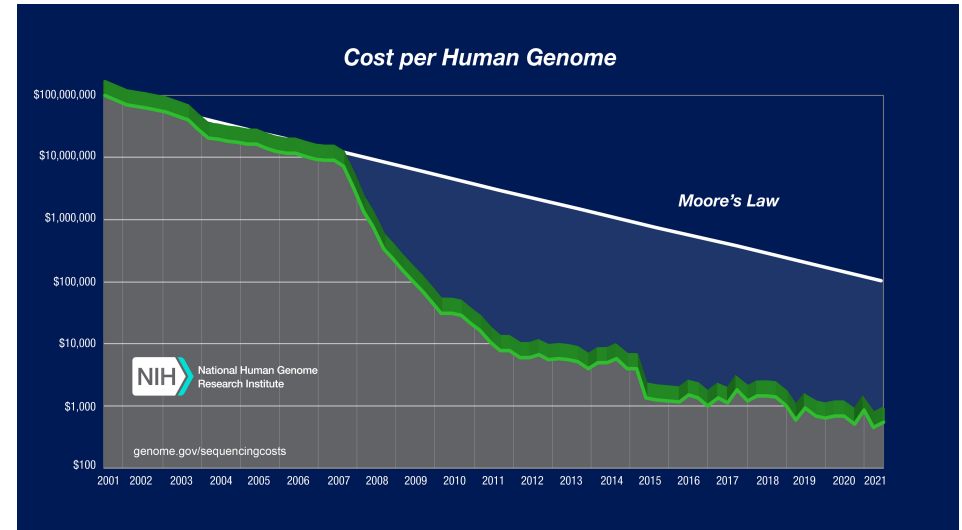




Mapping Sequencing Reads to a Reference Genome



High Throughput Sequencing

Example applications:

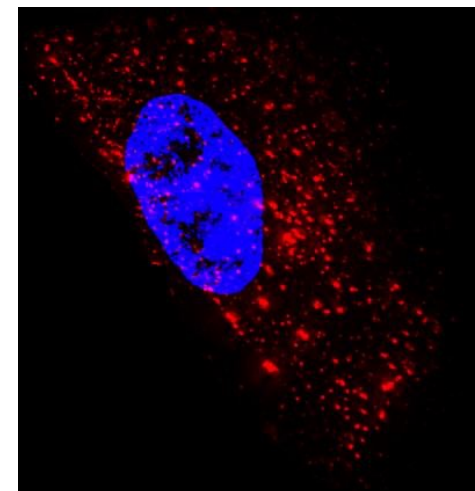
- Sequencing a genome (DNA)
- Sequencing a transcriptome and gene expression studies (RNA)
- ChIP (chromatin immunoprecipitation)

Example platforms:

