

# Read Mapping

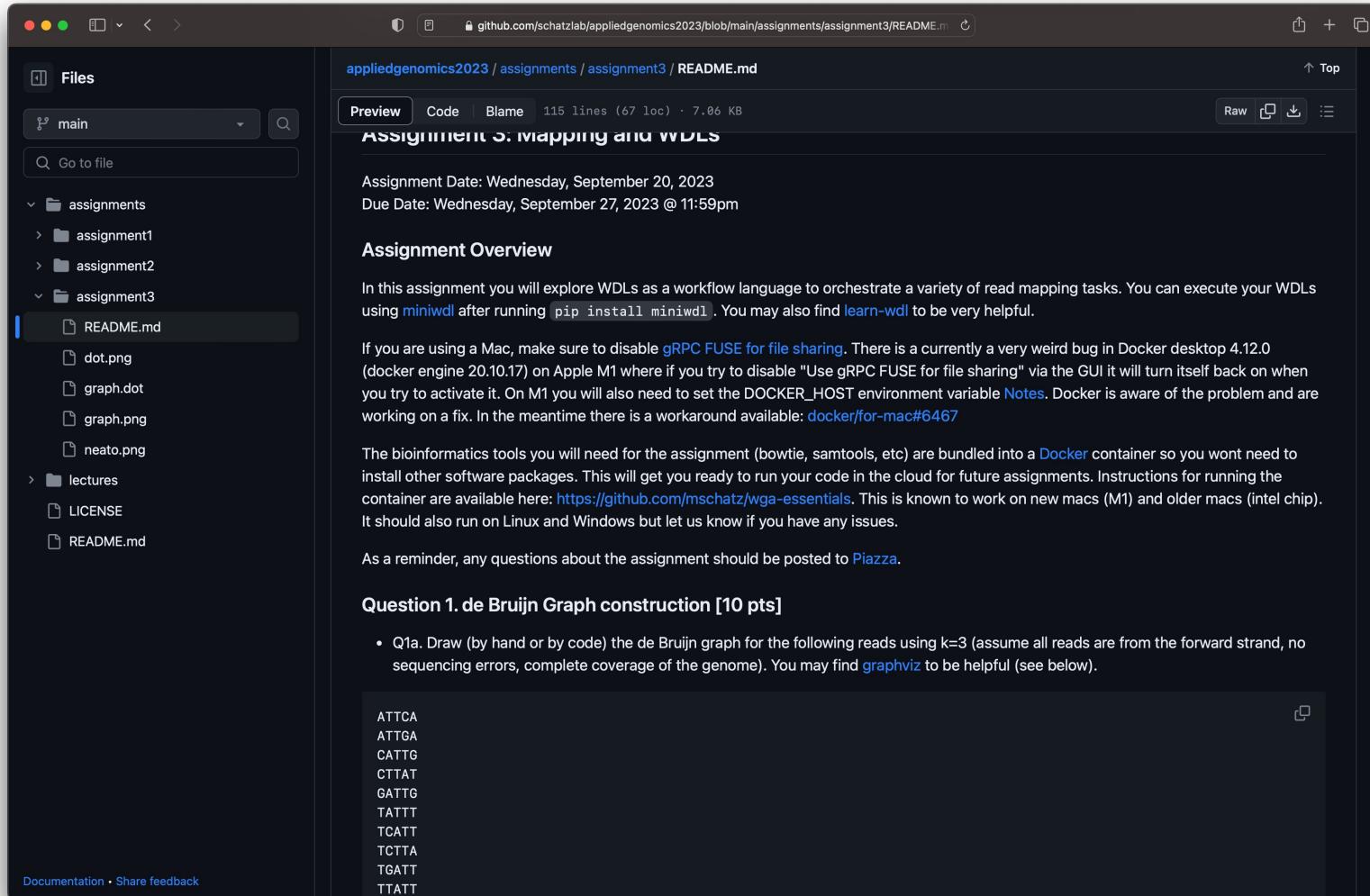
Michael Schatz

Sept 25, 2023  
Lecture 8: Applied Comparative Genomics



# Assignment 3: Genome Assembly

## Due Wednesday Sept 27 by 11:59pm



The screenshot shows a GitHub repository page for `appliedgenomics2023 / assignments / assignment3 / README.md`. The left sidebar displays a file tree with the following structure:

- main
- assignments
  - assignment1
  - assignment2
  - assignment3
    - README.md (selected)
    - dot.png
    - graph.dot
    - graph.png
    - neato.png
  - lectures
  - LICENSE
  - README.md

The README.md file content is as follows:

### Assignment 3: Mapping and WDLs

Assignment Date: Wednesday, September 20, 2023  
Due Date: Wednesday, September 27, 2023 @ 11:59pm

#### Assignment Overview

In this assignment you will explore WDLs as a workflow language to orchestrate a variety of read mapping tasks. You can execute your WDLs using `miniwdl` after running `pip install miniwdl`. You may also find `learn-wdl` to be very helpful.

If you are using a Mac, make sure to disable `gRPC FUSE for file sharing`. There is a currently a very weird bug in Docker desktop 4.12.0 (docker engine 20.10.17) on Apple M1 where if you try to disable "Use gRPC FUSE for file sharing" via the GUI it will turn itself back on when you try to activate it. On M1 you will also need to set the `DOCKER_HOST` environment variable [Notes](#). Docker is aware of the problem and are working on a fix. In the meantime there is a workaround available: [docker/for-mac#6467](#)

The bioinformatics tools you will need for the assignment (bowtie, samtools, etc) are bundled into a [Docker](#) container so you wont need to install other software packages. This will get you ready to run your code in the cloud for future assignments. Instructions for running the container are available here: <https://github.com/mschatz/wga-essentials>. This is known to work on new macs (M1) and older macs (intel chip). It should also run on Linux and Windows but let us know if you have any issues.

