

Bachelorarbeit

**Parallelisierung einer speichereffizienten  
Approximation der LZ77-Faktorisierung**

Gajann Sivarajah

Gutachter:

Prof. Dr. Johannes Fischer

M.Sc. Patrick Dinklage

Technische Universität Dortmund

Fakultät für Informatik

LS-11

<http://afe.cs.tu-dortmund.de>



# Inhaltsverzeichnis

<b>1</b>	<b>Einleitung</b>	<b>1</b>
1.1	Motivation und Hintergrund . . . . .	1
1.2	Aufbau der Arbeit . . . . .	1
<b>2</b>	<b>Grundlagen</b>	<b>3</b>
2.1	Eingabe . . . . .	3
2.2	Ausgabe $\rightarrow$ Faktorisierung . . . . .	3
2.3	Kompression . . . . .	4
2.3.1	Verlustfreie Kompression . . . . .	4
2.3.2	Dekompression . . . . .	4
2.3.3	String-Matching $\rightarrow$ Rabin-Karp . . . . .	4
2.3.4	Verlustbehaftete Kompression . . . . .	6
2.3.5	Binäre (De-)Kodierung . . . . .	6
2.3.6	Metriken . . . . .	7
2.4	Parallelität . . . . .	7
2.4.1	Shared-Memory-Modell . . . . .	7
2.4.2	Metriken . . . . .	7
<b>3</b>	<b>Kompressionsalgorithmen</b>	<b>9</b>
3.1	(exakte) LZ77-Kompression . . . . .	9
3.1.1	Konzept . . . . .	9
3.1.2	Theoretisches Laufzeit- und Speicherverhalten . . . . .	9
3.2	Approximation der LZ77-Faktorisierung(Approx. LZ77) . . . . .	10
3.2.1	Konzept . . . . .	10
3.2.2	Theoretisches Laufzeit- und Speicherverhalten . . . . .	12
3.3	Parallelisierung von Approx. LZ77(Approx. LZ77Par) . . . . .	14
3.3.1	Konzept . . . . .	14
3.3.2	Theoretisches Laufzeit- und Speicherverhalten . . . . .	15
3.4	Praktische Optimierungen . . . . .	15
3.4.1	Dynamische Endrunde(DynEnd) - Laufzeit vs. Qualität* . . . . .	16

3.4.2	Dynamische Startrunde(DynStart) - Laufzeit vs. Speicher . . . . .	16
3.4.3	Vorberechnete Runde(PreMatching) - Laufzeit vs. Speicher . . . . .	17
3.4.4	Minimale Tabellengröße(ScanSkip) - Laufzeit vs. Qualität . . . . .	18
<b>4</b>	<b>Praktische Evaluation</b>	<b>19</b>
4.1	Testumgebung . . . . .	19
4.2	Implementierung . . . . .	19
4.2.1	Klassenstruktur . . . . .	19
4.2.2	Code . . . . .	20
4.3	Messung . . . . .	20
4.3.1	Eingabedaten . . . . .	20
4.3.2	Messgrößen . . . . .	20
4.3.3	Messwerte . . . . .	22
4.4	Auswertung . . . . .	22
4.4.1	LZ77 . . . . .	22
4.4.2	Approx. LZ77 . . . . .	22
4.4.3	Approx. LZ77Par . . . . .	22
<b>A</b>	<b>Weitere Informationen</b>	<b>25</b>
A.1	Zusätzliche Messwerte . . . . .	25
	<b>Literaturverzeichnis</b>	<b>31</b>
	<b>Erklärung</b>	<b>31</b>

# Kapitel 1

## Einleitung

1.1 Motivation und Hintergrund

1.2 Aufbau der Arbeit



# Kapitel 2

## Grundlagen

Zunächst stellen wir die verwendete Terminologie und relevante Konzepte bzw. Phänomene dar.

### 2.1 Eingabe

Unsere Eingabe sei durch eine  $n$ -elementige Zeichenfolge  $S = e_1 \dots e_n$  über dem beschränkten numerischen Alphabet  $\Sigma$  mit  $e_i \in \Sigma \forall i = 1, \dots, n$  gegeben. Für jede beliebige Zeichenfolge  $S$  wird mit  $|S|$  dessen Länge, hier  $n$ , bezeichnet. Der Ausdruck  $S[i..j] \in \Sigma^{j-i+1}$  mit  $1 \leq i \leq j \leq n$  beschreibt die Teilfolge  $e_i \dots e_j$ , wobei im Falle, dass  $i = j$  ist, das einzelne Zeichen  $e_i$  referenziert wird. Alternativ kann ein einzelnes Zeichen  $e_i$  auch durch  $S[i]$  referenziert werden. Eine Teilfolge der Form  $S[1..k]$  mit  $1 \leq k \leq n$  wird als Präfix von  $S$  bezeichnet. Im Gegensatz dazu wird eine Teilfolge der Form  $S[k..n]$  als Suffix von  $S$  bezeichnet. Für zwei Teilfolgen  $S_1$  und  $S_2$  beschreibt der Ausdruck  $S_1 + S_2$  die Konkatenation der beiden Teilfolgen.

### 2.2 Ausgabe $\rightarrow$ Faktorisierung

Ein charakteristisches Merkmal der Familie der Lempel-Ziv-Kompressionsverfahren ist die Repräsentation der Ausgabe in Form einer Faktorisierung. Für eine Eingabe  $S = e_1 \dots e_n$  wird eine Faktorisierung  $F = f_1 \dots f_z$  mit  $z \leq n$  derart erzeugt, dass die Eingabe  $S$  durch die Faktorisierung in eine äquivalente Folge von nichtleeren Teilfolgen zerlegt werden. Dabei ist jeder Faktor  $f_i$  mit  $1 \leq i \leq z$  als Präfix von  $S[|f_1 \dots f_{i-1}| + 1..n]$  definiert, der bereits in  $S[1..|f_1 \dots f_i|]$  vorkommt oder als einzelnes Zeichen ohne vorheriges Vorkommen. Die im Folgenden betrachteten Algorithmen können speziell der Klasse der LZ77-Kompressionsverfahren zugeordnet werden, dessen Faktoren im Schema des Lempel-Ziv-Storer-Szymanski repräsentiert werden sollen.

$$F = f_1 \dots f_z \text{ mit } f_i = \begin{cases} (Length, Position) & \text{falls Referenz} \\ (0, Zeichen) & \text{sonst} \end{cases} \quad (2.1)$$

Zur Darstellung von Referenzen wird das Tupel aus der Position des vorherigen Vorkommens und der Länge des Faktors genutzt. Einzelne Zeichen können wiederum durch das Tupel aus dem Platzhalter 0 und dem entsprechenden Zeichen dargestellt werden. Das in 2.1 definierte Format beschreibt die gewünschte Ausgabe der im Folgenden betrachteten Algorithmen.

## 2.3 Kompression

### 2.3.1 Verlustfreie Kompression

Der Prozess der Kompression überführt eine Repräsentation einer finiten Datenmenge in eine möglichst kompaktere Form. Eine verlustfreie Kompression ist gegeben, falls die Abbildung zwischen der ursprünglichen und komprimierten Repräsentation bijektiv ist. Die Korrektheit einer verlustfreien Kompression kann daher durch die Angabe einer Dekompressionsfunktion nachgewiesen werden. Ist diese Voraussetzung nicht gegeben, so handelt es sich um eine verlustbehaftete Kompression, da eine Rekonstruktion der ursprünglichen Datenmenge nicht garantiert werden kann.

### 2.3.2 Dekompression

Die Dekompression beschreibt den Umkehrprozess der Kompression und erlaubt im Falle einer verlustfreien Kompression die Rekonstruktion der ursprünglichen Datenfolge. Im Falle von Verfahren der LZ77-Familie, kann die Dekompression durch die folgende Abbildung definiert werden,

$$DECOMP_{LZ77} : F(1..z) \rightarrow S(1..n). \quad (2.2)$$

Der dargestellte Algorithmus 2.1 beschreibt eine mögliche Implementierung der Dekompression für eine Faktorisierung  $F = f_1 \dots f_z$  zu der Eingabe  $S = e_1 \dots e_n$ . Der beschriebene Algorithmus iteriert durch alle Faktoren und fügt die referenzierten Zeichen einzeln in die Ausgabe  $S$  ein. Damit kann die Laufzeit des Algorithmus auf  $O(n)$  geschätzt werden.

### 2.3.3 String-Matching $\rightarrow$ Rabin-Karp

Im Rahmen des approximativen Algorithmus, welcher in dieser Arbeit beschrieben wird, werden Vergleiche von Zeichenfolgen mithilfe des Rabin-Karp-Fingerprints(RFP) durchge-



**Algorithmus 2.1** DECOMP<sub>LZ77</sub>


---

*Eingabe:*  $F = f_1 \dots f_z$  *Ausgabe:*  $S = e_1 \dots e_n$

---

```

 $S \leftarrow \emptyset$ 
for  $i = 1$  to  $z$  do
   $(len, ref) \leftarrow f_i$ 
  if  $len = 0$  then
     $S \leftarrow S + ref$ 
  else
    for  $j = 0$  to  $len - 1$  do
       $S \leftarrow S + S[ref + j]$ 
    end for
  end if
end for
return  $S$ 

```

---

führt. Sei  $p \in \mathbb{P}$  eine Primzahl und  $b \in \mathbb{N}$  eine Basis, so kann der RFP einer Zeichenfolge  $S$  der Länge  $n$  durch den Ausdruck

$$RFP(S) = \sum_{i=1}^n S[i]b^{n-i} \mod p \quad (2.3)$$

$$\in \{0, \dots, p-1\}$$

berechnet werden. Hierbei wird eine Zeichenfolge beliebiger Länge in eine Zahl aus dem Intervall  $[0, p-1]$  abgebildet. Der RFP erlaubt es, die Gleichheit zweier Zeichenfolgen zu widerlegen im Falle von unterschiedlichen Fingerprints. Im Falle von gleichen Fingerprints, kann die Gleichheit der Zeichenfolgen jedoch nicht garantiert werden. Die Wahrscheinlichkeit einer Kollision dieser Art bei Zeichenfolgen gleicher Größe ist jedoch beschränkt und praktisch gering. Insbesondere kann die Wahrscheinlichkeit durch die passende Wahl von  $p$  und  $b$  minimiert werden.

