

Suchen, Maps und Hashing



- Sie wissen, wie in Strings gesucht werden kann
- Sie kennen den Begriff der Invarianten
- Sie wissen, wie binäres Suchen funktioniert
- Sie wissen, wie Hashing funktioniert
- Sie kennen das Java Interface Map und die Implementation HashMap

Suchen

Beispiele, wo gesucht werden muss:

- Prüfen, ob ein Wort in einem Text vorkommt
- Zählen, wie oft ein Wort in einem Text vorkommt
- Überprüfen, ob alle Worte korrekt geschrieben sind
- Finden einer Telefonnummer in einem Telefonbuch
- Zählen, wie oft die richtige Zahlenkombination im Lotto gewählt wurde
- Prüfen, ob eine Kreditkartennummer gesperrt ist
- Prüfen, ob in zwei Listen die gleichen Elemente vorkommen

Finden eines Teilstrings in einem String

D r e i r e i n e r e i s e n d e

r e i s

r e i s

...

r e i s

r e i n

r e i n

...

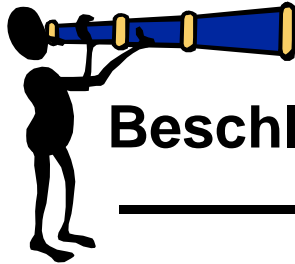
r e i n

- `int indexOf(String str, String pattern)`
 - liefert die Position, an der das Muster beginnt
 - -1 falls das Muster nicht vorkommt

Implementation in Java

```
int indexOf(String str, String pattern) {  
    for (int i = 0; i < str.length() - pattern.length() + 1; i++) {  
        int k;  
        for (k = 0;  
             k < pattern.length() && str.charAt(i+k) == pattern.charAt(k);  
             k++);  
        if (k == pattern.length()) return i;  
    }  
    return -1;  
}
```

- Muster wird an die Position i gesetzt
- Es wird mit dem String verglichen bis
 - Ende des Musters erreicht -> Erfolg
 - Nichtübereinstimmung
- Aufwand ist $O(\dots)$



Beschleunigte Stringsuche: Knuth-Morris-Pratt (KMP)

- Naive Stringsuche ist vor allem aufwändig, wenn grosse Teilübereinstimmungen vorkommen.
- Idee von KMP: bei der Suche wird Information gesammelt, die für die weitere Suche verwendet werden kann.
- Aufwand $O(m + n)$
 - aber: Muster muss vorher analysiert (compiliert) werden
 - Suche muss als Automat implementiert werden
- Optimierung von <5% der Fälle (wenn Teilübereinstimmung mit Muster)
- ➔ theoretisch interessant aber praktisch irrelevant

Beispiel

m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
P: ABCDABD
i: 0123456

- index in string
- string
- pattern
- index in pattern

Mismatch at S[3] → no 'A' found → set m=4

m: 01234567890123456789012
S: ABC ABCD**AB** ABCDABCDABDE
P: ABCD**ABD**
i: 0123456

Mismatch at S[10] → passed 'AB' →
set m=8 → P[0..1] need not to be checked

m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
P: ABCDABD
i: 01**23456**

Mismatch at S(10) → no 'C' found → set
m=11

m: 01234567890123456789012
S: ABC ABCDAB **ABCDAB**CDABDE
P: **ABCDABD**
i: 0123456

Mismatch at S[17] → passed 'AB' → set
m=15 → P[0..1] need not to be checked

m: 01234567890123456789012
S: ABC ABCDAB ABCD**ABCDAB**DE
P: **ABCDABD**
i: 01**23456**

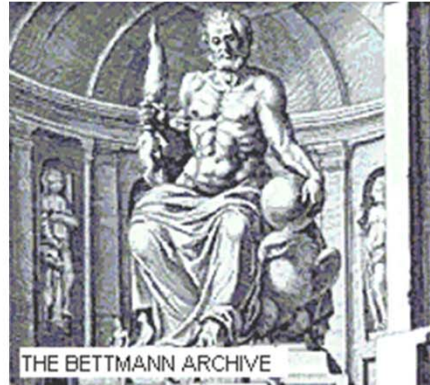
match

aus Wikipedia

Suche in zwei Arrays

Griechische Götter

Hermes
Hestia
Hypnos
Kronos
Poseidon
Rhea
Uranos
Zeus



Römische Götter

Amor
Apollo
Jupiter
Mars
Uranos
Venus
Vesta
Vulcanus

- Fragestellung: welchem Gott kann ich opfern, ohne mich mit einer der beiden Götterwelten anzulegen.
- Heutige Fragestellungen:
 - Mitgliedschaft in verschiedenen Institutionen
 - Rasterfahndung