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

#### Question 1. de Bruijn Graph construction [10 pts]

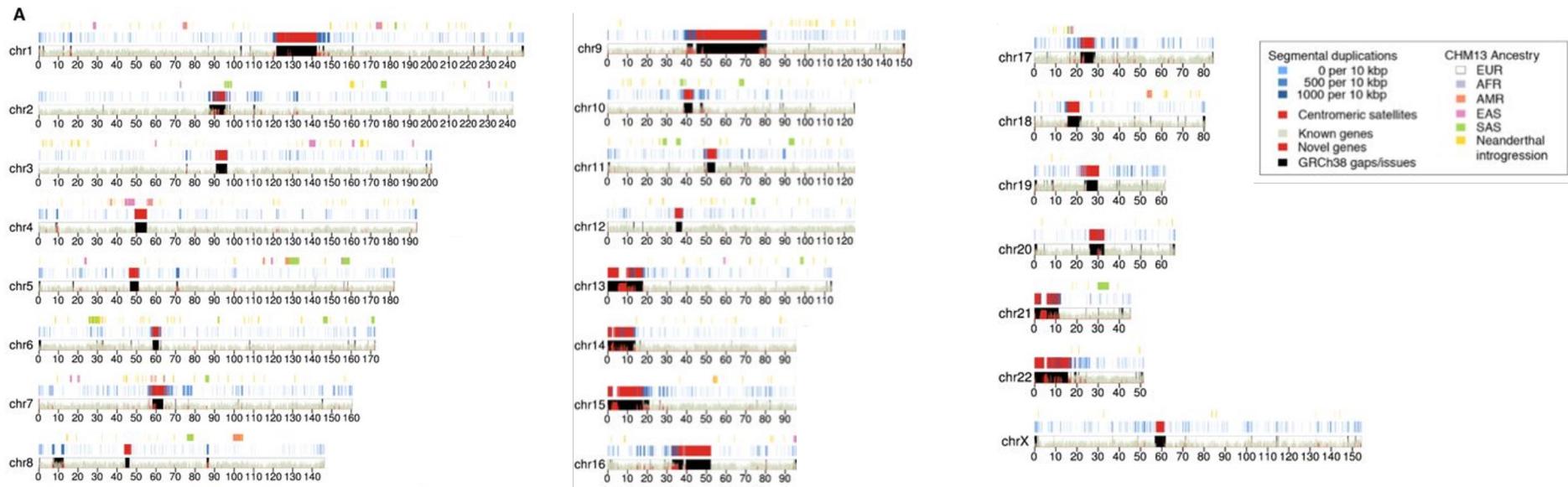
- Q1a. Draw (by hand or by code) the de Bruijn graph for the following reads using k=3 (assume all reads are from the forward strand, no sequencing errors, complete coverage of the genome). You may find [graphviz](#) to be helpful (see below).

ATTCA
ATTGA
CATTG
CTTAT
GATTG
TATTT
TCATT
TCTTA
TGATT
TTATT

<https://github.com/schatzlab/appliedgenomics2023/tree/main/assignments/assignment3>

Check Piazza for questions!

# The complete sequence of a human genome



CHM13v1.1 genome size is **3.057 Gbp with zero Ns**  
Every chromosome is telomere-to-telomere, quality estimated >Q70  
~190 Mbp (~8%) of new sequence vs. GRCh38, fixes thousands of errors

(Nurk et al. Science, 2022)

# T2T Variants: 1000 Genomes Project



26 populations from 5 superpopulations (continental regions)  
3202 samples (2504 core genomes + 698 offspring)

3202 samples x 30Gb = 96Tb input data | >5Pb of intermediate data | >>1M core hours

**Cell**

**High-coverage whole-genome sequencing of the expanded 1000 Genomes Project cohort including 602 trios**

**Graphical abstract**

The graphical abstract illustrates the expansion of the 1000 Genomes Project (1KGP) to include 602 trios. It highlights the increased power for variant discovery, including small variants, INDELs, and structural variants (SVs). It also shows the generation of an improved and accessible reference imputation panel.

**Authors**  
Marta Byrska-Bishop, Uday S. Evani, Xuefang Zhao, ..., Michael E. Takowski, Giuseppe Nardizi, Michael C. Zody

**Correspondence**  
mbyrska-bishop@nygenome.org (M.B.-B.), mczody@nygenome.org (M.C.Z.)

**In brief**  
High-coverage whole-genome sequencing (WGS) of the expanded 1000 Genomes Project (1KGP) cohort including 602 trios led to the discovery of additional rare non-coding single-nucleotide variants (SNVs), as well as coding and non-coding short insertions and deletions (INDELs) and structural variants (SVs) spanning the allele frequency spectrum compared to the original 1KGP resource based primarily on low-coverage WGS.

**Highlights**

- Expansion of the 1000 Genomes Project (1KGP) resource to include 602 trios
- High-coverage whole-genome sequencing of the expanded 1KGP cohort
- Discovery of more rare SNVs as well as INDELs and SVs across the frequency spectrum
- Generation of an improved and accessible reference imputation panel

