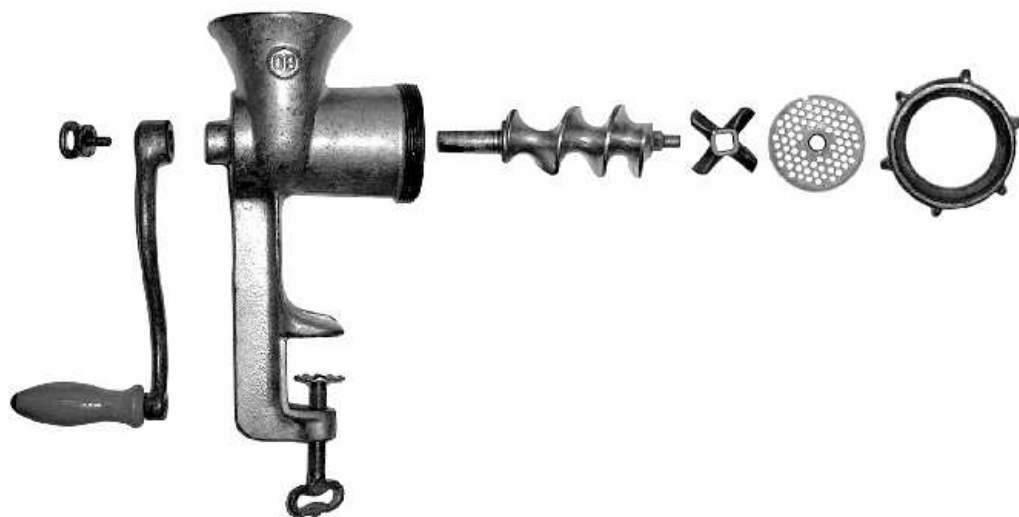
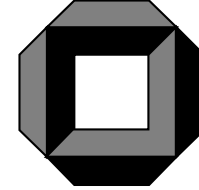


4 Hashing (Streuspeicherung)



“to **hash**” \approx “völlig **durcheinander** bringen”.

Paradoxerweise **hilft** das, Dinge wiederzufinden



Hashtabellen

speichere Menge $M \subseteq$ Element.

$\text{key}(e)$ ist eindeutig für $e \in M$.

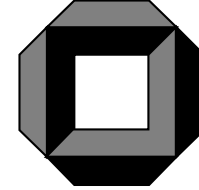
unterstütze **Wörterbuch**-Operationen in Zeit $O(1)$.

$M.\text{insert}(e : \text{Element})$: $M := M \cup \{e\}$

$M.\text{remove}(k : \text{Key})$: $M := M \setminus \{e\}, e = k$

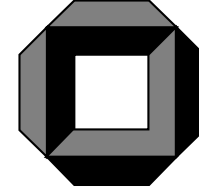
$M.\text{find}(k : \text{Key})$: return $e \in M$ with $e = k$; \perp falls nichts gefunden

(Konvention: key ist *implizit*), d.h. $e = k$ gdw $\text{key}(e) = k$)



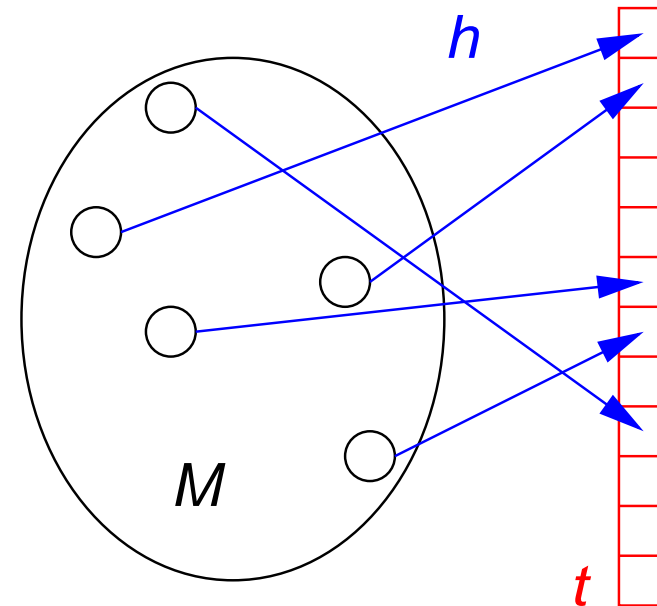
Hashing: Anwendungen

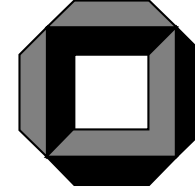
- Auslieferungsregale der UB Karlsruhe
- Entfernen exakter **Duplikate**
- Schach (oder andere kombinatorische Suchprogramme):
welche Stellungen wurden **bereits durchsucht** ?
- Symboltabelle** bei Compilern
- Assoziative Felder** bei Script-Sprachen wie perl oder awk
- Datenbank-Gleichheits-**Join**
(wenn eine Tabelle in den Speicher passt)
- Unsere Routenplaner: **Teilmengen** von Knoten,
z.B. Suchraum
- ...