Einfacher Algorithmus

```
int indexOf(String[] a, String[] b) {  
    for (int i = 0; i < a.length; i++) {  
        for (int j = 0; j < b.length; j++) {  
            if(a[i].equals(b[j])) return i;  
        }  
    }  
    return -1;  
}
```

- doppelt geschachtelte Schleife
- Aufwand $O(n*m)$
 - bei $2 * 10000$ Elementen $\rightarrow 10^8$

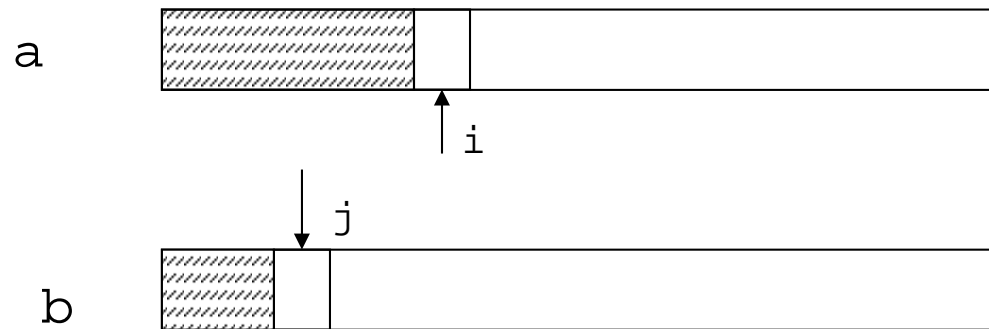
Griechische Götter

Hermes
Hestia
Hypnos
Kronos
Poseidon
Rhea
Uranos
Zeus

Römische Götter

Amor
Apollo
Jupiter
Mars
Uranos
Venus
Vesta
Vulcanus

Besserer Algorithmus wenn a und b sortiert



- **Invariante** $\forall k, n; k < j, n < i; b[k] \neq a[n]$
- Invariante ist eine Bedingung, die erhalten bleibt
- Idee: Bereich, für den die Invariante gilt, sukzessive erweitern
- $b[j] < a[i] \Rightarrow \forall k; k \leq j; b[k] < a[i] \Rightarrow j \text{ um } 1 \text{ erhöhen}$
- $b[j] > a[i] \Rightarrow \forall n; n \leq i; b[j] > a[n] \Rightarrow i \text{ um } 1 \text{ erhöhen}$

Schnelle Suche in zwei Arrays

```
int indexOf(String[] a, String[] b) {  
    int i = 0, j = 0; boolean found = false;  
    // {inv && i == 0 && j == 0}  
    while (!found && (i < a.length-1 || j < b.length-1)) {  
        int c = a[i].compareTo(b[j]);  
        if (c == 0) found = true;  
        else if (c < 0)  
            if (i < a.length-1) i++; else break;  
        else  
            if (j < b.length-1) j++; else break;  
    }  
    // {inv && (i == a.length-1 && j == b.length-1) ||  
    //   (found && a[i] == b[j])}  
    if (found) return i else return -1;  
}
```

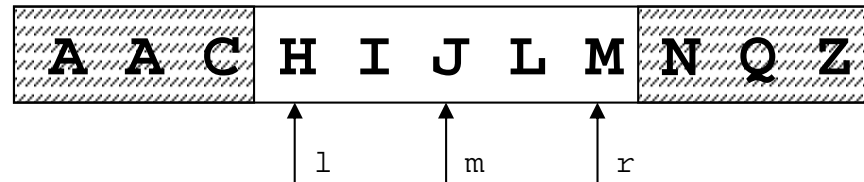
- i und j werden erhöht, so dass die Invariante erhalten bleibt
- Am Schluss gilt: Invariante **and** Abbruchbedingung
- Aufwand $O(n+m)$

Schnelle Suche in einem Array: Binäres Suchen



- Gegeben sei ein sortierter Array von Werten
- Wie kann in so einem Array effizient gesucht werden?
- Führe zwei Indizes ein: l und r
- Invariante: $\forall k, n; k < l, n > r; (a[k] < s) \wedge (a[n] > s)$

Binäres Suchen



nehme m als Index zwischen l und r

- falls $a[m] < s \rightarrow l = m+1$
- falls $a[m] > s \rightarrow r = m-1$
- falls $a[m] = s \rightarrow$ gefunden
- falls $l > r \rightarrow$ keine Elemente mehr zwischen l und $r \rightarrow$ nicht gefunden

