

Read Mapping

Michael Schatz

Feb 22, 2021

Lecture 9: Computational Biomedical Research



Assignment 3: de Bruijn Graphs

Due Feb 17 @ 11:59pm

Assignment Date: Wednesday, Feb. 10, 2020
Due Date: Wednesday, Feb. 17, 2020 @ 11:59pm

Assignment Overview

In this assignment you will take a closer look at coverage, build and analyze a simple de Bruijn graph, and write your own compacted de Bruijn graph generator that you will test on a small microbial genome using different length kmers.

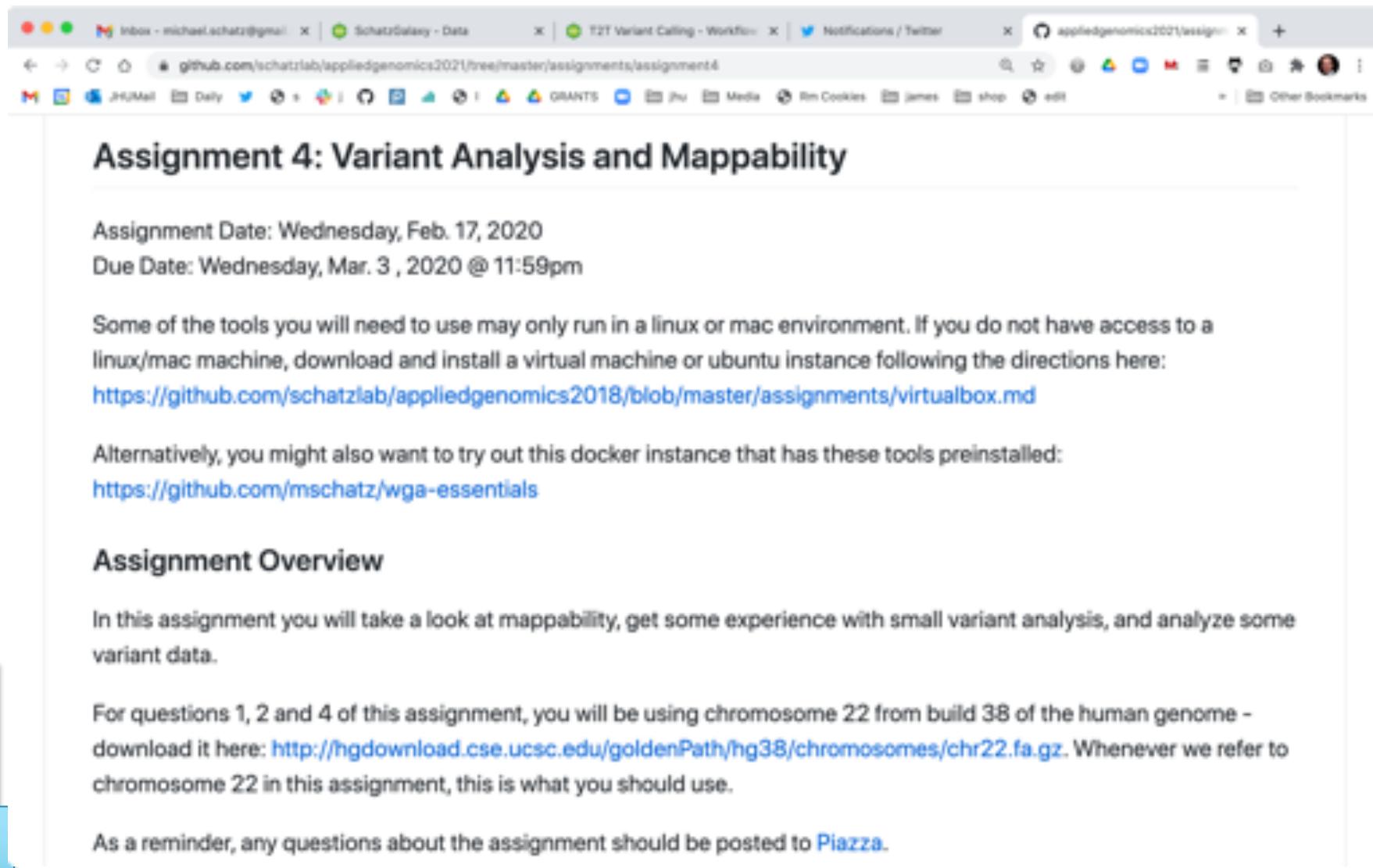
As a reminder, any questions about the assignment should be posted to [Piazza](#).

Question 1. Coverage simulator [20 pts]

- Q1a. How many 100bp reads are needed to sequence a 1Mbp genome to 5x coverage?
- Q1b. In the language of your choice, simulate sequencing 5x coverage of a 1Mbp genome with 100bp reads and plot the histogram of coverage. Note you do not need to actually output the sequences of the reads, you can just randomly sample positions in the genome and record the coverage. You do not need to consider the strand of each read. The start position of each read should have a uniform random probability at each possible starting position (1 through 999,901). You can record the coverage in an array of 1M positions. Overlay the histogram with a Poisson distribution with lambda=5
- Q1c. Using the histogram from 1b, how much of the genome has not been sequenced (has 0x coverage)? How well does this match Poisson expectations?
- Q1d. Now repeat the analysis with 15x coverage: 1. simulate the appropriate number of reads, 2. make a histogram, 3. overlay a Poisson

Assignment 4: Variant Analysis & Mappability

Due March 3 @ 11:59pm



The screenshot shows a web browser window with the following tabs:

- Inbox - michael.schatz@gmail.com
- SchatzGalaxy - Data
- T2T Variant Calling - Workflow
- Notifications / Twitter
- appliedgenomics2021/assignment4

The main content area displays the assignment details:

Assignment 4: Variant Analysis and Mappability

Assignment Date: Wednesday, Feb. 17, 2020
Due Date: Wednesday, Mar. 3 , 2020 @ 11:59pm

Some of the tools you will need to use may only run in a linux or mac environment. If you do not have access to a linux/mac machine, download and install a virtual machine or ubuntu instance following the directions here:
<https://github.com/schatzlab/appliedgenomics2018/blob/master/assignments/virtualbox.md>

Alternatively, you might also want to try out this docker instance that has these tools preinstalled:
<https://github.com/mschatz/wga-essentials>

Assignment Overview

In this assignment you will take a look at mappability, get some experience with small variant analysis, and analyze some variant data.

For questions 1, 2 and 4 of this assignment, you will be using chromosome 22 from build 38 of the human genome - download it here: <http://hgdownload.cse.ucsc.edu/goldenPath/hg38/chromosomes/chr22.fa.gz>. Whenever we refer to chromosome 22 in this assignment, this is what you should use.

As a reminder, any questions about the assignment should be posted to [Piazza](#).

<https://github.com/schatzlab/appliedgenomics2021>

Genomics Arsenal in the Year 2021

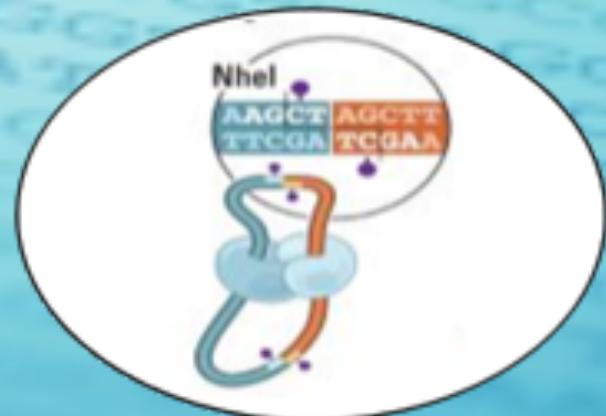
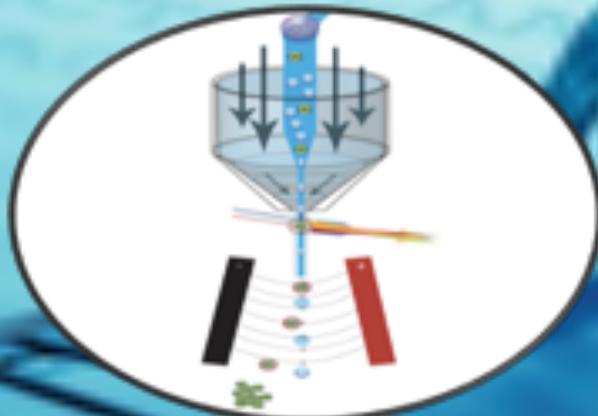
Sample Preparation



Sequencing

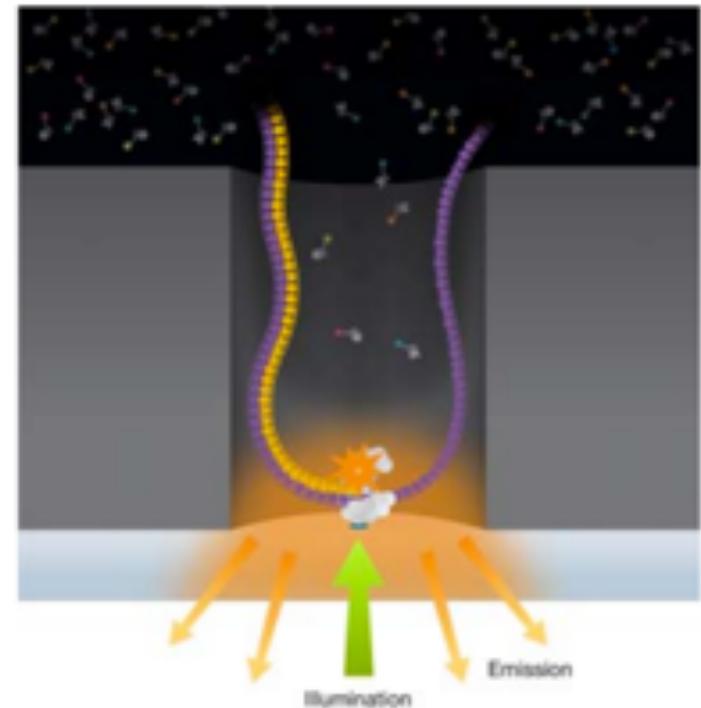
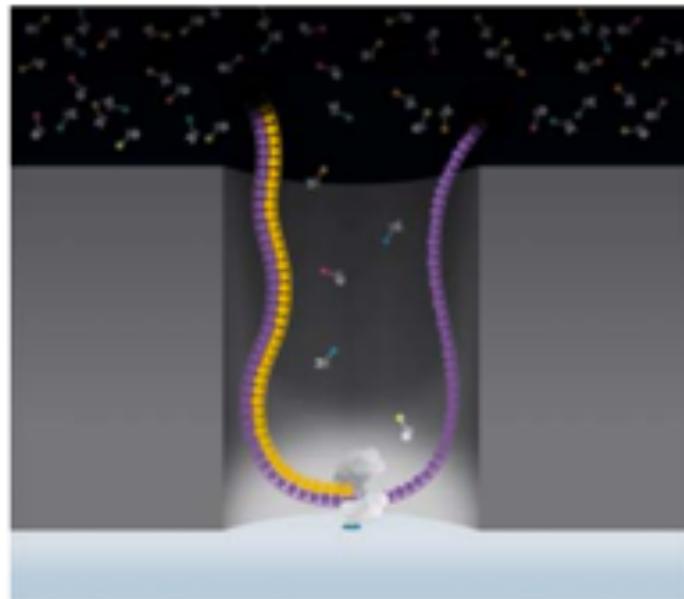


Chromosome Mapping

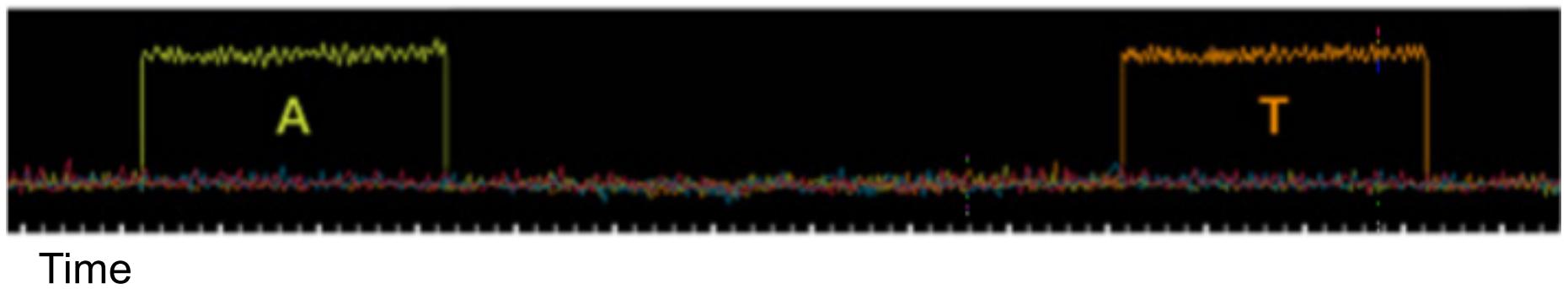


PacBio: SMRT Sequencing

Imaging of fluorescent phospholinked labeled nucleotides as they are incorporated by a polymerase anchored to a Zero-Mode Waveguide (ZMW).



