

Exkurs

Eine Algorithm-Engineering-Fallstudie

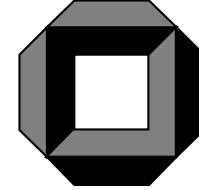
Volltextsuche mit invertiertem Index

Motivation:

Volltextsuchmaschinen wie Google haben die Welt verändert.

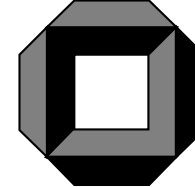
und machen das WWW erst nützlich.

Dazu sind Algorithmen unerlässlich



Algorithmen die Google Ausmachen

- Invertierter Index**. Zum Beispiel **Informatik Karlsruhe**
- Phrasensuche, z.B. **“to be or not to be”**
- Ranking. PageRank braucht im wesentlichen wiederholte Multiplikation dünn besetzter Matrizen mit Vektoren (iterative Eigenwertberechnung).
- Versteigerungsmechanismen für Werbung.



Konjunktive Volltextsuche

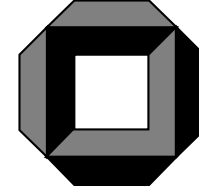
Gegeben:

Dokumente d_1, \dots, d_U

Schlüsselwörter t_1, \dots, t_k aus Wörterbuch mit W Einträgen.

Gesucht:

IDs ($\in 1..U$) der Dokumente, die t_1, \dots, t_k

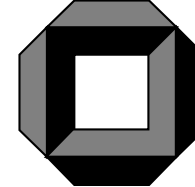


Invertierter Index

Für jedes Wort t berechne eine sortierte Liste $L(t)$ mit den docIDs der Dokumente, die t enthalten. (Entwurfstrategie Vorberechnung)

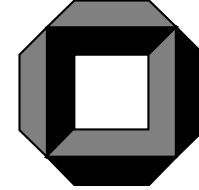
Ergebnis ist $L(t_1) \cap \dots \cap L(t_k)$

\rightsquigarrow Reduktion auf Mengenschnitt sortierter Folgen.



Indexkonstruktion

```
foreach Document  $d_i$  do  
    foreach term  $t \in d$  do  
         $L(t).pushBack(i)$ 
```



Mengenschnitt sortierter Folgen

Gut, wenn Folgenlängen sehr unterschiedlich:

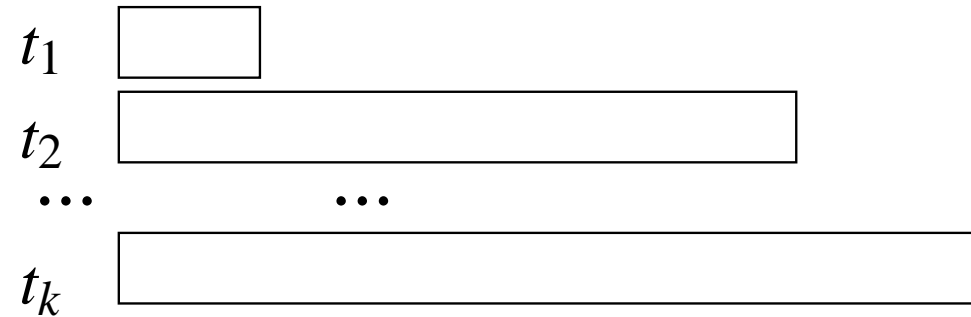
Sortiere Suchterme so um, dass

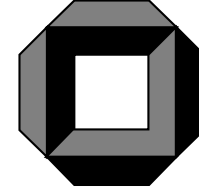
$$|L(t_1)| \leq \dots \leq |L(t_k)|$$

berechne

$$(\dots((L(t_1) \cap L(t_2)) \cap L(t_3)) \dots) \cap L(t_k)$$

↪ Reduktion auf Mengenschnitt **zweier** sortierter Folgen.

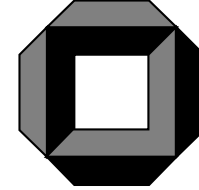




Mengenschnitt zweier sortierter Folgen

Function zipper($\langle \alpha_1, \beta_1, \dots, \omega_1 \rangle, \langle \alpha_2, \beta_2, \dots, \omega_2 \rangle$)
 if $\alpha_1 < \alpha_2$ **then return** zipper($\langle \beta_1, \dots, \omega_1 \rangle, \langle \alpha_2, \beta_2, \dots, \omega_2 \rangle$)
 if $\alpha_1 > \alpha_2$ **then return** zipper($\langle \alpha_1, \beta_1, \dots, \omega_1 \rangle, \langle \beta_2, \dots, \omega_2 \rangle$)
 else return $\langle \alpha_1 \rangle \circ$ zipper($\langle \beta_1, \dots, \omega_1 \rangle, \langle \beta_2, \dots, \omega_2 \rangle$)

Hochoptimierte, nichtrekursive, **feld**basierte Variante von zipper
funktioniert sehr gut, solange Längenverhältnis < 10



Mengenschnitt zweier ungleich langer Folgen

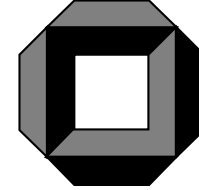
O.B.d.A. sei $m = |L_1| \leq |L_2| = n$ Zipper braucht $O(m + n)$ Zeit.

