



KAUST Academy - Bioinformatics Program 2024

Instructor: Azza Althagafi

Ph.D. in Computer Science (Bioinformatics)

azza.althagafi@kaust.edu.sa

TA: Ibtisam Saleh

ibtisamssaleh@outlook.com

Riyadh, 25-28 Feb, 2024



Course Overview



Course Overview



أكاديمية كاوهست
KAUST ACADEMY

Reference Book: Bioinformatics Algorithms, Phillip Compeau and Pavel Pevzner ([Link](#))

Day	Date	Topic
1	Sun 25 Feb	<u>Chapter 1: The Origins of Replication (Part#1)</u>
2	Mon 26 Feb	<u>Chapter 1: The Origins of Replication (Part#2)</u> <u>Introduction to coding (Python)</u>
3	Tues 27 Feb	<u>Chapter 3: Genome Assembly</u>
4	Wed 28 Feb	<u>Exam (9-12)</u> <u>Chapter 5: Sequence Alignment (1-5)</u>



Course Materials and Policies



- **Day 1 (25 Feb) – Day 4 (28 Feb):**
 - All day 8 am - 5 pm
 - Theoretical part
 - Practical part (Coding): Using [Rosalind](#)
- All the slides will be uploaded in:
<http://tinyurl.com/27chtrbb>



Course Material

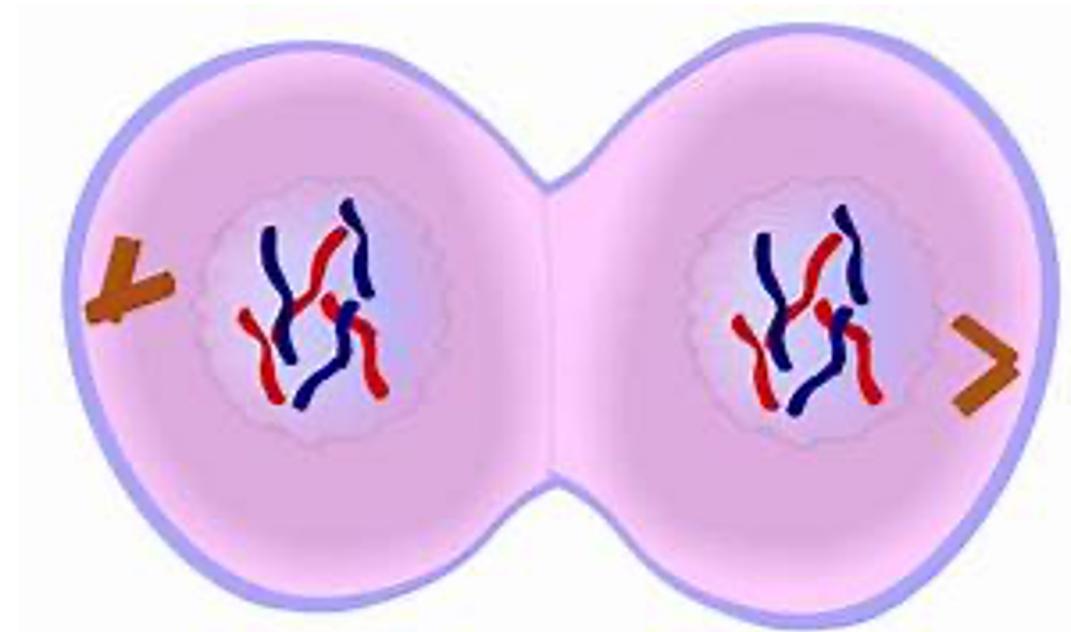


Where in a Genome Does DNA Replication Begin?

- A Journey of a Thousand Miles...

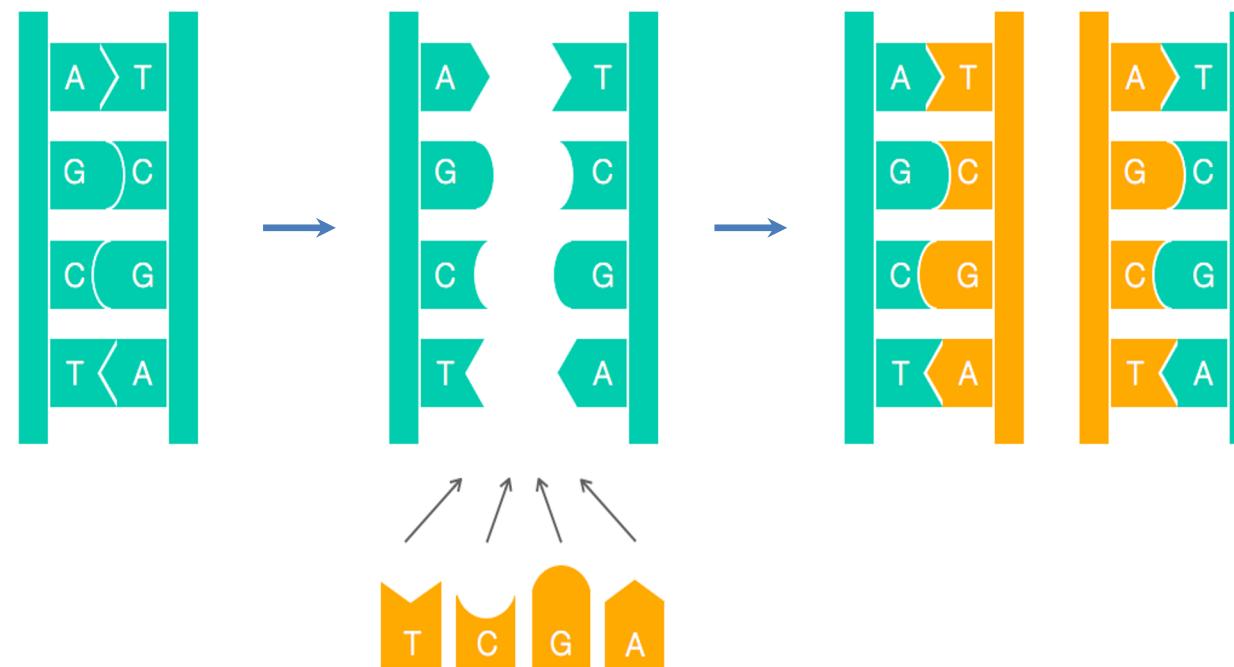
The majority of credit for these slides are attributed to Pavel and Phillip @ *Bioinformatics Algorithms: an Active Learning Approach*

**Genome replication is one of the most important tasks carried out in the cell!
Before a Cell Divides, it Must Replicate its Genome**



The two strands of the parent DNA molecule unwind during replication, and then each parent strand acts as a template for the synthesis of a new strand.

As a result, the replication process begins with a pair of complementary strands of DNA and ends with two pairs of complementary strands.



Although this figure successfully models DNA replication on a simple level, the details of replication turned out to be much more intricate than Watson and Crick imagined. An astounding amount of molecular logistics is required to ensure DNA replication.

