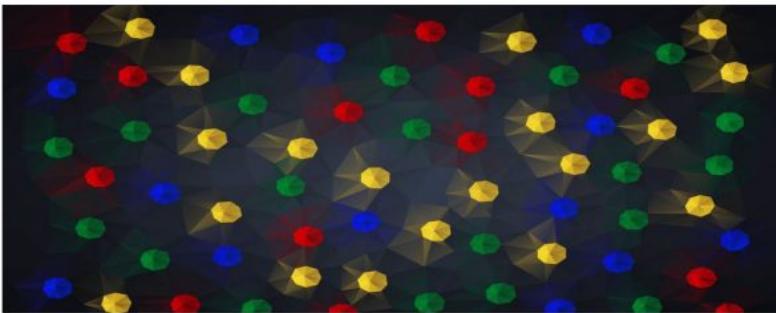
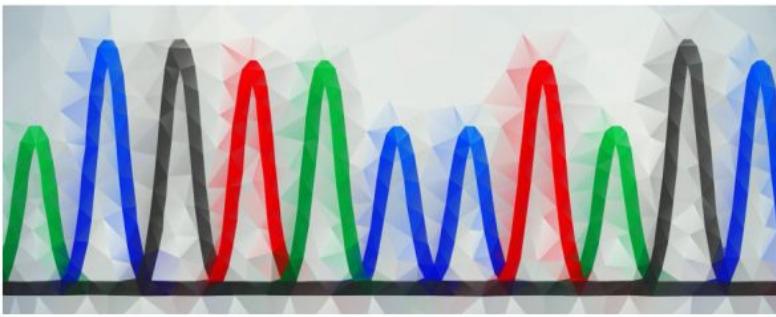
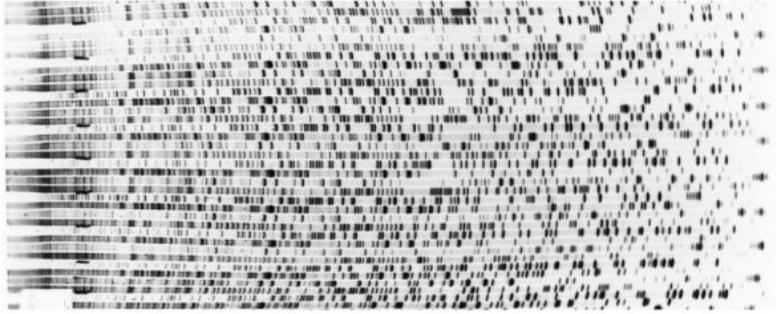


