



Universität Hamburg
Fakultät für Mathematik,
Informatik und Naturwissenschaften
Department Informatik

Bachelorarbeit

Speichereffiziente Methoden zur Repräsentation von paarweisen Sequenz-Alignments

Thorben Wiese

3wiese@informatik.uni-hamburg.de

Studiengang B.Sc. Informatik

Matr.-Nr. 6537204

Fachsemester 6

Erstgutachter Universität Hamburg:
Zweitgutachter Universität Hamburg:

Prof. Dr. Stefan Kurtz
Dr. Giorgio Gonnella

Inhaltsverzeichnis

1	Einleitung	1
2	Methoden	5
2.1	CIGAR-Strings	5
2.1.1	Kodierung eines CIGAR-Strings	5
2.2	Trace Point Konzept	8
2.2.1	Differenzen-Kodierung	10
2.2.2	Kodierung der Trace Point Differenzen	10
3	Resultate	13
3.1	Entropie der Repräsentationen	13
3.2	Testläufe CIGAR Kodierung	14
3.3	Testläufe Differenzen Kodierung	17
4	Diskussion	19
4.1	Bewertung CIGAR-Kodierung	19
4.2	Bewertung Kodierung der Differenzen der Trace Points	19
5	Programm	21
5.1	Aufbau	21
5.2	Funktionalität	22
5.2.1	Informationsverlust bei der encode()-Funktion	23
6	Fazit	25
	Literaturverzeichnis	27

Abbildungsverzeichnis

2.1	Huffman-Bäume der Kodierung des CIGAR-Strings	8
2.2	Huffman-Baum der Delta-Kodierung	12
3.1	Entropie des CIGAR-Strings	13
3.2	Entropie der Trace Point Differenzen	14
3.3	Größe der naiven binären Kodierung für einen CIGAR-String eines paarweisen Sequenz-Alignments in Bit	15
3.4	Größe der unären Kodierung für einen CIGAR-String eines paar- weisen Sequenz-Alignments in Bit	15
3.5	Größe der Huffman-Kodierung für einen CIGAR-String eines paarweisen Sequenz-Alignments in Bit	16
3.6	Größe der unären Kodierung für die Differenzen der Trace Points eines paarweisen Sequenz-Alignments in Bit	17
3.7	Größe der Huffman-Kodierung für die Differenzen der Trace Points eines paarweisen Sequenz-Alignments in Bit	18
5.1	UML-Diagramm	21

Tabellenverzeichnis

2.1	Relative Wahrscheinlichkeiten CIGAR-String	6
2.2	Unäre Kodierung des CIGAR-Strings	7
2.3	Huffmann-Kodierung des CIGAR-Strings	8
2.4	Relative Wahrscheinlichkeiten der Delta-Kodierung	11
2.5	Unäre Kodierung der Delta-Kodierung	11
2.6	Huffman-Kodierung der Delta-Kodierung	12

1 Einleitung

Ein Sequenzalignment wird in der Bioinformatik dazu verwendet, zwei oder mehrere Sequenzen von zum Beispiel DNA-Strängen oder Proteinsequenzen miteinander zu vergleichen und die Verwandtschaft zu bestimmen. Ein Alignment ist das Ergebnis eines solchen Vergleichs. Bei einem globalen Alignment wird jeweils die gesamte Sequenz betrachtet, bei einem lokalen Alignment lediglich Teilabschnitte der beiden Sequenzen. Um die verschiedenen Sequenzen vergleichen zu können, berechnet man einen Score oder die Kosten, um den Aufwand, den man betreiben muss, um die gegebene Sequenz in die Zielsequenz umzuwandeln, beschreiben zu können. Hierbei wird jeweils das Optimum, also entweder der maximale Score oder die minimalen Kosten gesucht. Die verschiedenen Schritte, um die Symbole der Strings zu verändern, sind bei Gleichheit ein 'match', bei der Substitution ein 'mismatch', bei der Löschung eine 'deletion' und bei der Einfügung eine 'insertion', welche je nach Verfahren unterschiedlich gewichtet werden können. Hierbei haben ähnliche Sequenzen einen hohen Score und geringe Kosten und unterschiedliche Sequenzen analog einen kleinen Score und hohe Kosten.

Ziel dieser Bachelorarbeit ist es, eine speichereffiziente Repräsentation von paarweisen Sequenzalignments zu implementieren und die Funktionsweise, sowie Vergleiche zu anderen Verfahren zu diskutieren.

Die Edit-Operationen

Die in diesem Kapitel eingeführten Begriffe werden in [Kurtz, S. 5-7, 14-16] definiert.

Sei \mathcal{A} eine endliche Menge von Buchstaben, die man Alphabet nennt. Für DNA-Sequenzen verwendet man üblicherweise die Menge der Basen, also $\mathcal{A} = \{a, c, g, t\}$. \mathcal{A}^i sei die Menge der Sequenzen der Länge i aus \mathcal{A} und ε sei die leere Sequenz. Formal ausgedrückt ist eine Edit-Operation ein Tupel

$$(\alpha, \beta) \in (\mathcal{A}^1 \cup \{\varepsilon\}) \times (\mathcal{A}^1 \cup \{\varepsilon\}) \setminus \{(\varepsilon, \varepsilon)\}.$$

Eine äquivalente Schreibweise von (α, β) ist $\alpha \rightarrow \beta$. Es gibt drei verschiedene Edit-Operationen

$a \rightarrow \varepsilon$ ist eine Deletion für alle $a \in \mathcal{A}$

$\varepsilon \rightarrow b$ ist eine Insertion für alle $b \in \mathcal{A}$

$a \rightarrow b$ ist eine Substitution für alle $a, b \in \mathcal{A}$

Dabei ist zu beachten, dass $\varepsilon \rightarrow \varepsilon$ keine Edit-Operation darstellt.

Ein Alignment von zwei Sequenzen u und v lässt sich nun als eine Sequenz $(\alpha_1 \rightarrow \beta_1, \dots, \alpha_h \rightarrow \beta_h)$ von Edit-Operationen definieren, sodass $u = \alpha_1 \dots \alpha_h$ und $v = \beta_1 \dots \beta_h$ gilt.

