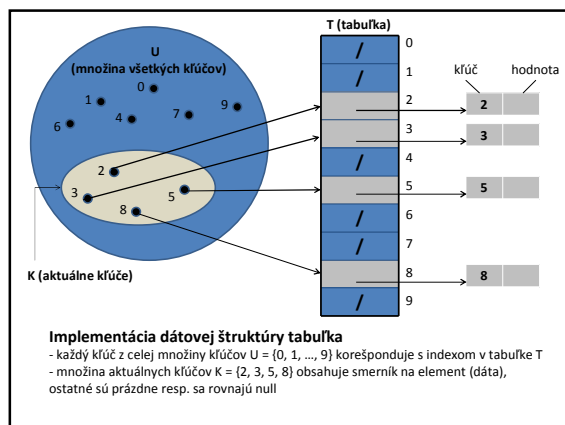


Tabuľka

Tabuľka

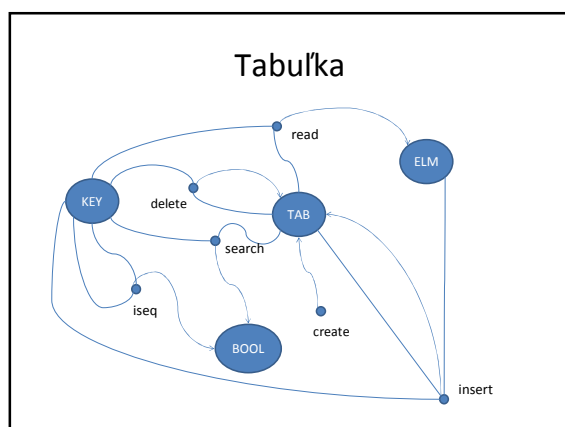
- Dátová štruktúra, ktorá asociuje hodnotu s kľúčom.
- V pamäti sa dáta uchovávajú ako dvojica kľúč – hodnota
- K hodnote sa pristupuje pomocou kľúča



Tabuľka – formálna špecifikácia

- Druhy: TAB, ELM, KEY, BOOL
- Operácie:
 - CREATE() → TAB vytvorenie prázdnej tabuľky
 - INSERT(tab, key, elem) → TAB vloženie prvku
 - READ(tab, key) → ELM výber prvku
 - DELETE(tab, key) → TAB vymazanie prvku
 - ISEQ(key, key) → BOOL porovnanie 2 prvkov
 - SEARCH(key, tab) → BOOL test, či sa v tabuľke nachádza prvok

Tabuľka



Slovník

0	AD	3
1	F8	4
2		
3	BC	5
4	E9	1
5	CB	0

kľúč Smerník na synonymum

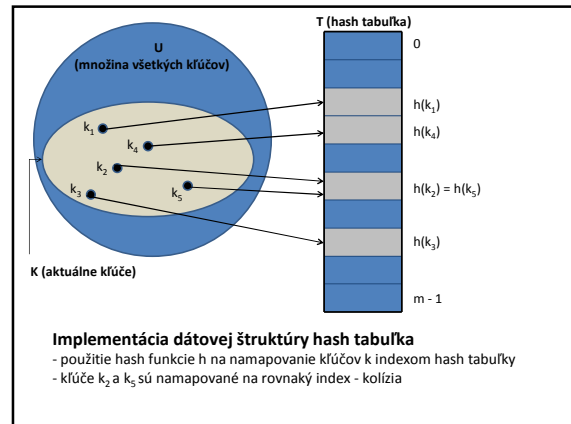
$$f(AD) = f(BC) = f(CB)$$

$$f(F8) = f(E9)$$

Každé slovo obsahuje aj smerník na synonymum. Všetky synonymá potom môžeme získať postupným prechádzaním vzniknutého spájaného zoznamu

Hash tabuľka

- Tabuľka, kde sa kľúč vypočíta pomocou špeciálnej (hešovacej) funkcie
- Príklad jednoduchkej hešovacej funkcie:
 - Kľúč = súčet ascii hodnôt jednotlivých znakov
 - Value = OKNO
$$\text{OKNO} \rightarrow \text{key} = 117(\text{O}) + 113(\text{K}) + 116(\text{K}) + 117(\text{O}) = 463$$



Hešovacia funkcia

- Základné vlastnosti hešovacej funkcie:
 - Výpočet by nemal byť náročný
 - Mala by byť navrhnutá tak, aby vznikalo čo najmenej kolízií

Aká má byť $h(k)$?