Binäres Suchen

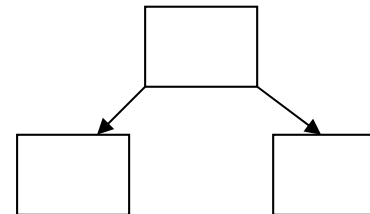
```
int binarySearch(int[] a, int v) {  
    int l = -1;  
    int r = a.length;  
    {inv && l == -1 && r == a.length}  
    while (l <= r) {  
        int m = (l + r) / 2;  
        if (a[m] == v) return m;  
        else if (a[m] < v) l = m+1;  
        else if (a[m] > v) r = m-1;  
    }  
    {inv && l > r}  
    return -1;  
}
```

Demo

- In jedem Durchgang wird $r - l$ halbiert $\rightarrow \log_2$ Schritte
- Aufwand: $O(\log n)$

Aufwand für Suchen und Einfügen

- **sortierter Array**
 - Einfügen $O(n)$
 - binäres Suchen $O(\log(n))$
- **lineare sortierte Liste**
 - Einfügen $O(n)$
 - Suchen $O(n)$
- **sortierter Binärbaum**
 - Einfügen $O(\log(n))$
 - Suchen $O(\log(n))$



Frage:

- Gibt es ein Verfahren, dessen Aufwand unabhängig von der Anzahl Elemente ist?

Maps und Sets

Maps (Abbildungen)

Gegeben sei eine Menge von Datensätzen der Form

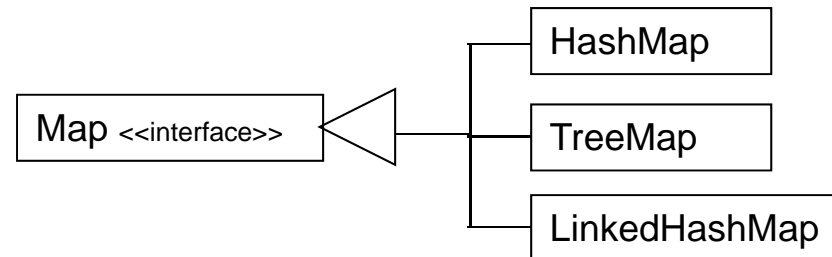


- Der Schlüssel kann ein Teil des Inhaltes sein
- Der Schlüssel besteht im einfachsten Fall aus einem String oder einem numerischen Wert.
- Mittels dem Schlüssel kann der Datensatz wiedergefunden werden.
- Bsp:
 - AHV-Nummer - Personen
 - Matrikel-Nummer- Studenten
- Aufgabe: Es sollen Daten (Objekte) in einen Behälter eingefügt und mittels ihrem Schlüssel wiedergefunden werden können.
- Eine solche Datenstruktur heisst *Map*. Sie bildet eindeutige Schlüssel auf Werte (Inhalte) ab. → Java Interface *Map*.

Maps (Abbildungen)

Eine Map kann auf unterschiedliche Arten implementiert werden.

- Im einfachsten Fall, wenn der Schlüssel ein *int* ist, ist dies ein
- Das Java Collection Framework kennt drei Implementierungen:
 - HashMap
 - TreeMap
 - LinkedHashMap

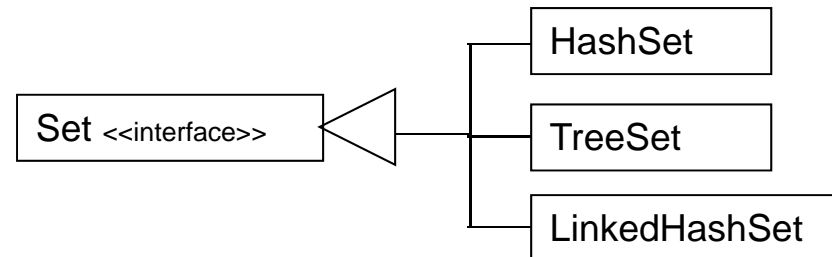


Map<K,V> Interface

<code>void clear()</code>	Löschen aller Elemente
<code>int size()</code>	Anzahl Elemente
<code>V put(K key, V value)</code>	Einfügen eines Elementes
<code>V get (K key)</code>	Finden eines Elementes
<code>V remove(K key)</code>	Löschen eines Elementes
<code>boolean containsKey(K key)</code>	ist Element mit Schlüssel in Table
<code>boolean containsValue(V value)</code>	hat ein Element den Wert
<code>Collection<V> values()</code>	alle Werte als Collection
<code>Set<K> keySet()</code>	alle Schlüssel als Set
<code>Set<Map.Entry<K,V>> entrySet()</code>	alle Schlüssel/Wertepaare als Set

Sets (Mengen)

- Ein *Set* ist eine Collection von Elementen, die keine Duplikate enthält.
- Gleich wie eine Map, aber statt Schlüssel/Wertepaare hat man nur Schlüssel.
- Java Interface *Set*, mit den Implementierungen
 - HashSet
 - TreeSet
 - LinkedHashSet