- 454
- [Illumina](#)
- SOLiD



RNA





Sequencing Output

- Hundreds of millions of sequencing reads, each ~200 nts in length
- We need to map each read to the genome, i.e., determine the region of the genome each read corresponds to

```

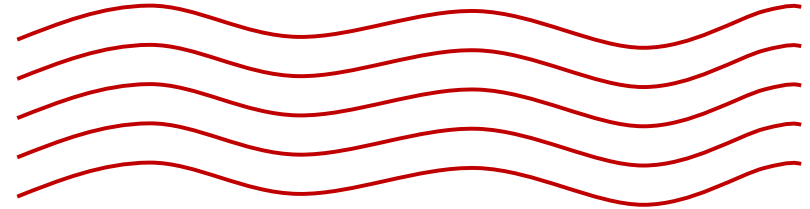
- ACGTAGTCGTAGTCGCTTACGATTTCGTATGCGTACGTGTAGTCTACGTGTA
- CCGCGCGTACTCTCGATGTACGTGTACGTACGATGTACTGTAGTACGTGT
- TTTGATCGTAGTGTCACTGAGCAACACCATTACTTACTATCTTGGACATC
- TGGGGCGATCGAGGATTCTAGTTATCGAGTGTCCGGGATTATCGGATCGAA
- GGCACTATACGTAGCGTATCGATTAGCACTGCGCGGCTATACGTCTGCGAT
- AGCGGGGCTGGACGACGGCTAGGCTCATCGTCGATCGATCGATCGTAAA
- TCGGTCGATCGAGTGTCTCGCGGCTCTCGAGAGGCTAGTAGAGAGCTG
- CTCCTCCGAGCGTAGTCCGATACGTATCGGATCTGACGATCGAGTCTGAT
- CCGTATCGTAGTCCGCGATCGATCGATTGAGTCTGATGATTACGTAGT
- TACGGCGACGTTACGCGTAGTCTGATGATGATGAGCGTCTAGTTCGATGCG
- GAAGCGTACATGCTAGTCTACGTCTCGCAGTTCGTATGGGATCGTATGACG
- ACGCAGCATCGATCTATGCATCGATGCTAGTCTGATGCTAGTCTATGCGAT
- ACGTAGTCGTAGTCGCTTACGATTTCGTATGCGTACGCTAGTCTACGTGTA
- CCGCGCGTACTCTCGATGTACGTGTACGTACGATGTACTGTAGTACGTGT
- TTTGATCGTAGTGTCACTGAGCAACACCATTACTTACTATCTTGGACATC
- TGGGGCGATCGAGGATTCTAGTTATCGAGTGTCCGGGATTATCGGATCGAA
- GGCACTATACGTAGCGTATCGATTAGCACTGCGCGGCTATACGTCTGCGAT
- AGCGGGGCTGGACGACGGCTAGGCTCATCGATCGGATCGGATCGGATAAA
- TCGGTCGATCGAGTGTCTCGCGGCTCTCGAGAGGCTAGTAGAGAGCTG
- CTCCTCCGAGCGTAGTCCGATACGTATCGGATCTGACGATCGAGTCTGAT
- CCGTATCGTAGTCCGCGATCGATCGATTGAGTCTGATGATTACGTAGT
- TACGGCGACGTTACGCGTAGTCTGATGATGAGCGTCTAGTTCGATGCG
- GAAGCGTACATGCTAGTCTACGTCTCGCAGTTCGTATGGGATCGTATGACG
- ACGCAGCATCGATCTATGCATCGATGCTAGTCTGATGCTAGTCTATGCGAT
- GTGCGTAGTCGTATATGCCTAGCATCGTTAGTCTAGCGTAGTCTACGTGTA
- ATCATCGGCGATAGTCTACGTAGTTATATCTACGCGCGCCCATCTCGCAA
- CGAGTAGGAGTCGTAGTCGTAGTCGATGCATCGAGTCTCGGATAGTCGTAG
- GACAGTCCGCGATGTATAGCAGCAGTTACGTAGCTAGTGTAGCTAGTA
- CAACTTTCGGCATCTTCGGTCTCTCTCTCTCTAGATAGAGACTTACGATCG
- TCCGGCATCGTAGTCCGGCGGACTATGGCGACACACGCTAGCATACC
- AGACTGAGTATATCGGCGGATGCGACTGTAGCTATATACGGCATCGCTC
- GGATCGAGTCACTCGGATCGAGTCTGAGCAGCAGTCTGATAGTATAGTCC
- GCGATCGTAGTCAGTTCGAGTCTGCGATCTGATGCGAGTACGTAGTCGTA
- TTATTCGCGCAGTGTGCTAGTCTGATGCTAGTCTGATGCTAGTCTGATGCT
- AGTACGTAGTATCTGAGCGTCTCTCTACGGACATCGATGCTACAGCTGA
- TTTATTACGACGATAGTGGCCATTGCGTATGAGTGTGCTAGTCTAGTCTAG
- GTGCGTAGTCGTATATGCCTAGCATCGTTAGTCTAGCTAGTCTAGTCTAGT
- ATCATCGGCGATAGTCTACGTAGTTATATCTACGCGCGCCCATCTCGCAA
- CGAGTAGGAGTCGTAGTCGTAGTCGATGCATCGAGTCTCGGATAGTCGTAG
- GACAGTCCGCGATGTATAGCAGCAGTTACGTAGCTAGTGTAGCTAGTA
- CAACTTTCGGCATCTTCGGTCTCTCTCTCTAGATAGAGACTTACGATCG
- TCCGGCATCGTAGTCCGGCGGACTATGGCGACACACGCTAGCATACC
- AGACTGAGTATATCGGCGGATGCGACTGTAGCTATATACGGCATCGCTC
- GGATCGACTCACTCGATCGAGTCTGAGCAGCAGTCTGATAGTATAGTCC
- GCGATCGTAGTCACTCGAGTCTGCGATCTGATGCGAGTACGTAGTCGTA
- TTATTCGCGCAGTGTGCTAGTCTGATGCTAGTCTGATGCTAGTCTGATGCT
- AGTACGTAGTCTGAGCGTCTCTCTACGGACATCGATGCTACAGCTGA
- TTTATTACGACGATAGTGGCCATTGCGTATGAGTGTGCTAGTCTAGTCTAG

