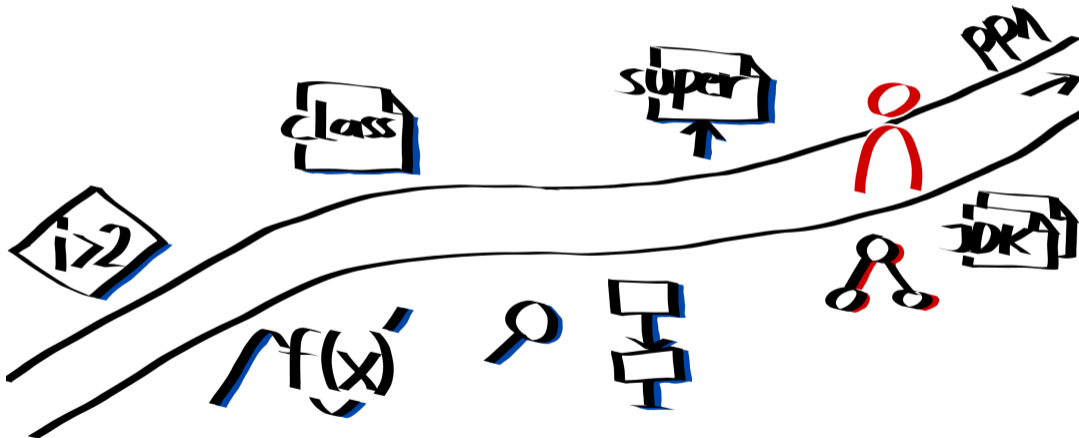


Kapitel 7: Datenstrukturen für effiziente Suche

VL 25: Hashing

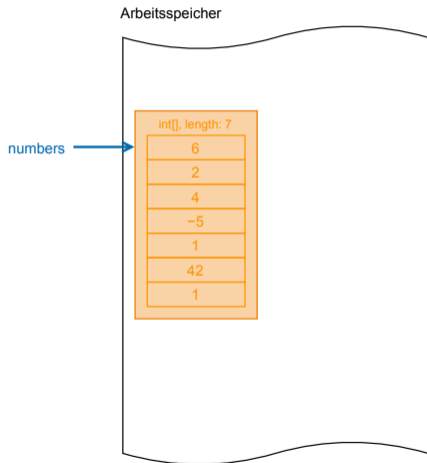


Wo stehen wir gerade?



- Stoff ab heute nicht mehr praktisch geübt
 - ⇒ Sie müssen zu diesem Stoff in Klausur keine Methoden schreiben
 - ! trotzdem kleine Verständnisaufgaben möglich
- Selbsttestfragen im Ilias (Material zu Woche 14 und 15)
 - Wer Fragen beantworten kann und weiß, *warum* welche Antwort richtig/falsch ist, ist gut für Klausur vorbereitet
- ! Stoff trotzdem für Praxis (und Propra) sehr(!) relevant

- wenn sehr oft bestimmte Daten gesucht:
auch binärer Suchbaum langsam
 - Beispiel: Bestimmte Kund:innen in großem Kundenstamm finden
- Ziel: Datensuche „in einem Schritt“
- Grundidee: Objekte in Array speichern, Objekteigenschaften (z. B. Kundennummer) geben Index vor
 - direkter Zugriff auf richtige Array-Position möglich



Definition

Eine Hashfunktion h ist eine Abbildung von Objekten auf Werte (typischerweise ganze Zahlen).

Anwendungsgebiete:

- Prüfsummen (MD5 (veraltet), SHA-512; Prüfung auf Übertragungsfehler/Manipulationen)
- sichere Speicherung von Passwörtern (kryptografische Hashfunktionen)
- Datenstrukturen, die schnelle Suche ermöglichen (u. a. in Datenbanken)

Run this command in your terminal in the directory the iso was downloaded to verify the SHA256 checksum:

```
echo  
"93bdab204067321ff131f560879db46bee3b994bf24836bb78538640f689  
e58f *ubuntu-20.04.2.0-desktop-amd64.iso" | shasum -a 256  
--check
```

You should get the following output:

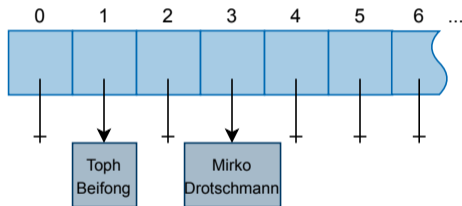
```
ubuntu-20.04.2.0-desktop-amd64.iso: OK
```

Quelle: ubuntu.com

- Datenstruktur, die wie Menge funktioniert
 - behält keine (sinnvolle) Sortierung
 - keine doppelten Elemente
- Finden von Daten „in einem Schritt“
- Grundidee: Objekte in großem Array speichern, Hashfunktion gibt Index vor
- Frage: Wie sieht die Hashfunktion aus?

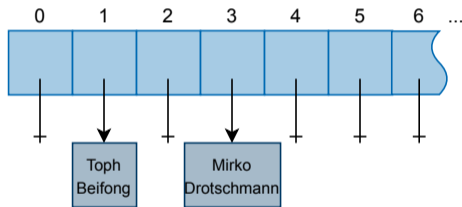
Personen sollen häufig und schnell über Nachnamen gefunden werden

o	$h(o)$
Beifong	1
Drotschmann	3
Nguyen-Kim	13
Schönfeld	18



Personen sollen häufig und schnell über Nachnamen gefunden werden

o	$h(o)$
Beifong	1
Drotschmann	3
Nguyen-Kim	13
Schönfeld	18



h : Position des ersten Buchstabens im Alphabet (minus 1)

o	$h(o)$
Beifong	1
Drotschmann	3
Nguyen-Kim	13
Schönfeld	18
Brenneis	1

- Problem: können nicht zwei Elemente an selbem Index speichern
 - Kollision: $h(o_1) = h(o_2) \wedge o_1 \neq o_2$
- Lösungsansätze?

- Hashfunktion sollte so viele Informationen wie möglich berücksichtigen
- Beispiel:
 - WAGNER → 230107140518
 - BACH → 02010308
- Problem: ergibt sehr große Indizes/Arrays
- Lösung?

- Hashfunktion sollte so viele Informationen wie möglich berücksichtigen
- Beispiel:
 - WAGNER → 230107140518
 - BACH → 02010308
- Problem: ergibt sehr große Indizes/Arrays
- Lösung:
 - Modulo der Array-Größe rechnen

- Hashfunktion sollte so viele Informationen wie möglich berücksichtigen
- Beispiel:
 - WAGNER → 230107140518
 - BACH → 02010308
- Problem: ergibt sehr große Indizes/Arrays
- Lösung:
 - Modulo der Array-Größe rechnen
 - ! kann aber weiterhin zu Kollisionen führen:
 - 230107140518 mod 10 = 8
 - 02010308 mod 10 = 8

- Möglichst wenig Kollisionen
 - möglichst viele Objekteigenschaften berücksichtigen
 - gute Streuung der Hashwerte
 - am besten surjektive Abbildung auf ausreichend große Wertemenge
 - Hashsets mit primen Array-Größen führen in Praxis zu weniger Kollisionen
- schnell berechenbar
- ! Aber: Auch mit sehr guten Hashfunktionen Kollisionen nie ganz vermeidbar (insbes. wenn Hashset schon sehr voll).
 - Lösungsstrategien für Kollisionen erforderlich

Lösung 1: Hashset dynamisch vergrößern¹

- Vergrößerung des Arrays, wenn zu stark gefüllt/zu viele Kollisionen
- Dabei meist direkt Verdopplung der Speicherkapazität
 - positiv für amortisierte Laufzeit → $O(\log n)$
- Nachteil?