- Pre tabuľku s m položkami:
- $\sum P(k) = 1/m$ pre $j=0,1,\dots,m-1$
 - suma cez k : $h(k)=j$
 - ale zvyčajne nepoznáme $P(k)$
- Ak poznáme $P(k)$, napr.:
 - Kľúče sú náhodné reálne čísla nezávislo rovnomerne rozdelené v intervale $0 \leq k < 1$,
 - $h(k) = \text{floor}(k.m)$

Ako navrhnuť $h(k)$?

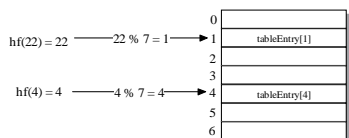
- Pomocou heuristík, opierajúc sa o kvalitatívne znalosti o P .
- Tabuľka symbolov v prekladači:
 - vieme, že často sa vyskytujú v programe symboly, ktoré sa len málo líšia, napr: refA, refB
 - Dobrá $h(k)$ by mala minimalizovať prípady, že takéto symboly pošle na rovnaké miesto v tabuľke.

Návrh $h(k)$ metódou delenia

- $h(k) = k \bmod m$
 - $m = 16, k = 36, h(36) = 4$.
- m nemá byť mocnina 2. Prečo?
 - $m = 2^p$: $h(k)$ závisí len od dolných p bitov kľúča
- m nemá byť mocnina 10. Prečo?
- m má byť prvočíslo nie blízke nejakej mocnine 2.

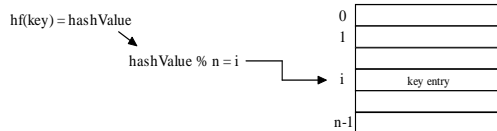
Hešovacia funkcia - príklad

- $hf(x) = x$, kde x je kladné číslo.
- Tabuľku predstavuje poľa s veľkosťou 7



Hešovanie

Hash Value: $hf(key) = hashValue$
 HashTable index: $hashValue \% n$



Hešovanie

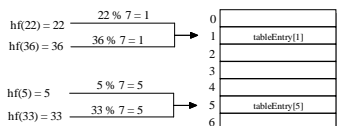
- Hešovacia tabuľka rozdeľuje elementy do série zlinkovaných listov označovaných ako sektory
- Hešovacia funkcia mapuje hodnotu na index v tabuľke. Funkcia poskytuje podobný prístup k elementu ako index v poli.

Hešovanie

- Hešovaciú tabuľku predstavuje pole referencií, kde
 - Hešovacia funkcia má kľúč ako argument a vráti integer hodnotu
 - Výpočtom zvyšku po delení (delí sa veľkosťou tabuľky) získame mapovanie kľúčov na indexy v tabuľke

Hešovacia funkcia - príklad

- Pri zvolenej hešovacej funkcii nastáva kolízia pri každých 2 kľúčoch, ktoré sa navzájom líšia o násobok veľkosti tabuľky
- $36 - 22 = 14 = 2 * 7 \rightarrow$ nastane kolízia



Jednoduchá hešovacia funkcia

- Častokrát je kľúčom reťazec (string)
 - Vo funkcii môžeme kombinovať sekvenciu znakov z reťazca

```
public int hashCode()
{
    int hash = 0;

    for (int i = 0; i < n; i++)
        hash = 31*hash + s[i];

    return hash;
}
```

Jednoduchá hešovacia funkcia

Výpočet hašovacej hodnoty pre 3 rôzne reťazce:

Hodnota pre strB je negatívna kôli pretečeniu

```
String strA = "and", strB = "uncharacteristically",  
      strC = "algorithm";  
  
hashValue = strA.hashCode();    // hashCode = 96727  
hashValue = strB.hashCode();    // hashCode = -2112884372  
hashValue = strC.hashCode();    // hashCode = 225490031
```

- ak nastane prípad, že hešovacia funkcia vráti záporné číslo, pribudne problém pri práci s poľom (záporný index neexistuje).