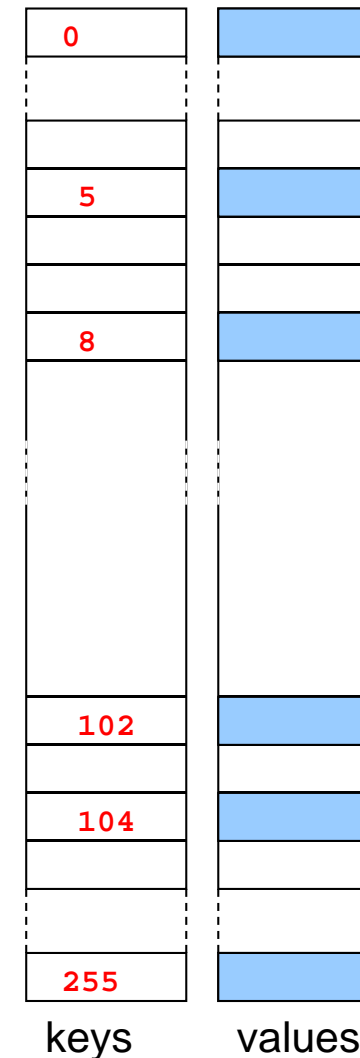


Implementationen einer HashMap: Hashtable

Menge der möglichen Schlüsselwerte klein

Werte in Array an ihrer Indexposition speichern

- Schlüsselbereich (char): 0..255
- Werte beliebig
- Einfach `Object[] h = new Object[256];`
- Einfügen `h[i] = o;`
- Suchen: `o = h[i];`
- Aufwand
 - Einfügen: $O(1)$
 - Suchen: $O(1)$



Idee : Hashing

Problem: der Array ist nur schwach belegt

- geht noch für Buchstaben
- was aber bei Zahlen? ($>2^{32}$) oder Strings

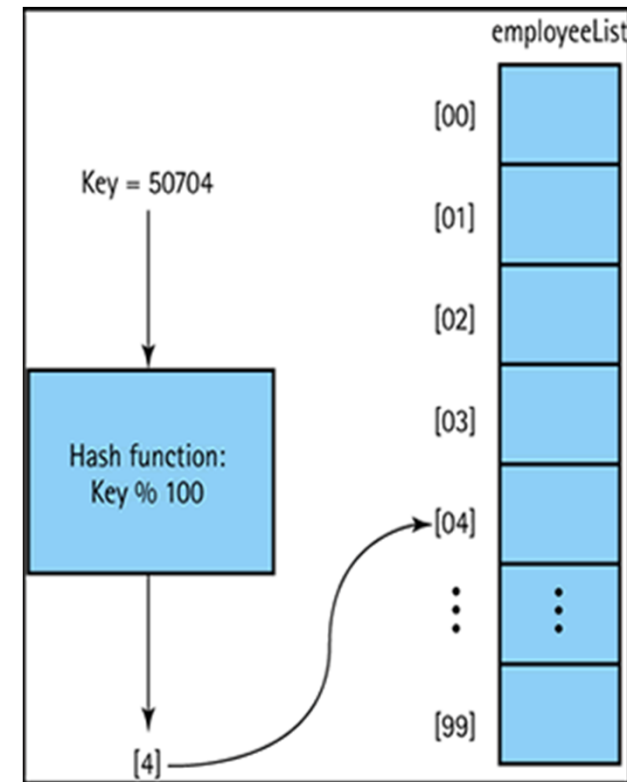
- Lösung:

- Es wird eine Funktion h verwendet, welche den grossen Wertebereich auf einen kleineren abbildet.

Einfachste Funktion: $X \text{ modulo } \text{tableSize}$

→ eine Zahl zwischen 0 und $\text{tableSize}-1$

- Eine solche Funktion nennt man **Hash Funktion**.



Hash Funktion

Problem 1:

- Hash Wert nicht eindeutig, i.e. $\nexists h^{-1}(k)$

Lösung

⇒ **Originalwert** in Tabelle (z.B. Array) speichern

Problem 2:

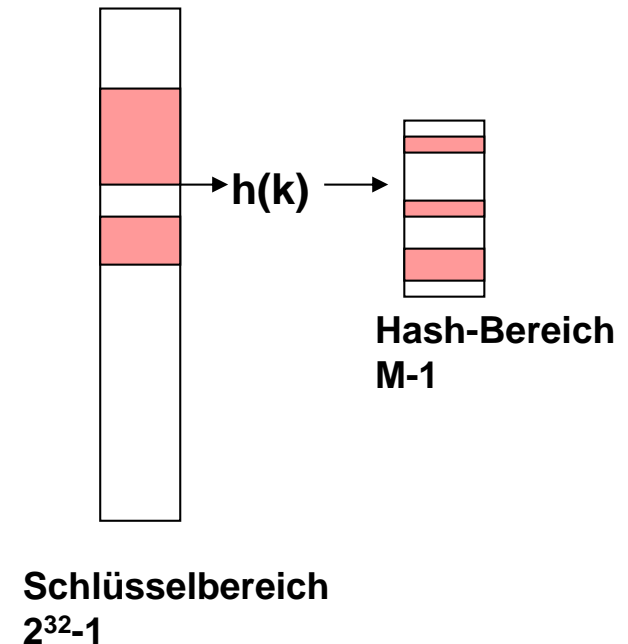
- Zwei unterschiedliche Objekte können den gleichen Hash Wert haben. d.h. sie müssten an der gleichen Stelle gespeichert werden ⇒ **Kollision**

Lösung:

- Kollisionen werden vermieden - oder verringert
- Kollisionen werden aufgelöst
- verschiedene Verfahren zur Auflösung (später):
 - linear/quadratic probing, ...

Hash Funktion

- Der grosse Schlüsselbereich wird mittels der Hash-Funktion auf einen kleinen Bereich abgebildet
- Problem
 - **Massierungen** im (Schlüssel-)Wertebereich
⇒ kann zu gleichen Hash-Werten führen ⇒
Kollisionen
- Hash = Durcheinander
 - die Hash-Funktion bringt den Schlüssel so "durcheinander", dass er möglichst **gleichmässig** auf den ganzen Hash-Bereich abgebildet wird
- Zwei gute Hashfunktionen
 - $h(k) = k \% M \mid M \in \text{Primzahl}$
 - $h(k) = (k * N) \% M \mid N, M \in \text{Primzahl}$



Hashtable ohne Kollisionen

```
public class Hashtable {  
    final int MAX = 100;  
    final int INVAL = 0;          // nicht das Gelbe vom Ei  
    int[] keys = new int[MAX];    // initialisiert mit INVAL  
    int[] vals = new int[MAX];  
  
    private int h(int key) {  
        return key % 97;  
    }  
  
    public void put(int key, int val) {  
        int h = h(key);  
        if (keys[h] == INVAL) {  
            keys[h] = key;  
            vals[h] = val;  
        }  
        else { /* COLLISION */ }  
    }  
  
    public int get(int key) {  
        int h = h(key);  
        if (keys[h] == key) {  
            return vals[h];  
        }  
        else { /* COLLISION or not found */ }  
    }  
}
```

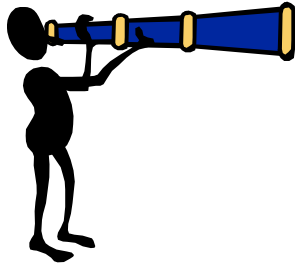
Hash Funktion

Überprüfe ob Feld frei

Speichere Zahl

Überprüfe ob
Schlüssel korrekt

Hole Zahl



Hashing von Strings

- Hashfunktion = ASCII Wert \times Position 128^n

*1 *128 *128² *128³

B	A	U	M
---	---	---	---

- es entstehen sehr grosse Zahlen:
 - Ein vier Zeichen langer String führt bereits zu einer Zahl in der Grössenordnung von $128^4 = 2^{28}$, was nur wenig kleiner als der grösste 32-Bit Integer ist.
- In der Praxis werden deshalb Polynome zur Umwandlung von Strings verwendet (Horner Schema): $A_3x^3 + A_2x^2 + A_1x^1 + A_0x^0$
 - kann als $((A_3)x + A_2)x + A_1)x + A_0$ gerechnet werden.
- Modulo-Arithmetik: Es gilt:
 - $(a+b) \bmod m = ((a \bmod m) + (b \bmod m)) \bmod m$
 - $(a*b) \bmod m = ((a \bmod m) * (b \bmod m)) \bmod m$

Hashing von Strings

Beispiel B A U M (ASCII-Werte B=66, A=65, U=85, M=77)

- $((((77*128+85)*128+65)*128+66) \bmod 19$
- $((((1*14+9)*14+8)*14+9) \bmod 19$
- $((4*14+8)*14+9) \bmod 19$
- $((7*14+9) \bmod 19$
- $107 \bmod 19$
- **12**

Java Methoden hashCode()

