CS 600.226: Data Structures Michael Schatz

Nov 16 2018 Lecture 33: BWT + Assembly



Assignment 8: Competitive Spelling Bee

Out on: November 9, 2018 Due by: November 16, 2018 before 10:00 pm Collaboration: None Grading:

Packaging 10%, Style 10% (where applicable), Testing 10% (where applicable), Performance 30% (where applicable), Functionality 40% (where applicable)

Overview

Your "one" task for this assignment is to take the simple spell checker we give you and to turn it into the fastest, most memory-efficient spell checker in the course, subject to the constraints detailed below. You are expected to do this by (once again) implementing the Map interface, this time using one of several hash table techniques (your choice, see below).

> Remember: javac –Xlint:all & checkstyle *.java & Junit & Jaybee BenchMarks

Assignment 9: StringOmics

Out on: November 16, 2018 Due by: November 30, 2018 before 10:00 pm Collaboration: None Grading:

Packaging 10%, Style 10% (where applicable), Testing 10% (where applicable), Performance 10% (where applicable), Functionality 60% (where applicable)

Overview

The ninth assignment focuses on data structures and operations on strings. In this assignment you will implement encoding and decoding using the Burrows Wheeler Transform as well as encoding and decoding in a simple form of run length encoding. In the final problem you will be asked to measure the space savings using run length encoding with and without applying the Burrows Wheeler Transform first.

Remember: javac –Xlint:all & checkstyle *.java & Junit (No JayBee)

Part I: Burrows Wheeler Transform

Personal Genomics

How does your genome compare to the reference?



Exact Matching Review & Overview

Where is GATTACA in the human genome?



*** These are general techniques applicable to any text search problem ***

Algorithmic challenge

How can we combine the speed of a suffix array O(m + lg(n)) (or even O(m)) with the size of a brute force analysis (n bytes)?

What would such an index look like?



Bowtie: Ultrafast and memory efficient alignment of short DNA sequences to the human genome

Slides Courtesy of Ben Langmead

• Reversible permutation of the characters in a text



A block sorting lossless data compression algorithm. Burrows M, Wheeler DJ (1994) *Digital Equipment Corporation*. Technical Report 124

• Permutation of the characters in a text



• BWT(T) is the index for T

A block sorting lossless data compression algorithm.

Burrows M, Wheeler DJ (1994) Digital Equipment Corporation. Technical Report 124

• Reversible permutation of the characters in a text



BWT(T) is the index for T

implicitly encodes Suffix Array

A block sorting lossless data compression algorithm. Burrows M, Wheeler DJ (1994) Digital Equipment Corporation. Technical Report 124

- Recreating T from BWT(T)
 - Start in the first row and apply LF repeatedly, accumulating predecessors along the way



[Decode this BWT string: ACTGA\$TTA]



ref[614]:

It_was_the_best_of_times,_it_was_the_worst_of_times,_it_was_the_age_ of_wisdom,_it_was_the_age_of_foolishness,_it_was_the_epoch_of_belief ,_it_was_the_epoch_of_incredulity,_it_was_the_season_of_Light,_it_wa s_the_season_of_Darkness,_it_was_the_spring_of_hope,_it_was_the_wint er_of_despair,_we_had_everything_before_us,_we_had_nothing_before_us ,_we_were_all_going_direct_to_Heaven,_we_were_all_going_direct_the_o ther_way_-_in_short,_the_period_was_so_far_like_the_present_period,_ that_some_of_its_noisiest_authorities_insisted_on_its_being_received ,_for_good_or_for_evil,_in_the_superlative_degree_of_comparison_only.\$

rle(ref)[614]:

It_was_the_best_of_times,_it_was_the_worst_of_times,_it_was_the_age_ of_wisdom,_it_was_the_age_of_fo2lishnes2,_it_was_the_epoch_of_belief ,_it_was_the_epoch_of_incredulity,_it_was_the_season_of_Light,_it_wa s_the_season_of_Darknes2,_it_was_the_spring_of_hope,_it_was_the_wint er_of_despair,_we_had_everything_before_us,_we_had_nothing_before_us ,_we_were_al2_going_direct_to_Heaven,_we_were_al2_going_direct_the_o ther_way___in_short,_the_period_was_so_far_like_the_present_period,_ that_some_of_its_noisiest_authorities_insisted_on_its_being_received ,_for_go2d_or_for_evil,_in_the_superlative_degre2_of_comparison_only.\$

ref[614]:

