

# Informatik II - SS 2014

## (Algorithmen & Datenstrukturen)

Vorlesung 7 (21.5.2014)

Binäre Suche, Hashtabellen I



**UNI  
FREIBURG**

Fabian Kuhn

Algorithmen und Komplexität

# Abstrakte Datentypen : Dictionary

**Dictionary:** (auch: Maps, assoziative Arrays, Symbol Table)

- Verwaltet eine ~~Kollektion~~<sup>Menge</sup> von Elementen, wo bei jedes Element durch einen eindeutigen Schlüssel (key) repräsentiert wird

(key, value)

## Operationen:

- create : erzeugt einen leeren Dictionary
- D.insert(key, value) : fügt neues (key,value)-Paar hinzu
  - falls schon ein Eintrag für key besteht, wird er ersetzt
- D.find(key) : gibt Eintrag zu Schlüssel key zurück
  - falls ein Eintrag vorhanden (gibt sonst einen Default-Wert zurück)
- D.delete(key) : löscht Eintrag zu Schlüssel key

- Wir kümmern uns in einer ersten Phase nur um die Basisoperationen *insert*, *find*, *delete* (und *create*)

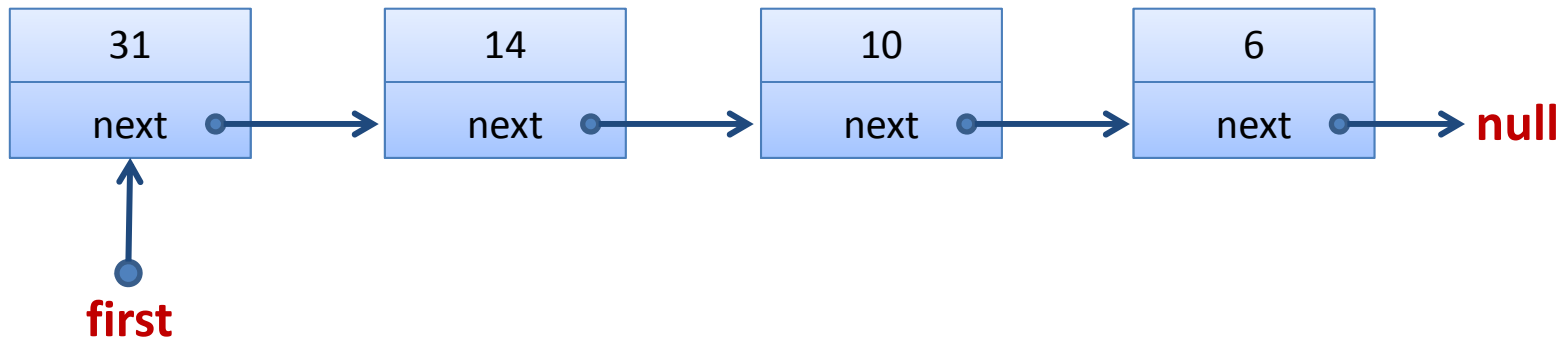
## Dictionary Beispiele:

- Wörterbuch (key: Wort, value: Definition / Übersetzung)
- Telefonbuch (key: Name, value: Telefonnummer)
- DNS Server (key: URL, value: IP-Adresse)
- Python Interpreter (key: Variablenname, value: Wert der Variable)  
Java/C++ Compiler (key: Variablenname, value: Typinformation)

