





### Team Members

Chang Heen Sunn U1920383H

Luo Wenyu U1922009G

Ooi Wei Chern U1920504J

Saw William Joseph U1921246F

Tan Ching Fhen U1920787D

### × 🗆 🗕

# String -Matching

Matching genomic pattern with COVID 19 genome











## Presentation Flow

- → Knuth-Morris-Pratt (KMP) Algorithm
- → Boyer-Moore-Horspool (BMH) Algorithm
- → Comparison and Conclusion





#### Knuth-Morris-Pratt

- String-searching algorithm by the means of repeated prefix which is also a suffix of the pattern





#### Pattern: AAGGAAAAG

Prefix ->	Suffix <-
A	G
AA	AG
AAG	AAG
AAGGAAAA	AGGAAAAG



#### Pi\_table / LPS array

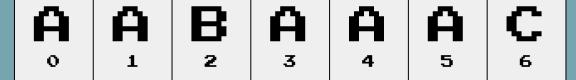
Longest proper

Prefix which is also a
Suffix array



```
def \pi_table(self, pattern):
   n = len(pattern)
   lps = [0]
   for i in range(1, n):
        j = lps[i-1]
        while j>0 and pattern[j] != pattern[i]:
            j = lps[j-1]
        if pattern[j] == pattern[i]:
            lps.append(j+1)
        else:
            lps.append(j)
    return lps
```

#### X ANALYSIS



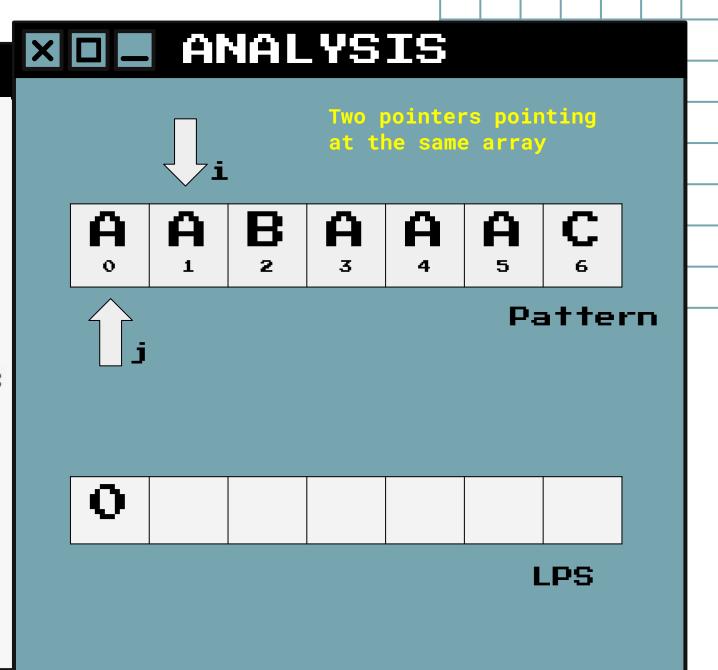
Pattern



LPS

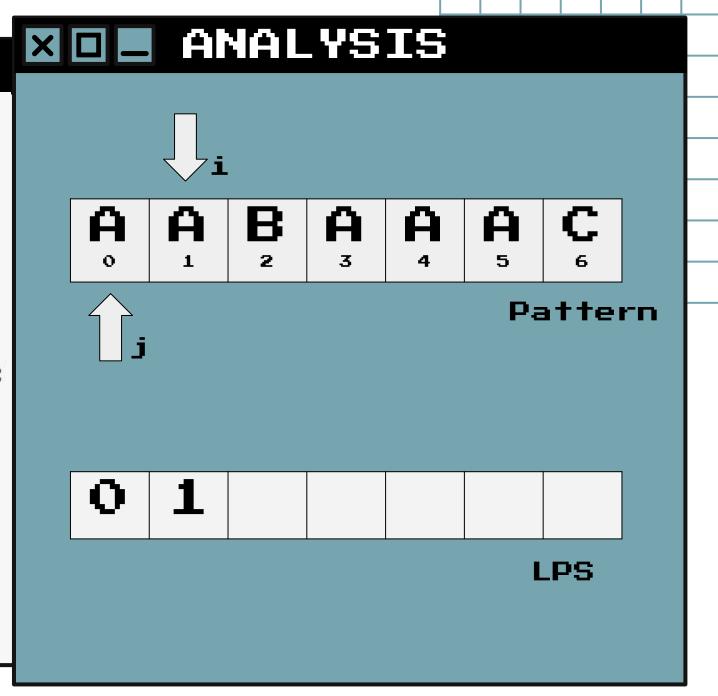
#### 

```
def \pi table(self, pattern):
   n = len(pattern)
   lps = [0]
    for i in range(1, n):
        j = lps[i-1]
        while j>0 and pattern[j] != pattern[i]:
            j = lps[j-1]
        if pattern[j] == pattern[i]:
            lps.append(j+1)
        else:
            lps.append(j)
    return lps
```

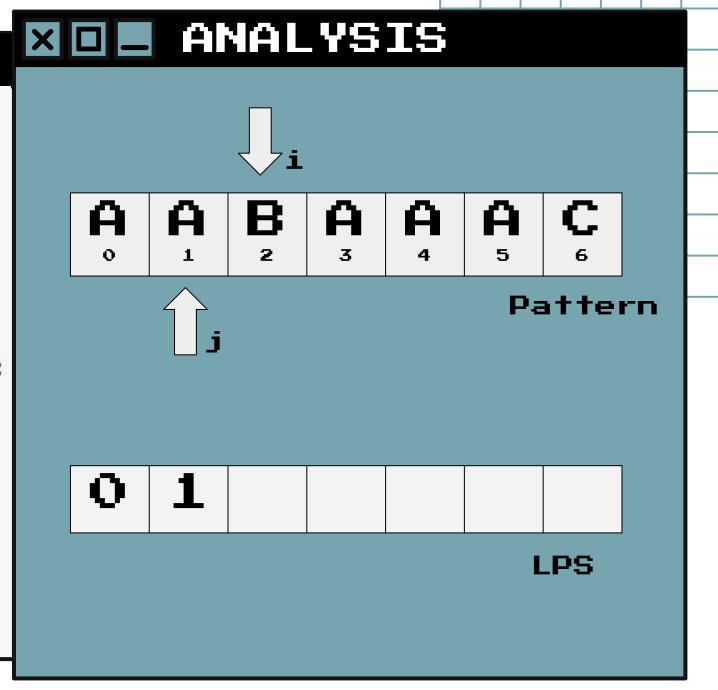


#### X D E CODES

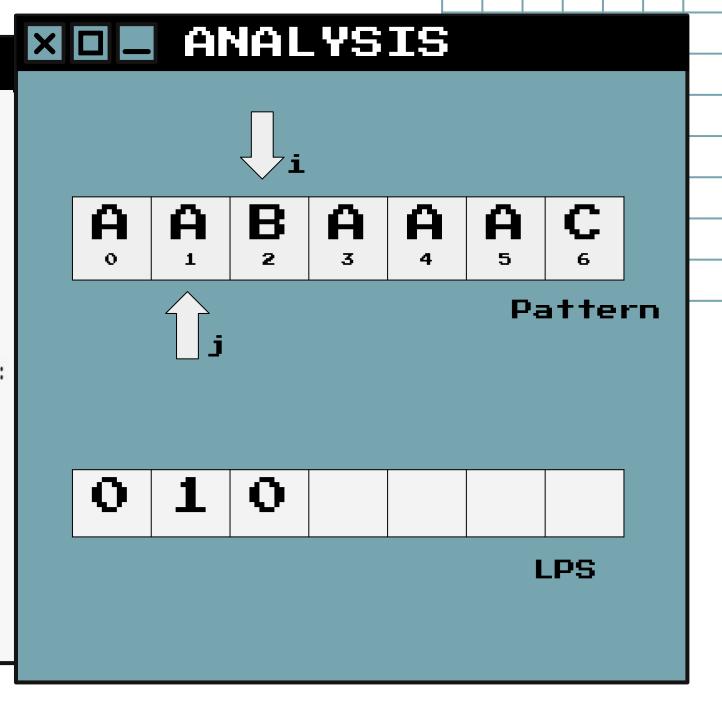
```
def \pi_table(self, pattern):
    n = len(pattern)
    lps = [0]
    for i in range(1, n):
        j = lps[i-1]
        while j>0 and pattern[j] != pattern[i]:
            j = lps[j-1]
        if pattern[j] == pattern[i]:
            lps.append(j+1)
        else:
            lps.append(j)
    return lps
```



```
def \pi_table(self, pattern):
    n = len(pattern)
    lps = [0]
    for i in range(1, n):
        j = lps[i-1]
        while j>0 and pattern[j] != pattern[i]:
            j = lps[j-1]
        if pattern[j] == pattern[i]:
            lps.append(j+1)
        else:
            lps.append(j)
    return lps
```