-nasledujúci výpočet však zabezpečí, že index bude vždy kladné číslo:
 $tableIndex = (hashValue \& Integer.MAX_VALUE) \% tableSize$

Hešovacia funkcia - príklad

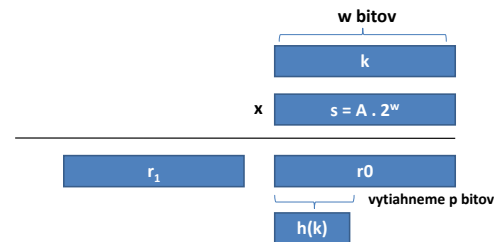
- Výpočet hešovacej funkcie pre produkt, ktorý je charakterizovaný sériovým číslom

```
public class Product  
{  
    private int serialNum;  
    public int hashCode()  
    {  
        long hashValue = serialNum;  
        hashValue *= hashValue;  
        return (int)(hashValue % Integer.MAX_VALUE);  
    }  
}
```

Návrh $h(k)$ metódou násobenia

1. Vynásobiť kľúč k zvolenou konštantou A , $0 < A < 1$ a vybrať zlomkovú časť z $k.A$.
 2. Vynásobiť túto hodnotu m a vziať celú časť.
 $h(k) = \text{floor}(m \cdot (k.A \bmod 1))$
kde $k.A \bmod 1$ označuje zlomkovú časť z $k.A$, t.j.
 $k.A - \text{floor}(k.A)$
- voľba A : podľa Knutha približne
 $\text{sqr}(5) - 1 = 0.6180339887$

Hešovacia funkcia - príklad



Násobenie ako metóda hešovania

- w -bitová reprezentácia kľúča sa vynásobí w -bitovou hodnotou $s = A \cdot 2^w$
- p najvyšších bitov nižšej w -bitovej polovice výsledku zvolíme ako hash hodnotu funkcie $h(k)$

kolízia

- Keď hešovacia hodnota dvoch (alebo viacerých) elementov ukazuje na rovnaké miesto v tabuľke nastáva kolízia. Dva elementy nemôžu byť uložené na rovnakej pozícii v tabuľke.
- Možnosti riešenia problému:
 - Umiestnenie jedného z kolidujúcich elementov na inú pozíciu v tabuľke (*linear probing*)
 - Navrhnutie takej štruktúry, ktorá bude schopná uchovávať viacero elementov s rovnakou hešovacou hodnotou (sekvencia spájaných zoznamov)

Otvorené adresovanie

- Všetky prvky sa ukladajú priamo v tabuľke
- Každá položka tabuľky obsahuje buď prvok reprezentovanej dynamickej množiny (slovníka) alebo NIL.
- Nič sa nezreťazuje, nič nie je mimo tabuľky.
- Miera naplnenia tabuľky je vždy najviac 1.