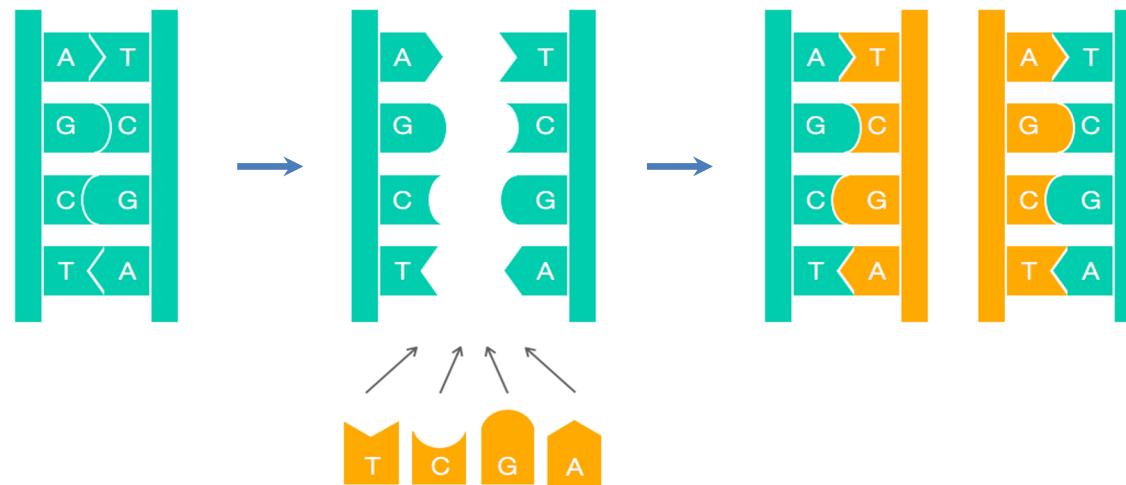


Figure: A naive view of replication. Nucleotides adenine (A) and thymine (T) are complements of each other, as are cytosine (C) and guanine (G). Complementary nucleotides bind to each other in DNA.

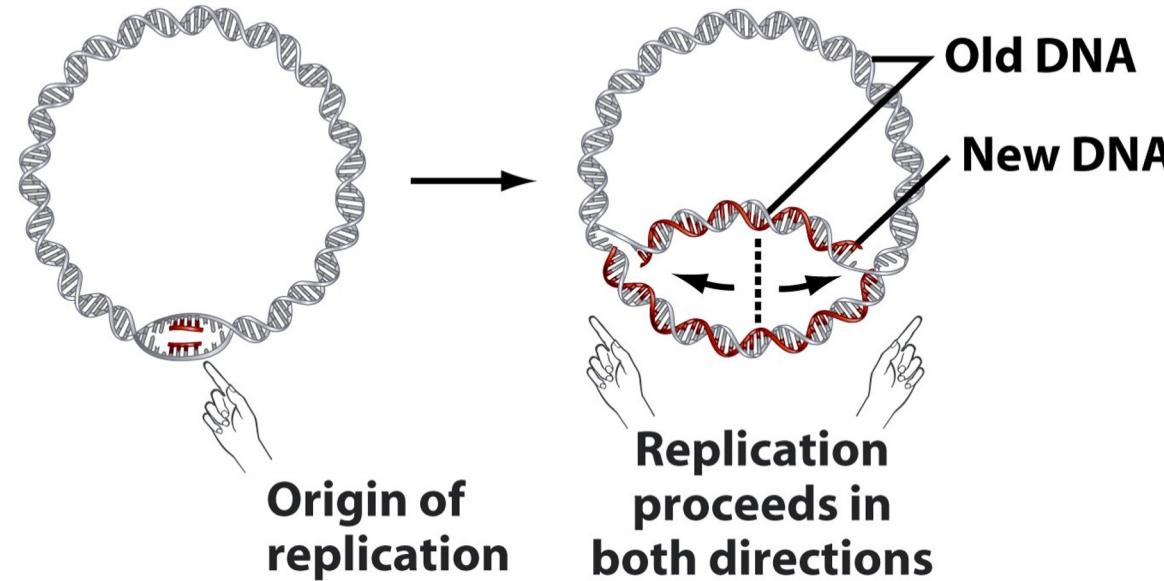
At first glance, a computer scientist might not imagine that these details have any computational relevance.

To mimic the process in the above figure algorithmically, we only need to take a string representing the genome and return a copy of it!

Where in a genome does it all begin?

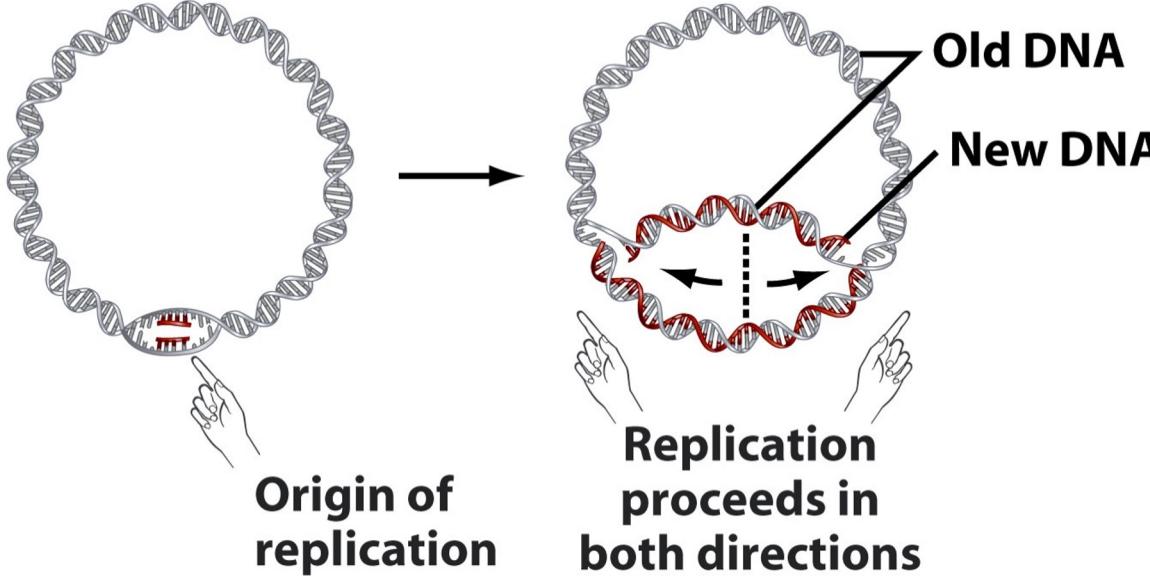


Replication begins in a region called the replication origin (*oriC*)



- Locating *ori* is important for understanding cell replication and has implications for biomedical problems.
- Gene therapy methods use viral vectors, genetically engineered mini-genomes that can penetrate cell walls, to carry artificial genes for various applications.
- Viral vectors have been used in agriculture.
- Gene therapy successfully saved many lives.

Replication begins in a region called the replication origin (*oriC*)



Where in a genome does it all begin?