Ein (über)optimistischer Ansatz

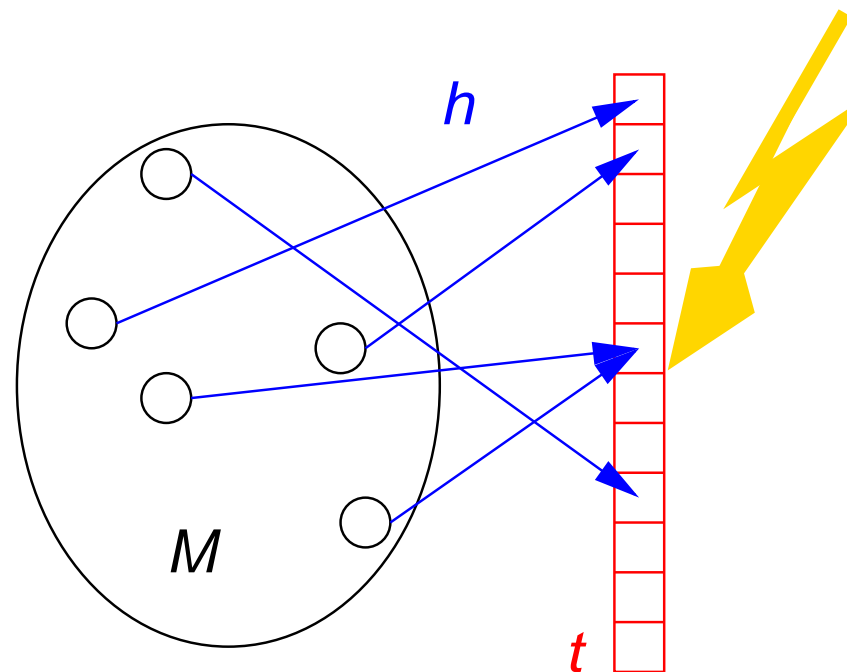
Eine perfekte **Hash-Funktion** h
bildet Elemente von M **injektiv**
auf eindeutige Einträge
der **Tabelle** $t[0..m-1]$ ab, d.h.,
 $t[h(\text{key}(e))] = e$



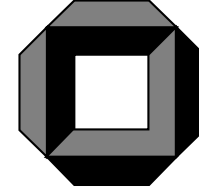


Kollisionen

Perfekte Hash-Funktionen sind schwer zu finden



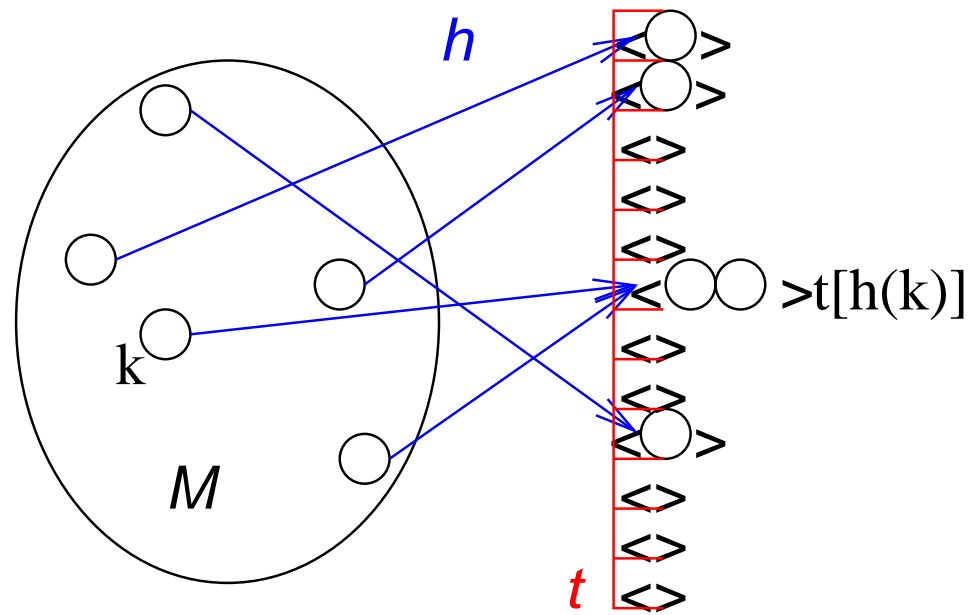
Beispiel: Geburtstagsparadox

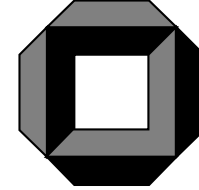


Kollisionsauflösung

Beispiel **geschlossenes Hashing**

Tabelleneinträge: Elemente \rightsquigarrow **Folgen** von Elementen





4.1 Hashing mit verketteten Listen

Implementiere die Folgen beim geschlossenen Hashing durch **einfach verkettete Listen**

insert(e): Füge e am Anfang von $t[h(e)]$ ein.