Intensity

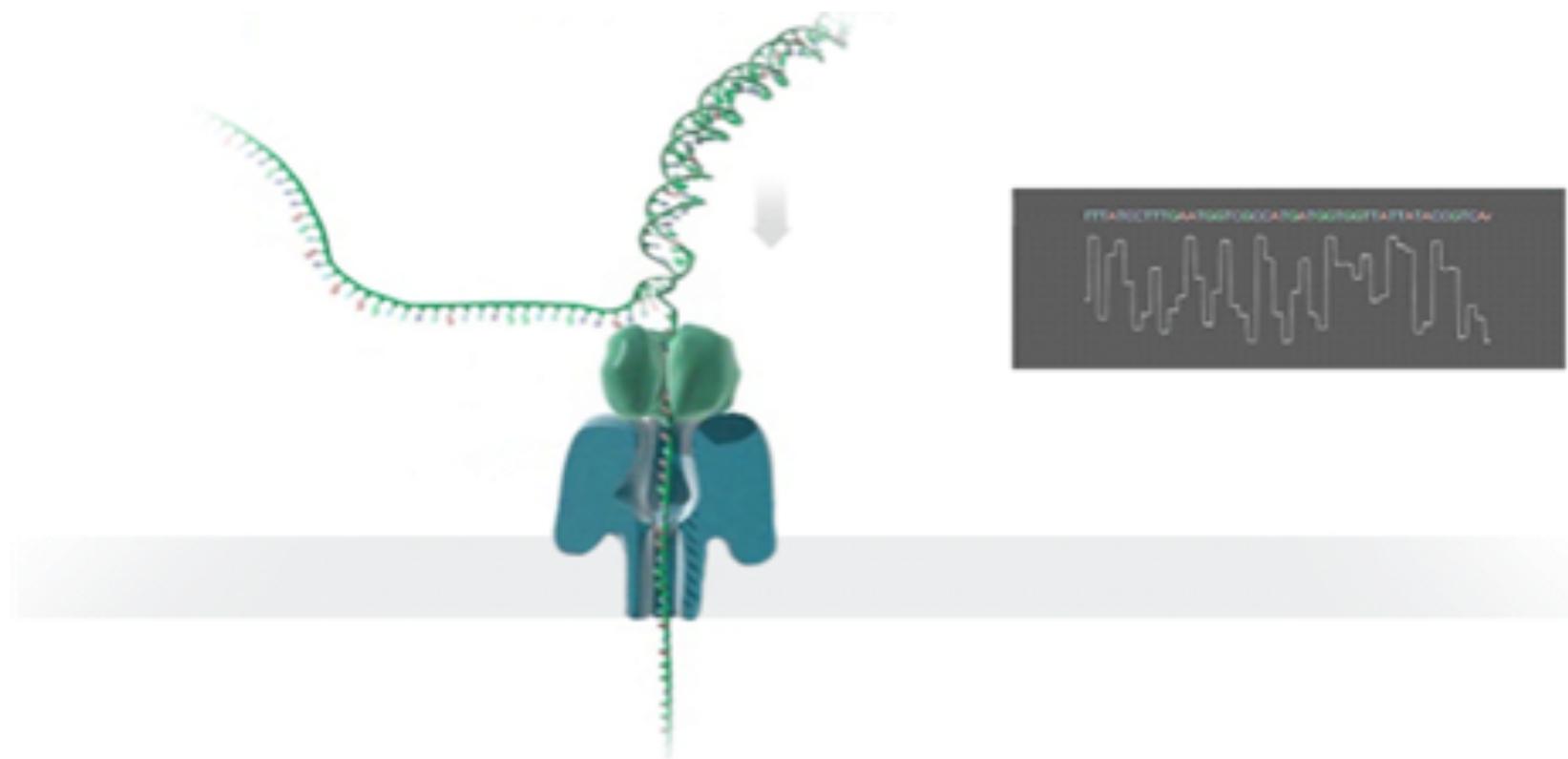


Oxford Nanopore Technologies (ONT)

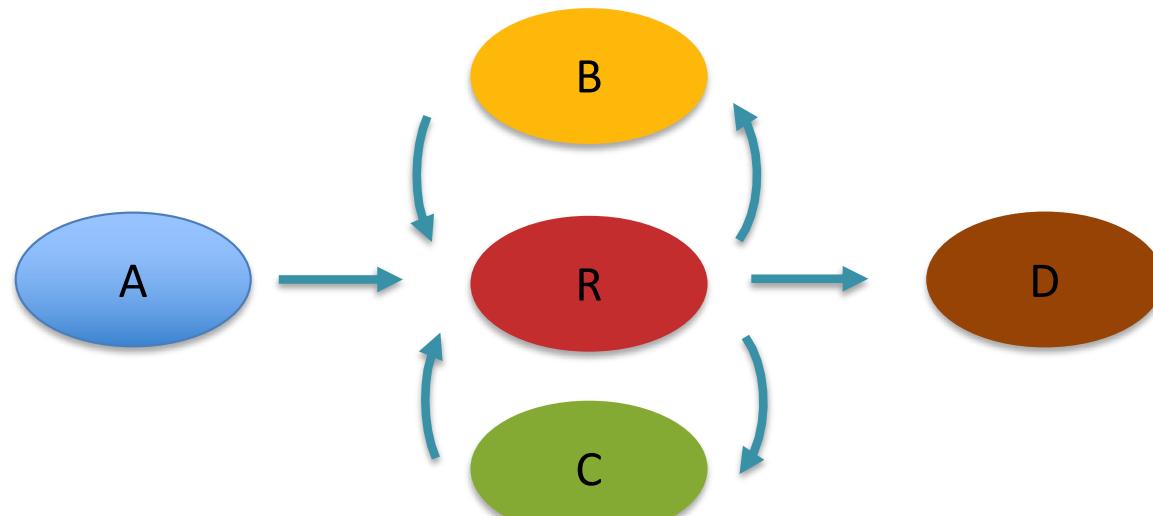
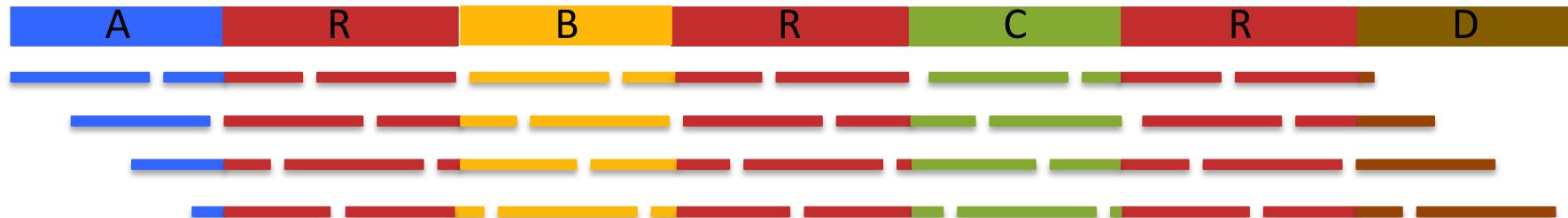


Nanopore Sequencing

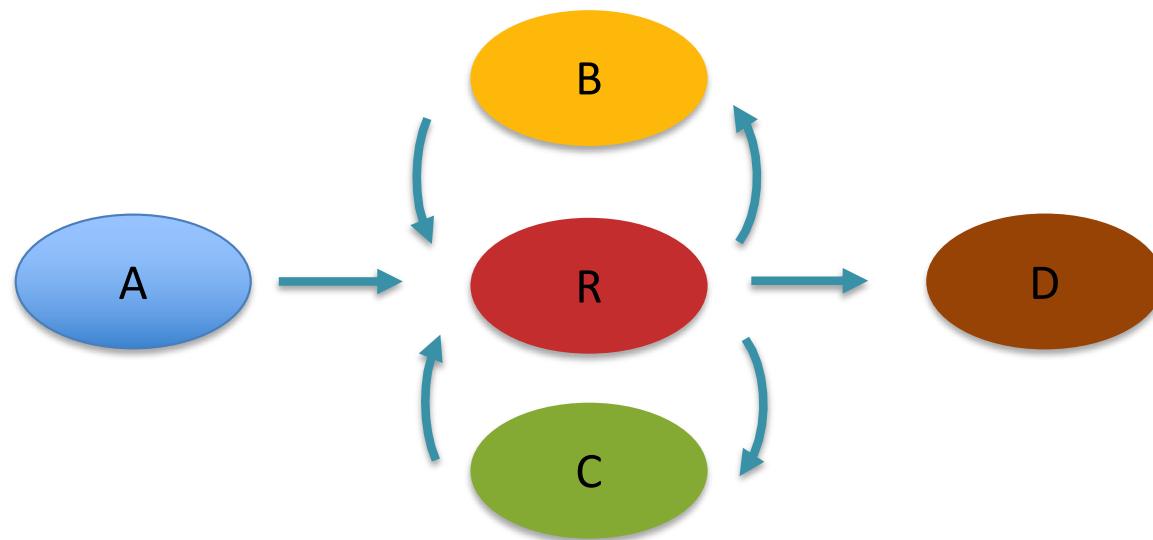
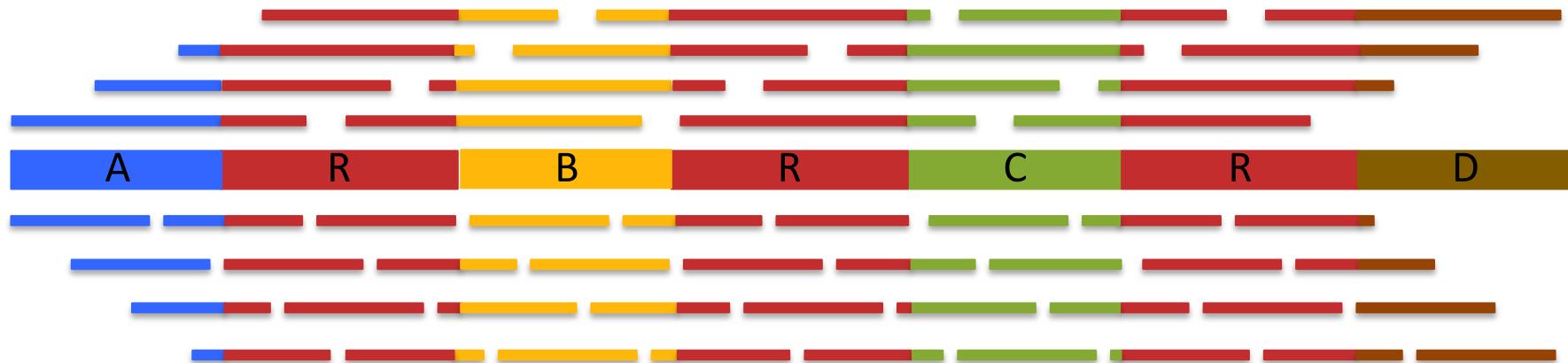
Sequences DNA/RNA by measuring changes in ionic current as nucleotide strand passes through a pore