# GENOME ANALYSIS & BIOINFORMATICS

## Week 2:

Bioinformatics and searching  
for patterns in sequence data



# The future of SEQUENCING

ever more  
MASSIVE

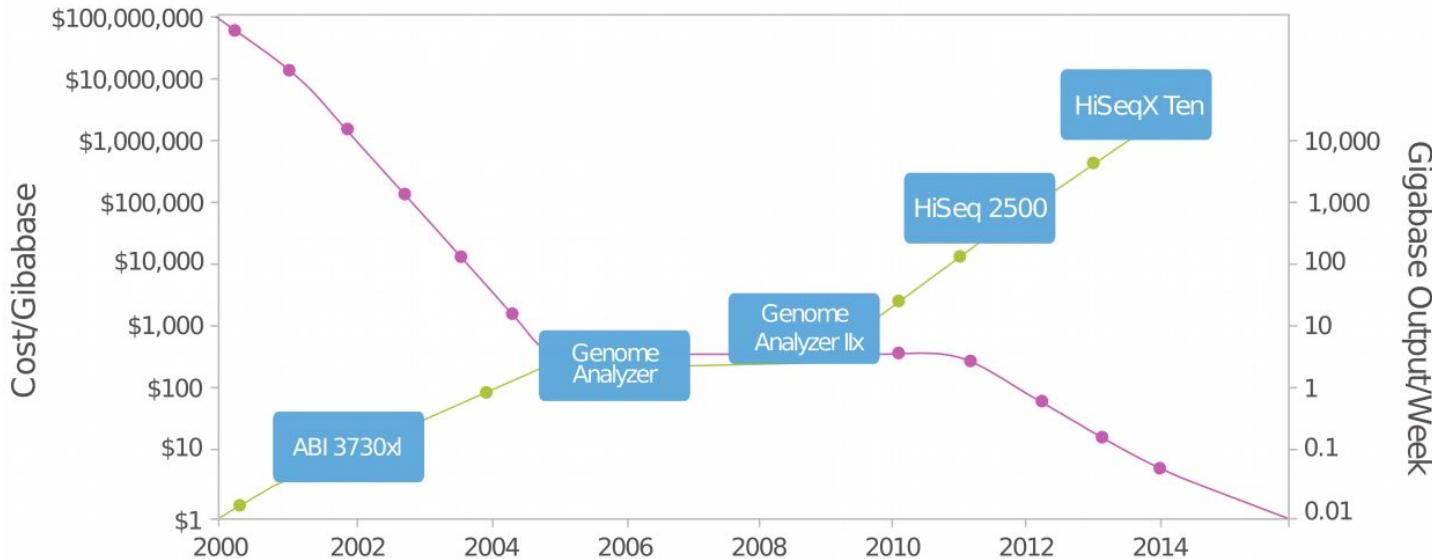
ever more  
PARALLEL

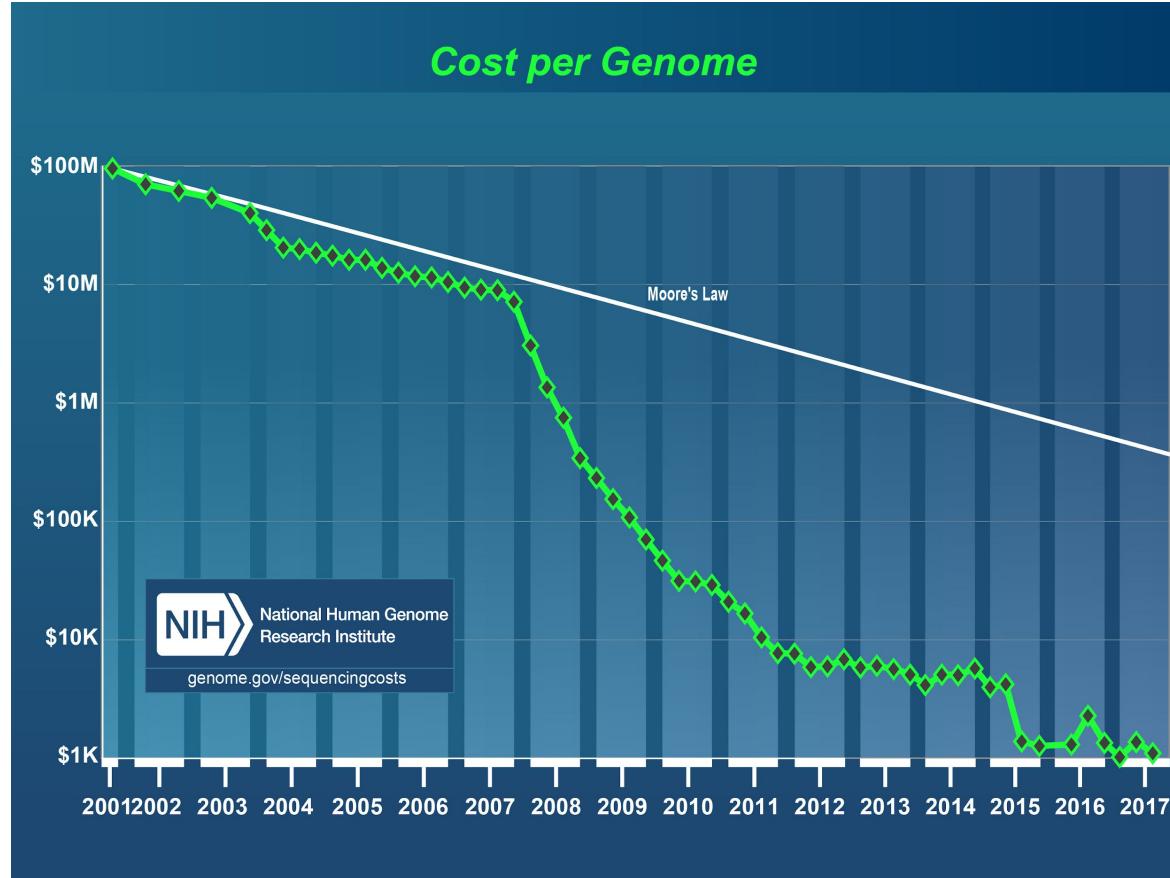
ever more  
DATA

Next generation sequencing  
has outpaced

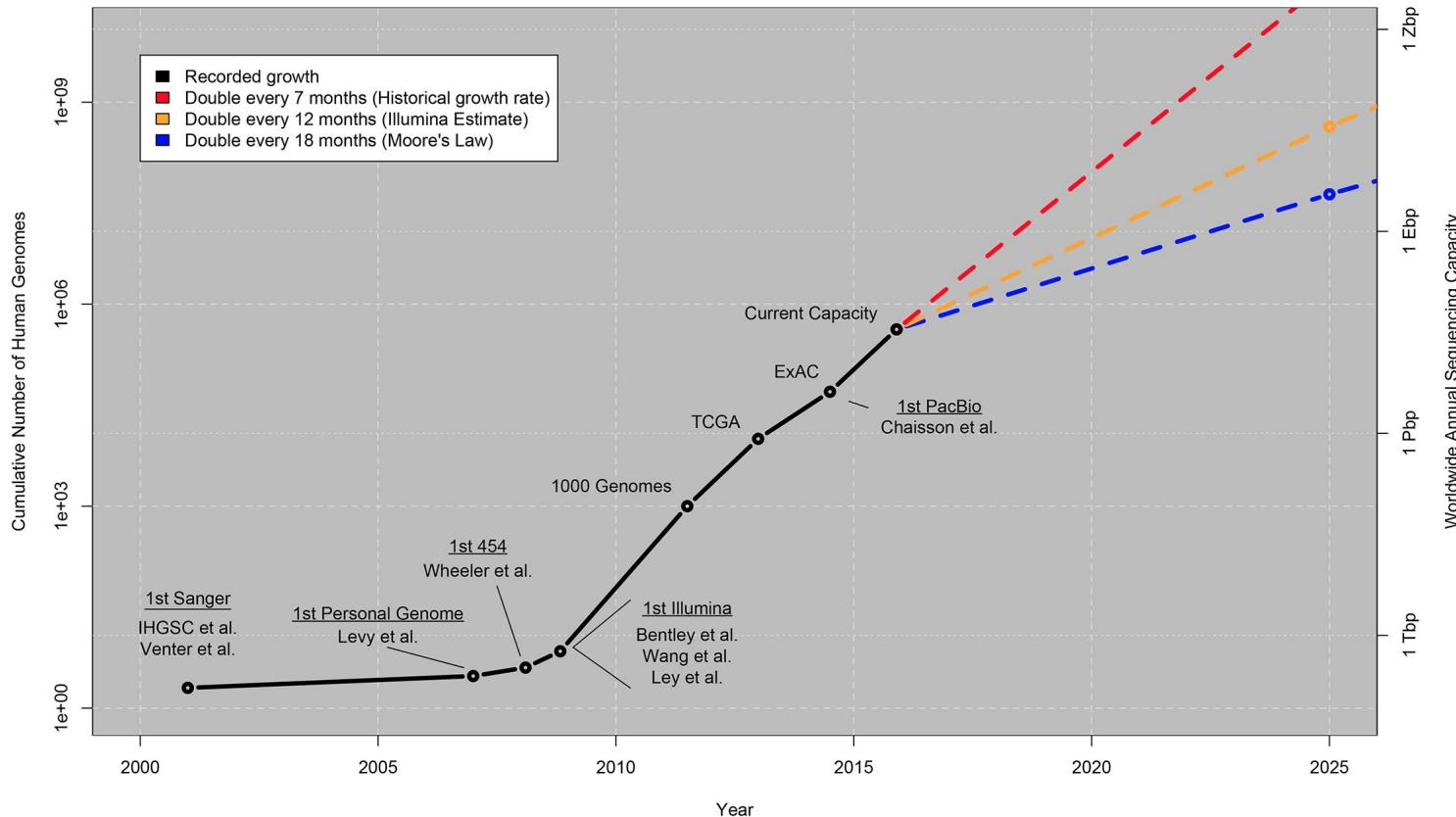
## MOORE'S LAW:

*“overall processing power of computers  
will double every two years”*

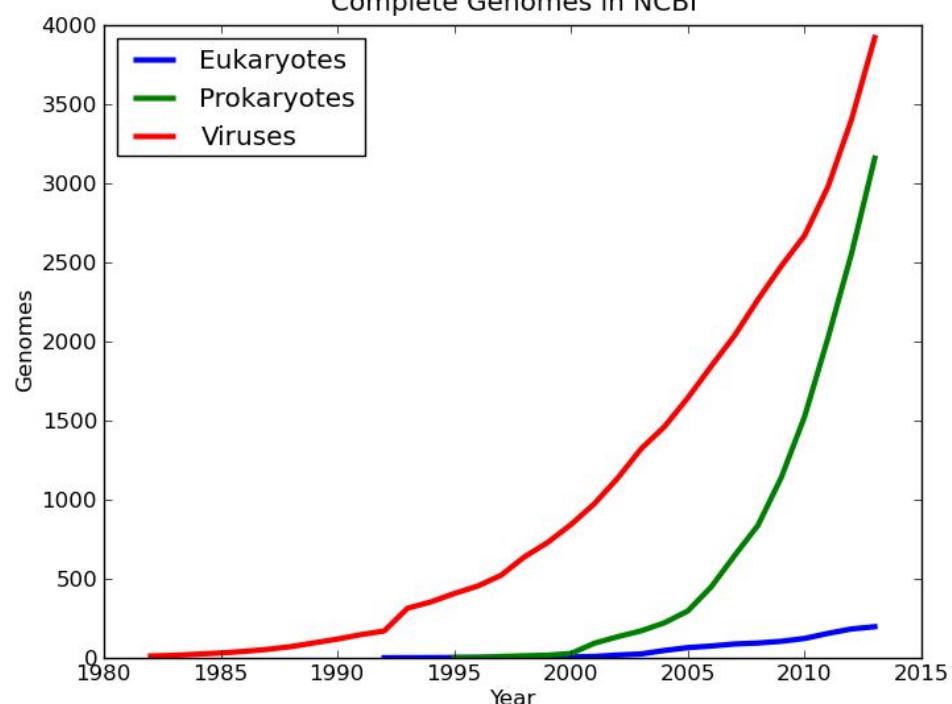




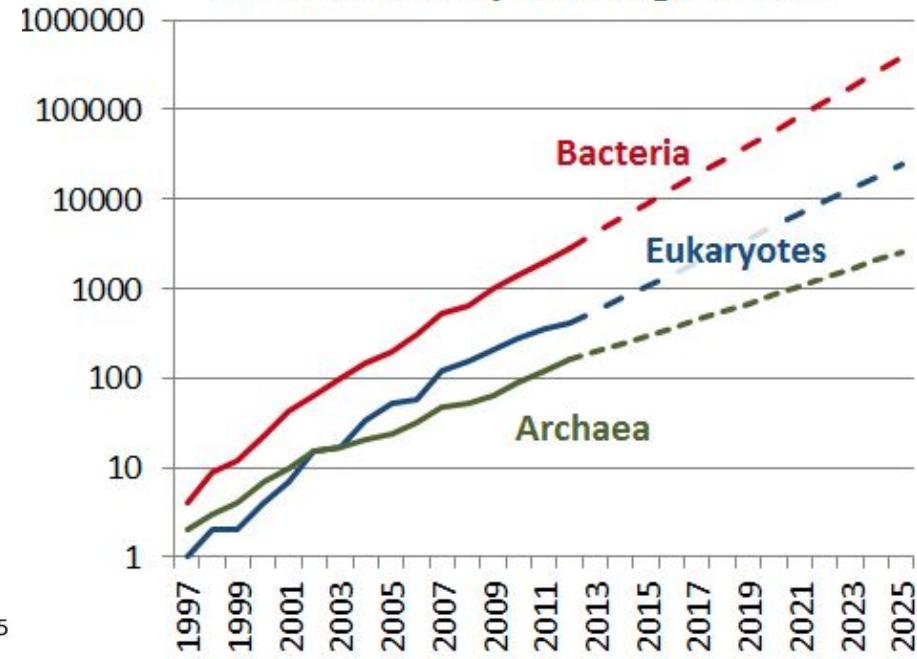
## Growth of DNA Sequencing



Complete Genomes in NCBI



Cumulative sequenced genomes



<b>Data Phase</b>	<b>Astronomy</b>	<b>Twitter</b>	<b>YouTube</b>	<b>Genomics</b>
<b>Acquisition</b>	25 zetta-bytes/year	0.5–15 billion tweets/year	500–900 million hours/year	1 zetta-bases/year
<b>Storage</b>	1 EB/year	1–17 PB/year	1–2 EB/year	2–40 EB/year
<b>Analysis</b>	In situ data reduction	Topic and sentiment mining	Limited requirements	Heterogeneous data and analysis
	Real-time processing	Metadata analysis		Variant calling, ~2 trillion central processing unit (CPU) hours
	Massive volumes			All-pairs genome alignments, ~10,000 trillion CPU hours
<b>Distribution</b>	Dedicated lines from antennae to server (600 TB/s)	Small units of distribution	Major component of modern user's bandwidth (10 MB/s)	Many small (10 MB/s) and fewer massive (10 TB/s) data movement

doi:10.1371/journal.pbio.1002195.t001



2017  
1 Human Genome  
per Hour

200 TB of raw data  
per day

# BIG DATA



VOLUME

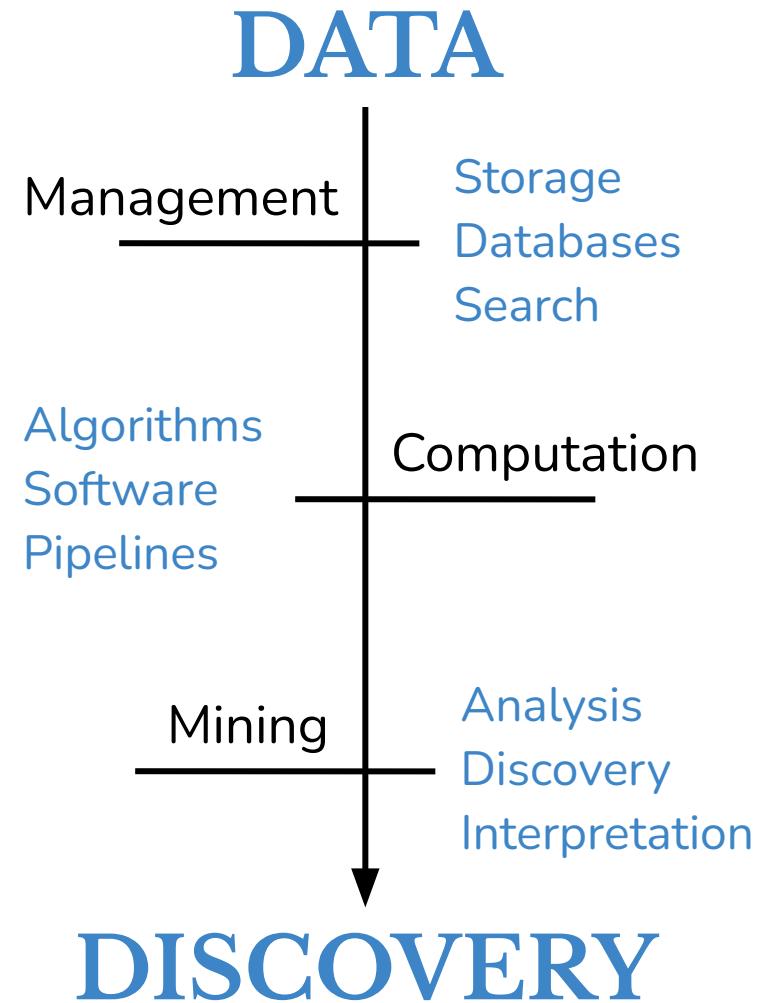
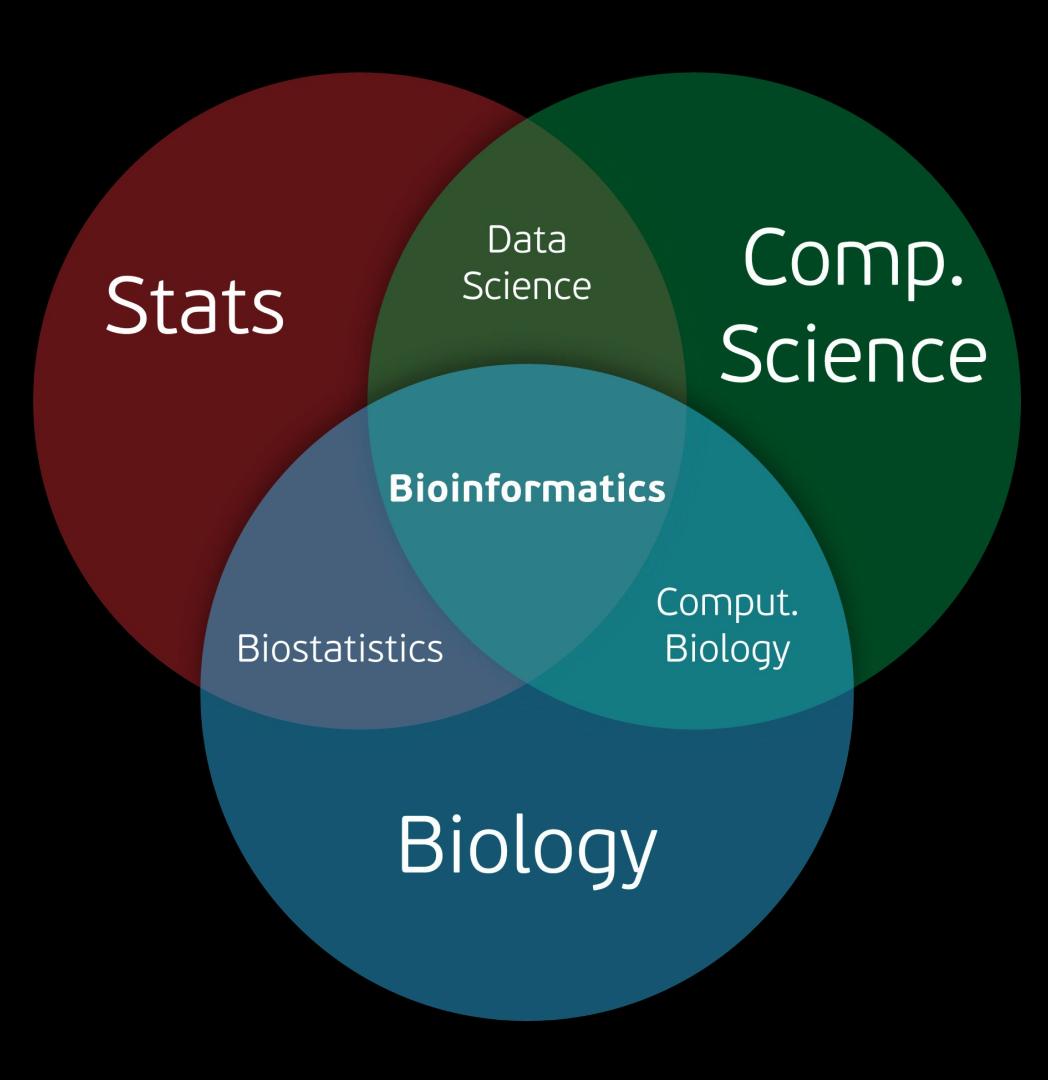
VELOCITY

