



BESCHLEUNIGUNG DES SABA-ALIGNMENTS

Diplomvortrag

Frank Hoffmann

Dresden, 27.08.2014



Problemstellung

Parallelisierung

Auswertung/Performance

Zusammenfassung

Problemstellung

Alignment zweier Sequenzen von Symbolen. Als Symbole dienen die Basen Adenin (A), Guanin (G), Cytosin (C) und Thymin (T).

<code>db</code>	<code>ACCGTA</code>
<code>read</code>	<code>A-CGTT</code>
<code>align(db, read)</code>	<code>MDMMX</code>

Problemstellung

Berechnet das optimale Alignment zweier Sequenzen (db,read).

Ein Alignment besteht aus Editieroperationen, welche mit Kosten behaftet sind. Für die Berechnung wird eine Matrix benutzt, in welcher die benötigten Kosten aufaddiert werden.

		db[x]																					
		T	G	C	A	G	G	T	A	T	A	C	A	G	C	C	A	T	A	T	A	G	C
read[y]	A	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□
	C	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□
	A	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□
	G	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□
	A	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□
	T	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□
	A	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□
	T	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□
A	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	

$$SI(x, y) = \min \begin{cases} SI(x, y - 1) + CIE \\ S(x, y - 1) + CIB \end{cases}$$

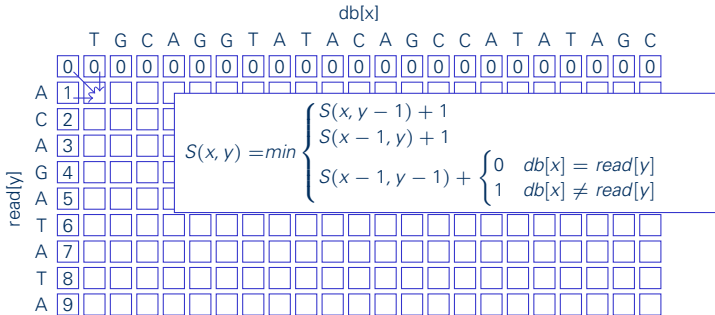
$$SD(x, y) = \min \begin{cases} SD(x - 1, y) + CDE \\ S(x - 1, y) + CDB \end{cases}$$

$$S(x, y) = \min \begin{cases} SI(x, y) \\ SD(x, y) \\ S(x - 1, y - 1) + \begin{cases} CM & db[x] = read[y] \\ CX & db[x] \neq read[y] \end{cases} \end{cases}$$

Problemstellung

Berechnet das optimale Alignment zweier Sequenzen (db,read).

Ein Alignment besteht aus Editieroperationen, welche mit Kosten behaftet sind. Für die Berechnung wird eine Matrix benutzt, in welcher die benötigten Kosten aufaddiert werden.



Problemstellung

Berechnet das optimale Alignment zweier Sequenzen (db,read).

Ein Alignment besteht aus Editieroperationen, welche mit Kosten behaftet sind. Für die Berechnung wird eine Matrix benutzt, in welcher die benötigten Kosten aufaddiert werden.

		db[x]																					
		T	G	C	A	G	G	T	A	T	A	C	A	G	C	C	A	T	A	T	A	G	C
		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
read[y]	A	1	1	1	1	0	1	1	1	0	1	0	1	1	1	1	0	1	0	1	0	1	1
	C	2	2	2	1																		
	A	3	3	3	2																		
	G	4	4	3	3	2																	
	A	5	5	4	4	3																	
	T	6	5	5	5	4	3	3	2	3	3	4	4	3	2	2	2	3	2	3	2	3	3
A	7	6	6	6	5	4	4	3	2	3	3	4	4	3	3	3	2	3	2	3	2	3	4
T	8	7	7	7	6	5	5	4	3	2	3	4	5	4	4	4	3	2	3	2	3	3	4
A	9	8	8	8	7	6	6	5	4	3	2	3	4	5	5	5	4	3	2	3	2	3	4

$$S(x, y) = \min \begin{cases} S(x, y - 1) + 1 \\ S(x - 1, y) + 1 \\ S(x - 1, y - 1) + \begin{cases} 0 & db[x] = read[y] \\ 1 & db[x] \neq read[y] \end{cases} \end{cases}$$

Problemstellung

Berechnet das optimale Alignment zweier Sequenzen (db,read).