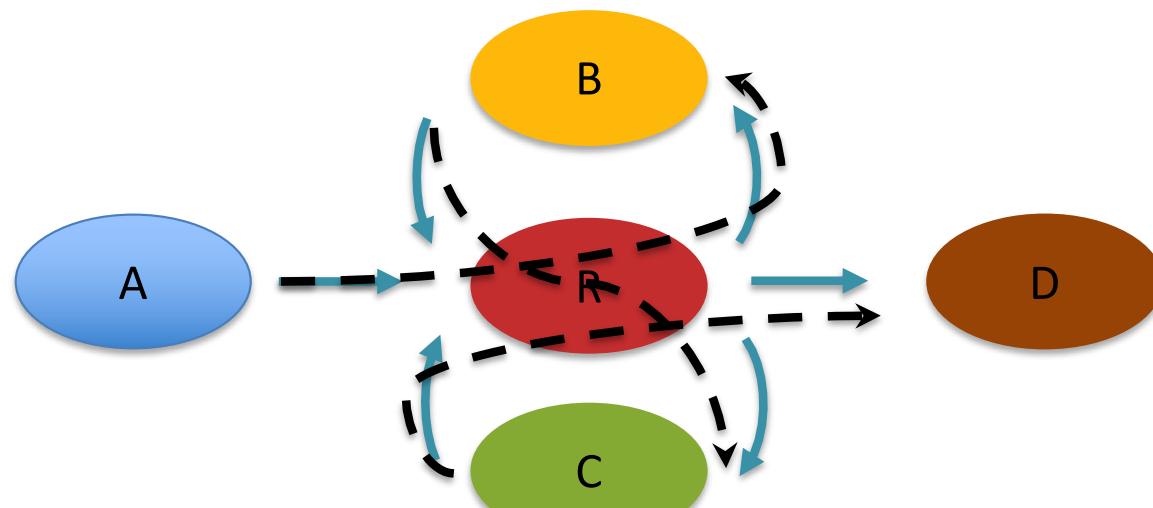
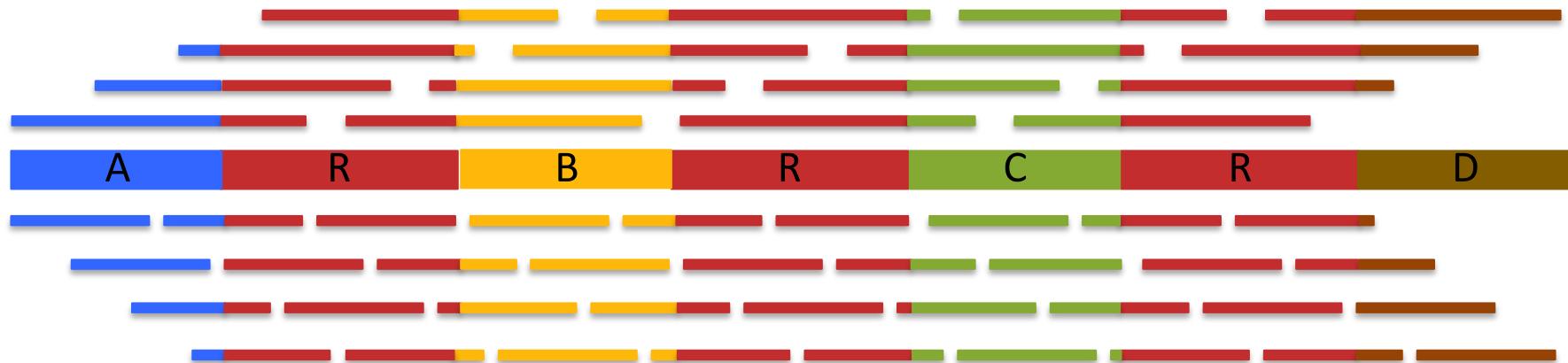
Assembly Complexity



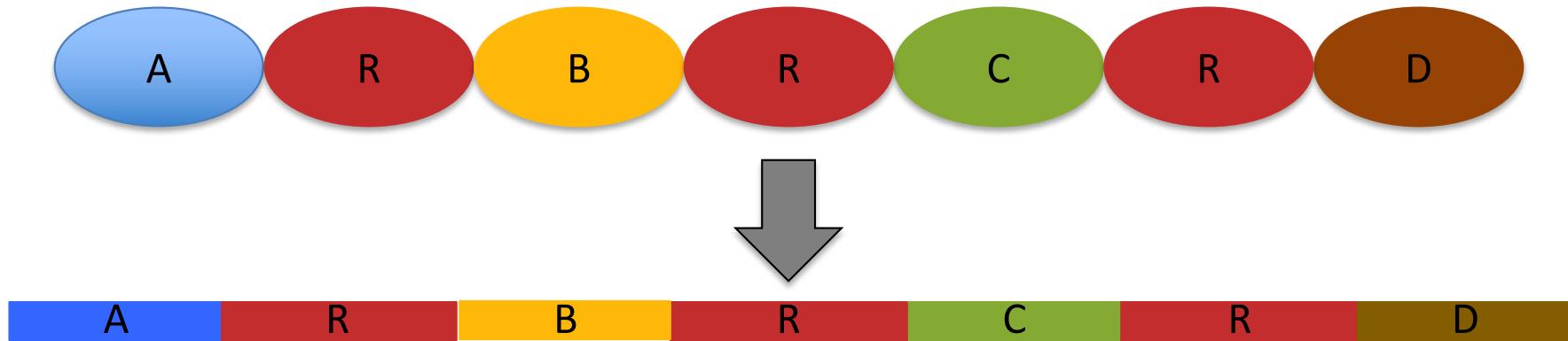
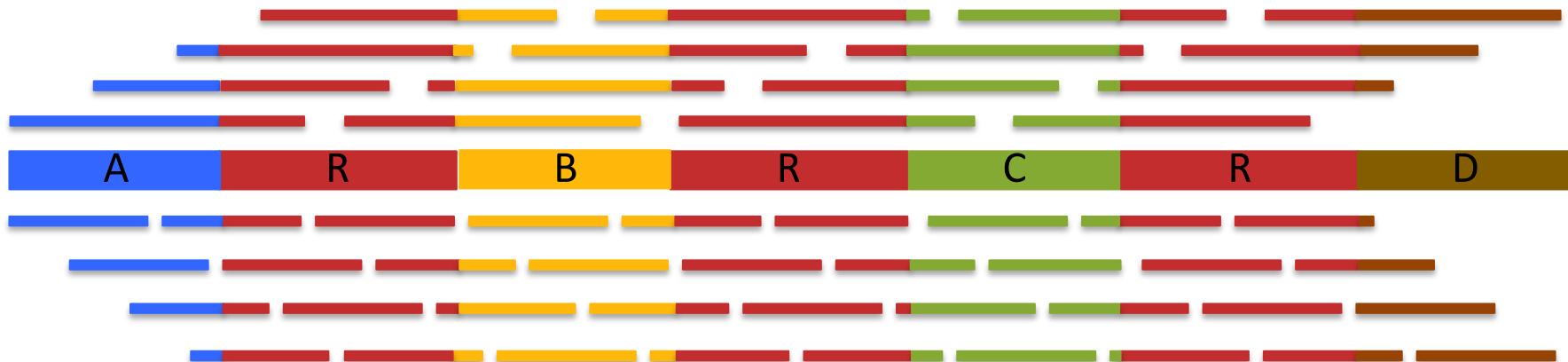
Assembly Complexity



Assembly Complexity



Assembly Complexity

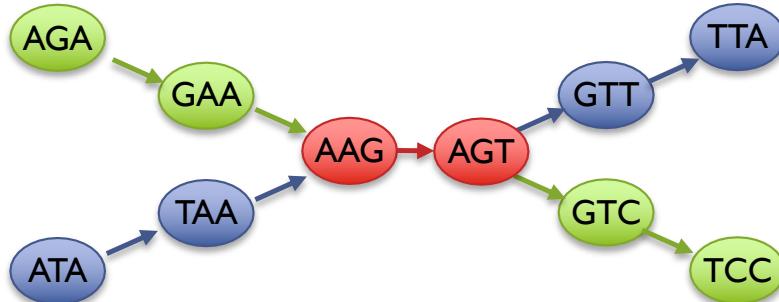


The advantages of SMRT sequencing

Roberts, RJ, Carneiro, MO, Schatz, MC (2013) *Genome Biology*. 14:405

Two Paradigms for Assembly

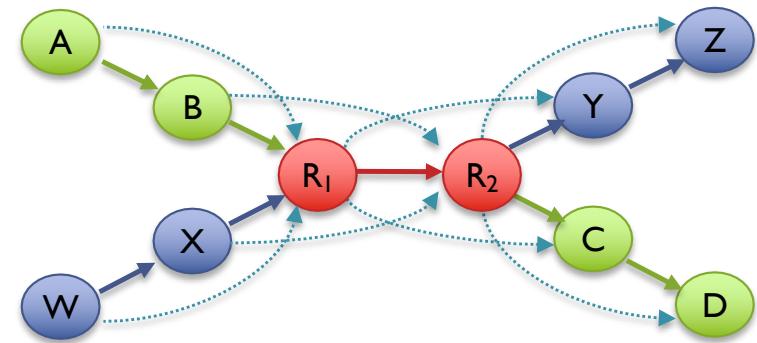
de Bruijn Graph



Short read assemblers

- Repeats depends on word length
- Read coherency, placements lost
- Robust to high coverage

Overlap Graph

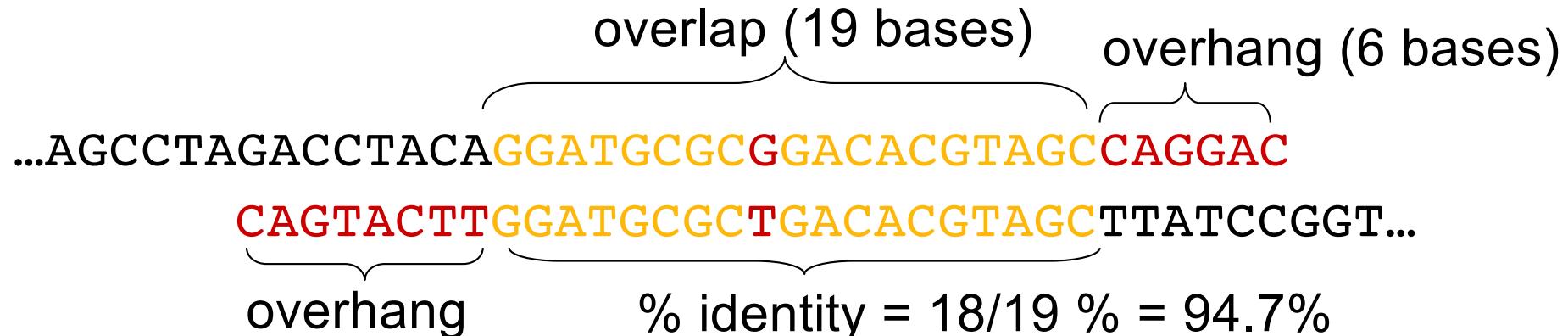


Long read assemblers

- Repeats depends on read length
- Read coherency, placements kept
- Tangled by high coverage

Assembly of Large Genomes using Second Generation Sequencing
Schatz MC, Delcher AL, Salzberg SL (2010) *Genome Research*. 20:1165-1173.

Overlap between two sequences



overlap - region of similarity between regions

overhang - un-aligned ends of the sequences

The assembler screens merges based on:

- length of overlap
- % identity in overlap region
- maximum overhang size.

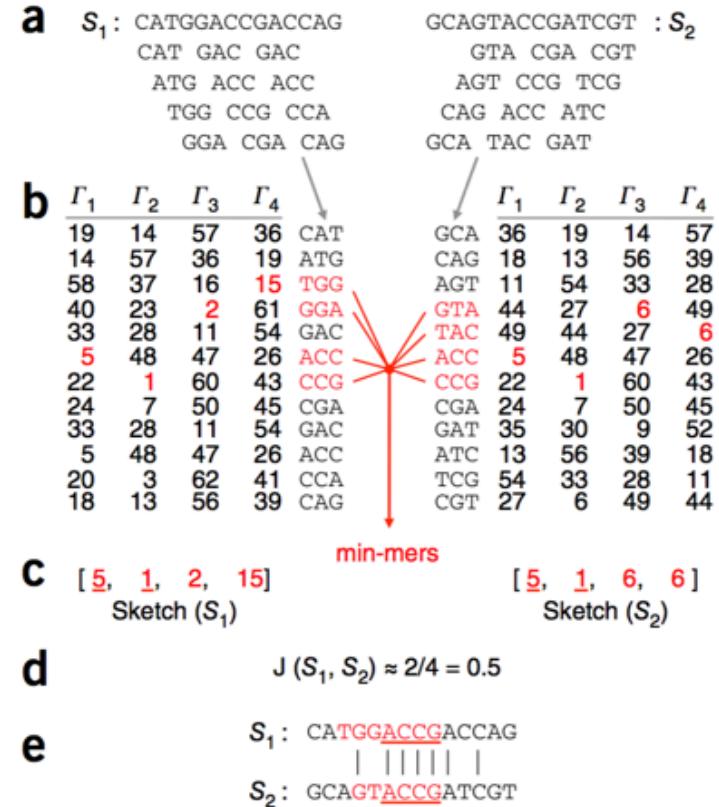
[How do we compute the overlap?]

[Do we really want to do all-vs-all?]

