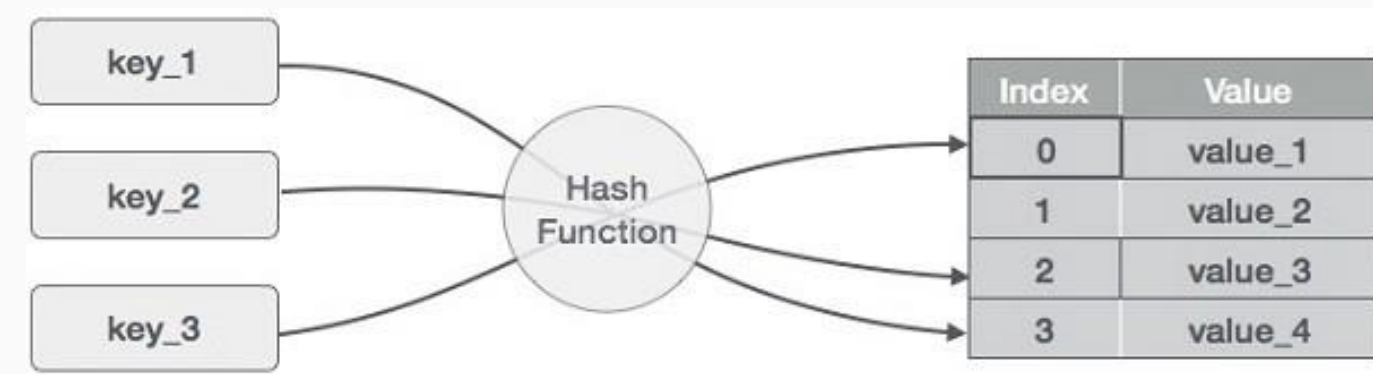
Agregando datos

Se inicia con un tamaño fijo del array, y se genera el hash code

```
int arrayLength = 10;  
int hashCode = getHashCode(key);  
int index = hashCode % arrayLength;  
array[index] = value
```

Qué problemas podrían aparecer?

Dos elementos diferentes con el mismo hash code, significa que se les asigna el mismo índice dentro del array

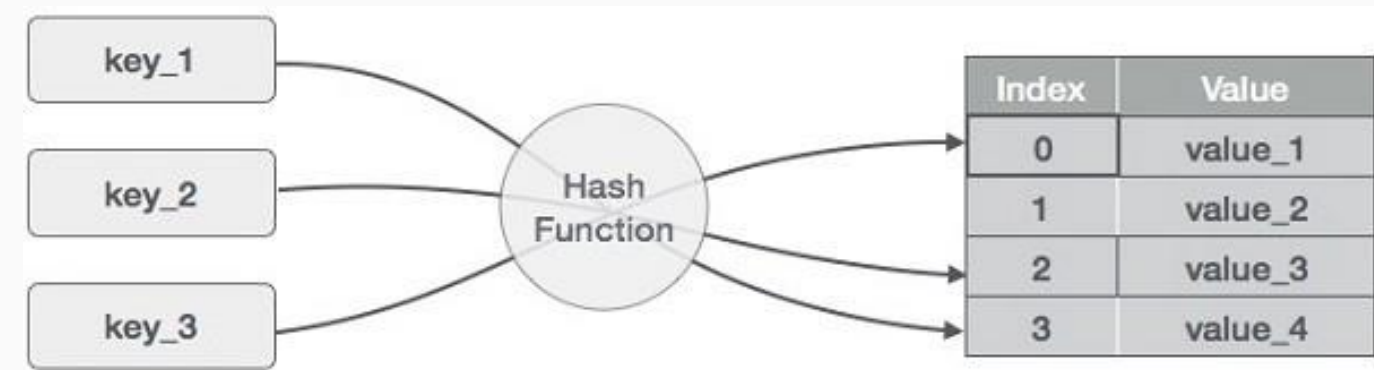


Manejando colisiones

Cómo podrían manejar colisiones?

Recuerden que el usuario de una tabla hash no debe saber que hubo una colisión

- Open addressing: Moverse al siguiente índice disponible
- Chaining: Almacenar los elementos en una lista enlazada



Manejando colisiones (open addressing)