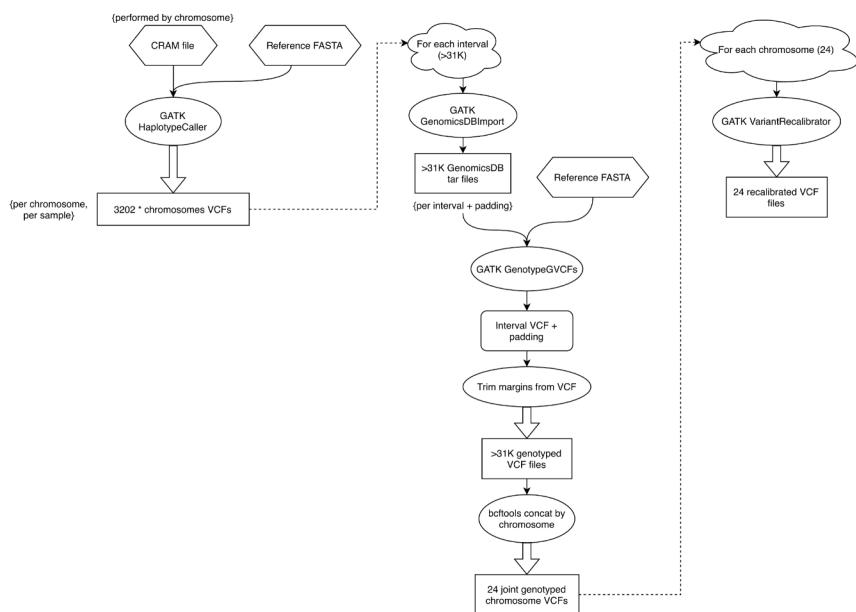
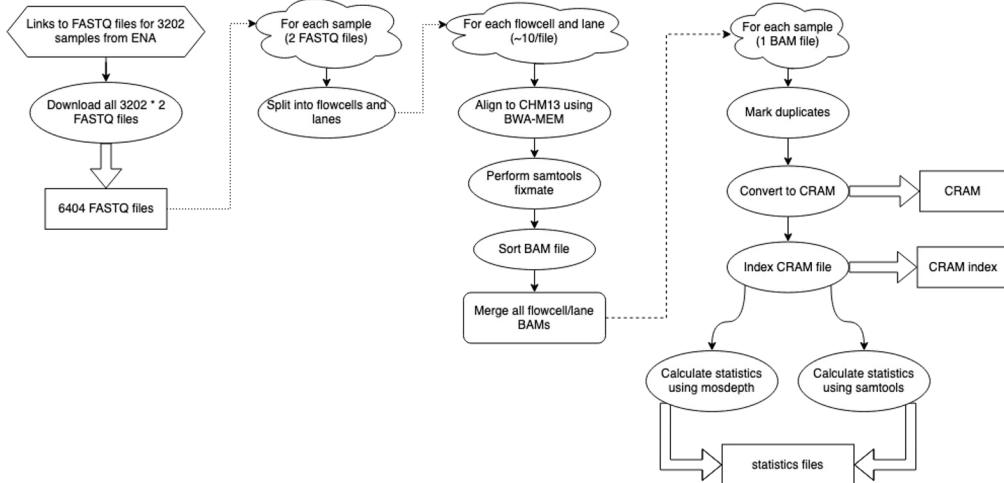
Byrska-Bishop et al., 2022, Cell 188, 3429–3440  
September 1, 2022 © 2022 The Authors. Published by Elsevier Inc.  
<https://doi.org/10.1016/j.cell.2022.08.004>

**CellPress**

# 3202 Genomes to go...



Samantha Zarate



```

task samtoolsStats {
    input {
        File inputCram
        File cramIndex
        File targetRef
        String sampleName
    }

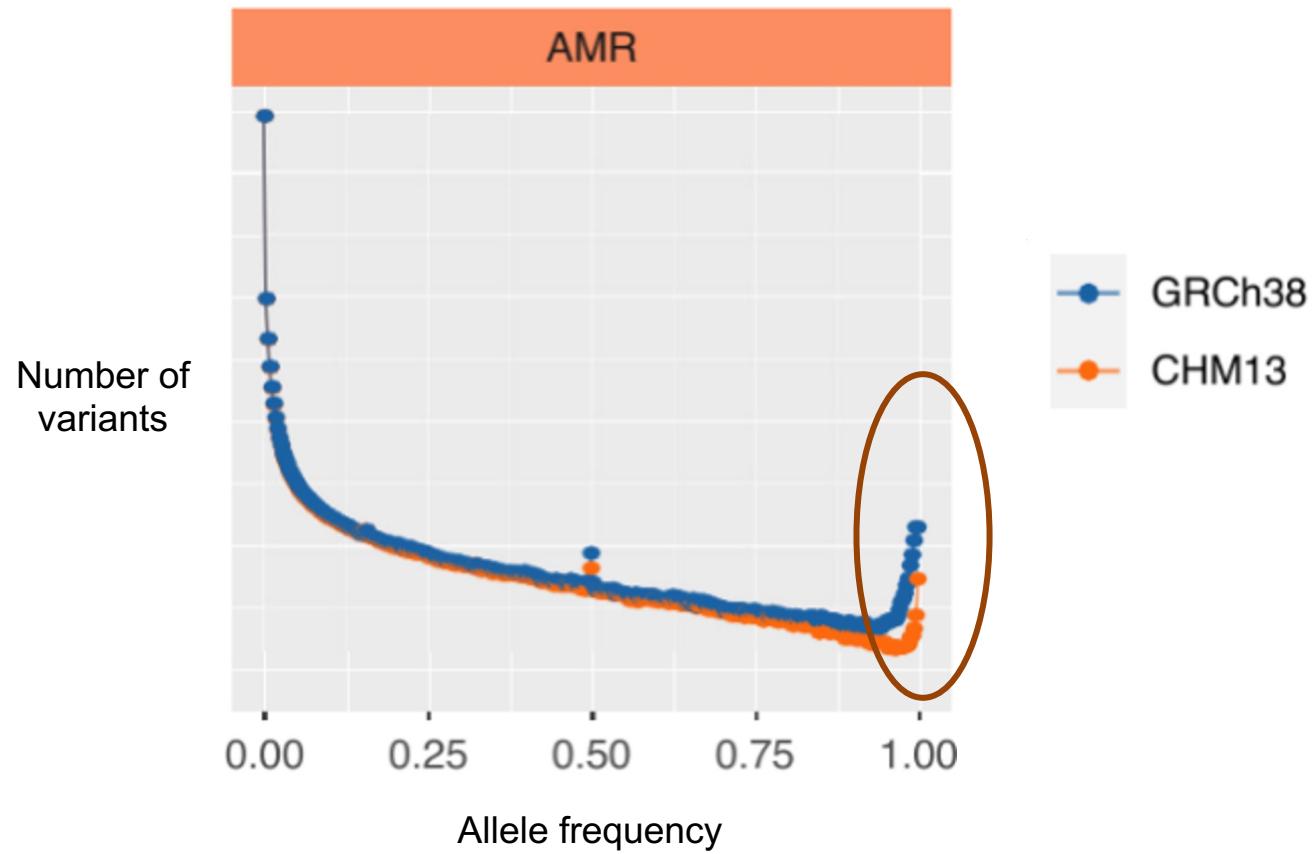
    command <<<
        samtools stats -r "~{targetRef}" \
            --reference "~{targetRef}" \
            -@ "$(nproc)" \
            "~{inputCram}" > "~{sampleName}.samtools.stats.txt"
    >>>

    Int diskGb = ceil(2.0 * size(inputCram, "G"))

    runtime {
        docker : "szarate/t2t_variants:v0.0.2"
        disks : "local-disk ${diskGb} SSD"
        memory: "12G"
        cpu : 16
        preemptible: 3
        maxRetries: 3
    }

    output {
        File stats = "~{sampleName}.samtools.stats.txt"
    }
}
  
```

# Explaining Decreased Per-Sample Count





NEWS CAREERS COMMENTARY JOURNALS ▾

**Science**

HOME > COLLECTIONS > COMPLETING THE HUMAN GENOME

**COMPLETING THE HUMAN GENOME**

A fully sequenced human genome was announced more than 20 years ago. However, owing to technological limitations, some genomic regions remained unresolved. Here, *Science* and other journals present research by the Telomere-to-Telomere (T2T) Consortium, reporting on the endeavor to complete a comprehensive human reference genome.