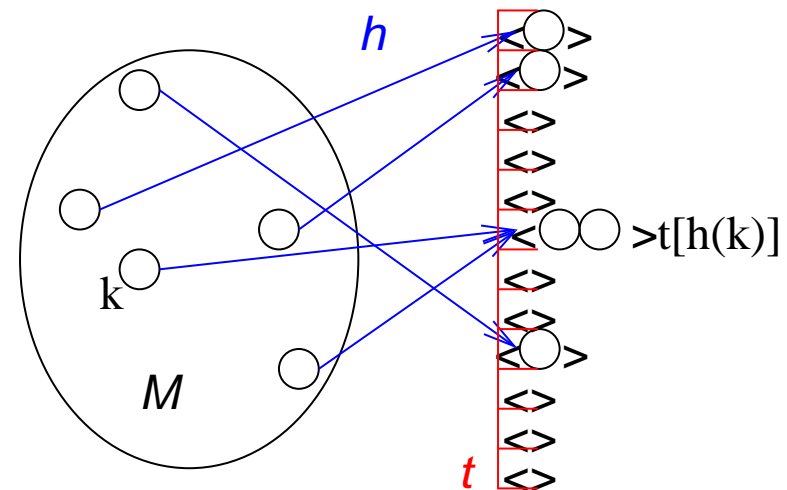
konstante Zeit

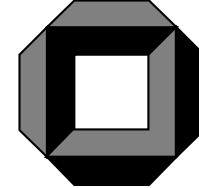
remove(k): Durchlaufe $t[h(k)]$. Element e mit $h(e) = k$ gefunden \rightsquigarrow löschen und zurückliefern.

find(k): Durchlaufe $t[h(k)]$. Element e mit $h(e) = k$ gefunden \rightsquigarrow zurückliefern.

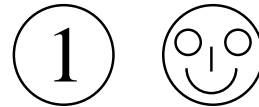
Sonst: \perp zurückgeben.

Zeit **$O(|M|)$** im schlechtesten Fall für remove und find





W.T. im Schnelldurchlauf



Hash-Beispiel

Elementarereignisse Ω zuf. Hash-Funktionen $\{0..m-1\}^{\text{Key}}$

Ereignisse: Teilmengen von Ω $\mathcal{E}_{42} = \{h \in \Omega : h(4) = h(2)\}$

$p_x =$ Wahrscheinlichkeit von $x \in \Omega$ Gleichverteilung $p_h = m^{-|\text{Key}|}$

$\mathbb{P}[\mathcal{E}] = \sum_{x \in E} p_x$ $\mathbb{P}[\mathcal{E}_{42}] = \frac{1}{m}$

Zufallsvariable $X_0 : \Omega \rightarrow \mathbb{R}$ $X = |\{e \in M : h(e) = 0\}|$

Erwartungswert $\mathbb{E}[X_0] = \sum_{y \in \Omega} p_y X(y)$ $\mathbb{E}[X] = \frac{|M|}{m}$ (*)

Linearität des Erwartungswerts: $\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$

Beweis von (*):

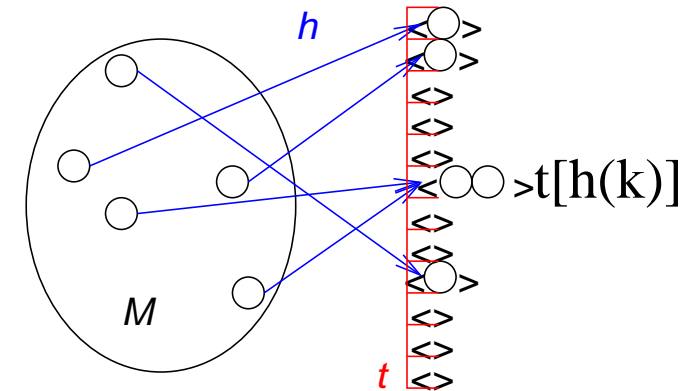
Betrachte die 0-1 ZV $X_e = 1$ für $h(e) = 0$, $e \in M$ und $X_e = 0$ sonst.

$$\mathbb{E}[X_0] = \mathbb{E}\left[\sum_{e \in M} X_e\right] = \sum_{e \in M} \mathbb{E}[X_e] = \sum_{e \in M} \mathbb{P}[X_e = 1] = |M| \cdot \frac{1}{m}$$



Analyse für zufällige Hash-Funktionen

Satz 1. Die erwartete *Ausführungszeit* von $\text{remove}(k)$ und $\text{find}(k)$ ist $\mathbf{O(1)}$ falls $|M| = \mathbf{O}(m)$.



Beweis. Konstante Zeit plus *Zeit für Durchlaufen von $t[h(k)]$* .

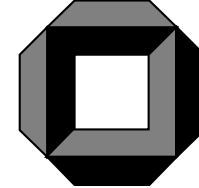
$$X := |t[h(k)]| = |\{e \in M : h(e) = h(k)\}|.$$

Betrachte die 0-1 ZV $X_e = 1$ für $h(e) = h(k)$, $e \in M$ und $X_e = 0$ sonst.

$$\begin{aligned} \mathbf{E}[X] &= \mathbf{E}\left[\sum_{e \in M} X_e\right] = \sum_{e \in M} \mathbf{E}[X_e] = \sum_{e \in M} \mathbb{P}[X_e = 1] = \frac{|M|}{m} \\ &= \mathbf{O}(1) \end{aligned}$$



Das gilt *unabhängig* von der Eingabe M .



