

Kap. 7 Sortierverfahren

Kap. 7.0 Darstellung, Vorüberlegungen

$\Sigma = \{(x, \alpha)\}$ mit $<$

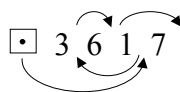
internes Sortieren
 Σ paßt

externes Sortieren
 Σ paßt nicht in AS

Darstellung von $<$:

3 6 1 7 ungeordnet, $<$ berechnen

1 3 6 7 $< \sim$ phys. Reihung



Liste verkettet

sortierter Baum mit Inordnung

- mit Zeigern
- Adreßrechnung

1

1. Einfache Feldorganisation:

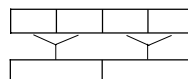
$S[1], \dots, S[n]$

Für Sortieren am Platz: Tausch-Operation!

2. Doppelfeld:

$S[1], \dots, S[n]$

$T[1], \dots, T[n]$



für Misch-Verfahren:

Speicher: $2n$;

Zeit: $O(n \log n)$

Sätze variabler Länge?

2

3. Listenorganisation:

(x, α) erweitern um Verweis v auf Nachfolger

Tauschoperation?

Anpassung von Quicksort?

4. Zeigerfeld: für lange Sätze und

Sätze variabler Länge. Sortiere Zeiger!

v_x  (x, α)

v_x gespeichert in Feld $V[1] \dots V[n]$

$v_x < \cdot v_y \Leftrightarrow x < y$

sortiere V nach $< \cdot$, kopiere S nach $< \cdot$

3

5. Feld von Paaren: (x, v_x)

falls α variabel lang

x feste Länge

speichere (x, v_x) in V ,

sortiere V nach $<$ auf X .

Auch o.k. falls Σ nicht in AS paßt,

aber $\Sigma_v = \{(x, v_x)\}$ paßt

Problem: gut für random Zugriff auf Σ , aber sequentielle Verarbeitung von Σ nach $<$ problematisch: 1 Plattenzugriff pro Element von Σ .

4

Kap. 7.1 Zu internen Sortierverfahren

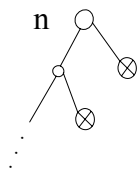
Einfache Verfahren mit $O(n^2)$

- Auswahl
- Einfüge (bubble)
- Austausch

Quicksort: Zeit?

Speicher: Kaskadenrekursion erfordert Rückstellung einer

Partitionsklasse für Laufzeitkeller



⊗ zurückgestellt, Kellertiefe?

? $2n$

? $n \log n$

5

Baumsortierung: Heapsort (Stadelsort)

Speicher: am Platz, schichtenweise sequentielle Darstellung
des Binär-Baums

Zeit: $O(n \log n)$ garantiert

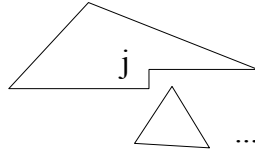
Zentral für externe Sortierverfahren

Umwandlung heap in sortiertes Feld:

$M[1:n]$ Baum $\Rightarrow M[1:i]$ mit $i \leq n$ ist Baum mit
denselben Vater-Sohn Beziehungen

6

Lemma: Falls $M[1:n]$ Stadel ist, ist auch $M[1:i]$ und $M[j:n]$ Stadel für alle i, j .



Sort – Merge intern:

Zeit: $O(n \log n)$

Speicher: $2n$ bzw. $1.5n$

7

Kap. 7.2 Externes Sortieren

Phase 1: Initialläufe: sortierte Teilmengen durch internes Sortieren

$$A = (a_1, a_2, \dots, a_{e_1}); \quad i < j \Rightarrow a_i \leq a_j$$

$$B = (b_1, b_2, \dots, b_{e_2})$$

$$C = \overline{AB} \text{ durch Mischung (merge) von A und B}$$

Seien a_j, b_i kleinste, ungemischte Elementen von A, B

$$c_{k+1} := \min\{a_j, b_i\}$$

Bei Mischung von m Initialläufen z.B. $m = 5000$

$$c_{k+1} := \min_{i=1 \dots 5000} \{a_i\}$$

verwende Heap!