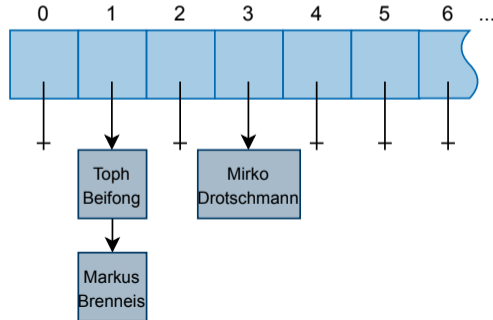
¹Dynamisches Hashing

- Vergrößerung des Arrays, wenn zu stark gefüllt/zu viele Kollisionen
- Dabei meist direkt Verdopplung der Speicherkapazität
 - positiv für amortisierte Laufzeit → AI_{Dat}
- Nachteil: Einträge müssen neu gehasht werden (Modulo neuer Speichergröße)
 - spezielle Hashfunktionen, die dieses Problem abschwächen → AI_{Dat}

¹Dynamisches Hashing

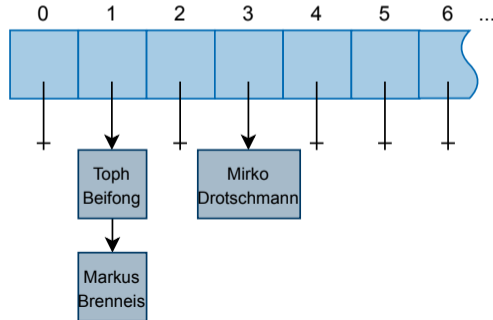
Lösung 2: Mehrerer Einträge am selben Index

- Einträge speichern Referenz auf Liste
- Nachteil?

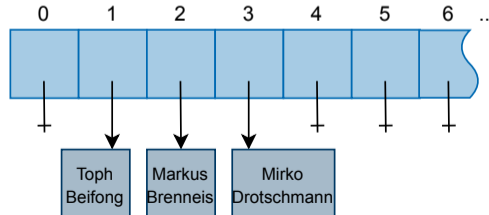


Lösung 2: Mehrerer Einträge am selben Index

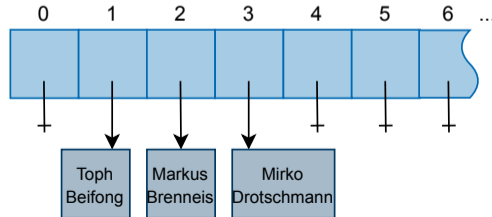
- Einträge speichern Referenz auf Liste
- Nachteil:
 - langsamere Suche, wenn Listen lang werden



- Alternativen Platz gemäß Sondierungsreihenfolge suchen
- Mögliche Strategien:
 - lineares Sondieren: $h(o) + 1, h(o) + 2, h(o) + 3, \dots$ durchprobieren
 - quadratisches Sondieren: $h(o) \pm 1, h(o) \pm 2, h(o) \pm 4, \dots$
 - komplexere Strategien \rightarrow AIDat



- Alternativen Platz gemäß Sondierungsreihenfolge suchen
 - Mögliche Strategien:
 - lineares Sondieren: $h(o) + 1, h(o) + 2, h(o) + 3, \dots$ durchprobieren
 - quadratisches Sondieren: $h(o) \pm 1, h(o) \pm 2, h(o) \pm 4, \dots$
 - komplexere Strategien \rightarrow AIDat
- ! Sondierung beim Suchen/Löschen/etc. zu berücksichtigen
- ! Erfordert eine Lösung für den Fall, dass alle Indices besetzt



- Alle Klassen erben Standardimplementierung von `hashCode`
 - sollte – ähnlich wie `equals` – überschrieben werden

! `equals` und `hashCode` müssen konsistent sein:

`a.equals(b)` \Rightarrow `a.hashCode() == b.hashCode(b)`

```
1  @Override
2  public int hashCode() {
3      return lastname.charAt(0) - 'A';
4  }
5
6  @Override
7  public boolean equals(Object other) {
8      if(this == other)
9          return true;
10     if(other == null || getClass() != other.getClass())
11         return false;
12     // Objektinhalte vergleichen
13     return ((Person) other).lastname.equals(this.lastname);
14 }
```