Vergleichsbasierte **untere Schranke** (Skizze):

Es gibt $\binom{n}{m}$ verschiedene Möglichkeiten, m Elemente in n Elementen zu positionieren.

\Rightarrow wir brauchen $\log \binom{n}{m}$ Bits Information (durch **Vergleiche**).

$$\log \binom{n}{m} \geq \log \left(\frac{n}{m} \right)^m = m \log \frac{n}{m}$$



Mengenschnitt zweier ungleich langer Folgen

Obere Schranke

Function exponentialSearchIntersect($\langle \alpha_1, \beta_1, \dots, \omega_1 \rangle, L_2$)

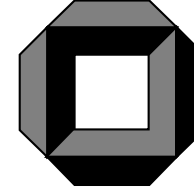
suche α_1 in L_2 , d.h., A_2 B_2
spalte $L_2 = A_2 \circ B_2$ mit

$\leq \alpha_1$	$> \alpha_1$
-----------------	--------------

$C = \text{exponentialSearchIntersect}(\langle \beta_1, \dots, \omega_1 \rangle, B_2)$

if $\alpha_1 = \text{last}(A_2)$ **then** $C := \langle \alpha_1 \rangle \circ C$

return C



Exponentielle Suche

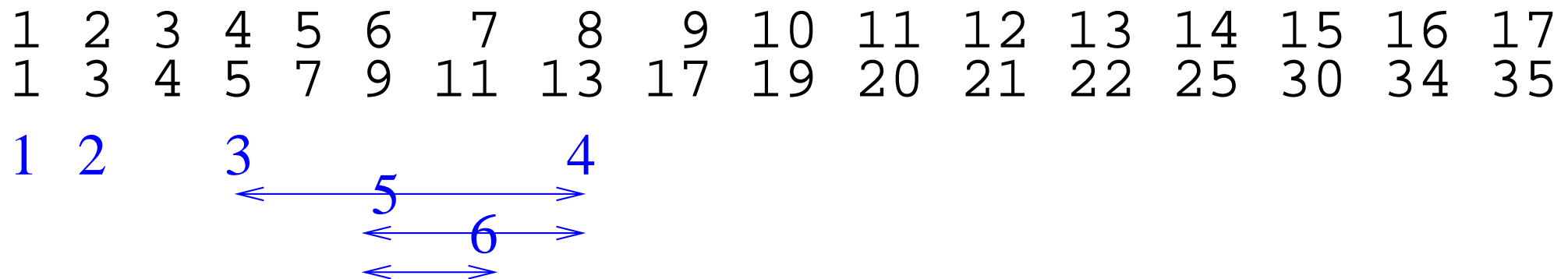
Function exponentialSearch($e, a : \mathbf{Array}$)

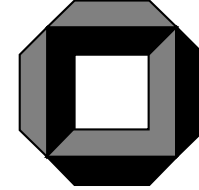
$i := 1$

while $e \leq a[i]$ **do** $i := 2i$

now binary search in $a[i/2..i] \rightsquigarrow$ **return** found pos $j \dots$

Zeitaufwand $O(1 + \log k)$





Mengenschnitt zweier ungleich langer Folgen

Analyse

exponentialSearchIntersect teilt L_2 in Teillisten der Längen $\Delta_1, \dots, \Delta_m$.

$$\text{Aufwand } \sum_{i=1}^m O(1 + \log \Delta_i) = O\left(m\left(1 + \log \frac{n}{m}\right)\right)$$

(Konvexität der Logarithmusfunktion)...

...

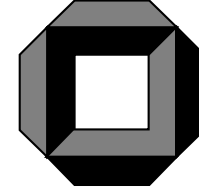
Leider sind diese Verfahren in der Praxis ziemlich langsam.

Idee: weitergehende Vorverarbeitung der Listen.

Idee: breche untere Schranke. Nutze aus, dass die Elemente ganze

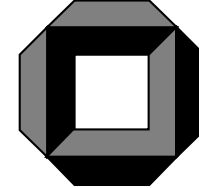
Zahlen sind

↪ **mittlere** Ausführungszeit $O(m)$.



Verfeinerungen und Verallgemeinerungen

- Parallele Konstruktion. Zum Beispiel mittels MCSTL?
`http://algo2.iti.uni-karlsruhe.de/singler/mcstl/`
- Phrasensuche mittels **Positional Index**
- Nur Paare bestimmter Suchwörter ermöglichen vernünftige Sucheingrenzung? Vorberechnung von Listen auch für solche Wortpaare.
- Datenkompression



Zusammenfassung

Algorithmenentwurfstechnik **Vorverarbeitung**: Hier Reduktion auf Manipulation von Begriffs- und Dokument-Nummern.

Algorithmenentwurfstechnik **Vorbereitung**: Hier Inverted Index

Algorithmenentwurfstechnik **Basic Toolbox**: Hier unbounded array, mischen, Mengenschnitt, exponentielle Suche, binäre Suche, Wörterbücher, sortieren, Datenkompression,...