String: $s[0]*31^{n-1} + s[1]*31^{n-2} + \dots s[n-1]$

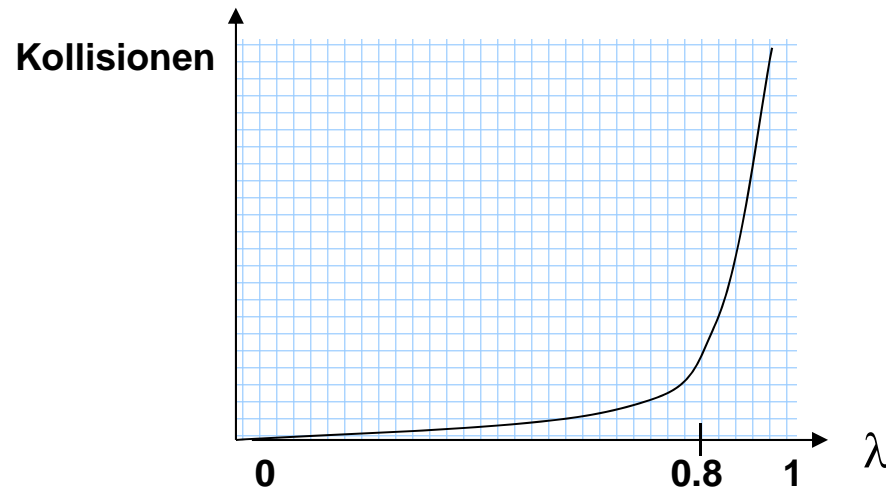
Float: Integer Bit Darstellung

Double: XOR der Integer Bit Darstellung der beiden Teile

Object: Adresse

Kollisionen

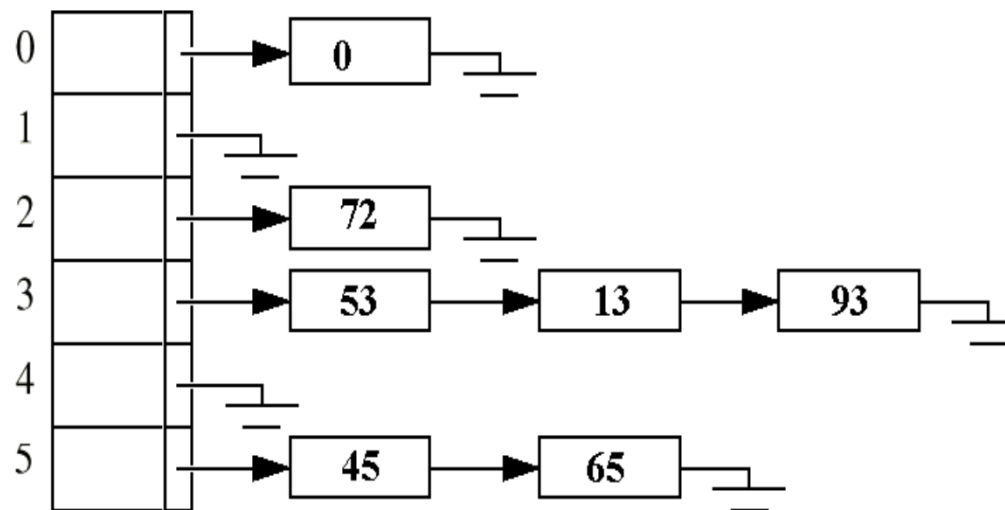
- Anzahl Kollisionen hängt von der Güte der Hash-Funktion und der Belegung der Zellen ab
- Der LoadFactor λ
 - sagt wie stark der Hash-Bereich belegt ist
 - bewegt sich zwischen 0 und 1.
 - Anzahl Kollisionen ist abhängig von λ und h : $f(h, \lambda)$



Kollisionsauflösung 1 : Separate Chaining

Hashtable lediglich als Ankerpunkt für Listen aller Objekte, die den gleichen HashWert haben : **Überlauflisten** (*Separate Chaining*)

Eine solche Hashtable wird auch als *offene* Hashtable bezeichnet.



- Overhead durch Verwendung einer weiteren Datenstruktur
- Verfahren ist gut bei Load-Faktor nahe (oder grösser) als 1

$$h(\text{key}) = \text{key} \% 10$$

Übung

Gegeben sind:

- eine Hashtabelle der Grösse 10
- eine Hash-Funktion $h(x) = x \bmod 10$
- Input: 4371, 1323, 6173, 4199, 4344, 9679, 1989.

Wie sieht die Tabelle aus, nachdem der Input unter Verwendung von Separate Chaining Hashing verarbeitet wurde?

Kollisionsauflösung 2: Open Addressing

Open Addressing: Techniken, wo bei Kollision eine freie Zelle sonstwo in der HashTable gesucht wird.

Wird auch als *geschlossene* Hashtable bezeichnet.