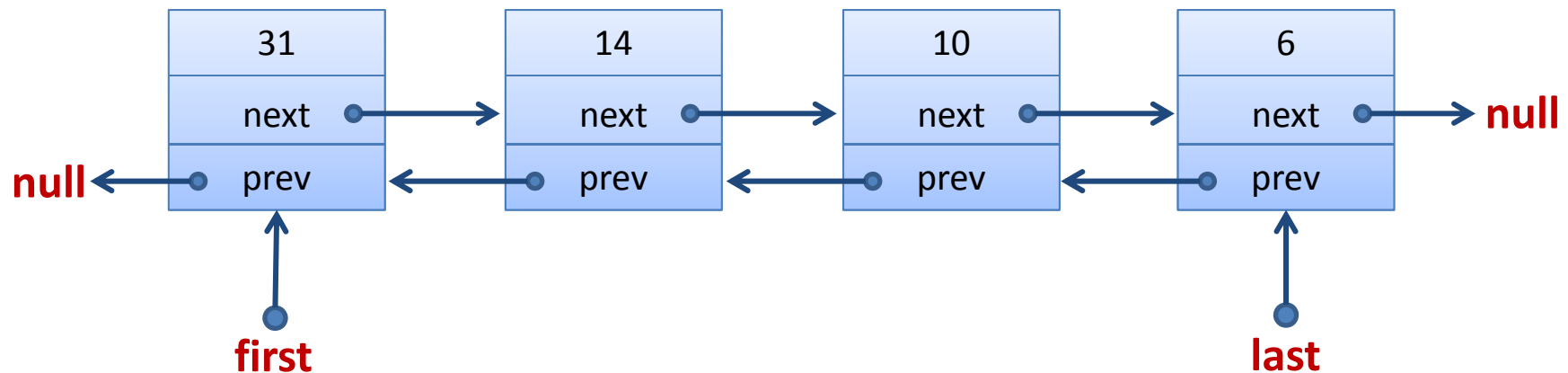
**In all diesen Fällen ist insbesondere eine schnelle find-Op. wichtig!**

# Verkettete Listen: Struktur

## Einfach verkettete Liste (Singly Linked List):





## Doppelt verkettete Liste (Doubly Linked List):



# Dictionary mit verketteten Listen

## Operationen:

- *create*:
  - lege neue leere Liste an
- *D.insert(key, value)*:
  - füge neues Element vorne ein 
  - Annahme: Es gibt noch keinen Eintrag mit dem Schlüssel *key* *(notwendig?)*
- *D.find(key)*:
  - gehe von vorne durch die Liste 
- *D.delete(key)*:
  - suche zuerst das Listenelement (wie in *find*)
  - lösche Element dann aus der Liste
  - Bei einfach verketteten Listen muss man stoppen, sobald *current.next.key == key* ist!



# Dictionary mit verketteten Listen

---

## Laufzeiten:

*create:*  $O(1)$

*insert:*  $O(1)$

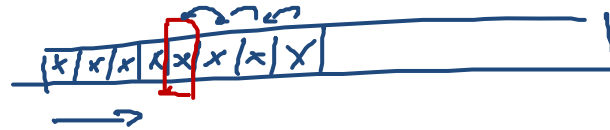
*find:*  $O(n)$

*delete:*  $O(n)$

Ist das gut?

# Dictionary mit Array

## Operationen:



- *create*:
  - lege neues Array der Länge NMAX an
- *D.insert(key, value)*:
  - füge neues Element hinten an (falls es noch Platz hat)
  - Annahme: Es gibt noch keinen Eintrag mit dem Schlüssel *key* ←
- *D.find(key)*:
  - gehe von vorne (oder hinten) durch die Elemente
- *D.delete(key)*:
  - suche zuerst nach dem *key*
  - lösche Element dann aus dem Array:

**Man muss alles dahinter um eins nach vorne schieben!**

# Dictionary mit Array

---

## Laufzeiten:

*create:*  $O(1)$

*insert:*  $O(1)$

*find:*  $O(n)$

*delete:*  $O(n)$

Bessere Ideen?



# Benutze sortiertes Array?

---

- **Teure Operation** bei Liste/Array, insbesondere find
- Falls (sobald) sich die Einträge nicht zu sehr ändern, ist *find* die wichtigste Operation!
- Kann man in einem (nach Schlüsseln) sortierten Array schneller nach einem bestimmten Schlüssel suchen?
  - Beispiel: Suche Tel.-Nr. einer Person im Telefonbuch...

## Ideen:

Gehe in die Mitte  
und dann entweder links oder rechts

# Binäre Suche

Benutze Divide and Conquer Idee!