Very fast approximate overlapping

Maybe we don't need to compute the exact identity of the overlap region, just approximate it

- If two reads overlap, they should share many of the same kmers: Their Jaccard coefficient should be high: $|\text{intersection}| / |\text{union}|$
- But tracking all of the kmers for a read is a lot of overhead
- Instead, compare the “sketch” of the reads: a small fraction of kmers carefully chosen
- LSH: Find the sketch by applying N hash functions to the kmers, and keeping the minimum hash values reported from each ($N=4$ in example)
- This forms a nice “random” sample of the reads, and the Jaccard coefficient is a good approximation of the sequence similarity



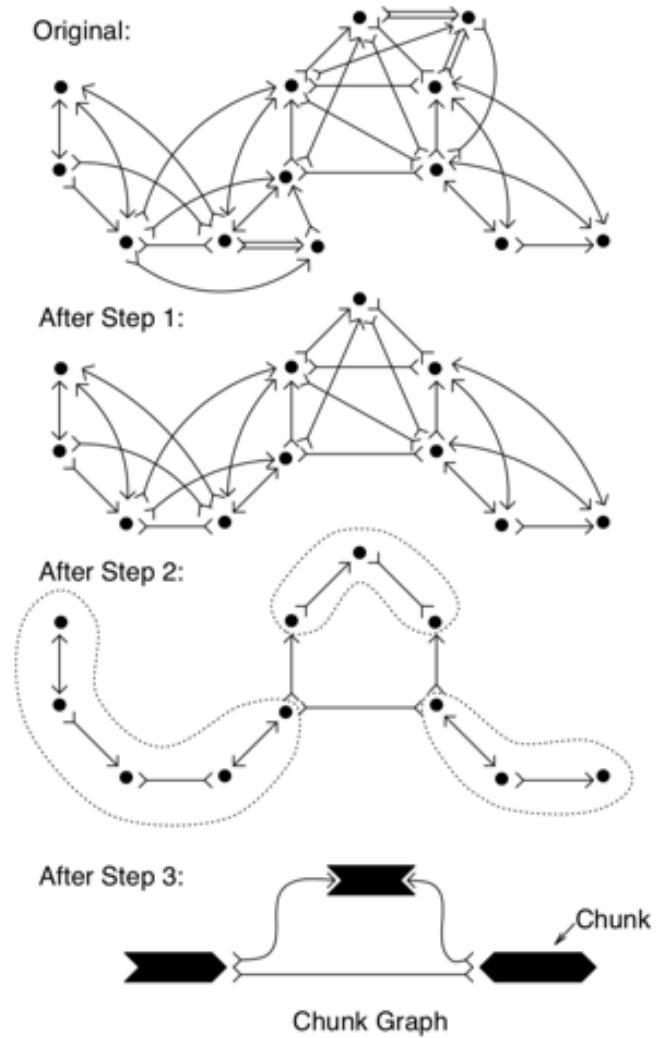
Unitigging: Pruning the Overlap Graph

The overlap graph has many redundant edges:

- If the average coverage is D, we should expect D overlaps at the beginning of the read, and D at the end

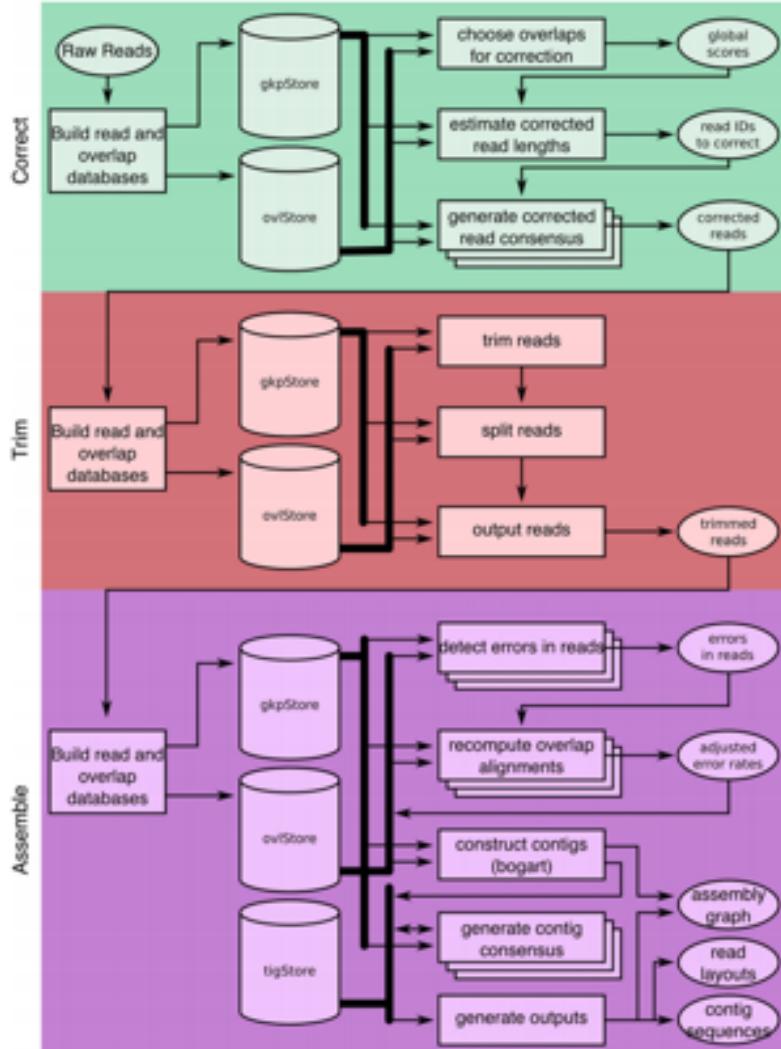
Transform the graph to simplify the assembly problem (without changing the valid solutions):

1. **Contained reads removal:** Short reads that are substrings of longer reads don't advance the assembly, remove those nodes and all of the edges
2. **Transitive edge removal:** If A \rightarrow B, and B \rightarrow C, remove the transitive edge A \rightarrow C
3. **“Chunkification”:** Linear subgraphs define uniquely assemblable segments: “unitigs”



Towards Simplifying and Accurately Formulating Fragment Assembly
Myers (1995) *J Comput Biol.* Summer;2(2):275-90.

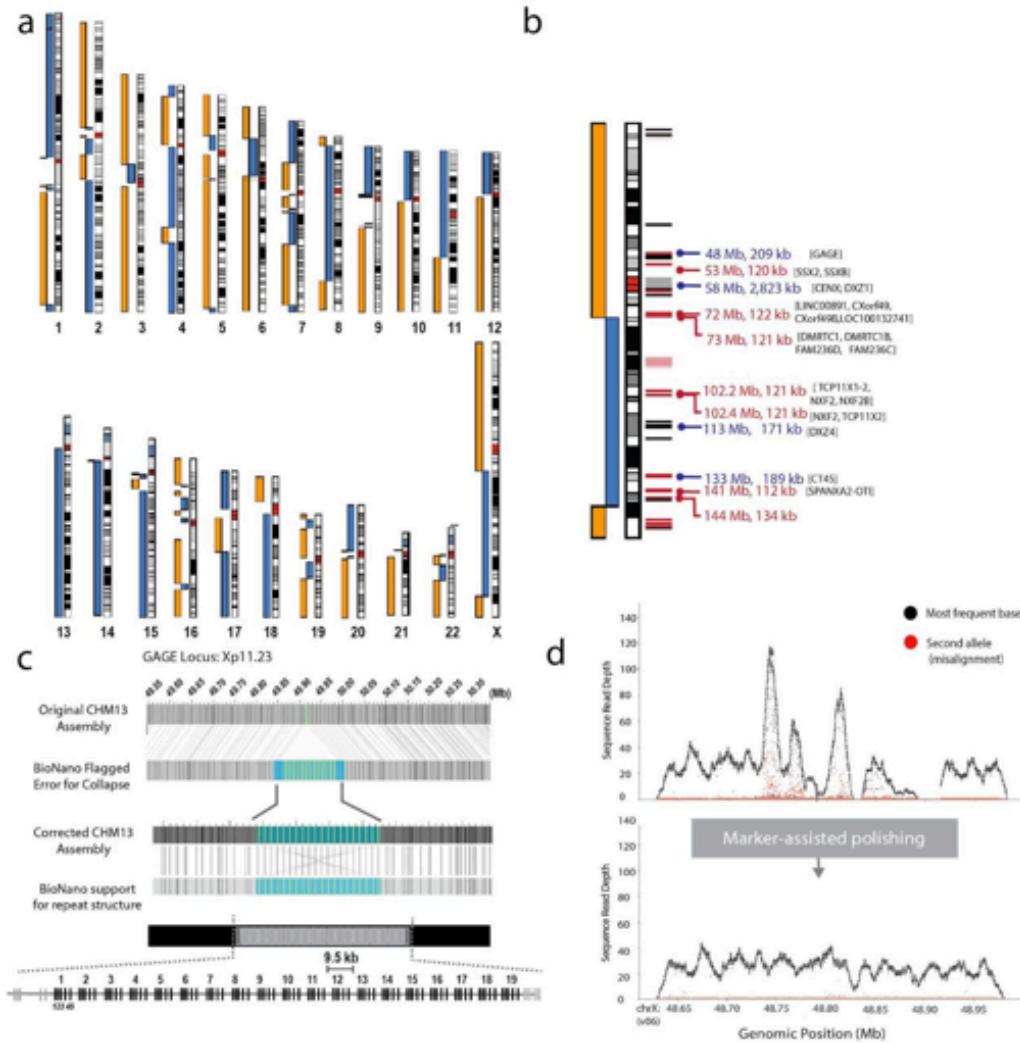
Canu Workflow



Three rounds of analysis:

- 1. Error Correction:** Use MHAP to overlap the reads, then compute a mini assembly centered around each read of good overlaps to error correct
- 2. Trim:** Use MHAP to recompute overlaps to find regions that are not well supported and discard
- 3. Unitigging:** Use Dynamic Programming to carefully overlap the error corrected reads, construct overlap graph, and then “unitig” those overlaps to build the contigs

First Telomere-to-Telomere Human Chromosome



Telomere-to-telomere assembly of a complete human X chromosome
Miga et al. (2020) Nature. <https://doi.org/10.1038/s41586-020-2547-7>

First Telomere-to-Telomere Human Genome



“We estimate that the quality of our polished assembly approaches Q70 (1 error per 10 million bases), with no known structural errors.”

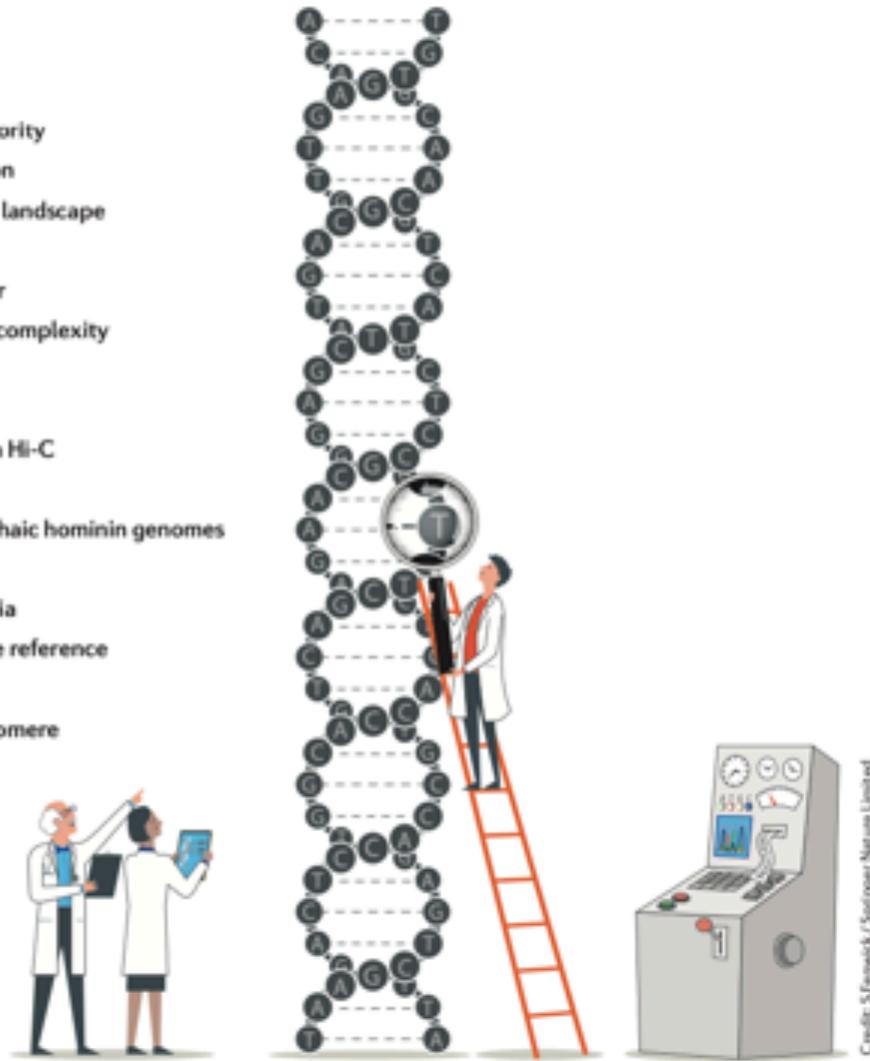
The (near) complete sequence of a human genome
<https://genomeinformatics.github.io/CHM13v1/>

nature

Genomic sequencing

MILESTONES

- S3 Foreword
- S4 Timeline
- S5 The Human Genome Project
- S6 Sequencing the unculturable majority
- S7 Sequencing — the next generation
- S8 ChIP-seq captures the chromatin landscape
- S9 The dawn of personal genomes
- S10 A sequencing revolution in cancer
- S11 Transcriptomes — a new layer of complexity
- S12 Long reads become a reality
- S13 Exploring whole exomes
- S14 Probing nuclear architecture with Hi-C
- S15 Sequencing one cell at a time
- S16 Waking the dead: sequencing archaic hominin genomes
- S17 Cataloguing a public genome
- S18 Our most elemental encyclopaedia
- S19 Pan-genomes: moving beyond the reference
- S20 Genomes go platinum
- S21 Filling in the gaps telomere to telomere