4.2 Universelles Hashing

Idee: nutze nur bestimmte “einfache” Hash-Funktionen

Definition 1. $\mathcal{H} \subseteq \{0..m-1\}^{\text{Key}}$ ist *universell*

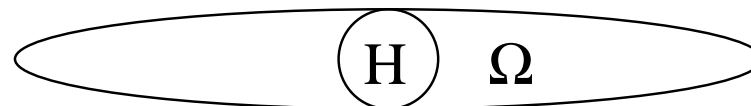
falls für alle x, y in Key mit $x \neq y$ und zufälligem $h \in \mathcal{H}$,

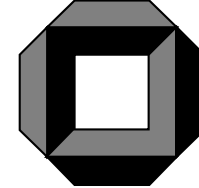
$$\mathbb{P}[h(x) = h(y)] = \frac{1}{m} .$$

Satz 2. *Theorem 1 gilt auch für universelle Familien von Hash-Funktionen.*

Beweis. Für $\Omega = \mathcal{H}$ haben wir immer noch $\mathbb{P}[X_e = 1] = \frac{1}{m}$.

Der Rest geht wie vorher. □





Eine einfache universelle Familie

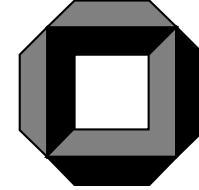
m sei eine Primzahl, $\text{Key} \subseteq \{0, \dots, m-1\}^k$

Satz 3. Für $\mathbf{a} = (a_1, \dots, a_k) \in \{0, \dots, m-1\}^k$ definiere

$$h_{\mathbf{a}}(\mathbf{x}) = \mathbf{a} \cdot \mathbf{x} \bmod m, H = \left\{ h_{\mathbf{a}} : \mathbf{a} \in \{0..m-1\}^k \right\}.$$

H ist eine universelle Familie von Hash-Funktionen

$$\left(\begin{array}{|c|c|c|} \hline x_1 & x_2 & x_3 \\ \hline * & * & * \\ \hline a_1 & a_2 & a_3 \\ \hline \end{array} \right) \text{mod } m = h_{\mathbf{a}}(\mathbf{x})$$



Beweis. Betrachte $\mathbf{x} = (x_1, \dots, x_k)$, $\mathbf{y} = (y_1, \dots, y_k)$ mit $x_j \neq y_j$
zähle \mathbf{a} -s mit $h_{\mathbf{a}}(\mathbf{x}) = h_{\mathbf{a}}(\mathbf{y})$.

Für jede Wahl von a_i s, $i \neq j$, \exists genau ein a_j mit $h_{\mathbf{a}}(\mathbf{x}) = h_{\mathbf{a}}(\mathbf{y})$:

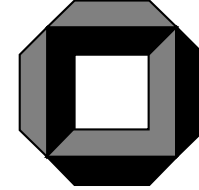
$$\begin{aligned}\sum_{1 \leq i \leq k} a_i x_i &\equiv \sum_{1 \leq i \leq k} a_i y_i \pmod{m} \\ \Leftrightarrow a_j (x_j - y_j) &\equiv \sum_{i \neq j, 1 \leq i \leq k} a_i (y_i - x_i) \pmod{m} \\ \Leftrightarrow a_j &\equiv (x_j - y_j)^{-1} \sum_{i \neq j, 1 \leq i \leq k} a_i (y_i - x_i) \pmod{m}\end{aligned}$$

m^{k-1} Möglichkeiten a_i auszuwählen (mit $i \neq j$).

m^k ist die Gesamtzahl \mathbf{a} s, d.h.,

$$\mathbb{P}[h_{\mathbf{a}}(\mathbf{x}) = h_{\mathbf{a}}(\mathbf{y})] = \frac{m^{k-1}}{m^k} = \frac{1}{m}.$$





Bit-basierte Universelle Familien

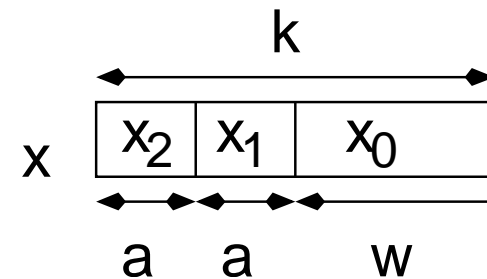
Sei $m = 2^w$, $\text{Key} = \{0, 1\}^k$

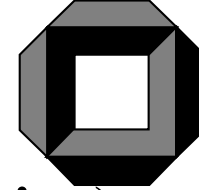
Bit-Matrix Multiplikation: $H^\oplus = \{h_{\mathbf{M}} : \mathbf{M} \in \{0, 1\}^{w \times k}\}$

wobei $h_{\mathbf{M}}(\mathbf{x}) = \mathbf{M}\mathbf{x}$ (Arithmetik mod 2, d.h., xor, and)

Tabellenzugriff: $H^{\oplus \square} = \{h_{(t_1, \dots, t_b)}^\oplus : t_i \in \{0..m-1\}^{\{0..w-1\}}\}$

wobei $h_{(t_1, \dots, t_b)}^\oplus((x_0, x_1, \dots, x_b)) = x_0 \oplus \bigoplus_{i=1}^b t_i[x_i]$





4.3 Hashing mit Linearer Suche (Linear Probing)

Offenes Hashing: zurück zur Ursprungsidee.

