

Parallelisierung einer speichereffizienten Approximation der LZ77-Faktorisierung

Gajann Sivarajah

LZ-Kompression - Konzept

Eingabe: $S = e_1 \dots e_n$

- $e_i \in \Sigma = \{0, \dots, 255\}$

Ausgabe: $F = (f_1, \dots, f_z)$

- $f_1 \dots f_z = S$
- $f_i = \begin{cases} (Länge, Position) & , \text{ falls Referenz} \\ (0, Zeichen) & , \text{ sonst} \end{cases}$

Algorithmus: $COMP_{LZ} : S \rightarrow F \iff DECOMP_{LZ} : F \rightarrow S$

LZ-Kompression - Gütemaße

Qualität:

- $FR = \frac{z}{n} \iff CR = \frac{|F|_{Bin}}{|S|_{Bin}}$

Zeit:

- $T(n, p)$, hier $n := 200MB$ und $p = 16 \vee 128$

Speicher:

- $Mem_{Peak} :=$ Spitze der allokierten Speichers

LZ77

Konzept:

- Scanne von links nach rechts
- Maximiere jeden Faktor $|f_i| \rightarrow \textit{Greedy}$

Zeit / Speicher:

- Zeit: $O(n)$
- Speicher: $O(n)$

Approx. LZ77 - Konzept

Ablauf:

- Rundenbasierter Algorithmus
- Runde $r \Rightarrow$ Extrahiere Faktoren der Länge $\frac{|S|}{2^r}$
- Letzte Runde $r_{End} = \log |S| \Rightarrow$ Alle Zeichen sind faktorisiert

Approx. LZ77 - Konzept

Runde:

- (Noch unverarbeitete) Zeichenfolge in Blöcke aufteilen
- Unter den Blöcken Duplikate/Referenzen finden(*InitTables*)
- Freie Suche nach Referenzen in S (*ReferenceScan*)
- Extrahiere Faktoren aus Referenzen

Approx. LZ77 - Konzept

InitTables

- Erzeuge *RFPTable* und *RefTable*:
 - $RFPTable(RFP) =$ Linkester Block mit RFP als Hash
 - $RefTable(Block) = \begin{cases} \text{Position einer Referenz zu } Block & \text{,falls bekannt} \\ \text{Position von } Block & \text{, sonst} \end{cases}$
- Blöcke, die nicht in *RFPTable* eingetragen werden \Rightarrow **Faktoren**

Approx. LZ77 - Konzept

ReferenceScan

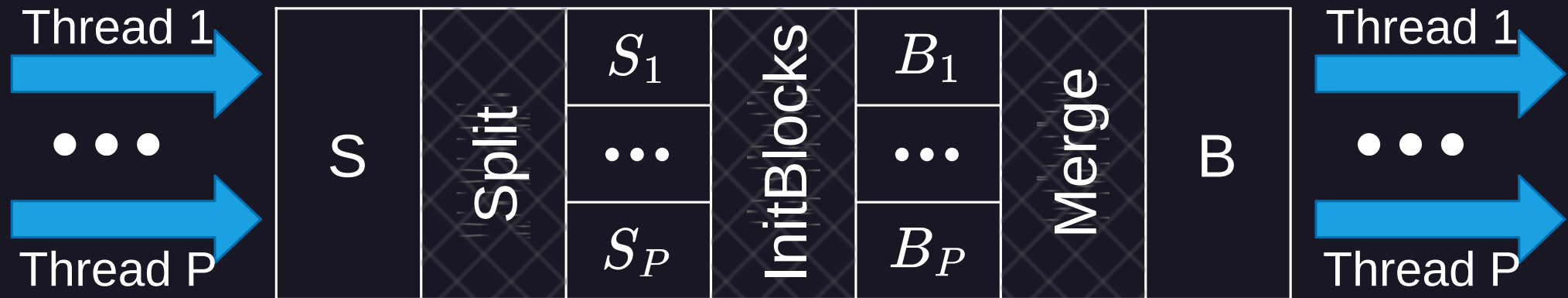
- Scan von links nach rechts \Rightarrow Bewege RFP-Fenster
- Treffer in RFPTable + Links von Eintrag in RefTable \Rightarrow **Faktor**