It_was_the_best_of_times,_it_was_the_worst_of_times,_it_was_the_age_ of_wisdom,_it_was_the_age_of_foolishness,_it_was_the_epoch_of_belief ,_it_was_the_epoch_of_incredulity,_it_was_the_season_of_Light,_it_wa s_the_season_of_Darkness,_it_was_the_spring_of_hope,_it_was_the_wint er_of_despair,_we_had_everything_before_us,_we_had_nothing_before_us ,_we_were_all_going_direct_to_Heaven,_we_were_all_going_direct_the_o ther_way__in_short,_the_period_was_so_far_like_the_present_period,_ that_some_of_its_noisiest_authorities_insisted_on_its_being_received ,_for_good_or_for_evil,_in_the_superlative_degree_of_comparison_only.\$

bwt[614]:

ref[614]:

It_was_the_best_of_times,_it_was_the_worst_of_times,_it_was_the_age_ of_wisdom,_it_was_the_age_of_foolishness,_it_was_the_epoch_of_belief ,_it_was_the_epoch_of_incredulity,_it_was_the_season_of_Light,_it_wa s_the_season_of_Darkness,_it_was_the_spring_of_hope,_it_was_the_wint er_of_despair,_we_had_everything_before_us,_we_had_nothing_before_us ,_we_were_all_going_direct_to_Heaven,_we_were_all_going_direct_the_o ther_way__in_short,_the_period_was_so_far_like_the_present_period,_ that_some_of_its_noisiest_authorities_insisted_on_its_being_received ,_for_good_or_for_evil,_in_the_superlative_degree_of_comparison_only.\$

bwt[614]:

bwt[614]:

rle(bwt)[464]:

.dlms2ftysesdtrsns_y_2\$_yfofe4tg2sfefefg2e2drofr,l2re2f-,fs,9nfrsdn2 hereghet2edndete2ge2nste2,s5t,es3ns2f2te2dt10r,4e3feh2_2p_2fpDw11e2h l_ew_5eo2_ne3oa2eo2_4seph2r2hvh2w2egmgh7kr2w2h2s2Hr3vtr2ib2dbcbvs_2t hw2p3vm2irdn2ib_2eo12_4e2n6a2i_3ec2_2t18s_tsgltsLlvt2_3h2o2re_wr2ad2 wlors_9r_2lteiril2re_oua2no2i2oeo4i3hki6o_2ieitsp2ioi_12g2nodsc_s3_g fhf_f3hwh_nsmo_2ue2_sio3ae4o2_i2cgp2e2aoaeo2e2s2eu2teta11i_2ei_in_2a 2ie_e3rei_hrs3nac2i2Ii7sn_15oyoui_2a_i3ds_2ai2ae2_21tlar

ref[614]:

It_was_the_best_of_times,_it_was_the_worst_of_times,_it_was_the_age_ of_wisdom,_it_was_the_age_of_foolishness,_it_was_the_epoch_of_belief ,_it_was_the_epoch_of_incredulity,_it_was_the_season_of_Light,_it_wa s_the_season_of_Darkness,_it_was_the_spring_of_hope,_it_was_the_wint er_of_despair,_we_had_everything_before_us,_we_had_nothing_before_us ,_we_were_all_going_direct_to_Heaven,_we_were_all_going_direct_the_o ther_way_-_in_short,_the_period_was_so_far_like_the_present_period,_ that_some_of_its_noisiest_authorities_insisted_on_its_being_received ,_for_good_or_for_evil,_in_the_superlative_degree_of_comparison_only.\$

rle(bwt)[464]:

.dlms2ftysesdtrsns_y_2\$_yfofe4tg2sfefefg2e2drofr,l2re2f-,fs,9nfrsdn2 hereghet2edndete2ge2nste2,s5t,es3ns2f2te2dt10r,4e3feh2_2p_2fpDw11e2h l_ew_5eo2_ne3oa2eo2_4seph2r2hvh2w2egmgh7kr2w2h2s2Hr3vtr2ib2dbcbvs_2t hw2p3vm2irdn2ib_2eo12_4e2n6a2i_3ec2_2t18s_tsgltsLlvt2_3h2o2re_wr2ad2 wlors_9r_2lteiril2re_oua2no2i2oeo4i3hki6o_2ieitsp2ioi_12g2nodsc_s3_g fhf_f3hwh_nsmo_2ue2_sio3ae4o2_i2cgp2e2aoaeo2e2s2eu2teta11i_2ei_in_2a

^{2ie_e3rei}. Saved 614-464 = 150 bytes (24%) with zero loss of information!

Common to save 50% to 90% on real world files with bzip2

BWT Exact Matching

 LFc(r, c) does the same thing as LF(r) but it ignores r's actual final character and "pretends" it's c:

> LFc(5, g) = 8 a c a a c g a a c g a c a c a a c g a c a c g a c a c a a c g a c a c a a c g a c a c g a c a a c gRank: 2 g a c a a c gF

BWT Exact Matching

 Start with a range, (top, bot) encompassing all rows and repeatedly apply LFc:

top = LFc(top, qc); bot = LFc(bot, qc)

qc = the next character to the left in the query



Ferragina P, Manzini G: Opportunistic data structures with applications. FOCS. IEEE Computer Society; 2000.

[Search for TTA this BWT string: ACTGA\$TTA]

Algorithm Overview



Genetic Associations



https://www.ebi.ac.uk/gwas/diagram

Part 3: Genome Assembly

Shredded Book Reconstruction

Dickens accidentally shreds the first printing of <u>A Tale of Two Cities</u>
– Text printed on 5 long spools

It v	vas	thev!	aesthæ	fbes	ti o fetsiņ	itesy	ais vi	aes telo	erstror	of tir	nes, i	t was t	he	a gg cb	fwixsidalon	niti twava	sthe alg	e ageo	ofictories	mess,	
It v	vas	thevi	aestt he	of	times,	it v	vas t	the n	e wors	t of t	imes,	it was	the	tinge age	woißatoised,	ointnyviats	thevag	ethê fa	geish hos	kshne	ss,
It v	vas	the va	.sb æ t	b£\$	inesiri	tew;	aist w	abeh	eowstro	f tifn	as,eit,	t was	the	ag e of	wisdom,	i it wa	as the a	age of	li stolis	sņess,	 []
It v	vas	t tik a	esbelse	bes	innesin	tes,	was	abeh	eo nst ro	f tift	ċs ,es	it was	the	age of	vi sciscio ,	nit, istassa	etshehæg	age f fo	ofistolis	sness,	
It	w	alt th	esbelset	b£\$	imés in	eist,	vitas	ahehw	owstro	f of t	imes,	it was	the	age of	ovfiewilsolo;	niț itravsa	atshtchæg	e a gto	ofistoliss	siness,	

- How can he reconstruct the text?
 - 5 copies x 138,656 words / 5 words per fragment = 138k fragments
 - The short fragments from every copy are mixed together
 - Some fragments are identical



Greedy Reconstruction



The repeated sequence make the correct reconstruction ambiguous

It was the best of times, it was the [worst/age]

Model the assembly problem as a graph problem

de Bruijn Graph Construction

- $D_k = (V, E)$
 - V = All length-k subfragments (k < l)
 - E = Directed edges between consecutive subfragments
 - Nodes overlap by k-I words



- Locally constructed graph reveals the global sequence structure
 - Overlaps between sequences implicitly computed

de Bruijn, 1946 Idury and Waterman, 1995 Pevzner, Tang, Waterman, 2001



de Bruijn Graph Assembly



Genomics Across the Tree of Life





http://technical.ly/baltimore/2016/11/01/johns-hopkins-genome-algorithm-wine/





Next Steps

- I. Reflect on the magic and power of Suffix Arrays and the BWT!
- I. Assignment 9 due Friday November 30 @ 10pm