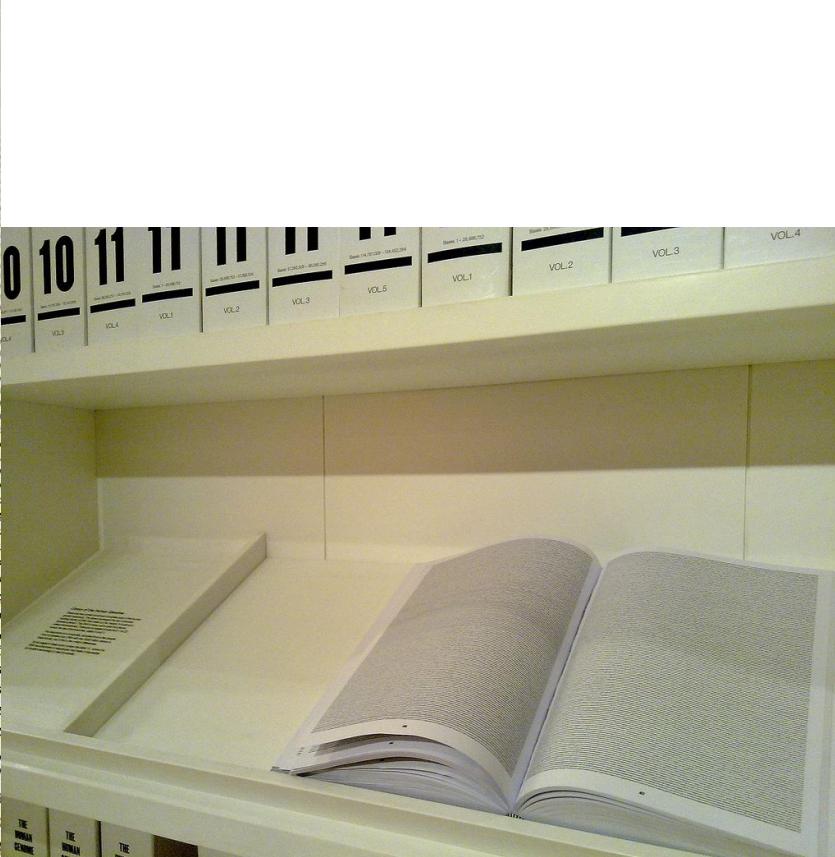
VARIETY

CAPTURING  
STORAGE  
ANALYSIS  
SHARING  
TRANSFER  
VISUALIZATION

# BIOINFORMATICS

An interdisciplinary field that develops and applies **computer technologies** and computational methods to **genome data** for understanding biological processes





CCACACGACAGCTAAGATCCAAACTGGGATTAGATACCCCCTATGCTTAGCCTTAAACATAAATAATTT  
ATTAAACAAAATTATTCGCCAGAGAACTACTAGCAACAGCTAAAGGACTTGGCGGTGCTTTA  
ACCCCCCTAGAGGAGCCTGTTCTATAATCGATAAAACCCGATAGACCTCACCCCTTGCTAATTCACT  
CTATATACGCCATCTTCAGCAAACCCCTAAAAGGAACAAGAGTAAGCACAATCATCTTGCTAAAAAAAG  
TTAGGTCAAGGTGTAACCCATGGGGTGGAAAGAAATGGGCTACATTTTCTATTCAAGAACAACTTACGAA  
AGTTTTTATGAAACTAAAACCTAAAGGTGGATTAGTAGTAATCAAGAATAGAGAGCTGATTGAATAA  
GGCAATGAAGCATGCACACACCAGCCGTACCCCTCTCAAGTGGCACAAGGCAAACACAACCTATTGAAA  
CCAATAAAACGCAAGAGGGAGGCAAGTCGTAACAAGGTAAGCATACTGGAAAGTGTGCTGGACGGATCA  
AAAGTGTAGCTTAAACAAAGCATCTAGCTTACCCAGAGGATTTCACGCATGTGACCACCTTGAAACCCAG  
AGCTAGCCCAGATAACAACCTAACCAACTACCCTAGACCAATTAAATAAAACATTCACTAGTAGTATAATTAA  
AGTATAGGAGATAGAAATTCTTTAATCGGAGCTATAGAGAGAGTACCGCAAGGGAATGATGAAAGATTA  
CTTAAAGTAACAAACAGCAAAGATTACCCCTCTACCTTGCTATAATGAGTTAGCCAGAAATAACCTAA  
CAAAGAGAACTTAAGCTAGGTTCCCGAACCGAGCAGCTACCTATGAACAATCCACTGGGATGAACTC  
ATCTATGTTGAAATAGT GAGAAGATTATAGGTAGAGGTGAAAGGCCTAACGACCTGGTGTAGCTG  
GTTGCCAGAATAGAATTCTAGTTCAACTTAAACTTGCTACAAAATTAAATTAAATGCAAGTTT  
AAAATATATTCTAAAAGGTACAGCTTTAGAATCAAGGATACAGCCTTACTTAGAGAGTAATACTGA  
TTAAATCATAGTAGGCCTAGAAGCAGCCATTAAGAAAGCGTTAAGCTAACACCCATATCAACTT  
GATAACAAAAATCTAATTAACTCTTAATATAAACTGGGCTAATCTATTAAATATAGAAGCAACAA  
TGCTAATATGAGTAACGAGAAGTACTTCTCAACGGCTAACAGCAACGGATAACCAACTGATA  
GTTAACACAACCGTAGAAATAATCCAACAATAAAACACCTACCAAAACCAATTGTTAACACACAGGT  
TGCAGACTAAGGAAAGATTAAAAGAAGTAAAAGGAACCTGGCAACACAAACCCGCTGTTACCAAAAAA  
CATCACCTCAGCATTCCCAGTATTGGAGGACTGCCTGCCCGTGACATCAGTTAACGGCCGGTAT  
CCTGACCGTCAAAGGTAGCATAATCATTGTTCTAAATAAGGACTTGTATGAAAGGCCACAGGAGGG  
TTTAACTGTCCTTACTTCAATCAGTGAATTGACCTTCCCGTGAAGAGGCGGAAATAAAATAAAGA  
CGAGAAGACCCATGGAGCTTCAATTAAATTAGCTAAAGGATTATTACAGACGGACAGGAACACA  
TATTCTTCCATGAGCTAGCAATTAGTTGGGGCGACCTGGAGTACAAAATAACCTCGAGTGTATT  
AATCTAGACGTACCACTGAAATGCTCACTTACTTATTGATCCTTGTGACGGAAACAAGT  
TACCCCTAGGGATAACAGCGCAATCCTATTAAAGAGTCCATCGACAATAGGTTACGACCTCGATGTT  
GGATCAGGACATCTTAATGGTCGACAGCTATTAAAGGGTCTTGTGACGGTAAAGTCCCTACGTGA  
TCTGAGTTCAAGCCGGAGCAATCCAGGTGGCTTCTATCTATTAAATGACCTCTCCAGTACGAAAGGA  
CAAGAGAAGTAAGGCCTCCCTACCAAAAGCCCTTAAGACCAATAGATGACTTGTCAAACCTAGTAAG  
TCTACCCCAACGTTGCCAAGAGACAGGGCTTGTAGGGTGGCAGAGCCGGTATTGCTACAAACTT  
AAACCTTAACTCAGAGGTTCAAATCCTCCCTAACACTATGTTATAATTAAACATTCTCACTAGT  
CGTACCCATTCTCTGCCGTAGCCTTAAACACTAGTGGAGCAGGAAACTGGCTACACAAACTT  
CGTAAAGGCCAACATTGAGGACCTATGGTCTCCCTACACCTATCGCAGATGCTACAAAATT  
CCAAAGAACCTTGCCTACATCGCCACAACTATTTATTAGCCCTATTCTAGCCCTAAC