²<https://docs.oracle.com/javase/8/docs/api/java/util/Set.html>

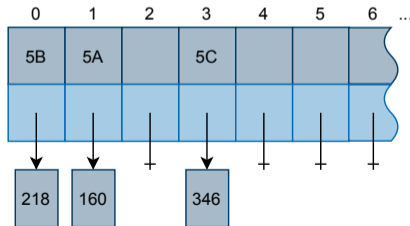
- Wichtig: Hashwert eines Objekts darf sich nach speichern im HashSet nicht mehr ändern.
- ⇒ Alle Objekteigenschaften, die für Hashwert relevant sind, sollten `final` sein
- Oracle²: „Great care must be exercised if mutable objects are used as set elements.“

```
1 HashSet<Student> hashSet = new HashSet<>();
2 Student student = new Student("Alex", 1234);
3
4 hashSet.add(student);
5 System.out.println(hashSet.contains(student)); // true
6
7 student.changeStudentNumber(1243);
8 System.out.println(hashSet.contains(student)); // false
```

²<https://docs.oracle.com/javase/8/docs/api/java/util/Set.html>

- Speichern Objekt (Value) zusammen mit selbst gewähltem Schlüssel (Key)
- Schlüssel-Objekt wird fürs Hashing verwendet
 - ⇒ Value-Objekte dürfen veränderlich sein
- Häufige Verwendung in nicht-OOP-Sprachen (z. B. Clojure)

```
1 HashMap<String,Integer> rooms = new HashMap<>();  
2 rooms.put("5A", 160);  
3 rooms.put("5B", 218);  
4 rooms.put("5C", 346);  
5 System.out.println(rooms.get("5A")); // 160
```



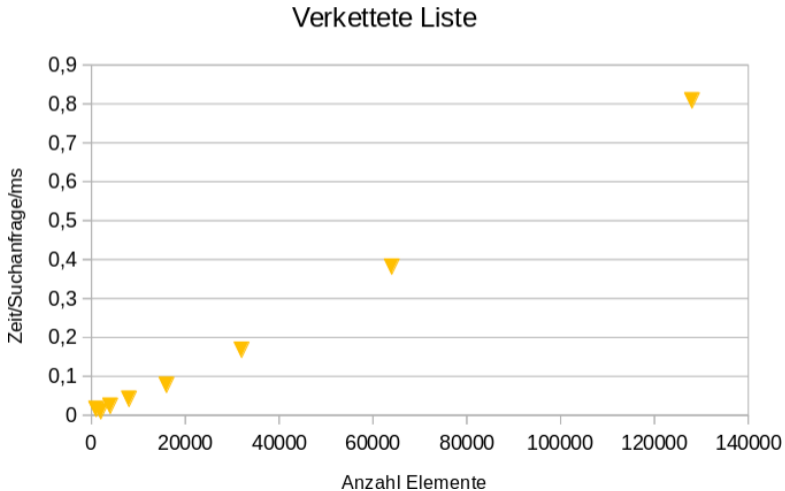
³Map, Dictionary, Assoziatives Array, Lookup-Tabelle

n : Anzahl der gespeicherten Elemente

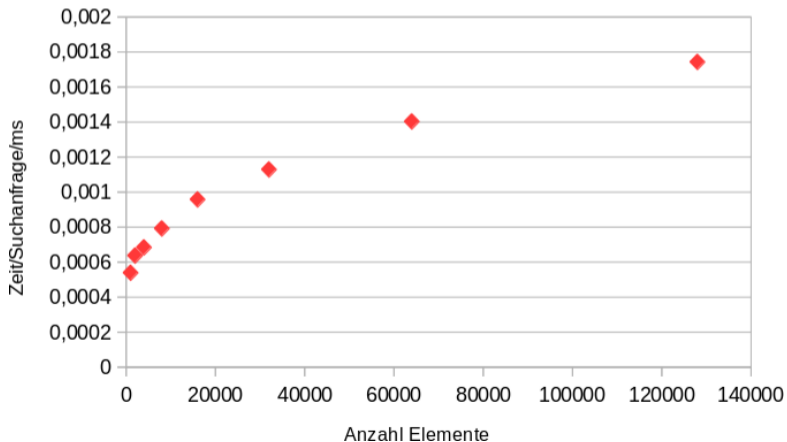
	Array	Anzahl der Schritte im „typischen“ Fall ⁴			Hashset
		Verk. Liste	Bin. Suchbaum ausgeglichen	entartet	
Hinzufügen (irgendwo)	–	1	$\log(n)$	n	1
Hinzufügen an Index i	–	n	–	–	–
Entfernen von Element	–	n	$\log(n)$	n	1
Suche von Element	n	n	$\log(n)$	n	1
Element an Index i	1	n	–	–	–
dynamischer Speicherbedarf	nein	ja		ja	nein ⁵
Duplikate möglich	ja	ja		nein	nein
Best. Reihenfolge speicherbar	ja	ja		nein	nein
Objekte müssen sortierbar sein	nein	nein		ja	nein
Objekte müssen hashbar sein	nein	nein		nein	ja