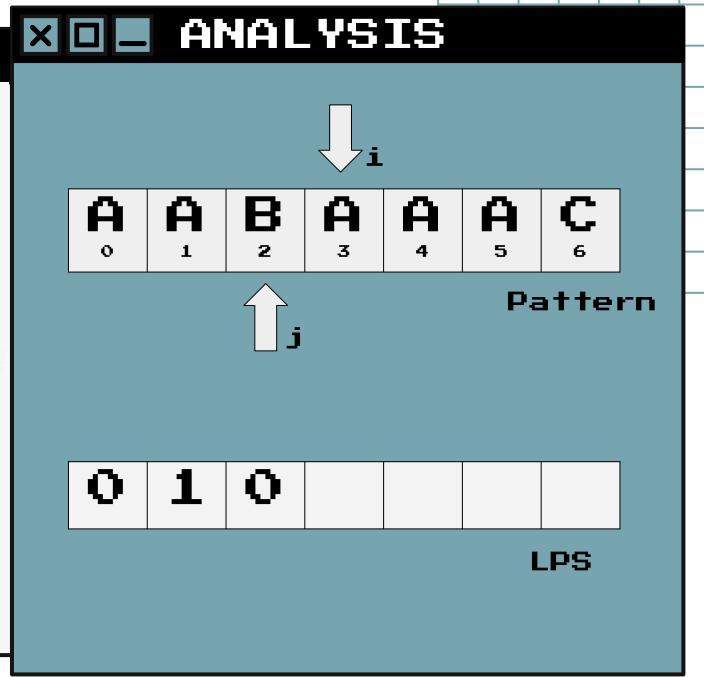


```
def \pi_table(self, pattern):
    n = len(pattern)
    lps = [0]
    for i in range(1, n):
        j = lps[i-1]
        while j>0 and pattern[j] != pattern[i]:
            j = lps[j-1]
        if pattern[j] == pattern[i]:
            lps.append(j+1)
        else:
            lps.append(j)
    return lps
```



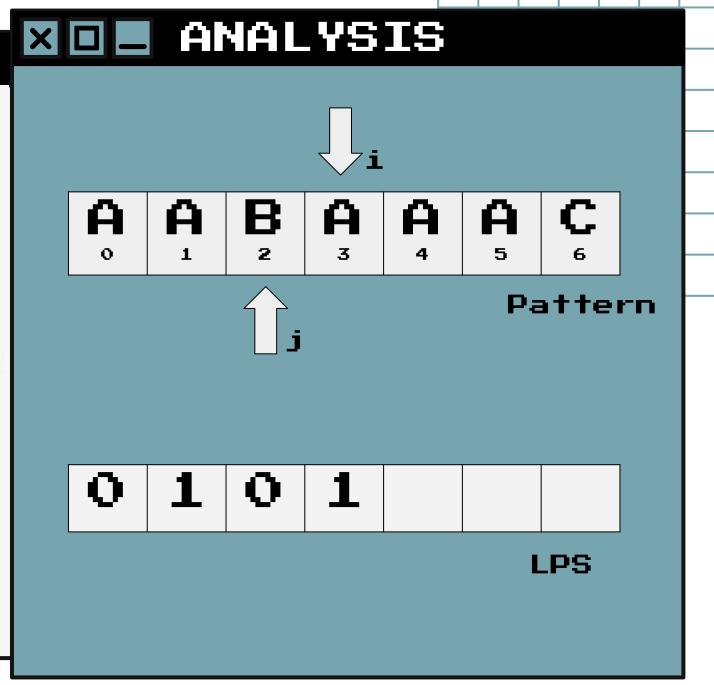
#### XDE CODES

```
def \pi_table(self, pattern):
    n = len(pattern)
    lps = [0]
    for i in range(1, n):
        j = lps[i-1]
        while j>0 and pattern[j] != pattern[i]:
            j = lps[j-1]
        if pattern[j] == pattern[i]:
            lps.append(j+1)
        else:
            lps.append(j)
    return lps
```

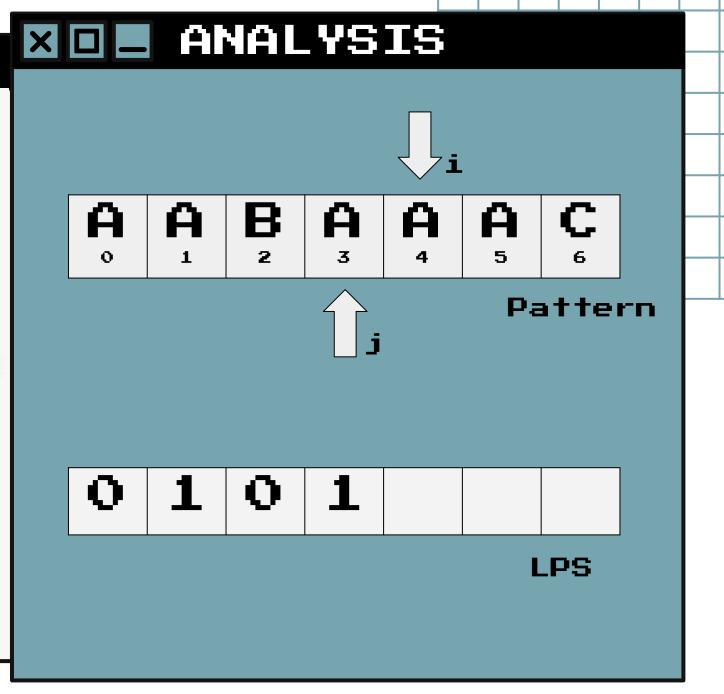


#### XDE CODES

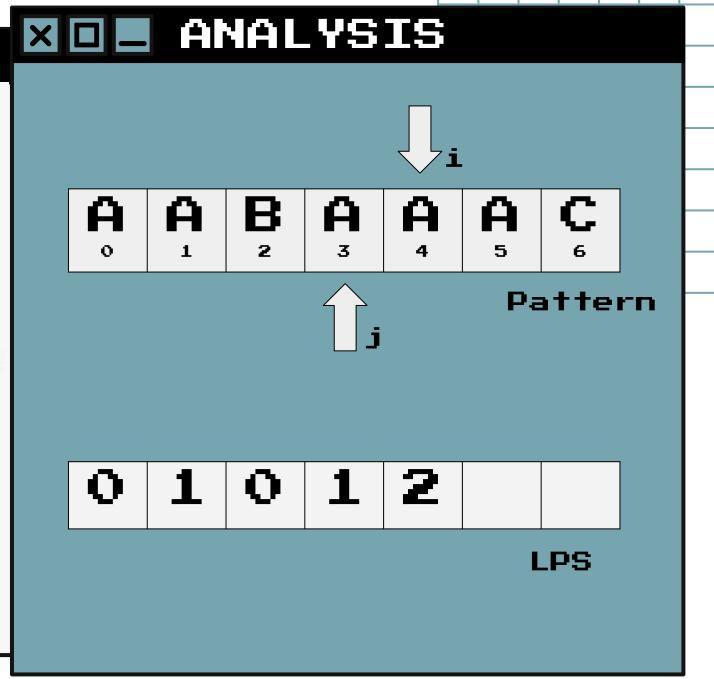
```
def \pi_table(self, pattern):
    n = len(pattern)
    lps = [0]
    for i in range(1, n):
        j = lps[i-1]
        while j>0 and pattern[j] != pattern[i]:
            j = lps[j-1]
        if pattern[j] == pattern[i]:
            lps.append(j+1)
        else:
            lps.append(j)
    return lps
```



