

# Hashverfahren

# Hashing

- Hashverfahren
  - Aufgabe
  - Realisierung
- Hashfunktionen
  - Anforderungen
  - Wahl einer Hashfunktion
- Kollisionen
  - Geburtstagsparadoxon
  - Strategien zur Kollisionsbehandlung
  - Hashverfahren mit Verkettung der Überläufer
  - Offene Hashverfahren
- Ausblick
  - Universelles Hashing

# Hashing

## Aufgabe

Dynamische Verwaltung von Daten, wobei jeder Datensatz eindeutig durch einen Schlüssel charakterisiert ist

Viele Anwendungen benötigen nur einfache Daten-Zugriffsmechanismen (**dictionary operations**):

- Suche nach Datensatz bei gegebenem Schlüssel  $x$   
 $\text{search}(x)$
- Einfügen eines neuen Datensatzes  $d$  mit Schlüssel  $x$   
 $\text{insert}(x, d)$  (abgekürzt  $\text{insert}(x)$ )
- Entfernen eines Datensatzes bei gegebenem Schlüssel  $x$   
 $\text{delete}(x)$

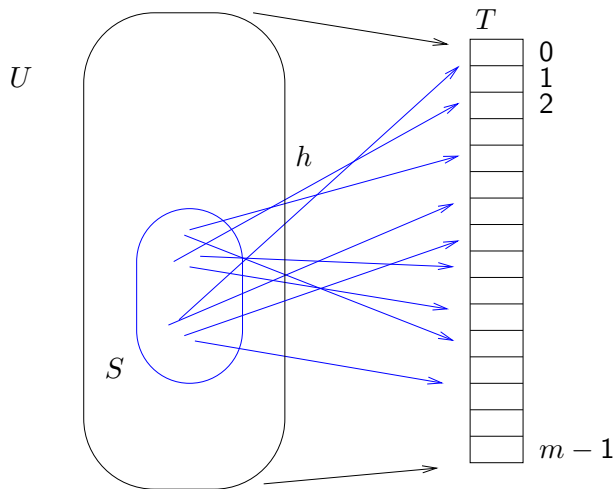
Menge potentieller Schlüssel (**Universum**) kann **sehr** groß sein!

# Hashing

## Hashing

- Menge  $U$  potentieller Schlüssel sehr groß, aktuelle Schlüsselmenge  $S$  jeweils nur kleine Teilmenge des Universums (im allgemeinen  $S$  nicht bekannt)
- **Idee:** durch Berechnung feststellen, wo Datensatz mit Schlüssel  $x$  gespeichert
- Abspeicherung der Datensätze in einem Array  $T$  mit Indizes  $\{0, 1, \dots, m - 1\}$ : **Hashtabelle**
- **Hashfunktion**  $h$  liefert für jeden Schlüssel  $x \in U$  eine Adresse in Hashtabelle, d.h.  $h : U \rightarrow \{0, 1, \dots, m - 1\}$ .

# Hashing



# Vorteil von Hashing

Sei  $n := |S|$ .

- Balancierte Suchbäume (AVL-Bäume, B-Bäume):  
dictionary operations haben Komplexität  $O(\log n)$
- Hashing:  
für alle Operationen **mittlere** Komplexität  $O(1)$

## Belegungsfaktor

Quotient  $\beta := n/m$  heißt Belegungsfaktor oder Auslastungsfaktor einer Hashtabelle der Größe  $m$ .

Mittlerer Aufwand für **dictionary operations** als Funktion in  $\beta$  abschätzbar

→ Anzahl aktueller Schlüssel geht nur indirekt in Aufwand ein

# Hashing

## Herausforderung:

Anzahl möglicher Schlüssel viel größer als Hashtabelle, also

$$|U| \gg m$$

- Hashfunktion muss verschiedene Schlüssel  $x_1$  und  $x_2$  auf gleiche Adresse abbilden.
- $x_1$  und  $x_2$  beide in aktueller Schlüsselmenge  
 → Adresskollision

## Hashverfahren gegeben durch:

- eine Hashtabelle,
- eine Hashfunktion, die Universum der möglichen Schlüssel auf Adressen einer Hashtabelle abbildet,
- eine Strategie zur Auflösung möglicher Adresskollisionen.