```

G-5



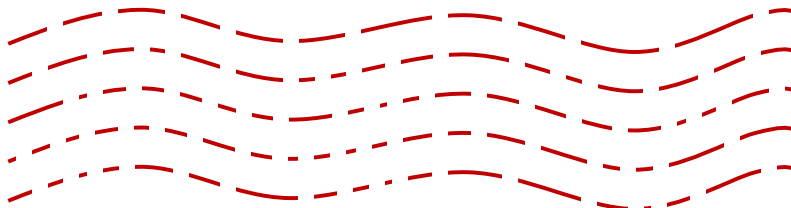
Nucleic Acid Sequencing



G-6



Nucleic Acid Sequencing

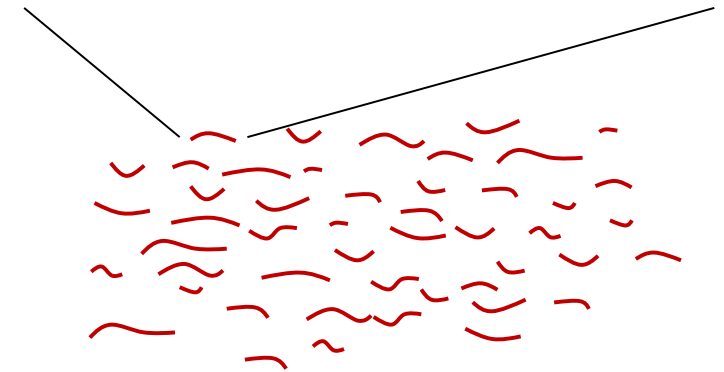


G-7



Nucleic Acid Sequencing

CGTAGTAGTCACAGTCTACGTATATGGGCTCAGCATATAGCGTATAGCGGACTTAGCCATCGTA



G-8



Nucleic Acid Sequencing

```

>CGTAGTAGTCACAGTCTACGTATATGGGCTCAGCATATAGCGTATAGCGGACTTAGCCATCG
>GCGTATAGTCTATATACGACTATCGGCTCGGTCGAGAGCAGATATATGCAGTTATATGCTAG
>CCTACGTTATATCGATACTACTAGTCTCGTCATGAGCGAGTAGATAGTATGACGAGCGCAGTAC
...
>CGATATTAGCCTAGCATCATTACGGCGAGACTCTCGGCTCGCTATATAGCGCTATAGCGAT
>CGGCTATAGCGCATATGCTCAGTAGCTATTAGCAGTATTACGATTATAGTCTCGGCGCATTAC
>TTTCGGGGATAAGTCTTCGCTTATGCGCAGATTATACGGCCGTATATTTGCATTTAGCATTT
>GGCGTATGGCGGATATCGGCGGTCATAGCAGCCGATTAGGCTACGCCGATGCATCG
>CGCGATCGCGCGGATCGCGTCAGTCGCGCAGTAGCGCGGCATAGTCGTATCGGCGCCG
>TGACAGAAGCTATAAGAGTCAGTAGATCTGAGTATTAGCATTATCGGCGCGATGCGCGATAACG
>GCGTATAGTCTATATACGACTTATCGGCTCGGTCGAGCAGATATATGCAGTTATATGCTA
>CGCGATCGCGCGGATCGCGTCAGTCGCGCAGTAGCGCGGCATAGTCGTATCGGCGCCGATCGC
>ATAGCAGCAGTGTAGGATATGCTGCTCGTTTCGACTATCATATCTCGCTGGTCTAGCA
>TGACAGAAGCTATAAGAGTCAGTAGATCTGAGTATTAGCATTATCGGCGCGATGCGCGA
>CCTACGTTATATCGATACTACTAGTCTCGTCATGAGCGAGTAGATAGTATGACGAGCGAGCATCG
>CGTAGTAGTCACAGTCTACGTATATGGGCTCAGCATATAGCGTATAGCGGACTTAGCCATCG
>TTTCGGGGATAAGTCTTCGCTTATGCGCAGATTATACGGCCGTATATTTGCATTTAGCATTT
>GGCGTATGGCGGATATCGGCGGTCATAGCAGCCGATTAGGCTACGCCGATGCATCGTCGAGTA