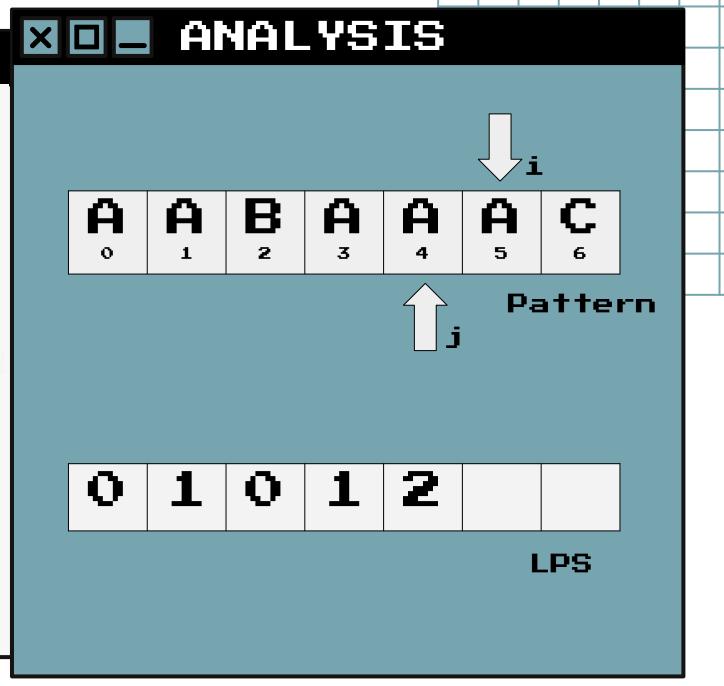
```
def \pi_table(self, pattern):
    n = len(pattern)
    lps = [0]
    for i in range(1, n):
        j = lps[i-1]
        while j>0 and pattern[j] != pattern[i]:
            j = lps[j-1]
        if pattern[j] == pattern[i]:
            lps.append(j+1)
        else:
            lps.append(j)
    return lps
```



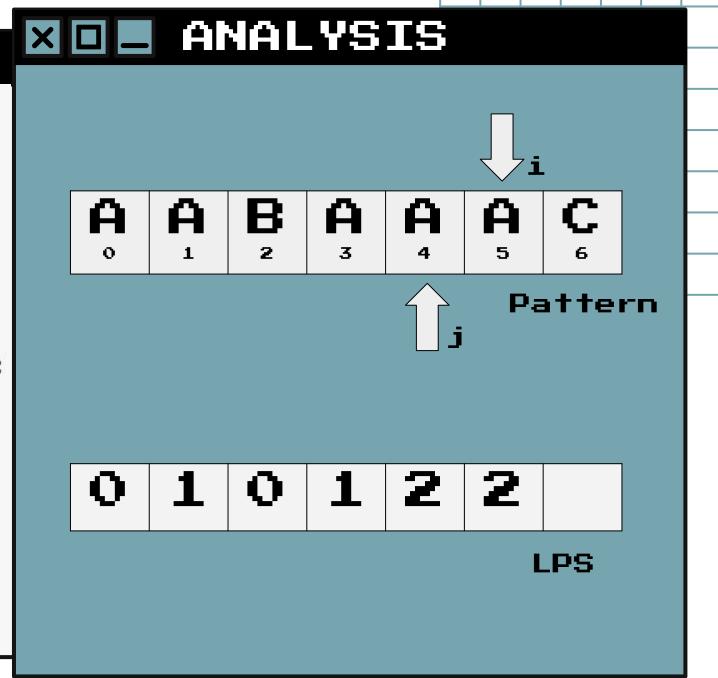
```
def \pi_table(self, pattern):
    n = len(pattern)
    lps = [0]
    for i in range(1, n):
        j = lps[i-1]
        while j>0 and pattern[j] != pattern[i]:
            j = lps[j-1]
        if pattern[j] == pattern[i]:
            lps.append(j+1)
        else:
            lps.append(j)
    return lps
```



```
def \pi_table(self, pattern):
    n = len(pattern)
    lps = [0]
    for i in range(1, n):
        j = lps[i-1]
        while j>0 and pattern[j] != pattern[i]:
            j = lps[j-1]
        if pattern[j] == pattern[i]:
            lps.append(j+1)
        else:
            lps.append(j)
    return lps
```

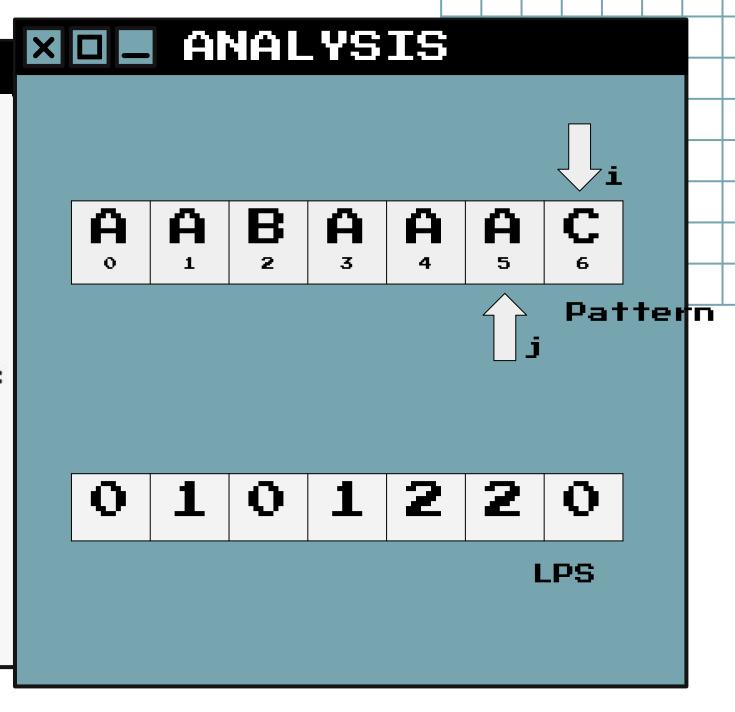


```
def \pi_table(self, pattern):
    n = len(pattern)
    lps = [0]
    for i in range(1, n):
        j = lps[i-1]
        while j>0 and pattern[j] != pattern[i]:
            j = lps[j-1]
        if pattern[j] == pattern[i]:
            lps.append(j+1)
        else:
            lps.append(j)
    return lps
```



#### X D = CODES

```
def \pi_table(self, pattern):
    n = len(pattern)
    lps = [0]
    for i in range(1, n):
        j = lps[i-1]
        while j>0 and pattern[j] != pattern[i]:
            j = lps[j-1]
        if pattern[j] == pattern[i]:
            lps.append(j+1)
        else:
            lps.append(j)
    return lps
```



```
def \pi_table(self, pattern):
   n = len(pattern)
   lps = [0]
    for i in range(1, n):
        j = lps[i-1]
        while j>0 and pattern[j] != pattern[i]:
            j = lps[j-1]
        if pattern[j] == pattern[i]:
            lps.append(j+1)
        else:
            lps.append(j)
    return lps
```

#### X - ANALYSIS



Pattern

0 1 0 1 2 2 0

LPS

#### X D = ANALYSIS

0										
A	В	A	В	A	В	A	В	В	A	С



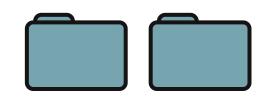
Two pointers pointing at two different arrays



Pat.		A	В	A	В	В	A	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0

#### 

```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
   m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```



#### X D = ANALYSIS

0										
A	В	A	В	A	В	A	В	В	A	С

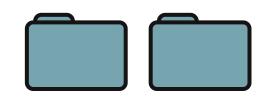


To check if: text[i] == pattern[j+1]



Pat.		A	В	A	В	В	A	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0

```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
   m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```





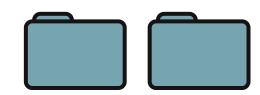
0										
Α	В	A	В	A	В	A	В	В	A	C





Pat.		Α	В	A	В	В	A	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0

```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
   m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```





0										
Α	В	A	В	A	В	A	В	В	A	С

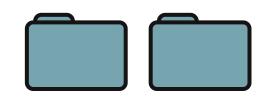




Pat.		Α	В	A	В	В	A	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0

#### XDE CODES

```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
   m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```