**FILTERS**

6 RESULTS FOUND

**SPECIAL ISSUE RESEARCH ARTICLE**  
**Segmental duplications and their variation in a complete human genome**  
BY MITCHELL R. VOLGLER, XAVI GUITART, PHILIP C. DISHUCK, LUDOVICA MERCURI, WILLIAM T. HARVEY, ARIEL GERSHMAN, MARK DIEKHANS, ARVIS SULOVARI, KATHERINE M. MUNSON, ALEXANDRA P. LEWIS, [...] EVAN E. EICHLER  
SCIENCE • VOL. 376, NO. 6588 • 01 APR 2022



**SPECIAL ISSUE RESEARCH ARTICLE**  
**Complete genomic and epigenetic maps of human centromeres**  
BY NICOLAS ALTEMOSSE, GLENNA A. LOGSDON, ANDREY V. BZIKADZE, PRAGYA SIDHWANI, SASHA A. LANGLEY, GINA V. CALDAS, SAVANNAH J. HOYT, LEV URALSKY, FEDOR D. RYABOV, COLIN J. SHEW, [...] KAREN H. MICA  
SCIENCE • VOL. 376, NO. 6588 • 01 APR 2022



**SPECIAL ISSUE RESEARCH ARTICLE**  
**From telomere to telomere: The transcriptional and epigenetic state of human repeat elements**  
BY SAVANNAH J. HOYT, JESSICA M. STORER, GABRIELLE A. HARTLEY, PATRICK G. S. GRADY, ARIEL GERSHMAN, LEONARDO G. DE LIMA, CHARLES LIMOUSE, REZA HALABIAN, LUKE WOJENSKI, MATIAS RODRIGUEZ, [...] RACHEL J.  
+16 authors • SCIENCE • VOL. 376, NO. 6588 • 01 APR 2022



**SPECIAL ISSUE RESEARCH ARTICLE**  
**A complete reference genome improves analysis of human genetic variation**  
BY SERGEY AGAEZOV, STEPHANIE M. YAN, DANIELA C. SOTO, MELANIE KIRSCH, SAMANTHA ZARATE, PAVEL AVDEYEV, DYLAN J. TAYLOR, KISHWAR SHAFIN, ALAINA SHUMATE, CHUNLIN XIAO, [...] MICHAEL C. SCHATZ  
SCIENCE • VOL. 376, NO. 6588 • 01 APR 2022

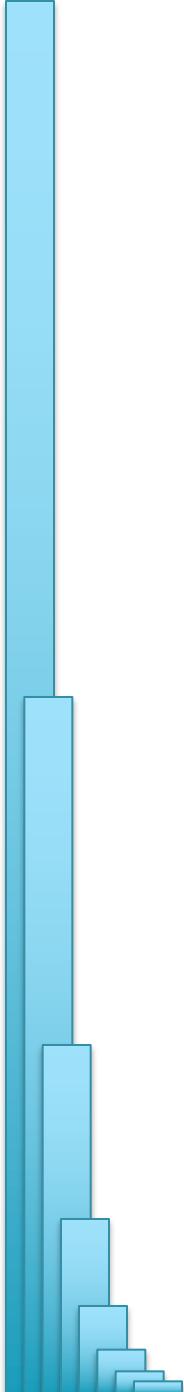


**SPECIAL ISSUE RESEARCH ARTICLE**  
**Epigenetic patterns in a complete human genome**  
BY ARIEL GERSHMAN, MICHAEL E. G. SAURIA, XAVI GUITART, MITCHELL R. VOLGLER, PAUL W. HOOK, SAVANNAH J. HOYT, MITEN JAIN, ALAINA SHUMATE, ROHAM RAZAGHI, SERGEY KOREN, [...] WINSTON TIMP  
VOL. 376, NO. 6588 • 01 APR 2022



**SPECIAL ISSUE RESEARCH ARTICLE**  
**The complete sequence of a human genome**  
BY SERGEY NURK, SERGEY KOREN, ARANG RHEE, MIKKO RAUTIAINEN, ANDREY V. BZIKADZE, ALLA MIKEHENKO, MITCHELL R. VOLGLER, NICOLAS ALTEMOSSE, LEV URALSKY, ARIEL GERSHMAN, [...] ADAM M. PHILLIPPI  
SCIENCE • VOL. 376, NO. 6588 • 31 MAR 2022 : 44-53