```

G-9



Mapping to Reference Genome

Reference Genome

```

CGTAGTAGTCACAGTCTACGTATATGGGCTCAGCATATAGCGTATAGCGGACTTAGCCATCGCGTG
TAGTCTAGTCAGTATAGCGATCAGTACTATGCAGTCTACGTAGTCTGTTATGACGTAGTCGATGTA
GCTAGTATCGTAGTACGGCATAAGTCGCGCATGGCTGCGTCTCATCATATCGGCACGACGCTCT
AGAGTAGTCACACTTGTGTGTATATAGCGCGGGAGGAGAGCTCTCTGAATAGCTGAGATGCGTA
AGTGCAGGAGAGATCTCTCTGAGAGAGTCTCGGGGATCTCTCTAGAAAGCTCTCGGAAGGATCTCG
AGAGGACTCTC6CGTAGAGAGCTTACAGAGACGATATATATGCGATTAGTACGTATGTC

```

Sequencing Read

```
TCATCATATCGGCACGACGCTCTAGAGTAGTCACACTTGTGTGTATATAGCGCGGGAGGAGA
```

G-10



Burrows-Wheeler Transform (BWT)

ATCATTAAATCATG\$

GTAASCCTTTTAAAA

ATCATTAAATCATG\$	0	\$ATCATTAAATCATG	14
TCATTAAATCATG\$A	1	AAATCATG\$ATCATT	6
CATTAAATCATG\$AT	2	AATCATG\$ATCATTA	7
ATTAAATCATG\$ATC	3	ATCATG\$ATCATTTAA	8
TTAAATCATG\$ATCA	4	ATCATTAAATCATG\$	0
TAAATCATG\$ATCAT	5	ATG\$ATCATTAAATC	11
AAATCATG\$ATCATT	6	ATTAAATCATG\$ATC	3
AATCATG\$ATCATTAA	7	CATG\$ATCATTAAAT	10
ATCATG\$ATCATTAA	8	CATTAAATCATG\$AT	2
TCATG\$ATCATTAAA	9	G\$ATCATTAAATCAT	13
CATG\$ATCATTAAAT	10	TAAATCATG\$ATCAT	5
ATG\$ATCATTAAATC	11	TCATG\$ATCATTAAA	9
TG\$ATCATTAAATCA	12	TCATTAAATCATG\$A	1
G\$ATCATTAAATCAT	13	TG\$ATCATTAAATCA	12
\$ATCATTAAATCATG	14	TTAAATCATG\$ATCA	4

G-11



Efficient Substring Search

ATCATTAAATCATG\$

GTAASCCTTTTAAAA

TCA

\$ATCATTAAATCATG	14
AAATCATG\$ATCATT	6
AATCATG\$ATCATTA	7
ATCATG\$ATCATTTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
CATG\$ATCATTAAAT	10
CATTAAATCATG\$AT	2
G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTAAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

G-12



Efficient Substring Search

ATCATTAAATCATG\$

GTAASCCTTTTAAAA

TCA

\$ATCATTAAATCATG	14
AAATCATG\$ATCAIT	6
AATCATG\$ATCATTA	7
ATCATG\$ATCATTTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
CATG\$ATCATTTAAAT	10
CATTAAATCATG\$AT	2
G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTTAAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

G-13



Efficient Substring Search

ATCATTAAATCATG\$

GTAASCCTTTTAAAA

TCA

Range of rows starting with
(1,7)

Thus, the substring
is in the reference sequence at indices
{6,7,8,0,11,3}

\$ATCATTAAATCATG	14
AAATCATG\$ATCAIT	6
AATCATG\$ATCATTA	7
ATCATG\$ATCATTTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
...CATG\$ATCATTTAAAT	10
CATTAAATCATG\$AT	2
G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTTAAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

G-14



Efficient Substring Search

ATCATTAAATCATG\$

GTAASCCTTTTAAAA

TCA

Range of rows starting with
(1,7)