Approx. LZ77 - Güte

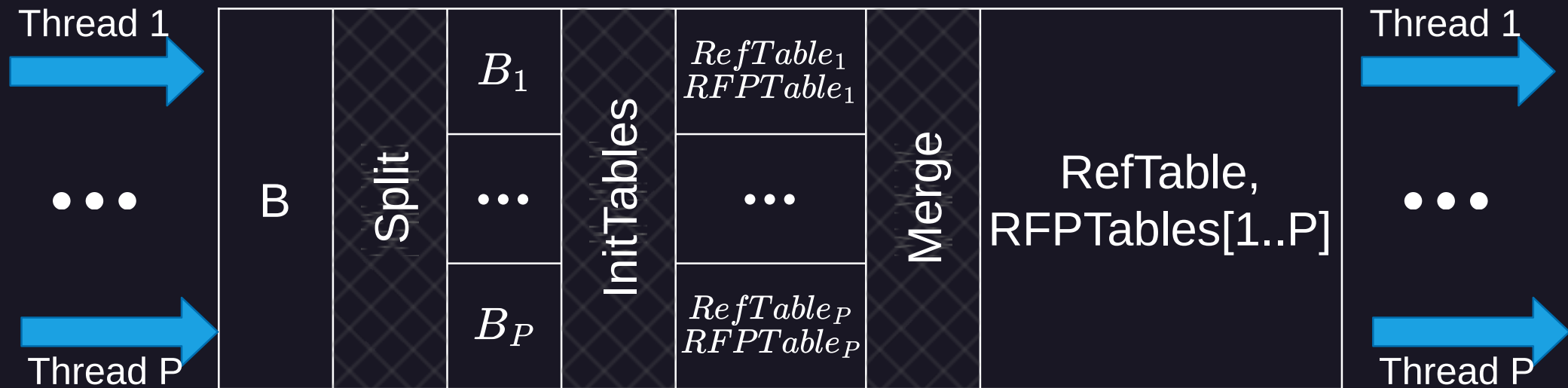
Zeit: $O(n \log n)$

Speicher: $O(z)$

Approx. LZ77Par - $S \Rightarrow$ Blöcke



Approx. LZ77Par - InitTables



Approx. LZ77Par - ReferenceScan



Optimierungen - DynStart

Optimierungen - DynEnd

Optimierungen - PreMatching

Optimierungen - ScanSkip

- $|F_{ReferenceScan}| \leq |RFPTable| = |Blocks| - |F_{InitTables}|$
- $k = \frac{|RFPTable|}{|Blocks|}$
- **Führe ReferenceScan nur bei $k \geq k_{min} \in [0, 1]$ durch**

