

Informatik II: Algorithmen und Datenstrukturen SS 2015

Vorlesung 2b, Mittwoch, 29. April 2015
(Sortieren in Linearzeit, Untere Schranke $n \cdot \log n$)

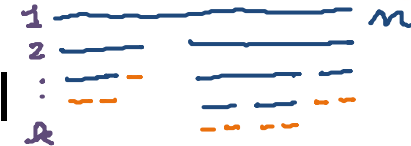
Prof. Dr. Hannah Bast
Lehrstuhl für Algorithmen und Datenstrukturen
Institut für Informatik
Universität Freiburg

Blick über die Vorlesung heute

- Laufzeitanalyse Sortieren, Teil 2
 - Quicksort: Laufzeitanalyse im "average case"
Gestern war: "worst case" (n^2) und "best case" ($n \cdot \log n$)
 - Andere Sortierverfahren: MergeSort, HeapSort
Kurz die Idee ... keine Details und keine Programme
 - Sortieren in Linearzeit: 0-1-Sort und CountingSort
Das geht aber nur für bestimmte Eingaben
 - Untere Schranke: vergleichsbasiert nicht besser als $n \cdot \log n$
 - **Übungsblatt 2:** Rechenaufgabe zum Logarithmus; Rekursion mit anderer Aufteilung; CountingSort Permutation als Ausgabe

QuickSort im "Average Case" 1/3

- im schlechtesten Fall*
 $\ell_2 = n$
- Anzahl Operationen im allgemeinen Fall



- Sei k die maximale Rekursionstiefe und sei die Größe der Teilfelder auf Stufe j : $n_{j,1}, n_{j,2}, \dots, n_{j,\ell_j}$ (Summe $\leq n$)
- Dann gibt es eine Konstante $A > 0$ so dass gilt:

$$T(n) \leq T(n_{1,1}) + T(n_{1,2}) + A \cdot n$$

$$T(n_{1,1}) + T(n_{1,2}) \leq T(n_{2,1}) + \dots + T(n_{2,\ell_1}) + A \cdot n$$

$$T(n_{2,1}) + \dots + T(n_{2,\ell_2}) \leq T(n_{3,1}) + \dots + T(n_{3,\ell_3}) + A \cdot n \quad \text{usw.}$$

- Dann kann man leicht zeigen (durch Induktion)

$$T(n) \leq T(n_{k,1}) + \dots + T(n_{k,\ell_k}) + k \cdot A \cdot n \leq (k + 1) \cdot A \cdot n$$

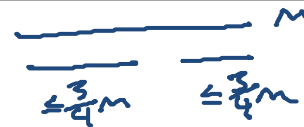
$= 1$ *$= 1$*

Intuitiv: auf jeder Rekursionsstufe linear viel Arbeit, und insgesamt höchstens k Stufen

QuickSort im "Average Case" 2/3

■ Gute Pivots

n \dots $m/4$ $|$ $m/2$ $|$ $m/4$ \dots
(sonst, nur zum Verständnis)
die wären alle gut



$(\frac{3}{4})^2 = \frac{9}{16} \approx \frac{1}{2}$
also ung. doppelt so "klein"

- Sagen wir, ein Pivot sei gut, wenn die resultierenden Teilfelder beide $\leq c \cdot n$ groß sind, für $c = \frac{3}{4}$ best case: $c = \frac{1}{2}$

Dann ist ein zufälliger Pivot gut mit Wahrscheinlichkeit $\frac{1}{2}$

- Betrachten wir eine Sequenz von ineinander geschachtelten rekursiven Aufrufen: dann ist man nach spätestens $\log_{1/c} n$ guten Pivots bei einem Feld der Größe 1
- Daraus folgt: die maximale Rekursionstiefe ist mit hoher Wahrscheinlichkeit $\leq C \cdot \log_{1/c} n$ für eine Konstante $C > 0$

$c = \frac{1}{2} : \log_2 n, c = \frac{3}{4} : \log_{4/3} n$

Das lässt sich mit elementarer Wahrscheinlichkeitsrechnung zeigen, aber lassen wir in dieser (Grund-)Vorlesung weg