Suche nach der Zahl (dem Key) 19:

2	3	4	6	9	12	15	16	17	18	19	20	24	27	29
---	---	---	---	---	----	----	----	----	----	----	----	----	----	----

16 <sup>?</sup> > 19

20 > 19

# Binäre Suche

2	3	4	6	9	12	15	16	17	18	19	20	24	27	29
---	---	---	---	---	----	----	----	----	----	----	----	----	----	----

Algorithmus (Array  $A$  der Länge  $n$ , Suche nach Schlüssel  $x$ ):

$l = 0; r = n - 1$

while  $r > l$  do

$m = \lfloor (l + r) / 2 \rfloor$

if  $A[m] < x$  then (x ist rechts von m)

$l = m + 1$

else if  $A[m] > x$  then (x ist links von m)

$r = m - 1$

else

$l = m; r = m$

$A[l] \stackrel{?}{=} x$

# Binäre Suche

2	3	4	6	9	12	15	16	17	18	19	20	24	27	29
---	---	---	---	---	----	----	----	----	----	----	----	----	----	----

Algorithmus (Array  $A$  der Länge  $n$ , Suche nach Schlüssel  $x$ ):

```
l = 0; r = n - 1;
while r > l do
    m = (l + r) / 2;
    if A[m] < x then
        l = m + 1
    else if A[m] > x then
        r = m - 1
    else
        l = m; r = m
```

Falls Schlüssel  $x$  im Array ist, dann gilt am Schluss  $A[l] = x$

# Ist der Algorithmus korrekt?

---

## Wie überprüft man das?

- Empirisch: Unit Test oder auch systematischere Tests...
- **Formal?**
  - Korrektheit ist (meistens) noch wichtiger als Performance!
- **Vorbedingung**
  - Bedingung, welche am Anfang (der Methode / Schleife / ...) gilt
- **Nachbedingung**
  - Bedingung, welche am Schluss (der Methode / Schleife / ...) gilt
- **Schleifeninvariante**
  - Bedingung welche am Anfang / Ende jedes Schleifendurchlaufs gilt

# Ist der Algorithmus korrekt?

```
l = 0; r = n - 1;
while r > l do
    m = (l + r) / 2;
    if A[m] < x then l = m + 1
    else if A[m] > x then r = m - 1
    else l = m; r = m
```

*Handwritten annotations: 'l', 'x', 'r' above the while loop; an arrow pointing to the 'do' keyword.*

## Vorbedingung

- *Array ist am Anfang sortiert, Array hat Länge  $n$*

## Nachbedingung

- *Falls  $x$  im Array ist, dann gilt  $A[l] = x$*

## Schleifeninvariante

- *Falls  $x$  im Array ist, dann gilt  $A[l] \leq x \leq A[r]$*

# Ist der Algorithmus korrekt?

## Vorbedingung

while  $(r > l)$  do ...

- Array ist am Anfang sortiert, Array hat Länge  $n$

$l = 0; r = n - 1;$

## Schleifeninvariante

$A[l] \leq x \leq A[r]$

- Falls  $x$  im Array ist, dann gilt  $A[l] \leq x \leq A[r]$
- Vorbedingung und Zuweisung zu  $l$  und  $r \rightarrow$  Schleifeninvariante
  - Invariante gilt am Anfang des ersten Schleifendurchlaufs

## Nachbedingung

$A[r] \leq A[l] \leq x \leq A[r]$

- Falls  $x$  im Array ist, dann gilt  $A[l] = x$
- Abbruchbedingung while-Schleife  $\rightarrow$   $l \geq r$  und damit  $A[l] \geq A[r]$
- Falls  $x$  im Array ist, dann folgt aus der Schleifeninvariante und da  $A$  sortiert ist, dass  $A[l] = A[r]$  und damit  $A[l] = x$

$A$  ist sortiert  
↙

# Ist der Algorithmus korrekt?

```

l = 0; r = n - 1;
while r > l do
  m = (l + r) / 2;
  if A[m] < x then l = m + 1
  else if A[m] > x then r = m - 1
  else l = m; r = m