Die Edit-Distanz

Sei eine Kostenfunktion δ mit $\delta(a \rightarrow b) \geq 0$ für alle Substitutionen $a \rightarrow b$ und $\delta(\alpha \rightarrow \beta) > 0$ für alle Einfügungen und Löschungen $\alpha \rightarrow \beta$ gegeben. Die Kosten für ein Alignment $A = (\alpha_1 \rightarrow \beta_1, \dots, \alpha_h \rightarrow \beta_h)$ ist die Summe der Kosten aller Edit-Operationen des Alignments.

$$\delta(A) = \sum_{i=1}^h \delta(\alpha_i \rightarrow \beta_i)$$

Ein Beispiel einer Kostenfunktion ist die Einheitskostenfunktion

$$\delta(\alpha \rightarrow \beta) = \begin{cases} 0, & \text{wenn } \alpha, \beta \in \mathcal{A} \text{ und } \alpha = \beta \\ 1, & \text{sonst.} \end{cases}$$

Die Edit-Distanz von zwei Sequenzen ist wie folgt definiert:

$$\text{edist}_\delta(u, v) = \min\{\delta(A) \mid A \text{ ist Alignment von } u \text{ und } v\}$$

Ein Alignment A ist optimal, wenn $\delta(A) = \text{edist}_\delta(u, v)$ gilt.

Wenn δ die Einheitskostenfunktion ist, so ist $\text{edist}_\delta(u, v)$ die Levenshtein Distanz [Kurtz, S. 19-21].

Ein Alignment kann für eine Edit-Distanz e mit der Einheitskostenfunktion in

$O(e)$ Zeit berechnet werden [Kurtz, S. 41-42].

2 Methoden

2.1 CIGAR-Strings

Ein Dateiformat, welches zur Speicherung von Alignments verwendet wird, ist das SAM-Format oder die binär komprimierte Version BAM. Dieses codiert ein Alignment in einem sogenannten CIGAR-String, der aus einzelnen Zeichen besteht, die jeweils eine Edit-Operation bezeichnen, also M für eine Substitution, I für eine Insertion und D für eine Deletion. Gleiche aufeinanderfolgende Operationen werden als Kombination von Quantität und Symbol geschrieben.

Beispiel 1. Sei $u = \text{actgaact}$, $v = \text{actagaat}$ und das Alignment $A = (a \rightarrow a, c \rightarrow c, t \rightarrow t, \dots)$ gegeben.

a	c	t	-	g	a	a	c	t
a	c	t	a	g	a	a	-	t

Ein Alignment wird üblicherweise in drei Zeilen geschrieben, wobei in der ersten Zeile die Sequenz u und in der dritten Zeile die Sequenz v geschrieben wird. In der mittleren Zeile symbolisiert das Zeichen '|' eine Substitution, wobei üblicherweise nur ein Match markiert wird. Außerdem wird ein ε aus der Edit-Operation in diesem Fall durch das Zeichen '-' dargestellt.

Dieses Alignment wird durch den CIGAR-String `3M1I3M1D1M` repräsentiert [The SAM/BAM Format Specification Group 2015].

2.1.1 Kodierung eines CIGAR-Strings

Beispiel 2. Sei das Alignment A

0	5	0	5	0	5	0	5	0
g	a	g	c	-	a	-	t	-
g	a	g	c	c	c	t	t	t
g	c	c	t	t	t	g	c	t
a	g	t	a	g	t	a	c	t
a	g	a	-	a	-	t	-	a
a	-	a	-	a	-	t	-	a

```

| | | | | | | | | | | | | | | |
gaccaagtag--g-cgtggacctt-gctcggg-ctgtaagaga
0      5      0      5      0      5      0      5      0

```

gegeben, welches durch den CIGAR-String 4M1I1M1I1M1I1M2D1M1D1M1I8M1D7M1D5M1I4M repräsentiert werden kann.

Im Folgenden vergleiche ich für diese Art der Alignment-Repräsentation die Verfahren der naiven binären Kodierung, der unären Kodierung und der Huffman-Kodierung.

Sei das Alphabet $\mathcal{A}_1 = \{M, I, D\}$, welches alle Symbole aus dem CIGAR-String enthält, das Alphabet $\mathcal{A}_2 = \{1, 2, 4, 5, 7, 8\}$, welches alle Zahlen aus dem CIGAR-String enthält, sowie die relativen Wahrscheinlichkeiten $p(c_1)$ jeden Symbols $c_1 \in \mathcal{A}_1$ und $p(c_2)$ jeden Symbols $c_2 \in \mathcal{A}_2$ gegeben.

c_1	Häufigkeit	$p(c_1)$	c_2	Häufigkeit	$p(c_2)$
M	10	$\frac{10}{38}$	1	13	$\frac{13}{38}$
I	5	$\frac{5}{38}$	4	2	$\frac{2}{38}$
D	4	$\frac{4}{38}$	2	1	$\frac{1}{38}$
			5	1	$\frac{1}{38}$
			7	1	$\frac{1}{38}$
			8	1	$\frac{1}{38}$

Tabelle 2.1: Relative Wahrscheinlichkeiten CIGAR-String

Bei einer naiven binären Kodierung wird jedes Symbol $c_i \in \mathcal{A}_i$ mit $\lceil \log_2 n \rceil$, $n = |\mathcal{A}|$ Bit kodiert, also $\lceil \log_2 3 \rceil + \lceil \log_2 6 \rceil = 2 + 3 = 5$ Bit pro Symbol. Insgesamt ergibt das somit $19 \cdot 2 + 19 \cdot 3 = 95$ Bit.

Die unäre Kodierung kodiert jedes Symbol nach der Häufigkeit des Auftretens im Alphabet mit $i - 1$ '1'-Bits, gefolgt von einem '0'-Bit, wobei i die Position des Symbols in einer nach der Häufigkeit absteigend sortierten Liste ist. Das am Häufigsten auftretende Symbol des Alphabets wird also mit '0', das zweithäufigste mit '10', das dritthäufigste mit '110' usw. kodiert [Moffat u. Turpin 2002, S. 29-30].

Der oben genannte CIGAR-String wird demnach wie folgt unär kodiert:

Insgesamt benötigt die unäre Kodierung also $32 + 35 = 67$ Bit.