Credit: S. Emrich / Springer Nature Limited

Produced by:
Nature, Nature Genetics and
Nature Review

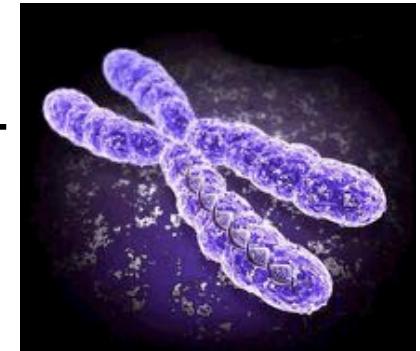
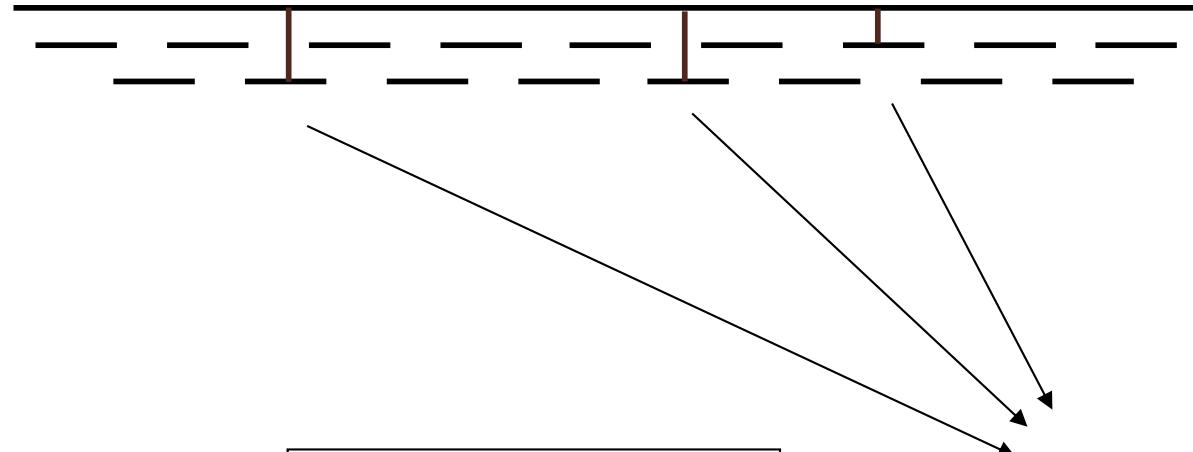
• II •

<https://www.nature.com/immersive/d42859-020-00099-0/index.html>

Read Mapping

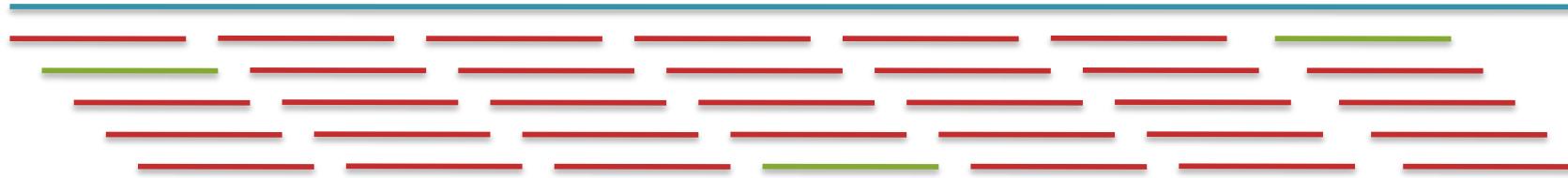
Personal Genomics

How does your genome compare to the reference?



Heart Disease —
Cancer —
Presidential smile —

Brute Force Analysis



- Brute Force:
 - At every possible offset in the genome:
 - Do all of the characters of the query match?
- Analysis
 - Simple, easy to understand
 - Genome length = n [3B]
 - Query length = m [7]
 - Comparisons: $(n-m+1) * m$ [21B]
- Overall runtime: $O(nm)$
 - [How long would it take if we double the genome size, read length?]
 - [How long would it take if we double both?]

Brute Force Reflections

Why check every position?

- GATTACA can't possibly start at position 15

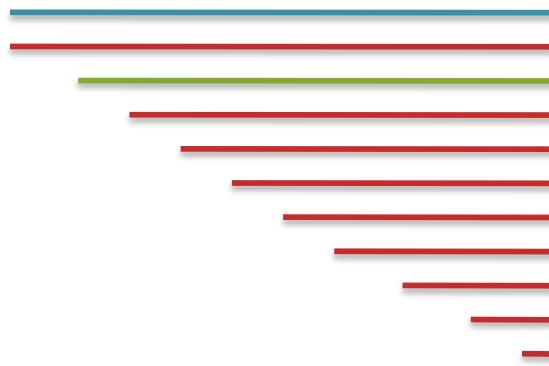
[WHY?]

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	...
T	G	A	T	T	A	C	A	G	A	T	T	A	C	C	...
								G	A	T	T	A	C	A	

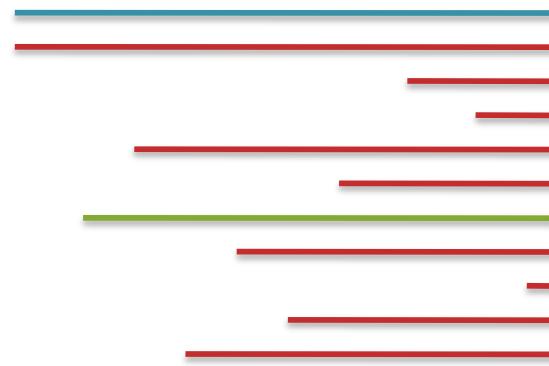
- Improve runtime to $O(n + m)$ [3B + 7]
 - If we double both, it just takes twice as long
 - Knuth-Morris-Pratt, 1977
 - Boyer-Moyer, 1977, 1991
- For one-off scans, this is the best we can do (optimal performance)
 - We have to read every character of the genome, and every character of the query
 - For short queries, runtime is dominated by the length of the genome

Suffix Arrays: Searching the Phone Book

- What if we need to check many queries?
 - We don't need to check every page of the phone book to find 'Schatz'
 - Sorting alphabetically lets us immediately skip 96% (25/26) of the book *without any loss in accuracy*
- Sorting the genome: Suffix Array (Manber & Myers, 1991)
 - Sort every suffix of the genome



Split into n suffixes



Sort suffixes alphabetically

[Challenge Question: How else could we split the genome?]

Searching the Index

- Strategy 2: Binary search
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - Lo = 1; Hi = 15;

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Searching the Index

- Strategy 2: Binary search
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - $\text{Lo} = 1; \text{Hi} = 15; \text{Mid} = (1+15)/2 = 8$
 - Middle = Suffix[8] = CC

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Searching the Index

- Strategy 2: Binary search
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - $\text{Lo} = 1; \text{Hi} = 15; \text{Mid} = (1+15)/2 = 8$
 - Middle = Suffix[8] = CC
=> Higher: Lo = Mid + 1

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Lo →

Hi →

Searching the Index

- Strategy 2: Binary search
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - Lo = 1; Hi = 15; Mid = $(1+15)/2 = 8$
 - Middle = Suffix[8] = CC
=> Higher: Lo = Mid + 1
 - Lo = 9; Hi = 15;

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Lo
→

Hi
→

Searching the Index

- Strategy 2: Binary search
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - Lo = 1; Hi = 15; Mid = $(1+15)/2 = 8$
 - Middle = Suffix[8] = CC
=> Higher: Lo = Mid + 1
 - Lo = 9; Hi = 15; Mid = $(9+15)/2 = 12$
 - Middle = Suffix[12] = TACC

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Lo
→

Hi
→

Searching the Index

- Strategy 2: Binary search
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - Lo = 1; Hi = 15; Mid = $(1+15)/2 = 8$
 - Middle = Suffix[8] = CC
=> Higher: Lo = Mid + 1
 - Lo = 9; Hi = 15; Mid = $(9+15)/2 = 12$
 - Middle = Suffix[12] = TACC
=> Lower: Hi = Mid - 1
 - Lo = 9; Hi = 11;

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Lo
→

Hi
→

Searching the Index

- Strategy 2: Binary search
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - Lo = 1; Hi = 15; Mid = $(1+15)/2 = 8$
 - Middle = Suffix[8] = CC
=> Higher: Lo = Mid + 1
 - Lo = 9; Hi = 15; Mid = $(9+15)/2 = 12$
 - Middle = Suffix[12] = TACC
=> Lower: Hi = Mid - 1
 - Lo = 9; Hi = 11; Mid = $(9+11)/2 = 10$
 - Middle = Suffix[10] = GATTACC

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Lo

Hi

Searching the Index

- Strategy 2: Binary search
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - Lo = 1; Hi = 15; Mid = $(1+15)/2 = 8$
 - Middle = Suffix[8] = CC
=> Higher: Lo = Mid + 1
 - Lo = 9; Hi = 15; Mid = $(9+15)/2 = 12$
 - Middle = Suffix[12] = TACC
=> Lower: Hi = Mid - 1
 - Lo = 9; Hi = 11; Mid = $(9+11)/2 = 10$
 - Middle = Suffix[10] = GATTACC
=> Lower: Hi = Mid - 1
 - Lo = 9; Hi = 9;



#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Searching the Index

- Strategy 2: Binary search
 - Compare to the middle, refine as higher or lower
- Searching for GATTACA
 - Lo = 1; Hi = 15; Mid = $(1+15)/2 = 8$
 - Middle = Suffix[8] = CC
=> Higher: Lo = Mid + 1
 - Lo = 9; Hi = 15; Mid = $(9+15)/2 = 12$
 - Middle = Suffix[12] = TACC
=> Lower: Hi = Mid - 1
 - Lo = 9; Hi = 11; Mid = $(9+11)/2 = 10$
 - Middle = Suffix[10] = GATTACC
=> Lower: Hi = Mid - 1
 - Lo = 9; Hi = 9; Mid = $(9+9)/2 = 9$
 - Middle = Suffix[9] = GATTACA...
=> Match at position 2!

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Lo Hi

Binary Search Analysis

- Binary Search

 Initialize search range to entire list

$\text{mid} = (\text{hi}+\text{lo})/2$; $\text{middle} = \text{suffix}[\text{mid}]$

 if query matches middle: done

 else if query < middle: pick low range

 else if query > middle: pick hi range

 Repeat until done or empty range

[WHEN?]

- Analysis

- More complicated method

- How many times do we repeat?

- How many times can it cut the range in half?

- Find smallest x such that: $n/(2^x) \leq 1$; $x = \lg_2(n)$

[32]

- Total Runtime: $O(m \lg n)$

- More complicated, but much faster!

- Looking up a query loops 32 times instead of 3B

[How long does it take to search 6B or 24B nucleotides?]

Binary Search Analysis

- Binary Search

Initialize search range to entire list

$mid = (hi+lo)/2$; $middle = suffix[mid]$

if query matches middle: done

else if query < middle: pick low range

else if query > middle: pick hi range

Repeat until done or empty range

[WHEN?]

- Analysis

- More complicated method

- How many times do we repeat?

- How many times can it cut the range in half?