```

$A[l] \leq x \leq A[r]$   
INV  
INV

## Schleifeninvariante

- Falls  $x$  im Array ist, dann gilt  $A[l] \leq x \leq A[r]$

$l = m + 1$ :  $l \quad x \quad m \quad r$   
 $A[m] < x$

$r = m - 1$ : analog



# Terminiert der Algorithmus?

$l = 0; r = n - 1;$

while  $r > l$  do

$m = (l + r) / 2;$

if  $A[m] < x$  then  $l = m + 1$  ←

else if  $A[m] > x$  then  $r = m - 1$  ←

else  $l = m; r = m$  ←

$$\frac{l+r}{2} - \frac{1}{2} \leq m \leq \frac{l+r}{2}$$

- Veränderung der Anz. Elemente ( $r - l + 1$ ) pro Schleifendurchlauf?

⊖  $l = m + 1$ : *wenes l*

$$\underline{r - (m + 1) + 1} \leq r - \left( \frac{l+r}{2} + \frac{1}{2} \right) + 1 = \underline{\frac{r-l+1}{2}}$$

–  $r = m - 1$ :

$$\underline{(m - 1) - l + 1} \leq \frac{l+r}{2} - 1 - l + 1 = \frac{r-l}{2} < \underline{\frac{r-l+1}{2}}$$

– Sonst wird  $x$  gefunden und  $r - l + 1$  wird 1

## Terminiert der Algorithmus?

- In jedem Schleifendurchlauf wird die Anzahl der Elemente mindestens halbiert.
- Der Algorithmus terminiert!

## Laufzeit?

$$T(n) \leq T(\lfloor n/2 \rfloor) + c, \quad \underline{\underline{T(1) \leq c}}$$

$$\begin{aligned} T(n) &\leq T(n/2) + c \\ &\leq T(n/4) + 2c \\ &\leq T(n/8) + 3c \\ &\vdots \\ &\leq T(n/2^k) + k \cdot c \\ &\leq T(1) + c \cdot \log_2 n \leq c(1 + \log n) \end{aligned}$$

# Laufzeit Binäre Suche

Der Algorithmus terminiert in Zeit  $O(\log n)$ .

$$T(n) \leq T(\lfloor n/2 \rfloor) + c$$

$$T(n) \leq c(1 + \log n)$$

Verankerung:  $T(1) \leq c(1 + \log 1) = c \quad \checkmark$

Schritt:  $T(n) \leq T(\lfloor n/2 \rfloor) + c$   
 $\leq c(1 + \log \lfloor n/2 \rfloor) + c$   
 $\leq c(1 + \log_2(n/2)) + c$   
 $= \underline{c(1 + \log_2 n)} \quad \checkmark$

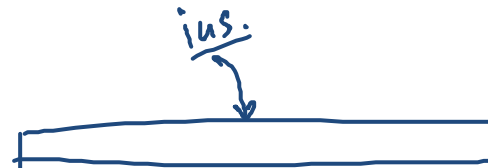
$$\log_2(n/2) = \log_2 n - 1$$

find in Zeit  $O(\log n)$

# Dictionary mit sortiertem Array

## Operationen:

- *create*:
  - lege neues Array der Länge  $NMAX$  an
- *D.find(key)*:
  - Suche nach *key* mit binärer Suche
- *D.insert(key, value)*:
  - suche nach *key* und füge neues Element an der richtigen Stelle ein
  - Einfügen: alles dahinter muss um eins nach hinten geschoben werden!
- *D.delete(key)*:
  - suche zuerst nach dem *key* und lösche den Eintrag
  - Löschen: alles dahinter muss um eins nach vorne geschoben werden!



# Dictionary mit sortiertem Array

---

## Laufzeiten:

*create:*  $O(n)$

*insert:*  $O(n)$

*find:*  $O(\log n)$

*delete:*  $O(n)$

Können wir alle Operationen schnell machen?

- und das *find* noch schneller?

# Direkte Adressierung

Mit einem Array können alles schnell machen,  
...falls das Array gross genug ist.