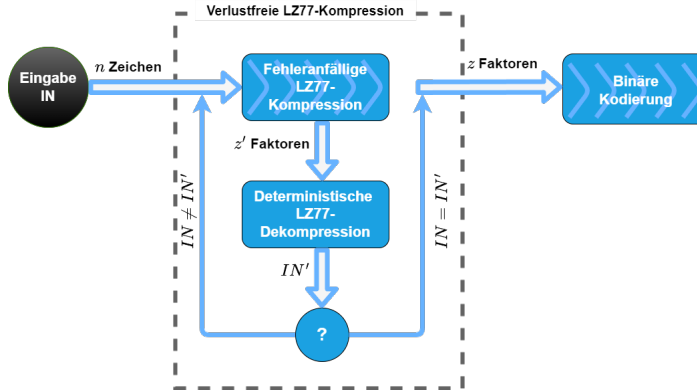
Rabin-Karp-Fingerprints erlauben verschiedene Operationen auf Zeichenfolgen, die im Rahmen der approximativen LZ77-Faktorisierung effizient genutzt werden können. Zum Einen kann ein beschränktes Fenster  $S_W = S(j..j+w)$  der Länge  $w < n$  leicht verschoben werden. Sei  $RFP(S_W)$  der Fingerabdruck des Fensters, so kann der Fingerabdruck durch die Verschiebung um ein Zeichen nach rechts durch den Ausdruck,

$$RFP(S(j+1..j+w+1)) = (RFP(S_W) - S[j]b^{w-1})b + S[j+w+1] \mod p, \quad (2.4)$$

beschrieben. Desweiteren seien zwei Teilfolgen  $S_1$  und  $S_2$  der gleichen Länge  $n$  gegeben. Der RFP der Konkatenation  $S_1 + S_2$  der beiden Teilfolgen kann durch den Ausdruck,

$$RFP(S_1 + S_2) = (RFP(S_1)b^n + RFP(S_2)) \mod p, \quad (2.5)$$

Abbildung 2.1: Las-Vegas-Algorithmus



berechnet werden. Ähnlich zur Konkatenation von Zeichenfolgen ist die Operation 2.5 ebenfalls assoziativ, jedoch nicht kommutativ.

### 2.3.4 Verlustbehaftete Kompression

Im Rahmen dieser Arbeit werden wir einen Approximationsalgorithmus betrachten, der aufgrund der verwendeten RFP-Technik für Vergleiche von Zeichenfolgen eine fehlerhafte Faktorisierung mit einer beschränkten Wahrscheinlichkeit erzeugen kann. Die Korrektheit der Dekompression kann intern und extern durch explizite Vergleiche der Zeichenfolgen erkannt werden. Da der Kompressionsprozess in diesem Fall mit anderen Parametern wiederholt werden kann, können wir einen verlustfreien Las-Vegas-Algorithmus konstruieren.

In Abbildung 2.1 wird eine Regelsteuerung illustriert. Der Algorithmus wird solange wiederholt, bis eine korrekte Faktorisierung erzeugt wurde. Dass die Anzahl der Wiederholungen beschränkt ist, werden wir in der Analyse des Approximationsalgorithmus und der praktischen Evaluation zeigen.

### 2.3.5 Binäre (De-)Kodierung

Die Kodierung  $K_{IN} : \Sigma^* \rightarrow \{0, 1\}^*$  überführt unsere Eingabe aus dem Alphabet  $\Sigma$  in eine binäre Repräsentation. Die Umkehrabbildung  $K_{IN}^{-1} : \{0, 1\}^* \rightarrow \Sigma^*$  definiert die Dekodierung und überführt eine binäre Repräsentation in eine Zeichenfolge aus dem Alphabet  $\Sigma$ . Im Rahmen dieser Arbeit gehen wir davon aus, dass unsere Eingabe  $S$  über dem Alphabet  $\Sigma = \{1, \dots, 255\}$  erzeugt wurde und jedes Zeichen durch 8 Bits, oder 1 Byte, dargestellt wird. Für die Länge,  $|S|_{Bin}$ , der binären Repräsentation folgt,

$$|S|_{Bin} = |K_{IN}(S)| = 8 * |S|. \quad (2.6)$$

Die eingelesene Eingabefolge wird durch den Kompressionsalgorithmus in die Faktorisierung  $F = f_1 \dots f_z$  überführt. Die bijektive Abbildung  $K_{OUT} : F \rightarrow \{0, 1\}^*$  definiert die Kodierung der Faktoren in eine binäre Repräsentation. Analog dazu wird die Dekodierung

$K_{OUT}^{-1} : \{0,1\}^* \rightarrow F$  definiert. Im Gegensatz zur Eingabe, werden wir keine Kodierung bzw. Dekodierung der Faktoren vorgeben, da diese durch den Kompressionsalgorithmus nicht beschränkt wird. Für eine beliebige lineare Kodierung  $K_{OUT}$  ergibt sich die binäre Ausgabegröße  $|F|_{Bin}$  durch

$$|F|_{Bin} = \sum_{i=1}^z |K(f_i)|. \quad (2.7)$$

### 2.3.6 Metriken

Die Qualität einer Kompression kann durch verschiedene Metriken quantifiziert werden. Zum Einen beschreibt die Kompressionsrate  $CR$  den Grad der Kompression und ist durch den Ausdruck,

$$CR = \frac{|F|_{Bin}}{|S|_{Bin}} \quad (2.8)$$

, definiert. Da die Kodierung der Faktoren nicht eindeutig aus der Wahl des Kompressionsalgorithmus eingegrenzt wird, ist stattdessen die Anzahl der erzeugten Faktoren ein weiteres geeignetes Gütemaß. Für die Eingabe  $S$  der Länge  $n$  und der Ausgabe  $f_1 \dots f_z$  sei die Faktorrare durch

$$FR = \frac{z}{n} \quad (2.9)$$

gegeben. In beiden Fällen wird ein niedriger Wert bevorzugt, da dieser auf eine bessere Extraktion von Redundanzen hinweist.

## 2.4 Parallelität

Das Ziel dieser Arbeit ist die Entwicklung und Evaluation eines parallel Kompressionsalgorithmus. Im Folgenden definieren wir die Rahmenbedingungen und Konzepte der Parallelität.

### 2.4.1 Shared-Memory-Modell

Unser Algorithmus agiere auf einem Shared-Memory-Modell mit  $P$  Ausführungseinheiten, welches im Gegensatz zum Distributed-Memory-Modell allen beteiligten Ausführungseinheiten bzw. Prozessoren einen gemeinsamen Zugriff auf den Speicher ermöglicht. Im Rahmen der Arbeit und Kommunikation unter den Prozessoren wird man jedoch auf Konflikte bei gleichzeitigen Speicherzugriffen zustoßen. Ein parallel modellierter Algorithmus muss explizit hinsichtlich der Korrektheit und Effizienz Mechanismen zur Synchronisation implementieren.

### 2.4.2 Metriken

Das Ziel der Parallelisierung eines Algorithmus liegt hauptsächlich in einer Verbesserung der Laufzeit, insbesondere unter Berücksichtigung von Ressourcenkonflikten. Die zeitliche

Beschleunigung der Laufzeit kann durch den Speedup  $SP$  bemessen werden. Für eine Eingabe  $S$  der Länge  $n$  brauche ein sequenzieller Durchlauf  $T(n, p = 1)$  Zeit, während ein paralleler Algorithmus mit  $P$  Prozessoren  $T(n, p = P)$  an Zeit benötigt. Der Speedup ist dabei definiert durch

$$SP(n, P) = \frac{T(n, 1)}{T(n, P)}. \quad (2.10)$$

Ein idealer Speedup ist gegeben durch  $SP(n, P) = P$ . Verschiedene Effekte im Rahmen des Speicherzugriffs, der Synchronisation und der Kommunikation über mehrere Prozessoren können jedoch die Effizienz der Parallelisierung stark beeinträchtigen. Insbesondere können sequenzielle Abschnitte im Algorithmus aufgrund des Amdahlschen Gesetzes eine obere Schranke für den Speedup setzen.

## Kapitel 3

# Kompressionsalgorithmen

### 3.1 (exakte) LZ77-Kompression

Der im Folgenden beschriebene Algorithmus für die Generierung einer exakten LZ77-Faktorisierung dient als Referenz für die Evaluation der approximativen Algorithmen.

#### 3.1.1 Konzept

Wie bereits in Kapitel xx beschrieben, erzeugen Algorithmen der LZ77 - Familie eine Faktorisierung einer Eingabezeichenfolge  $S$ , wobei die Faktoren entweder Referenzen zu vorherigen Zeichenfolgen oder einzelne Zeichen sein können. Im Rahmen der exakten LZ77 - Faktorisierung wird ein Greedy - Ansatz verwendet, um von links nach rechts stets die längste Zeichenfolge zu referenzieren, die bereits links von der aktuellen Position vorkommt. In 3.1 wird der Algorithmus zur Generierung einer exakten LZ77-Faktorisierung beschrieben. Der Algorithmus erzeugt zunächst ein SuffixArray, welches allen Suffixen der Eingabe eine lexikographische Ordnung zuweist. Mithilfe der lexikographischen Ordnung können Kandidaten für Referenzen effizient gefunden werden. Hierfür werden mit Hilfe des SuffixArrays zwei Arrays, das Next Smaller Value(NSV) und das Previous Smaller Value(PSV) erzeugt. Sei die aktuelle Position in der Eingabe  $k$ , so muss aufgrund von positionellen und lexikographischen Einschränkungen die Position  $ref$  der längsten vorherigen Referenz  $NSV[k]$  oder  $PSV[k]$  sein. Die maximale Länge der übereinstimmenden Präfixe zwischen  $S(NSV[k]..n)$  und  $S(k..n)$  bzw.  $S(PSV[k]..n)$  und  $S(k..n)$  wird durch die Funktion  $LCP$  berechnet. Das Ergebnis dieser Berechnung bestimmt den Faktor  $(len, ref)$ , welcher in der Eingabe an Position  $k$  beginnt. Der Algorithmus terminiert, wenn die gesamte Eingabe abgearbeitet wurde.