- Find smallest x such that: $n/(2^x) \leq 1$; $x = \lg_2(n)$

[32]

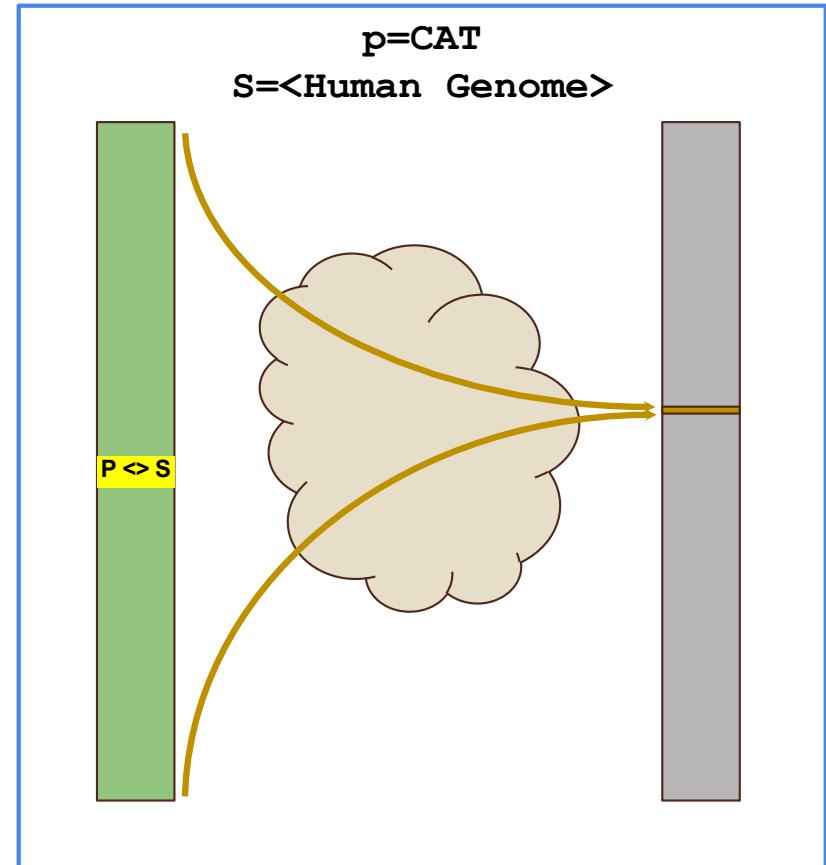
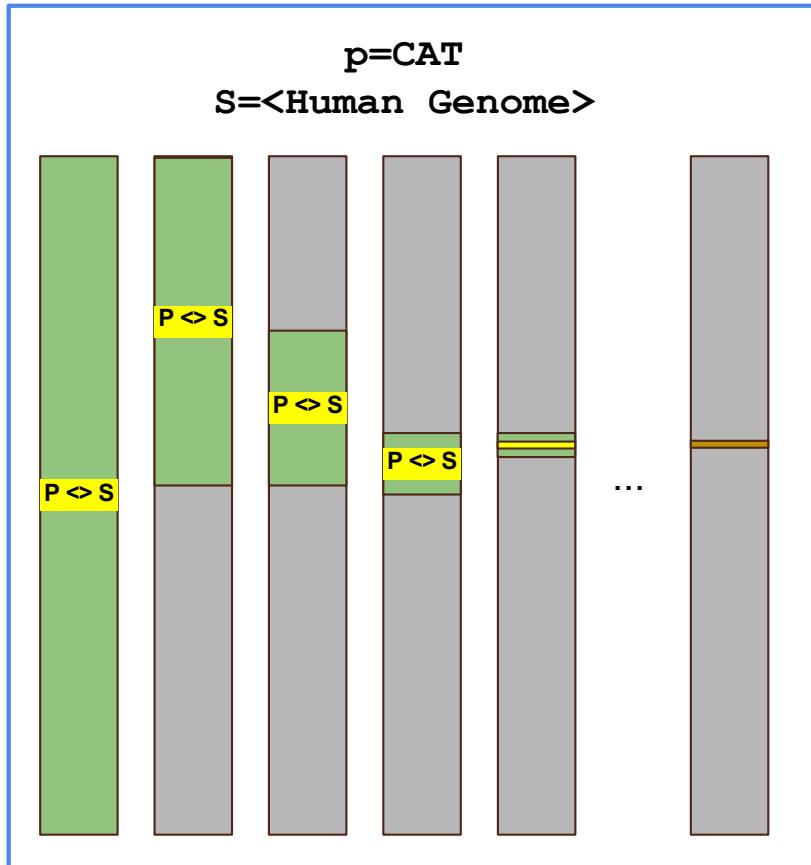
- Total Runtime: $O(m \lg n)$

- More complicated, but much faster!

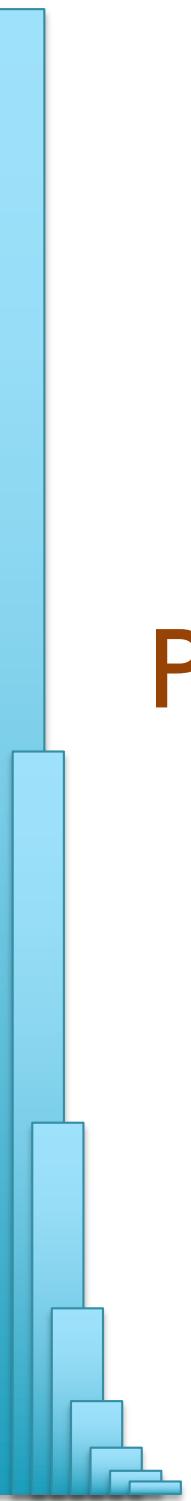
- Looking up a query loops 32 times instead of 3B

Can be reduced to $O(m + \lg n)$
using an auxiliary data structure called the LCP array

Sapling: Accelerating Suffix Array Queries with Learned Data Models



*What if instead of a slow algorithmic approach to find the correct rows,
we could somehow quickly guess/predict the correct rows?*

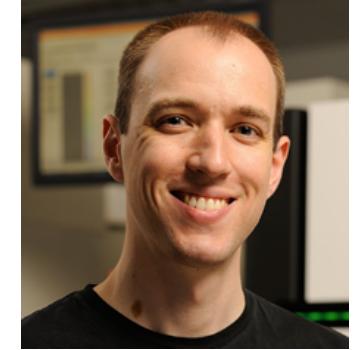


Part 2: Burrows Wheeler Transform

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?

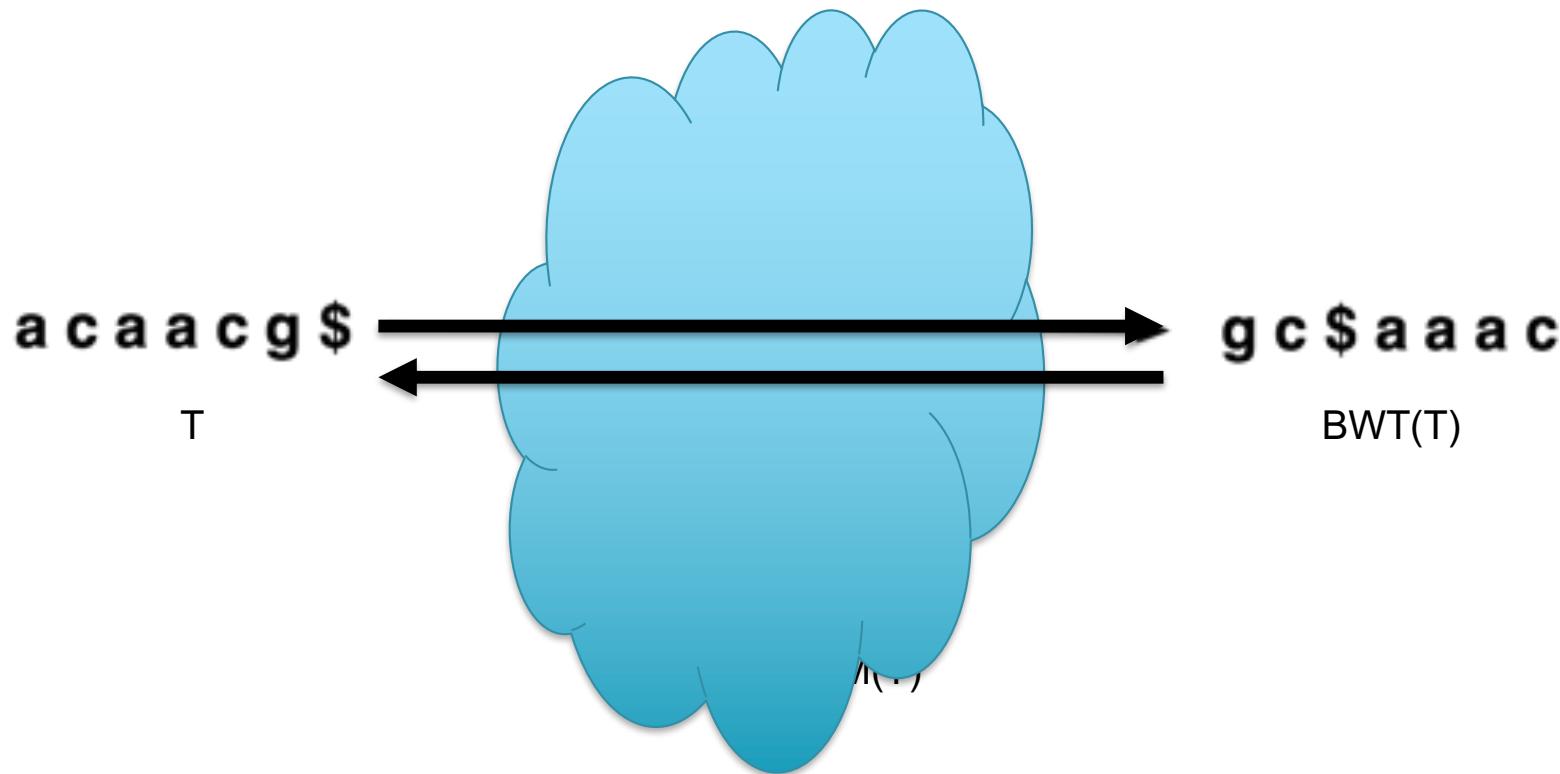


Bowtie: Ultrafast and memory efficient alignment of short DNA sequences to the human genome

Slides Courtesy of Ben Langmead

Burrows-Wheeler Transform

- 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

Burrows-Wheeler Transform

- Permutation of the characters in a text

a c a a c g \$ →

T

A block sorting lossless data compression algorithm.

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

Burrows-Wheeler Transform

- Permutation of the characters in a text

a c a a c g \$
c a a c g \$ a
a a c g \$ a c
a c a a c g \$ → a c g \$ a c a
c g \$ a c a a
g \$ a c a a c
\$ a c a a c g

T

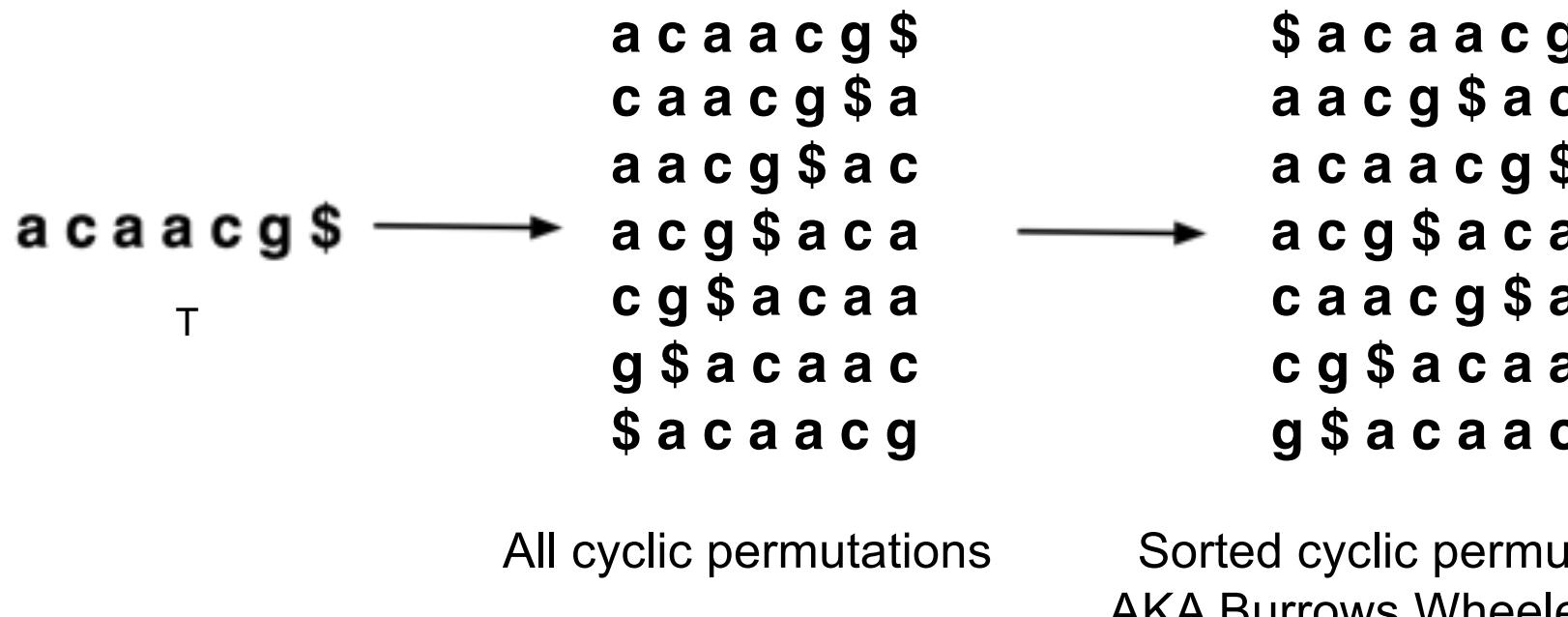
All cyclic permutations

A block sorting lossless data compression algorithm.

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

Burrows-Wheeler Transform

- 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*

Burrows-Wheeler Transform

- Permutation of the characters in a text

a c a a c g \$ →
T

\$ a c a a c g
a a c g \$ a c
a c a a c g \$
a c g \$ a c a
c a a c g \$ a
c g \$ a c a a
g \$ a c a a c

Sorted cyclic permutations
AKA Burrows Wheeler Matrix

A block sorting lossless data compression algorithm.

