Bowtie Langmead, Trapnell, Pop, Salzberg 2009

Software

Ultrafast and memory-efficient alignment of short DNA sequences to the human genome Ben Langmead, Cole Trapnell, Mihai Pop and Steven L Salzberg

Address: Center for Bioinformatics and Computational Biology, Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20742, USA.

Correspondence: Ben Langmead. Email: langmead@cs.umd.edu

Published: 4 March 2009

Genome **Biology** 2009, **10:**R25 (doi:10.1186/gb-2009-10-3-r25)

The electronic version of this article is the complete one and can be found online at http://genomebiology.com/2009/10/3/R25

© 2009 Langmead et al.; licensee BioMed Central Ltd. This is an open access article distributed under the terms of the Creative Commons Attribution License (http://creativecommons.org/licenses/by/2.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

CS 4390/5390 Fall 2019

Open Access

Received: 21 October 2008 Revised: 19 December 2008 Accepted: 4 March 2009

The "BWT Index" discussed previously is also called the "FM Index" Originally defined by Ferragina and Manzini in 2000/2005

Reminder that the BWT/FM-index is:

- A data structure for a sting T containing
 - BWT_{T\$} encoded as a wavelet tree

FM-Index

• an integer array C continuing the counts of each character from Σ in T

Refresher on BWTs

Can be constructed using the		
last character of the	(a)	¢
lexicographic order of all	(a)	э а
cyclic rotations of the text	acaacg\$-	a ≁a c

\$acaacg aacg\$ac acaacg\$ acg\$aca → gc\$aaac caacg\$a cg\$acaa g\$acaac

Refresher on BWTs

- Can be constructed using the last character of the lexicographic order of all cyclic rotations of the text
- Encodes the original text, which can be recovered by a walk in the sequence



g\$acaac

- (a) \$acaacg aacg\$ac acaacg\$ acaacg\$ acaacg\$ → acg\$aca → gc\$aaac caacg\$a cg\$acaa g\$acaac
 - cg \$acaacg aacg\$aca acg\$aca caacg\$a caacg\$a caacg\$a caacg\$a caacg\$a caacg\$a caacg\$a caacg\$a

				a	С	g
5	а	C	a	а	С	g
a	a	С	g	\$	a	С
a	С	а	a	C	g	\$
a	С	g	\$	a	С	а
2	а	a	С	g	\$	а
2	9	\$	a	C	CI.	a
a	\$	a	C	a	a	С

			a	a	C
\$	а	С	a	a	С
а	a	C	g	\$	a
а	C	а	a	C	g
a	R	g	\$	a	•
С	a	a	0	g	\$
С	g	\$	а	C	a
a	\$	a	C	a	a

а

а

		С	а	а	С	g
\$	a	c	a	a	С	g
a	-	C	g	\$	t	С
а	C	a	8	C	g	\$
а	С	g	\$	a	C	а
С	a	a	С	g	\$	а
С	g	\$	а	C	a	а
g	\$	a	C	a	a	С



Refresher on BWTs

(a)

(b)

g

Sacaacg

aacg\$ac

acaacg\$

acg\$aca

caacg\$a

cg\$acaa

gSacaac

\$ac

aac

aca

acq

cag

Can be constructed using the last character of the lexicographic order of all cyclic rotations of the text

Encodes the original text, which can be recovered by a walk in the sequence

Searching for patterns is done back to front using similar techniques to sequence recovery





сg	ac
aacg	\$acaac
g\$c	aacg\$a
a 🤈 g \$	acaacg
aca	acg\$aca
cg\$a	caacg\$;
a c a a	Cytaca

		a	a	С	g
\$ a					g
a a					С
a c					\$
a	g	\$	a	•	a
c a	a	0	q	\$	a
cg				a	a
0.0					0

	саа	сg
\$ a		сg
a	cg\$	t C
ac	aac	g \$
ac	g\$a	ca
ca		\$ a
cg		a a
a s		ac

	a	С	a	a	С	g
\$						g
а					3	С
а			a	K	g	\$
а	9	Ø	\$			a
С	a	a	C	g	8	a
С						a
g						С



Using a BWT to Align Reads

allow for errors

Previously mentioned some method to overcome this

- Bowtie assumes all changes are single point changes (i.e. mismatches only) They use a backtracking search to find matching locations • The quality scores are used to prioritize alignments
- - Other speed-ups are included to ensure all matching locations are found