Outline

- Search for Hidden Messages in Replication Origin
 - What is a Hidden Message in Replication Origin?
 - Some Hidden Messages are More Surprising than Others
 - Clumps of Hidden Messages
- From a Biological Insight toward an Algorithm for Finding Replication Origin
 - Asymmetry of Replication
 - Why would a computer scientist care about assymetry of replication?
 - Skew Diagrams
 - Finding Frequent Words with Mismatches
 - Open Problems

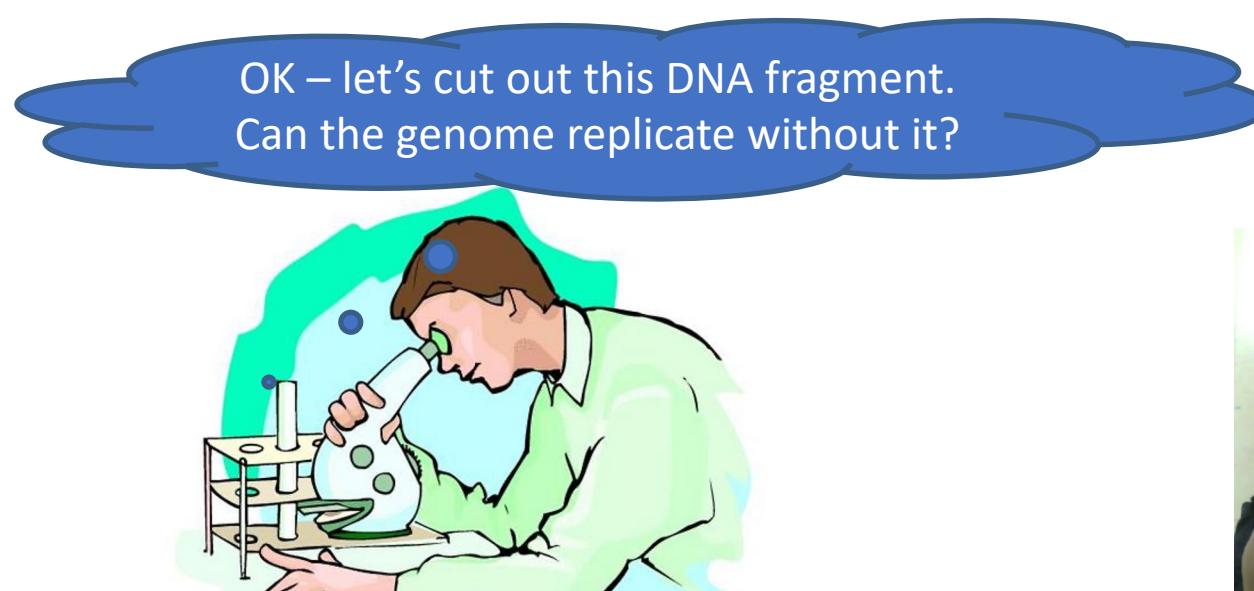


Finding Origin of Replication

In the following problem, we assume that a genome has a single *ori* and is represented as a **DNA string**, or a string of nucleotides from the four-letter alphabet {A, C, G, T}.

Finding *oriC* Problem: Finding *oriC* in a genome.

- **Input.** A genome.
- **Output.** The location of *oriC* in the genome.

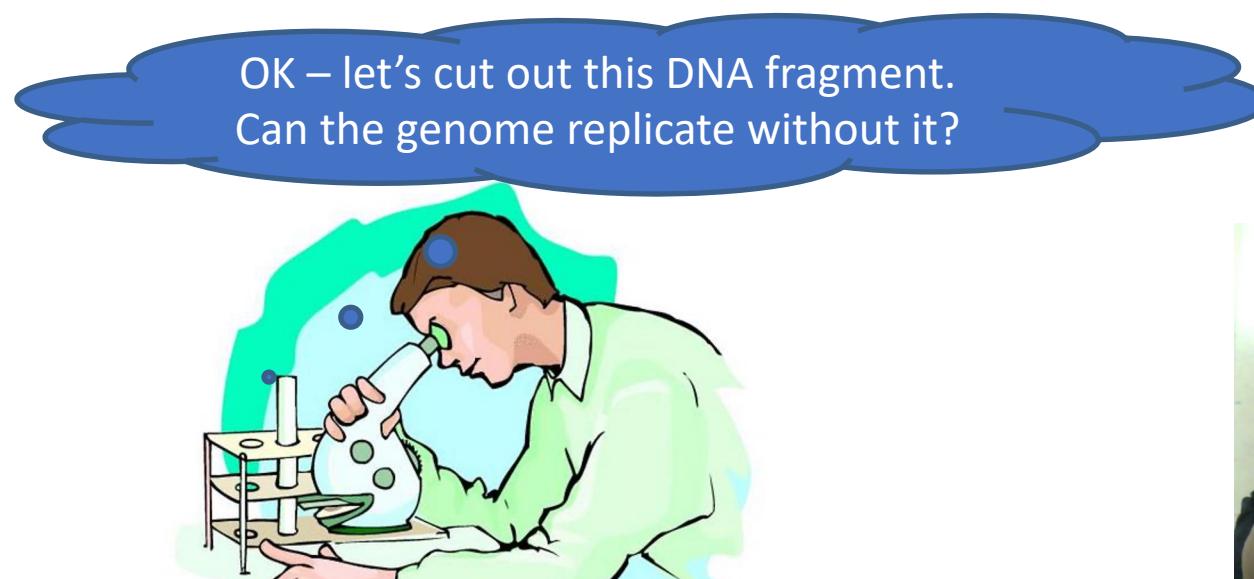


Finding Origin of Replication

In the following problem, we assume that a genome has a single *ori* and is represented as a **DNA string**, or a string of nucleotides from the four-letter alphabet {A, C, G, T}.

Finding *oriC* Problem: Finding *oriC* in a genome.

- **Input.** A genome.
- **Output.** The location of *oriC* in the genome.



Finding Origin of Replication

Finding *oriC* Problem: Finding *oriC* in a genome.

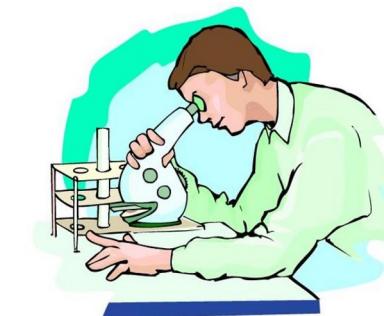
- **Input.** A genome.
- **Output.** The location of *oriC* in the genome.

First, these methods are much faster than experimental approaches.

Second, the results of many experiments cannot be interpreted without computational analysis.

In particular, existing experimental approaches to *ori* prediction are rather time consuming. As a result, *ori* has only been experimentally located in a handful of species.

Thus, we would like to design a computational approach to find *ori* so that biologists are free to spend their time and money on other tasks.



How Does the Cell Know to Begin Replication in Short *oriC*?

In this chapter, we will focus on the relatively easy case of finding *ori* in bacterial genomes, most of which consist of a single circular chromosome. Research has shown that the region of the bacterial genome encoding *ori* is typically a few hundred nucleotides long. Our plan is to begin with a bacterium in which *ori* is known, and then determine what makes this genomic region special in order to design a computational approach for finding *ori* in other bacteria.

Replication origin of *Vibrio cholerae* (≈ 500 nucleotides):



```
atcaatgatcaacgtaaagcttctaaggcatgatcaagggtgctcacacagtttatccacaac  
ctgagtggatgacatcaagataggtcggttatctccttcctcgactctcatgacca  
cgaaaagatgatcaagagaggatgattttgcgcattatcgcaatgaataacttgtgactt  
gtgcttccaattgacatcttcagcgccatattgcgcgtggccaaaggtagcgaggcgggatt  
acgaaagcatgatcatggctgtttctgtttatcttgactgagacttgttagga  
tagacggttttcatcactgactgccaaggcctactctgcgtacatcgaccgtaaat  
tgataatgaatttacatgctccgcgacgatttacctcttgcattttactgatccgattgaag  
atcttcaattgttaattcttgcctcgactcatggcatgatgagctttgatcatgtt  
tccttaaccctctattttacggaagaatgatcaagctgctttgatcatcgttc
```

There must be a **hidden message** telling the cell to start replication here.