⇒ Setzt einen LoadFactor < ~0.8 voraus.

- **lineares Sondieren (*Linear Probing*):**

sequentiell nach nächster
freier Zelle suchen
(mit *Wrap around*).

- **quadratisches Sondieren (*Quadratic Probing*):**

in wachsenden Schritten der
Reihe nach $F+1$, $F+4$, $F+9$, . .
. , $F+i^2$ prüfen
(mit *Wrap around*).

Linear Probing 1

```
hash( 89, 10 ) = 9
hash( 18, 10 ) = 8
hash( 49, 10 ) = 9
hash( 58, 10 ) = 8
hash(  9, 10 ) = 9
```

in eine Hash-Tabelle mit 10 Feldern werden der Reihe nach 89, 18, 49, 58 und 9 eingefügt.

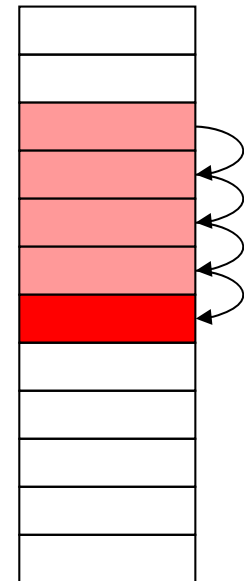
After Insert 89 After Insert 18 After Insert 49 After Insert 58 After Insert 9

0			49	49	49
1				58	58
2					9
3					
4					
5					
6					
7					
8		18	18	18	18
9	89	89	89	89	89

- Hash-Funktion Input modulo Tabellengrösse
- Bei zunehmendem Load Faktor dauert es immer länger bis eine Zelle gefunden wird. (Einfügen und Suchen)
- find funktioniert wie insert: Element wird in Tabelle ausgehend vom HashWert gesucht bis Wert gefunden oder leere Zelle.

Linear Probing 2

```
int findPos(Object x ) {  
    int currentPos = hash(x);  
  
    while( array[ currentPos ] != null &&  
        !array[currentPos].equals( x ) ) {  
        currentPos = (currentPos + 1) % array.length ;  
    }  
    return currentPos;  
}
```



zur Bestimmung einer
Ausweichzelle wird einfach die
nächste genommen → primary
Clustering Phänomen

Linear Probing 3

Performance

ziemlich schwierig abzuschätzen, da der Aufwand nicht nur vom Load Faktor, sondern auch von der Verteilung der belegten Zellen abhängt.

Phänomen des Primary Clustering:

mussten einmal freie Zellen neben dem Hash-Wert belegt werden, steigt die Wahrscheinlichkeit, dass weiter gesucht werden muss für:

- alle Ausgangswerte mit gleichem Hash-Wert
- all jene, deren Hash-Wert in eine der nachfolgenden Zellen verweist.

Folge:

- Verlängerung des durchschnittlichen Zeitaufwandes zum Sondieren
- erhöhte Wahrscheinlichkeit, dass weiteres Sondieren nötig wird

⇒ Bei hohem Load Faktor/ungünstigen Daten bricht die Performance ein!

Übung

Gegeben sind:

- eine Hashtabelle der Grösse 10
- eine Hash-Funktion $H(X) = X \bmod 10$
- Input: 4371, 1323, 6173, 4199, 4344, 9679, 1989.

Wie sieht die Tabelle aus, nachdem der Input unter Verwendung einer linearen Sondiermethode verarbeitet wurde?

Quadratic Probing 1

```
hash( 89, 10 ) = 9
hash( 18, 10 ) = 8
hash( 49, 10 ) = 9
hash( 58, 10 ) = 8
hash(  9, 10 ) = 9
```

in eine Hash-Tabelle mit 10 Feldern werden der Reihe nach 89, 18, 49, 58 und 9 eingefügt.

After Insert 89 After Insert 18 After Insert 49 After Insert 58 After Insert 9

0			49	49	49
1					
2				58	58
3					9
4					
5					
6					
7					
8		18	18	18	18
9	89	89	89	89	89

Hash-Funktion Input
modulo Tabellengrösse

jetzt bleiben Lücken in
Hash-Tabelle offen

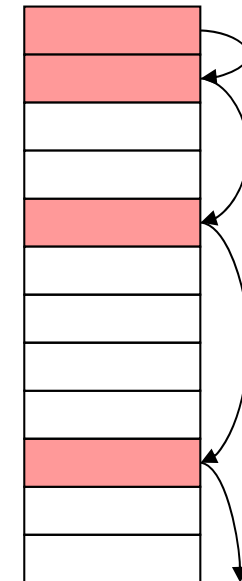
die zuletzt eingefügte 9
findet ihren Platz
unbeeinflusst von der
zuvor eingefügten 58.