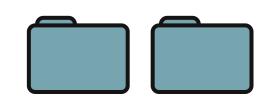
0										
Α	В	A	В	A	В	A	В	В	A	С





Pat.		Α	В	A	В	В	A	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0

```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
   m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```





						6				
Α	В	A	В	A	В	A	В	В	A	C

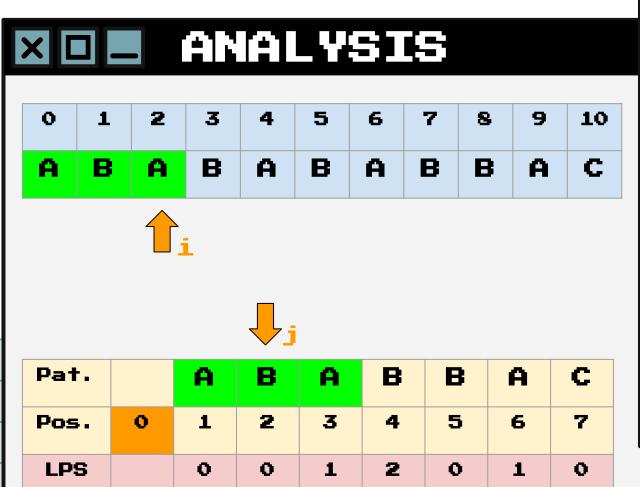




Pat.		Α	В	A	В	В	A	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0

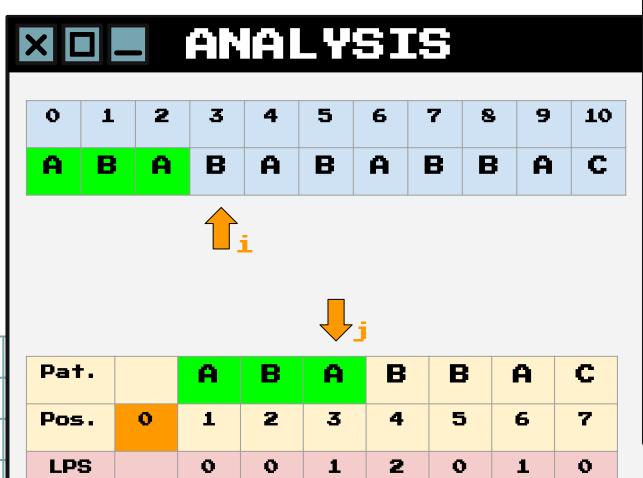
```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
   m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```



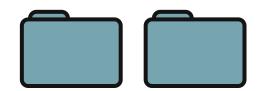


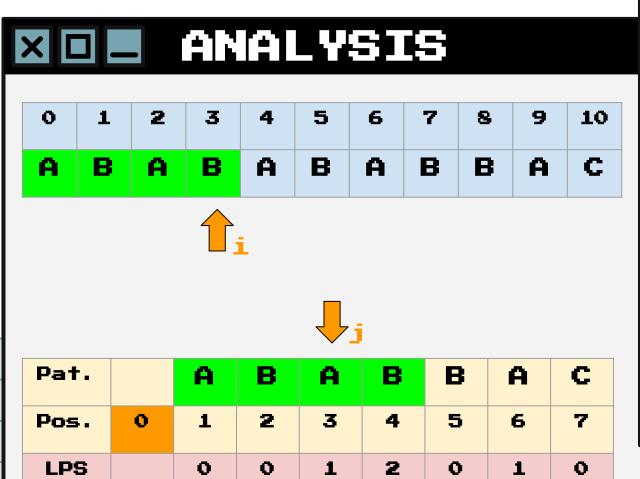
```
X D CODES
def knuth_morris_pratt(self, text, pattern):
   lps = self.\pi_table(pattern)
   result, j = [], 0
   m, n = len(text), len(pattern)
   for i in range(m):
       while j>0 and text[i] != pattern[j]:
           j = lps[j-1]
       if text[i] == pattern[j]:
           j += 1
       if j == n:
           result.append((i-j)+1)
           j = 0
```



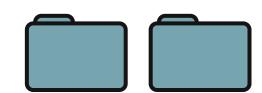


```
X D CODES
def knuth_morris_pratt(self, text, pattern):
   lps = self.\pi_table(pattern)
   result, j = [], 0
   m, n = len(text), len(pattern)
   for i in range(m):
       while j>0 and text[i] != pattern[j]:
           j = lps[j-1]
       if text[i] == pattern[j]:
           j += 1
       if j == n:
           result.append((i-j)+1)
           j = 0
```





```
X D CODES
def knuth_morris_pratt(self, text, pattern):
   lps = self.\pi_table(pattern)
   result, j = [], 0
   m, n = len(text), len(pattern)
   for i in range(m):
       while j>0 and text[i] != pattern[j]:
           j = lps[j-1]
       if text[i] == pattern[j]:
           j += 1
       if j == n:
           result.append((i-j)+1)
           j = 0
```





			3							
Α	В	Α	В	Α	В	A	В	В	A	С

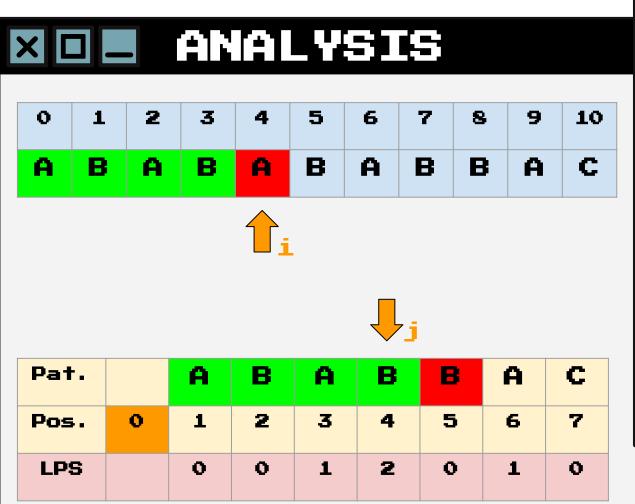




Pat.		Α	В	Α	В	В	A	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0

```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
   m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```





```
X D CODES
def knuth_morris_pratt(self, text, pattern):
   lps = self.\pi_table(pattern)
   result, j = [], 0
   m, n = len(text), len(pattern)
   for i in range(m):
       while j>0 and text[i] != pattern[j]:
           j = lps[j-1]
       if text[i] == pattern[j]:
           j += 1
       if j == n:
           result.append((i-j)+1)
           j = 0
```



0										
A	В	A	В	Α	В	A	В	В	A	С

i moves 2 steps forwards from where it initially was

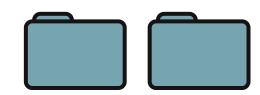
j moves 2 steps backwards



Pat.		Α	В	Α	В	В	A	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0

### 

```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
    m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```





						6				
Α	В	A	В	A	В	A	В	В	A	С



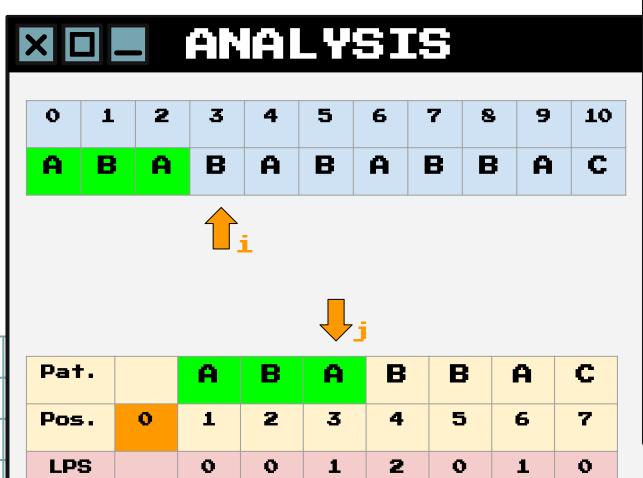


Pat.		Α	В	A	В	В	A	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0

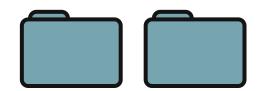
#### 

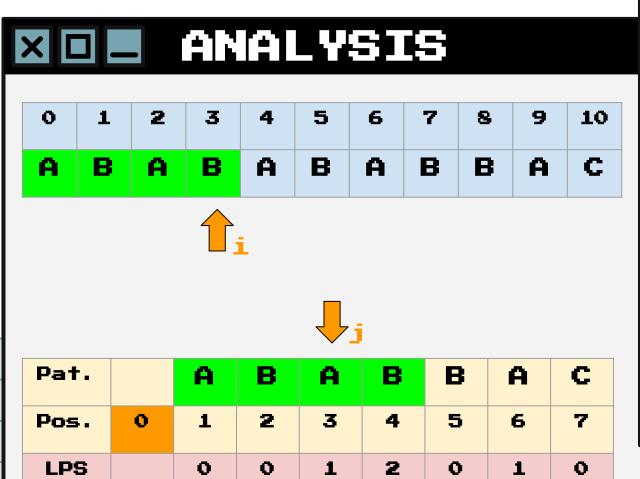
```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
   m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```



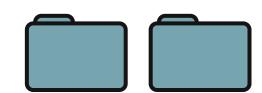


```
X D CODES
def knuth_morris_pratt(self, text, pattern):
   lps = self.\pi_table(pattern)
   result, j = [], 0
   m, n = len(text), len(pattern)
   for i in range(m):
       while j>0 and text[i] != pattern[j]:
           j = lps[j-1]
       if text[i] == pattern[j]:
           j += 1
       if j == n:
           result.append((i-j)+1)
           j = 0
```