How Does the Cell Know to Begin Replication in Short *oriC*?



There must be a **hidden message** telling the cell to start replication here.

How does the bacterial cell know to begin replication exactly in this short region within the much larger *Vibrio cholerae* chromosome, which consists of 1,108,250 nucleotides?

There must be some “hidden message” in the *ori* region ordering the cell to begin replication here.

Indeed, we know that the initiation of replication is mediated by *DnaA*, a protein that binds to a short segment within the *ori* known as a ***DnaA* box**.

You can think of the *DnaA* box as a message within the DNA sequence telling the *DnaA* protein: “bind here!”

The question is how to find this hidden message without knowing what it looks like in advance—can you find it? In other words, can you find something that stands out in *ori*? This discussion motivates the following problem.



The Hidden Message Problem



STOP and Think: Does the Hidden Message Problem represent a clearly stated computational problem?

أكاديمية KAUST
KAUST ACADEMY

Hidden Message Problem. Finding a hidden message in a string.

- **Input.** A string *Text* (representing replication origin).
- **Output.** A hidden message in *Text*.

This is not a computational problem either!



The notion of “**hidden message**” is not precisely defined.

The *ori* region of *Vibrio cholerae* is currently just as puzzling as the parchment discovered by William Legrand in Edgar Allan Poe's story "The Gold-Bug". Written on the parchment was



“The Gold-Bug” Problem



53++!305))6*;4826)4+.)4+);806*;4
8!8`60))85;]8*:+*8!83(88)5*!;46(
;88*96*?;8)*+(;485);5*!2:*+(;495
6*2(5*4)8`8*;4069285);)6!8)4++;1
(+9;48081;8:8+1;48!85;4)485!5288
06*81(+9;48;(88;4(+?34;48)4+;161
;:188;+?;

A secret message left by pirates
("The Gold-Bug" by Edgar Allan Poe)

“The Gold-Bug” Problem



53++!305))6*;4826)4+.)4+);806*;4
8!8`60))85;]8*:+*8!83(88)5*!;46(
;88*96*?;8)*+(;485);5*!2:*+(;495
6*2(5*4)8`8*;4069285);)6!8)4++;1
(+9;48081;8:8+1;48!85;4)485!5288
06*81(+9;48;(88;4(+?34;48)4+;161
;:188;+?;

Upon seeing the parchment, the narrator remarks, "Were all the jewels of Golconda awaiting me upon my solution of this enigma, I am quite sure that I should be unable to earn them

Why is “;48” so Frequent?

Hint: The message is in English

```
53++!305))6*;4826)4+. )4+);806*48
!8`60))85;]8*:+*8!83(88)5*!46(88
*96*?;8)*+(;485);5*!2:*+(;4856*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;+?;
```

“THE” is the Most Frequent English Word



53++!305))6***THE**26)4+.)4+)806***THE**
!8`60))85;]8*:+*8!83(88)5*!;46(;
88*96*?;8)*+(**THE**5);5*!2:*+(;4956
*2(5*4)8`8*;4069285);)6!8)4++;1(
+9**THE**081;8:8+1**THE**!85;4)485!52880
6*81(+9**THE**; (88;4(+?34**THE**)4+;161;
:188;+?;



Could you Complete Decoding the Message?

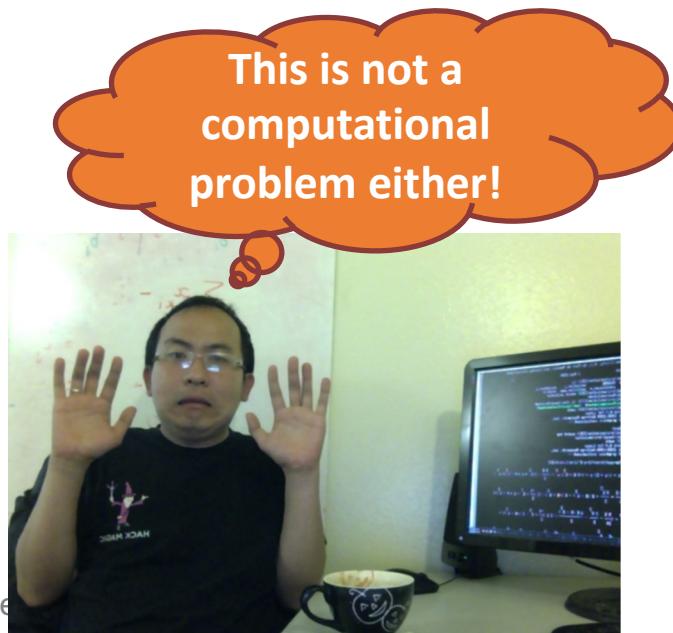
53++!305)) 6***THE**26) **H**+.) **H**+) 806***THE**
!**E**`60))**E**5;]**E***:+***E**!**E**3(**EE**)5*!**TH**6(T
EE*96*?;**E**)*+(**THE**5)**T**5*!2:/*+ (**TH**956
*2(5***H**)**E`E*****TH**0692**E**5)**T**)6!**E**)**H**++**T**1(
+9**THE**0**E**1**TE**:**E**+1**THE**!**E**5**T**4)**HE**5!52**88**0
6***E**1(+9**THE****T**(**EETH**(+?34**THE**)**H**+**T**161**T**
:1**EET**+?**T**

Operating under the assumption that DNA is a language of its own, let's borrow Legrand's method and see if we can find any surprisingly frequent "words" within the *ori* of *Vibrio cholerae*

The Hidden Message Problem Revisited

Hidden Message Problem. Finding a hidden message in a string.

- **Input.** A string *Text* (representing *oriC*).
- **Output.** A hidden message in *Text*.

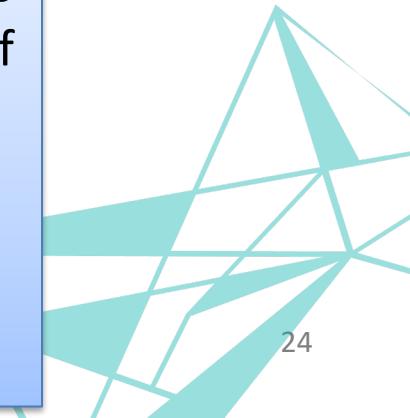


The notion of “**hidden message**” is not precisely defined.

Hint: For various biological signals, certain words appear surprisingly frequently in small regions of the genome.

AATTT is a surprisingly frequent 5-mer in:

ACA**AATTTGCAT****AATTT**CGGGAA**AATTCCT**



The Frequent Words Problem

Frequent Words Problem. Finding most frequent k -mers in a string.

- **Input.** A string $Text$ and an integer k .
- **Output.** All **most frequent k -mers** in $Text$.

This is better, but where is
the definition of “a most
frequent k -mer?”



We will 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 .

The Frequent Words Problem

Frequent Words Problem. Finding most frequent k -mers in a string.

- **Input.** A string $Text$ and an integer k .
- **Output.** All **most frequent k -mers** in $Text$.



A k -mer **Pattern** is a **most frequent k -mer** in a text if no other k -mer is more frequent than **Pattern**.

AATTT is a most frequent 5-mer in:

ACA**AATTT**GCAT**AATTT**CGGGA**AATTT**CCT

Count(ACA**ACTAT**GC**ACTAT**CGGGA**ACTAT**CCT,
ACTAT) = 3

Does the Frequent Words Problem Make Sense to Biologists?