QuickSort im "Average Case" 3/3

■ Fazit

- Auch im "average case" gibt es also eine Konstante $C > 0$, so dass gilt:

$$T(n) \leq C \cdot n \cdot (1 + \log_2 n)$$

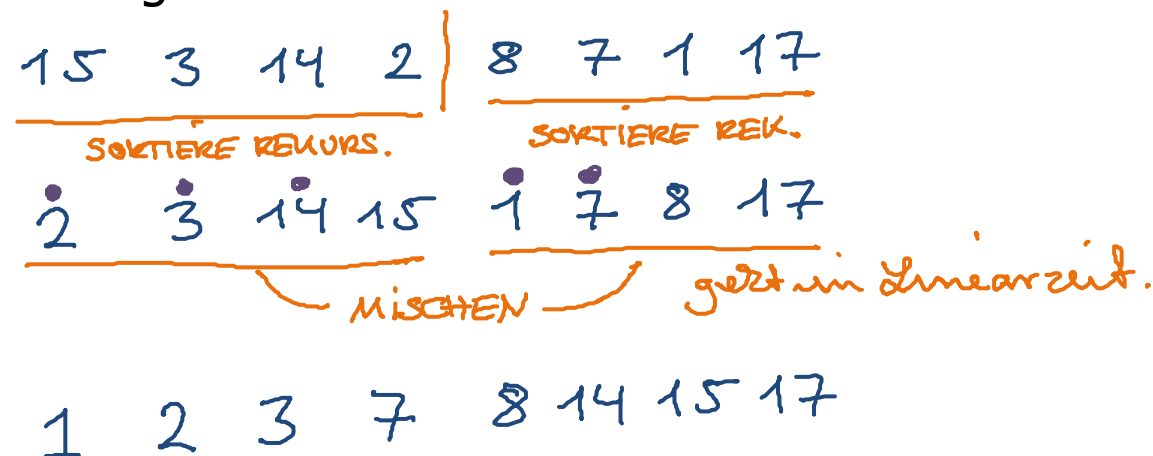
Quadratische Laufzeit also nur bei ungeschickt gewählten Pivots, und auch dann nur bei bestimmten Eingaben

Mit dem "Median-of-Three" Pivot bekommt man übrigens auch eine Laufzeit proportional zu $n \cdot \log n$

*nur das C ist etwas größer
als beim "best case"*

■ MergeSort

- Auch MergeSort teilt das Feld rekursiv auf, aber immer in zwei möglichst gleich große Hälften ... $\lfloor n/2 \rfloor$ und $\lceil n/2 \rceil$
- Es gibt aber keinen "divide" Schritt ... dafür müssen die beiden (rekursiv) sortieren Hälften nach den rekursiven Aufrufen noch "gemischt" werden



■ MergeSort, Fortsetzung

- Die Laufzeit ist dann **in jedem Fall** wie bei QuickSort im "best case" und im "average case", nämlich

$$T(n) \leq C \cdot n \cdot (1 + \log_2 n) \text{ für eine Konstante } C > 0$$

- Trotzdem wird man in der Praxis QuickSort bevorzugen:

Bei geeigneter Wahl des Pivots hat auch QuickSort praktisch immer $T(n) \leq C \cdot n \cdot (1 + \log_2 n)$

QuickSort macht alles "in place" in dem Eingabefeld

Für das "Mischen" bei MergeSort braucht man dagegen zusätzlichen Speicher ... das kostet Zeit, und führt zu einem höheren Faktor C in der Laufzeit

■ HeapSort

- HeapSort benutzt einen **binären Heap** zum Sortieren

Das ist eine Datenstruktur, die aus einer Menge von n Elementen mit $\leq C \cdot \log n$ Operationen das kleinste Element extrahieren und entfernen kann

Dazu in einer späteren Vorlesung mehr !