```
X D CODES
def knuth_morris_pratt(self, text, pattern):
   lps = self.\pi_table(pattern)
   result, j = [], 0
   m, n = len(text), len(pattern)
   for i in range(m):
       while j>0 and text[i] != pattern[j]:
           j = lps[j-1]
       if text[i] == pattern[j]:
           j += 1
       if j == n:
           result.append((i-j)+1)
           j = 0
```





			3							
Α	В	Α	В	Α	В	A	В	В	A	С

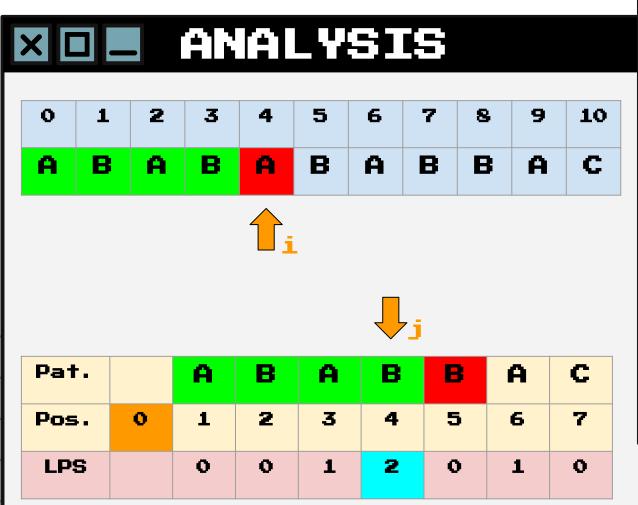




Pat.		Α	В	Α	В	В	A	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0

```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
   m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```





```
X D CODES
def knuth_morris_pratt(self, text, pattern):
   lps = self.\pi_table(pattern)
   result, j = [], 0
   m, n = len(text), len(pattern)
   for i in range(m):
       while j>0 and text[i] != pattern[j]:
           j = lps[j-1]
       if text[i] == pattern[j]:
           j += 1
       if j == n:
           result.append((i-j)+1)
           j = 0
```





				4						
A	В	Α	В	A	В	A	В	В	A	C



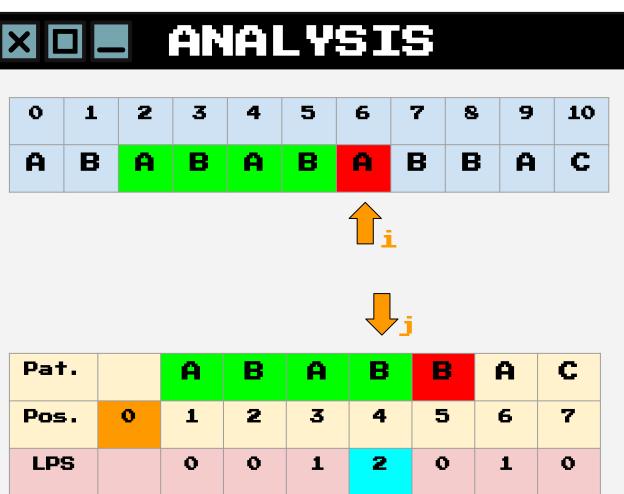


Pat.		Α	В	A	В	В	A	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0

## XDE CODES

```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
   m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```





```
X D CODES
def knuth_morris_pratt(self, text, pattern):
   lps = self.\pi_table(pattern)
   result, j = [], 0
   m, n = len(text), len(pattern)
   for i in range(m):
       while j>0 and text[i] != pattern[j]:
           j = lps[j-1]
       if text[i] == pattern[j]:
           j += 1
       if j == n:
           result.append((i-j)+1)
           j = 0
```





				4						
A	В	A	В	Α	В	A	В	В	A	С





Pat.		Α	В	A	В	В	A	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0

## 

```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
   m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```





					5					
Α	В	A	В	Α	В	Α	В	В	A	С

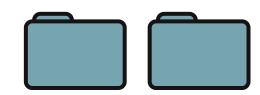




Pat.		Α	В	Α	В	В	A	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0



```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
   m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```





0										
A	В	A	В	Α	В	А	В	В	Α	С

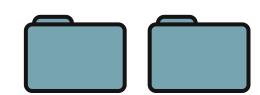




Pat.		Α	В	Α	В	В	Α	С
Pos.	0	1	2	3	4	5	6	7
LPS		0	0	1	2	0	1	0

### XDE CODES

```
def knuth_morris_pratt(self, text, pattern):
    lps = self.\pi_table(pattern)
    result, j = [], 0
   m, n = len(text), len(pattern)
    for i in range(m):
        while j>0 and text[i] != pattern[j]:
            j = lps[j-1]
        if text[i] == pattern[j]:
            j += 1
        if j == n:
            result.append((i-j)+1)
            j = 0
```



### X D TIME COMPLEXITY

#### **Building LPS array:**

time complexity = O(m)

#### Searching:

comparison is a multiple of n, time complexity = O(n)

Total time complexity: O(n+m)



- lps array is <u>NOT actually</u> linear
- the pattern is usually short

Overall (searching)=> O(n+m)

$$=> 0(n)$$

for n>>m





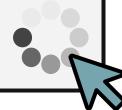
### **▼□■ TIME COMPLEXITY**

Best case: each time it shifts a suffix could be found

Virus genome: A C T G G A C T G G A C T ...

Target genome: A C T G G A C G

Best case time complexity = O(n)





### **▼□■** TIME COMPLEXITY

Worst case: All the characters in the pattern are distinct. BUT...

pat	Т	A	С	G
Pos	0	1	2	3
lps	0	0	0	0

Worst case time complexity = O(n)



 Search genome from back to front

```
current index=len(self.genome)-len(target)
while(current index>0):
   flag=1
    for i in range(len(target)):
        genome character=self.genome[current index+i]
       target character=target[i]
        if genome character!=target character:
            flag=0
            break
    if flag==1:
       results.append(current index)
    skip=skipper.get(self.genome[current index],len(target))
    current index-=skip
```

- Search genome from back to front
- Check pattern is from front to back

```
current index=len(self.genome)-len(target)
while(current index>0):
    flag=1
    for i in range(len(target)):
        genome character=self.genome[current index+i]
        target character=target[i]
        if genome character!=target character:
            flag=0
            break
    if flag==1:
        results.append(current index)
    skip=skipper.get(self.genome[current index],len(target))
    current index-=skip
```

- Search genome from back to front
- Check pattern is from front to back
- Use 'bad match table'

```
E.g. ('A':3,'T':2)
```

```
results=[]
skipper={c:len(target)-i-1 for i,c in enumerate(target[:0:-1])}
current index=len(self.genome)-len(target)
while(current index>0):
   flag=1
   for i in range(len(target)):
        genome character=self.genome[current_index+i]
       target character=target[i]
        if genome character!=target character:
           flag=0
            break
   if flag==1:
        results.append(current index)
    skip=skipper.get(self.genome[current_index],len(target))
    current index-=skip
```

- Search genome from back to front
- Check pattern is from front to back
- Use 'bad match table'