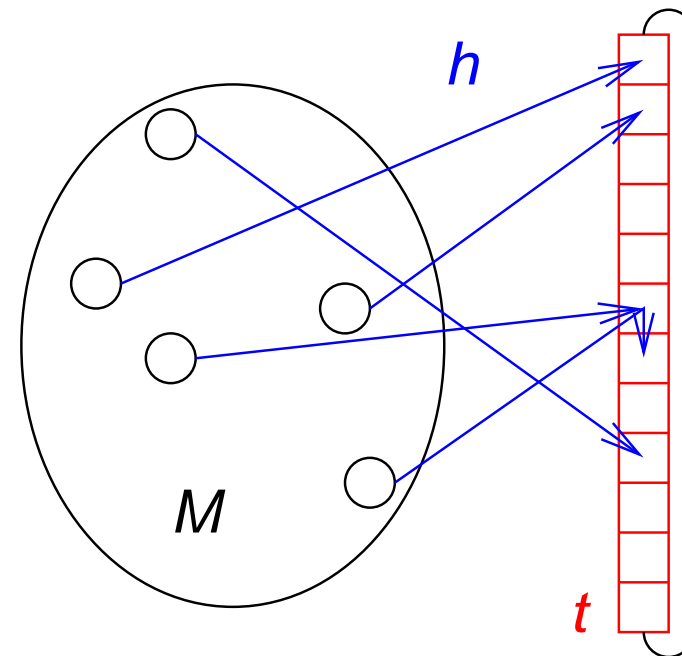
Elemente werden direkt in der Tabelle gespeichert.

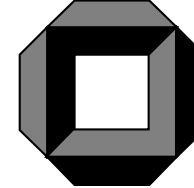
Kollisionen werden durch Finden anderer Stellen aufgelöst.

linear probing: Suche nächsten freien Platz.

Am Ende fange von vorn an.

- einfach
- Platzeffizient
- Cache-effizient





Der einfache Teil

Class BoundedLinearProbing($m, m' : \mathbb{N}; h : \text{Key} \rightarrow 0..m - 1$)

$t = [\perp, \dots, \perp] : \mathbf{Array} [0..m + m' - 1]$ **of** Element

invariant $\forall i : t[i] \neq \perp \Rightarrow \forall j \in \{h(t[i])..i - 1\} : t[j] = \perp$

Procedure insert($e : \text{Element}$)

for $i := h(e)$ **to** ∞ **while** $t[i] \neq \perp$ **do** ;

assert $i < m + m' - 1$

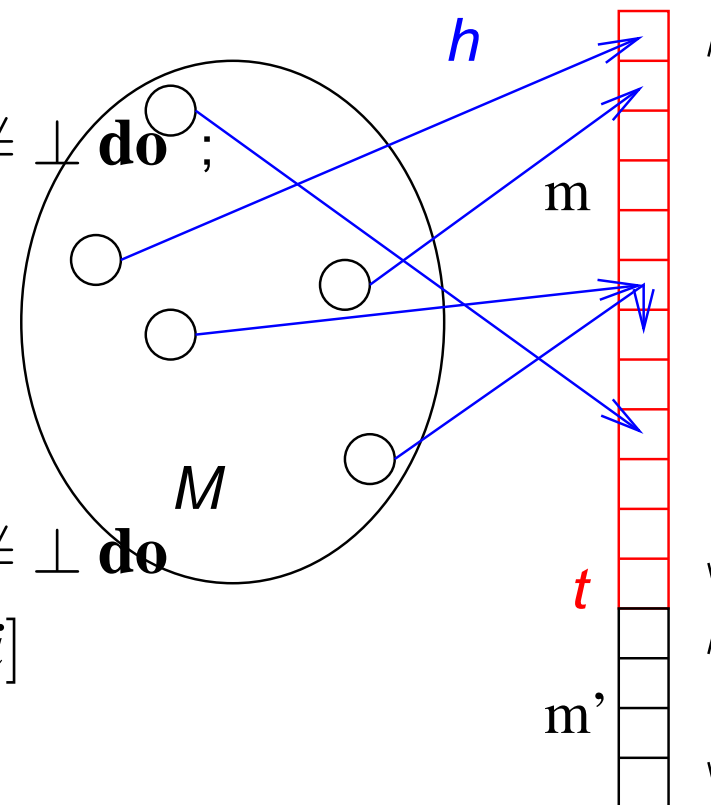
$t[i] := e$

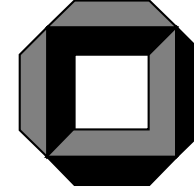
Function find($k : \text{Key}$) : Element

for $i := h(k)$ **to** ∞ **while** $t[i] \neq \perp$ **do**

if $t[i] = k$ **then return** $t[i]$

return \perp

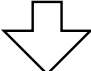




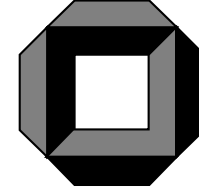
insert : axe, chop, clip, cube, dice, fell, hack, hash, lop, slash