Frequent Words Problem. Finding most frequent k -mers in a string.

- **Input.** A string $Text$ and an integer k .
- **Output.** All **most frequent k -mers** in $Text$.

Replication is performed by **DNA polymerase** and the initiation of replication is mediated by a protein called ***DnaA***.

DnaA binds to short (typically 9 nucleotides long) segments within the replication origin known as a ***DnaA box***.

A *DnaA* box is a hidden message telling *DnaA*: “**bind here!**” And *DnaA* wants to see multiple *DnaA* boxes.

What is the Runtime of Your Algorithm?

Frequent Words Problem. Finding most frequent k -mers in a string.

- **Input.** A string $Text$ and an integer k .
- **Output.** All **most frequent k -mers** in $Text$.

- $|Text|^2 \cdot k$
- $4^k + |Text| \cdot k$
- $|Text| \cdot k \cdot \log(|Text|)$???
- $|Text|$

You will later see how a **naive and slow** algorithm with $|Text|^2 \cdot k$ runtime can be turned into a **fast** algorithm with $|Text|$ runtime ($|Text|$ stands for the length of string $Text$)

Code Challenge: Implement PatternCount()



```
PatternCount(Text, Pattern)
    count ← 0
    for i ← 0 to |Text| - |Pattern|
        if Text(i, |Pattern|) = Pattern
            count ← count + 1
    return count
```

Let's refresh our foundational knowledge of Python.



Python



- Python Syntax
- String and Console Output
- Conditionals and control flow
- Functions
- Lists and Dictionaries

www.codecademy.com



Code Challenge:

Implement PatternCount()



```
PatternCount(Text, Pattern)
    count ← 0
    for i ← 0 to |Text| - |Pattern|
        if Text(i, |Pattern|) = Pattern
            count ← count + 1
    return count
```



Code Challenge:

Implement PatternCount()



```
PatternCount(Text, Pattern)
    count ← 0
    for i ← 0 to |Text| - |Pattern|
        if Text(i, |Pattern|) = Pattern
            count ← count + 1
    return count
```

```
1 import sys
2
3 # Please do not remove package declarations because these are used by the autograder.
4
5 # Insert your PatternCount function here, along with any subroutines you need
6 def pattern_count(text: str, pattern: str) -> int:
7     count = 0
8     overlap = len(text) - len(pattern) + 1
9     for i in range(overlap):
10         start = i
11         end = i + len(pattern)
12         if text[start:end] == pattern:
13             count += 1
14
15     return count
```



The Frequent Words Problem



We say that *Pattern* is a **most frequent k -mer** in *Text* if it maximizes $\text{Count}(\text{Text}, \text{Pattern})$ among all k -mers.

- You can see that **ACTAT** is a most frequent 5-mer of ACAACTATGCATACTATCGGGAACTATCCT,
- and **ATA** is a most frequent 3-mer of CGATATATCCATAG.



The Frequent Words Problem

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 *Count(i)* stores *Count(Text, Pattern)* for *Pattern* = *Text(i, k)*

<i>Text</i>	A	C	T	G	A	C	T	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: The array *Count* for *Text* = ACTGACTCCCACCCCC and *k* = 3. For example, *Count(0)* = *Count(4)* = 2 because ACT (shown in boldface) appears twice in *Text*.

The Frequent Words Problem

```
FrequentWords(Text, k)
    FrequentPatterns ← an empty set
    for i ← 0 to |Text| - k
        Pattern ← the k-mer Text(i, k)
        Count(i) ← PatternCount(Text, Pattern)
    maxCount ← maximum value in array Count
    for i ← 0 to |Text| - k
        if Count(i) = maxCount
            add Text(i, k) to FrequentPatterns
    remove duplicates from FrequentPatterns
    return FrequentPatterns
```

STOP and Think: Take another look at the pseudocode for **FrequentWords()**. How fast is this algorithm?

The Frequent Words Problem

```
FrequentWords(Text, k)
    FrequentPatterns ← an empty set
    for i ← 0 to |Text| - k
        Pattern ← the k-mer Text(i, k)
        Count(i) ← PatternCount(Text, Pattern)
        maxCount ← maximum value in array Count
    for i ← 0 to |Text| - k
        if Count(i) = maxCount
            add Text(i, k) to FrequentPatterns
    remove duplicates from FrequentPatterns
    return FrequentPatterns
```

Although **FrequentWords()** finds most frequent k -mers, it is not very efficient. Each call to **PatternCount($Text, Pattern$)** checks whether the k -mer $Pattern$ appears in position 0 of $Text$, position 1 of $Text$, and so on. Since each k -mer requires $|Text| - k + 1$ such checks, each one requiring as many as k comparisons, the overall number of steps of **PatternCount($Text, Pattern$)** is $(|Text| - k + 1) \cdot k$.

Furthermore, **FrequentWords()** must call **PatternCount()** $|Text| - k + 1$ times (once for each k -mer of $Text$), so that its overall number of steps is $(|Text| - k + 1) \cdot (|Text| - k + 1) \cdot k$. To simplify the matter, computer scientists often say that the runtime of **FrequentWords()** has an upper bound of $|Text|^2 \cdot k$ steps and refer to the **complexity** of this algorithm as $O(|Text|^2 \cdot k)$.

A Faster Frequent Words Approach

If you were to solve the Frequent Words Problem by hand for a small example, you would probably form a table like the one in the figure below for *Text* equal to "ACGTTTCACGTTTACGG" and k equal to 3.

You would slide a length- k window *Text*, and if the current k -mer substring of *text* does not occur in the table, then you would create a new entry for it.

Otherwise, you would add 1 to the entry corresponding to the current k -mer substring of *Text*. We call this table the **frequency table** for *Text* and k .

ACG	CGT	GTT	TTT	TTC	TCA	CAC	TTA	TAC	CGG
3	2	2	3	1	1	1	1	1	1

A Faster Frequent Words Approach



أكاديمية كاوهست
KAUST ACADEMY

```
FrequencyTable(Text, k)
    freqMap ← empty map
    n ← |Text|
    for i ← 0 to n - k
        Pattern ← Text(i, k)
        if freqMap[Pattern] doesn't exist
            freqMap[Pattern] ← 1
        else
            freqMap[Pattern] ← freqMap[Pattern]+1
    return freqMap
```

```
BetterFrequentWords(Text, k)
    FrequentPatterns ← an array of strings of length 0
    freqMap ← FrequencyTable(Text, k)
    max ← MaxMap(freqMap)
    for all strings Pattern in freqMap
        if freqMap[Pattern] = max
            append Pattern to frequentPatterns
    return frequentPatterns
```



Code Challenge: Solve the Frequent Words Problem.