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

Burrows-Wheeler Transform

- Permutation of the characters in a text

a c a a c g \$ →
T

\$	a	c	a	a	c	g
a	a	c	g	\$	a	c
a	c	a	a	c	g	\$
a	c	g	\$	a	c	a
c	a	a	c	g	\$	a
c	g	\$	a	c	a	a
g	\$	a	c	a	a	c

Sorted cyclic permutations
AKA Burrows Wheeler Matrix

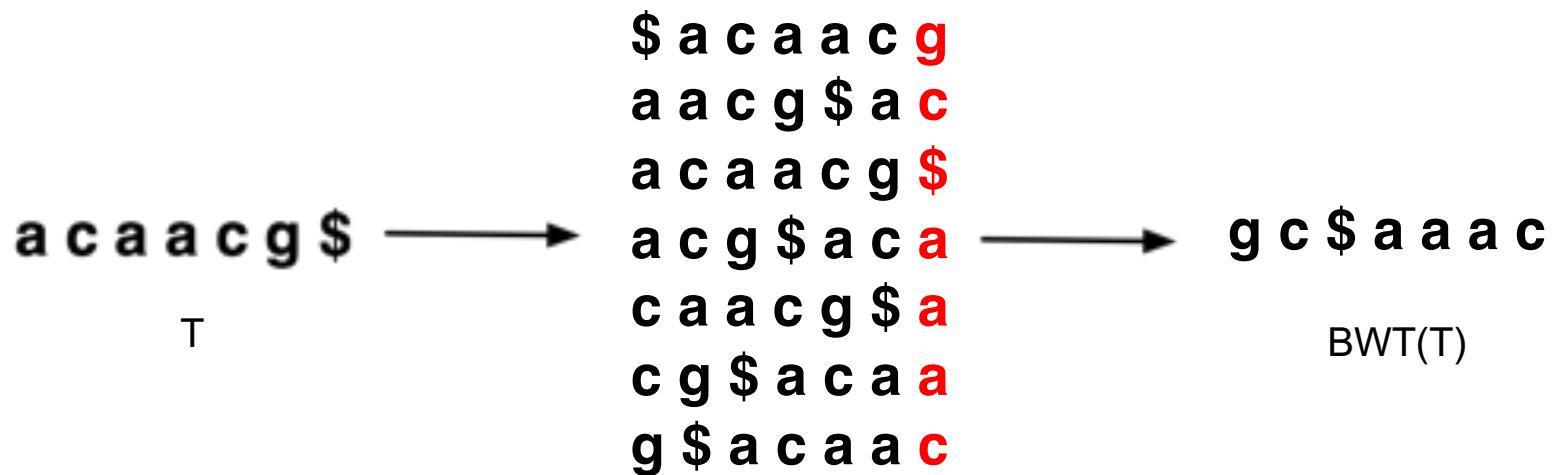
Last Column = Burrows Wheeler Transform

A block sorting lossless data compression algorithm.

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

Burrows-Wheeler Transform

- Permutation of the characters in a text



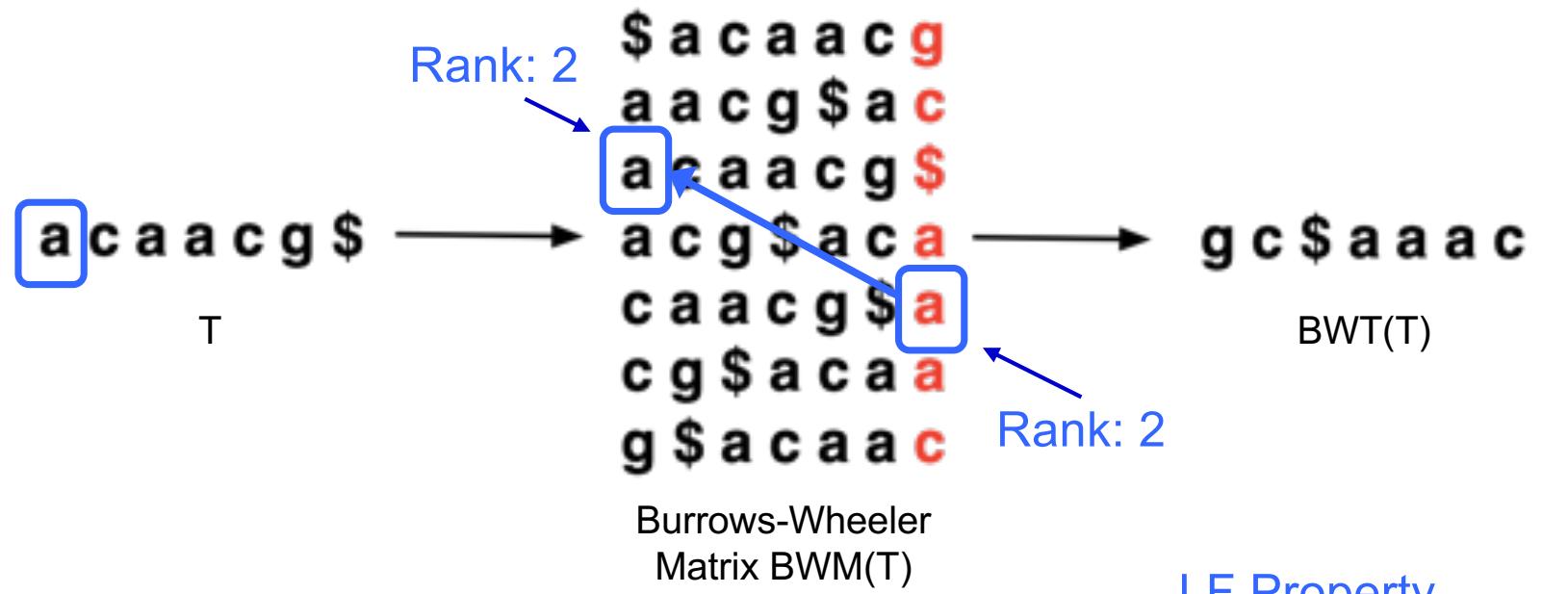
$\text{BWT}(T)$ is the index for T

A block sorting lossless data compression algorithm.

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

Burrows-Wheeler Transform

- Reversible permutation of the characters in a text



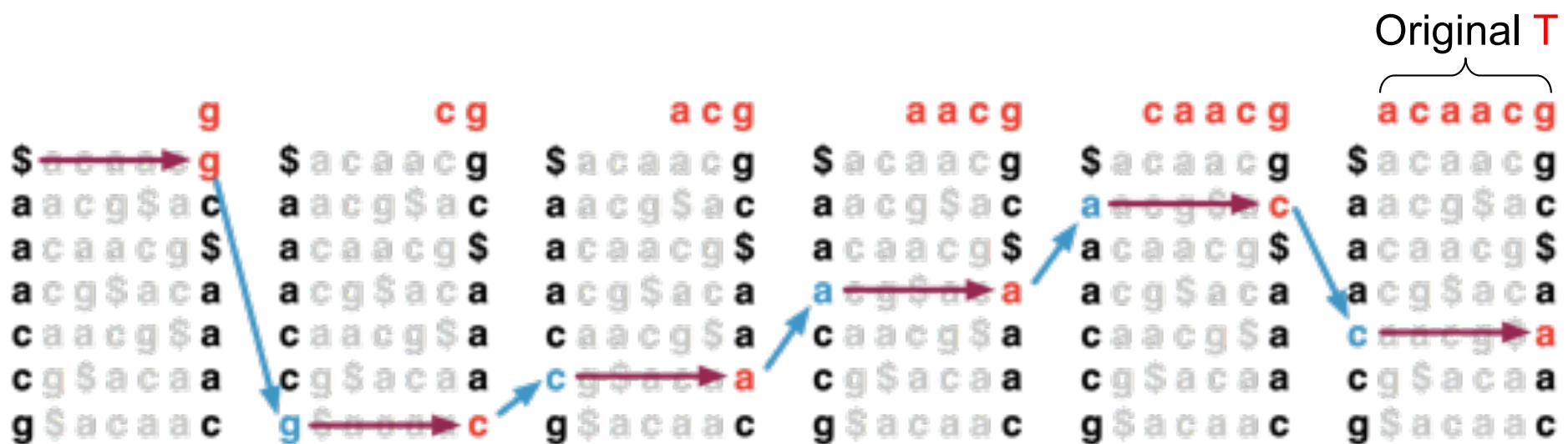
- $\text{BWT}(T)$ is the index for T

A block sorting lossless data compression algorithm.

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

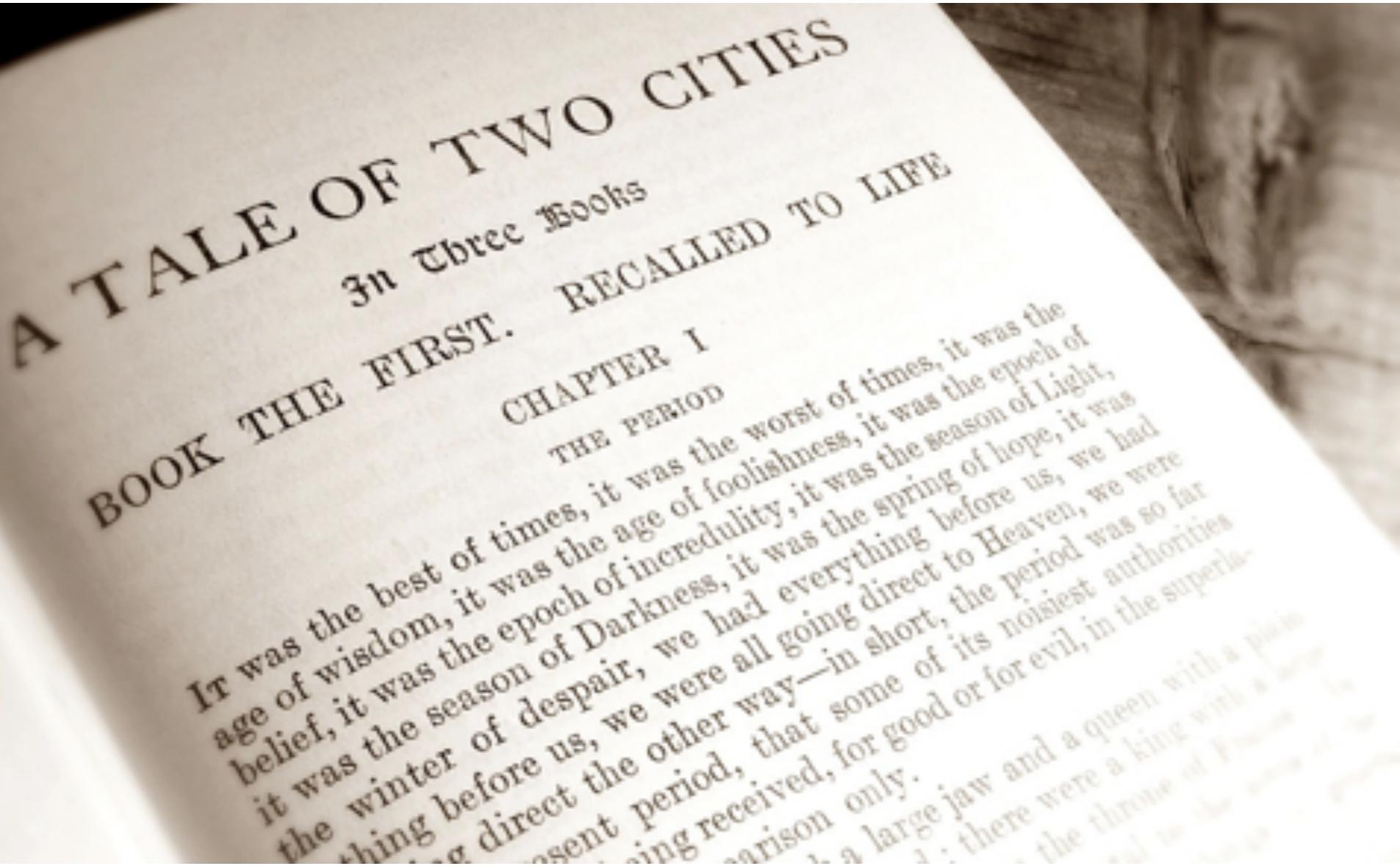
Burrows-Wheeler Transform

- Recreating T from $\text{BWT}(T)$
 - Start in the first row and apply **LF** repeatedly, accumulating predecessors along the way



[Decode this BWT string: ACTGA\\$TTA]

Run Length Encoding



Run Length Encoding

`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.$
```

Run Length Encoding:

- Replace a “run” of a character X with a single X followed by the length of the run
- GAAAAAAAATTACA => GA8T2ACA (reverse is also easy to implement)
- If your text contains numbers, then you will need to use a (slightly) more sophisticated encoding

Run Length Encoding