an bo cp dq er fs gt hu iv jw kx ly mz
 tt 0 1 2 3 4 5 6 7 8 9 10 11 12

⊥	⊥	⊥	⊥	axe	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥
⊥	⊥	chop	⊥	axe	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥
⊥	⊥	chop	clip	axe	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥
⊥	⊥	chop	clip	axe	cube	⊥	⊥	⊥	⊥	⊥	⊥	⊥
⊥	⊥	chop	clip	axe	cube	dice	⊥	⊥	⊥	⊥	⊥	⊥
⊥	⊥	chop	clip	axe	cube	dice	⊥	⊥	⊥	⊥	fell	⊥
⊥	⊥	chop	clip	axe	cube	dice	⊥	⊥	⊥	hack	fell	⊥
⊥	⊥	chop	clip	axe	cube	dice	hash	⊥	⊥	⊥	fell	⊥
⊥	⊥	chop	clip	axe	cube	dice	hash	lop	⊥	hack	fell	⊥
⊥	⊥	chop	clip	axe	cube	dice	hash	lop	slash	hack	fell	⊥

remove  clip

⊥	⊥	chop	clip	axe	cube	dice	hash	lop	slash	hack	fell	⊥
⊥	⊥	chop	lop	axe	cube	dice	hash	lop	slash	hack	fell	⊥
⊥	⊥	chop	lop	axe	cube	dice	hash	slash	slash	hack	fell	⊥
⊥	⊥	chop	lop	axe	cube	dice	hash	slash	⊥	hack	fell	⊥



4.4 Verketteten \leftrightarrow Lineare Suche

Volllaufen: Verketteten weniger empfindlich.

Unbeschränktes **offenes** Hashing hat nur amortisiert konst.

Einfügezeit

Cache: Lineare Suche besser. Vor allem für **doall**

Platz/Zeit Abwägung: Kompliziert! Abhängig von n , **Füllgrad**,

Elementgröße, Implementierungsdetails bei Verketteten

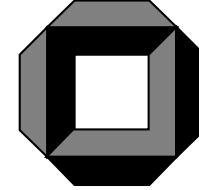
(shared dummy!, t speichert Zeiger oder item),

Speicherverwaltung bei Verketteten, beschränkt oder nicht,...

Referentielle Integrität: Nur bei Verketteten !

Leistungsgarantien: Universelles Hashing funktioniert so nur mit

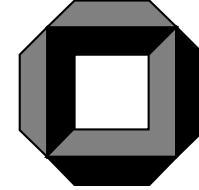
Verketteten



4.5 Perfektes Hashing

Idee: mache h injektiv.

Braucht $\Omega(n)$ bits Platz !



Here: Fast Space Efficient Hashing

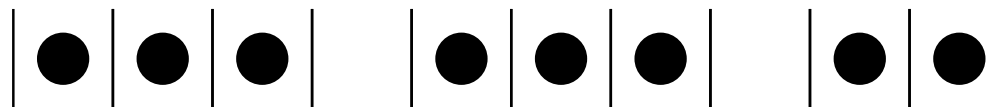
[Fotakis, Pagh, Sanders, Spirakis]

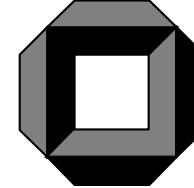
Represent a set of n elements (with associated information) using space $(1 + \epsilon)n$.

Support operations **insert**, **delete**, **lookup**, (doall) efficiently.

Assume a truly random hash function h

(Trick [Dietzfelbinger, Weidling 2005]: läßt sich rechtfertigen.)

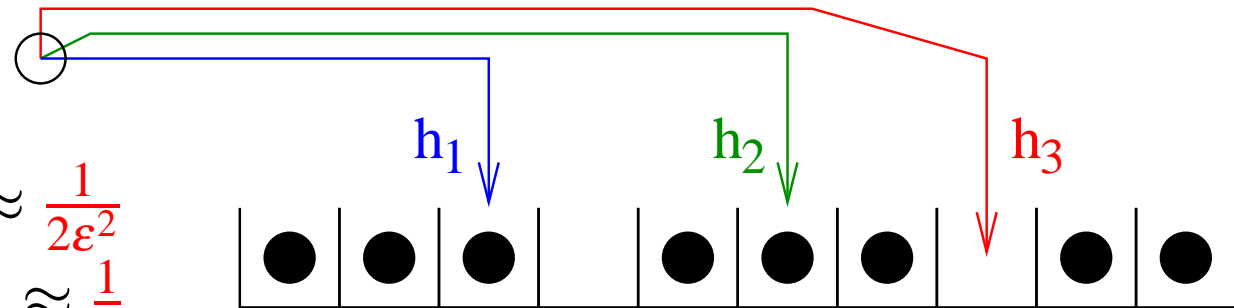




Related Work

Linear probing: $E[T_{\text{find}}] \approx \frac{1}{2\varepsilon^2}$

Uniform hashing: $E[T_{\text{find}}] \approx \frac{1}{\varepsilon}$

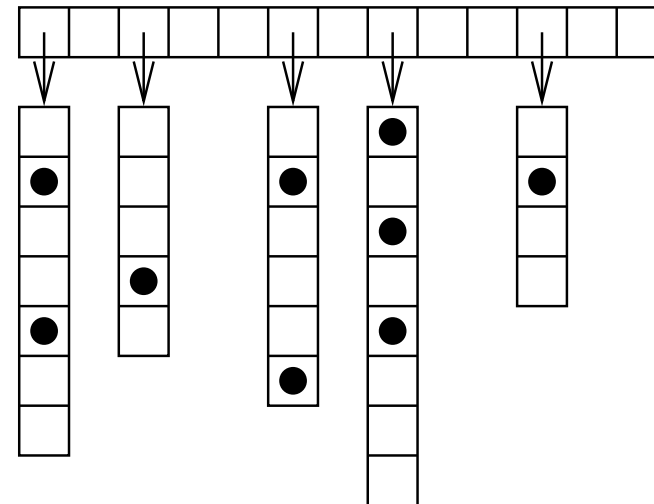


Dynamic Perfect Hashing,

[Dietzfelbinger et al. 94]

Worst case constant time

for lookup but ε is not small.