```
E.g. {'A':3,'T':2}
```

 Skips by len(pattern) otherwise

```
results=[]
skipper={c:len(target)-i-1 for i,c in enumerate(target[:0:-1])}
current index=len(self.genome)-len(target)
while(current index>0):
    flag=1
    for i in range(len(target)):
        genome character=self.genome[current index+i]
        target character=target[i]
        if genome character!=target character:
            flag=0
            break
    if flag==1:
        results.append(current index)
    skip=skipper.get(self.genome[current_index],len(target))
    current index-=skip
```

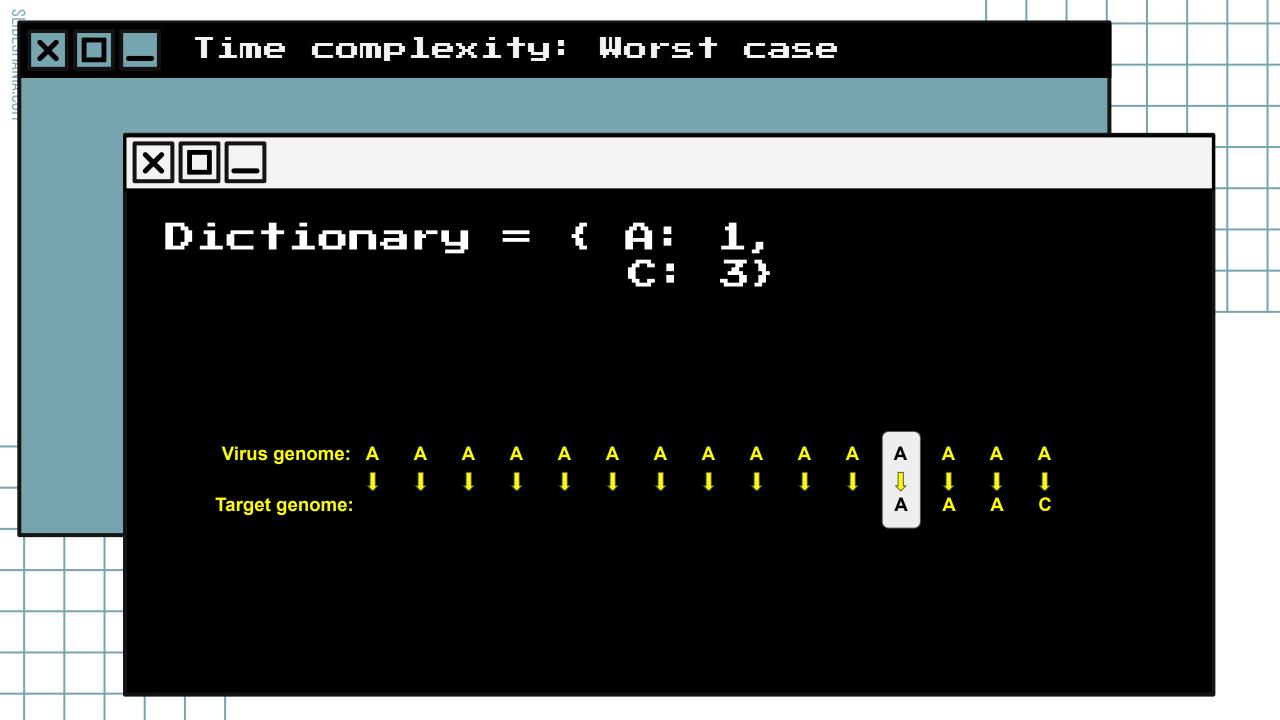


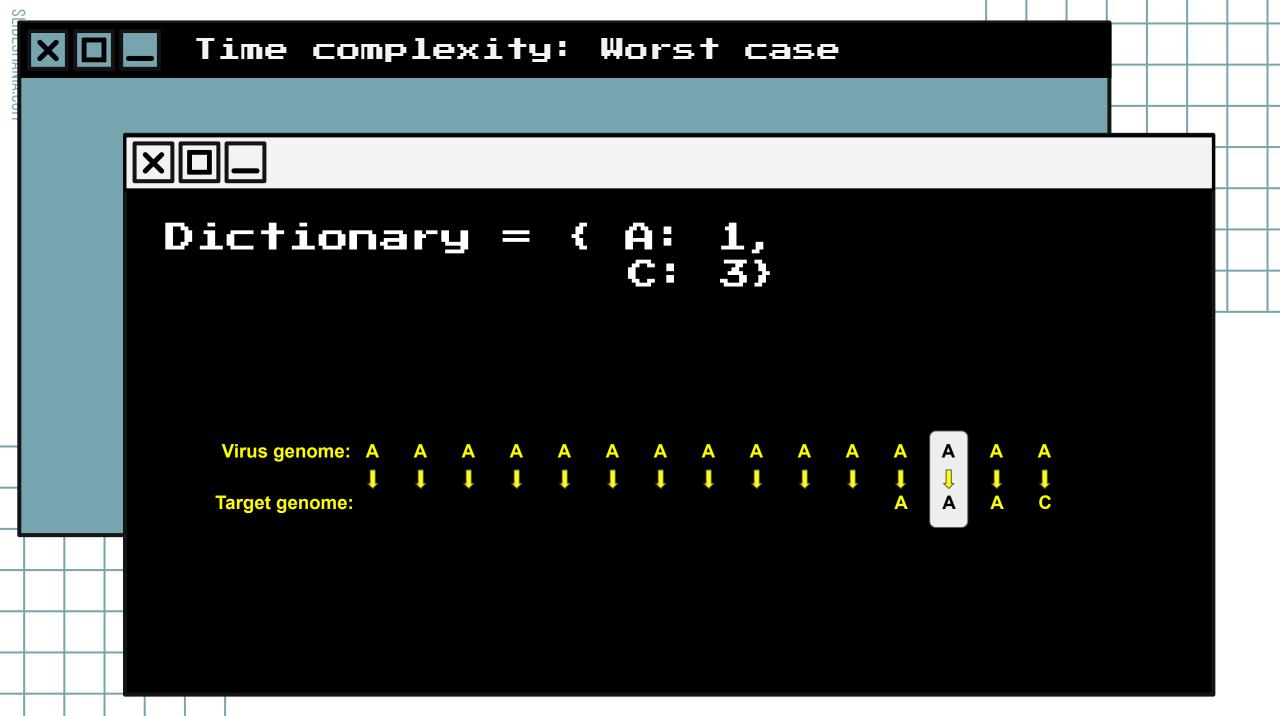
# How the Bad Match Table is made?

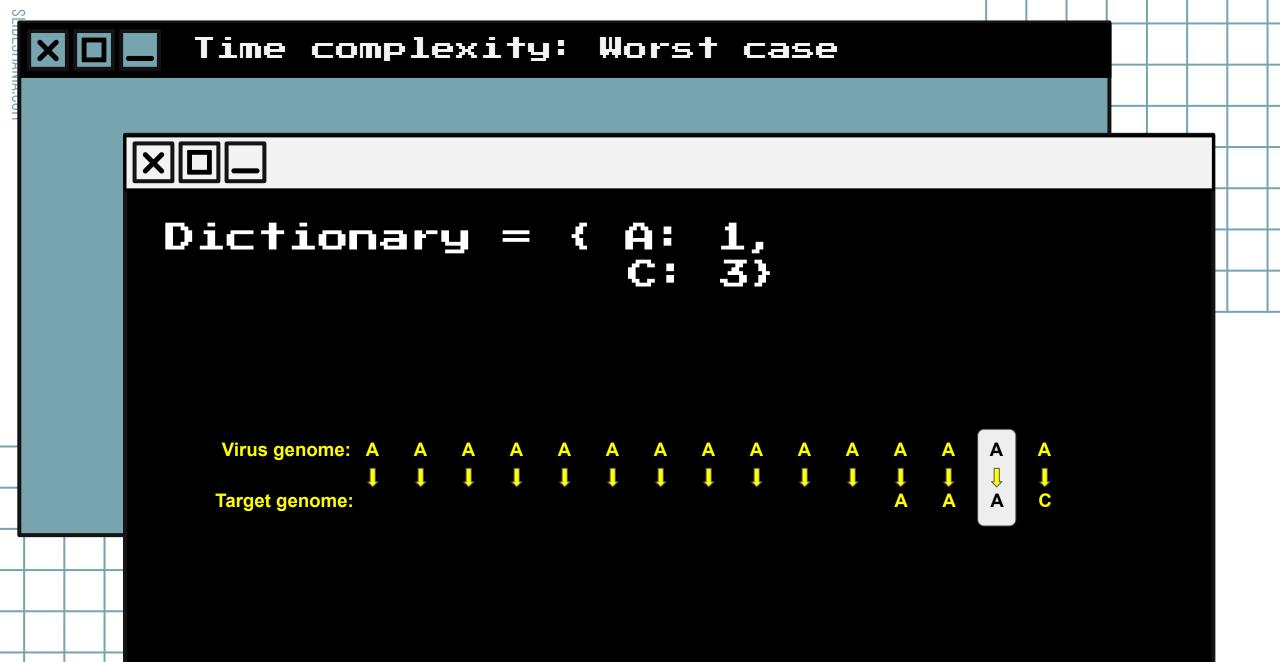


- Ignore first character
- Based on first occurrence of a character