Symbol	Kodierung	Anzahl Bits	Symbol	Kodierung	Anzahl Bits
M	0	$10 \cdot 1 = 10$	1	0	$13 \cdot 1 = 13$
D	10	$5 \cdot 2 = 10$	4	10	$2 \cdot 2 = 4$
I	110	$4 \cdot 3 = 12$	2	110	$1 \cdot 3 = 3$
			5	1110	$1 \cdot 4 = 4$
			7	11110	$1 \cdot 5 = 5$
			8	111110	$1 \cdot 6 = 6$
Gesamtanzahl:		32			35

Tabelle 2.2: Unäre Kodierung des CIGAR-Strings

Der durchschnittliche Bedarf für ein Symbol beträgt

$$\frac{67}{38} \approx 1.76 \frac{\text{Bit}}{\text{Symbol}}$$

Bei einer *minimalen* binären Kodierung, wie sie auch im Huffman-Algorithmus verwendet wird, werden die Längen der Codewörter anhand der relativen Wahrscheinlichkeit des Symbols im Alphabet angepasst. Somit lässt sich eine Kodierung ermöglichen, welche im Durchschnitt weniger Bit pro Symbol beansprucht [Moffat u. Turpin 2002, S. 53-57]. Eine sparsamere Variante des regulären Huffman-Algorithmus ist der Kanonische Huffman-Algorithmus, welcher im Gegensatz zu der ursprünglichen Variante eine eindeutige Menge von Codewörtern liefert und keinen vollständigen Huffman-Baum, sondern lediglich die Anzahl der Codewörter für jede vorhandene Codewortlänge, sowie die sortierten Symbole benötigt, um die Informationen zu dekodieren.

Der kanonische Huffman-Algorithmus würde bei dem oben genannten Beispiel des CIGAR-Strings nach [Moffat u. Turpin 2002, S. 54] die Symbole wie in Tabelle 2.3 beschrieben kodieren. Die Huffman-Bäume beider Alphabete sind in Abbildung 2.1 dargestellt. Für die Dekodierung sind somit zusätzlich die Listen (1, 2), (M,D,I) und (0, 2, 4), (1, 4, 2, 5, 7, 8) zu speichern.

Die Größe dieser Kodierung ist demnach $28 + 42 = 70$ Bit und der durchschnittliche Bedarf für ein Symbol beträgt

$$\frac{70}{38} \approx 1.84 \frac{\text{Bit}}{\text{Symbol}}.$$

Symbol	Kodierung	Anzahl Bits	Symbol	Kodierung	Anzahl Bits
M	0	$10 \cdot 1 = 10$	1	00	$13 \cdot 2 = 26$
D	10	$5 \cdot 2 = 10$	4	01	$2 \cdot 2 = 4$
I	11	$4 \cdot 2 = 8$	2	100	$1 \cdot 3 = 3$
			5	101	$1 \cdot 3 = 3$
			7	110	$1 \cdot 3 = 3$
			8	111	$1 \cdot 3 = 3$
Gesamtanzahl:		28			42

Tabelle 2.3: Huffman-Kodierung des CIGAR-Strings

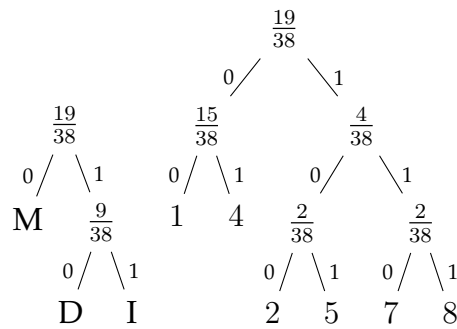


Abbildung 2.1: Huffman-Bäume der Kodierung des CIGAR-Strings

2.2 Trace Point Konzept

Ein neuer Ansatz der speichereffizienten Repräsentation von Alignments wurde von Gene Myers in [Myers 2015] beschrieben und basiert auf dem Konzept der Trace Points.

Sei A ein Alignment von $u[i...j]$ und $v[k...l]$ mit $i < j$ und $k < l$ und sei $\Delta \in \mathbb{N}$. Sei $p = \lceil \frac{i}{\Delta} \rceil$. Man unterteilt $u[i...j]$ in $\tau = \lceil \frac{j}{\Delta} \rceil - \lceil \frac{i}{\Delta} \rceil$ Substrings $u_0, u_1, \dots, u_{\tau-1}$ mit

$$u_q = \begin{cases} u[i...p \cdot \Delta] & \text{falls } q = 0 \\ u[(p + q - 1) \cdot \Delta + 1... (p + q) \cdot \Delta] & \text{falls } 0 < q < \tau - 1 \\ u[(p + \tau - 2) \cdot \Delta...j] & \text{falls } q = \tau - 1 \end{cases}$$

Für alle q mit $0 \leq q < \tau - 1$ sei t_q der letzte Index des Substrings von v , der in A mit u_q aligniert. t_q nennt man Trace Point. Für $q = 0$ aligniert u_0 mit $v_0 = v[k...t_0]$. Für alle q mit $0 < q < \tau - 1$ aligniert u_q mit $v_q = v[t_{q-1} + 1...t_q]$.

Seien i, j, k, ℓ, Δ und die Trace-Points eines Alignments von u und v gegeben. Dann kann ein Alignment A' von u und v mit $\delta(A') \leq \delta(A)$ konstruiert werden. Danach bestimmt man aus den Trace-Points die Substring-Paare u_q und v_q , berechnet hierfür ein optimales Alignment und konkateniert die Alignments von den aufeinanderfolgenden Substring-Paaren zu A' .

Beispiel 3.

Sequenz 1: gagcatgttgccctggctcctttgctaggtactgtagaga

Sequenz 2: gaccaagtaggcgtggaccttgctcggctctgtaagaga

Delta: 15

Gesamtalignment:

```

0      5      0      5      0      5      0      5      0
gagc-a-t-gttgcc-tggcctttgctaggtactgta-gaga
|| | | | | | | ||| |||| ||| ||| ||||| ||||
gaccaagtag--g-cgtggacctt-gctcggctctgtaagaga
0      5      0      5      0      5      0      5      0

```

seq1[0...14] aligniert mit seq2[0...15]

```

gagc-a-t-gttgcc-tgg
|| | | | | | | |||
gaccaagtag--g-cgtgg

```

seq1[15...29] aligniert mit seq2[16...28]

```

tcctttgctaggtac
|||| ||| ||| |
acctt-gctcggct-c

```

seq1[30...37] aligniert mit seq2[29...37]

```

tgta-gaga
|||| ||||
tgtaagaga

```

Trace Points: [15, 28]

2.2.1 Differenzen-Kodierung

Gegeben sei eine Liste $L = (a_1, a_2, \dots, a_n)$ mit $a_i < a_{i+1}, 0 < i \leq n$.