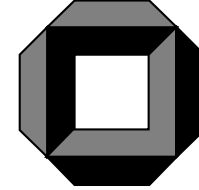


Approaching the Information Theoretic Lower Bound:

[Brodnik Munro 99, Raman Rao 02]

Space $(1 + o(1)) \times$ lower bound without associated information

[Botelho Pagh Ziviani 2007] static case.



Cuckoo Hashing

[Pagh Rodler 01]

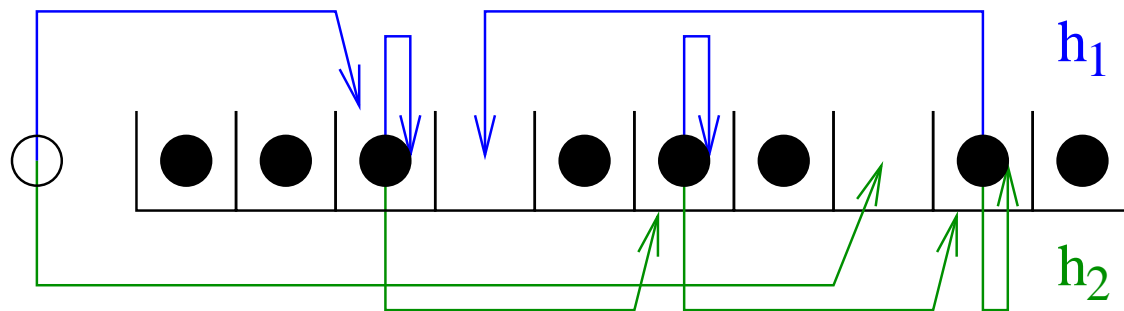
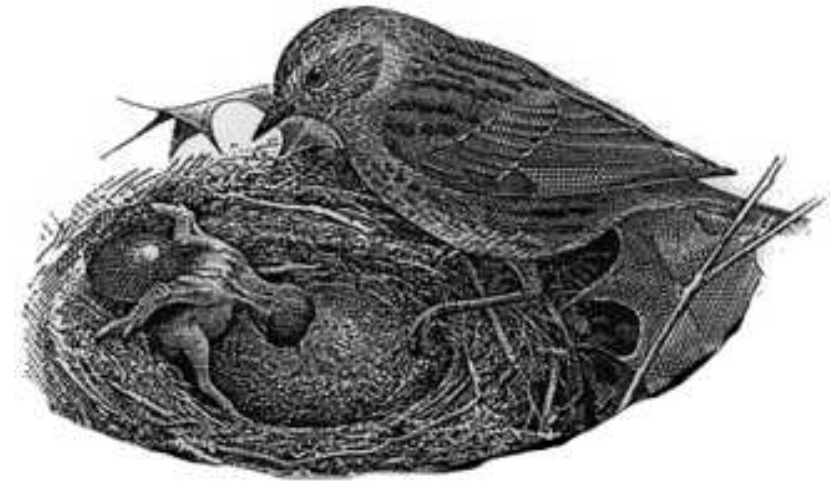
Table of size $(2 + \epsilon)n$.

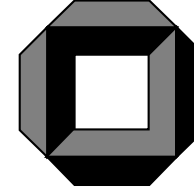
Two choices for each element.

Insert moves elements;
rebuild if necessary.

Very fast lookup and insert.

Expected constant insertion time.





d -ary Cuckoo Hashing

d choices for each element.

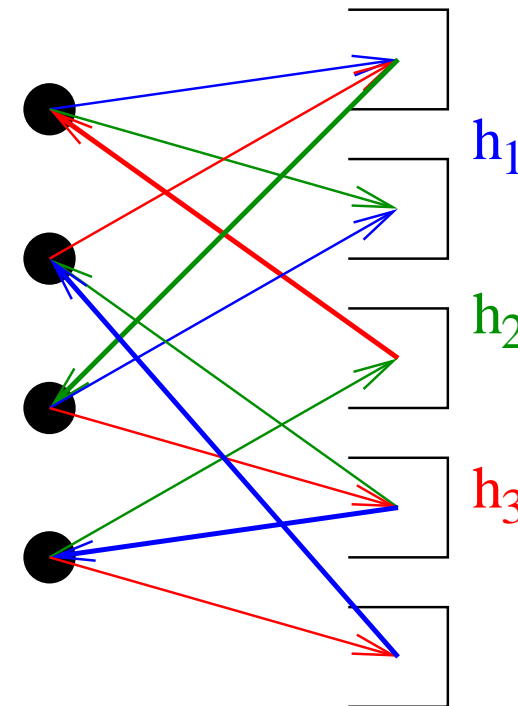
Worst case d probes for **delete** and **lookup**.