# Anforderungen an Hashfunktionen

## Gute Hashfunktionen sollten:

- surjektiv sein, d.h. den ganzen Wertebereich umfassen,
- die zu speichernden Schlüssel (möglichst) gleichmäßig verteilen, d.h. für alle Speicherplätze  $i$  und  $j$  sollte gelten  $|h^{-1}(i)| \approx |h^{-1}(j)|$ ,
- effizient berechenbar sein.

Voraussetzung:

$$U = \{0, 1, \dots, N - 1\} \text{ und}$$

$$h : U \rightarrow \{0, 1, \dots, m - 1\}$$

Häufig: Divisions- oder Kongruenzmethode

$$h(x) := x \bmod m \rightarrow |h^{-1}(i)| \in \{\lfloor N/m \rfloor, \lceil N/m \rceil\}$$



# Wahl einer Hashfunktion

$$h(x) := x \bmod m$$

Problem: Daten oft nicht gleichverteilt!

Beispiel: Texte in Zahlen übertragen, oft viele Leerzeichen,  
 bestimmte Wörter häufiger etc.

Wichtig: geeignete Wahl von  $m$

- $m$  Zweierpotenz:  $x \bmod m$  wählt nur **die letzten  $\log m$  Bits**
- $m$  Primzahl:  $x \bmod m$  **beeinflusst alle Bits**

Beispiel:  $m = 11$  und  $S = \{49, 22, 6, 52, 76, 34, 13, 29\}$

Hashwerte:  $h(49) = 49 \bmod 11 = 5$ ,  $h(22) = 22 \bmod 11 = 0$ ,  
 6, 8, 10, 1, 2, 7

## Zur Erinnerung:

Im Allgemeinen unvermeidbar, dass Kollisionen auftreten, denn aus  $N \gg m$  folgt Existenz eines Speicherplatzes  $i$  mit  $|h^{-1}(i)| \geq N/m$ .

## Frage:

Sei  $n := |S|$ . Wie wahrscheinlich sind Kollisionen bei  $n \ll m$ ?

## Geburtstagsparadoxon

Bei wie vielen (zufällig gewählten) Personen ist es *wahrscheinlich*, dass hiervon zwei am selben Datum (Tag und Monat) Geburtstag haben?

## Annahme:

- Daten unabhängig
- $\text{Prob}(h(x) = j) = 1/m$

$\text{Prob}(i\text{-tes Datum kollidiert nicht mit den ersten } i - 1 \text{ Daten, wenn diese kollisionsfrei sind}) = \frac{m - (i - 1)}{m}$

## Intuition:

Egal welche Speicherplätze die ersten  $i - 1$  Daten belegen,  $m - i + 1$  der  $m$  Möglichkeiten sind *gut*.

$$\text{Prob}(n \text{ Daten kollisionsfrei}) = \frac{m-1}{m} \cdot \frac{m-2}{m} \dots \frac{m-n+1}{m}$$

**Beispiel:**  $m = 365$

$$\text{Prob}(23 \text{ Daten kollisionsfrei}) \approx 0.49$$

$$\text{Prob}(50 \text{ Daten kollisionsfrei}) \approx 0.03$$

Prob( $2m^{1/2}$  Daten kollisionsfrei) =

$$\frac{m-1}{m} \dots \frac{m-m^{1/2}}{m} \dots \frac{m-2m^{1/2}+1}{m}$$

$$\leq 1 \cdot \left( \frac{m-m^{1/2}}{m} \right)^{m^{1/2}} = \left( 1 - \frac{1}{m^{1/2}} \right)^{m^{1/2}} \approx \frac{1}{e}$$

Hashing muss mit Kollisionen leben und benötigt Strategien zur Kollisionsbehandlung!

