

---

---

# Perfekte Hashfunktionen

—  
so kompakt wie es geht

---

---

Yaojia Hou & Christian Schwarz

*Fachgebiet Komplexitätstheorie,  
Institut für Theoretische Informatik,*

Technische Universität Ilmenau

18. November 2021

# Das Hashing-Problem

---

Gegeben:  $U$  .. Universum,  
 $S \subseteq U$  .. Schlüssel  
 $m = (1 + \epsilon)|S|$

Konstruiere: Abbildung  $H : U \rightarrow [m]$ , mit  $H|_S$  injektiv  
Datenstruktur  $DS$  zum Speichern von  $H$

Ziel: Speicherplatz für  $DS$  klein  
Auswerteszeit für  $H$  konstant  
auch für große  $|S|$

# Roadmap

---

Einführung

Konstruktion

Partitionierung

(große Schlüsselmengen)

Auswertung & Fazit

Einführung

Konstruktion

Partitionierung

(große Schlüsselmengen)

Auswertung & Fazit

# Konstruktion

---

für jeden Schlüssel  $x \in S$ , berechne  $k$  Zahlen  $h_x^1, \dots, h_x^k \in [m]$

$$B_x = (h_x^1, \dots, h_x^k)$$

gesucht

Abbildung  $\sigma : S \rightarrow [k]$  so, dass

$$S \rightarrow [m] : x \mapsto h_x^{\sigma(x)} \text{ injektiv ist}$$

## Konstruktion ('2)

---

Wähle abhängig von  $\epsilon$  einen von zwei Algorithmen:

$$m = (1 + \epsilon)n \text{ mit } \epsilon \in \{1, 0.23, 0.04, 0\}$$

$\epsilon = 0.23 \longrightarrow$  *peeling* eines Hypergraphen Botelho et al. [BPZ13]

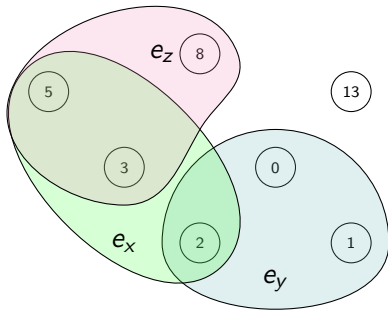
sonst  $\longrightarrow$  *Cuckoo Insertion*  
+ Lösen eines LGS VL "Hashing"

$\epsilon = 0 \longrightarrow$  arbeite mit  $m = 1.04n$  & konstruiere Rang DS

# Peeling

geg.: für jeden Schlüssel  $x \in S$ ,  $k$  Zahlen  $h_x^1, \dots, h_x^k \in [m]$

die Tupel  $e_x := (h_x^1, \dots, h_x^k)$  bilden Kanten eines ( $k$ -uniformen) Hypergraphen  $H$



## Peeling ('2)

---

### Ergebnis

Aufzählung  $\underline{e} = (e_1, e_2, \dots)$  der Kanten von  $H$  so, dass:

$e_{i+1}$  hat Ecke vom Grad 1 in  $H - \{e_1, \dots, e_i\}$

kann scheitern, falls  $H$  'Kreise' hat

Algorithmus:

entferne sukzessive Kanten  $e$  mit einer Ecke  $v$  vom Grad 1

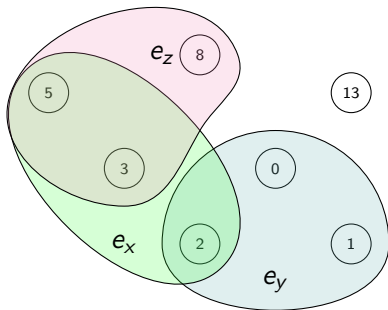


## Peeling (Beispiel)

---

Algorithmus:

entferne sukzessive Kanten  $e$  mit einer Ecke  $v$  vom Grad 1



$$\rightarrow \underline{e} = (e_z, e_x, e_y)$$

$$\underline{v} = (8, 3, 1)$$

# Darstellung von $\sigma$

---

$\sigma(x)$  – Auswahl eines Hashwertes  $h_x^1, \dots, h_x^k$

peeling liefert Aufzählung  $\underline{e}$  und Folge  $\underline{v}$

$\underline{v}_i$  pw. verschieden  
' $\underline{v}_i \in \underline{e}_i$ '

für  $x \in S$  sei  $j$  der Index von  $e_x = (h_x^1, \dots, h_x^k)$  in  $\underline{e}$   
wähle  $\sigma(x)$  so, dass  $\underline{v}_j = h_x^{\sigma(x)}$

Problem: Wie ist  $\sigma$  zu speichern?