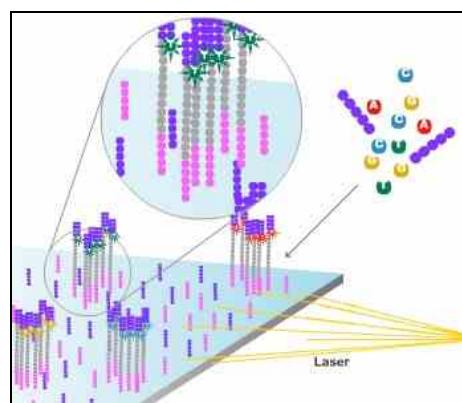
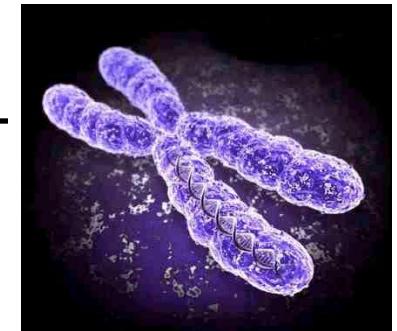
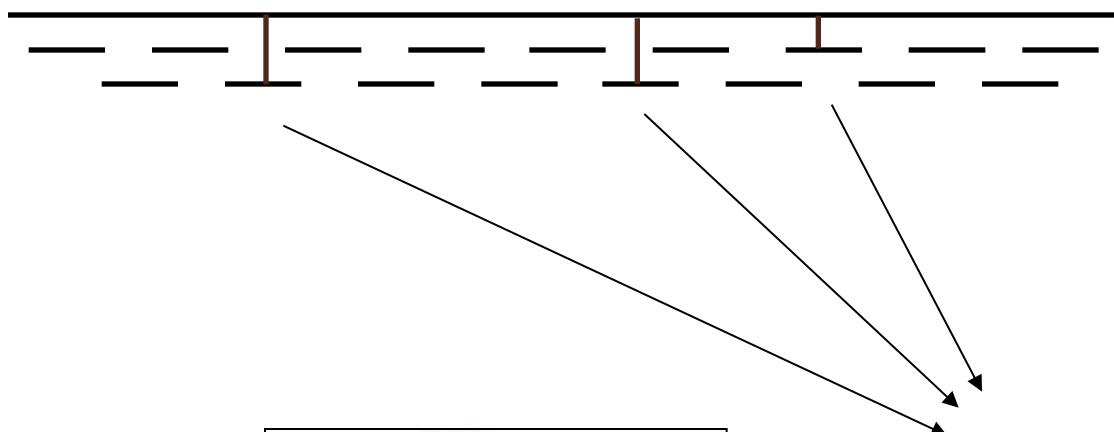




# Read Mapping

# Personal Genomics

How does your genome compare to the reference?



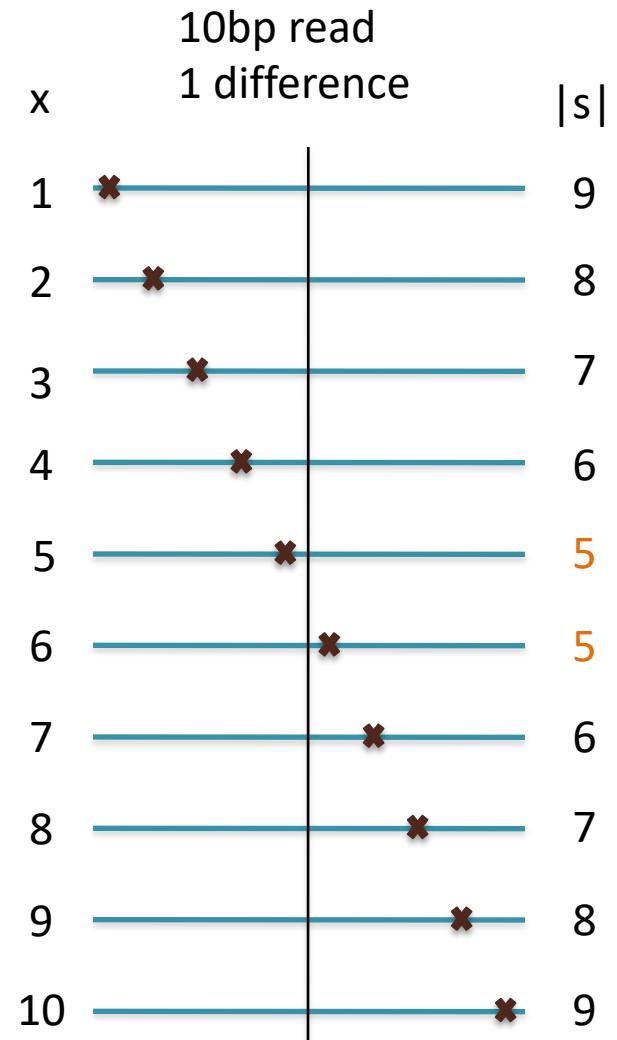
Heart Disease  
Cancer  
Presidential smile

— — —  
— — —  
— — —

# Seed-and-Extend Alignment

Theorem: An alignment of a sequence of length  $m$  with at most  $k$  differences **must** contain an exact match at least  $s=m/(k+1)$  bp long  
(Baeza-Yates and Perleberg, 1996)

- Proof: Pigeonhole principle
  - 1 pigeon can't fill 2 holes
- Seed-and-extend search
  - Use an index to rapidly find short exact alignments to seed longer in-exact alignments
    - BLAST, MUMmer, Bowtie, BWA, SOAP, ...
  - Specificity of the depends on seed length
    - Guaranteed sensitivity for  $k$  differences
    - Also finds some (but not all) lower quality alignments <- heuristic



# SAM / BAM Files

```
Coor      12345678901234 5678901234567890123456789012345
ref       AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGGCCAT

+r001/1    TTAGATAAAGGATA*CTG
+r002    aaaAGATAA*GGATA
+r003    gcctaAGCTAA
+r004    ATAGCT.....TCAGC
-r003    ttagctTAGGC
-r001/2    CAGCGGCAT
```

```
@HD VN:1.6 SO:coordinate
@SQ SN:ref LN:45
r001 99 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003 0 ref 9 30 5S6M * 0 0 GCCTAACGCTAA * SA:Z:ref,29,-,6H5M,17,0;
r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
r003 2064 ref 29 17 6H5M * 0 0 TAGGC * SA:Z:ref,9,+,5S6M,30,1;
r001 147 ref 37 30 9M = 7 -39 CAGCGGCAT * NM:i:1
```

# SAM / BAM Files

## 1.4 The alignment section: mandatory fields

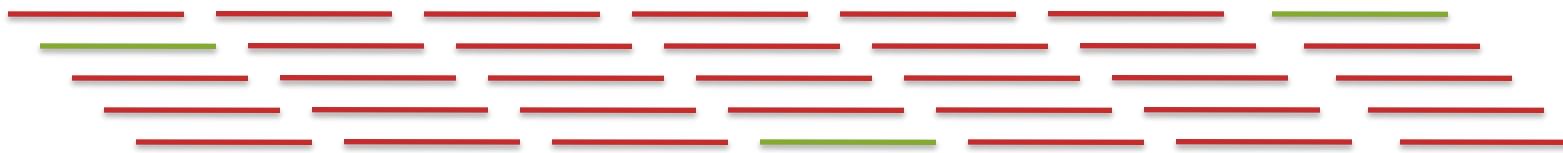
In the SAM format, each alignment line typically represents the linear alignment of a segment. Each line consists of 11 or more TAB-separated fields. The first eleven fields are always present and in the order shown below; if the information represented by any of these fields is unavailable, that field's value will be a placeholder, either '0' or '\*' as determined by the field's type. The following table gives an overview of these mandatory fields in the SAM format:

Col	Field	Type	Regexp/Range	Brief description
1	QNAME	String	[!-?A-~]{1,254}	Query template NAME
2	FLAG	Int	[0, 2 <sup>16</sup> - 1]	bitwise FLAG
3	RNAME	String	\* [:rname:^*=:] [:rname:]*	Reference sequence NAME <sup>11</sup>
4	POS	Int	[0, 2 <sup>31</sup> - 1]	1-based leftmost mapping POSition
5	MAPQ	Int	[0, 2 <sup>8</sup> - 1]	MAPping Quality
6	CIGAR	String	\* ([0-9]+[MIDNSHPX=])+	CIGAR string
7	RNEXT	String	\* = [:rname:^*=:] [:rname:]*	Reference name of the mate/next read
8	PNEXT	Int	[0, 2 <sup>31</sup> - 1]	Position of the mate/next read
9	TLEN	Int	[-2 <sup>31</sup> + 1, 2 <sup>31</sup> - 1]	observed Template LENgth
10	SEQ	String	\* [A-Za-z.=]+	segment SEQuence
11	QUAL	String	[!-~]+	ASCII of Phred-scaled base QUALity+33

<sup>11</sup> Reference sequence names may contain any printable ASCII characters with the exception of certain punctuation characters, and may not start with '\*' or '='. See Section 1.2.1 for details and an explanation of the [:rname:] notation.

# 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 =  $n$
  - Query length =  $m$
  - Comparisons:  $(n-m+1) * m$
- 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?

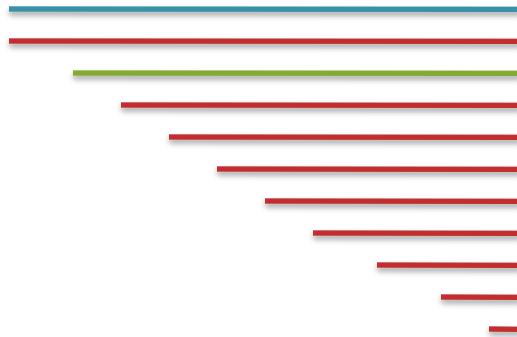
- GATTACA can't possibly start at position 15 [WHY?]

I	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	...
								<b>G</b>	<b>A</b>	<b>T</b>	<b>T</b>	<b>A</b>	<b>C</b>	<b>A</b>	

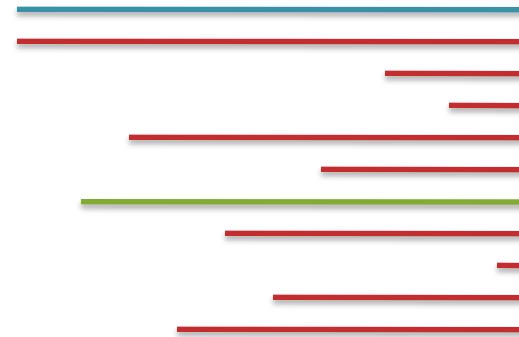
- Improve runtime to  $O(n + m)$  [3B + 7]
  - If we double both, it just takes twice as long
  - Knuth-Morris-Pratt, 1977
  - Boyer-Moyer, 1977, 1991
- For one-off scans, this is the best we can do (optimal performance)
  - We have to read every character of the genome, and every character of the query
  - 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
  - $Lo = 1; Hi = 15; Mid = (1+15)/2 = 8$
  - 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
  - $Lo = 1; Hi = 15; Mid = (1+15)/2 = 8$
  - 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

# 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;

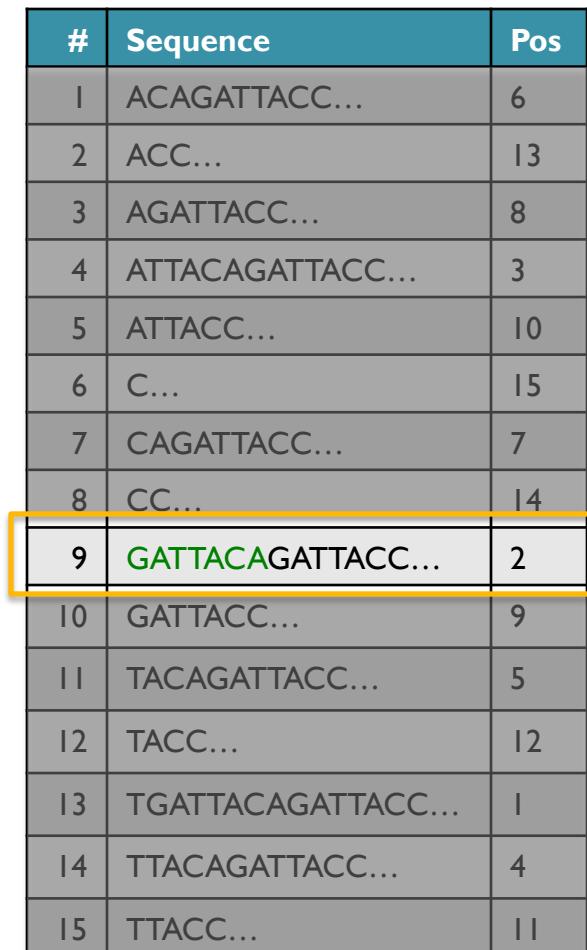
#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

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; 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	GATTACA <b>GATTACC...</b>	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

# 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

[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]