Vkladanie pri otvorenom adresovaní

- Postupne sa skúša nájsť prázdne miesto
- Nie postupne ďalšieho a ďalšieho suseda $O(m)$
- Postupnosť skúšaných miest závisí od kľúča, t.j. hešovací funkcia dostane ďalší parameter $h(k,0), h(k,1), \dots, h(k,m-1)$
- Táto postupnosť m miest/adries musí byť permutáciou
- $0,1,\dots,m-1$

Hash insert(T, k)

```
insert(T table, k key)
{
    i ← 0
    repeat j ← h(k, i)
        if T[j] = NIL
            then T[j] ← k
            return j
        else i ← i + 1
    until i = m
    error "hash table overflow"
```

Hash search(T, k)

```
search(T table, k key)
{
    i ← 0
    repeat j ← h(k, i)
        if T[j] = k
            then return j
        else i ← i + 1
    until T[j] = NIL or i = m
    return NIL
}
```

Lineárne skúšanie (linear probing)

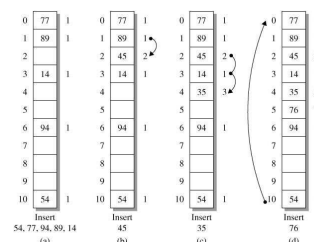
- Vloženie prvku:
 - Na začiatku sa všetky bunky tabuľky označia ako prázdne
 - Aplikuje sa hešovací funkcia a zvyšok po delení tejto hodnoty veľkosťou tabuľky predstavuje index v tabuľke. Ak je bunka prázdna, vloží sa do nej element
 - Inak sa postupne prehľadávajú ďalšie bunky v poradí a element sa vloží do prvej voľnej.

Lineárne skúšanie (linear probing)

- Uvažujme bežnú hešovaciu funkciu $h'(k): U \rightarrow \{0,1,\dots,m-1\}$
- $h(k,i) = (h'(k) + i) \bmod m$ pre $i = 0,1,\dots,m-1$
- Problém: strapce

Linear Probing

tableIndex = $x \% 11$



- vloženie prvkov na ich príslušné pozície bez posunov
- prvok 45 sme museli posunúť o jedno miesto kvôli obsadenému miestu 1
- posunutie 35 až o dve miesta (4 namiesto 2)
- posunutie prvku 76 až na miesto 5 namiesto pôvodného miesta 10

Linear Probing - algoritmus

```
//výpočet indexu tabuľky
int index = (item.hashCode() & Integer.MAX_VALUE) % n

// uloženie pôvodného indexu
int origIndex = index;

//cyklické prehľadávanie tabuľky a hľadanie voľnej pozície
// nájde miesto alebo sa tabuľka zaplní (origIndex == index).
do
{
    if table[index] is empty
        insert item in table at table[index] and return
    else if table[index] matches item
        return
    // posunutie v tabuľke
    index = (index+1) % n;
}
while (index != origIndex);
throw new BufferOverflowException();
```

Linear Probing

- Táto metóda je vhodná v prípade, ak veľkosť tabuľky je rádovo väčšia ako počet elementov, ktoré sa do nej budú vkladať.
- Dobrá hešovací funkcia minimalizuje kolízie a ak aj nastanú, tak v tabuľke bude dostatok voľných miest pre vyhľadanie náhradnej pozície

kvadratické skúšanie (quadratic probing)

- Uvažujme bežnú hešovaciu funkciu $h'(k): U \rightarrow \{0, 1, \dots, m-1\}$
- $h(k, i) = (h'(k) + c_1 \cdot i + c_2 \cdot i^2) \bmod m$ pre $i = 0, 1, \dots, m-1$
- Problém: ako zvoliť c_1 a c_2
- Strapce vznikajú sekundárne, menej

dvojité hešovanie

- Uvažujme bežné hešovacie funkcie $h_1(k): U \rightarrow \{0, 1, \dots, m-1\}$
 $h_2(k): U \rightarrow \{0, 1, \dots, m-1\}$
- $h(k, i) = (h_1(k) + i \cdot h_2(k)) \bmod m$ pre $i = 0, 1, \dots, m-1$

Hešovacia funkcia - príklad

Vkladanie pomocou dvojitého hešovania

- máme hešovaciu tabuľku s veľkosťou 13 s $h_1(k) = k \bmod 13$ a $h_2(k) = 1 + (k \bmod 11)$