A

\$ATCATTAAATCATG	14
AAATCATG\$ATCAIT	6
AATCATG\$ATCATTA	7
ATCATG\$ATCATTTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
...CATG\$ATCATTTAAAT	10
CATTAAATCATG\$AT	2
G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTTAAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

G-15



Efficient Substring Search

ATCATTAAATCATG\$

GTAASCCTTTTAAAA

TCA

Range of rows starting with
(1,7)

A

\$ATCATTAAATCATG	14
AAATCATG\$ATCAIT	6
AATCATG\$ATCATTA	7
ATCATG\$ATCATTTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
...CATG\$ATCATTTAAAT	10
CATTAAATCATG\$AT	2
G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTTAAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

G-16



Efficient Substring Search

ATCATTAAATCATG\$

GTAA\$CCTTTTAAAA

TCA

Range of rows starting with
(1,7)

A

\$ATCATTAAATCATG	14
AAATCATG\$ATCATT	6
AATCATG\$ATCATT	7
ATCATG\$ATCATTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
CATG\$ATCATTAAAT	10
CATTAAATCATG\$AT	2
...G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTAAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

Range of rows starting with
(7,9)

CA

Thus, the substring

CA

is in the reference sequence at indices
{10,2}



Efficient Substring Search

ATCATTAAATCATG\$

GTAA\$CCTTTTAAAA

TCA

Range of rows starting with
(1,7)

A

\$ATCATTAAATCATG	14
AAATCATG\$ATCATT	6
AATCATG\$ATCATT	7
ATCATG\$ATCATTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
CATG\$ATCATTAAAT	10
CATTAAATCATG\$AT	2
...G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTAAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

Range of rows starting with
(7,9)

CA



Efficient Substring Search

ATCATTAAATCATG\$

GTAA\$CCTTTTAAAA

TCA

Range of rows starting with
(1,7)

A

\$ATCATTAAATCATG	14
AAATCATG\$ATCATT	6
AATCATG\$ATCATT	7
ATCATG\$ATCATTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
CATG\$ATCATTAAAT	10
CATTAAATCATG\$AT	2
...G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTAAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

Range of rows starting with
(7,9)

CA



Efficient Substring Search

ATCATTAAATCATG\$

GTAA\$CCTTTTAAAA

TCA

Range of rows starting with
(1,7)

A

\$ATCATTAAATCATG	14
AAATCATG\$ATCATT	6
AATCATG\$ATCATT	7
ATCATG\$ATCATTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
CATG\$ATCATTAAAT	10
CATTAAATCATG\$AT	2
...G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTAAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

Range of rows starting with
(7,9)

CA

Range of rows starting with
(11,13)

TCA

Thus, the substring

TCA

is in the reference sequence at indices
(9,1)



Efficient Substring Search

ATCATTAAATCATG\$

GTAA\$CCTTTTAAAA

TCA

Range of rows starting with
A
[1,7]

A

\$ATCATTAAATCATG	14
AAATCATG\$ATCATT	6
AATCATG\$ATCATTA	7
ATCATG\$ATCATTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
CATG\$ATCATTAAAT	10
CATTAAATCATG\$AT	2
G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

G-21



Efficient Substring Search

ATCATTAAATCATG\$

GTAA\$CCTTTTAAAA

TCA

Range of rows starting with
A
[1,7]

A

Range of rows starting with
CA
[7,9]

CA

Range of rows starting with
TCA
[11,13]

TCA

Each step of the search *must be fast!*
O(1) time.

\$ATCATTAAATCATG	14
AAATCATG\$ATCATT	6
AATCATG\$ATCATTA	7
ATCATG\$ATCATTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
CATG\$ATCATTAAAT	10
CATTAAATCATG\$AT	2
G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

G-22



Precompute Helper Information

GTAA\$CCTTTTAAAA

Returns the number of nucleotide characters in the BWT that are lexicographically less than c.

```
public int getNumberCharactersLessThan(char c);
    getNumberCharactersLessThan('$'); returns 0
    getNumberCharactersLessThan('C'); returns 7
    getNumberCharactersLessThan('T'); returns 10
```

Returns the number of occurrences of nucleotide character c in the BWT up to but not including index i.

```
public int getNumberOccurrencesPriorToIndex(char c, int i);
    getNumberOccurrencesPriorToIndex('A', 5); returns 2
    getNumberOccurrencesPriorToIndex('A', 13); returns 4
    getNumberOccurrencesPriorToIndex('T', 8); returns 2
```

G-23



Efficient Substring Search

ATCATTAAATCATG\$

GTAA\$CCTTTTAAAA

TCA

How can we compute this?

