

Wie kann das LCP-Array als Maß für die Komprimierbarkeit genutzt werden?

Überlegungen:

- Welche Texte sind gut komprimierbar?

↳ All-a-Text: $a^na \xrightarrow{LZ77} (0, a)(1, n-1)$

↳ Texte mit vielen Wdh.: $\underbrace{abc \dots xyz}_a \alpha a d e \alpha x y \alpha$

- Welche Texte lassen sich nicht gut komprimieren?

↳ Texte mit wenigen Wdh. oder nur sehr kurzen Wdh.

$abc \underbrace{ab} d \underbrace{abc} e \underbrace{ab} f \dots$ Abspeichern der Faktoren kann teuer werden

→ Was bedeutet das für das LCP-Array

↳ All-a-text: LCP: $n-1 \mid n-2 \mid n-3 \mid \dots \mid 0$

↳ Texten mit vielen Wdh.: $1 \mid 0 \mid 2 \mid 0 \mid \alpha \mid \alpha \mid \alpha \mid \alpha + 4 \mid \dots$

↳ Wenige Wdh. $1 \mid 0 \mid 2 \mid 2 \mid 2 \mid \dots$

Also etwas konkreter

↳ "Je höher die Summe der LCP-Werte ist, desto komprimierbarer sollte ein Text sein" (evtl + Normierung durch die Länge des Textes)

↳ Was ist z.B. mit dem Maximum

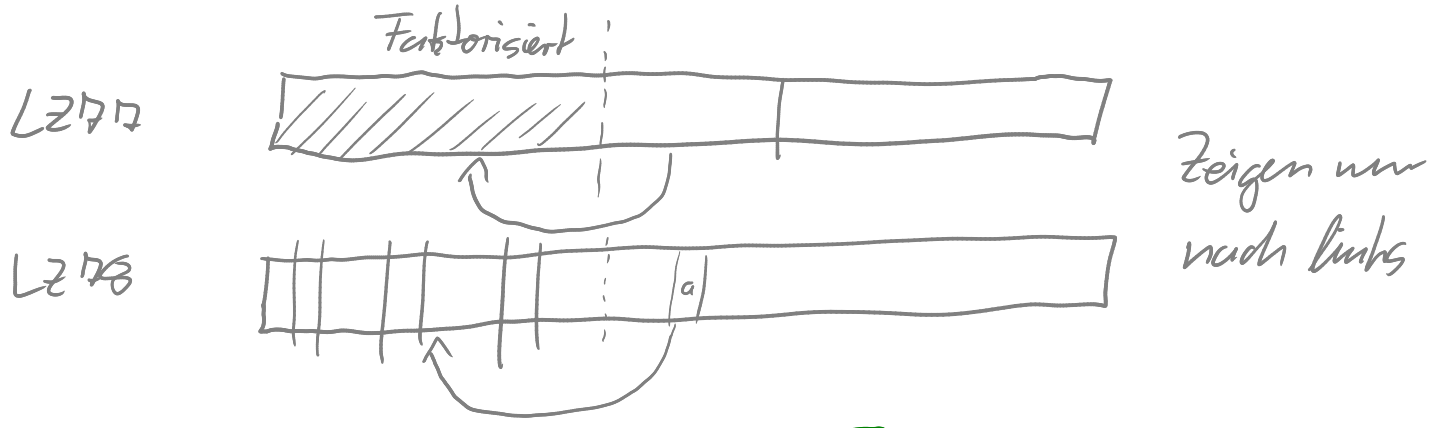
↳ $\underbrace{a \dots a}_{n/2} \underbrace{b \dots b}_{n/2}$

Das LCP-Array zum Komprimieren nutzen

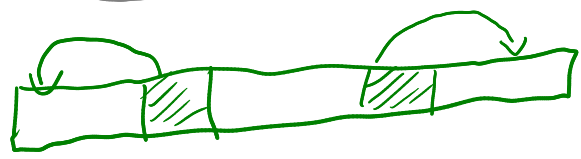
Basierend auf „Bidirectional Text Compression in External Memory“, Dinklage et al. ESA '19.