**Annahme:** Schlüssel sind ganze Zahlen zwischen 0 und  $M - 1$

<u>0</u>	1	<u>null</u>
	2	null
<u>→</u>	3	<u>Value 1</u>
	4	null
	5	null
	6	null
<u>→</u>	7	<u>Sebastian</u>
	8	<u>Value 3</u>
	9	null
	:	:
<u>M-1</u>	<u>M</u>	null

*find(3) → "Value 1"*

*insert(7, "Sebastian")*

*delete(5)*

*alles O(1)*

## 1. Direkte Adressierung benötigt zu viel Platz!

- Falls Schlüssel ein beliebiger *int* (32 bit) sein kann:  
Wir benötigen ein Array der Grösse  $2^{32} \approx 4 \cdot 10^9$ .  
Bei 64 bit Integers sind's sogar schon mehr als  $10^{19}$ ...

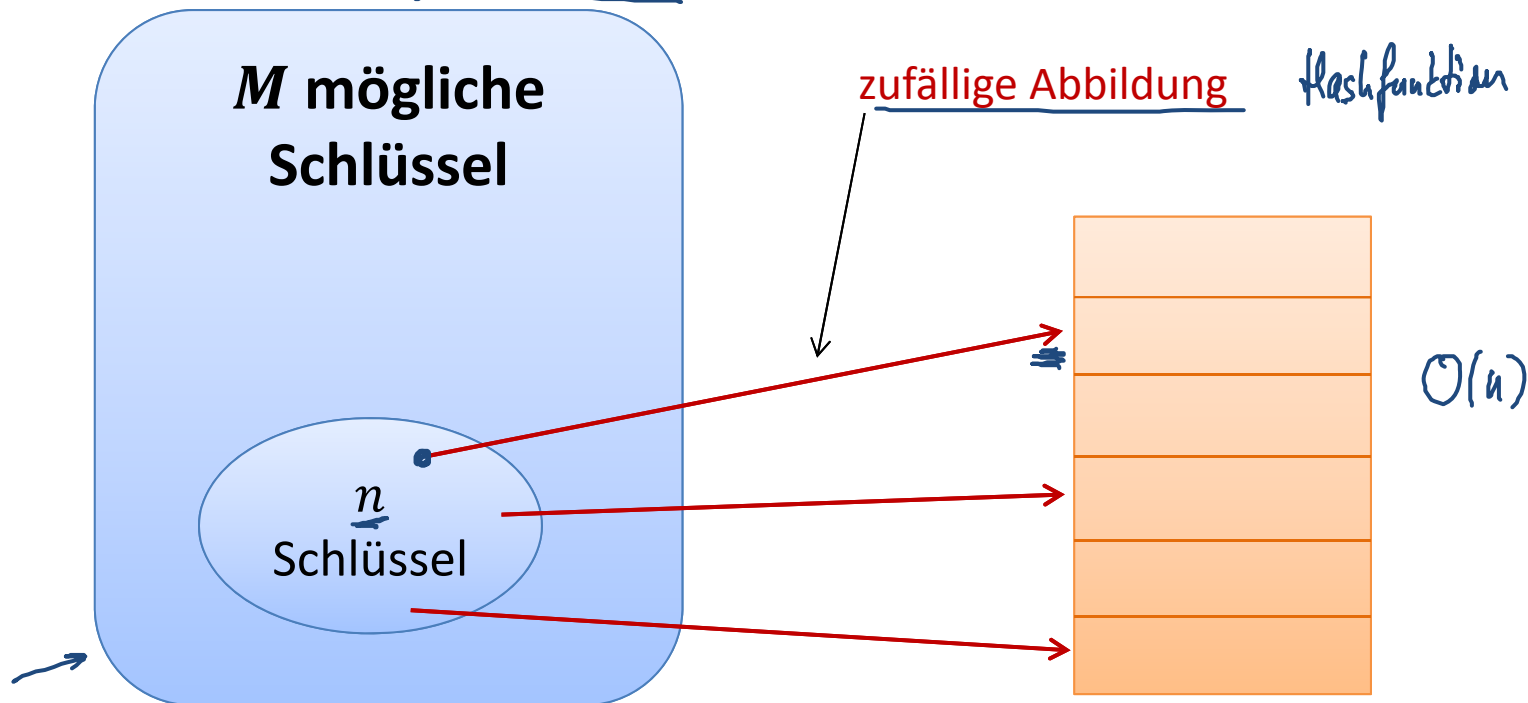
## 2. Was tun, wenn die Schlüssel keine ganzen Zahlen sind?