Range of rows starting with
A

A

Range of rows starting with
CA

CA

Range of rows starting with
TCA

TCA

\$ATCATTAAATCATG	14
AAATCATG\$ATCATT	6
AATCATG\$ATCATTA	7
ATCATG\$ATCATTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
CATG\$ATCATTAAAT	10
CATTAAATCATG\$AT	2
G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

G-24



Efficient Substring Search

ATCATTAAATCATG\$

TCA

Range of rows starting with **A**

A

Range of rows starting with **CA**

CA

Range of rows starting with **TCA**

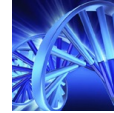
TCA

How can we compute these?

GTAA\$CCTTTTAAAA

\$ATCATTAAATCATG	14
AAATCATG\$ATCATT	6
AATCATG\$ATCATT	7
ATCATG\$ATCATTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
CATG\$ATCATTAAAT	10
CATTAAATCATG\$AT	2
G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTAAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

G-25



Another Example

ATCATTAAATCATG\$

ATT

Range of rows starting with **T**

T

Range of rows starting with **TT**

TT

Range of rows starting with **ATT**

ATT

GTAA\$CCTTTTAAAA

\$ATCATTAAATCATG	14
AAATCATG\$ATCATT	6
AATCATG\$ATCATT	7
ATCATG\$ATCATTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
CATG\$ATCATTAAAT	10
CATTAAATCATG\$AT	2
G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTAAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

G-26



A Substring Not in the Reference

ATCATTAAATCATG\$

GTC

Range of rows starting with **C**

C

Range of rows starting with **TC**

TC

Range of rows starting with **GTC**

GTC

GTAA\$CCTTTTAAAA

\$ATCATTAAATCATG	14
AAATCATG\$ATCATT	6
AATCATG\$ATCATT	7
ATCATG\$ATCATTAA	8
ATCATTAAATCATG\$	0
ATG\$ATCATTAAATC	11
ATTAAATCATG\$ATC	3
CATG\$ATCATTAAAT	10
CATTAAATCATG\$AT	2
G\$ATCATTAAATCAT	13
TAAATCATG\$ATCAT	5
TCATG\$ATCATTAAA	9
TCATTAAATCATG\$A	1
TG\$ATCATTAAATCA	12
TTAAATCATG\$ATCA	4

G-27



Algorithm: Mapping read to genome

← Read is processed one character (NT) at a time, from right to left.

TCA

Base Case: Last character (NT) in read, **ch**

The range of rows where **ch** appears in the BWT is **[start, end)** where

start = `getNumberCharactersLessThan(ch)`

if (**ch** is not 'T')

end = `getNumberCharactersLessThan(NT after ch alphabetically)`

if (**ch** is 'T')

end = length of BWT

G-28



Algorithm: Mapping read to genome

← Read is processed one character (NT) at a time, from right to left.
 TCA

Iterative Case: Character (NT) at index i in read, ch

The range of rows where the substring from i onward in the read appears in the BWT is [$start$, end) where

$start = \text{getNumberCharactersLessThan}(ch) + \text{getNumberOccurrencesPriorToIndex}(ch, start \text{ for character at index } i+1)$

$end = \text{getNumberCharactersLessThan}(ch) + \text{getNumberOccurrencesPriorToIndex}(ch, end \text{ for character at index } i+1)$

G-29



Efficient Substring Search with Errors

GCTAGAGTATATCGGCGCTATGCGTTGTCTCTCTAGAGAGACGATTCTAGTCTGCCTGTC

Suppose we tried to map a read of length 60 nts to a reference genome and found that the read did not map. Perhaps the read contains one or two errors from the sequencing process.

- Break the read up into three pieces and map each of the three pieces to the reference.
- If none of the three pieces map, then there are at least three errors in the read.
- If one or more of the pieces map, then we use that mapping to perform a *fast* alignment.

G-30