The BWT and FM-Index are insufficient for aligning reads since it doesn't

Start by matching the exact sequence

If the algorithm reaches a point with no matches swap out characters already matched and restart search from that there

When ties occur, start with the character with the lowest quality score, keep the rest in a stack





Start by matching the exact sequence

If the algorithm reaches a point with no matches swap out characters already matched and restart search from that there

When ties occur, start with the character with the lowest quality score, keep the rest in a stack



Start by matching the exact sequence

If the algorithm reaches a point with no matches swap out characters already matched and restart search from that there

When ties occur, start with the character with the lowest quality score, keep the rest in a stack



Start by matching the exact sequence

If the algorithm reaches a point with no matches swap out characters already matched and restart search from that there

When ties occur, start with the character with the lowest quality score, keep the rest in a stack



Start by matching the exact sequence

If the algorithm reaches a point with no matches swap out characters already matched and restart search from that there

When ties occur, start with the character with the lowest quality score, keep the rest in a stack



Start by matching the exact sequence

If the algorithm reaches a point with no matches swap out characters already matched and restart search from that there

When ties occur, start with the character with the lowest quality score, keep the rest in a stack

Keep track of how many changes are made

"Bowtie conducts a quality-aware, greedy, randomized, depth-first search through the space of possible alignments."



Backtracking Options

The user specifies the sum of the quality scores that can be changed • one medium quality change

- Bowtie outputs the first valid alignment by default (within the specified constraints) can be modified to complete the backtracking and return the "best" alignment
 - 2x-3x slower to do this

User can specify a number of alignments to consider

- default is to use only one
- might want the two best alignments
- 2 alignments is ~2x slower than using only 1

this means that a mapping can have lots of low quality replacements, or

In low quality reads, lots of time may be spent backtracking since there are many possible changes at low quality positions.

They mitigate this by creating two indexes (as we saw previously), one for the forward and one for the reverse of the string • the backtracking is performed somewhat simultaneously on both index

as we will see next

One other step they take is to concentrate on the "high-quality" end of a read (the first 28 characters read) which is most reliable

Excessive Backtracking

Phased Search

Split the seed (first 28 bases) into two parts, hi-half and lo-half

Assume we're allowing 2 changes in the seed, a good alignment will have either:

- 1. no mismatches
- 2. no mismatches in hi-half, 1 or 2 mismatches in lo-half
- 3.1 or 2 mismatches in hi-half, no mismatches in lo-half
- 4.1 mismatch in hi-half, 1 mismatch in lo-half

Phase 1 uses the mirror index and invokes the aligner to find alignments for cases 1 & 2.
Phases 2 and 3 cooperate to find alignments for case 3:
Phase 2 finds partial alignments with mismatches only in the hi-half, and phase 3 attempts to extend those partial alignments into full alignments.
Finally, phase 3 invokes the aligner to find alignments for case 4.



Phased search with reverse strand

Since both the read and its reverse complement are possibilities for ma match, need to consider both.

Phases 2-4 here map to phases 1-3 previously.



Performance

Maq (Li, Ruan, Durban 2008), SOAP (Li, Li, Kristiansen, Wang 2008) the leading competitors at the time

Both used hashing to find potential mapping locations

Performance

	Platform	CPU time	Wall clock time	Reads mapped per hour (millions)	Peak virtual memory footprint (megabytes)	Bowtie speed-up	Read aligne
Bowtie -v 2	Sorvor	15 m 7 s	15 m 41 s	33.8	1,149	251~	
SOAP	Server	91 h 57 m 35 s	91 h 47 m 46 s	0.10	13,619	JJIX	
Bowtie		16 m 41 s	17 m 57 s	29.5	1,353	50.0.	
MAQ	PC	17 h 46 m 35 s	17 h 53 m 7 s	0.49	804	09.0X	
Bowtie	Samor	17 m 58 s	18 m 26 s	28.8	1,353	107.	
MAQ	Server	32 h 56 m 53 s	32 h 58 m 39 s	0.27	804	107 X	



Bowtie was (at the time) the fastest short read aligniner

time)

Is able to run on a standard PC

When first published didn't use mate-pair information



- Used a one-time index based on a BWT that could be reused (novel at the

Bowtie2



Bowtie2

eds

ŝ

e

















Bowtie2