#### 3.1.2 Theoretisches Laufzeit- und Speicherverhalten

Die Berechnung des SuffixArrays und die folgende Berechnung der NSV- und PSV-Arrays können mithilfe von Algorithmen aus der Literatur(siehe xx) in  $O(n)$  Laufzeit durchgeführt

**Algorithmus 3.1**  $\text{COMP}_{\text{LZ77}}$ 


---

Eingabe:  $S = e_1 \dots e_n$  Ausgabe:  $F = f_1 \dots f_z$   
 $SA \leftarrow \text{SuffixArray}(S)$   
 $(NSV, PSV) \leftarrow (\text{NSVArray}(S, SA), \text{PSVArray}(S, SA))$   
 $F \leftarrow \emptyset$   
 $k \leftarrow 1$   
**while**  $k \leq n$  **do**  
      $(len, ref) \leftarrow (0, 0)$   
      $l_{nsv} \leftarrow \text{LCP}(S(NSV[k]..n), S(k..n))$   
      $l_{psv} \leftarrow \text{LCP}(S(PSV[k]..n), S(k..n))$   
     **if**  $l_{nsv} > l_{psv}$  **then**  
          $(len, ref) \leftarrow (l_{nsv}, NSV[k])$   
     **else if**  $l_{nsv} < l_{psv}$  **then**  
          $(len, ref) \leftarrow (l_{psv}, PSV[k])$   
     **else**  
          $(len, ref) \leftarrow (0, S[k])$   
     **end if**  
      $F \leftarrow F + (len, ref)$   
      $k \leftarrow k + len + 1$   
**end while**  
**return**  $F$

---

werden. In der abschließenden Schleife repräsentiert die  $k$ -te Iteration den  $k$ -ten Faktor, wobei die Iteration für die Berechnung der Faktorlänge  $O(|f_k|)$  Laufzeit benötigt. Damit ergibt sich eine Gesamtlaufzeit von  $O(n + \underbrace{\sum_{i=1}^z |f_i|}_n) = O(n)$  für die Generierung der exak-

ten LZ77-Faktorisierung. Der Speicherbedarf des Algorithmus beträgt  $O(n)$ , da sich die Größe des SuffixArrays und der NSV- und PSV-Arrays linear zur Eingabelänge verhalten. Es sollte jedoch angemerkt werden, dass die Linearität des Speicherbedarfs einen hohen konstanten Faktor hat und unabhängig von der Beschaffenheit der Eingabe und der Anzahl der Faktoren ist.

## 3.2 Approximation der LZ77-Faktorisierung(Approx. LZ77)

### 3.2.1 Konzept

Im Rahmen dieser Arbeit wird die erste Phase einer speichereffizienten Approximation des LZ77-Algorithmus betrachtet, welche auch eine verlustfreie Kompression darstellt. Wie in [1] beschrieben, kann die Kombination aller drei Phasen des Algorithmus eine

2-Approximation bezüglich der Faktorraten ermöglichen. Die resultierenden Faktoren entsprechen jedoch dem LZ77-Schema, sodass eine verlustfreie Dekompression mit 2.1 möglich ist. Im weiteren Verlauf werden wir die sequenzielle erste Phase des Algorithmus als Approx.LZ77 bezeichnen. Im Gegensatz zur exakten LZ77-Faktorisierung werden Referenzen nicht durch einen Greedy-Ansatz mit einem Scan von links nach rechts gefunden. Stattdessen wird eine Approximation der exakten LZ77-Faktorisierung erzeugt, die einen Tradeoff zwischen der Faktorraten und der Performanz, hier der Speicherverbrauch, des Algorithmus darstellt. Die resultierenden Faktoren sind insbesondere dadurch definiert, dass ihre Länge einer Zweierpotenz entspricht. Analog dazu gehen wir ohne Beschränkung der Allgemeinheit davon aus, dass die Länge der Eingabe ebenfalls eine Zweierpotenz ist. Eine abweichende Eingabelänge kann stets durch entsprechendes Padding erreicht werden. Die Approximation wird in mehreren Runden durchgeführt, wobei in jeder Runde die noch unverarbeitete Eingabe in Blöcke gleicher Größe eingeteilt wird. Innerhalb einer Runde werden für Zeichenfolgen, die durch die jeweiligen Blöcke repräsentiert werden, Referenzen, also vorherige Vorkommen, gesucht. Im Erfolgsfall wird ein entsprechender Faktor extrahiert und der Block gilt als markiert. Die verbleibenden Blöcke werden im Übergang zur nächsten Runde in der Hälfte geteilt. Der Algorithmus terminiert, wenn die Eingabe vollständig verarbeitet wurde, was spätestens nach  $\log_2(|S|)$  Runden der Fall ist, da Blöcke der Größe 1 nicht weiter geteilt werden können. Die Zeichenfolgen der Blöcke, die in der letzten Runde nicht markiert wurden, werden als einzelne Zeichen interpretiert und faktorisiert. Der Ablauf des Algorithmus wird in 3.2 dargestellt. Initial wird die Eingabe in  $2^r$  Blöcke gleicher Größe eingeteilt, wobei  $r$  die initiale Runde beschreibt. Daraufhin wird eine Schleife über alle Runden gestartet, wobei in jeder Iteration die aktuelle Menge der Blöcke auf Referenzen untersucht und ggf. markiert werden. Die Menge der markierten Blöcke, die auch zur Menge der Faktoren  $F$  hinzugefügt werden, wird in der Routine ProcessRound bestimmt. Nachdem die Menge der Blöcke um die markierten Blöcke bereinigt wurde, werden die verbleibenden Blöcke halbiert und die nächste Runde initiiert.

---

**Algorithmus 3.2**  $\text{COMP}_{\text{ApproxLZ77}}$ 


---

*Eingabe:*  $S = e_1 \dots e_n$  *Ausgabe:*  $F = f_1 \dots f_z$

---

```

1:  $F \leftarrow \emptyset$ 
2:  $r \leftarrow 1$ 
3:  $\text{Blocks}[1..2^r] \leftarrow \text{InitBlocks}(S, 2^r)$  // Split S into  $2^r$  equal blocks
4: while  $r \leq \log_2(|S|)$  do
5:    $\text{markedBlocks}[1..z_r] \leftarrow \text{ProcessRound}(r, S, \text{Blocks}, F)$ 
6:    $\text{Blocks} \leftarrow \text{NextNodes}(\text{Blocks} \setminus \text{markedBlock}[1..z_r])$  // Halve unmarked blocks
7:    $r \leftarrow r + 1$ 
8: end while
9: return  $F$ 

```

---

Die Funktionsweise der ProcessRound-Routine wird in 3.3 beschrieben. Die Ausführung kann in drei Schritte unterteilt werden. Im ersten Schritt werden die Tabellen RFPTable und RefTable mithilfe der Routine InitTables initialisiert. Die RFPTable ordnet jedem einzigartigen Rabin-Karp-Fingerprint(RFP), der über die Menge aller Blöcke erzeugt wird, den linken Block zu, welcher diesen Wert aufweist. In diesem Rahmen vernachlässigen wird explizit Kollisionfälle des RFP und gehen davon aus, dass die Gleichheit der RFPs zweier Zeichenfolgen ebenfalls die Gleichheit der Zeichenfolgen impliziert. Die RefTable speichert für jeden Block die Position des linken Vorkommens seiner Zeichenfolge in der Eingabe, die zum aktuellen Zeitpunkt bekannt ist. Die Routine InitTables, wie in 3.4 beschrieben, erzeugt die RFPTable, indem alle Blöcke von links nach rechts durchlaufen werden und das Paar aus einem Block und dessen RFP in die Tabelle nur dann eingefügt wird, wenn der RFP noch nicht vorhanden ist. Als Konsequenz wird für jeden RFP, der eine einzigartige Zeichenfolge repräsentiert, der entsprechende linke Block gespeichert. Die RefTable wird analog initialisiert. Falls der RFP eines Blocks bereits in der RFPTable vorhanden ist, so kann dem Block ein vorheriges Vorkommen zugeordnet werden. Andernfalls ist der Block das linke Vorkommen seiner Zeichenfolge und dessen Position wird entsprechend gespeichert. Es ist zu betonen, dass die InitTables-Routine durch Aktualisierungen der RefTable bereits markierte Blöcke bzw. Faktoren erzeugt.

Im zweiten Schritt werden zusätzliche Referenzen innerhalb der gesamten Eingabe  $S$  gesucht. Dazu wird der RFP eines Fensters der Größe  $\frac{|S|}{2^r}$  über die Eingabe berechnet und mit den Einträgen der RFPTable verglichen. Im Falle eines Treffers wird die Position des Fensters ggf. in der RefTable eingetragen. Im Rahmen des Scans wird ein Rolling-Hash verwendet, um die RFPs der Fenster effizient um eine Position zu verschieben. Schließlich wird im dritten Schritt die RefTable genutzt, um die markierten Blöcke zu bestimmen. Ein Block wird markiert, wenn die Position des linken Vorkommens seiner Zeichenfolge in der RefTable kleiner als die aktuelle Position des Blocks ist. Die markierten Blöcke werden in die Menge der Faktoren  $F$  eingefügt, wobei die Reihenfolge der Faktoren entweder während des Einfügens oder durch eine nachfolgende Sortierung sichergestellt werden muss.