## Darstellung von $\sigma$ ('2)

---

Array  $\text{sigma}[m]$  von Werten aus  $\{0, \dots, k\}$

$\sim$  Ecken-Label in  $H$

Idee

$$\sigma(x) - 1 \stackrel{!}{=} (\text{sigma}[h_x^1] + \dots + \text{sigma}[h_x^k]) \text{ MOD } k$$

$\sigma(x) \in \{1, \dots, k\}$

initial: alle Einträge =  $k$

durchlaufe  $\underline{e}$  in *umgekehrter Reihenfolge* ( $j = n..1$ )

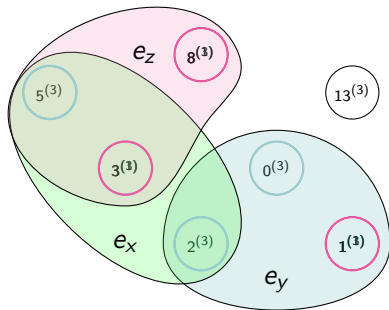
und passe  $\text{sigma}[\underline{v}_j]$  an

# Darstellung von $\sigma$ (Beispiel)

$$\underline{e} = (e_z, e_x, e_y) \text{ und } \underline{v} = (8, 3, 1)$$

$$\mathfrak{g} = h^2 \sigma_{\underline{v}}[h^3 \underline{e}] = [3, 1, 3, 1, \dots]$$

$$2 - 1 \equiv 3 \quad 3 + 3 + 3$$



Ecken in  $e$  seien der Größe nach sortiert,  
also  $h_x^1 < \dots < h_x^3$

## Variante aus der VL (Übersicht)

geg.: für jeden Schlüssel  $x \in S$ ,  $k$  Zahlen  $h_x^1, \dots, h_x^k \in [m]$

berechne  $\sigma$  mittels *Cuckoo Hashing*

konstruiere Matrix  $A \in \{0, 1\}^{S \times [m]}$  mit Zeilen der Form

$$\boxed{00 \dots 0 \quad \underbrace{1 * \dots *}_{64} \quad 00 \dots 0} \quad * \in \{0, 1\}$$

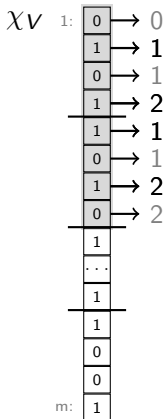
löse das LGS

$$A \bar{z} = (\sigma(x))_{x \in S} \quad \bar{z} \in \{0, 1\}^{m \times \lceil \log k \rceil}$$

# Rang Datenstruktur

Für Menge  $V \subseteq [m]$  sei  $\text{rank}_V(x) := |\{v \in V : v \leq x\}|$

Idee: speichere Char. Vektor von  $V$



→ ersetze 0101 und 1010  
durch 'Abbildung' [1, 1, 2, 2]

teile  $\{0000, \dots, 1111\}$  in Klassen auf  
→ 8 Klassen  
→ 3 Bit (statt 4)  
→ *0.75m Bit* insgesamt

Einführung

Konstruktion

Partitionierung

(große Schlüsselmenen)

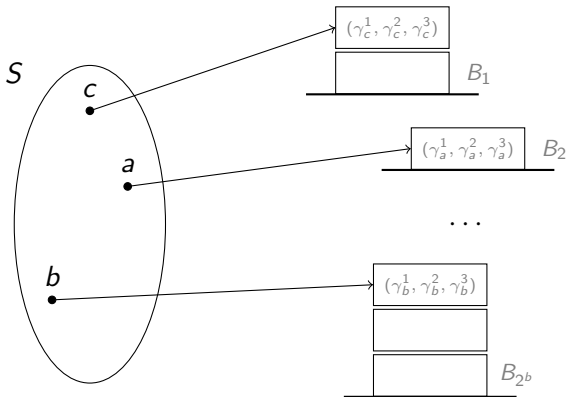
Auswertung & Fazit

# große Schlüsselmengen

Teile  $S$  auf in  $2^b$  *buckets*  $B_1, B_2, \dots$

Berechne für jedes  $x \in S$  Werte  $\gamma_x^1, \gamma_x^2, \gamma_x^3$  aus  $[2^{16}]$