# Kollisionsbehandlung

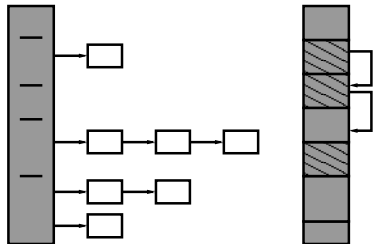
Hashverfahren unterscheiden sich in Strategien zur Kollisionsbehandlung:

- mittels verketteter Listen:  
Jede Komponente der Hashtabelle enthält Zeiger auf Überlaufliste.
- mittels offener Adressierung:  
Im Kollisionsfall nach fester Regel alternativen freien Platz in Hashtabelle suchen. Für jeden Schlüssel Reihenfolge, in der Speicherplätze betrachtet werden, vorgegeben:  
**Sondierungsfolge**

# Kollisionsbehandlung

Verschiedene Arten der Kollisionsbehandlung:

- mittels verketteter Listen  
(links)
- mittels offener Adressierung  
(rechts)



# Hashing mit Verkettung

## Realisierung:

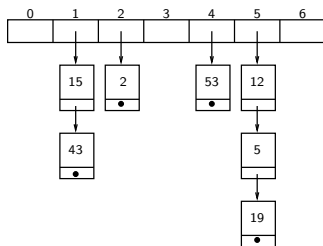
Jede Komponente der Hashtabelle enthält Zeiger auf paarweise disjunkte lineare Listen. Die  $i$ -te Liste  $L(i)$  enthält alle Schlüssel  $x \in S$  mit  $h(x) = i$ .

Vorteil: Alle Operationen werden unterstützt und  $n > m$  ist möglich. (Für  $n \gg m$  jedoch **Rehashing** ratsam) Nachteil: Speicherplatzbedarf für Zeiger

- **search**( $x$ ): Berechne  $h(x)$  und suche in Liste  $L(h(x))$ .
- **insert**( $x$ ) (nach erfolgloser Suche): Berechne  $h(x)$  und füge  $x$  in Liste  $L(h(x))$  ein.
- **delete**( $x$ ) (nach erfolgreicher Suche): Berechne  $h(x)$ , suche  $x$  in Liste  $L(h(x))$  und entferne  $x$ .

# Beispiel Verkettung der Überläufer

Beispiel:  $m = 7$  und  $h(x) = x \bmod m$   
 $S = \{2, 5, 12, 15, 19, 43, 53\}$





# Analyse

Bei zufälligen Daten und ideal streuenden Hashfunktion gilt für

$$X_{ij} := \begin{cases} 1 & i\text{-tes Datum kommt in Liste } L(j) \\ 0 & \text{sonst} \end{cases}$$

$$\text{Prob}(X_{ij} = 1) = \frac{1}{m}$$

$$\rightarrow E(X_{ij}) = 1 \cdot \frac{1}{m} + 0 \cdot \frac{m-1}{m} = \frac{1}{m}$$

$X_j = X_{1j} + \dots + X_{nj}$  zählt Anzahl Daten in Liste  $L(j)$ .

$$E(X_j) = E(X_{1j} + \dots + X_{nj}) = E(X_{1j}) + \dots + E(X_{nj}) = \frac{n}{m}$$

## Quadratisches Sondieren

Frage: Ist  $t(\cdot, j)$  für alle  $j$  und  $m$  bijektiv?

**Nein**, aber immer wenn  $m \equiv 3 \pmod{4}$  und  $m$  eine Primzahl ist  
(Beweis: Zahlentheorie)

Besser als **lineares Sondieren**, aber für großes  $m$  sind die ersten  
Werte noch *nah* an  $j$