

Parallelisierung einer speichereffizienten Approximation der LZ77-Faktorisierung

Gajann Sivarajah

Eingabe: $S = e_1 \dots e_n$

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

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

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

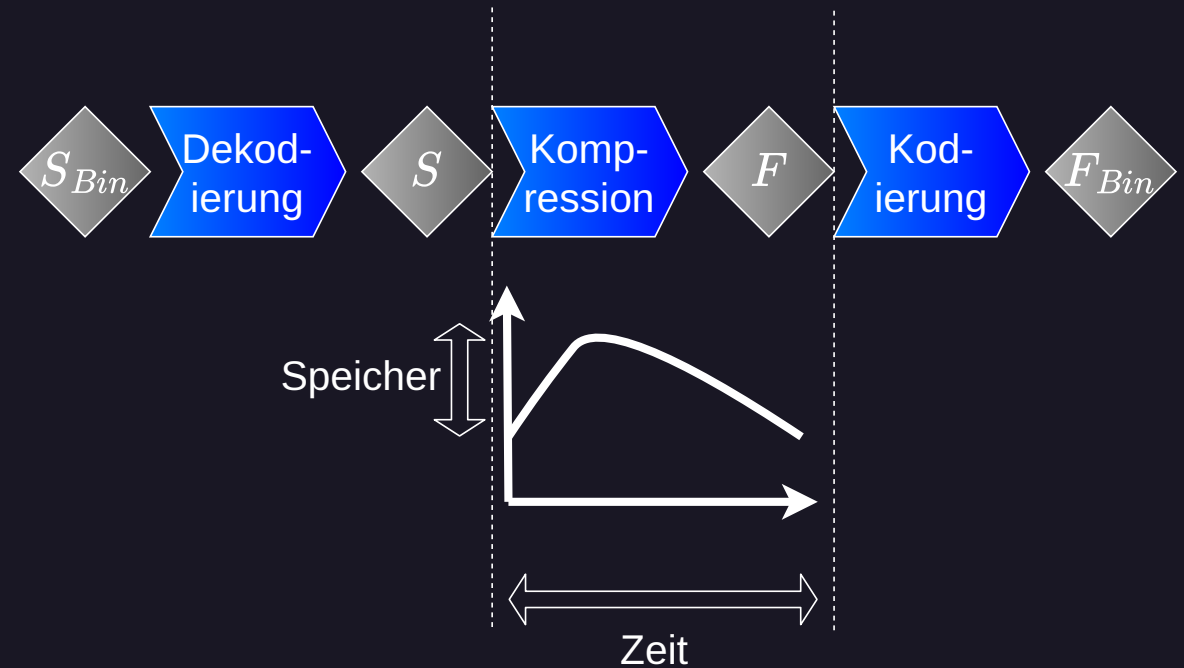
Algorithmus: $COMP_{LZ} : S \rightarrow F \iff DECOMP_{LZ} : F \rightarrow S$

Qualität:

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

Perfomanz:

- Speicher: Mem_{Peak}
- Zeit: $T(n, p)$



Konzept:

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

Zeit / Speicher:

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

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

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

InitTables

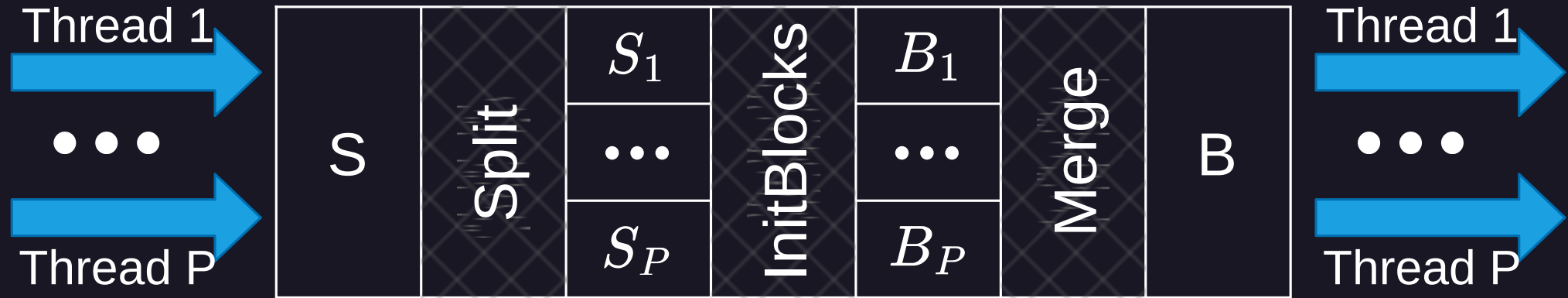
- Erzeuge *RFPTable* und *RefTable*:
 - $RFPTable(RFP) =$ Linkster Block mit RFP als Hash
 - $RefTable(Block) = \begin{cases} \text{Position einer Referenz zu } Block & \text{, falls bekannt} \\ \text{Position von } Block & \text{, sonst} \end{cases}$
- Blöcke, die nicht in *RFPTable* eingetragen werden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

ReferenceScan

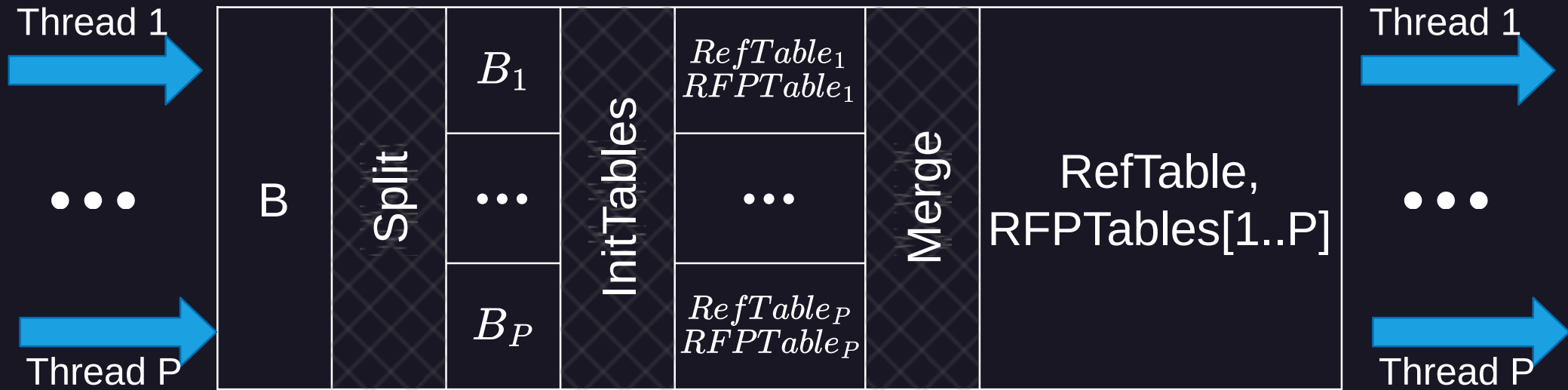
- Scan von links nach rechts \Rightarrow Bewege RFP-Fenster
- Treffer in RFPTable + Links von Eintrag in RefTable \Rightarrow **Faktor**