8

Phase 2: Reines Mischen: Datei mit N Elementen

0. Läufe der Länge $1 = 2^0$ vorhanden

1. Läufe der Länge $2 = 2^1$

⋮

d. Läufe der Länge 2^d

für kleinstes d mit $2^d \geq N \Rightarrow d = \lceil \log_2 N \rceil$

Ziel: bei 1. Durchgang möglichst lange Initialläufe durch internes Sortieren

Bei S Speicher:

Länge S: in place Verfahren, feste Satzlänge

Länge S/2: internes Mischen

Länge 2S: Stadelort, erspart bis zu 2 externe Mischdurchgänge 9

Mit Auswahlverfahren:

M [1] ... M [s]

1. min. suchen i

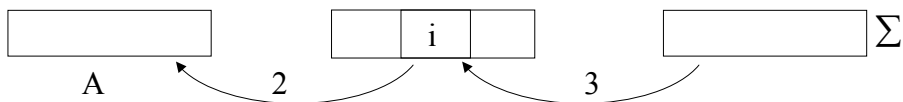


$\min_i \{M[i], M[i] \geq y\}$

2. Ausgabe: $y := M[i]$

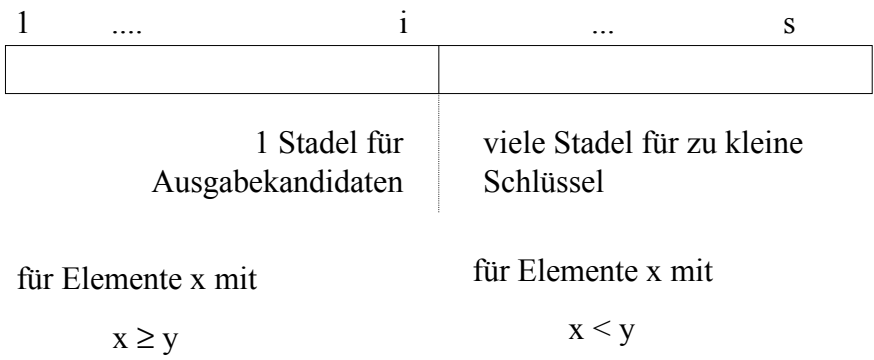
3. Eingabe: $M[i] := x$

x kann an nächster min-Suche teilnehmen, falls $x \geq y$



Problem: min-Suche in M von $O(S)$

Mit Heap:



11

Algorithmus für Initialläufe:

```

var L: {für letzten ausgegebenen Schlüssel}
H: {für neue Eingabe von  $\Sigma$ }
{Stadel in M[1..S]}
starte neuen Initiallauf:
i := S;
while i > 0 do {Ausgabestadel M[1..i] ≠  $\emptyset$ }
begin L := M[1]; gebe L aus;
H := Eingabe vom  $\Sigma$ ;
if H ≥ L then {H noch zu Initiallauf}
begin M[1] := H; SENKE (1,S);
else begin M[1] = M[i];
M[i] = H;
SENKE (i,S);
i := i-1;
SENKE (1,i) Ver- besserung?
end
end {keine E/A Stockung!!!}
    
```

12

Externes Mischen mit Bändern

2 Wege, 3 Bänder:

$$U = A_1, A_2, A_3$$

$$V = B_1, B_2, B_3$$

$$W = C_1 = \overline{A_1 B_1}, C_2 = \overline{A_2 B_2}, C_3 = \overline{A_3 B_3}$$

Merge von C_1, C_2, C_3

1. Symmetrisches Mischen, vollständiges Kopieren

$$\{W = C_1, C_2, C_3\}$$

$$U = C_1, C_3$$

$$V = C_2, \phi$$

$$W = \overline{C_1 C_2}, C_3$$

$$U = \overline{C_1 C_2}$$

$$V = C_3$$

$$W = \overline{\overline{C_1 C_2} C_3}$$

=> 6 Durchgänge

Kopiervorgänge für C_3 ?

