

# Hashing



$U$   
"Universum"

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

Problem:  $|U|$  sehr groß, insbesondere  
 $m \ll |U|$

d.h.  $h$  nicht injektiv, also  $h(k_1) = h(k_2)$   
für  $k_1 \neq k_2$ .

  
Kollision

Umgang mit Kollisionen:

1) Hashing mit Überlauf („chaining“)

2) Hashing mit Ersatzadresse  
(„open addressing“)

Idee:

- speichere alle Datensätze in der Hashtabelle (d.h. keine Überlauf-Listen)
- $\Rightarrow m \geq n$
- Bei Kollision muss Ersatzadresse gesucht werden.

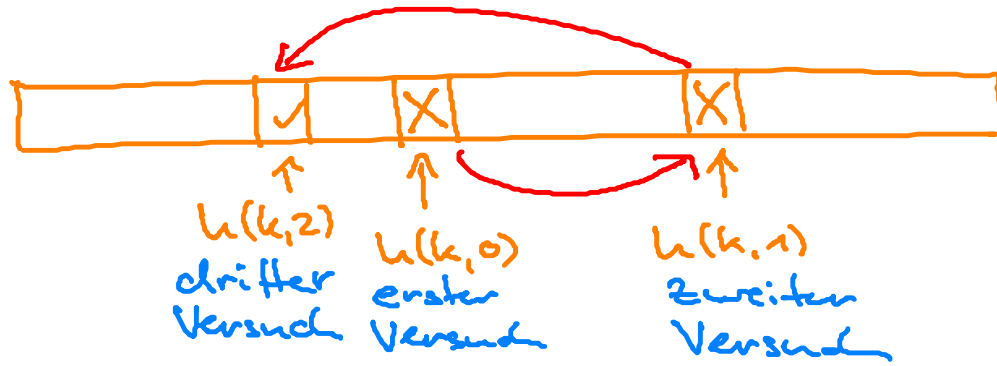
$$h: U \times \{0, 1, \dots, m-1\} \rightarrow \{0, 1, \dots, m-1\}$$

↑  
Schlüssel

↑  
# erfolgloser  
Versuche bei  
Suche nach  
freiem Eintrag

↑  
Adresse

$h(k, i) \hat{=}$   $(i+1)$ -ter Versuch, eine freie Adresse für  $k$  zu finden.



Dieses Verfahren berechnet also nacheinander die Adressen

$h(k,0), h(k,1), \dots, h(k,i), \dots$

bis eine freie Adresse gefunden wird.

Damit auf diese Art immer eine freie Adresse gefunden wird (unter der Annahme, dass  $n \leq m$ ), muss gelten:

$h(k,0), h(k,1), \dots, h(k, m-1)$   
ist Permutation von  $0, 1, \dots, m-1$

„Permutationsbedingung“

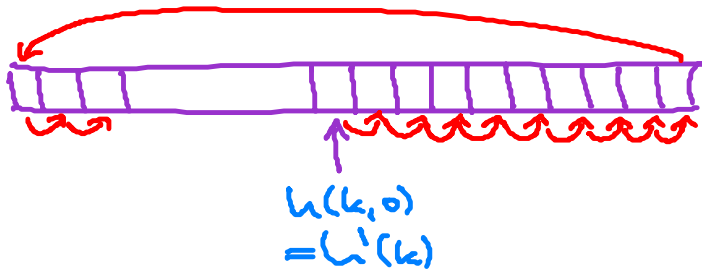
Achtung: Dies erschwert das Löschen!

Beispiele für open addressing:

- Lineares Sondieren

$$h(k, i) := (h'(k) + i) \bmod m$$

$\uparrow$   
 gewöhnliche  
 Hashfkt.  $\hat{=} h(k, 0)$



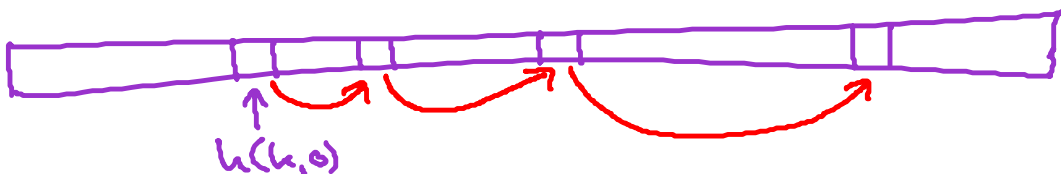
Permutationsbedingung offensichtlich erfüllt!

Nachteil: „primäres Clustering“, d.h. es bilden sich immer größere Blöcke konsequenter belegter Adressen, denn:

Falls die Adressen  $a, a+1, a+2, \dots, a+i-1$  bereits belegt sind und als nächstes ein neues Datenelement mit zufälliger Schlüssel eingetüft wird, so landet dieses mit W'keit  $\frac{i+1}{m}$  an der Stelle  $a+i$ .

Quadratisches Sondieren:

$$h(k, i) = (h'(k) + c_1 \cdot i + c_2 \cdot i^2) \bmod m$$



Sprungweite nimmt quadratisch zu!

Beachte: Permutationsbedingung hängt von  $m, c_1, c_2$  ab und muss im Einzelfall überprüft werden.