[ROSLIND](#)



KAUST Academy

Code Challenge: Solve the Frequent Words Problem.

```
3 # Please do not remove package declarations because these are used by the autograder.
4 def pattern_count(text: str, pattern: str) -> int:
5     count = 0
6     overlap = len(text) - len(pattern) + 1
7     for i in range(overlap):
8         start = i
9         end = i + len(pattern)
10        if text[start:end] == pattern:
11            count += 1
12
13    return count
14 # Insert your frequent_words function here, along with any subroutines you need
15 def frequent_words(text: str, k: int) -> list[str]:
16     """Find the most frequent k-mers in a given text."""
17     frequent_patterns = set()
18     count = {}
19
20     # iterate through DNA Text and count kmers
21     for i in range(len(text) - k):
22         pattern = text[i : i + k]
23         # add pattern to dictionary
24         # key is i, start position of kmer
25         # value is count, using PatternCount function from 1A
26         count[i] = pattern_count(text, pattern)
27
28         # find maximum count value in dictionary
29         max_count = max(count.values())
30
31     # iterate through text again, if count at that position is max count, slice that kmer and add to set
32     for position in range(len(text) - k):
33         if count[position] == max_count:
34             frequent_patterns.add(text[position : position + k])
35
36 return frequent_patterns
```



Frequent Words in *Vibrio cholerae*

k	3	4	5	6	7	8	9
count	25	12	8	8	5	4	3
k-mers	tga	atga	gatca	tgatca	atgatca	atgatcaa	atgatcaag
			tgatc			cttgatcat	
						tcttgcata	
						ctcttgatc	

The table above reveals the most frequent k -mers in the *ori* region from *Vibrio cholerae*, along with the number of times that each k -mer occurs.

STOP and Think: Do any of the counts in the above table seem surprisingly large?

Frequent Words in *Vibrio cholerae*

-For example, the 9-mer **ATGATCAAG** appears three times in the *ori* region of *Vibrio cholerae*—is it surprising?

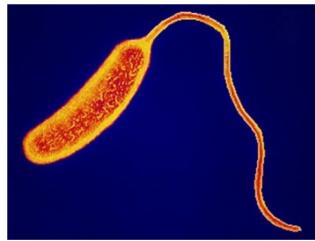
```
atcaatgatcaacgtaagcttctaagcATGATCAAGgtgctcacacagtttatccacaac  
ctgagtgatgacatcaagataggtagtcgttatctccttcctcgtaactctcatgacca  
cgaaagATGATCAAGagaggatgattttggccatatcgcaatgaataacttgtgactt  
gtgcttccaattgacatcttcagcgccatattgcgctggccaaggtagcgagcggatt  
acgaaagcatgatcatggctgttctgttatcttgcgtttgactgagacttgttagga  
tagacggttttcatcaactgacttagccaaagcctactctgcctgacatcgaccgtaaat  
tgataatgaatttacatgctccgcacgattacaccttgcgtttgactgatccgattgaag  
atcttcaattgttaattcttcgcactcatagccatgatgagcttgcgtttgactgatgtt  
tccttaaccctctatTTTACGGAAGAATGATCAAGctgctgcttgcgtttgactcatcgttc
```

We highlight a most frequent 9-mer instead of using some other value of k because experiments have revealed that bacterial *DnaA* boxes are usually nine nucleotides long.

Outline

- **Search for Hidden Messages in Replication Origin**
 - What is a Hidden Message in Replication Origin?
 - **Some Hidden Messages are More Surprising than Others**
 - Clumps of Hidden Messages
- **From a Biological Insight toward an Algorithm for Finding Replication Origin**
 - Asymmetry of Replication
 - Why would a computer scientist care about assymetry of replication?
 - Skew Diagrams
 - Finding Frequent Words with Mismatches
 - Open Problems

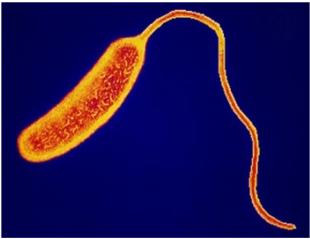
oriC of Vibrio cholerae



```
atcaatgatcaacgtAACGTTCTAAGCATGATCAAGGTGCTCACACAGTTATCCACAACCTGAGTGG
atgacatcaagatAGGTGCGTGTATCTCCTCCTCGTACTCTCATGACCACGAAAGATGATCAAG
agaggatgattCTTGGCCATATCGCAATGAATACTTGTGACTTGTGCTTCCAATTGACATCTTCAGC
GCCATATTGCGCTGGCCAAGGTGACGGAGCGGGATTACGAAAGCATGATCATGGCTGTTGTTCTGTT
ATCTTGTGTTGACTGAGACTTGTGTTAGGATAGACGGTTTTCATCACTGACTAGCCAAGCCTTACTCT
GCCTGACATCGACCCTAAATTGATAATGAATTACATGCTTCCGCACGATTACCTCTGATCATCG
ATCCGATTGAAGATCTTCAATTGTTAATTCTCTTGCCCTCGACTCATAGCCATGATGAGCTTGTCA
TGTTCTTAACCCTCTATTACGGAAGAATGATCAAGCTGCTGATCATCGTTTC
```



Too Many Frequent Words – Which One is a Hidden Message?



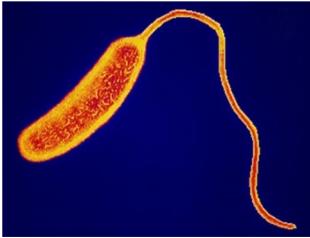
```
atcaatgatcaacgtAACGTTCTAAGCATGATCAAGgtgctcacacagtttatccacaacacctgagg  
atgacatcaagataggTCGGTGTatctcCTTCCtCTCGTactctcatgaccacggaaagATGATCAAG  
agaggatgattCTTGGCCatATCGCAATGAATACTTGTgacttGTgCTTCCAATTGACATCTTCAGC  
GCCATATTGCGCTGGCCAAGGTGACGGAGCGGGATTACGAAAGCATGATCATGGCTGTTGTTCTGTT  
ATCTTGTGTTGACTGAGACTTGTAGGATAGACGGTTTCTCACTGACTAGCCAAGCCTTACTCT  
GCCTGACATCGACCCTAAATTGATAATTACATGCTTCCGCGACGATTACCTCTTGATCATC  
ATCCGATTGAAGATCTTCAATTGTTAATTCTCTTGCCTCGACTCATGCCATGATGAGCTCTTGATCA  
TGTTTCTTAACCCCTCTATTACGGAGAATGATCAAGCTGCTGCTCTTGATCATCgtttc
```

Most frequent 9-mers in this *oriC* (all appear 3 times):

ATGATCAAG, CTTGATCAT, TCTTGGATCA, CTCTTGATC

Is it **STATISTICALLY** surprising to find a 9-mer appearing **3 or more** times within ≈ 500 nucleotides?

Hidden Message Found!



atcaatgatcaacgtaagcttctaagc**ATGATCAAG**gtgctcacacagtttatccacaacacctgagtgg
atgacatcaagatagtcgttgtatctccttcgtactctcatgaccacggaaag**ATGATCAAG**
agaggatgattcttggccatatcgcaatgaataacttgtgacttgcattccaattgacatcttcagc
gccatattgcgtggccaaggtgacggagcgggattacgaaagcatgatcatggctgttctgtt
atcttggggactgagacttgttaggatagacggttttcatcactgacttagccaaagccttactct
gcctgacatcgaccgtaaattgataatgaatttacatgcttccgcacgatttacct**CTTGATCAT**cg
atccgattgaagatcttcaattgttaattcttgcctcgactcatgatgatgagct**CTTGATCA**
Tgtttccttaaccctctattttacggaaga**ATGATCAAG**ctgctgct**CTTGATCAT**cgttcc

ATGATCAAG

||||||| are **reverse complements** and likely **DnaA** boxes
TACTAGTTC (**DnaA** does not care what strand to bind to)

It is **VERY SURPRISING** to find a 9-mer appearing **6 or more** times
(counting reverse complements) within a short ≈ 500 nucleotides.