```
E.g. TAAAG
0 1 2 3 4
('A':1, 'G':4)
```



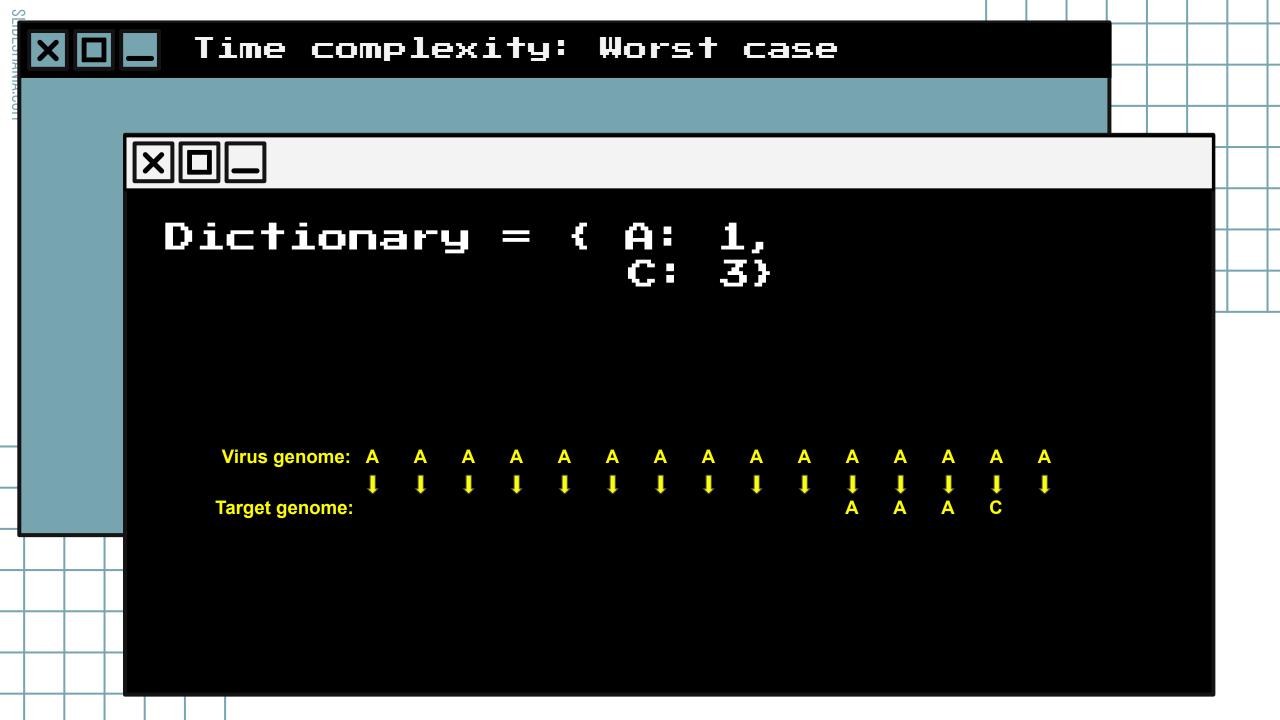


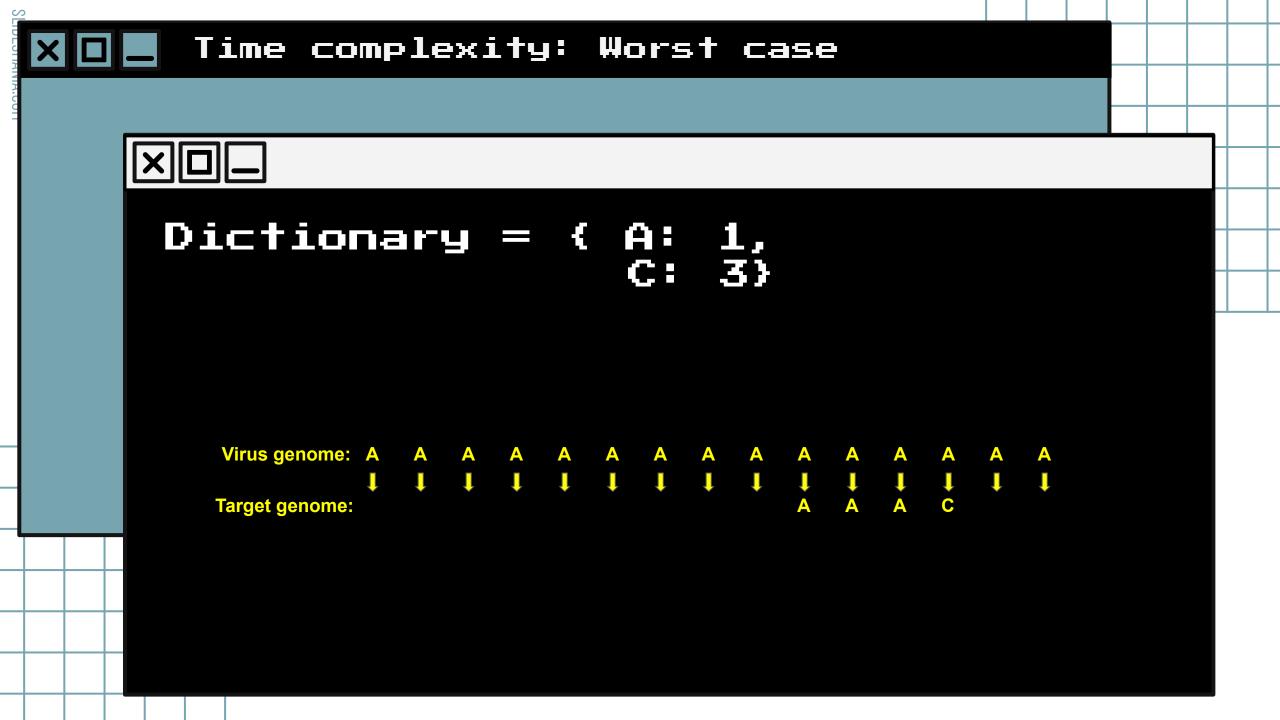


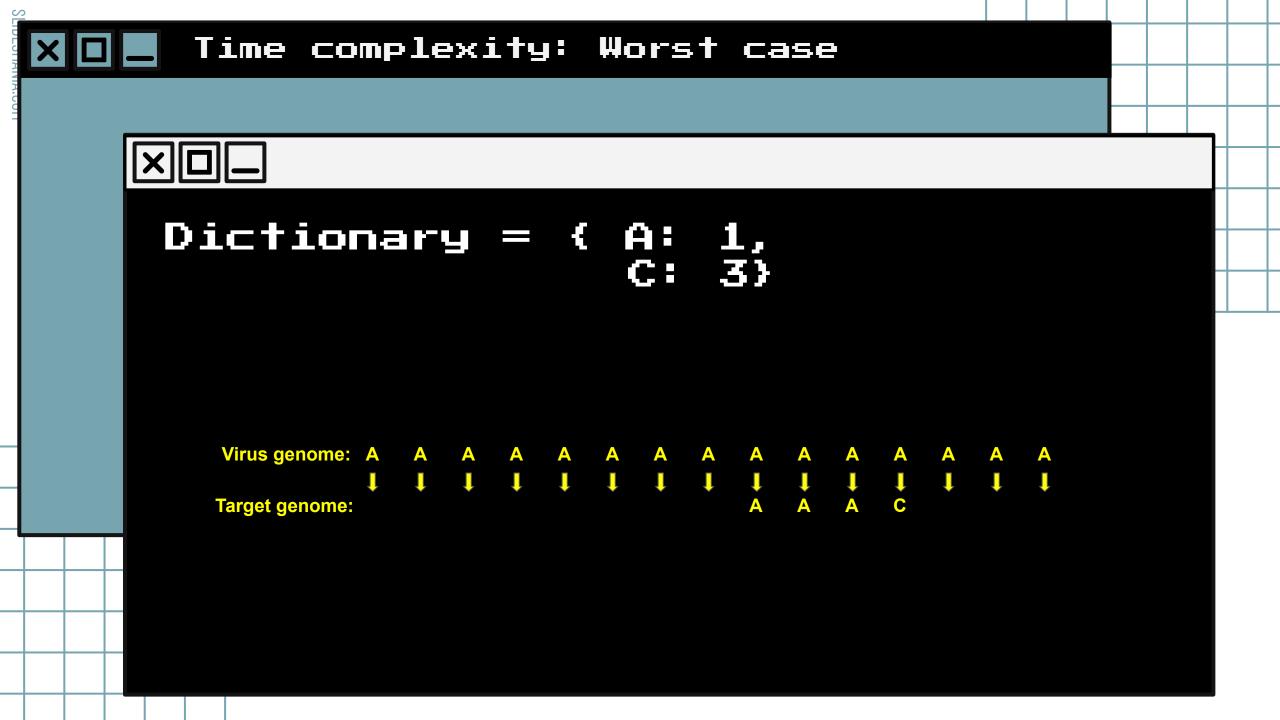