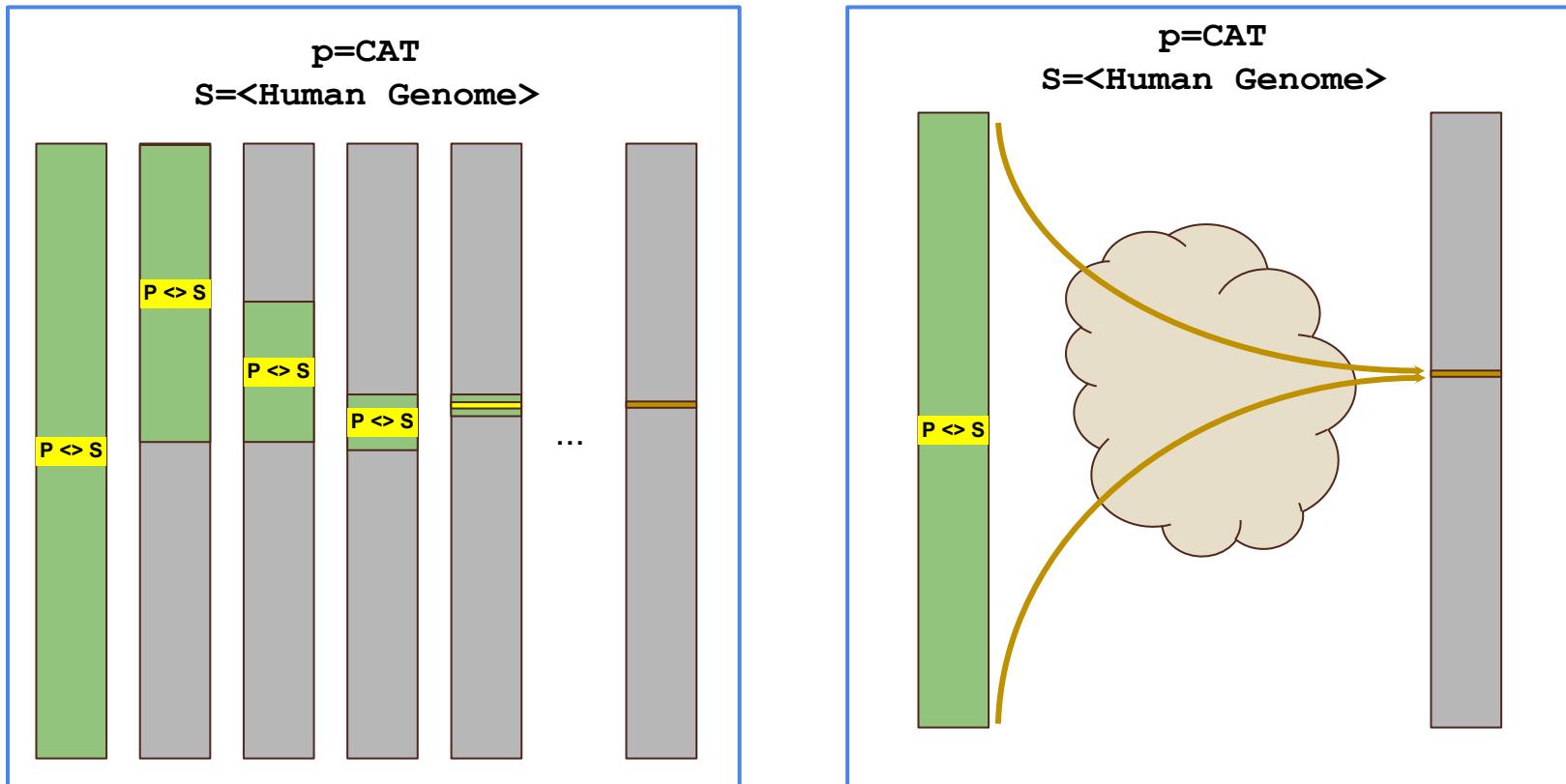
- Total Runtime:  $O(m \lg n)$

- More complicated, but **much** faster!

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

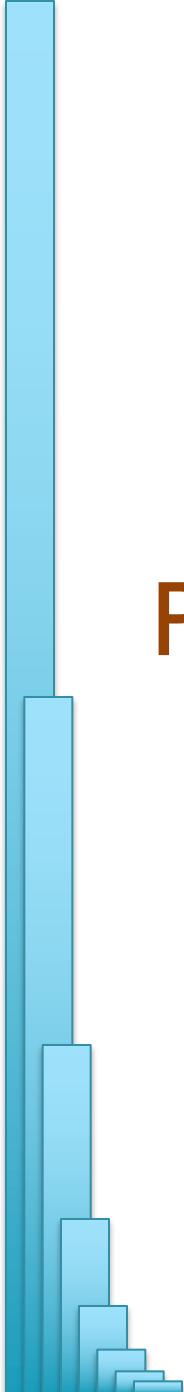
Can be reduced to  $O(m + \lg n)$   
using an auxiliary data structure called the LCP array

# 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?*

Kirsche, M, Das, A, Schatz, MC (2020) Bioinformatics  
doi: <https://doi.org/10.1093/bioinformatics/btaa911>