Complementary strand on a template strand

- The figure below shows a template strand **AGTCGCGATAGT** and its complementary strand **ACTATGCGACT**. The beginning and end of a DNA strand are denoted 5' (pronounced “five prime”) and 3' (pronoun

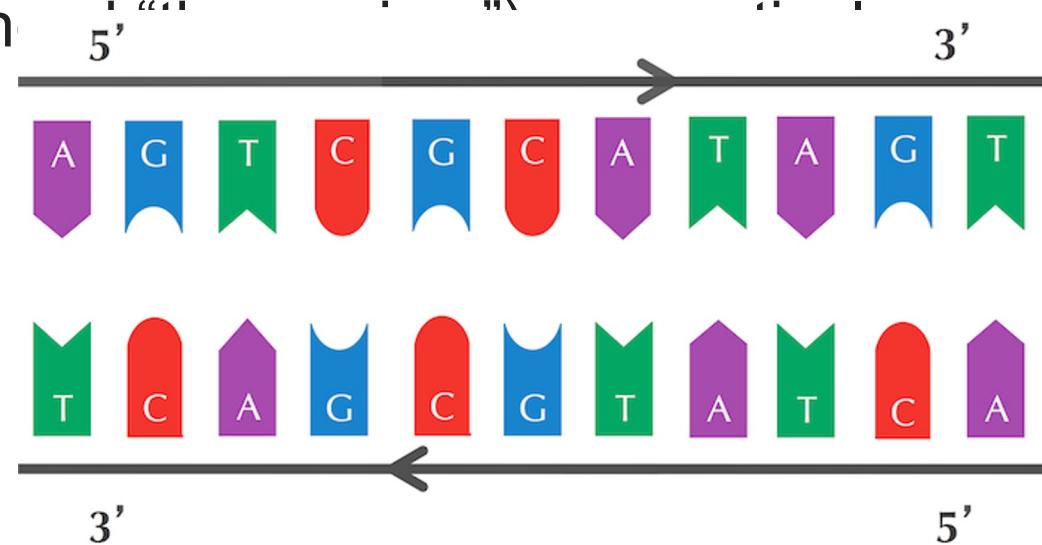
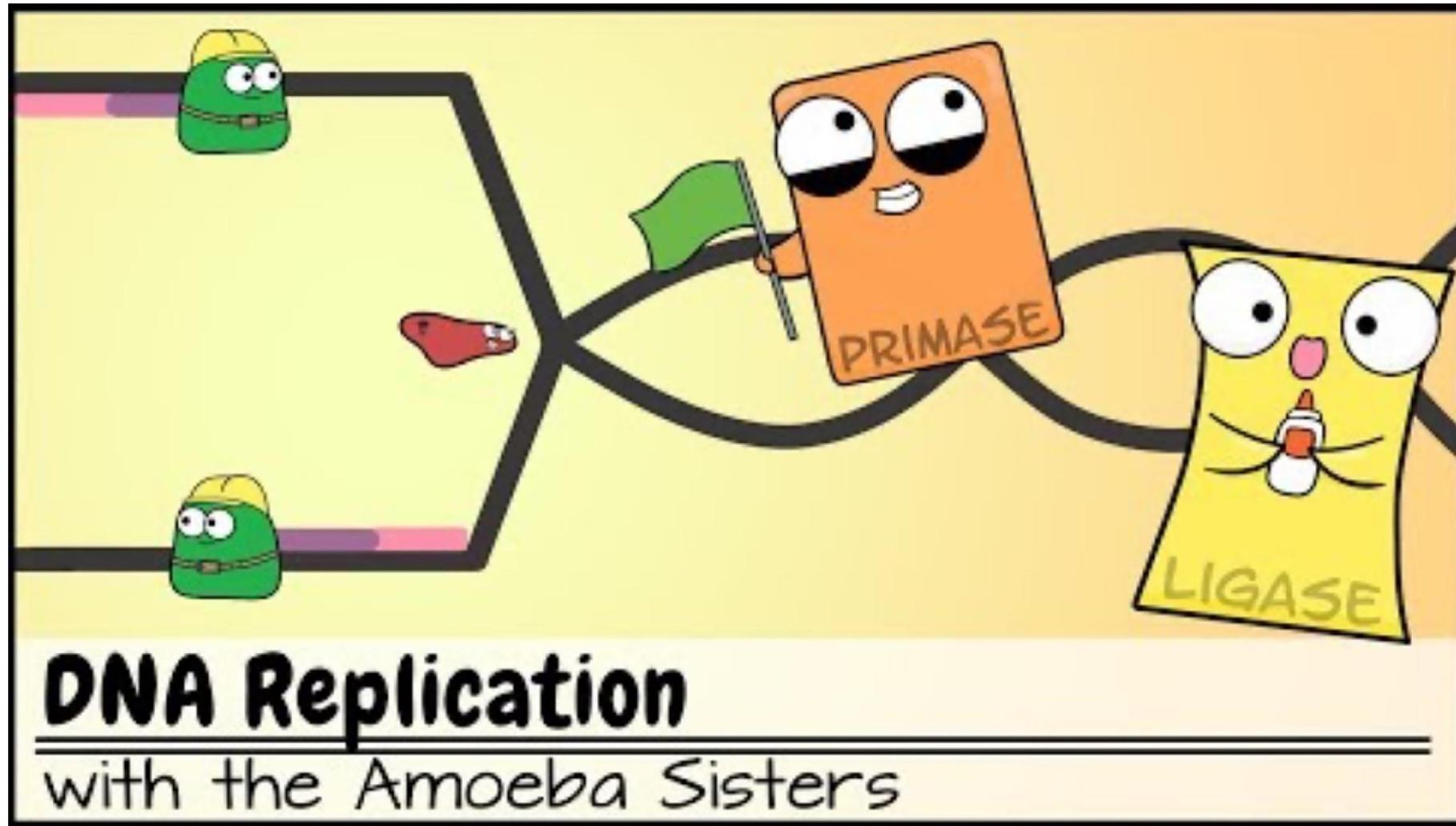


Figure: Complementary strands run in opposite directions. Each strand is read in the $5' \rightarrow 3'$ direction.



Code Challenge: Solve the Reverse Complement Problem.



Reverse Complement Problem: Find the reverse complement of a DNA string.

Input: A DNA string Pattern.

Output: Pattern_{rc} , the reverse complement of Pattern.



Code Challenge: Solve the Reverse Complement Problem.



Reverse Complement Problem: Find the reverse complement of a DNA string.

Input: A DNA string Pattern.

Output: Pattern_{rc} , the reverse complement of Pattern.

```
1 import sys
2
3 # Please do not remove package declarations because these are used by the autograder.
4
5 # Insert your reverse_complement function here, along with any subroutines you need
6 def reverse_complement(Text: str) -> str:
7     """Calculate the reverse complement of a DNA pattern."""
8     complement = []
9     for i in Text:
10         if i == "A":
11             complement.append("T")
12         elif i == "T":
13             complement.append("A")
14         elif i == "G":
15             complement.append("C")
16         elif i == "C":
17             complement.append("G")
18         else:
19             print("Not ACTG!")
20
21     reverse_complement = list(reversed(complement))
22
23     return "".join(reverse_complement)
24
```



Exercise Break



Return a space-separated list of starting positions (in increasing order) where **CTTGATCAT** appears as a substring in the *Vibrio cholerae* genome.