- Wo kommt das (*key,value*)-Paar ("Sebastian", "Assistent") hin?
- Wo soll der Schlüssel 3.14159 gespeichert werden?
- Pythagoras: "Alles ist Zahl"  
"Alles" kann als Folge von Bits abgespeichert werden:  
**Interpretiere Bit-Folge als ganze Zahl**
- **Verschärft das Platz-Problem noch zusätzlich!**

# Hashing : Idee

## Problem

- Riesiger Raum  $S$  an möglichen Schlüsseln
- Anzahl  $n$  der wirklich benutzten Schlüssel ist **viel** kleiner
  - Wir möchten nur Arrays der Grösse  $\approx n$  (resp.  $O(n)$ ) verwenden...
- Wie können wir  $M$  Schlüssel auf  $O(n)$  Array-Positionen abbilden?





Schlüsselraum  $S$ ,  $|S| = M$  (alle möglichen Schlüssel)

Arraygrösse  $m$  ( $\approx$  Anz. Schlüssel, welche wir max. speichern wollen)

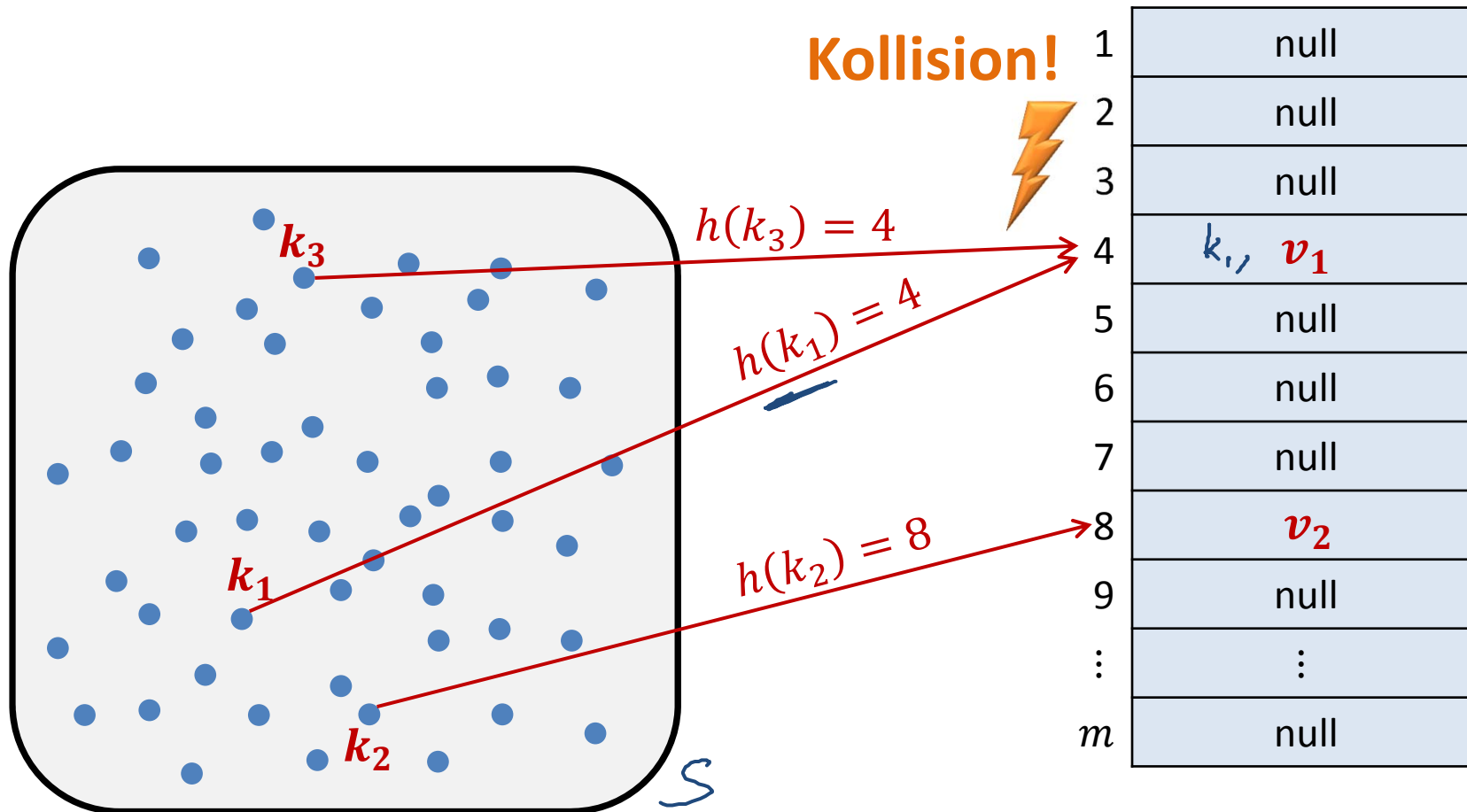
## Hashfunktion

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

- Bildet Schlüssel vom Schlüsselraum  $S$  in Arraypositionen ab