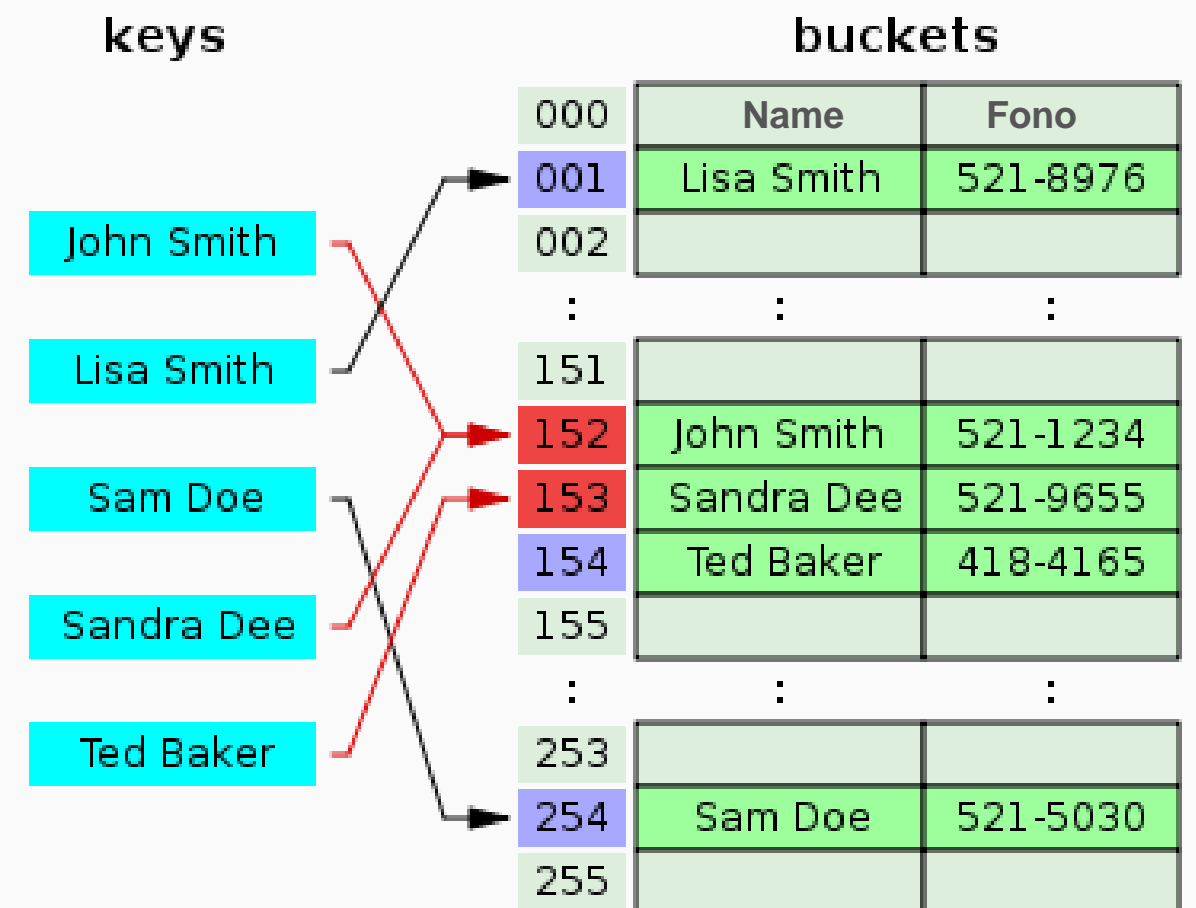
Si el espacio ya está ocupado, entonces se va a buscar el siguiente espacio disponible

```
int hashCode = hashCode(key);
int index = hashCode % arrayLength
while (array[index] != null){ // int/float/double is null??
    index++;
}
array[index] = value
```

Al momento de buscar el elemento se empieza a comparar el hash code hasta encontrar el elemento o estar seguros que no está en la tabla hash

Cómo impacta los agrupamientos (muchos hashcode iguales) ?