# BIOINFORMATICS

Searching for patterns in  
biological data

## **Finding Origin of Replication Problem:**

**Input:** A DNA string *Genome*.

**Output:** The location of *oriC* in *Genome*.

## The OriC region of *Vibrio cholerae*

atcaatgatcaacgtaagcttctaaggatcatgatcaagggtgctcacacagttatccacaac  
ctgagtggatgacatcaagataggtcgttgtatctccttcctctcgtaactctcatgacca  
cgaaaagatgatcaagagaggatgattcttggccatatcgcaatgaatacttgtgactt  
gtgcttccaattgacatcttcagcgccatattgcgctggccaagggtgacggagcgggatt  
acgaaagcatgatcatggctgttctgttatcttggtttgactgagacttgttagga  
tagacggttttcatcactgacttagccaaagcctactctgcctgacatcgaccgtaaat  
tgataatgaatttacatgcttccgcgacgattacacttgcattgatcatcgatccgattgaag  
atcttcaattgttaattctttgcctcgactcatagccatgatgagctttgatcatgtt  
tccttaaccctctatTTTACGGAAGAATGATCAAGCTGCTTGCATTGATCATCGTTCTTGCCTCGACTCATAGCCATGATGAGCTTGTGATCATGTTCTTCGATTGAAG

## The OriC region of *Vibrio cholerae*

```
atcaatgatcaacgtaagcttctaaggatcatgatcaagggtgctcacacagttatccacaac  
ctgagtggatgacatcaagataggcggtgtatctccttcctctcgtaactctcatgacca  
cgaaaaagatgatcaagagaggatgattcttggccatatcgcaatgaatacttgtgactt  
gtgcttccaattgacatcttcagcgccatattgcgctggccaagggtgacggagcgggatt  
acgaaagcatgatcatggctgttctgttatcttggtttgactgagacttgttagga  
tagacggttttcatcactgactagccaaagcctactctgcctgacatcgaccgtaaat  
tgataatgaatttacatgcttccgcgacgattacacttgcattgatcatcgatccgattgaag  
atcttcaattgttaattctttgcctcgactcatagccatgatgagctttgatcatgtt  
tccttaaccctctatTTTACGGAAGAATGATCAAGCTGCTGCTGATCATCGTTCTTCA
```

### Hidden Message Problem:

*Find a “hidden message” in the replication origin.*

**Input:** A string *Text* (representing the replication origin of a genome).

**Output:** A hidden message in *Text*.

## Legrand's parchment in The Gold-Bug story

53++!305))6\*;4826)4+.)4+);806\*;48!8'60))85;1+(;:+\*8  
!83(88)5\*!;46(;88\*96\*?;8)\*+(;485);5\*!2:\*+(;4956\*2(5  
\*-4)8'8\*;4069285);)6!8)4++;1(+9;48081;8:8+1;48!85:4  
)485!528806\*81(+9;48;(88;4(+?34;48)4+;161;:188;+?;

## Legrand's parchment in The Gold-Bug story

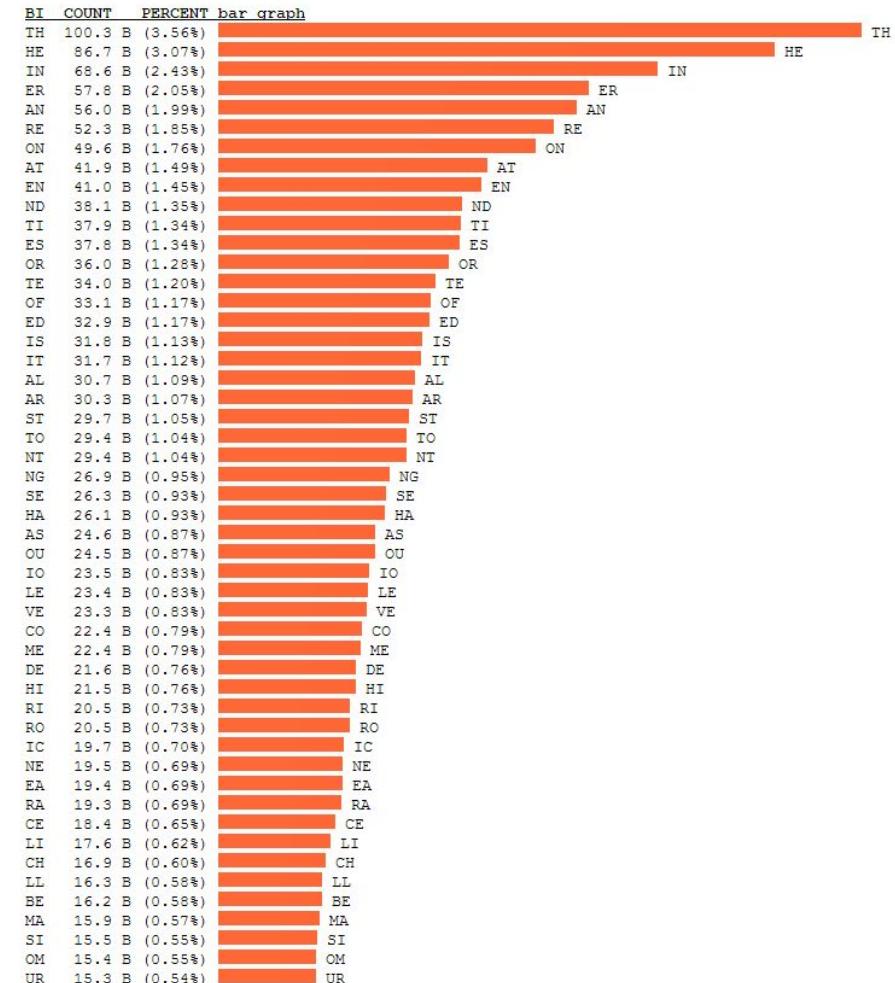
53++!305))6\*;4826)4+. )4+);806\*;48!8‘60))85;1+(;:+\*8  
!83(88)5\*!;46(;88\*96\*?;8)\*+(;485);5\*!2:/\*+(;4956\*2(5  
\*-4)8‘8\*;4069285);)6!8)4++;1(+9;48081;8:8+1;48!85:4  
)485!528806\*81(+9;48;(88;4(+?34;48)4+;161;:188;+?;