### 3.2.2 Theoretisches Laufzeit- und Speicherverhalten

Die Laufzeit des Algorithmus wird durch die Anzahl der Runden und der Extraktion von Referenzen in jeder Runde bestimmt. Die Anzahl der Runden beträgt maximal  $\log_2(|S|) = \log_2(n)$ . In jeder Runde werden Referenzen innerhalb der Blöcke durch die InitTables-Routine bestimmt. Die Menge aller Blöcke, die in dieser Routine bearbeitet werden, decken maximal die gesamte Eingabe  $S$  ab. Die Laufzeit der Routine erhält durch die Berechnung des RFPs über dieser Zeichenfolge eine Laufzeitschätzung von  $O(n)$ . Die Referenzsuche über die gesamte Eingabe durch die ReferenceScan-Routine benötigt ebenfalls  $O(n)$  Laufzeit, da das berücksichtigte Fenster in linearer Zeit über die Eingabe verschoben wer-



**Algorithmus 3.3** ProcessRound*Eingabe:*  $r, S, \text{Blocks}$  *Ausgabe:*  $\text{markedBlocks}$ 


---

```

1: ( $\text{RFPTable}, \text{RefTable}$ )  $\leftarrow \text{InitTables}(\text{Blocks})$ 
2:  $\text{ReferenceScan}(S, \text{Blocks}, \text{RFPTable}, \text{RefTable})$ 
3: for  $i \leftarrow 1$  to  $|\text{Blocks}|$  do
4:   if  $\text{RefTable}[i] < \text{Blocks}[i].\text{Pos}$  then
5:      $\text{markedBlocks.insert}(i)$ 
6:      $F.\text{insert}(\frac{|S|}{2^r}, \text{RefTable}[i])$  // Ordered Insert
7:   end if
8: end for
9: return  $\text{markedBlocks}$ 

```

---

**Algorithmus 3.4** InitTables*Eingabe:*  $\text{Blocks}$  *Ausgabe:*  $\text{RFPTable}, \text{RefTable}$ 


---

```

1:  $\text{RFPTable} \leftarrow \text{HashTable}[\text{RFP}, \text{LeftMostBlock}]$ 
2:  $\text{RefTable}[1..|S|] \leftarrow (0, \dots, 0)$ 
3: for  $i \leftarrow 1$  to  $|\text{Blocks}|$  do
4:    $\text{RFP} \leftarrow \text{RFP}(\text{Blocks}[i])$ 
5:   if  $\text{RFP}$  in  $\text{RFPTable}$  then
6:      $\text{RefTable}[i] \leftarrow \text{RFPTable}[\text{RFP}].\text{Pos}$ 
7:   else
8:      $\text{RFPTable.insert}(\text{RFP}, \text{Blocks}[i])$ 
9:      $\text{RefTable}[i] \leftarrow \text{Blocks}[i].\text{Pos}$ 
10:  end if
11: end for
12: return  $\text{RFPTable}, \text{RefTable}$ 

```

---

den kann. Die Laufzeit der ProcessRound-Routine beträgt somit  $O(n)$ . Insgesamt ergibt sich eine Gesamtlaufzeit von  $O(n \cdot \log_2(n))$  für die Generierung der approximativen LZ77-Faktorisierung. Der Speicherbedarf des Algorithmus wird durch die Größe der RFPTable und der RefTable bestimmt, die wiederum durch die Anzahl der Blöcke bestimmt wird. In jeder Runde repräsentiert die Menge der Blöcke die noch unverarbeitete Eingabe, sodass aus jedem Block im Laufe des Algorithmus mindestens ein Faktor extrahiert wird. Die Anzahl der Blöcke in der Runde übersteigt damit nie die Anzahl der endgültigen Faktoren. Der Speicherbedarf kann damit konservativ auf  $O(z)$  abgeschätzt werden.

**Algorithmus 3.5** ReferenceScan*Eingabe:*  $S, Blocks, RFPTable, RefTable$ 


---

```

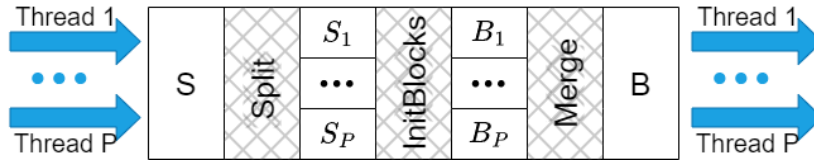
blockSize  $\leftarrow \frac{|S|}{|Blocks|}$ 
 $RFP \leftarrow RFP(S[1..blockSize])$ 
for  $i \leftarrow 1$  to  $|S| - blockSize$  do
  if  $RFP \text{ in } RFPTable$  and  $i < RefTable[RFPTable[RFP]]$  then
     $RefTable[RFPTable[RFP]] \leftarrow i$ 
  end if
   $RFP \leftarrow RFP(S[i + 1..i + blockSize])$  // Rolling Hash
end for

```

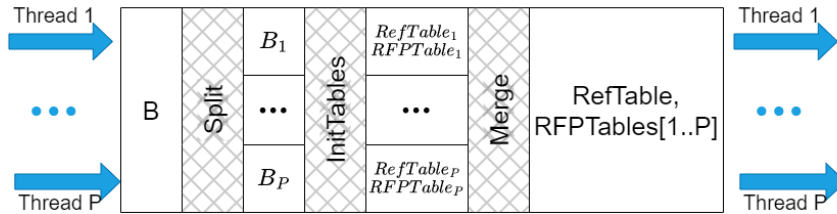
---

**3.3 Parallelisierung von Approx. LZ77(Approx. LZ77Par)****3.3.1 Konzept**

Im Folgenden beschreiben wir das verwendete Schema zur Parallelisierung des Approx. LZ77-Algorithmus. Die Parallelisierung zielt auf eine schrittweise Bearbeitung der Eingabe, die jeweils auf die verschiedenen Prozessoren bzw. Threads verteilt werden. In 3.1 wird die

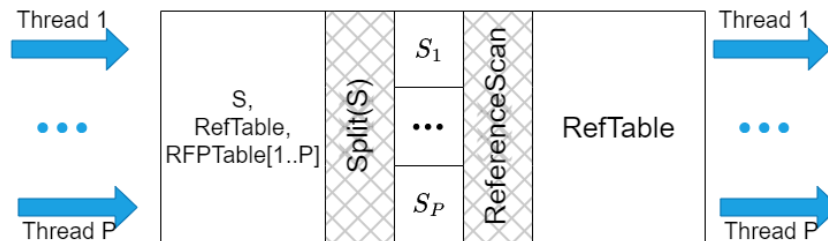
**Abbildung 3.1:** Parallele Generierung der initialen Blöcke

parallele Generierung der initialen Blöcke dargestellt. Hierbei wird die Eingabe  $S$  in  $P$  Teile aufgebrochen, sodass die Routine *InitBlocks* auf jedem Prozessor die zugehörigen Blöcke erzeugen kann. Da die chronologische Ordnung erhalten bleibt, können die erzeugten Blöcke ohne zusätzlichen Aufwand kombiniert werden.

**Abbildung 3.2:** Parallele Initialisierung der RFP-Tabelle und der Referenztafel

In 3.2 wird die parallele Initialisierung der RFP-Tabelle und der Referenztafel dargestellt. Die Menge aller Blöcke wird in  $P$  Teile aufgeteilt, wobei als Kriterium für die Aufteilung der RFP jedes Blocks verwendet wird. Als Konsequenz werden identische Zeichenfolgen nicht auf verschiedene Prozessoren verteilt. Jeder Prozessor erzeugt eine eige-

ne RFP-Tabelle, wobei die Referenztabelle durch alle Prozessoren gemeinsam aktualisiert wird. Auch hier wird durch den Aufteilungsschlüssel eine Überlappung der Zugriffe vermieden. Das Ergebnis ist eine konsistente Referenztabelle und eine Menge von  $P$  RFP-Tabellen, die jeweils eine klar definierte Teilmenge der RFPs enthalten.



**Abbildung 3.3:** Paralleler Scan der Eingabe nach zusätzlichen Referenzen

In 3.3 wird der parallele Scan der Eingabe  $S$  nach zusätzlichen Referenzen dargestellt. Die Eingabe  $S$  wird in  $P$  Teile aufgeteilt. Auf jede Teilmenge wird die sequenzielle Routine `ReferenceScan` angewendet, wobei die Referenztabelle durch alle Prozessoren gemeinsam aktualisiert wird und ein paralleler Lesezugriff auf die RFP-Tabellen stattfindet.

Schließlich müssen die implizit erzeugten Faktoren in die Menge der Faktoren  $F$  eingefügt werden. Da die Reihenfolge der Faktoren erhalten bleiben muss, wird eine nachträgliche parallele Sortierung verwendet. Die Implementierung der parallelen Sortierung ist jedoch nicht Gegenstand dieser Arbeit und wird daher nicht weiter betrachtet. Eine mögliche Implementierung ist in [] beschrieben.

### 3.3.2 Theoretisches Laufzeit- und Speicherverhalten

Eine theoretische Laufzeit von  $O(\frac{n \log n}{p})$  kann erreicht werden, wobei  $p$  die Anzahl der Prozessoren ist. Der Speicherbedarf des Algorithmus beträgt  $O(z)$ . Dies stellt jedoch eine ideale Abschätzung dar, die in der Praxis nicht erreicht werden kann. Insbesondere die Interaktion mit dem Speicher und die Kommunikation zwischen den Prozessoren führen zu einer oberen Schranke des Speedups.

## 3.4 Praktische Optimierungen

Im Folgenden betrachten wir optionale Optimierungen, die die durchschnittliche Laufzeit von Approx. LZ77 verbessern können auf Kosten von anderen Metriken. Jede einzelne Technik ist unabhängig von den anderen nutzbar, wobei eine positive Korrelation zu erwarten ist.