$\text{index} = (\text{hashCode} + \text{salto}) \% \text{arrayLength}$

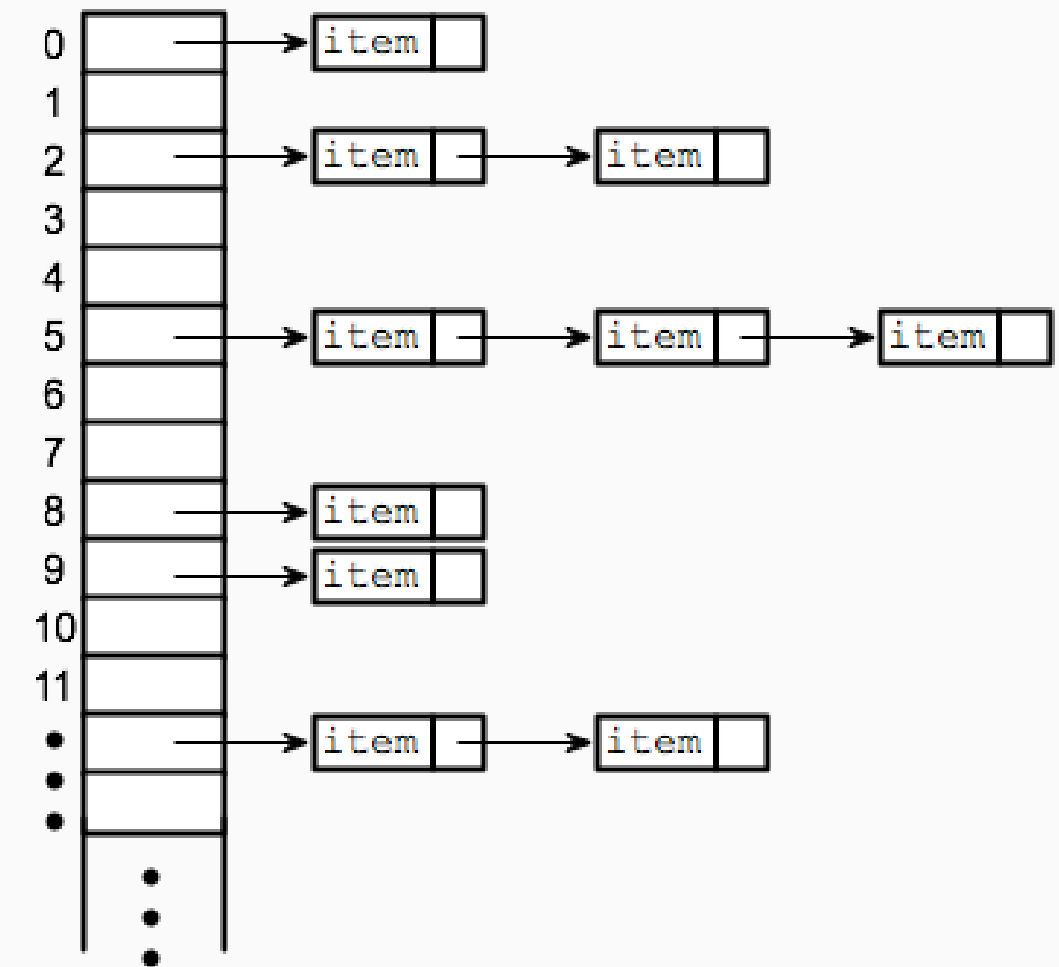


Manejando colisiones (chaining)

Si el espacio ya está ocupado, entonces se agrega el elemento a la lista enlazada

```
int hashCode = getHashCode(key);  
int index = hashCode % arrayLength  
array[index].push_back(value)
```

Para encontrar el elemento, solo se debe buscar dentro de la lista enlazada



Cómo impacta los agrupamientos (muchos hashcode iguales) ?

Haciendo crecer la estructura

La probabilidad de que hayan colisiones, va a ser proporcional a la cantidad de espacios disponibles del array

Para controlar esto, se debe tener un factor de llenado, que nos indicará cuán lleno está nuestro hash table.

fillFactor = # of elements / hash table length

Al agregar un elemento verificar:

```
if (fillFactor >= maxFillFactor)
    rehashing();
```

¿Cuánto cuesta el rehashing?

¿Rehashing con chaining?

Original Hash Table	
0	6
1	15
2	
3	24
4	
5	
6	13

After Inserting 23	
0	6
1	15
2	23
3	24
4	
5	
6	13

After Rehashing	
0	
1	
2	
3	
4	
5	
6	6
7	23
8	24
9	
10	
11	
12	
13	13
14	
15	15
16	

Manejando colisiones (open addressing)

Ejemplo:

Insertar los siguientes keys:

2, 3, 6, 8, 10, 11, 15, 19, 20, 22, 25, 30

Tamaño de la tabla hash: 5

FillFactor = 0,8

HashCode = Key

Manejando colisiones (open addressing)

Ejemplo:

Insertar los siguientes keys:

2, 3, 6, 8, 10, 11, 15, 19, 20, 22, 25, 30

Tamaño de la tabla hash: 5

maxFillFactor = 0,8

HashCode = Key

0	
1	6
2	2
3	3
4	8

N=10
fillFactor: 8/10

0	10
1	11
2	2
3	3
4	
5	15
6	6
7	
8	8
9	19

Manejando colisiones (chaining)

Ejemplo:

Insertar los siguientes keys:

2, 3, 6, 8, 10, 11, 15, 19, 20, 22, 25, 30

Tamaño de la tabla hash: 5

FillFactor = 0,8

HashCode = Key

Manejando colisiones (chaining)

Ejemplo:

Insertar los siguientes keys:

2, 3, 6, 8, 10, 11, 15, 19, 20, 22, 25, 30

Tamaño de la tabla hash: 5

maxFillFactor = 0,5

HashCode = Key

maxColision = 3

0	10,15
1	6,11
2	2
3	3,8,7
4	19

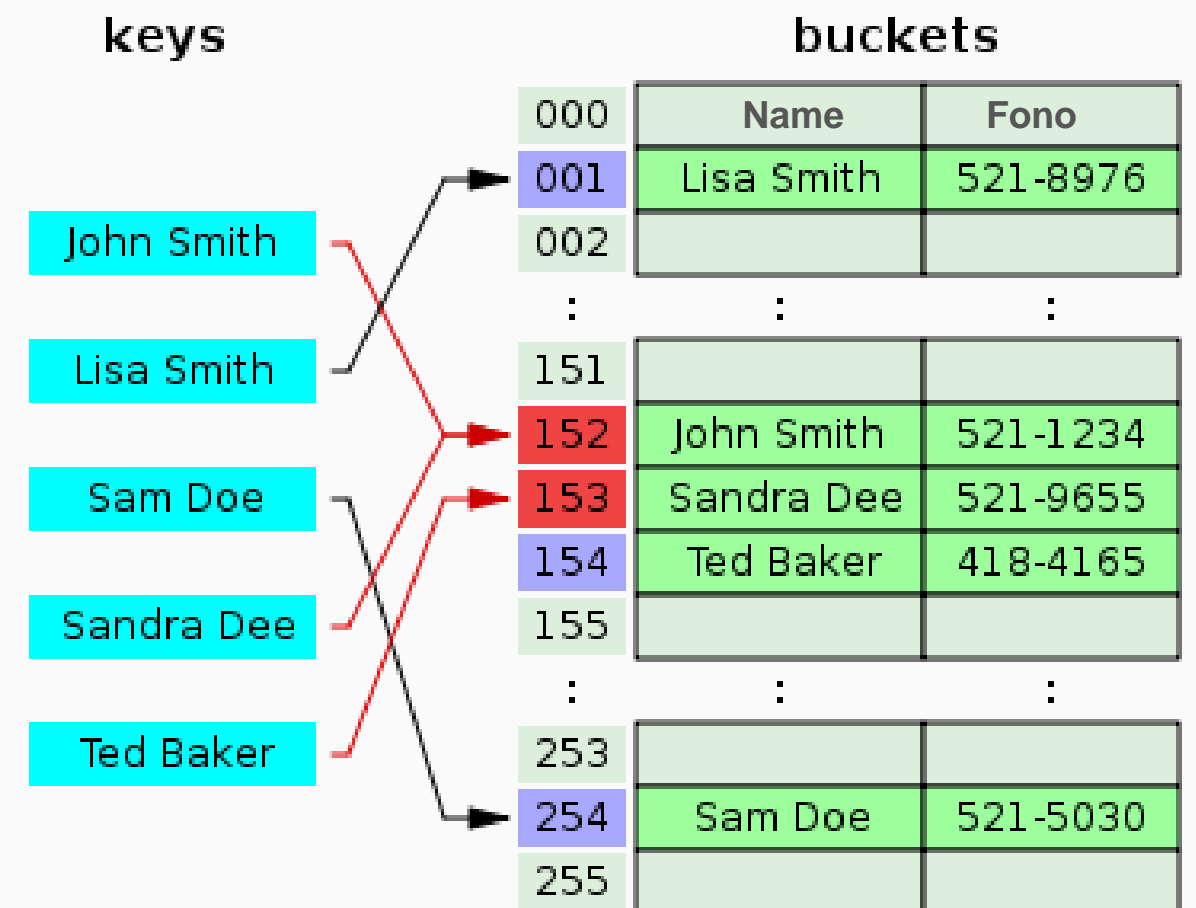
N=10
fillFactor: $10/(10*3)$

0	10,20
1	11
2	2,22
3	3
4	
5	15
6	6
7	
8	8
9	19

Encontrar elementos (open addressing)

Al buscar un elemento va a depender de como estemos manejando las colisiones:

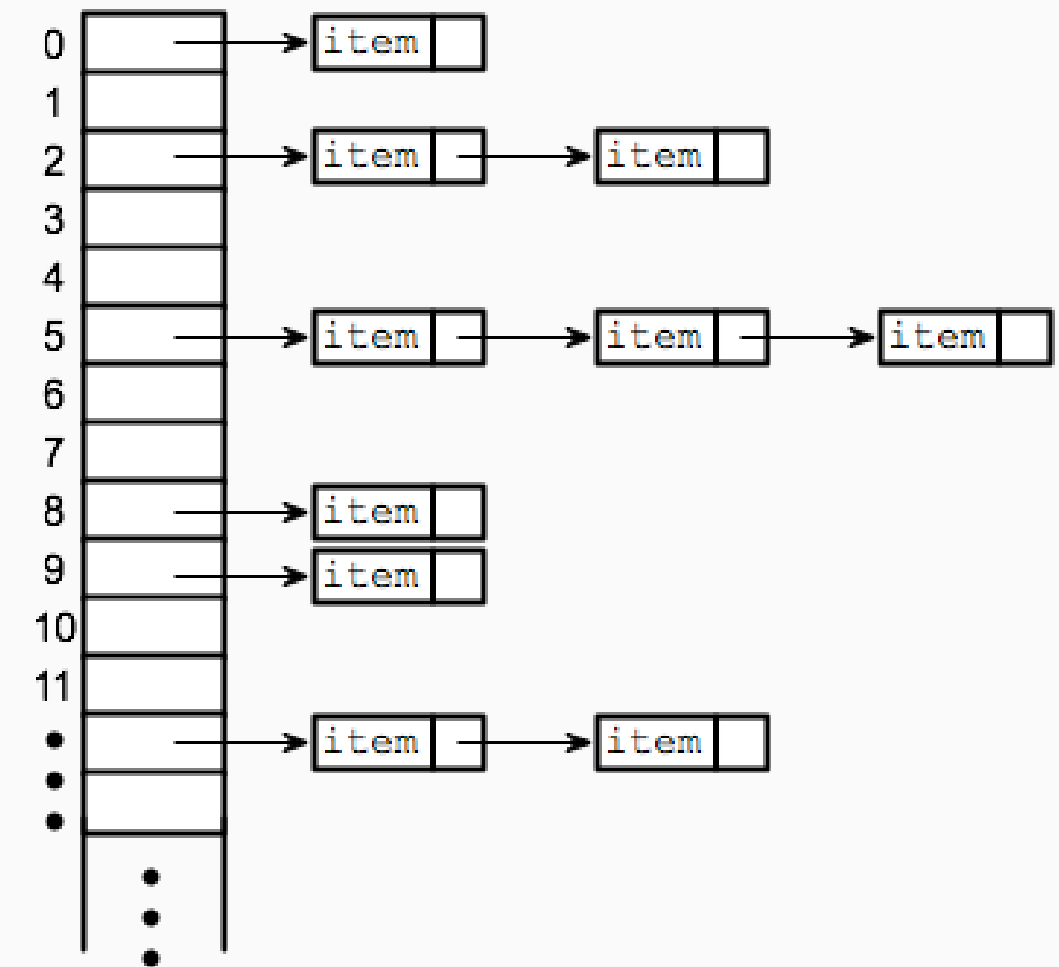
- Open addressing:
 - a. Se obtiene el índice de la key
 - b. Se verifica que no sea nulo
 - c. Si hay colisión, se sigue avanzando hasta encontrar el elemento



Encontrar elementos (chaining)

Al buscar un elemento va a depender de como estemos manejando las colisiones:

- Chaining:
 - a. Se obtiene el índice de la key
 - b. Se busca en la lista de la posición



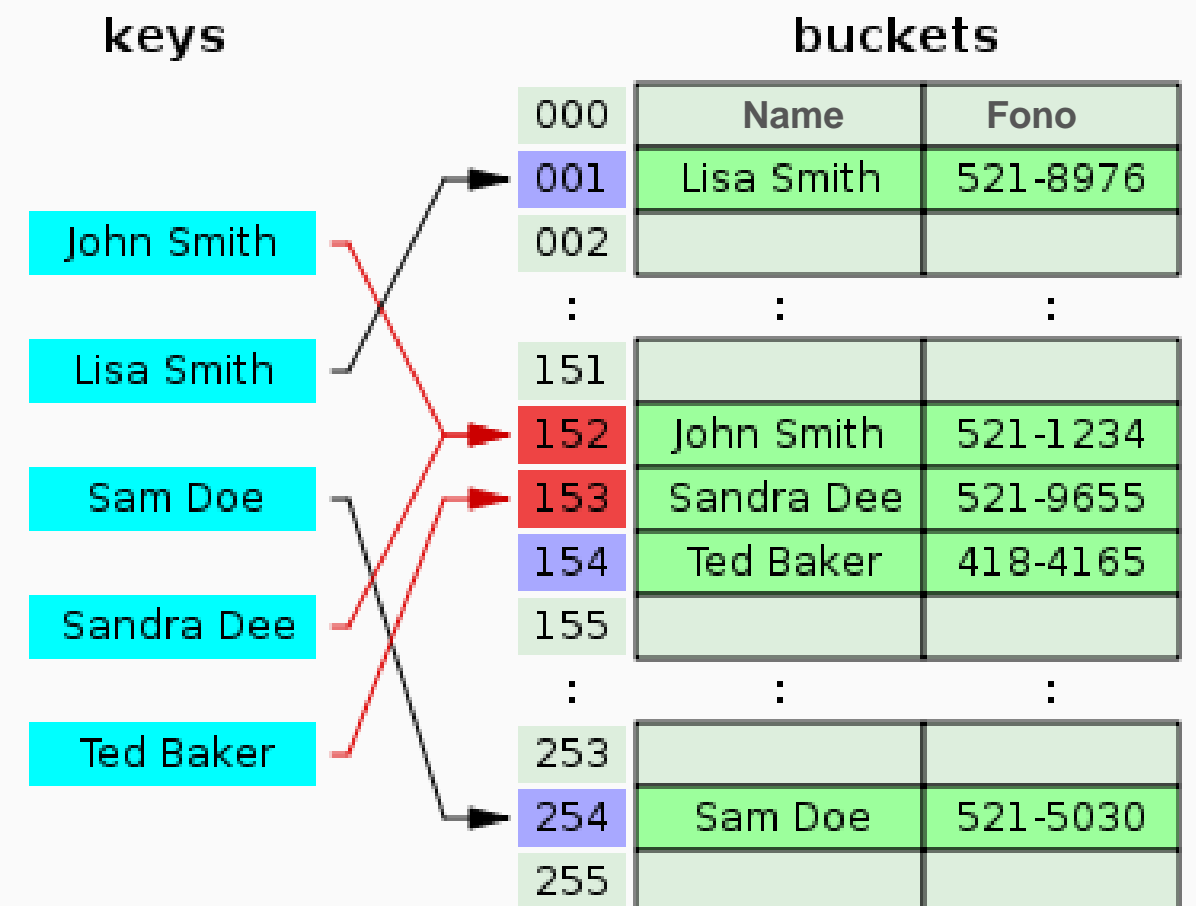
Remover elementos (open addressing)