13

Aufwand für N Initialläufe:

$\lceil \log_2 N \rceil$	Mischdurchgänge	$\lceil \log_m N \rceil$
$\lceil \log_2 N \rceil - 1$	Kopierdurchgänge	$\lceil \log_m N \rceil - 1$
1	Initiallauf-Durchgang	1

gesamt: $2 \cdot \lceil \log_m N \rceil$ Durchgänge

1 Durchg. = 1 x lesen ganz Σ

1 x schreiben

für m-Wege Mischen, m+1 Bänder

14

2. Symmetrisches Mischen, unvollständiges Kopieren

nach 1. Mischdurchgang:

$$\{W = C_1, C_2, C_3\}$$

$$U = C_1$$

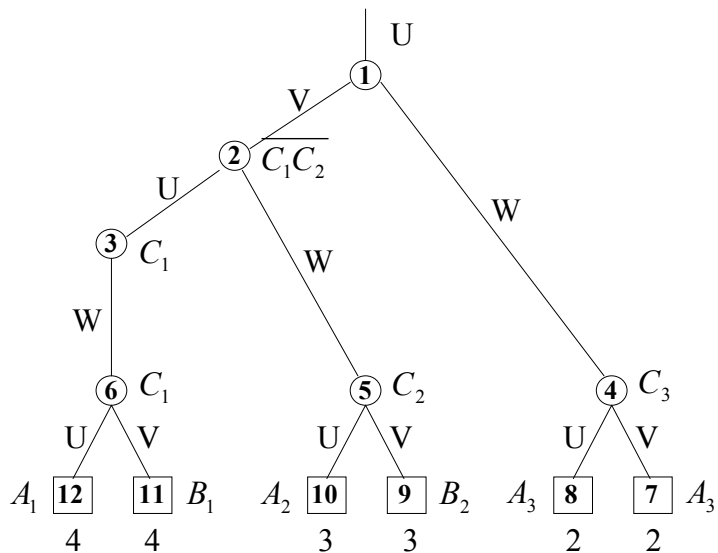
$$V = \overline{C_1 C_2}$$

$$U = \overline{\overline{C_1 C_2 C_3}}$$

=> 4 Durchgänge statt 6

15

Mischbaum nach Knuth



16

Symmetrisches Mischen, 2m Bänder

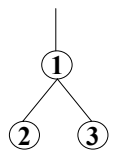
fast ohne Kopieren,

je m Bänder alternierend als Eingabe und Ausgabe

Baum für vorherige Initialläufe, 4 Bänder, symm. Mischen Kosten?

Mehrphasen Mischen (polyphase):

Sortieren ohne zu kopieren?



U	V	W
0	0	<u>1</u>
1	<u>1</u>	0
<u>2</u>	0	1
0	2	3

letzter Mischdurchgang, Zustände vorher?

Rekursionsgleichung für Anzahl Läufe auf Band

$$0 \quad C_i \quad C_{i+1}$$

$$C_{i+1} \quad (C_{i+1} + C_i) \quad 0$$

$$= C_{i+2}$$

i.e Fibonacci

17

Problem: Länge der entstehenden Zwischenläufe:

z.B. 5-Wege, 65 Läufe

0	0	0	0	0	<u>1₆₅</u>	65
1 ₁	1 ₅	1 ₉	1 ₁₇	<u>1₃₃</u>	0	33
2 ₁	2 ₅	2 ₉	<u>2₁₇</u>	0	1 ₁	34
4 ₁	4 ₅	<u>4₉</u>	0	2 ₁	3 ₁	36
8 ₁	<u>8₅</u>	0	4 ₁	6 ₁	7 ₁	40
16 ₁	0	8 ₁	12 ₁	14 ₁	15 ₁	65
						273

18

Ohne Rückspulung:

Bänder vorwärts schreiben, \uparrow
 rückwärts lesen: \downarrow

0	0	1 \uparrow
1 \downarrow	1 \downarrow	0
2 $\downarrow\uparrow$	0	1 \uparrow
0	2 $\downarrow\uparrow$	3 $\uparrow\downarrow\uparrow$