Anstatt jeden Wert $a \in L$ als solchen abzuspeichern, kann alternativ die Differenz eines Wertes a_i zu dem nachfolgenden Wert a_{i+1} abgespeichert werden. Lediglich der erste (oder letzte) Wert aus L wird benötigt, um später sukzessive die ursprüngliche Liste rekonstruieren zu können.

$$L_{diff} = (a_1, (a_2 - a_1), (a_3 - a_2), \dots, (a_n - a_{n-1}))$$

Bei gleichmäßig ansteigenden Werten ist die Abweichung der Differenzen zweier aufeinanderfolgender Werte in der Liste untereinander gering und die Menge der zu kodierenden Symbole verringert sich.

2.2.2 Kodierung der Trace Point Differenzen

Beispiel 4. Sei $\Delta = 5$ und das Alignment A

```

0      5      0      5      0      5      0      5      0
gagc-a-t-gttgcc-tggtcctttgctaggtactgta-gaga
|| | | | | | | ||| ||| ||| ||| ||| |||
gaccaagtag--g-cgtggacctt-gctcgggt-ctgtaagaga
0      5      0      5      0      5      0      5      0

```

wie in Abschnitt 2.1.1 mit den dazugehörigen TracePoints 5, 10, 15, 20, 24, 28 und 34 gegeben. Es ergibt sich somit das Alphabet $\mathcal{A} = \{5, 10, 15, 20, 24, 28, 34\}$ mit ausschließlich positiven und aufsteigenden Werten.

Für die Trace Point Darstellung ist somit eine Differenzen-Kodierung möglich. Als neue Liste zu kodierender Werte ergibt sich nach 2.2.1 $L_{diff} = (5, 5, 5, 5, 4, 4, 6)$.

Um aus den Trace Points ein neues Alignment rekonstruieren zu können, benötigt man zusätzlich mindestens den Δ -Wert, damit die Grenzen der Substrings beider Sequenzen berechnet werden können. Hierfür muss also der Δ -Wert zu L_{diff} hinzugefügt werden.

Für das oben genannten Beispiel ergibt sich somit

$$L_{diff} = (\Delta, a_1, (a_2 - a_1), (a_3 - a_2), \dots, (a_n - a_{n-1})) = (5, 5, 5, 5, 5, 4, 4, 6).$$

Sei das Alphabet $\mathcal{A} = \{4, 5, 6\}$ für L_{diff} , sowie die relativen Wahrscheinlichkeiten $p(c)$ jeden Symbols $c \in \mathcal{A}$ gegeben.

c	Häufigkeit	$p(c)$
5	5	$\frac{5}{8}$
4	2	$\frac{2}{8}$
6	1	$\frac{1}{8}$

Tabelle 2.4: Relative Wahrscheinlichkeiten der Delta-Kodierung

Bei der naiven binären Kodierung ergibt sich analog zu 2.1.1 ein Bedarf von $\lceil \log_2 8 \rceil = 3$ Bit pro Symbol, also $8 \cdot 3 = 24$ Bit insgesamt.

Die unäre Kodierung ergibt für dieses Beispiel die folgende Kodierung:

Symbol	Kodierung	Anzahl Bits
5	0	$5 \cdot 1 = 5$
4	10	$2 \cdot 2 = 4$
6	110	$1 \cdot 3 = 3$
Gesamtanzahl:		12

Tabelle 2.5: Unäre Kodierung der Delta-Kodierung

Die Größe dieser Kodierung ist demnach 12 Bit mit einem durchschnittlichen Bedarf pro Symbol von

$$\frac{12}{8} \approx 1.5 \frac{\text{Bit}}{\text{Symbol}}$$

Die Ausführung des Huffman-Algorithmus kodiert nach [Moffat u. Turpin 2002, S. 54] die Symbole wie in Tabelle 2.6 aufgelistet. Der dazugehörige Huffman-Baum aus 2.2 verdeutlicht die Kodierung der einzelnen Symbole, muss aber für den kanonischen Huffman-Algorithmus, wie in 2.1.1 beschrieben, nicht komplett gespeichert werden. Aufgrund der Beschaffenheit der Codewörter des kanonischen Huffman-Algorithmus ist hier lediglich die Speicherung der Listen $(1, 2)$, $(5, 4, 6)$ nötig.

Symbol	Kodierung	Anzahl Bits
5	0	$5 \cdot 1 = 5$
4	10	$2 \cdot 2 = 4$
6	11	$1 \cdot 2 = 2$
Gesamtanzahl:		11

Tabelle 2.6: Huffman-Kodierung der Delta-Kodierung

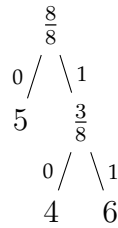


Abbildung 2.2: Huffman-Baum der Delta-Kodierung

Die Größe der Huffman-Kodierung ist demnach 11 Bit mit einem durchschnittlichen Bedarf pro Symbol von

$$\frac{11}{8} \approx 1.38 \frac{\text{Bit}}{\text{Symbol}}$$

3 Resultate

3.1 Entropie der Repräsentationen

Die folgenden Grafiken zeigen unabhängig vom Kodierungsverfahren die Entropie der CIGAR-Strings und Trace Point Differenzen von 10.000 zufällig generierten Sequenzpaaren mit je etwa 1.000 Basen, einer Fehlerrate von 15% und einem Δ -Wert von 100.