53++!305))6\*;**48**26)4+. )4+);806\*;**48**!8‘60))85;1+(;:+\*8  
!83(88)5\*!;46(;88\*96\*?;8)\*+(;**48**5);5\*!2:/\*+(;4956\*2(5  
\*-4)8‘8\*;4069285);)6!8)4++;1(+9;**48**081;8:8+1;**48**!85;4  
)485!528806\*81(+9;**48**;(88;4(+?34;**48**)4+;161;:188;+?;

## Legrand's parchment in The Gold-Bug story

53++!305))6\*;4826)4+.)4+);806\*;48!8'60))85;1+(;:+\*8  
!83(88)5\*!;46(;88\*96\*?;8)\*+(;485);5\*!2:/\*+(;4956\*2(5  
\*-4)8'8\*;4069285);)6!8)4++;1(+9;48081;8:8+1;48!85:4  
)485!528806\*81(+9;48;(88;4(+?34;48)4+;161;:188;+?;

53++!305))6\*THE26)H+. )H+)TE06\*THE!E'60))E5T1+(T:+\*E  
!E3(EE)5\*!TH6(TEE\*96\*?TE)\*+(THE5)T5\*!2:/\*+(TH956\*2(5  
\*-H)E'E\*TH0692E5)T)6!E)H++T1(+9THE0E1TE:E+1THE!E5TH  
)HE5!52EE06\*E1(+9THET(EETH(+?3HTHE)H+T161T:1EET+?T



**LET COUNT PERCENT bar graph**

E	445.2 B	12.49%	E
T	330.5 B	9.28%	T
A	286.5 B	8.04%	A
O	272.3 B	7.64%	O
I	269.7 B	7.57%	I
N	257.8 B	7.23%	N
S	232.1 B	6.51%	S
R	223.8 B	6.28%	R
H	180.1 B	5.05%	H
L	145.0 B	4.07%	L
D	136.0 B	3.82%	D
C	119.2 B	3.34%	C
U	97.3 B	2.73%	U
M	89.5 B	2.51%	M
F	85.6 B	2.40%	F
P	76.1 B	2.14%	P
G	66.6 B	1.87%	G
W	59.7 B	1.68%	W
Y	59.3 B	1.66%	Y
B	52.9 B	1.48%	B
V	37.5 B	1.05%	V
K	19.3 B	0.54%	K
X	8.4 B	0.23%	X
J	5.7 B	0.16%	J
Q	4.3 B	0.12%	Q
Z	3.2 B	0.09%	Z

## Legrand's parchment in The Gold-Bug story

53++!305))6\*;4826)4+.)4+);806\*;48!8'60))85;1+(;:+\*8  
!83(88)5\*!;46(;88\*96\*?;8)\*+(;485);5\*!2:\*+(;4956\*2(5  
\*-4)8'8\*;4069285);)6!8)4++;1(+9;48081;8:8+1;48!85:4  
)485!528806\*81(+9;48;(88;4(+?34;48)4+;161;:188;+?;

A good **glass** in the bishop's hostel in the devil's seat forty-one degrees and thirteen **minutes** northeast and by north main branch seventh limb east side shoot from the left eye of the **death's-head** a bee line from the tree through the **shot** fifty feet out.

## Counting words in a text

ACAACTATGCATACTATCGGGAACTATCCT

We use the term ***k*-mer** to refer to a string of length *k* and define  $\text{COUNT}(\text{Text}, \text{Pattern})$  as the number of times that a *k*-mer *Pattern* appears as a substring of *Text*. Following the above example,

$$\text{COUNT}(\text{ACAACTATGCATACTATCGGGAACTATCCT}, \text{ACTAT}) = 3.$$

Note that  $\text{COUNT}(\text{CGATATAATCCATAG}, \text{ATA})$  is equal to 3 (not 2)

A simple pseudocode for counting patterns in text

```
PATTERNCOUNT(Text, Pattern)
    count  $\leftarrow$  0
    for i  $\leftarrow$  0 to  $|Text| - |Pattern|$ 
        if Text(i,  $|Pattern|$ ) = Pattern
            count  $\leftarrow$  count + 1
    return count
```

What is the most frequent  $k$ -mer in  $Text$ ?

We say that  $Pattern$  is a **most frequent  $k$ -mer** in  $Text$  if it maximizes  $\text{COUNT}(Text, Pattern)$  among all  $k$ -mers. You can see that **ACTAT** is a most frequent 5-mer for  $Text = \text{ACA}\mathbf{ACTAT}GCAT\mathbf{ACTAT}CGGG\mathbf{ACTAT}CCT$ , and **ATA** is a most frequent 3-mer for  $Text = \text{CGA}\mathbf{TATA}\mathbf{TCC}\mathbf{ATAG}$ .

### Frequent Words Problem:

*Find the most frequent  $k$ -mers in a string.*

**Input:** A string  $Text$  and an integer  $k$ .

**Output:** All most frequent  $k$ -mers in  $Text$ .

## What is the most frequent $k$ -mer in $Text$ ?

A straightforward algorithm for finding the most frequent  $k$ -mers in a string  $Text$  checks all  $k$ -mers appearing in this string (there are  $|Text| - k + 1$  such  $k$ -mers) and then computes how many times each  $k$ -mer appears in  $Text$ . To implement this algorithm, called **FREQUENTWORDS**, we will need to generate an array **COUNT**, where  $\text{COUNT}(i)$  stores  $\text{COUNT}(Text, Pattern)$  for  $Pattern = Text(i, k)$  (see Figure 1.2).

<i>Text</i>	<b>A</b>	C	<b>T</b>	G	<b>A</b>	C	<b>T</b>	C	C	C	A	C	C	C	C
COUNT	2	1	1	1	2	1	1	3	1	1	1	3	3	3	3

**FIGURE 1.2** The array **COUNT** for  $Text = \text{ACTGACTCCCACCCCC}$  and  $k = 3$ . For example,  $\text{COUNT}(0) = \text{COUNT}(4) = 2$  because **ACT** (shown in boldface) appears twice in  $Text$  at positions 0 and 4.