- $h$  sollte möglichst nahe bei einer zufälligen Funktion sein
  - alle Elemente in  $\{0, \dots, m - 1\}$  etwa gleich vielen Schlüsseln zugewiesen sein
  - ähnliche Schlüssel sollten auf verschiedene Positionen abgebildet
- $h$  sollte möglichst schnell berechnet werden können
  - Wenn möglich in Zeit  $O(1)$
  - Wir betrachten es im folgenden als Grundoperation (Kosten = 1)

# Funktionsweise Hashtabellen

1.  $insert(k_1, v_1)$
2.  $insert(k_2, v_2)$
3.  $insert(k_3, v_3)$



# Hashtabellen : Kollisionen

## Kollision:

Zwei Schlüssel  $k_1, k_2$  kollidieren, falls  $h(k_1) = h(k_2)$ .

## Was tun bei einer Kollision?

- Können wir Hashfunktionen wählen, bei welchen es keine Kollisionen gibt?

*geht nicht  
wenn wir die Schlüssel kennen geht's*

- Eine andere Hashfunktion nehmen?

*müssten alles neu einfügen*

Größe  $m$

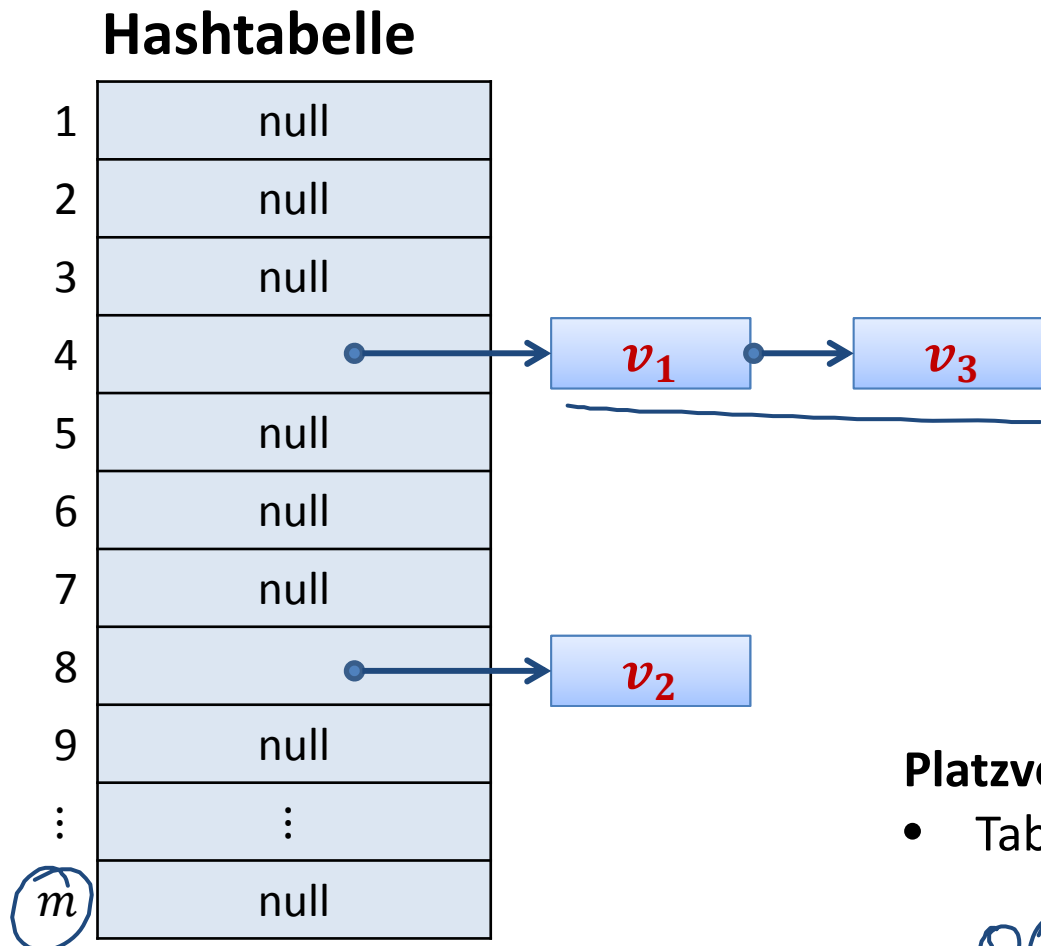
- Weitere Ideen?