- HeapSort schafft damit auch **in jedem Fall**

$T(n) \leq C \cdot n \cdot (1 + \log_2 n)$ für eine Konstante $C > 0$

- In der Praxis:

HeapSort etwas besser als MergeSort ... **kleineres C**

HeapSort etwas schlechter als QuickSort ... **größeres C**

Sortieren in Linearzeit 1/3

■ ZeroOneSort

- Geht Sortieren auch schneller als $n \cdot \log n$?
- Ja, zum Beispiel wenn alle Elemente nur 0 oder 1 sind
- Dann kann man einfach die 0en und 1en zählen

Dazu schreiben wir gerade ein Programm ZeroOneSort

0, 1, 1, 0, 0, 1, 1

#0 = 3, #1 = 4

(trivial in Linearzeit
z.B. einfach Summe
berechnen)

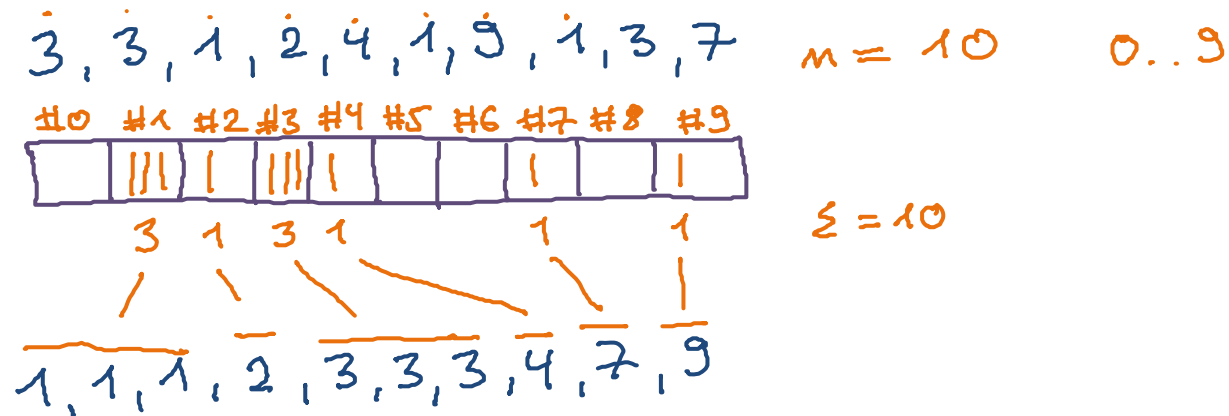
0 0 0 1 1 1 1
└───┘ └───┘
3 mal 4 mal

Sortieren in Linearzeit 2/3

■ CountingSort

- Die Idee klappt auch noch, wenn die Elemente aus dem Bereich $0 .. n - 1$ sind

Dazu schreiben wir gerade ein Programm CountingSort



■ Laufzeit

- Für beide Verfahren gilt offensichtlich

$$T(n) \leq C \cdot n \dots \text{für ein } C > 0$$

Man geht einmal über das Eingabefeld zum "Zählen",
und dann gibt man die (gleich große) Ausgabe aus

- Ist das vielleicht sogar für beliebige Eingaben möglich?

CountingSort braucht ein Feld der Größe m für Zahlen
aus dem Bereich $0..m-1$

Für $m \gg n$ ist die Laufzeit (und der Platzverbrauch)
dann proportional zu m , und nicht zu n

- Vergleichsbasiertes Sortieren
 - ZeroOneSort und CountingSort sortieren die Elemente nicht durch "Umsortieren", sondern durch Zählen
 - Wir wollen jetzt zeigen: wenn man "nur Umsortieren" zulässt, geht es tatsächlich nicht schneller als $n \cdot \log n$
 - Dazu müssen wir erst mal genauer fassen, was es heißt, dass ein Algorithmus "nur umsortiert"

■ Vorbetrachtung 1

- Wir werden uns bei unserem Beweis auf Algorithmen **von einer bestimmten Art** beschränken

Wir werden sehen: weil das die Argumentation erleichtert

- Nehmen wir an, ein Algorithmus **A** ist nicht von dieser Art, aber es gibt einen Algorithmus **A'** von der Art für den gilt:

Er läuft genauso schnell wie **A**, bis auf $\leq C_1 \cdot n$ Operationen

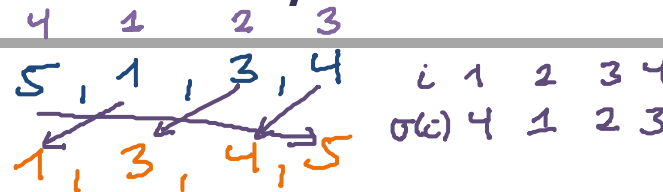
Die Ausgabe von **A** kann mit $\leq C_2 \cdot n$ Operationen in die Ausgabe von **A'** überführt werden