### 3.4.1 Dynamische Endrunde(DynEnd) - Laufzeit vs. Qualität\*

Sei eine Kodierung  $K$  für die Übersetzung der erzeugten Faktorenfolge  $F$  gegeben. Der Wert,

$$Min_{Bin}^{Ref} = \min\{|K(f)| \mid f \in F, f \text{ ist Referenz}\} \quad (3.1)$$

gibt die minimale Anzahl an Bits an, die für die Kodierung einer Referenz benötigt wird. Analog dazu beschreibt

$$Max_{Bin}^{Lit} = \max\{|K(f)| \mid f \in F, f \text{ ist Zeichen}\} \quad (3.2)$$

die maximale Anzahl an Bits, die für die Kodierung eines einzelnen Zeichens benötigt wird. Sei  $f_{ref}$  ein beliebiger referenzierender Faktor, welcher  $|f_{ref}| \leq \frac{Min_{Bin}^{Ref}}{Max_{Bin}^{Lit}}$  Zeichen referenziert. Die referenzierte Zeichenfolge von  $f_{ref}$  wird im Folgenden als  $S_{ref}$  mit  $|S_{ref}| = |f_{ref}|$  bezeichnet. Dann gilt für die Länge der kodierten Repräsentation von  $f_{ref}$ :

$$\begin{aligned} |K(f_{ref})| &\geq Min_{Bin}^{Ref} \\ &\geq |f_{ref}| \cdot Max_{Bin}^{Lit} \\ &\geq \sum_{i=1}^{|f_{ref}|} |K((0, S_{ref}(i)))|. \end{aligned} \quad (3.3)$$

Es folgt, dass ein referenzierender Faktor, dessen Länge eine obere Schranke von  $\frac{Min_{Bin}^{Ref}}{Max_{Bin}^{Lit}}$  Zeichen nicht überschreitet, nicht effizient kodiert werden kann. Stattdessen sollten die referenzierten Zeichen einzeln kodiert werden. Die Technik der dynamischen Endrunde greift diese Idee auf, indem Referenzen unterhalb einer Grenzlänge nicht berechnet werden. Gibt uns die Kodierung eine Grenzlänge  $l_{min}^{ref}$  vor, so kann der Algorithmus in Runde  $r = \lceil \log_2 |S| - \log_2 l_{min}^{ref} \rceil$  terminieren. Da potenziell referenzierende Faktoren aufgebrochen werden, kann die Qualität der Faktorisierung sinken, wobei das binäre Endprodukt kleiner wird. Es ergibt sich also eine steigende Faktorraten bei sinkender Kompressionsrate.

$$CR_{DynEnd}^{Approx.LZ77} \leq CR^{Approx.LZ77} \quad (3.4)$$

$$FR_{DynEnd}^{Approx.LZ77} \geq FR^{Approx.LZ77} \quad (3.5)$$

### 3.4.2 Dynamische Startrunde(DynStart) - Laufzeit vs. Speicher

Gegeben seien zwei initiale Runden  $r_{init1}$  und  $r_{init2}$  mit  $1 \leq r_{init1} < r_{init2} \leq \log_2 |S|$ , die auf der gesamten Eingabe  $S$  angewendet werden, so wird die Eingabe jeweils in  $2^{r_{init1}}$  bzw.  $2^{r_{init2}}$  Blöcke gleicher Größe eingeteilt. Die Menge der Blöcke werde im Folgenden als  $B_{init1}$  bzw.  $B_{init2}$  bezeichnet. Im Rahmen der Bearbeitung der Runden wird eine Menge von markierten Blöcken  $B_{init1}^{marked} \subset B_{init1}$  bzw.  $B_{init2}^{marked} \subset B_{init2}$  erzeugt, für die ein vorheriges Vorkommen bestimmt wurde. Aufgrund der Natur der Blockhalbierung in jeder Runde, kann jedem Block in  $B_{init1}$  eine Gruppe von  $2^{r_{init2}-r_{init1}}$  Blöcken in  $B_{init2}$  zugeordnet

werden, die die gleiche Zeichenfolge repräsentieren. Die Folgerung lässt sich insbesondere auch auf die markierten Blöcke anwenden, sodass die folgende Beziehung hergeleitet werden kann:

$$|B_{marked}^{init2}| \geq 2^{r_{init2}-r_{init1}} \cdot |B_{marked}^{init1}|. \quad (3.6)$$

Weiterhin folgt, dass die Existenz eines markierten Blocks in  $B^{init1}$  die Existenz von  $2^{r_{init2}-r_{init1}}$  benachbarten markierten Blöcken in  $B^{init2}$  impliziert. Die Umkehrung dieser Aussage liefert,

$$longestChain(B_{marked}^{init2}) < 2^{r_{init2}-r_{init1}} \Rightarrow B_{marked}^{init1} = \emptyset, \quad (3.7)$$

wobei  $LongestChain(B_{marked}^{init2})$  die längste Kette von benachbarten markierten Blöcken in  $B_{marked}^{init2}$  bezeichnet. Die Technik der dynamischen Startrunde greift diese Beziehung auf, indem initial die Runde  $r_{init} = \log_2|S|/2$  auf die gesamte Eingabe  $S$  angewendet wird. Im Anschluss kann der Wert  $longestChain(B_{marked}^{init})$  mithilfe eines Scans über die markierten Blöcke bestimmt werden. Der errechnete Wert impliziert eine Runde  $r_{Start}$  derart, dass vorherige Runden garantiert keine markierten Blöcke erzeugen und damit ausgelassen werden können. Der Wert  $r_{Start}$  ergibt sich wie folgt,

$$r_{Start} = r_{init} - \begin{cases} -1, & \text{falls } longestChain(B_{marked}^{init}) = 0 \\ \lfloor \log_2 longestChain(B_{marked}^{init}) \rfloor, & \text{sonst} \end{cases} \quad (3.8)$$

In Abhängigkeit von der Beschaffenheit der Eingabe, können maximal die Hälfte aller Runden ausgelassen werden, ohne eine Veränderung der Ergebnisse zu verursachen. In Runde  $r_{init} = \log_2|S|/2$  werden jedoch  $2^{\log_2|S|/2} = \sqrt{|S|}$  Blöcke erzeugt. Dies führt zu einer weiteren unteren Schranke für den Speicheraufwand des Algorithmus. Falls diese Technik angewandt wird, kann der Speicheraufwand mit  $O(\max\{\sqrt{n}, z\})$  abgeschätzt werden.

### 3.4.3 Vorberechnete Runde(PreMatching) - Laufzeit vs. Speicher

Analog zu der dynamischen Startrunde kann eine vorberechnete Runde ebenfalls genutzt werden, um den Arbeitsaufwand vorheriger Runden zu reduzieren. Sei  $r_{prematch} \leq \log_2|S|$  eine Runde, die auf die gesamte Eingabe angewendet wird. Als Ergebnis erhalten wir die Menge der markierten Blöcke  $B_{prematch}^{marked}$ . Weiterhin speichern wir uns den RFP aller Blöcke, die im Rahmen der Runde erzeugt werden. Wie in 3.4.2 gezeigt, kann jedem Block in einer vorherigen Runde einer Gruppe von Blöcken in einer späteren Runde zugeordnet werden, die die gleiche Zeichenfolge repräsentieren. Die Konkatenation von Zeichenfolgen kann entsprechend 2.5 in eine konstante Operation auf der Basis des RFP übersetzt werden. Gegeben sei eine Runde  $r_m$  mit  $1 \leq r_m \leq \log_2|S|$ . Für einen beliebigen Block  $b \in B_m$  können  $2^{r_{prematch}-r_m}$  viele Böcke  $(b_1, b_2, \dots, b_{2^{r_{prematch}-r_m}}) \in B_{prematch}$  gefunden werden, die die gleiche Zeichenfolge repräsentieren. So ergibt sich für den zugehörigen RFP,

$$RFP(b) = RFP(b_1) \oplus RFP(b_2) \oplus \dots \oplus RFP(b_{2^{r_{prematch}-r_m}}) \quad (3.9)$$

, wobei jede Operation in konstanter Zeit durchgeführt werden kann. Die Anzahl der Rechenschritte für die Berechnung des RFP eines Blockes hängt nun nicht mehr von der Länge der repräsentierten Zeichenfolge ab, sondern Rundendistanz zur vorberechneten Runde. Weiterhin kann die Menge der unmarkierten Blöcke  $B_{prematch}^{unmarked} = B_{prematch} \setminus B_{prematch}^{marked}$  genutzt werden, um die Menge der Blöcke  $B_m$  in Runde  $r_m$  zu reduzieren. Ein Block  $b \in B_m$  kann nur dann markiert werden, wenn die äquivalente Sequenz von Blocken  $(b_1, b_2, \dots, b_{2^{r_{prematch}-r_m}}) \in B_{prematch}$  markiert ist. Die Umkehrung dieser Aussage liefert einen Filter für alle vorherigen Runden. Die zusätzlich gespeicherten Daten erhöhen den Speicherbedarf des Algorithmus. Falls die vorberechnete Runde auf den Wert  $k \in \mathbb{N}$  festgelegt wird, so kann der Speicherbedarf mit  $O(\max\{2^k, z\})$  abgeschätzt werden.

#### 3.4.4 Minimale Tabellengröße(ScanSkip) - Laufzeit vs. Qualität

ToDo: Grundlagen + ApproxLZ77 ausweiten

# Kapitel 4

## Praktische Evaluation

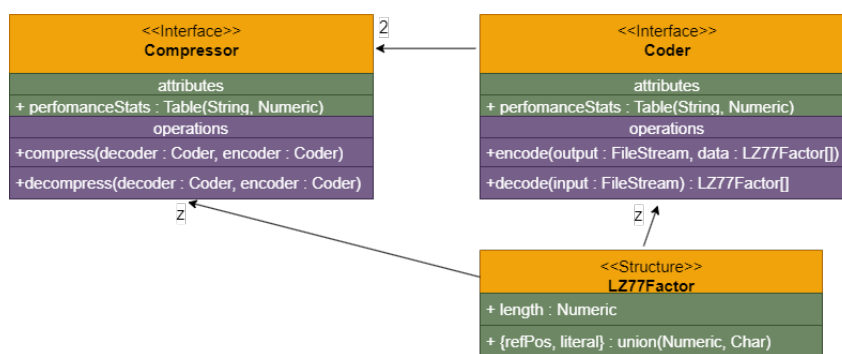
### 4.1 Testumgebung

Die folgenden Experimente wurden mithilfe einer AMD EPYC 7763 64-Core CPU mit 16 nutzbaren Hardwarethreads und 64GB Arbeitsspeicher durchgeführt. Das System verwendet Ubuntu 24.04 als Betriebssystem und GCC in der Version 13.2.0 als Compiler. Die Ausführung der Algorithmen mit einer spezifischen Anzahl von Threads wurde softwareseitig über OpenMP-Instruktionen realisiert.