## Kollisionen Lösungsansätze

- Annahme: Schlüssel  $k_1$  und  $k_2$  kollidieren
  1. Speichere beide (key,value)-Paare an die gleiche Stelle
    - Die Hashtabelle muss an jeder Position Platz für mehrere Elemente bieten
    - Wir wollen die Hashtabelle aber nicht einfach vergrößern (dann könnten wir gleich mit einer grösseren Tabelle starten...)
    - **Lösung: Verwende verkettete Listen**
  2. Speichere zweiten Schlüssel an eine andere Stelle
    - Kann man zum Beispiel mit einer zweiten Hashfunktion erreichen
    - Problem: An der alternativen Stelle könnte wieder eine Kollision auftreten
    - Es gibt mehrere Lösungen
    - Eine Lösung: Verwende **viele mögliche neue Stellen** (Man sollte sicherstellen, dass man die meistens nicht braucht...)

# Hashtabellen mit Chaining

- Jede Stelle in der Hashtabelle zeigt auf eine verkettete Liste



insert: key x  
gehe zur Stelle  $h(x)$

find:  
gehe zur Stelle  $h(x)$   
suche in Liste

delete:  
analog

## Platzverbrauch:

- Tabellengröße  $m$ , Anz. Elemente  $n$

$$O(m + n)$$

# Laufzeit Hashtabellen-Operationen

---

Zuerst, um's einfach zu machen, für den Fall ohne Kollisionen...

*create:*  $O(1)$

*insert:*  $O(1)$

*find:*  $O(1)$

*delete:*  $O(1)$

- Solange keine Kollisionen auftreten, sind Hashtabellen extrem schnell (falls die Hashfunktion schnell ausgewertet werden kann)
- Wir werden sehen, dass dies auch mit Kollisionen gilt...

# Laufzeit mit Chaining

---

Verkettete Listen an allen Positionen der Hashtabelle

*create:*  $O(1)$

*insert:*  $O(1)$

*find:*  $O(1 + \text{"Länge der Liste an Stelle } h(x)\text{"})$

Worst case:  $O(n)$

*delete:*  $O(1)$

# Funktionsweise Hashtabellen

Schlechtester Fall bei Hashing mit Chaining

- Alle Schlüssel, welche vorkommen, haben den gleichen Hashwert
- Ergibt eine verkettete Liste der Länge  $n$
- Wahrscheinlichkeit bei zufälligem  $h$ :

$$\left(\frac{1}{m}\right)^{n-1}$$