- Dann gilt $T_{A'}(n) \geq n \cdot \log n \Rightarrow T_A(n) \geq n \cdot \log n$

Wäre **A** schneller, könnten wir auch **A'** schneller machen

Untere Schranke Sortieren 3/12

■ Vorbetrachtung 2



- Bisher haben unsere Algorithmen für Eingabe x_1, \dots, x_n diese Zahlen in sortierter Reihenfolge ausgegeben
- Wir wollen im Folgenden Algorithmen betrachten, die stattdessen ausgeben: eine Permutation σ der Zahlen $1, \dots, n$ so dass $x_{\sigma(1)} \leq x_{\sigma(2)} \leq \dots \leq x_{\sigma(n)}$
- **Beobachtung:** alle unseren bisherigen Algorithmen können dahingehend abgewandelt werden, ohne dass sich die Anzahl Operationen um mehr als $C \cdot n$ ändert

Ü2: das gilt auch für ZeroOneSort und CountingSort

■ Vorbetrachtung 3

- In einem Programm (Python, Java, C++) können an diversen Stellen Verzweigungen auftreten

```
while ( ... ) { ... }
```

```
for ( ... ) { ... }
```

```
if (...) { ... } else { ... }
```

- Ohne Beschränken der Allgemeinheit seien all diese Verzweigungen von der Form `if (...) { ... } else { ... }`

Zum Beispiel ist `while (EXPR) { ... }` äquivalent zu:
`while (true) { if (EXPR) { ... } else { break; } }`

■ Vorbetrachtung 4

- Betrachten wir die Folge von Entscheidungen in den `if (...) { ... } else { ... }` Teilen im Ablauf eines Programms
- Dann entspricht jeder Ablauf einer Folge **I E E I I I E I I I E ...**
 - = `if`-Teil wird ausgeführt, **E** = `else`-Teil wird ausgeführt
- Ein Algorithmus heißt **vergleichsbasiert**, wenn die I/E Folge die Ergebnispermutation **eindeutig** bestimmt
- `MinSort`, `QuickSort`, `MergeSort`, `HeapSort` sind alle vergleichsbasiert bzw. können leicht dazu gemacht werden
- `ZeroOneSort` und `CountingSort` sind es nicht, und können auch nicht dahingehend abgeändert werden

- ZeroOneSort ist nicht "vergleichsbasiert"

- Hier sind zwei Eingaben mit gleicher I/E Folge aber verschiedenen Ergebnispermutationen

Wir lassen dabei die I und E von den "for" Schleifen weg, die hängen ja sowieso nicht von der Eingabe ab

σ_1 1 4 5 2 3 : IIII EE
0, 1, 1, 0, 0

σ_2 4 1 2 3 5 : IIII EE
1, 0, 0, 0, 1

Untere Schranke Sortieren 7/12

■ MinSort ist "vergleichsbasiert"

- Beispiel: zwei Eingaben mit gleicher I/E Folge und gleicher Permutation, und dritte Eingabe mit anderer I/E Folge und anderer Permutation

Wieder ohne die Is und Es von den "for" Schleifen

σ_1	² 5, ³ 8, ¹ 3, ⁵ 17, ⁴ 12	:	<u>E I E E</u> ...
			Iteration 1
σ_2	² 4, ³ 12, ¹ 2, ⁵ 35, ⁴ 17	:	<u>E I E E</u> ...
H			
σ_3	⁴ 4, ³ 3, ² 2, ⁵ 35, ¹ 1	:	<u>I I E I</u> ...