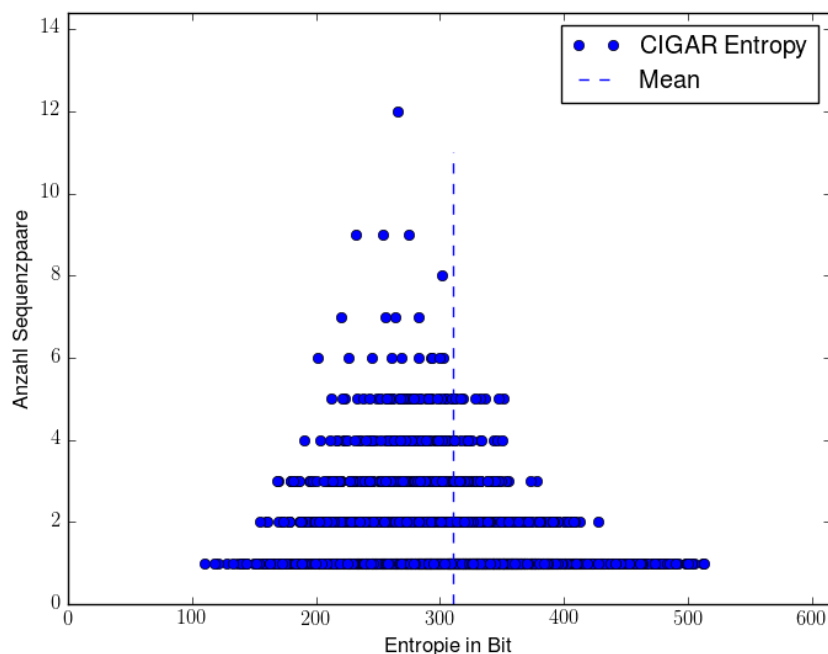


Abbildung 3.1: Entropie des CIGAR-Strings

Die Entropie für die CIGAR-String Repräsentation liegt für die gegebenen Sequenzpaare, wie in Abbildung 3.1 zu erkennen ist, zwischen 110 und 520 Bit, wobei der Durchschnittswert bei 310.17 Bit liegt. Es fällt auf, dass die meisten Werte nur einmal vorkommen und im gesamten Intervall zwischen 110 und 520 Bit auftreten. Für die weiteren Sequenzpaare, welche 2- bis 9-fach auftreten, wird das Intervall pyramidisch kleiner, bis das Maximum bei 12 Sequenzpaaren, welche mit etwa 280 Bit kodiert wurden, erreicht ist.

Die Differenzen der Trace Points weisen hingegen wie in Abbildung 3.2 verdeut-

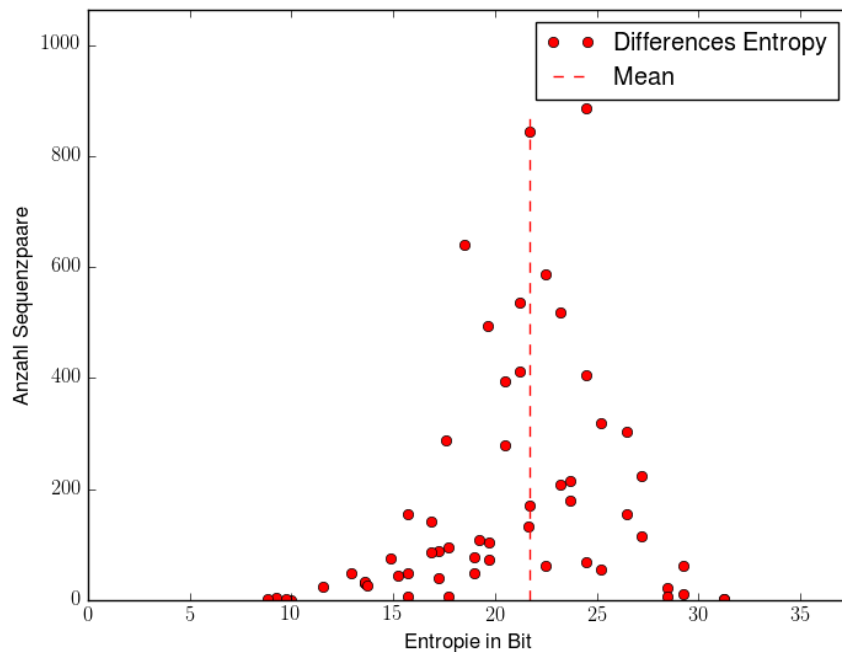


Abbildung 3.2: Entropie der Trace Point Differenzen

licht eine Entropie im Bereich von 8 bis 32 Bit auf, wobei der Durchschnittswert von 21.70 Bit nur etwa 7% des Durchschnittswertes der CIGAR-String Repräsentation ausmacht. Das Maximum der nahezu normalverteilten Entropie-Werte liegt hier bei 24.80 Bit, mit welchem etwa 850 Sequenzpaare kodiert wurden.

3.2 Testläufe CIGAR Kodierung

Die folgenden Grafiken wurden mit jeweils 10.000 zufällig generierte Sequenzpaaren mit je etwa 1.000 Basen, einer Fehlerrate von 15% und einem Δ -Wert von 100 berechnet und auf 10 Bit gerundet.

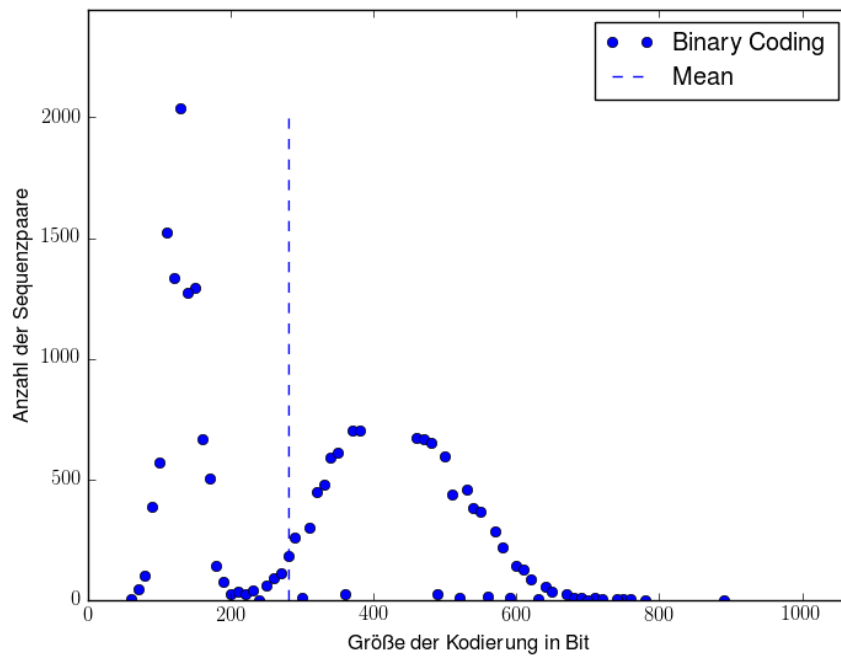


Abbildung 3.3: Größe der naiven binären Kodierung für einen CIGAR-String eines paarweisen Sequenz-Alignments in Bit