Ein Alignment besteht aus Editieroperationen, welche mit Kosten behaftet sind. Für die Berechnung wird eine Matrix benutzt, in welcher die benötigten Kosten aufaddiert werden.

		db[x]																					
		T	G	C	A	G	G	T	A	T	A	C	A	G	C	C	A	T	A	T	A	G	C
read[y]	A	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	C	1	1	1	0	1	1	1	0	1	0	1	0	1	1	1	0	1	0	1	0	1	1
	A	2	2	2	1	1	2	2	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1
	G	3	3	3	2	1	2	2	3	2	2	1	1	0	1	2	2	1	2	1	2	1	2
	A	4	4	3	3	2	1	2	3	3	3	2	2	1	0	1	2	2	2	2	2	2	1
	T	5	5	4	4	3	2	2	3	3	4	3	3	2	1	1	2	2	3	2	3	2	2
	A	6	5	5	5	4	3	3	2	3	3	4	4	3	2	2	2	3	2	3	2	3	3
	T	7	6	6	6	5	4	4	3	2	3	3	4	4	3	3	3	2	3	2	3	2	3
	A	8	7	7	7	6	5	5	4	3	2	3	4	5	4	4	4	3	2	3	2	3	3
	A	9	8	8	8	7	6	6	5	4	3	2	3	4	5	5	5	4	3	2	3	2	3

Problemstellung

Berechnet das optimale Alignment zweier Sequenzen (db,read).

Ein Alignment besteht aus Editieroperationen, welche mit Kosten behaftet sind. Für die Berechnung wird eine Matrix benutzt, in welcher die benötigten Kosten aufaddiert werden.

		db[x]																					
		T	G	C	A	G	G	T	A	T	A	C	A	G	C	C	A	T	A	T	A	G	C
read[y]	A	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	C	1	1	1	1	0	1	1	1	0	1	0	1	1	1	1	0	1	0	1	0	1	1
	A	2	2	2	1	1	1	2	2	1	1	1	0	1	1	1	1	1	1	1	1	1	1
	G	3	3	3	2	1	2	2	3	2	2	1	1	0	1	2	2	1	2	1	2	1	2
	A	4	4	3	3	2	1	2	3	3	3	2	2	1	0	1	2	2	2	2	2	2	1
	T	5	5	4	4	3	2	2	3	3	4	3	3	2	1	1	2	2	3	2	3	2	2
	A	6	5	5	5	4	3	3	2	3	3	4	4	3	2	2	2	3	2	3	2	3	3
	T	7	6	6	6	5	4	4	3	2	3	3	4	4	3	3	3	2	3	2	3	2	3
	A	8	7	7	7	6	5	5	4	3	2	3	4	5	4	4	4	3	2	3	2	3	4
	A	9	8	8	8	7	6	6	5	4	3	2	3	4	5	5	5	4	3	2	3	2	3

Problemstellung

Berechnet das optimale Alignment zweier Sequenzen (db,read).

Ein Alignment besteht aus Editieroperationen, welche mit Kosten behaftet sind. Für die Berechnung wird eine Matrix benutzt, in welcher die benötigten Kosten aufaddiert werden.

		db[x]																							
		T	G	C	A	G	G	T	A	T	A	C	A	G	C	C	A	T	A	T	A	G	C		
		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
read[y]	A	1	1	1	1	0	1	1	1	0	1	0	1	0	1	1	1	0	1	0	1	0	1	1	
	C	2	2	2	1	1	1	2	2	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	
	A	3	3	3	2	1	2	2	3	2	2	1	1	0	1	2	2	1	2	1	2	1	2	2	
	G	4	4	3	3	2	1	2	3	3	3	2	2	1	0	1	2	2	2	2	2	2	2	1	2
	A	5	5	4	4	3	2	2	3	3	4	3	3	2	1	1	2	2	3	2	3	2	2	2	
	T	6	5	5	5	4	3	3	2	3	3	4	4	3	2	2	2	3	2	3	2	3	3	3	
	A	7	6	6	6	5	4	4	3	2	3	3	4	4	3	3	2	3	2	3	2	3	4	4	
	T	8	7	7	7	6	5	5	4	3	2	3	4	5	4	4	4	3	2	3	2	3	3	4	
	A	9	8	8	8	7	6	6	5	4	3	2	3	4	5	5	5	4	3	2	3	2	3	4	

Problemstellung

Berechnet das optimale Alignment zweier Sequenzen (db,read).