Can we Now Find Hidden Messages in *Thermotoga petrophila*?

```
aactctatacctcctttgtcgatttgtgattatagagaaaatcttattaactgaaactaa  
aatggtagggttggtaggtttgtgtacatttgtagtatctgatttttaattacataccgta  
tattgtattaaattgacgaacaattgcatgaaattgaatatatgcacaaacaaacctaccaccaaac  
tctgtattgaccatTTtaggacaacttcagggtggtaggtttctgaagctctcatcaatagactat  
tttagtcttacaaacaatattaccgttcagattcaagattctacaacgctgtttatggcggt  
gcagaaaaacttaccacccatccagtatccaagccgatttcagagaaacctaccacttac  
cacttacctaccaccGGGTggtaagtgcagacattattaaaaacctcatcagaagcttgc  
aaatttcaataactcgaaacctaccacctgcgtcccattatttactactaataatagcagta  
taattgatctgaaaagaggtggtaaaaaaa
```

No single occurrence of **ATGATCAAG** or **CTTGATCAT** from
Vibrio Cholerae!!!

Applying the Frequent Words Problem to this replication origin:
AACCTACCA, ACCTACCAC, GGTAGGTTT, TGGTAGGTT,
AAACCTACC, CCTACCACC

Hidden Messages in *Thermotoga petrophila*



أكاديمية كاوهست
KAUST ACADEMY

```
aactctatacctccctttgtcgattgtgatggatagaaaaatcttattaactgaaactaa  
aatggtaggtttGGTGGTAGGtttgtgtacatttgtagtatctgatttttaattacataccgt  
tattgtattaaattgacgaacaattgcattggaaattgaatataatgcacaaaacaaaCCTACCAACCaaac  
tctgtattgaccattttaggacaacttcagGGTGGTAGGtttctgaagctctcatcaatagactat  
tttagtcttacaaacaatattaccgttcagattcaagattctacaacgctgtttatggcggt  
gcagaaaaacttaccacctaataccagtatccaaggccattcagagaaaacctaccacttac  
cacttaCCTACCAACCcgggtggtaagttgcagacattattaaaaacctcatcagaagcttgc  
aaattcaataactcgaaaCCTACCAACCtgcgtcccattatttactactaataatagcagta  
taattgatctgaaaagaggtggtaaaaaaa
```

Ori-Finder software confirms that

CCTACCAACC

||||||| are candidate hidden messages.

GGATGGTGG

We learned how to find hidden messages **IF oriC is given**. But we have no clue **WHERE oriC** is located in a (long) genome.



Outline

- **Search for Hidden Messages in Replication Origin**
 - What is a Hidden Message in Replication Origin?
 - Some Hidden Messages are More Surprising than Others
 - **Clumps of Hidden Messages**
- **From a Biological Insight toward an Algorithm for Finding Replication Origin**
 - Asymmetry of Replication
 - Why would a computer scientist care about assymetry of replication?
 - Skew Diagrams
 - Finding Frequent Words with Mismatches
 - Open Problems



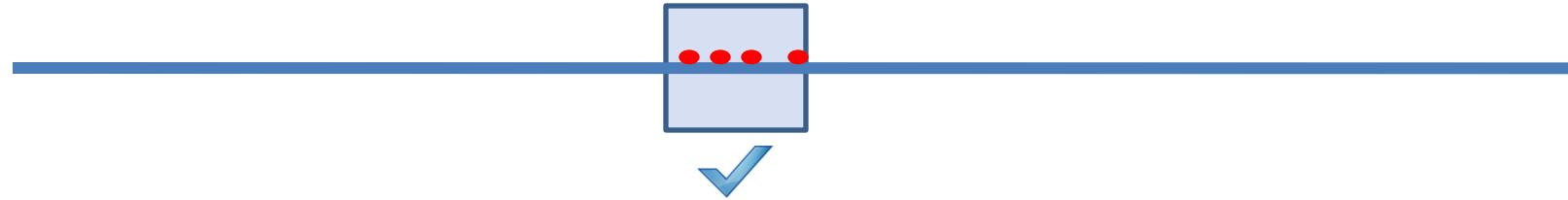
Finding Replication Origin



أكاديمية كاوهست
KAUST ACADEMY

Our strategy **BEFORE**: given a previously **known** *oriC* (a 500-nucleotide window), find **frequent words** (clumps) in *oriC* as candidate *DnaA* boxes.

replication origin → **frequent words**



Finding Replication Origin

Our strategy **BEFORE**: given previously **known** *oriC* (a 500-nucleotide window), find **frequent words** (clumps) in *oriC* as candidate *DnaA* boxes.

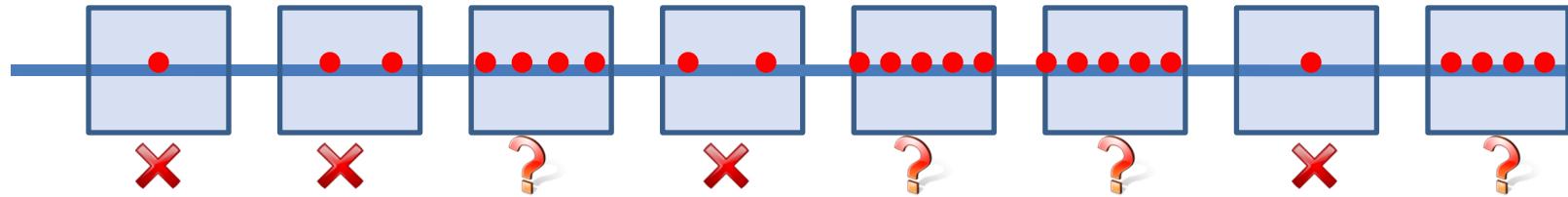
replication origin → **frequent words**

But what if the position of the replication origin within a genome is **unknown!**

Finding Replication Origin

Our strategy **BEFORE**: given previously **known** *oriC* (a 500-nucleotide window), find **frequent words** (clumps) in *oriC* as candidate *DnaA* boxes.

replication origin → **frequent words**



NEW strategy: find frequent words in **ALL** windows within a genome. Windows with **clumps** of frequent words are candidate replication origins.

frequent words → **replication origin**

What is a Clump?

Formal: A k -mer forms an (L, t) -clump inside $Genome$ if there is a **short** (length L) interval of $Genome$ in which it appears **many** (at least t) times.

Clump Finding Problem. Find patterns forming clumps in a string.

- **Input.** A string $Genome$ and integers k (length of a pattern), L (window length), and t (number of patterns in a clump).
- **Output.** All k -mers forming (L, t) -clumps in $Genome$.

There exist **1904 different** 9-mers forming $(500, 3)$ -clumps in $E. coli$ genome. It is absolutely unclear which of them point to the replication origin...

The Clump Finding Problem

```
FindClumps(Text, k, L, t)
    Patterns ← an array of strings of length 0
    n ← |Text|
    for every integer i between 0 and n - L
        Window ← Text(i, L)
        freqMap ← FrequencyTable(Window, k)
        for every key s in freqMap
            if freqMap[s] ≥ t
                append s to Patterns
    remove duplicates from Patterns
    return Patterns
```



Code Challenge:

Solve the Clump Finding Problem



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*.



Exercise Break



How many *different* 9-mers form (500,3)-clumps in the *E. coli* genome? (In other words, do not count a 9-mer more than once.)