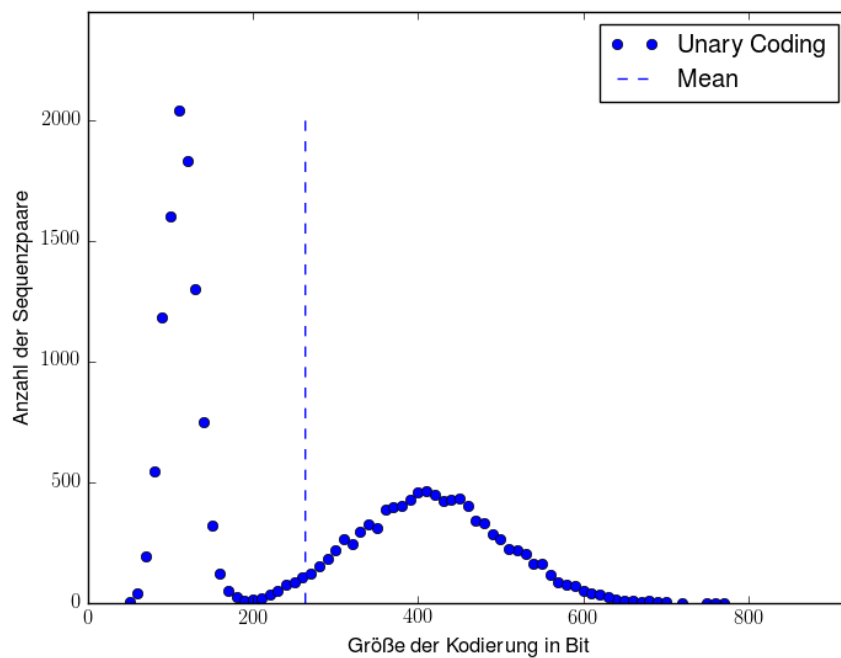


Abbildung 3.4: Größe der unären Kodierung für einen CIGAR-String eines paarweisen Sequenz-Alignments in Bit

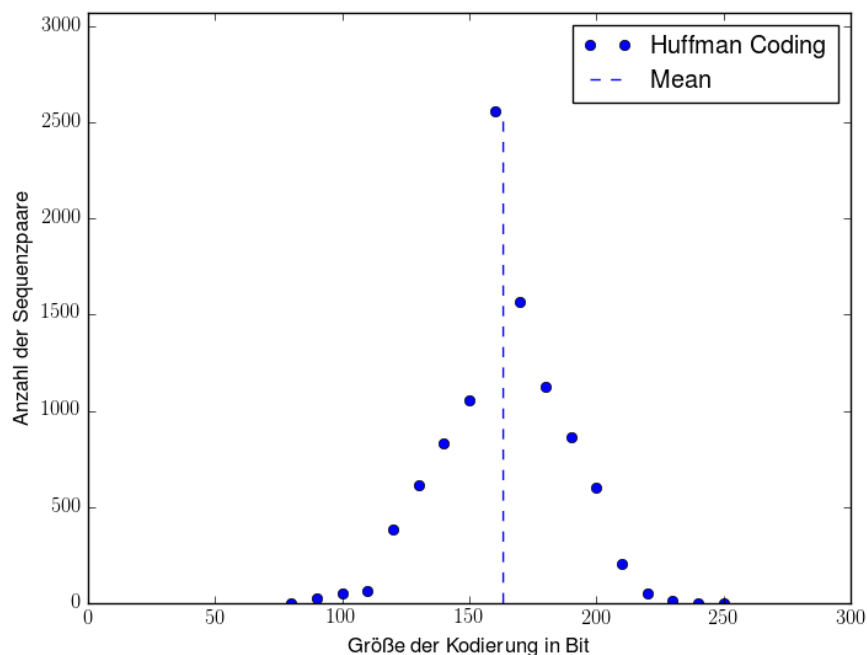


Abbildung 3.5: Größe der Huffman-Kodierung für einen CIGAR-String eines paarweisen Sequenz-Alignments in Bit

Die naive binäre Kodierung der CIGAR-Strings benötigt wie in Abbildung 3.3 dargestellt im Mittel 282.88 Bit, wobei Werte zwischen 50 und 900 Bit erreicht werden. Diese bilden zwei Maxima bei 2050 Sequenzpaaren mit jeweils etwa 120 Bit als globales Maximum und 700 Sequenzpaaren mit jeweils etwa 420 Bit als lokales Maximum.

Wie in Abbildung 3.4 zu erkennen ist, ähnelt die Anzahl der Bits für die unäre Kodierung der CIGAR-Strings der der naiven binären Kodierung. Hier werden ebenfalls zwei Maxima erreicht, wobei hier das globale Maximum bei 2050 Sequenzpaaren mit etwa 120 Bit und das lokale Maximum bei 500 Sequenzpaaren mit etwa 420 Bit liegt. Für die unäre Kodierung wird im Mittel 263.39 Bit für die Kodierung eines CIGAR-Strings und damit nur etwa 7% weniger als bei der naiven binären Kodierung benötigt.

Die Huffman-Kodierung kodiert wie in Abbildung 3.5 zu erkennen ist zwischen 70 und 250 Bit für die Kodierung der CIGAR-Strings. Die Bitanzahl ist hier nahezu normalverteilt mit einem Maximum von 2550 Sequenzpaaren für etwa 160 Bit, was in etwa dem Durchschnitt von 162.79 Bit entspricht. Die Huffman-Kodierung benötigt somit für die Kodierung im Schnitt nur 61.81% der Bitanzahl der unären

und 57.55% der naiven binären Kodierung.

3.3 Testläufe Differenzen Kodierung

Die folgenden Grafiken wurden mit jeweils 10.000 zufällig generierte Sequenzpaaren mit je etwa 1.000 Basen, einer Fehlerrate von 15% und einem Δ -Wert von 100 berechnet und auf 3 Bit gerundet.

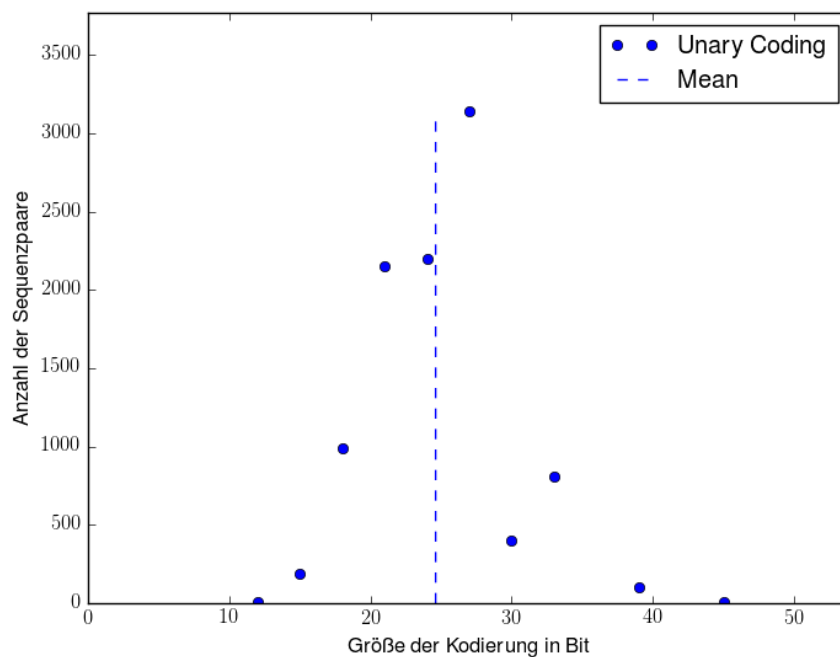


Abbildung 3.6: Größe der unären Kodierung für die Differenzen der Trace Points eines paarweisen Sequenz-Alignments in Bit