EM: Ist uns egal (für diese Übung)

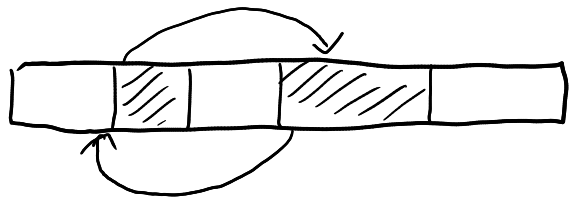
Bidirectional: Faktoren dürfen in beide Richtungen zeigen



Bidirectional



Achtung: Kreise!



Die Idee von PLCP-Comp

$$PLCP[SA[i]] = LCP[i]$$

↳ „LCP-Array in Textreihenfolge“

↳ Das PLCP-Array wird beim „engineered linear time“-Algorithmus konstruiert



① Suche maximalen Wert im PLCP-Array
 ↳ Konstruiere aus diesem Wert einen Faktor
 ↳ Faktorisier ab Textpos. i : $(\phi[i], \text{PLCP}[i])$

Mit wem möchte ich vergleichen werden
 $\phi[i] = \text{SA}[\text{ISA}[i]-1]$

Hier kommen Faktoren in beide Richtungen zeigen
 Mit welchem Suffix kommt der LCP-Wert zustande

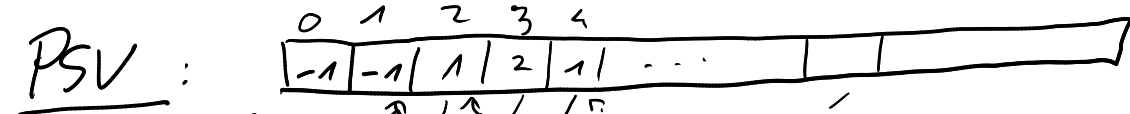
② PLCP-Array aktualisieren
 a) Alles was überdeckt ist löschen
 b) Alles was ein Präfix von dem Faktor ist (*)
 muss reduziert werden (*)

③ Weiter bei ①, so lange sinnvolle Faktoren existieren

Die Implementierungen von LZ77 / LZ78

LZ77

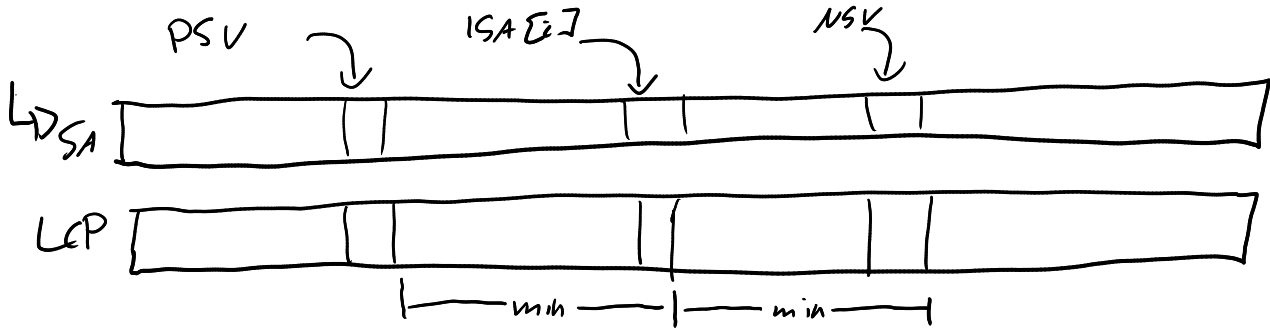
Faktoren: Länge, Text Position (wenn Länge = 1, dann steht hier das Zeichen)



(NSV analog)

Berechnung der Faktoren

↳ i ist Textpos. ab der der nächste Faktor beg.