Ein Alignment besteht aus Editieroperationen, welche mit Kosten behaftet sind. Für die Berechnung wird eine Matrix benutzt, in welcher die benötigten Kosten aufaddiert werden.

		db[x]																						
		T	G	C	A	G	G	T	A	T	A	C	A	G	C	C	A	T	A	T	A	G	C	
read[y]	A	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	C	1	1	1	0	1	1	1	0	1	0	1	0	1	1	1	0	1	0	1	0	1	1	
	A	2	2	2	1	1	2	2	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	
	G	3	3	3	2	1	2	2	3	2	2	1	1	0	1	2	2	1	2	1	2	1	2	2
	A	4	4	3	3	2	1	2	3	3	3	2	2	1	0	1	2	2	2	2	2	2	1	2
	T	5	5	4	4	3	2	2	3	3	4	3	3	2	1	1	2	2	3	2	3	2	2	2
	A	6	5	5	5	4	3	3	2	3	3	4	4	3	2	2	2	3	2	3	2	3	3	3
	T	7	6	6	6	5	4	4	3	2	3	3	4	4	3	3	3	2	3	2	3	2	3	4
	T	8	7	7	7	6	5	5	4	3	2	3	4	5	4	4	4	3	2	3	2	3	3	4
	A	9	8	8	8	7	6	6	5	4	3	2	3	4	5	5	5	4	3	2	3	2	3	4

Parallelisierung

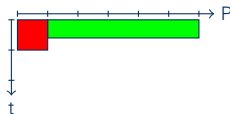
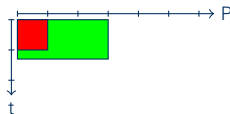
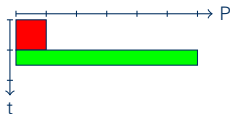
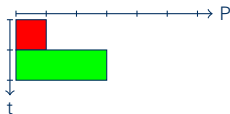
Das Einlesen der Reads benötigt 3,5 % der Zeit und das Ausgeben der Alignments 2,32 %.

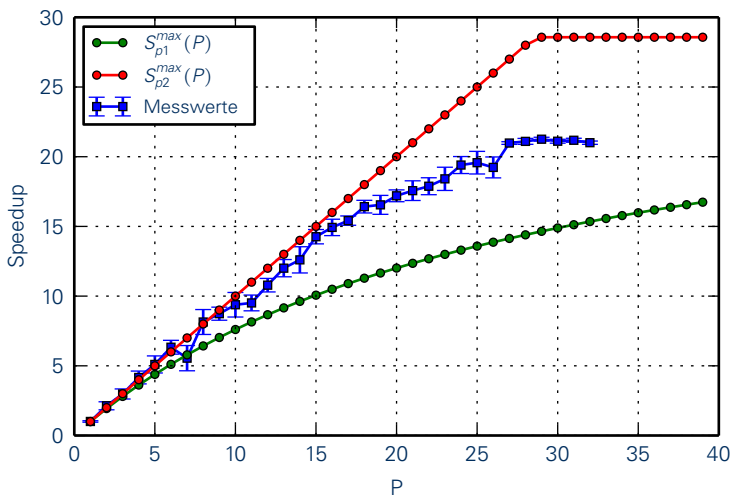
$$S_{p1}^{max}(P) = \frac{1}{f + \frac{(1-f)}{P}}$$

$$f = 0,035$$

(Amdahl's Law)