### 4.2 Implementierung

#### 4.2.1 Klassenstruktur

Abbildung 4.1: Klassenstruktur der Implementierung



Die in 4.1 dargestellte Klassenstruktur illustriert die grundlegende Abstraktion, die für die Implementierung der Algorithmen verwendet wurde. Compressor und Coder beschreiben jeweils ein Interface bzw. ein Template, welches durch ein konkretes Kompressionsverfahren und einer Kodierung spezialisiert werden kann. Jegliche Spezialisierungen teilen sich jedoch eine gemeinsame Definition eines Faktors im LZ77-Schema.

### 4.2.2 Code

Die konkrete Implementierung erfolgte in der Programmiersprache C++ im C++20 - Standard und dem Build - System CMake in der Version 3.28. Die Parallelisierung wurde über Präprozessor - Instruktionen von OpenMP realisiert.

## 4.3 Messung

### 4.3.1 Eingabedaten

Die folgenden Algorithmen wurden auf verschiedenen Dateien aus dem Pizza & Chili-Corpus getestet. Die verwendeten Dateien decken verschiedene Kontexte und damit Kompressionspotentiale ab. In der Tabelle 4.2 sind die verwendeten Dateien aufgelistet. Die

**Abbildung 4.2:** Auflistung der verwendeten Eingabedaten

Datei	Größe	Alphabetgröße	Beschreibung
dna	200MB	4	DNA-Sequenzen
english	200MB	256	Englische Texte
proteins	200MB	20	Proteinsequenzen
sources	200MB	256	Quellcode
xml	200MB	256	XML-Dateien

Größe der Dateien wurde auf 100MB beschränkt, um einen angemessenen Rahmen für die Laufzeitmessung zu erhalten.

### 4.3.2 Messgrößen

#### Laufzeit

Die Laufzeit der Algorithmen wurde innerhalb der Ausführung gemessen. Dabei wird die Zeitmessung nach dem Laden der Eingabedatei gestartet und mit dem vollständigen Auffüllen der Faktorfolge beendet. Damit wird das Einlesen der Eingabe und eine eventuelle Kodierung der Ausgabe nicht in die Laufzeitmessung einbezogen. Diese Strategie hat ihren Hintergrund in der Tatsache, dass die konkrete Ausprägung des Eingabe- und Ausgabestroms keine Aussagekraft über die Qualität der Kompression hat.

#### Speicher

Der Speicherverbrauch der Algorithmen wurde auch intern mithilfe einer externen Bibliothek gemessen. Dabei wurden Speicherallokationen auf dem Heap überwacht und gemessen. Im Rahmen dieser Arbeit wurde die Spitze des allokierten Speichers im Zeitraum nach dem Einlesen der Eingabedatei und nach dem vollständigen Auffüllen der Faktorfolge gemessen.



**Kompressionsrate CR\***

Die Kompressionsrate wird neben der Anzahl der Faktoren zum Großteil von der verwendeten Kodierung bestimmt. Wie bereits erwähnt, sind wir in der Wahl der Kodierung nicht beschränkt, sodass die Aussagekraft bezüglich der Qualität der Kompression nicht eingeschränkt ist. Es ist jedoch zu beachten, dass Faktoren, die durch Approx. LZ77 erzeugt werden stets eine Zweierpotenz als Länge annehmen. Die binäre Repräsentation dieser Längen kann daher in Abhängigkeit von der gewählten Kodierung kompakt konstruiert werden. Um dieses Phänomen zu illustrieren gehen wir im Folgenden eine naive Kodierung vor, auf dessen Grundlage wir die Kompressionsrate  $CR^*$  definieren.

$$|K_{LZ77}(f)| = 1 + \begin{cases} 2 \log_2(n) & \text{falls } |f| > 1 \\ 9 & \text{sonst} \end{cases} \quad (4.1)$$

Im Falle von LZ77 bestimmt ein Bit, ob es sich um eine Referenz oder ein einzelnes Zeichen handelt. Im Falle einer Referenz wird die Länge und die Position der Referenz mithilfe von  $\log_2(n)$  Bits kodiert. Im Falle eines einzelnen Zeichens wird die ASCII-Kodierung des Zeichens verwendet, die 8 Bits benötigt.

$$|K_{Approx.LZ77}| = 1 + \begin{cases} \log_2(n) + \log_2(\log_2(n)) & \text{falls } |f| > 1 \\ 9 & \text{sonst} \end{cases} \quad (4.2)$$

Im Falle von Approx. LZ77 kann die Länge mithilfe von  $\log_2(\log_2(n))$  Bits kodiert werden, da die Länge anhand einer einzelnen Bitposition bestimmt wird.

### 4.3.3 Messwerte

**Tabelle 4.1:** Messgrößen der Algorithmen

Eingabe	Algorithmus	Laufzeit[s]	Speicher[Byte]	FR	CR*
proteins	LZ77	15.40	20.00	9.98%	70.83%
	Approx.LZ77	44.06	9.94	15.34%	63.95%
	Approx.LZ77Par	6.99	10.21	15.34%	63.95%
sources	LZ77	13.47	20.00	5.60%	38.75%
	Approx.LZ77	40.43	6.42	10.05%	40.14%
	Approx.LZ77Par	6.49	5.90	10.05%	40.14%
english	LZ77	15.36	20.00	6.70%	47.24%
	Approx.LZ77	51.00	7.06	10.42%	43.39%
	Approx.LZ77Par	7.46	6.16	10.42%	43.39%
dna	LZ77	14.89	20.00	6.66%	47.46%
	Approx.LZ77	30.10	8.38	10.71%	45.53%
	Approx.LZ77Par	4.80	6.66	10.71%	45.53%
xml	LZ77	13.02	20.00	3.42%	23.61%
	Approx.LZ77	29.38	3.46	6.62%	26.78%
	Approx.LZ77Par	4.91	3.46	6.62%	26.78%