Evaluation - Qualität

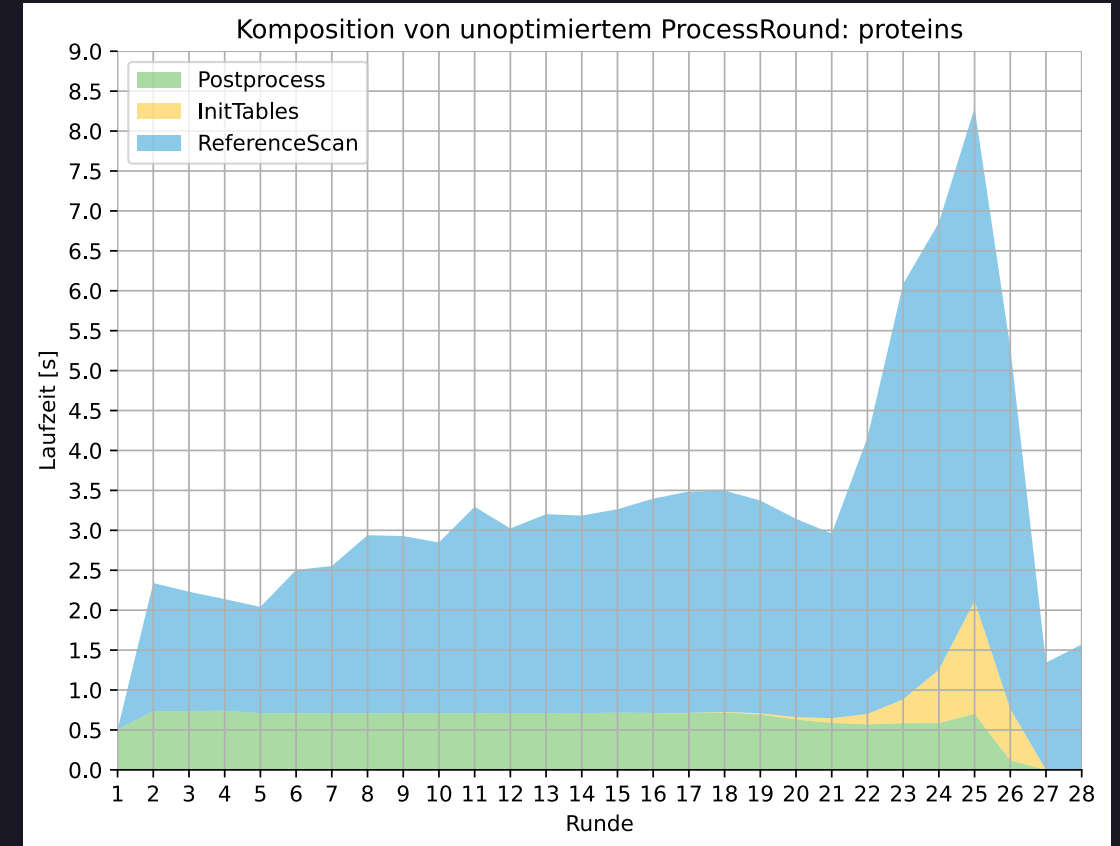
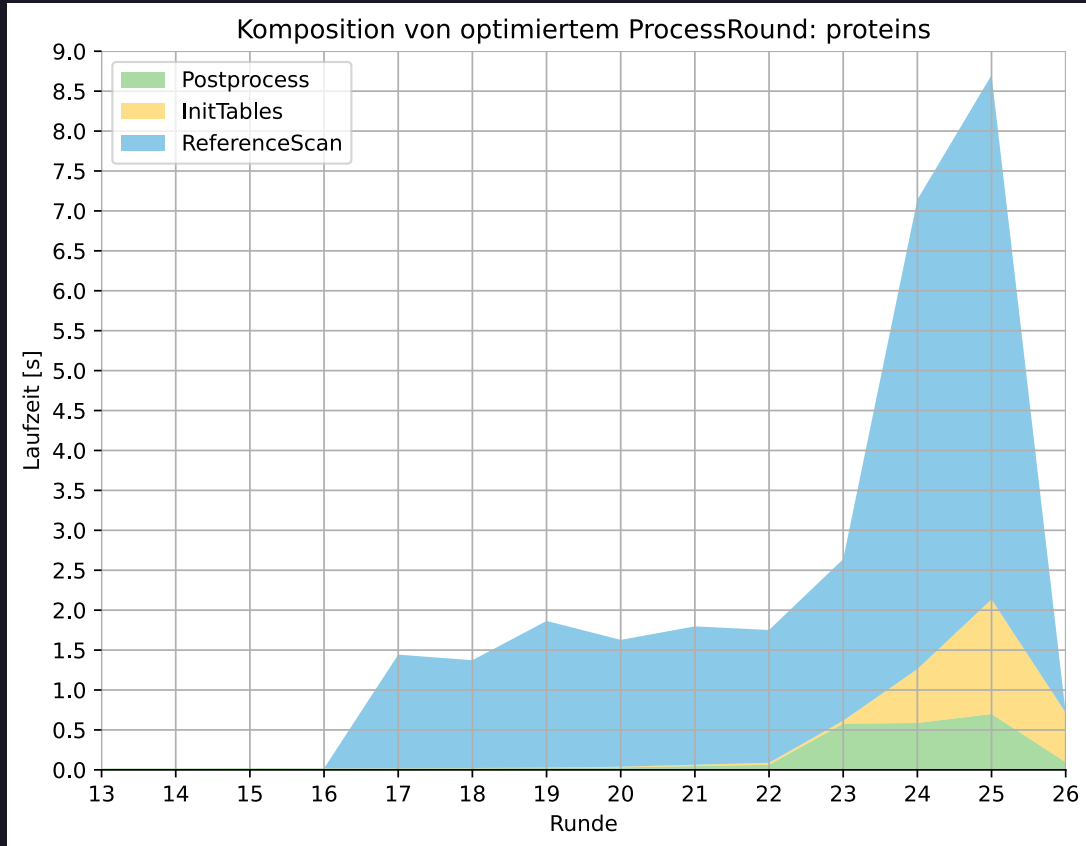
COMP	proteins	sources	dna	xml	english
LZ77	help	help	help	help	help
Approx. LZ77	help	help	help	help	help
Approx. LZ77Par	help	help	help	help	help

Evaluation - Speicher

COMP	proteins	sources	dna	xml	english
LZ77	help	help	help	help	help
Approx. LZ77	help	help	help	help	help
Approx. LZ77Par	help	help	help	help	help

Evaluation - Zeit

Evaluation - Optimierungen



Zusammenfassung

- Approx. LZ77 \rightarrow Approx. LZ77Par : Korrektheit nachgewiesen
- Zeitersparnis durch Optimierungen nachgewiesen
- $\text{Zeit}(\text{Approx. LZ77Par}) < \text{Zeit}(\text{LZ77}) < \text{Zeit}(\text{Approx. LZ77})$
- $\text{Speicher}(\text{Approx. LZ77Par}) \approx \text{Speicher}(\text{Approx. LZ77}) < \text{Speicher}(\text{LZ77})$

Offene Punkte

- Alternative Techniken (Hashtabelle, Bloom-Filter,...)
- Dynamische Generierung der Parameter $r_{PreMatch}$ und k_{min}
- Zweite und Dritte Phase des Approximationsalgorithmus