(Identität der Schlüssel wird danach 'vergessen')





## große Schlüsselmengen ('2)

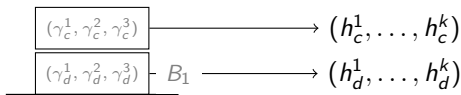
---

Konstruiere (M)PHF  $H_i$  für jeden bucket  $B_i$

$$(n = |B_i|, m = (1 + \epsilon) \max |B_j|)$$

$h_x^1, \dots, h_x^k$  werden aus den  $\gamma_x^i$  berechnet

(benötigt große Tabellen!)



verschiebe  $H_i(x)$  um  $(i - 1) \times m$  bzw.  $\sum_{j < i} |B_j|$

Einführung

Konstruktion

Partitionierung

(große Schlüsselmengen)

Auswertung & Fazit

# Übersicht der Ergebnisse

$n = 3\text{Mio.}$  Schlüssel (Zeichenketten der Länge 10)

Konfiguration		Speicherplatz [Bit/Schlüssel]	Auswertezeit
$m=1.23n$	EM-T	1.98	$24\mu\text{s}$
	RAM	1.97	30ns
$m=1.04n$	EM-T	2.09	$33\mu\text{s}$
	RAM	—	—
$m=n$	EM-T	2.98	$35\mu\text{s}$
	RAM	—	—
	[BPZ13]	2.67	$\approx 30\text{ns}$

# Übersicht der Ergebnisse ('2)

---

$n = 3\text{Mio.}$  Schlüssel (Zeichenketten der Länge 10)

Konfiguration		Speicherplatz [Bit/Schlüssel]	Auswertezeit
$m=1.23n$	EM-T	1.98	$24\mu\text{s}$
	EM+T	14.56	$1\mu\text{s}$
$m=1.04n$	EM-T	2.09	$33\mu\text{s}$
	EM+T	18.87	$2\mu\text{s}$
$m=n$	EM-T	2.98	$35\mu\text{s}$
	EM+T	19.75	$2\mu\text{s}$

# Übersicht der Ergebnisse ('3)

---

Schlüssel: Wörter mit 20..40 Buchstaben

Konfiguration		Speicherplatz		
n	b	m=1.23n	m=1.04n	m=n
1M	4	2.03	2.09	2.98
100M	11	2.00	2.11	2.99
1.75G	15	—	—	—

## Perfekte Hashfunktionen

Hinsichtlich *Speicherplatz* und *Auswertezeit*:

Algo. von Botelho et Al. besser (1.97 B/S,  $< 0.1\mu\text{s}$ )

aber: 23% größerer Bildbereich

Algo. aus VL nur 4..5%, dafür 2.1 B/S und  $\approx 30\mu\text{s}$

→ +0.13 Bit/Schlüssel

→ **aber!** 1000× langsamer

## Fazit ('2)

---

### Minimal Perfekte Hashfunktionen

Algo. von Botelho et Al.: 2.62 B/S und  $< 1\mu s$  (hier nicht getestet)

Algo. aus VL: 2.98 B/S und  $\approx 30\mu s$

→ +0.36 B/S

→  $\approx 100\times$  langsamer

## Fazit ('3)

---

zwei Algorithmen implementiert und verglichen

Speicherplatz: ähnlich ( $< 0.4$  B/S Differenz)

Auswertzeit: Algo. aus VL zu langsam

Problem: nur für  $n < 100\,000$  ist  $A \bar{z} = (\sigma(x))_x$  lösbar

bessere Ergebnisse bei anderer Wahl von  $A$ ?

Es geht noch besser ..

(RecSplit / PTHash)