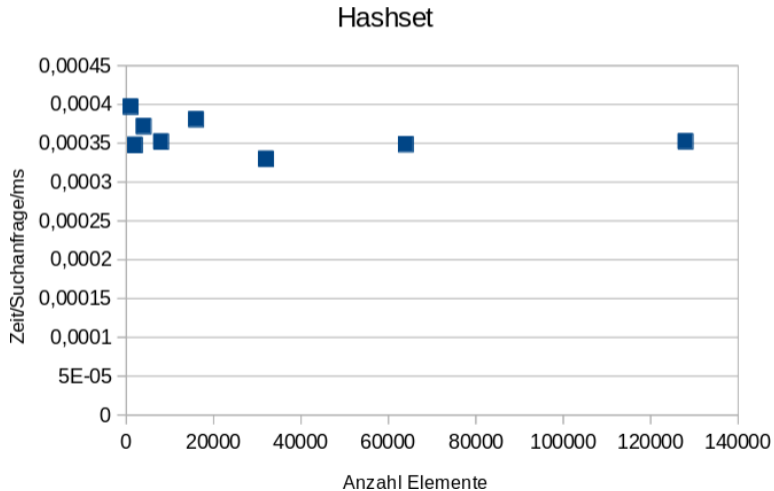
⁴genauere Analyse der *amortisierten Laufzeit* in AIDat

⁵in der von uns umgesetzten Variante

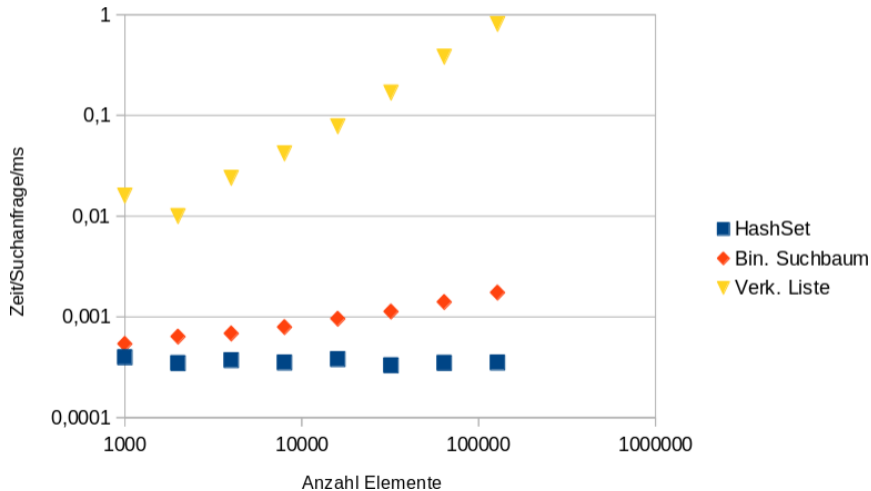


Binärer Suchbaum





Suchgeschwindigkeit im Vergleich



Sie können am Ende der Woche ...

- **erklären**, welche Vor- und Nachteile von Hashsets gegenüber anderen Datenstrukturen haben.

Hashing
Werte

Hashfunktion
Hashset

Schlüssel
Kollision

- im Laufe des Tages an der Stelle, wo Übungsblätter hochgeladen werden
- Empfehlung: unter realistischen Bedingungen (ausgedruckt, 120 Minuten, ohne Störung, nur Spicker) rechnen
- PK etwas anspruchsvoller konzipiert als die beiden echten Klausuren, letztere ungefähr gleich schwierig
- Inhalt: alles bis einschließlich Suchbäume (echte Klausuren decken auch letzte VLs noch ab!)
- Besprechung in den letzten Tutorien, evtl. zusätzlich Lösungsvideo