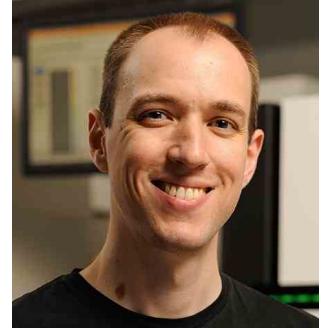


## Part 2: Burrows Wheeler Transform

# 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?

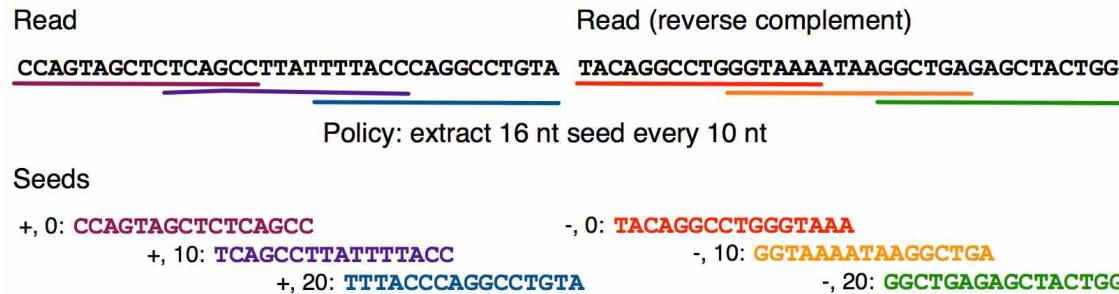


# Bowtie: Ultrafast and memory efficient alignment of short DNA sequences to the human genome

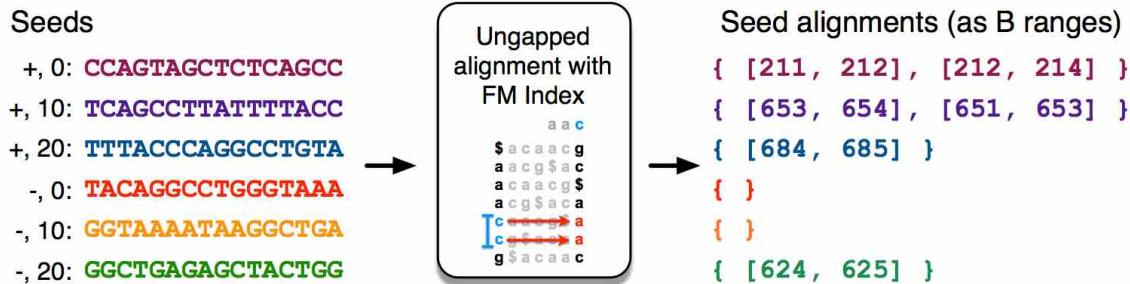
Slides Courtesy of Ben Langmead

# Algorithm Overview

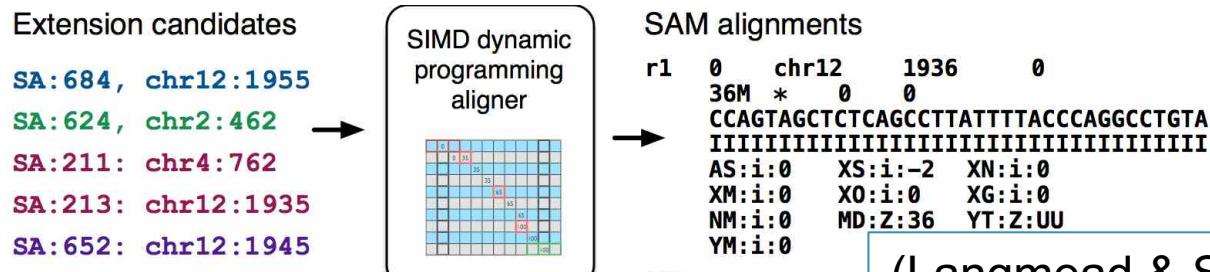
## 1. Split read into segments



## 2. Lookup each segment and prioritize



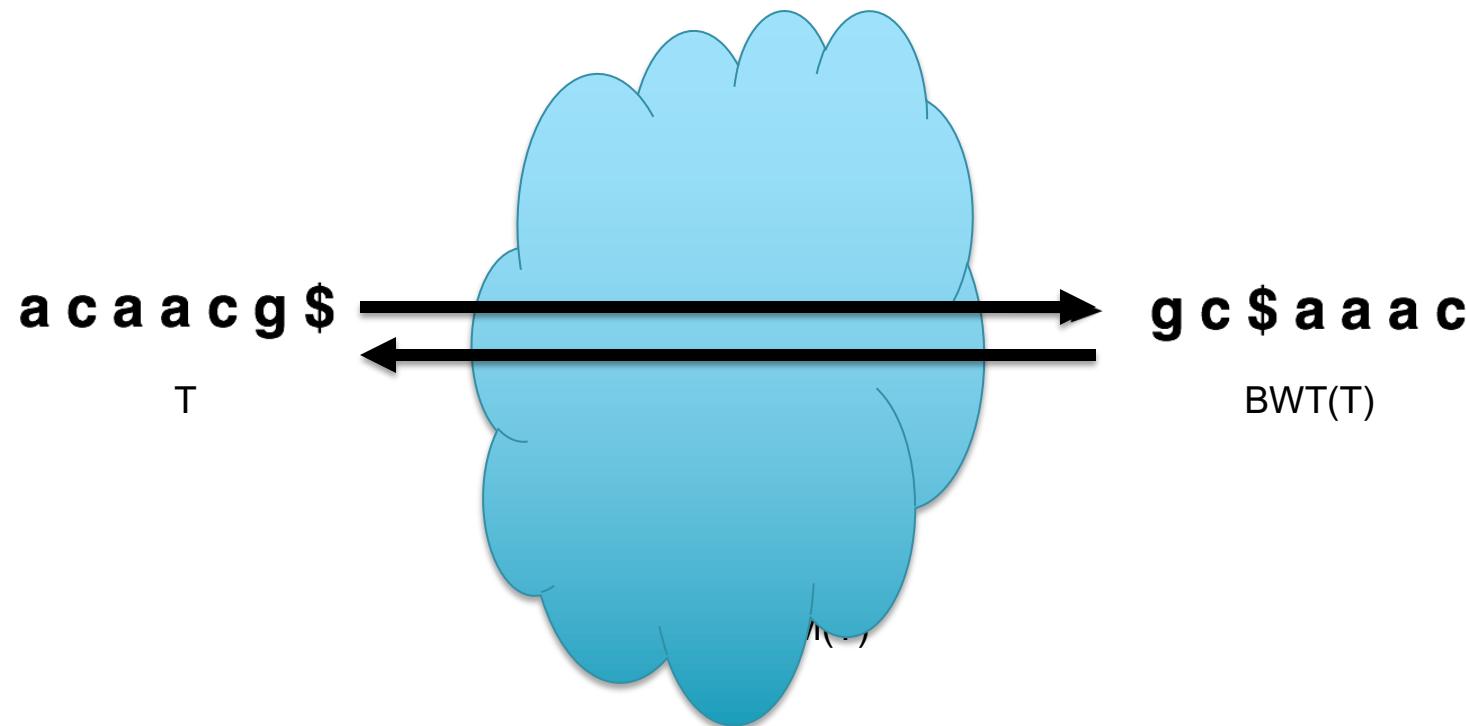
## 3. Evaluate end-to-end match



(Langmead & Salzberg, 2012)

# 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