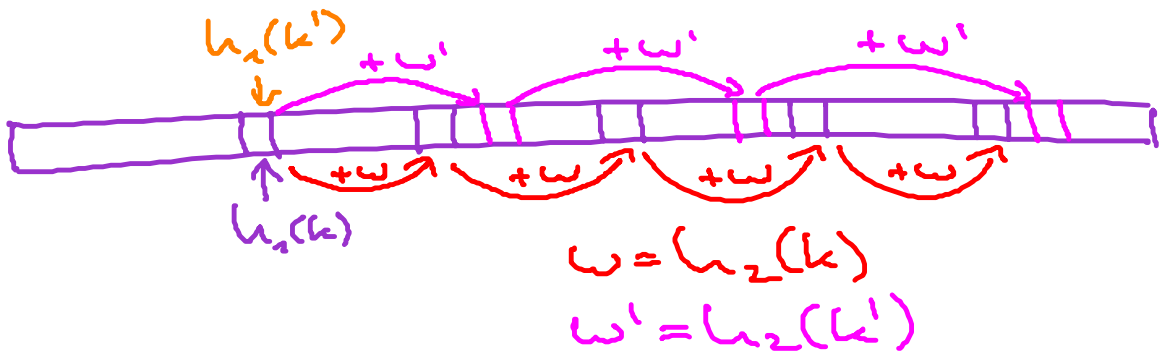
Vorteil: kein primäres Clustering!

Nachteil: sekundäres Clustering:  
Schlüssel, die auf dieselbe erste Adresse abgebildet werden haben auch dieselbe Folge von Ersatzadressen.

Doppel-Hash:

$$h(k, i) = (h_1(k) + i \cdot h_2(k)) \bmod m$$

$h_1, h_2$  herkömmliche Hashfkt.



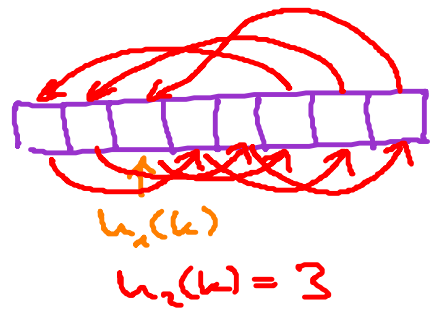
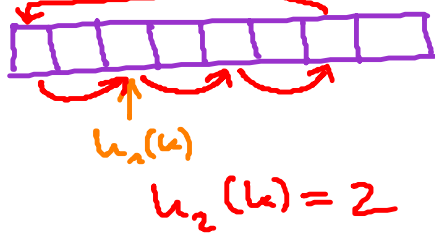
vermeidet primäres und sekundäres Clustering, d.h.

$h_1(k) = h_1(k')$  führt zu unterschiedl.

Ersatzadressen, falls  $h_2(k) \neq h_2(k')$

Frage: Wie erfüllt man die Permutationsbedingung?

Bsp:  $m = 8$



Permut. Bed. verletzt

Permut. Bed. erfüllt

Satz: Doppelhash erfüllt die Permutationsbedingung  $\Leftrightarrow$

$\forall k$  sind  $m$  und  $h_2(k)$  relativ prim  
d.h.  $\text{ggT}(m, h_2(k)) = 1$ .

Bemerkung:

- Diese Bedingung ist immer erfüllt, wenn  $m$  Primzahl ist.
- Sie ist z.B. auch erfüllt, wenn  $m = 2^q$  und  $h_2(k)$  ungerade für alle  $k$ .

Beweis: " $\Leftarrow$ ": Ann:  $\exists k$ , so dass

$$\text{ggT}(m, h_2(k)) = d > 1$$

$$\Rightarrow m = p \cdot d, \quad h_2(k) = q \cdot d$$

$$\Rightarrow p \cdot h_2(k) = q \cdot p \cdot d = q \cdot m \equiv 0 \pmod{m}$$

d.h. die  $p$ -te Ersatzadresse für  $k$  stimmt überein mit  $h_2(k)$ . Beachte, dass  $p < m$   
 $\Rightarrow$  nicht alle Adressen besucht

$\Rightarrow$  Perm. bed. verletzt  $\Downarrow$

$\Leftarrow$  Ann: Perm. bed. nicht erfüllt, d.h.:

$\exists k$  und  $0 \leq i < j \leq m-1$  mit

$$h_2(k) + i \cdot h_2(k) \equiv h_2(k) + j \cdot h_2(k) \pmod{m}$$

$$\Leftrightarrow (j-i) \cdot h_2(k) \equiv 0 \pmod{m}$$

$$\Leftrightarrow (j-i) \cdot h_2(k) = q \cdot m, \quad q \in \mathbb{N}$$

$$\Leftrightarrow h_2(k) = \frac{q \cdot m}{j-i}, \quad q \in \mathbb{N}$$

Da  $j-i < m$  gilt also  $\text{ggT}(h_2(k), m) > 1$   $\Downarrow$   
 $\square$

## Analyse des „open addressing“

Einfügen: Wie viele Sondierungen müssen durchgeführt werden bis eine freie Adresse gefunden wird?

Beantworte diese Frage unter der  
Gleichverteilungsannahme

d.h.: nächste Sondierung wählt immer gleichverteilt unter den  $m$  Adressen, also wird jede Adresse mit W'keit  $\frac{1}{m}$  gewählt.

Satz: Bei Auslastung  $\alpha = \frac{\lambda}{\mu} < 1$  ist die erwartete Anzahl Sondierungen beim Einfügen höchstens  $\frac{1}{1-\alpha}$ .

Bsp:  $\alpha = \frac{1}{2} \Rightarrow \leq 2$  Sondierungen im Erw. wert

$\alpha = 0,8 \Rightarrow \leq 5$  " " " "

Beweis beruht auf etwas W'-Theorie...