## A pseudocode for counting all frequent words

```
FREQUENTWORDS(Text, k)
    FrequentPatterns ← an empty set
    for  $i \leftarrow 0$  to  $|Text| - k$ 
        Pattern ← the  $k$ -mer  $Text(i, k)$ 
        COUNT( $i$ ) ← PATTERNCOUNT(Text, Pattern)
        maxCount ← maximum value in array COUNT
        for  $i \leftarrow 0$  to  $|Text| - k$ 
            if COUNT( $i$ ) = maxCount
                add  $Text(i, k)$  to FrequentPatterns
        remove duplicates from FrequentPatterns
    return FrequentPatterns
```

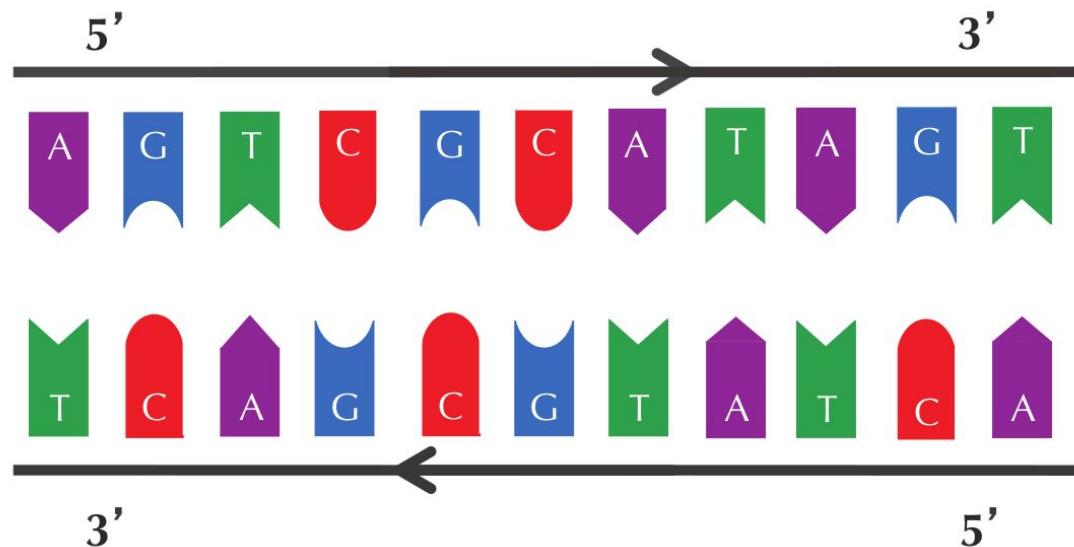
# Frequent words in *Vibro chlorea*

Do any of the counts in *Vibrio cholerae* seem exceptionally large?

atcaatgatcaacgtaagcttctaagc**ATGATCAAG**gtgctcacacagttatccacaac  
ctgagtgatgacatcaagataggcggtgttatctccttcctctcgtaactctcatgacca  
cgaaag**ATGATCAAG**agaggatgatttctggccatatcgaaatgaatacttgtgactt  
gtgcttccaattgacatcttcagcgccatattgcgtggccaagggtgacggagcgggatt  
acgaaagcatgatcatggctgttctgttatcttgcgtttgactgagacttgttagga  
tagacggttttcatcactgactgccaaggcctactctgcctgacatcgaccgtaaat  
tgataatgaatttacatgctccgcgacgattacacttgcattgtatcatcgatccgattgaag  
atcttcaattgttaattcttgcctcgactcatagccatgatgagctttgatcatgtt  
tccttaaccctctatTTTACGGAAGA**ATGATCAAG**ctgctgctttgatcatcgttc

$$P = 0.000769$$

DNApoly can read the complementary strand too



Do any of the counts in *Vibrio cholerae* seem exceptionally large?

atcaatgatcaacgtaagcttctaagc**ATGATCAAG**gtgctcacacagttatccacaac  
ctgagtggatgacatcaagataggcggtgttatctccttcgtactctcatgacca  
cgaaaa**ATGATCAAG**agaggatgattcttgccatatcgcaatgaatacttgtgactt  
gtgcttccaattgacatcttcagcgccatattgcgtggccaagggtgacggagcgggatt  
acgaaagcatgatcatggctgttgttatcttgcgtttgactgagacttgttagga  
tagacggttttcatcactgacttagccaaagcctactctgcctgacatcgaccgtaaat  
tgataatgaatttacatgctccgcgacgattacct**CTTGATCAT**cgatccgattgaag  
atcttcaattgttaattcttgcctcgactcatagccatgatgagct**CTTGATCAT**gtt  
tccttaaccctctatTTTacggaaga**ATGATCAAG**ctgctgct**CTTGATCAT**cgtttc

## **Pattern Matching Problem:**

*Find all occurrences of a pattern in a string.*

**Input:** Strings *Pattern* and *Genome*.

**Output:** All starting positions in *Genome* where *Pattern* appears as a sub-string.

**ATGATCAAG** appears 17 times in the genome

116556, 149355, **151913**, **152013**, **152394**, 186189, 194276, 200076, 224527,  
307692, 479770, 610980, 653338, 679985, 768828, 878903, 985368

Let's check if the proposed message can be found  
in the same region of *Thermotoga petrophila*?

atcaatgatcaacgtaagcttctaaggatcatgatcaagggtgctcacacagttatccacaac  
ctgagtggatgacatcaagataggcggttatctccttcctcgtaactctcatgacca  
cgaaagatgatcaagagaggatgatttttgcgcattatcgcaatgaataacttgtgactt  
gtgcttccaattgacatcttcagcgccatattgcgcgtggccaaaggtagcgagcgaggatt  
acgaaagcatgatcatggctgttctgttatcttggtttgactgagacttgttagga  
tagacggttttcatcactgacttagccaaagcctactctgcctgacatcgaccgtaat  
tgataatgaatttacatgctccgcgacgattacctcttgatcatcgatccgattgaag  
atcttcaattgttaattctttgcctcgactcatagccatgatgagctttgatcatgtt  
tccttaaccctctattttacggaagaatgatcaagctgctgcttgcattcgatcatcgttc

Let's check if the proposed message can be found  
in the same region of *Thermotoga petrophila*?

atcaatgatcaacgtaagcttctaaggatcatgatcaagggtgctcacacagttatccacaac  
ctgagtggatgacatcaagataggcggttatctccttcctcgtaactctcatgacca  
cgaaagatgatcaagagaggatgatttttgcgcattatcgcaatgaataacttgtgactt  
gtgcttccaattgacatcttcagcgccatattgcgcgtggccaaaggtagcgagcgaggatt  
acgaaagcatgatcatggctgttctgttatcttggtttgactgagacttggtagga  
tagacggttttcatcactgacttagccaaagccttactctgcctgacatcgaccgtaaat  
tgataatgaatttacatgctccgcgacgattacctcttgatcatcgatccgattgaag  
atcttcaattgttaattctttgcctcgactcatagccatgatgagctttgatcatgtt  
tccttaacccttattttacggaagaatgatcaagctgctgatcatcgtttc

This region does not contain a single occurrence of **ATGATCAAG** or **CTTGATCAT**!