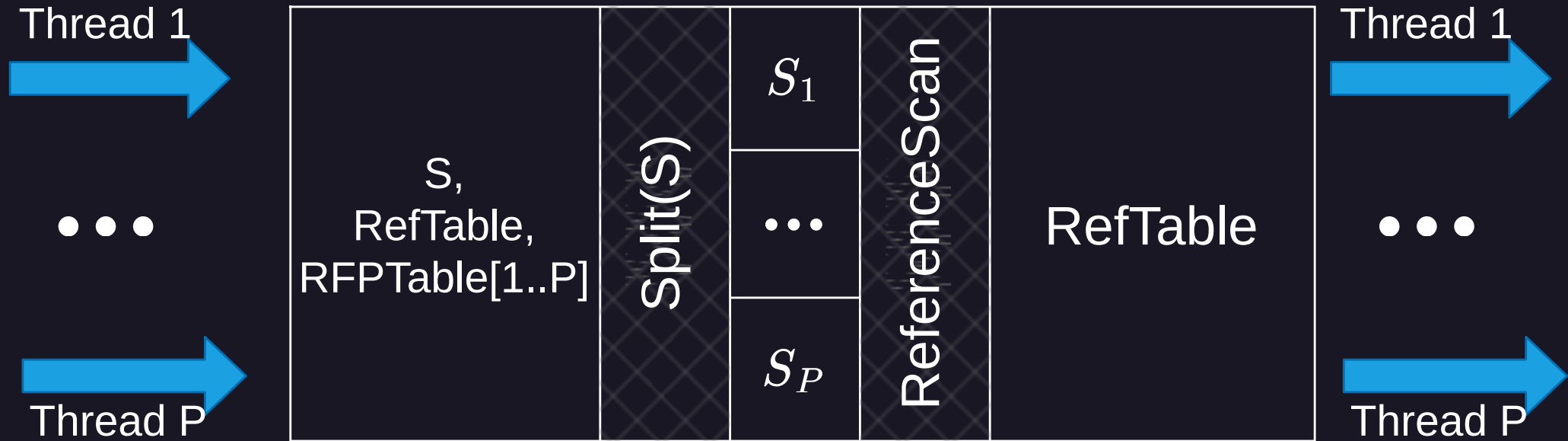
Zeit: $O(n \log n)$

Speicher: $O(z)$



Approx. LZ77Par - InitTables











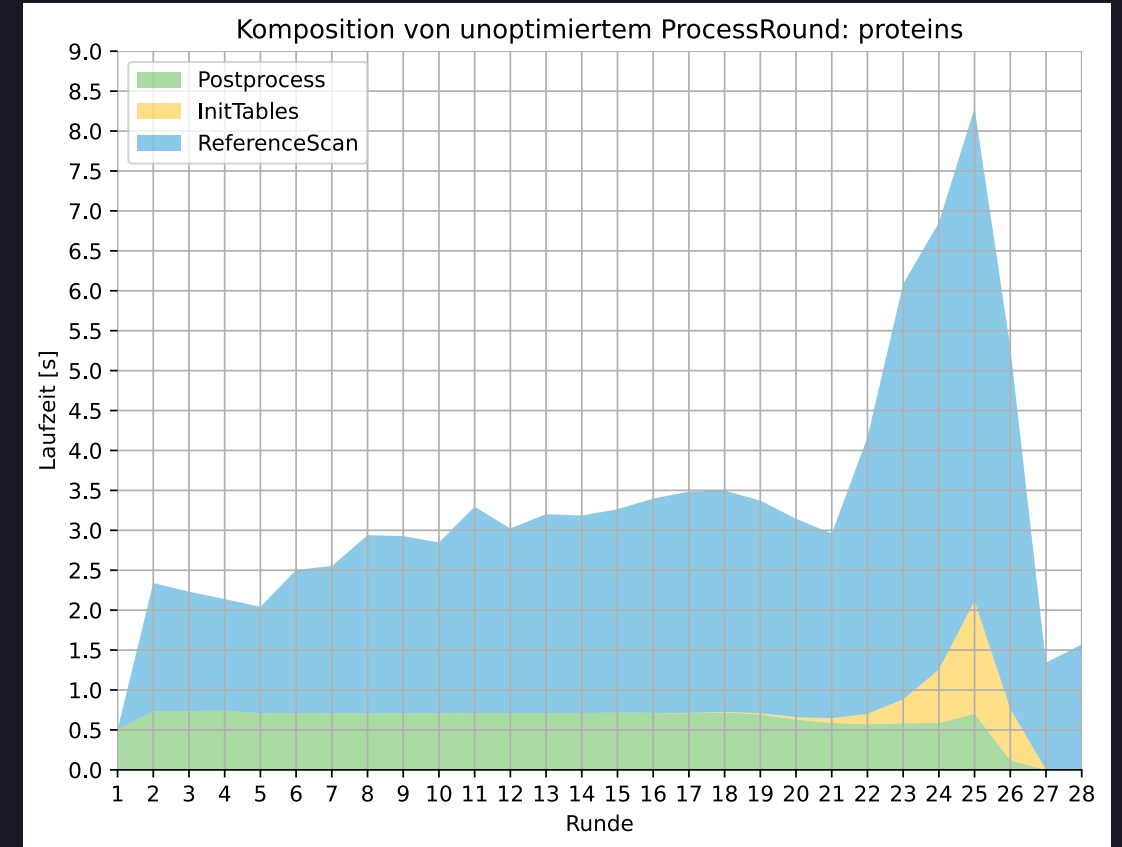
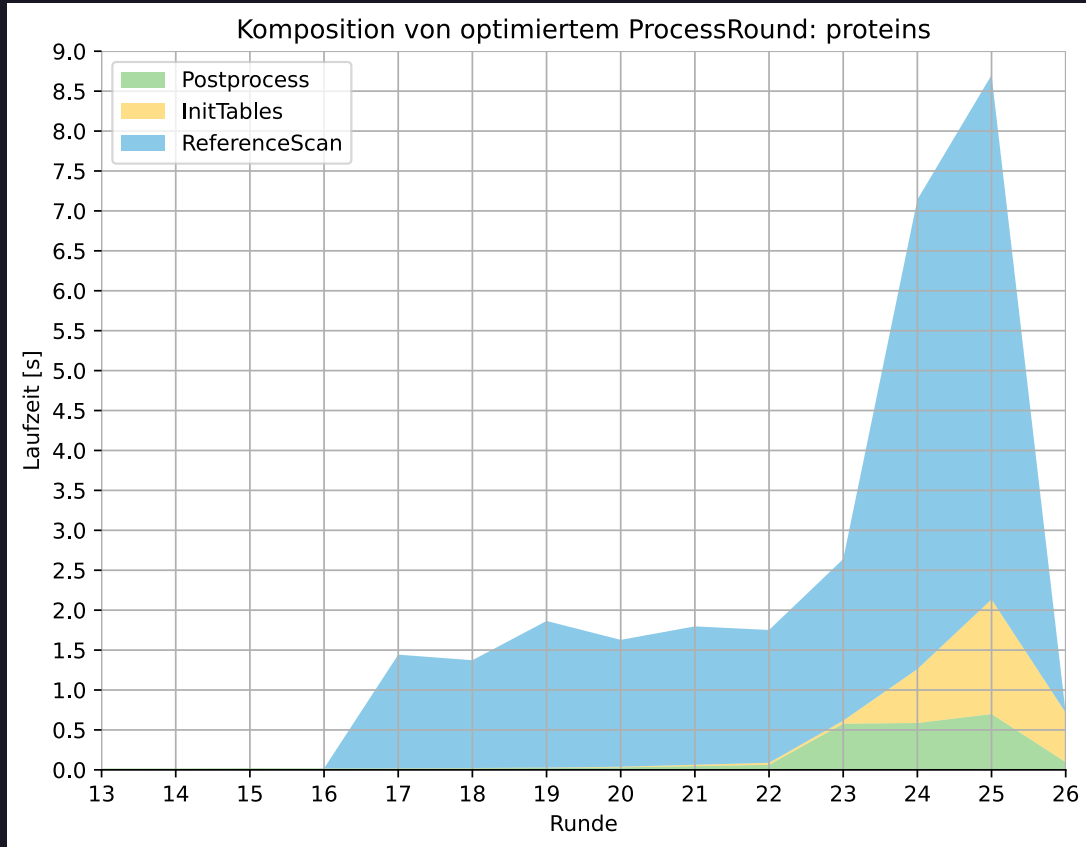
Optimierungen - PreMatching

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

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

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





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