- keď $14 \equiv 1 \pmod{13}$ a $14 \equiv 3 \pmod{11}$, tak potom kľúč 14 je vložený na prázdne miesto 9, potom ako sa zistilo, že miesta 1 a 5 sú obsadené

0	
1	79
2	
3	
4	69
5	98
6	
7	72
8	
9	14
10	
11	50
12	

Hešovacia funkcia - príklad

Použitie perfektného hešovania (perfect hashing) na uloženie množiny $K = \{10, 22, 37, 40, 60, 70, 75\}$

T	m_0	a_0	b_0	S_0
0	1	0	0	10
1	/			
2	4	10	18	60 75 / /
3	/			
4	/			
5	1	0	0	70
6	/			
7	9	23	88	40 37 / / / / /
8	/			

Hešovacia funkcia - príklad

Použitie perfektného hešovania (perfect hashing) na uloženie množiny $K = \{10, 22, 37, 40, 60, 70, 75\}$

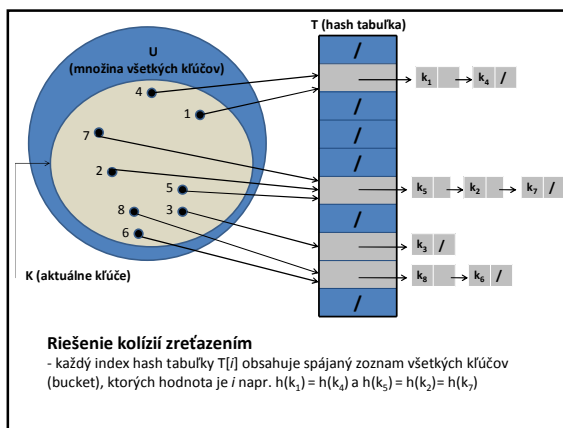
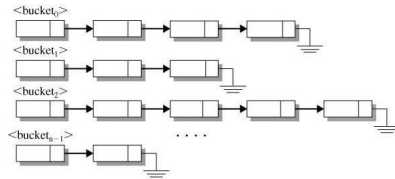
- Základná hešovací funkcia je $h(k) = ((ak + b) \bmod p) \bmod m$, kde $a \neq 0$, $b = 3$, 42 , $p = 101$ a $m = 9$
- Príklad: $h(75) = 2$, takže objekt 75 sa uloží na miesto s kľúčom 2
- Sekundárna hešovacia tabuľka S_j obsahuje všetky kľúče hešujúce index j
 - jej veľkosť je m_j
 - a hešovacia funkcia je $h_j(k) = ((a_jk + b_j) \bmod p) \bmod m_j$

-Keď $h_2(75) = 1$, kľúč 75 je uložený na miesto 1 v sekundárnej hešovacej tabuľke S2

-Takto nie sú žiadne kolízie v sekundárnych hešovacích tabuľkách
a vyhľadávanie trvá v najhoršom prípade konštantný čas

Sekvencia spájaných zoznamov

- V tomto prípade definujeme hešovaciú tabuľku ako indexovanú sekvenciu voľne spájaných zoznamov.
- Každý spájaný zoznam sa nazýva bucket a uchováva elementy s rovnakým indexom (rovnakou hešovacou hodnotou)

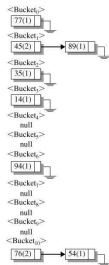


Sekvencia spájaných zoznamov

- Vloženie prvku
 - Hešovacou funkciou sa zistí index bucketu
 - Ak je bucket prázdny vloží sa prvok na prvú pozíciu
 - Ak nie je bucket prázdny, najskôr sa celý prehľadá, či sa v ňom element už nenachádza. Ak nie, tak sa vloží na začiatok.

Sekvencia spájaných zoznamov

Vloženie prvkov: {54, 77, 94, 89, 14, 45, 35, 76}
Veľkosť tabuľky je 11.