El remover elementos va a depender de como estemos manejando las colisiones:

Open addressing:

1. Se obtiene el índice del elemento
2. Se verifica que no sea nulo en el array
3. Si las keys coinciden, remover el elemento
4. Si las keys no coinciden (probablemente haya una colisión), revisar el siguiente elemento

Es bastante difícil de mantener

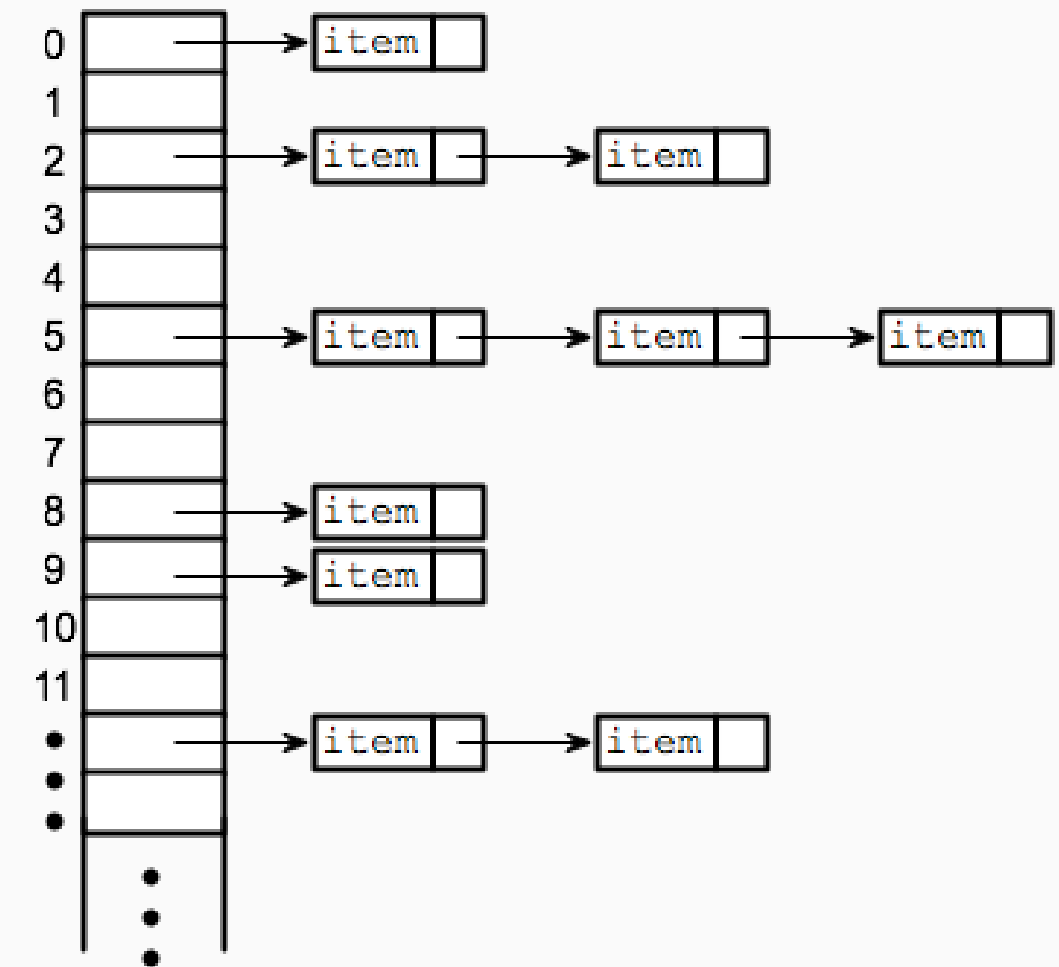


Remove elements (chaining)