$$S_{p2}^{max}(P) = \min\left(P, \frac{1}{\max(0,035, 0,0232)}\right)$$



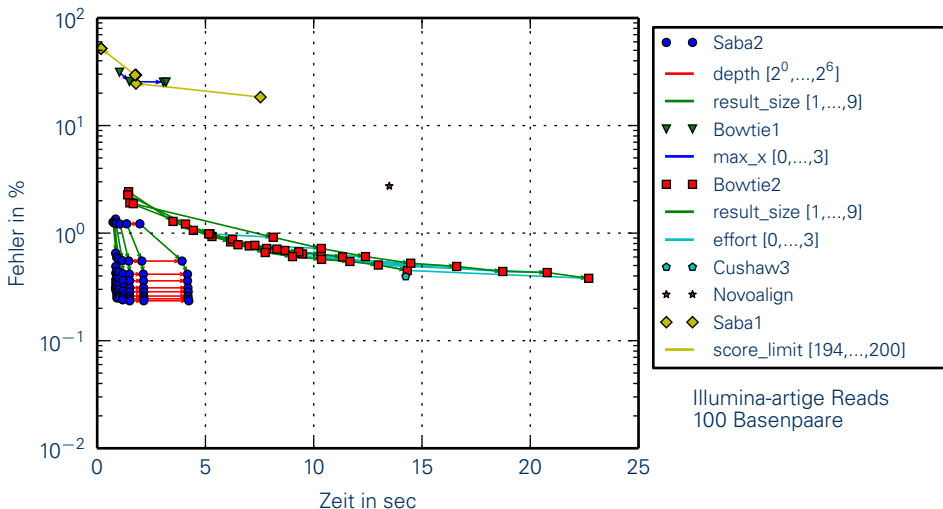


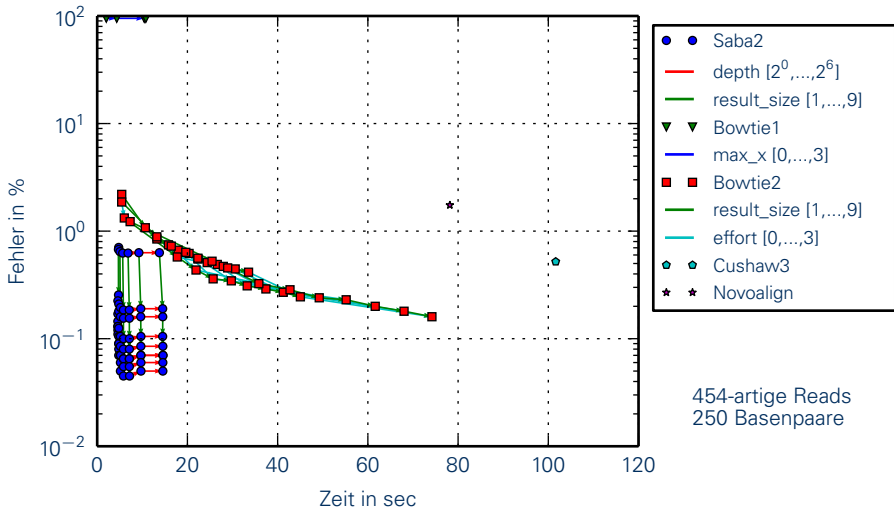
Auswertung/Performance

Analyse von Reads, welche mit verschiedenen Sequenzierungsarten simuliert wurden.

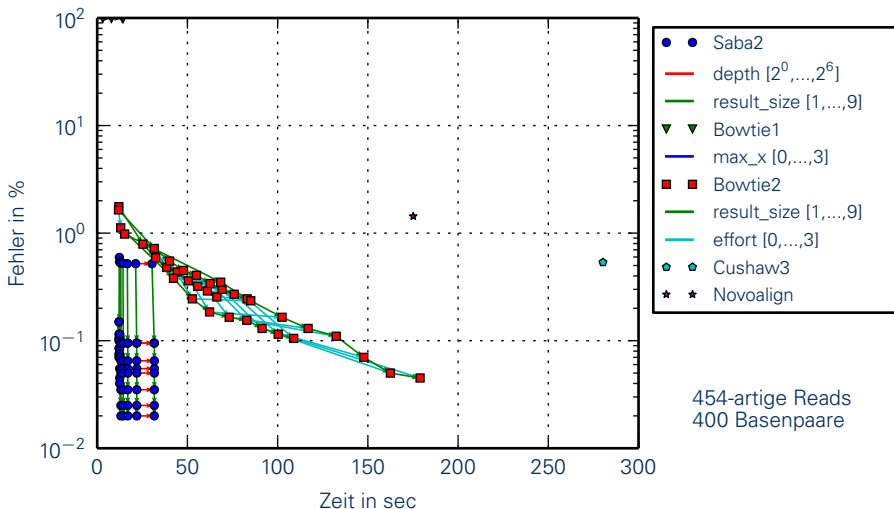
Algorithmus	Jahr	Basenpaare	Entwickler
Bowtie1	2013	35-100	Johns Hopkins University
Bowtie2	2014	50-1000	Johns Hopkins University
Cushaw3	2014	~100	Uni Mainz
Novoalign	2014	short and long	novocraft.com
Saba1	2013		Frank Hoffmann