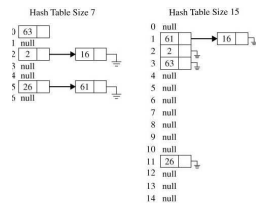


Sekvencia spájaných zoznamov

- V prípade väčšieho počtu kolízií je to výhodnejšia metóda
- Nie je obmedzená pevným počtom prvkov (resp. tak ako je obmedzené pole)

Rehašovanie

- So zvyšujúcim počtom prvkov v hešovacej tabuľke (a zároveň so zvyšujúcim počtom kolízií) sa efektivita vyhľadávania znižuje.
- Rehašovanie predstavuje zväčšenie veľkosti tabuľky, ak súčasná je zaplnená do určitej úrovne.



Hešovacia tabuľka – implementácia

```

Bucket:

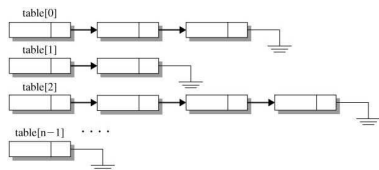
private static class Entry<T>
{
    // hodnota v hešovacej tabuľke
    T value;
    int hashCode;
    Entry<T> next;

    Entry(T value, int hashCode, Entry<T> next)
    {
        this.value = value;
        this.hashCode = hashCode;
        this.next = next;
    }
}

```

Hash Class Konštruktor

- Konštruktor vytvorí 12 prázdnych bucketov.



Hash Class

```

public class Hash<T> implements Collection<T>
{
    private Entry[] table;
    private int hashCode;
    private final double MAX_LOAD_FACTOR = .75;
    private int tableThreshold;

    private int modCount = 0;

    public Hash()
    {
        table = new Entry[17];
        hashCode = 0;
        tableThreshold =
            (int)(table.length * MAX_LOAD_FACTOR);
    }
    ...
}

```

Hash Class add()

- Pri vložení sa inkrementuje hashCode a modCount
- If hashCode ≥ tableThreshold nastane rehash().
- Veľkosť novej tabuľky bude:

$$2 * \text{table.length} + 1$$

Hash add() - 1

```

public boolean add(T item)
{
    int hashCode = item.hashCode() &
        Integer.MAX_VALUE,
        index = hashCode % table.length;
    Entry<T> entry;

    entry = table[index];

    // zistenie, či sa hodnota už v zozname nenachádza
    while (entry != null)
    {
        if (entry.value.equals(item))
            return false;

        entry = entry.next;
    }
}

```


Hash add() - 2

```
modCount++;

entry = new Entry<T>(item, hashValue,
    (Entry<T>)table[index]);

// vloží hodnotu
table[index] = entry;
hashTableSize++;

if (hashTableSize >= tableThreshold)
    rehash(2*table.length + 1);

return true;
}
```

Hash rehash() - 1

```
private void rehash(int newTableSize)
{
    Entry[] newTable = new Entry[newTableSize],
        oldTable = table;
    Entry<T> entry, nextEntry;
    int index;

    for (int i=0; i < table.length; i++)
    {
        entry = table[i];
```

Hash rehash() - 2

```
if (entry != null)
{
    do
    {
        nextEntry = entry.next;

        index = entry.hashValue % newTableSize;

        entry.next = newTable[index];
        newTable[index] = entry;
    }
```

Hash rehash() - 3

```
        entry = nextEntry;
    } while (entry != null);
}

table = newTable;

tableThreshold =
    (int)(table.length * MAX_LOAD_FACTOR);

oldTable = null;
}
```

Hash remove() - 1

```
public boolean remove(Object item)
{
    int index = (item.hashCode() &
        Integer.MAX_VALUE) % table.length;
    Entry<T> curr, prev;

    curr = table[index];
    prev = null;

    while (curr != null)
        if (curr.value.equals(item))
        {
            modCount++;
```

Hash remove() - 2

```
if (prev != null)
    prev.next = curr.next;
else
    table[index] = curr.next;

hashTableSize--;

return true;
}
else
{
    prev = curr;
    curr = curr.next;
}
return false;
}
```