Chaining:

1. Se obtiene el índice del elemento
2. Se elimina el elemento de la lista simplemente enlazada

Es fácil de mantener e implementar



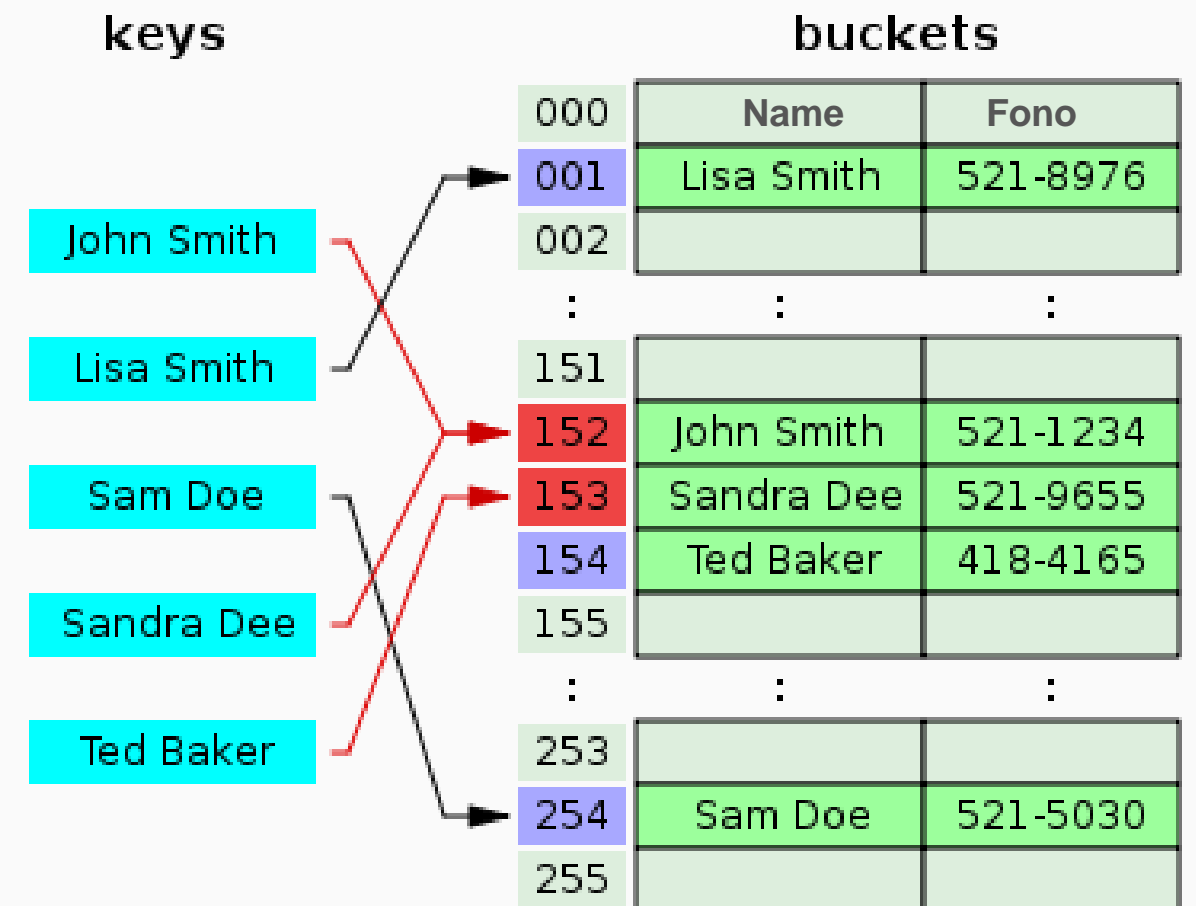
Enumerar elementos (open addressing)

Cuántas maneras de enumerar se deben tener?