Quadratic Probing 2

```
int findPos( Object x )
{
    int collisionNum = 0;
    int currentPos = hash(x);

    while( array[currentPos] != null &&
        !array[currentPos].equals( x ) ) {
        currentPos += 2 * ++collisionNum - 1;
        currentPos = currentPos % array.length;
    }

    return currentPos;
}
```



bessere Performance als
lineares Sondieren weil
primary Clustering
weniger auftritt.

Übung

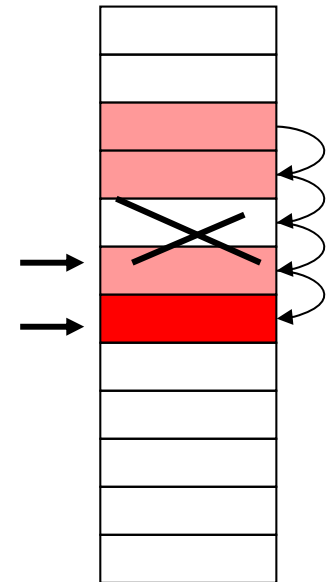
Gegeben sind:

- eine Hashtabelle der Grösse 10
- eine Hash-Funktion $H(X) = X \bmod 10$
- Input: 4371, 1323, 6173, 4199, 4344, 9679, 1989.

Wie sieht die Tabelle aus, nachdem der Input unter Verwendung einer quadratischen Sondiermethode verarbeitet wurde?

Löschen in Hashtabellen

- Werte können nicht einfach gelöscht werden, da sie die Folge der Ausweichzellen unterbrechen.
- Wenn ein Wert gelöscht wird, müssen alle Werte, die potentielle Ausweichzellen sind, gelöscht und wieder eingefügt werden (rehashing).
- Zweite Möglichkeit: gelöschte Zelle lediglich als "gelöscht" markieren → Tombstones



Vor- und Nachteile von Hashing

Vorzüge

- Suchen und Einfügen in Hash-Tabellen sehr effizient
- binäre Bäume können je nach Inputdaten degenerieren, Hash-Tabellen kaum.
- Der Implementationsaufwand für Hash-Tabellen ist geringer als derjenige für ausgeglichene binäre Bäume.

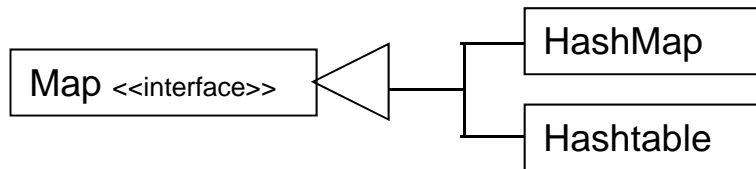
Einschränkungen

- das kleinste oder grösste Element lässt sich nicht einfach finden
- Geordnete Ausgabe nicht möglich
- die Suche nach Werten in einem bestimmten Bereich oder das Finden z.B. eines Strings, wenn nur der Anfang bekannt ist, ist nicht möglich

Hash-Tabellen sind geeignet wenn:

- die **Reihenfolge** nicht von Bedeutung ist
- nicht nach **Bereichen** gesucht werden muss
- die **ungefähre (maximale) Anzahl** bekannt ist

Hashtable in Java



Hashtable	↔	HashMap
seit java 1.0		seit java 1.2
synchronized		nicht synchronized
keine null-Werte		Collections.synchronizedMap() verwenden
zusätzliche Methoden		null-Werte erlaubt

Konstruktoren

```

HashMap<K,V>()
HashMap<K,V>(int initialCapacity)
HashMap<K,V>(int initialCapacity,
    float loadFactor)
  
```

Konstruktor

Konstruktor mit Grösse
Konstruktor mit Grösse und
load Faktor

Methoden

```

void clear()

int size()
V put(K key, V value)
V get (Object key)
V remove(Object key)
boolean containsKey(Object key)
boolean containsValue(Object value)
Collection<V> values()
Set<K> keySet()
  
```

Löschen aller Elemente

Anzahl Elemente
Einfügen eines Elementes
Finden eines Elementes
Löschen eines Elementes
ist Element mit Schlüssel in Table
hat ein Element den Wert
alle Werte als Collection
alle Schlüssel als Set

Zusammenfassung

- Suche
 - Einfache Stringsuche
 - Suche in zwei Sammlungen
 - Binäres Suchen
- Maps und Sets
- Hashing
 - Idee
 - Hashfunktion
 - gute Hashfunktionen
 - Kollisionsauflösung
 - Überlauflisten
 - lineares Sondieren
 - quadratisches Sondieren
 - Vor- und Nachteile
 - Hashtabellen in Java