Allg. Fall.: $m+1$ Bänder
 m Wege
 $m-1$ Fib. Zahlen für nächste Phase

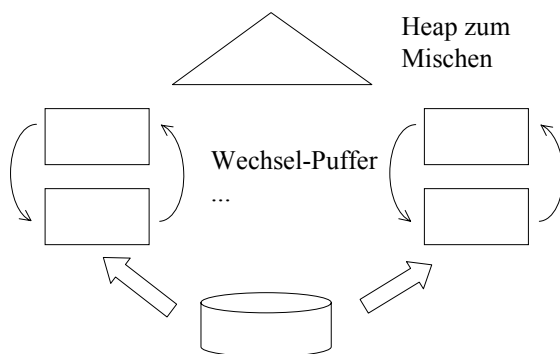
$$\begin{array}{ccccccc}
 C_{k-m+1} & \dots & C_{k-m+i} & \dots & \underline{C_k} & 0 \\
 (C_{k-m+1} + C_k) \dots & & (C_{k-m+i} + C_k) \dots & & 0 & C_k
 \end{array}$$

i.e. $C_{k-m+i} + C_k = C_{k+i}$ für $i = 1, \dots, m-1$

$m-1$ neue Fib.-Zahlen der Ordnung $m-1$

Kap. 7.3 Sortieren, Platten, Parallelität

Annahme: AS mit 10^n Bytes,
 Seitengröße 10^3 Bytes
 i.e. 10^{n-3} Seiten
 Länge Initilläufe: $2 \cdot 10^n$ Bytes



Mischgrad:
 $\frac{1}{2} \cdot 10^{n-3} = 5000$
 bei AS = 10MB
 1 Mischlauf:
 $\frac{1}{2} \cdot 10^{n-3} \cdot 2 \cdot 10^n$
 $= 10^{2n-3}$ Bytes

21

Workstations: $n=7$

d.h. 10^{11} B = 100 GB in 1 Lauf sortiert
 Problem: UNIX File-System

Mainframes: $n=8$

d.h. 10^{13} B = 10 TB in 1 Lauf sortiert
 Fazit: Bei heutiger Rechner-Technologie ist Sortieren
 ein linearer Algorithmus!

CPU-Last: AS = 10^n B

1 Element sei z.B. 100 Bytes
 \Rightarrow Heap hat 10^{n-2} Knoten,
 z.B. Bei $n=7$: 10^5 Knoten

22

Heap Höhe: $\log_2 10^5 \approx 17$ SENKE Aufrufe pro Element
mit je 4 Vergleichen:

string-Vergleich:

≈ 500 Befehle $\left\{ \begin{array}{l} - \text{extrahiere Attribut aus Tupel, Additionen,} \\ \text{Multiplikationen} \\ - \text{Adressrechnung für } M[k], M[2k], M[2k+1] \end{array} \right.$

\Rightarrow Pro El. $\approx 10^4$ Befehle?

Festplatte: 500 KB/s Nutzleistung $\sim \sim$ 5000 El/s

CPU: 5000 El/s * 10^4 Befehle/El = 50 MIPS oder mehr

\Rightarrow Potential für paralleles Sortieren!

1991: CPU bremst um Faktor 5!!

2000: CPU bremst bei RAID

23

Beispiel: Sortiere 100 MB auf SUN 4

Initialläufe: 1 x lesen	≈ 200 s
1 x schreiben	$\approx \underline{200}$ s
	400 s

Beim Einlesen Bremsfaktor 5
wegen langsamer CPU: **1000 s**
Prozessor-Engpass!!

Mischlauf: 1 x lesen	} random Zugriff zur Platte
1 x schreiben	

≤ 50 Seiten/s, d.h. Plattenleistung:
 $\leq \leq 50$ KB/s bei 1KB/Seite
 $\leq \leq 200$ KB/s bei 4 KB/Seite

24

1 KB/Seite: 1 x lesen ~ 2000 s
1 x schreiben ~ 2000 s
für Mischen 4000 s

4 KB/Seite: 1 x lesen ~ 500 s
1 x schreiben ~ 500 s
für Mischen 1000 s

Mischgrad 5 bis 20, kleiner Heap, CPU bremst nicht!
Mischphase ist Plattenengpass!!

Gesamtzeit: 2000 s mit 4 KB/Seite
5000 s mit 1 KB/Seite
5000 s = 83 min

Hinweis: bei niedrigem Mischgrad, z.B. 10, verwende
maximale E/A-Pufferung, mehrere Seiten pro
Wechselpuffer:

25