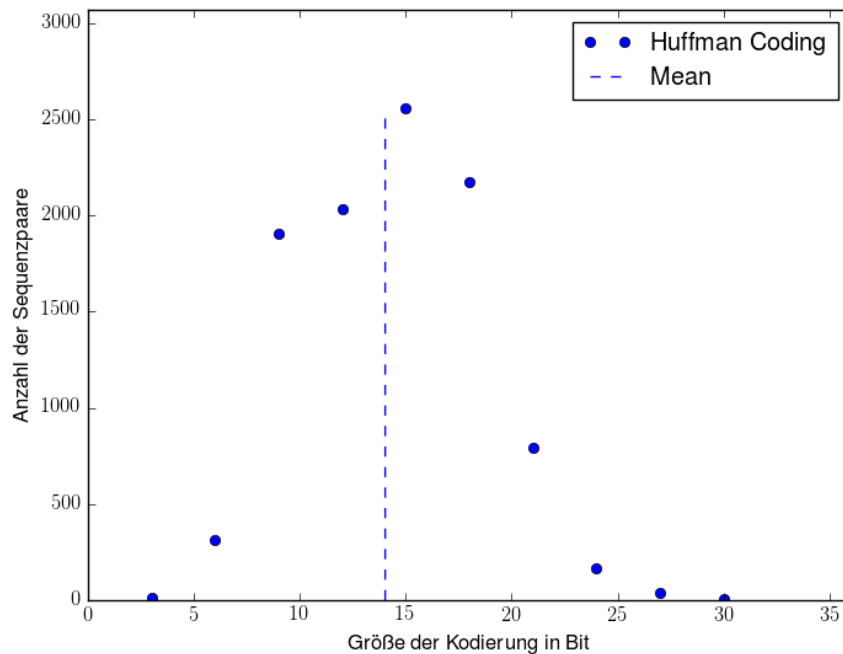


Abbildung 3.7: Größe der Huffman-Kodierung für die Differenzen der Trace Points eines paarweisen Sequenz-Alignments in Bit

Die naive binäre Kodierung benötigt in jedem Durchlauf und somit auch im Mittel konstant 40 Bit, da für jeden Durchlauf 9 Trace Points und der Δ -Wert gespeichert werden für $\lceil \log_2 10 \rceil \cdot 10 = 40$ Bit.

Für die unäre Kodierung der Differenzen der Trace Points wird, wie in Abbildung 3.6 verdeutlicht, zwischen 12 und 45 Bit benötigt. Die Bitanzahl ist in etwa normalverteilt, wobei das Maximum bei 3100 Sequenzpaaren mit 26 Bit liegt. Im Mittel wird für diese Kodierung 24.62 Bit und damit nur etwa 61.54% der binären Kodierung benötigt.

Die Huffman-Kodierung kodiert, wie in Abbildung 3.7 zu sehen ist, die Differenzen der Trace Points mit 4 bis 30 Bit, wobei hier das Maximum bei 2550 Sequenzpaaren und 15 Bit liegt. Sie verbraucht durchschnittlich nur 14.03 Bit und damit nur etwa 56.99% des Speicherbedarfs der unären und 35.10% der naiven binären Kodierung.