```
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.$
```

Run Length Encoding

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_noisiestAuthorities_insisted_on_its_being_received
,_for_good_or_for_evil,_in_the_superlative_degree_of_comparison_only.$
```

bwt[614]:

```
.dlmssftysesdtrsns_y__$yfofeeetggsfefefggeedrofr,llreef-,fs,,,
,,nfrsdnnhereghettedndeteegreenstee,ssssst,esssnssffteedtttttttr,,,
,,eeeefehh_p_fpDwwwwwwwwwwweehl_ew_____eo_neeeoaaeoo____sephrrvh
hwwegmghhhhhhkrrwwhssHrrrvtrribdbcbvs_thwwpppvmmirdnnib_eoooooo
oooooo____eennnnnnaai____ecc_ttts_tsgltsLlvtt____hhoor
e_wrraddwlors_____r_lteirillre_ouaanooioeooooiiihkiiiiiio_iei
tsppioi_____gnodsc_sss_gfhf_fffhwh_nsmo_uee_sioooaeeeeoo_ii
cgppeeaoaeooeesseutetaaaaaaaaai_ei_in_aaie_eeerei_hrsssnacciII
iiiiisn_____oyoui_a_iids_aiaee_____tlar
```

Run Length Encoding

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_noisiestAuthorities_insisted_on_its_being_received
,_for_good_or_for_evil,_in_the_superlative_degree_of_comparison_only.$
```

bwt[614]:

```
.dlmssftysesdtrsns_y__$yfofeeeetggsfefefggeedrofr,llreef-,fs,,,
,,nfrsdnnhereghettedndeteegreenstee,ssssst,esssnssffteedtttttttr,,,
,,eeeefehh_p_fpDwwwwwwwwwwweehl_ew_____eoo_neeeoaaeoo____sephrrvh
hwwegmghhhhhhkkrrwwhssHrrrvtrribdbcbvs_thwwpppvmmirdnnib_eoooooo
oooooooo____eennnnnnaai____ecc_tttttttttttts_tsgltsLlvtt__hhoor
e_wrraddwlors_____r_lteirillre_ouaanooioeooooiiihkiiiiio_iei
tsppioi_____gnodsc_sss_gfhf_fffhwh_nsmo_uee_sioooaeeeeoo_ii
cgppeeaoaeooeesseutetaaaaaaaaai_ei_in_aaie_eeerei_hrsssnacciII
iiiiisn_____ Why does the BWT tend to make runs in english text? _____tlar
```

Run Length Encoding

bwt[614]:

```
.dlmssftysesdtrsns_y$_yfofeeeetggsfefefggeedrofr,llreef-,fs,,,,,
,,nfrsdnnhereghettedndeteegreenstee,ssssst,esssnssffteedtttttttr.,
,,eefehh_p_fpDwwwwwwwwwwweehl_ew_____eo_neeeoaaeoo_____sephrrvh
hwwegmghhhhhhkrrwwhssHrrrvtrribdbcbvs__thwwpppvmmirdnnib__oooooooo
oooooooo____eennnnnnaai____ecc____tttttttttttts_tsgltsLlvtt____hhoor
e_wrraddwlors_____r_lteirillre_ouaanooioeooooiiihkiiiiio_iei
tsppioi_____gnodsc_sss_gfhf_fffhwh_nsмо_uee_siоооаеееоо_ii
cgppeeaoaeooeesseutetaaaaaaaaai_ei_in_aaie_eeerei_hrsssnacciII
iiiiisn_____oyoui_a_iids_aiaee_____tlar
```

rle(bwt)[464]:

```
.dlms2ftysesdtrsns_y_2$yfofe4tg2sfefefg2e2drofr,l2re2f-,fs,9nfrsdn2
hereghet2edndete2ge2nste2,s5t,es3ns2f2te2dt10r,4e3feh2_2p_2fpDw11e2h
1_ew_5eo2_ne3oa2eo2_4seph2r2hvh2w2egmgh7kr2w2h2s2Hr3vtr2ib2dbcbvs_2t
hw2p3vm2irdn2ib_2eo12_4e2n6a2i_3ec2_2t18s_tsgltsLlv2_3h2o2re_wr2ad2
wlors_9r_2lteiril2re_oua2no2i2oeo4i3hki6o_2ieitsp2ioi_12g2nodsc_s3_g
fhf_f3hwh_nsмо_2ue2_siо3ae4o2_i2cgp2e2aoaeo2e2s2eu2tetallи_2ei_in_2a
2ie_e3rei_hrs3nac2i2Ii7sn_15oyoui_2a_i3ds_2ai2ae2_21tlar
```

Run Length Encoding

bwt[614]:

```
.dlmssftysesdtrsns_y__$yfofeeeetggsfefefggeedrofr,1lreef-,fs,,,,,
,,nfrsdnnhereghettedndeteegreenstee,ssssst,esssnssffteedtttttttttr.,
,,eefehh_p_fpDWWWWWWWWWWWWWeehl_ew_____eo_neeeoaaeoo_____sephrrvh
hwwegmghhhhhhkrrwwhssHrrrvtrribdbcbvs__thwwpppvmmirdnnib__oooooooo
oooooooo____eennnnnnaai____ecc__ttttttttttttts_tsgltsLlvtt____hhoor
e_wrraddwlors_____r_lteirillre_ouaanooioeooooiiihkiiiiio_iei
tsppioi_____gnodsc_sss_gfhf_fffhwh_nsмо_uee_siоооаеееоо_ii
cgppeeaoaeooeesseuutetaaaaaaaaai__ei_in__aaie_eeerei_hrsssnacciII
iiiiisn_____oyoui_a_iids_aiaee_____tlar
```

rle(bwt)[464]:

```
.dlms2ftysesdtrsns_y_2$yfofe4tg2sfefefg2e2drofr,l2re2f-,fs,9nfrsdn2
hereghet2edndete2ge2nste2,s5t,es3ns2f2te2dt10r,4e3feh2_2p_2fpDw11e2h
1_ew_5eo2_ne3oa2eo2_4seph2r2hvh2w2egmgh7kr2w2h2s2Hr3vtr2ib2dbcbvs_2t
hw2p3vm2irdn2ib_2eo12_4e2n6a2i_3ec2_2t18s_tsgltsLlv2_3h2o2re_wr2ad2
wlors_9r_2lteiril2re_oua2no2i2oeo4i3hki6o_2ieitsp2ioi_12g2nodsc_s3_g
fhf_f3hwh_nsмо_2ue2_siо3ae4o2_i2cgp2e2aoaeo2e2s2eu2tet11i_2ei_in_2a
2ie_e3rei_hrs3nac2i2Ii7sn_15oyoui_2a_i3ds_2ai2ae2_21tlar
```

Run Length Encoding

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_noisiestAuthorities_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_i2cgp2e2aoaeo2e2s2eu2tet11i_2ei_in_2a
2ie_e3rei_hrs3nac2i2Ii7sn_15oyoui_2a_i3ds_2ai2ae2_21tlar
```

Run Length Encoding

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_noisiestAuthorities_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_nsмо_2ue2_sio3ae4o2_i2cgp2e2aoaeo2e2s2eu2tet11i_2ei_in_2a
2ie_e3rei
```

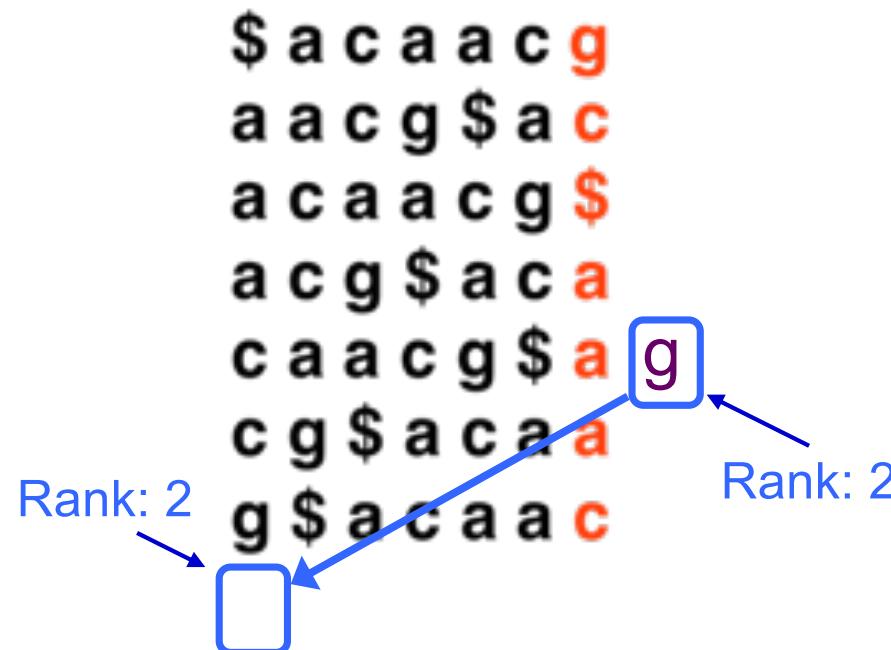
Saved 614-464 = 150 bytes (24%) with zero loss of information!

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

BWT Exact Matching

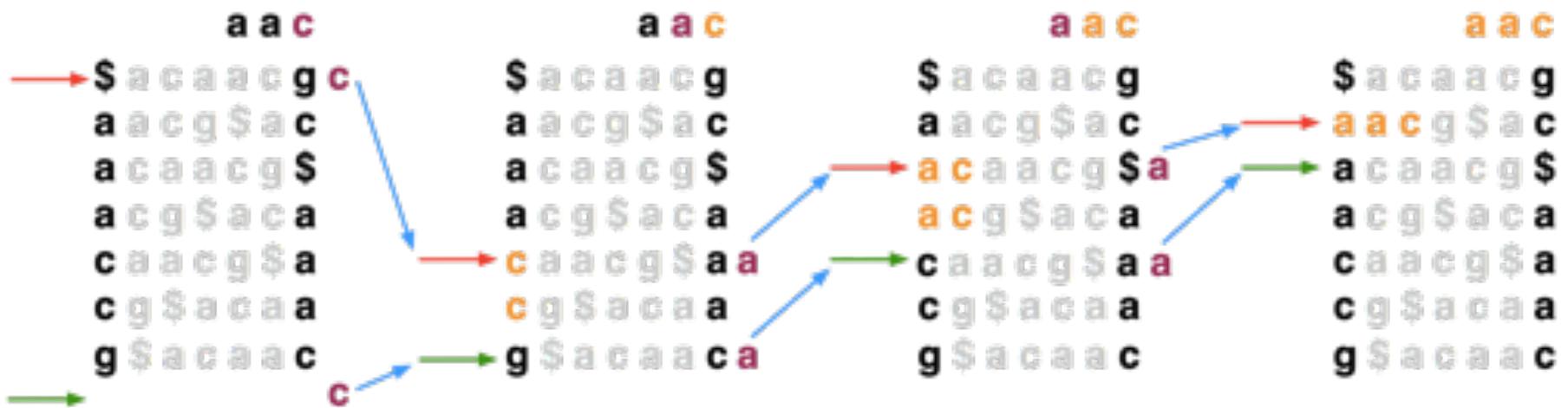
- $\text{LF}_c(r, c)$ does the same thing as $\text{LF}(r)$ but it ignores r 's actual final character and “pretends” it's c :

$$\text{LF}_c(5, g) = 8$$



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

1. Split read into segments

Read
CCAGTAGCTCTCAGCCTTATTTACCCAGGCCTGTA Read (reverse complement)
TACAGGCCTGGGTAAAATAAGGCTGAGAGCTACTGG

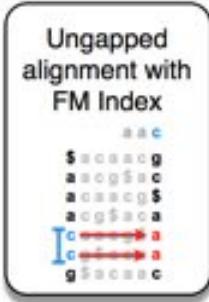
Policy: extract 16 nt seed every 10 nt

Seeds

+ , 0: CCAGTAGCTCTCAGCC	- , 0: TACAGGCCTGGGTAAA
+ , 10: TCAGCCTTATTTACC	- , 10: GGTAAAATAAGGCTGA
+ , 20: TTTACCCAGGCCTGTA	- , 20: GGCTGAGAGCTACTGG

2. Lookup each segment and prioritize

Seeds

+ , 0: CCAGTAGCTCTCAGCC	→	Ungapped alignment with FM Index	→	Seed alignments (as B ranges)
+ , 10: TCAGCCTTATTTACC				{ [211, 212], [212, 214] }
+ , 20: TTTACCCAGGCCTGTA				{ [653, 654], [651, 653] }
- , 0: TACAGGCCTGGGTAAA				{ [684, 685] }
- , 10: GGTAAAATAAGGCTGA				{ }
- , 20: GGCTGAGAGCTACTGG				{ }

3. Evaluate end-to-end match

Extension candidates

SA:684, chr12:1955	→	SIMD dynamic programming aligner	→	SAM alignments
SA:624, chr2:462				r1 0 chr12 1936 0
SA:211: chr4:762				36M * 0 0
SA:213: chr12:1935				CCAGTAGCTCTCAGCCTTATTTACCCAGGCCTGTA
SA:652: chr12:1945				IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII

(Langmead & Salzberg, 2012)