**Laufzeiten**

**Speicherverbrauch**

**FR und CR\***

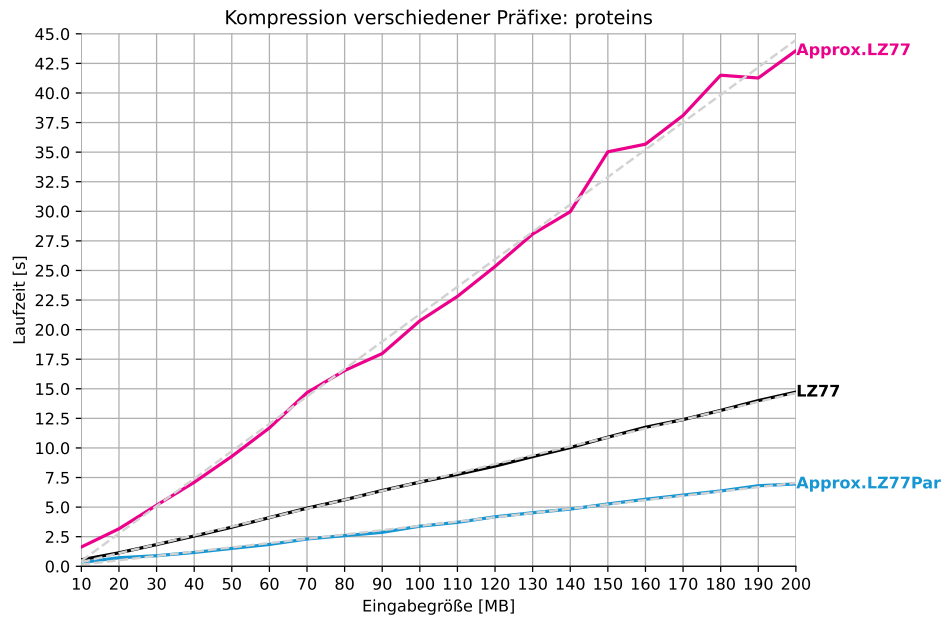
## 4.4 Auswertung

### 4.4.1 LZ77

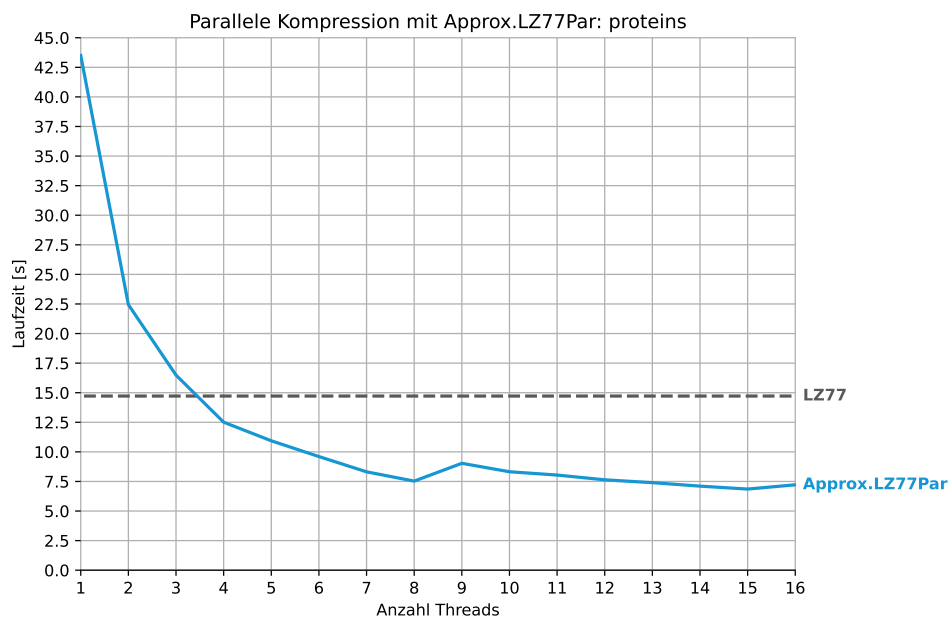
### 4.4.2 Approx. LZ77

### 4.4.3 Approx. LZ77Par

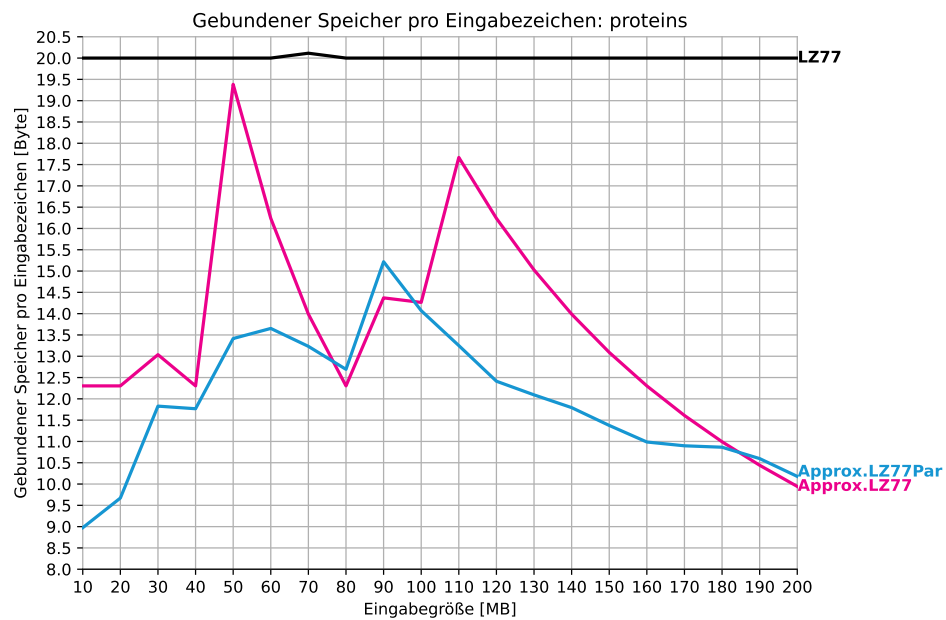
**Abbildung 4.3:** Laufzeitmessung von LZ77, Approx.LZ77 und Approx.LZ77Par(16 Threads) auf verschiedenen Präfixen von proteins. Als Vergleichsmaß wurde die lineare Regression der Kurven gestrichelt eingezeichnet.



**Abbildung 4.4:** Laufzeitmessung von Approx.LZ77Par mit verschiedener Anzahl an Threads für proteins



**Abbildung 4.5:** Speicherverbrauch von LZ77, Approx.LZ77 und Approx.LZ77Par(16 Threads) auf verschiedenen Präfixen von proteins. Aufgezeichnet wurde das Verhältnis von allokiertem Speicher zur Eingabegröße.

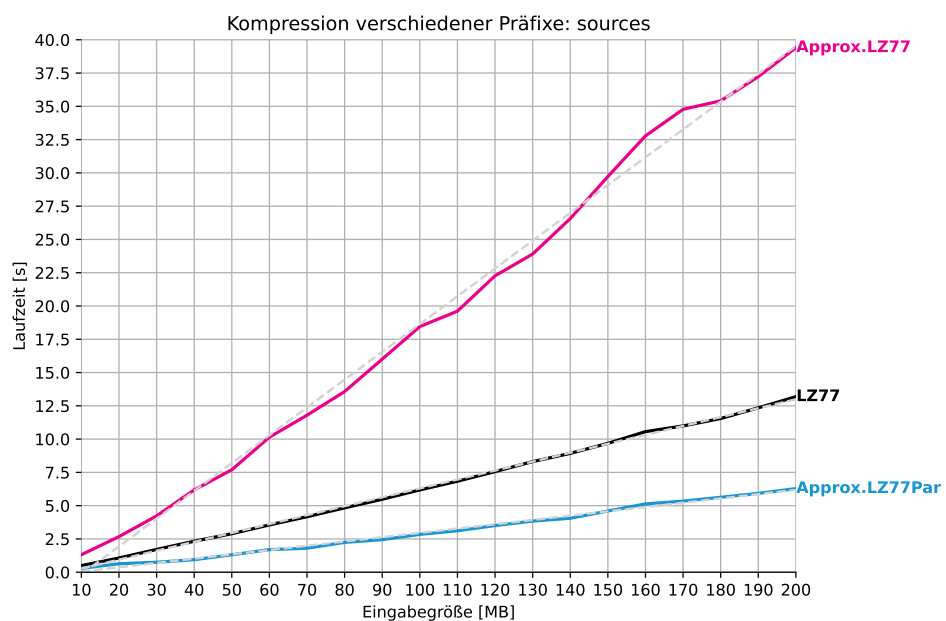


# Anhang A

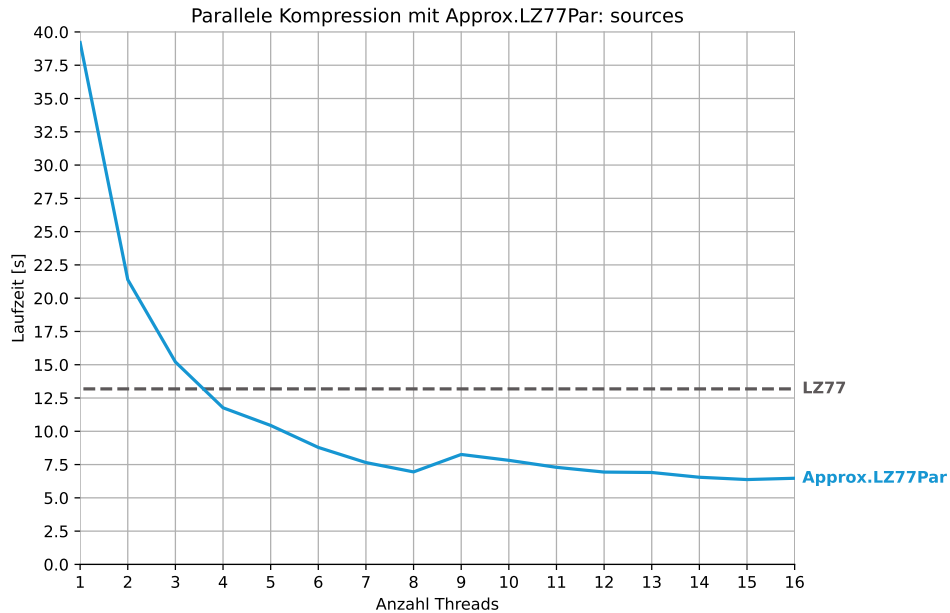
## Weitere Informationen

### A.1 Zusätzliche Messwerte

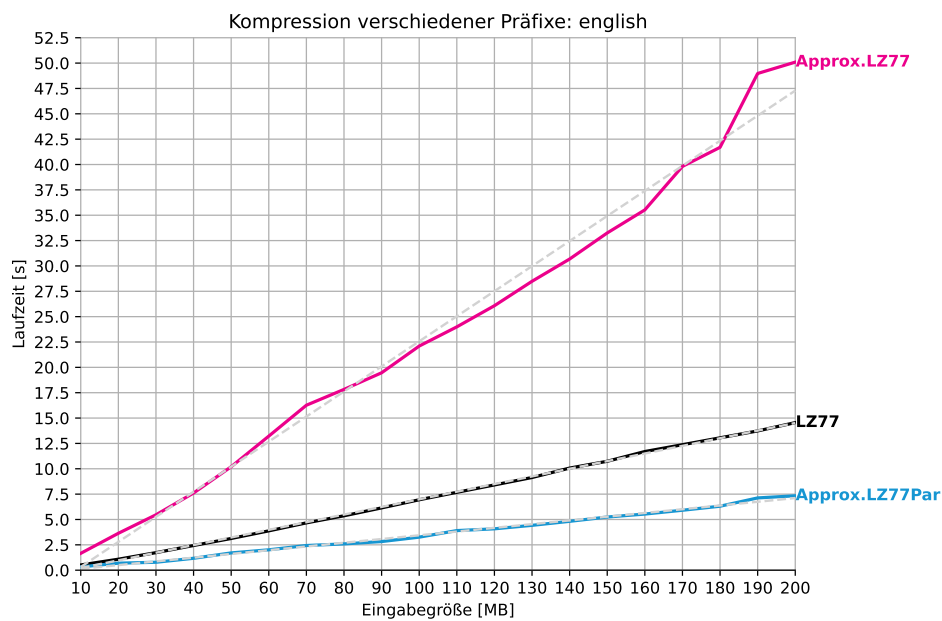
**Abbildung A.1:** Laufzeitmessung von LZ77, Approx.LZ77 und Approx.LZ77Par(16 Threads) auf verschiedenen Präfixen von sources. Als Vergleichsmaß wurde die lineare Regression der Kurven gestrichelt eingezeichnet.