- Illumina-artig mit 100 Basenpaaren
- 454-artig mit 250 Basenpaaren
- 454-artig mit 400 Basenpaaren

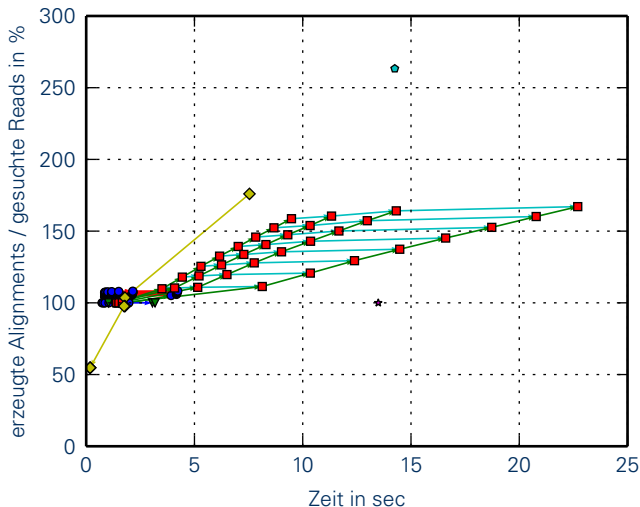




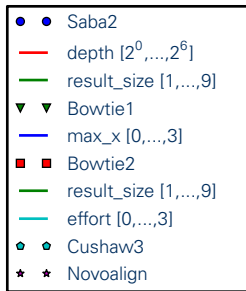
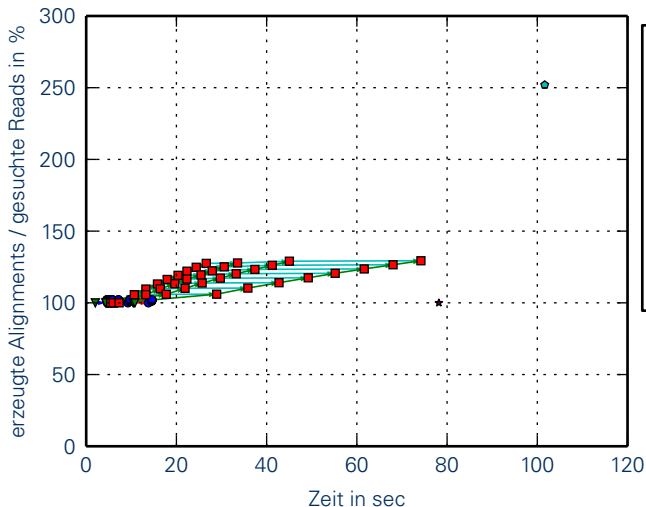
454-artige Reads
250 Basenpaare



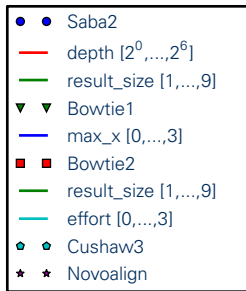
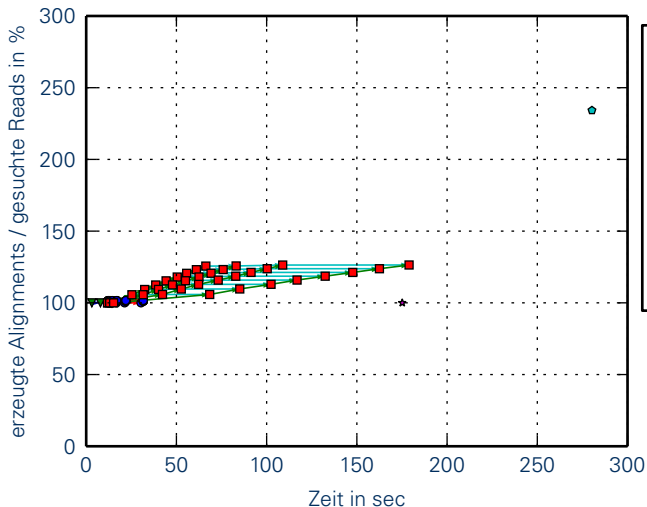
454-artige Reads
400 Basenpaare



Illumina-artige Reads
100 Basenpaare



454-artige Reads
250 Basenpaare



454-artige Reads
400 Basenpaare



Zusammenfassung

- Optimierung des Alignments
- Parallelisierung
- Vergleich mit anderen Verfahren



»Wissen schafft Brücken.«