Application of the Frequent Words Problem to the *oriC* region above reveals that the following six 9-mers appear in this region three or more times:

AACCTACCA

AAACCTACC

ACCTACCAC

CCTACCACC

GGTAGGTTT

TGGTAGGTT

aactctatacctcctttgtcgaatttgtgtgatttatagagaaaatcttattaactga  
aactaaaatggtagggtt **GGTGGTAGG**ttttgtgtacattttgttagtatctgattttaa  
ttacataccgtatattgtattaaattgacgaacaattgcattgaaattgaatatatgcaa  
acaaa**CCTACCACC**aaactctgtattgaccattttaggacaacttcag**GGTGGTAGG**ttt  
ctgaagctctcatcaatagactattttagtctttacaaacaatattaccgttcagattca  
agattctacaacgctgtttaatggcggtgcagaaaacttaccacctaattccagtat  
ccaagccgatttcagagaaacctaccacttacccactta**CCTACCACC**cgggtggta  
agttgcagacattattaaaaacctcatcagaagctgttcaaaaattcaataactcgaaa  
**CCTACCACC**tgcgtcccattatttactactaataatgcgtataattgatctga

## How can we find *oriC* in a newly sequenced genome?

Our plan is to slide a window of fixed length  $L$  along the genome, looking for a region where a  $k$ -mer appears several times in short succession. The parameter value  $L = 500$  reflects the typical length of *oriC* in bacterial genomes.

### **Clump Finding Problem:**

*Find patterns forming clumps in a string.*

**Input:** A string *Genome*, and integers  $k$ ,  $L$ , and  $t$ .

**Output:** All distinct  $k$ -mers forming  $(L, t)$ -clumps in *Genome*.

## How can we find *oriC* in a newly sequenced genome?

Our plan is to slide a window of fixed length  $L$  along the genome, looking for a region where a  $k$ -mer appears several times in short succession. The parameter value  $L = 500$  reflects the typical length of *oriC* in bacterial genomes.

### Clump Finding Problem:

*Find patterns forming clumps in a string.*

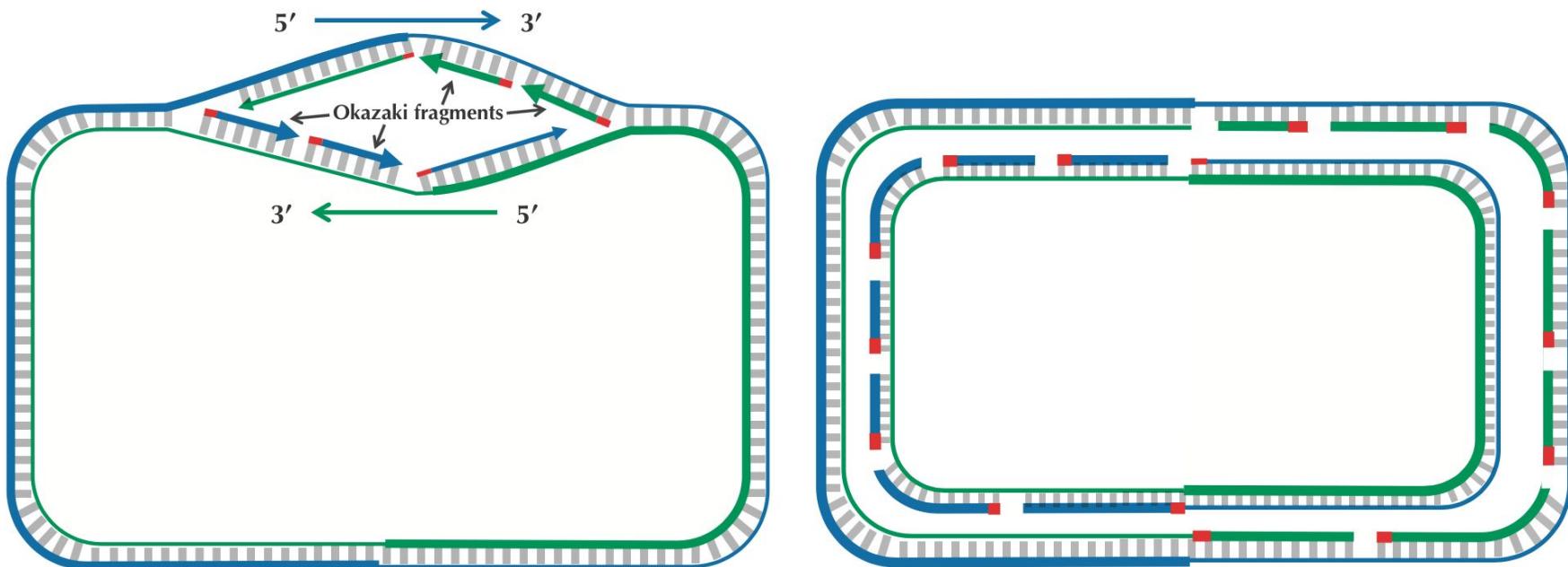
**Input:** A string *Genome*, and integers  $k$ ,  $L$ , and  $t$ .

**Output:** All distinct  $k$ -mers forming  $(L, t)$ -clumps in *Genome*.

*Escherichia coli (E. coli)*

We find hundreds of different 9-mers  
forming  $(500, 3)$ -clumps

# How can our knowledge about DNA replication help?

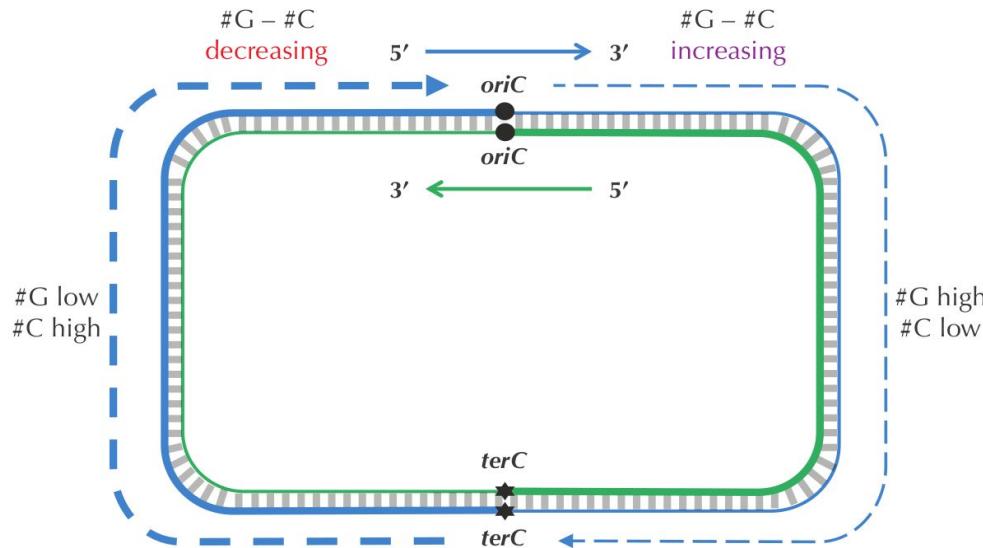


Let's compare the nucleotide counts of the F and R strands

	#C	#G	#A	#T
<b>Entire strand</b>	427419	413241	491488	491363
<b>Reverse half-strand</b>	219518	201634	243963	246641
<b>Forward half-strand</b>	207901	211607	247525	244722
<b>Difference</b>	+11617	-9973	-3562	+1919

# The skew diagram

	#C	#G	#G - #C
<b>Entire strand</b>	427419	413241	
<b>Reverse half-strand</b>	219518	201634	-17884
<b>Forward half-strand</b>	207901	211607	3706
<b>Difference</b>	+11617	-9973	



# The skew diagram