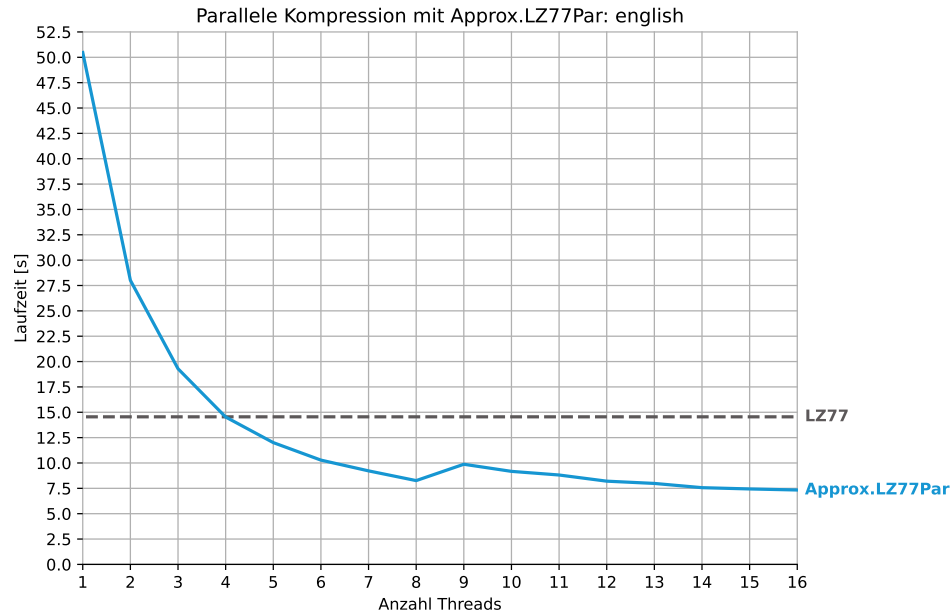
**Abbildung A.2:** Laufzeitmessung von Approx.LZ77Par mit verschiedener Anzahl an Threads für sources



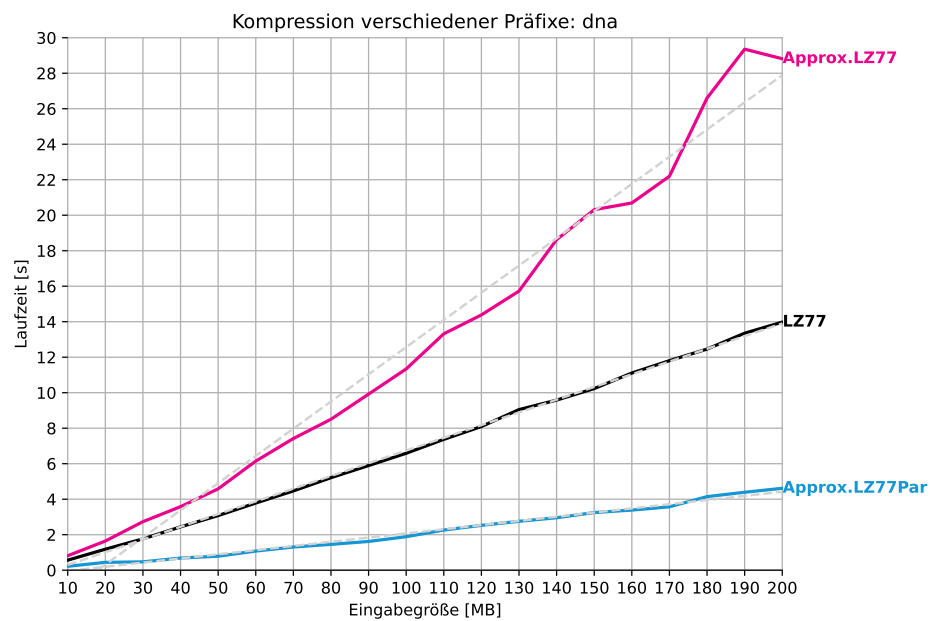
**Abbildung A.3:** Laufzeitmessung von LZ77, Approx.LZ77 und Approx.LZ77Par(16 Threads) auf verschiedenen Präfixen von english. Als Vergleichsmaß wurde die lineare Regression der Kurven gestrichelt eingezeichnet.



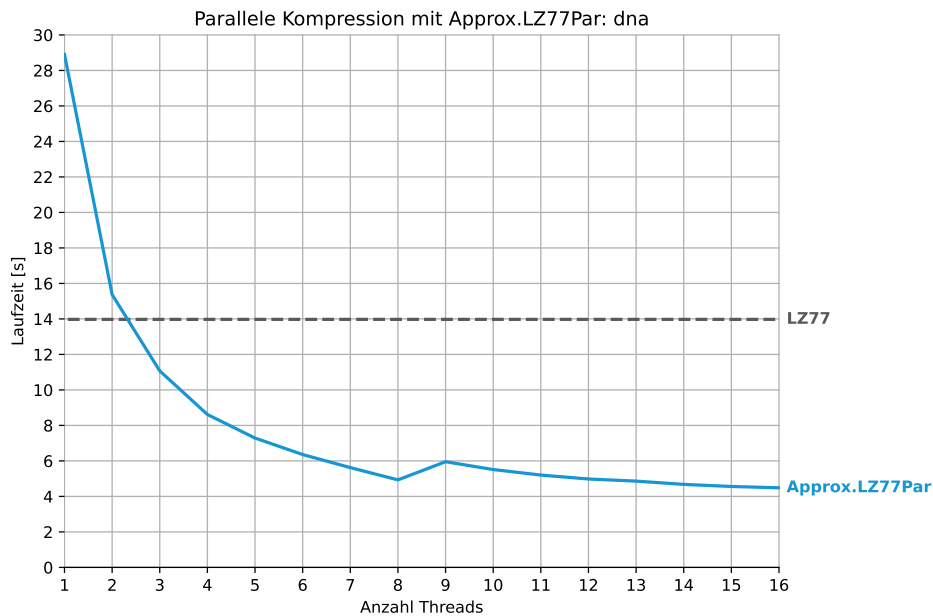
**Abbildung A.4:** Laufzeitmessung von Approx.LZ77Par mit verschiedener Anzahl an Threads für english



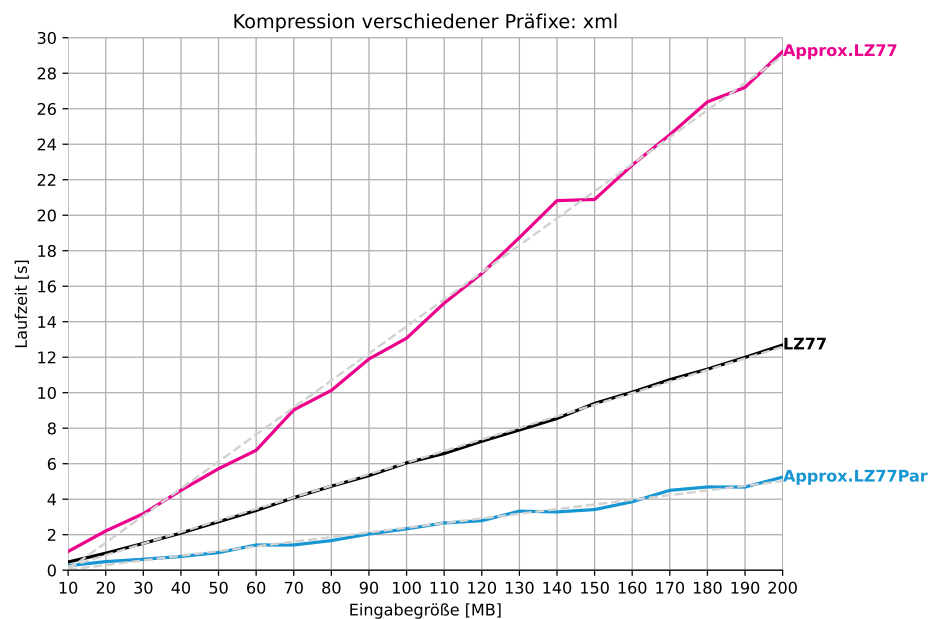
**Abbildung A.5:** Laufzeitmessung von LZ77, Approx.LZ77 und Approx.LZ77Par(16 Threads) auf verschiedenen Präfixen von dna. Als Vergleichsmaß wurde die lineare Regression der Kurven gestrichelt eingezeichnet.



**Abbildung A.6:** Laufzeitmessung von Approx.LZ77Par mit verschiedener Anzahl an Threads für dna

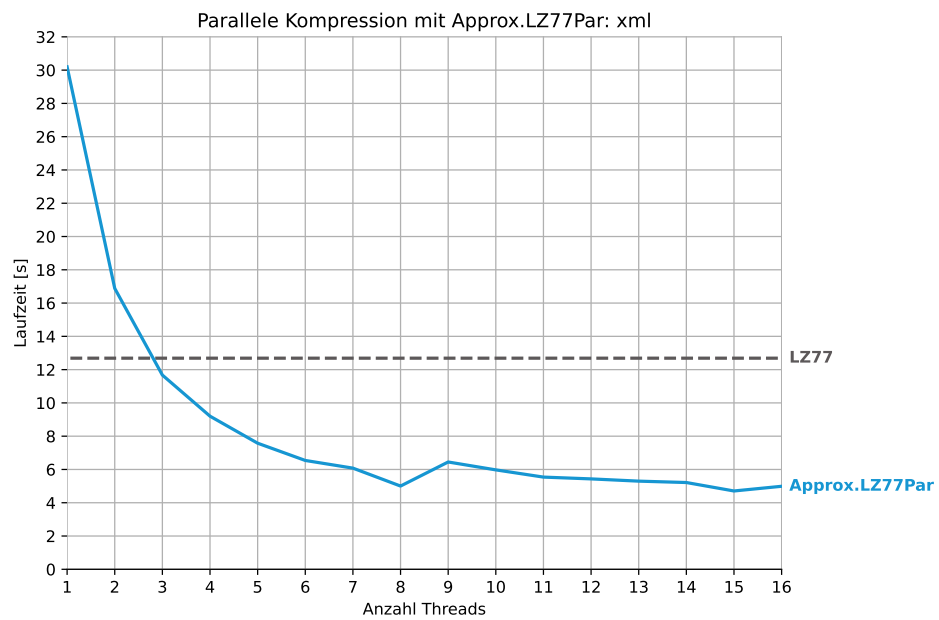


**Abbildung A.7:** Laufzeitmessung von LZ77, Approx.LZ77 und Approx.LZ77Par(16 Threads) auf verschiedenen Präfixen von xml. Als Vergleichsmaß wurde die lineare Regression der Kurven gestrichelt eingezeichnet.





**Abbildung A.8:** Laufzeitmessung von Approx.LZ77Par mit verschiedener Anzahl an Threads für xml





# Literaturverzeichnis

- [1] FISCHER, JOHANNES, TRAVIS GAGIE, PAWEŁ GAWRYCHOWSKI und TOMASZ KOCI-  
UMAKA: *Approximating LZ77 via small-space multiple-pattern matching*. In: *23rd Eu-  
ropean Symposium on Algorithms (ESA)*, Band 9294, Seiten 533–544. Springer, 2015.



Hiermit versichere ich, dass ich die vorliegende Arbeit selbstständig verfasst habe und keine anderen als die angegebenen Quellen und Hilfsmittel verwendet sowie Zitate kenntlich gemacht habe.

Dortmund, den 8. August 2024

Muster Mustermann