4 Diskussion

4.1 Bewertung CIGAR-Kodierung

4.2 Bewertung Kodierung der Differenzen der Trace Points

5 Programm

5.1 Aufbau

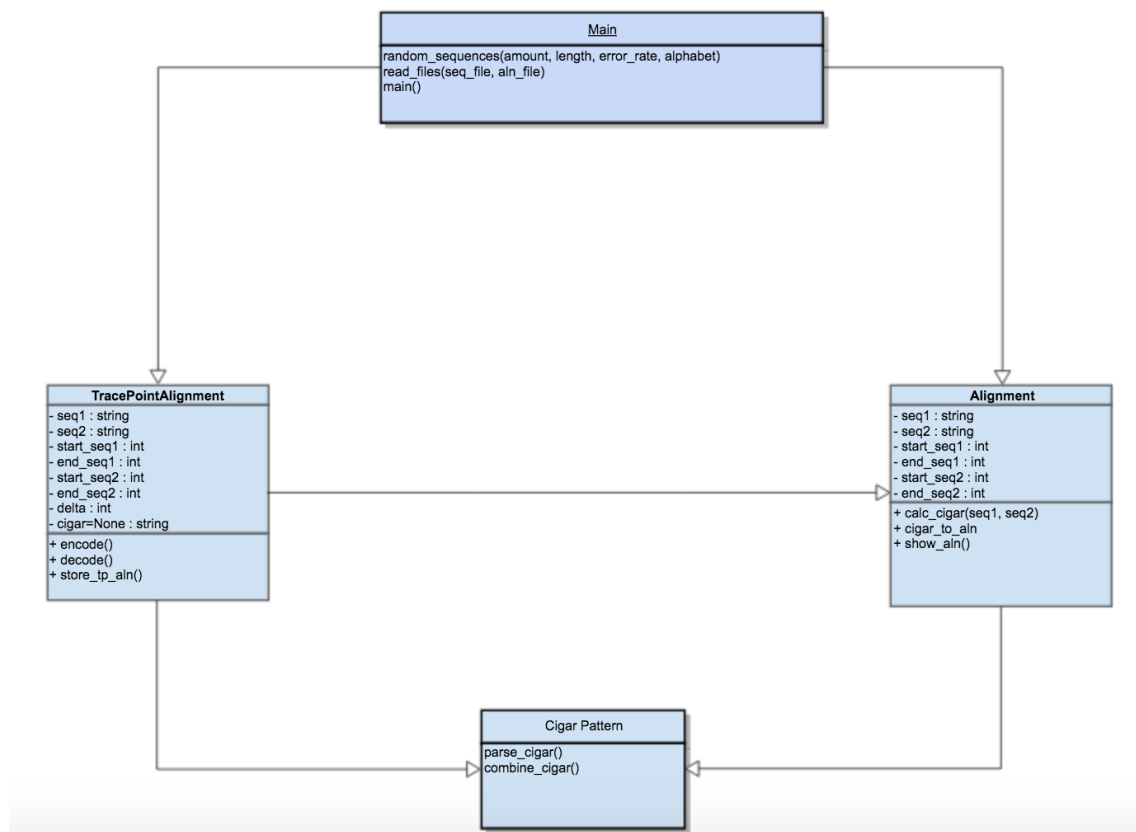


Abbildung 5.1: UML-Diagramm

5.2 Funktionalität

Algorithm 1 Computation of Trace Points from a given CIGAR-String

Input: $seq1, seq2, start_seq1, end_seq1, start_seq2, \Delta, cigar$ mit
 $|seq1|, |seq2|, |cigar| > 0$;
 $start_seq1, start_seq2 \geq 0$;
 $start_seq1 < end_seq1$ und
 $\Delta > 0$

Output: Array TP of Trace Points

```

1: function encode(seq1, seq2, start_seq1, end_seq1, start_seq2, Δ, cigar)
2:    $p \leftarrow MAX(1, \lceil start\_seq1/\Delta \rceil)$ 
3:    $\tau \leftarrow \lceil end\_seq1/\Delta \rceil - \lfloor start\_seq2/\Delta \rfloor$ 
4:    $uTP \leftarrow$  Array for interval termini in the first sequence
5:   for  $i \leftarrow 0$  upto  $|\tau|$  do
6:      $uTP[i] \leftarrow (p + i) \cdot (\Delta - 1)$ 
7:   end for
8:    $\#uChars, \#vChars, count \leftarrow 0$ 
9:    $TP \leftarrow$  Array for Trace Points
10:  for each ( $cig\_count, cig\_symbol$ ) in  $cigar$  do
11:    for  $i \leftarrow 0$  upto  $cig\_count$  do
12:      if  $cig\_symbol = 'T'$  then
13:        increment  $\#uChars$ 
14:      else if  $cig\_symbol = 'D'$  then
15:        increment  $\#vChars$ 
16:      else
17:        increment  $\#uChars, \#vChars$ 
18:      end if
19:      if  $\#uChars = uTP[count]$  then
20:        append ( $\#vChars$  to  $TP$ )
21:      end if
22:      if  $count \neq |uTP| - 1$  then
23:        return  $TP$ 
24:      else
25:        increment  $count$ 
26:      end if
27:    end for
28:  end for
29: end function

```

5.2.1 Informationsverlust bei der encode()-Funktion

Die encode-Funktion extrahiert aus dem gegebenen CIGAR-String die Trace Points, welche dann zusammen mit dem Δ -Wert und den Start- und Endpositionen der Sequenzabschnitte gespeichert werden. Hierbei geht die Information, wie die jeweiligen Intervalle zwischen den Trace Points zu den komplementären Intervallen in der Ursprungssequenz aligniert werden, verloren. Für die Rückgewinnung dieser Information muss in der decode()-Funktion zunächst ein neues Alignment des jeweiligen Intervall-Paares errechnet werden und alle Teilalignments zu einem Gesamtalignment konkatiniert werden.

Algorithm 2 Computation of a CIGAR-String from a given Trace Point Array

Input: $seq1, seq2, \Delta, TP$ mit
 $|seq1|, |seq2|, \Delta, |TP| > 0$

Output: CIGAR-String

```

1: function decode( $seq1, seq2, \Delta, TP$ )
2:    $cig \leftarrow$  empty String
3:   for  $i \leftarrow 0$  upto  $|TP|$  do
4:     if  $i = 0$  then
5:       append cigar( $seq1[0 \dots \Delta], seq2[0 \dots TP[i] + 1]$ ) to  $cig$ 
6:     else if  $i = |TP| - 1$  then
7:       append cigar( $seq1[i \cdot \Delta \dots |seq1|], seq2[TP[i - 1] + 1 \dots |seq2|]$ ) to  $cig$ 
8:     else
9:       append cigar ( $seq1[i \cdot \Delta \dots (i + 1) \cdot \Delta],$ 
10:                     $seq2[TP[i - 1] + 1 \dots TP[i] + 1]$ ) to  $cig$ 
11:   end if
12: end for
13:    $cig \leftarrow$  combine( $cig$ )
14:   return  $cig$ 
15: end function
16: function combine( $cigar$ )
17:    $cig \leftarrow$  empty String
18:    $tmp \leftarrow 0$ 
19:   for each  $cig\_count, cig\_symbol$  in  $cigar$  do
20:      $tmp \leftarrow tmp + previous\_cig\_count$ 
21:     if  $cig\_symbol = previous\_cig\_symbol$  then
22:       if not last element in  $cigar$  then
23:          $tmp \leftarrow 0$ 
24:       end if
25:     end if
26:     if last element in  $cigar$  then
27:       append  $tmp + cig\_count, cig\_symbol$  to  $cig$ 
28:     end if
29:   end for
30:   return  $cig$ 
31: end function

```

6 Fazit

Je größer der vorher definierte positive Parameter Δ ist, desto weniger Trace Points werden gespeichert und umso länger dauert die Berechnung, um die Teil-Alignments zu rekonstruieren. Bei einem kleinen Δ werden analog mehr Trace Points gespeichert, aber die Rekonstruktionszeit der Teil-Alignments ist geringer.

Mithilfe von Δ lässt sich somit ein Trade-Off zwischen dem Speicherplatzverbrauch und dem Zeitbedarf für die Rekonstruktion der Teil-Alignments einstellen.

Literaturverzeichnis

[Kurtz] KURTZ, Stefan: *Foundations of Sequence Analysis*. – Lecture notes for a course in the Wintersemester 2015/2016

[Moffat u. Turpin 2002] MOFFAT, Alistair ; TURPIN, Andrew: *Compression and Coding Algorithms*. Kluwer Academic Publishers, 2002

[Myers 2015] MYERS, Eugene: *Recording Alignments with Trace Points*. <https://dazzlerblog.wordpress.com/2015/11/05/trace-points/>.
Version: November 2015

[The SAM/BAM Format Specification Group 2015] THE SAM/BAM FORMAT SPECIFICATION GROUP: *Sequence Alignment/Map Format Specification*. <https://samtools.github.io/hts-specs/SAMv1.pdf>. Version: November 2015