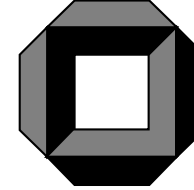
Task: maintain **perfect matching**

in the **bipartite graph**

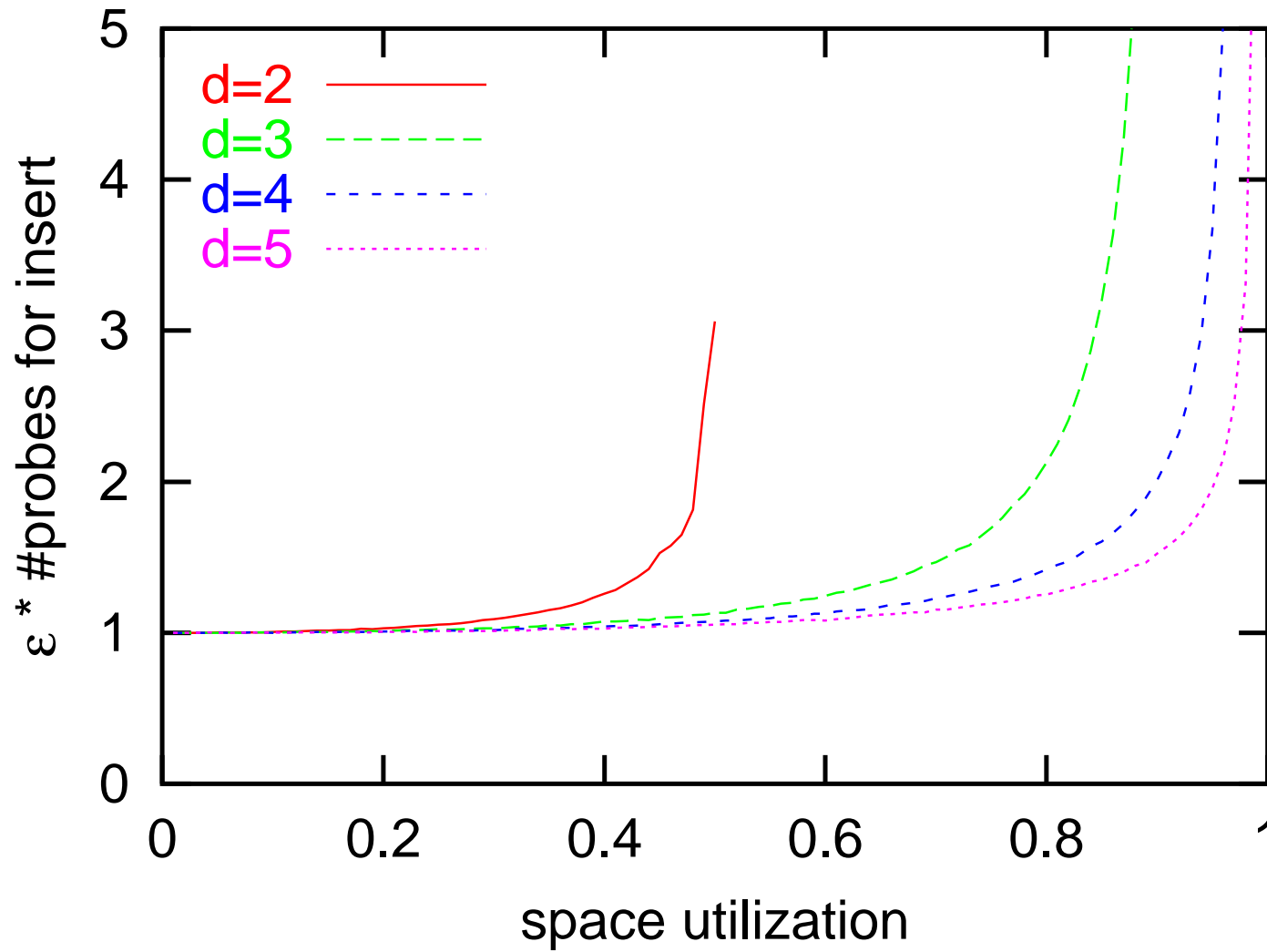
($L = \text{Elements}$, $R = \text{Cells}$, $E = \text{Choices}$),

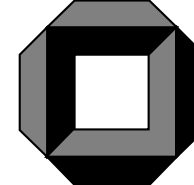
e.g., **insert** by **BFS**.





Experiments



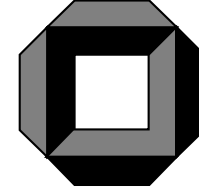


Analysis

Lookup and Delete: $d = O\left(\log \frac{1}{\varepsilon}\right)$ probes

Insert: $\left(\frac{1}{\varepsilon}\right)^{O(\log(1/\varepsilon))}$, (experiments) $\longrightarrow O(1/\varepsilon)$?

Approach: Properties of random bipartite graphs / hypergraphs.



4.6 Mehr Hashing

- Hohe Wahrscheinlichkeit und Garantien für den **schlechtesten Fall**
 \rightsquigarrow höhere Anforderungen an die Hash-Funktionen
- Hashing als Mittel zur Lastverteilung z.B., storage servers,
 (peer to peer Netze, . . .)
- Verschiedene Plattengrößen und Geschwindigkeiten
- Platten hinzufügen / ersetzen
- $O(1)$ find / perfektes Hashing