Dos, una para las llaves y otra para los valores

- Open addressing:

```
foreach (item in array) {  
    if (item != null){  
        ....  
    }  
}
```



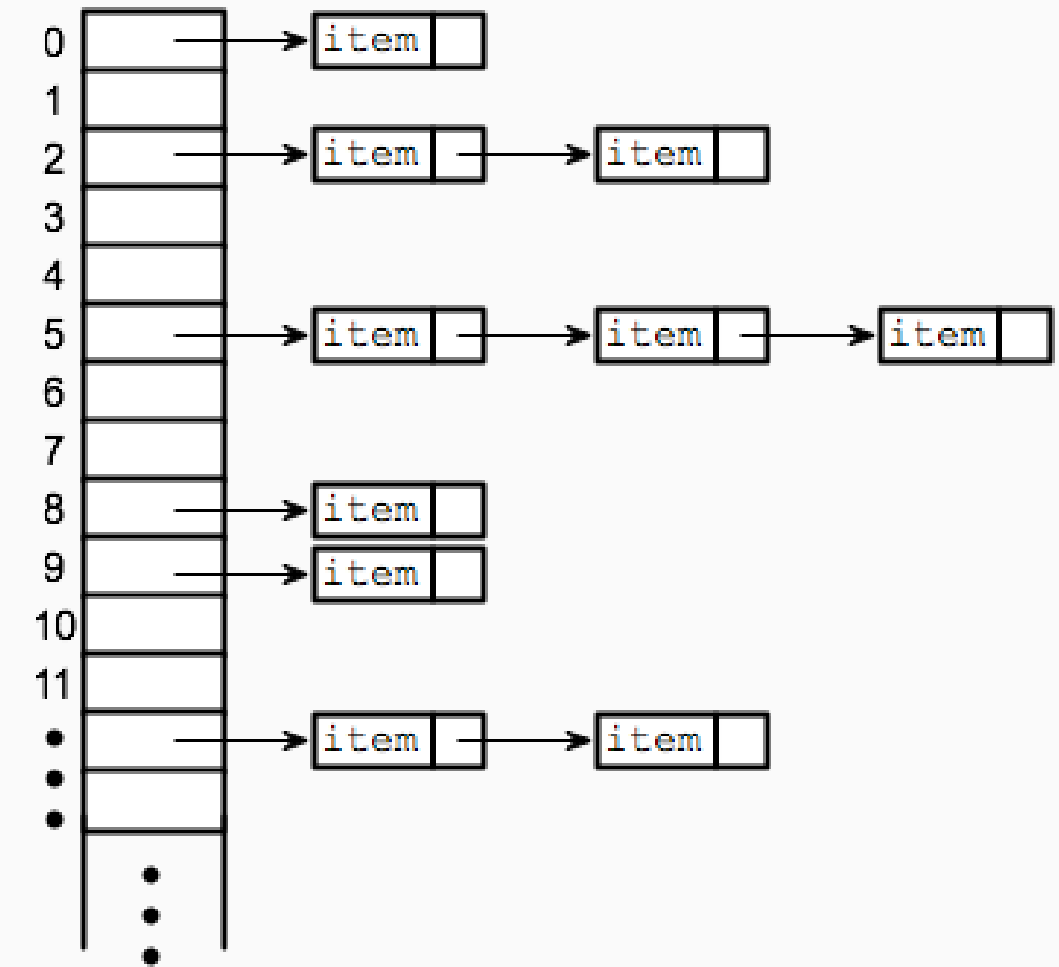
Enumerar elementos (chaining)

Cuántas maneras de enumerar se deben tener?

Dos, una para las llaves y otra para los valores

- Chaining:
foreach(list in array) {
 foreach(item in list) {

 }
}



Hashing (Complejidad Open Addressing)

Operación	Caso Mejor/promedio	Peor Caso
Inserción	$\Omega(1)$	$O(N)$
Búsqueda	$\Omega(1)$	$O(N)$
Eliminación	$\Omega(1)$	$O(N)$
Rehashing	$\Omega(N)$	$O(N)$
Traer Ordenado	$\Omega(N * \log(N))$	$O(N * \log(N))$

N: total de elementos

K: promedio de elementos en la colision $K \leq N$

Hashing (Complejidad Chaining)

Operación	Caso Mejor/promedio	Peor Caso
Inserción	$\Omega(1)$	$O(K)$
Búsqueda	$\Omega(1)$	$O(K)$
Eliminación	$\Omega(1)$	$O(K)$
Rehashing	$\Omega(N)$	$O(N) * O(k)$
Ordenar	$\Omega(N * \log(N))$	$O(N * \log(N))$

N: total de elementos

K: promedio de elementos en la colision $K \leq N$

Hashing (Complejidad)

¿Y si usamos un árbol binario de búsqueda balanceado?

```
Class NodeBTS{  
    TK key;  
    TV value;  
    NodeBTS *left;  
    NodeBTS* right;  
};
```

Operación	Caso Mejor/promedio	Peor Caso
Inserción	$\Omega(\log(N))$	$O(\log(N))$
Búsqueda	$\Omega(\log(N))$	$O(\log(N))$
Eliminación	$\Omega(\log(N))$	$O(\log(N))$
Rehashing	--	--
Ordenar	$\Omega(N)$	$O(N)$

Welcome to Algorithms and Data Structures! - CS2100