■ Beweis untere Schranke, Teil 1

– Wir betrachten jetzt einen beliebigen vergleichsbasierten Algorithmus **A**, der für eine Eingabe der Größe **n** eine sortierende Permutation σ der Zahlen **1, ..., n** ausgibt

– Sei $T(n)$ eine obere Schranke für die Anzahl der von **A** benötigten Operationen auf einer Eingabe der Größe **n**

Dann höchstens $T(n)$ Verzweigungs-Anweisungen

– Der Algorithmus gibt also für Eingabegröße **n** höchstens $2^{T(n)}$ verschiedene Permutationen aus

IIIEIEE...
x viele
 $\leq T(n)$

2^x Möglichkeiten
 $\leq 2^{T(n)}$

■ Beweis untere Schranke, Teil 2

- Ein korrekter Algorithmus muss für Eingabegröße n alle möglichen Permutationen erzeugen können, das sind $n!$

Wenn er eine Permutation nicht erzeugen könnte, würde er für die Eingabe, die genau diese Permutation zum Sortieren benötigt, nicht das richtige Ergebnis liefern

- Auf der vorherigen Folie hatten wir gesehen, dass bei Laufzeit $\leq T(n)$ höchstens $2^{T(n)}$ Permutationen erzeugt werden können
- Wäre $2^{T(n)} < n!$, würden nicht alle Permutationen erzeugt werden können und der Algorithmus wäre nicht korrekt
- Es muss also $2^{T(n)} \geq n!$ sein, oder äquivalent $T(n) \geq \log_2(n!)$

Untere Schranke Sortieren 10/12

■ Abschätzung von $\log_2(n!)$

$$(*) \log_2 x = \frac{\ln x}{\ln 2}$$

– Dafür wird oft die Stirling-Formel benutzt:

$$n! \geq \sqrt{2 \cdot \pi \cdot n} \cdot (n/e)^n$$

$$(**) \text{ Stammfunktion von } \ln x : x \cdot \ln x - x$$

– Wir können das aber auch (etwas weniger genau, aber ausreichend) elementar-mathematisch abschätzen:

$$\log_2(n!) = \log_2(n \cdot (n-1) \cdot \dots \cdot 1) \quad \underbrace{\hspace{10em}}_{=0}$$

$$= \log_2 n + \log_2(n-1) + \dots + \log_2 1$$

$$(*) = \frac{1}{\ln 2} \sum_{i=1}^n \ln i$$

$$\geq \frac{1}{\ln 2} \int_1^n \ln x \, dx$$

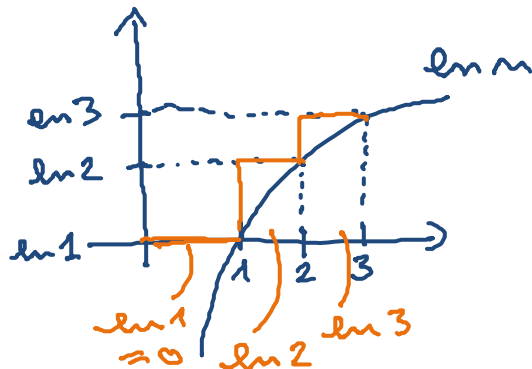
$$\geq \frac{1}{\ln 2} \left. x \cdot \ln x - x \right|_1^n$$

$$\underbrace{n \cdot \ln n - n + 1}_{\geq n \cdot (\ln n - 1)} \geq n \cdot \frac{\ln n}{2}$$

$\ln n \geq 2 \Leftrightarrow n \geq e^2$

Fazit:

$$\log(n!) \sim n \cdot \log n$$



- Beweis untere Schranke, Zusammenfassung

- Wir haben gezeigt:

- Sei $T(n)$ eine obere Schranke für einen Algorithmus, der n Elemente vergleichsbasiert sortiert

- Dann ist $T(n) \geq \log_2(n!) \geq \frac{1}{2} \cdot n \cdot \log_2 n$ für $n \geq 9$

■ Fazit

– Unter den vergleichsbasierten Algorithmen sind also QuickSort, MergeSort, HeapSort alle optimal !

– Aber in der Praxis zählen auch (unter anderem)

Konstante Faktoren ... zum Beispiel: " $n \log n$ " ist 10 mal schneller als " $10 n \log n$ "

Cache-Effizienz ... Zugriff auf aufeinanderfolgende Elemente im Speicher ist billiger als wenn "verstreut"

– Zu diesen Aspekten in einer späteren Vorlesung mehr !

- Weiterführende Literatur (bei Interesse)

- Untere Schranke für vergleichsbasiertes Sortieren

Mehlhorn/Sanders: [5.3 A Lower Bound \[for Sorting\]](#)

Wikipedia: [wiki/QuickSort#Formal_analysis](https://de.wikipedia.org/wiki/QuickSort#Formal_analysis)