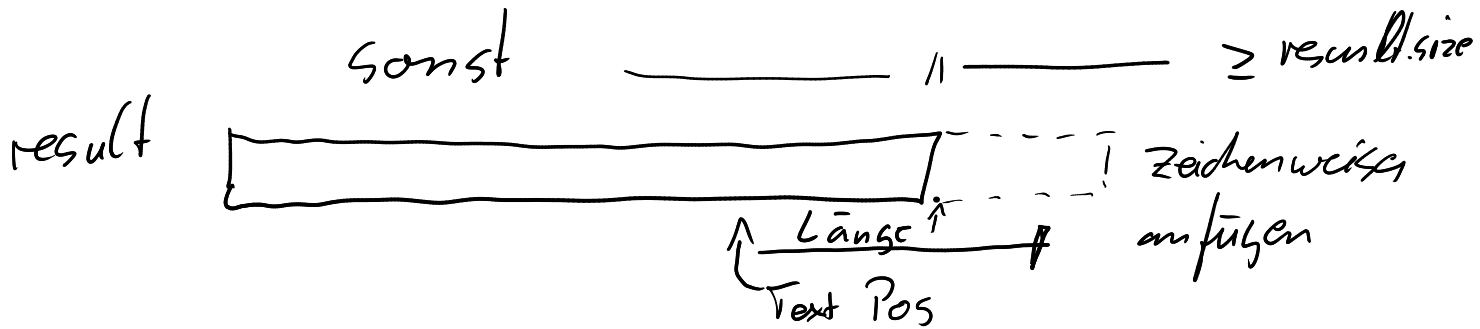
Struct LZ77 mit Konstruktor erlaubt Verwendung von emplace_back

↳ konstruiert an der richtigen Stelle im Speicher

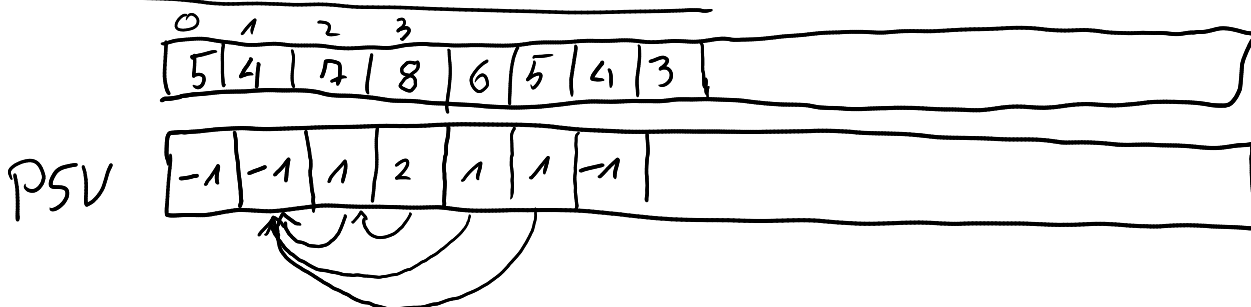
↳ kann eine Kopie sparen

Beim Dekomprimieren

↳ Länge > 1 und $factor \cdot text_pos + factor \cdot length < result.size$



PSV in Linearzeit

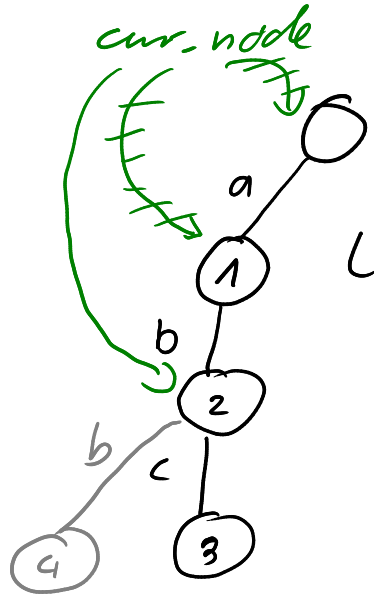


L778

(als Zahl mit 1 beg.)

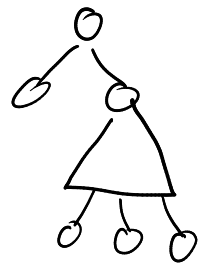
Faktoren: Referenz auf anderen Faktor, Zeichen

Ablauf:



a|ab|c|abc|abb

L778: (0, a)(1, b)(2, c)(2, b)



Dehkomprimieren

↳ Erzeuge Dictionary für alle Faktoren

Eingabe: Faktor - Nummer

Ausgabe: "String, den der Faktor kodiert"

↳ Pfad von der Wurzel bis zum Knoten des Faktors

In der Implementierung geben wir Textpos. + Länge zurück

Dict: $\{ \}$ (0, 1) (1, 2) (3, 3) (6, 3) ... () ()

result: 0 1 2 3 4 5 6 7 8
 a a b a b c a b b