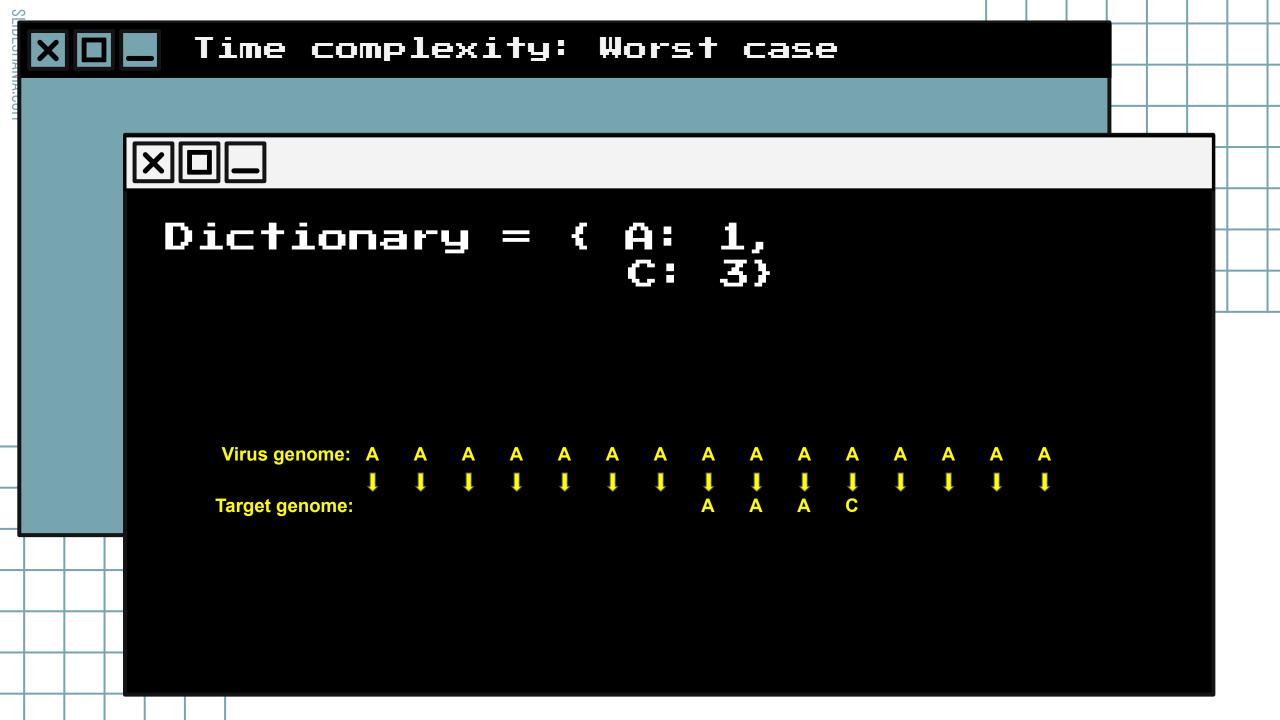


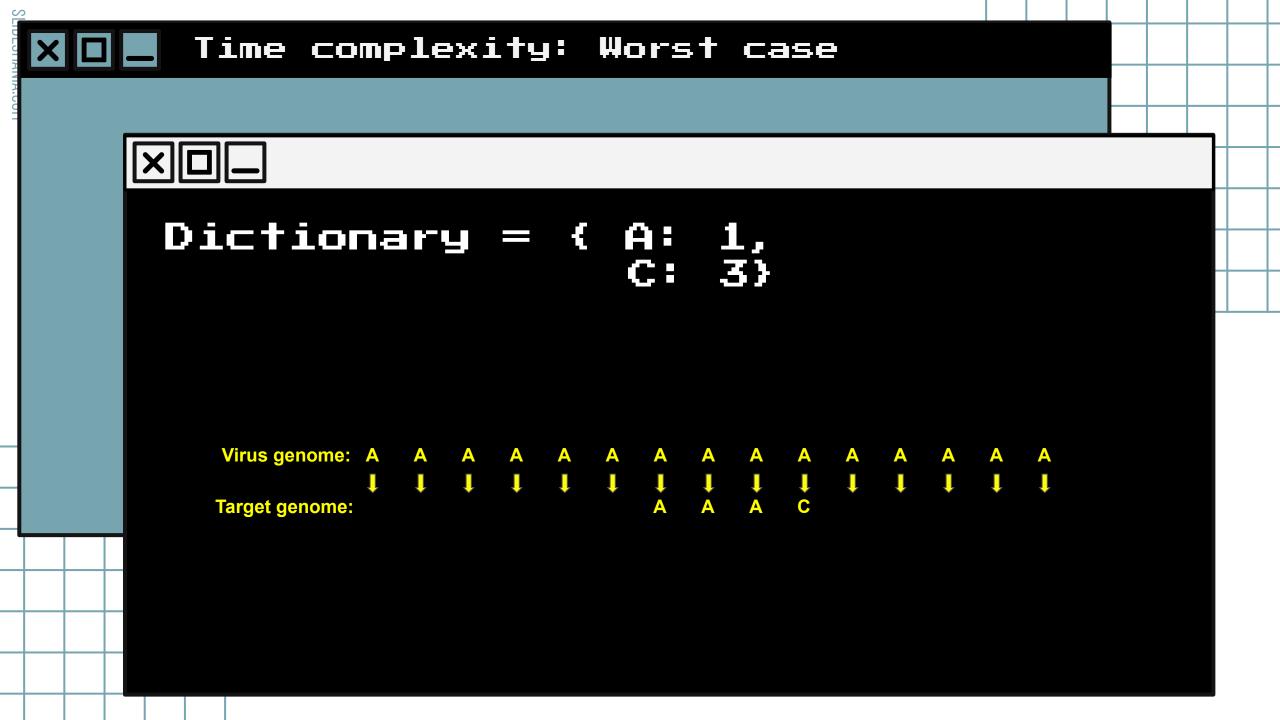
MISMATCH!!

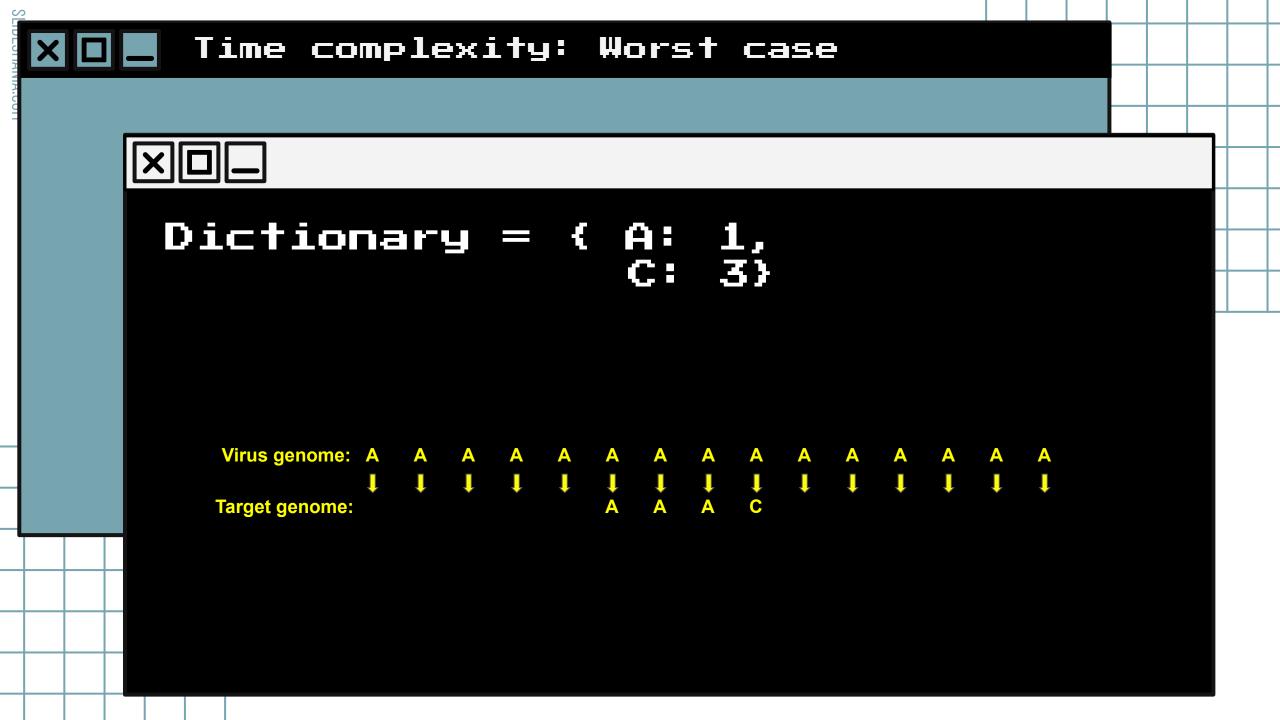


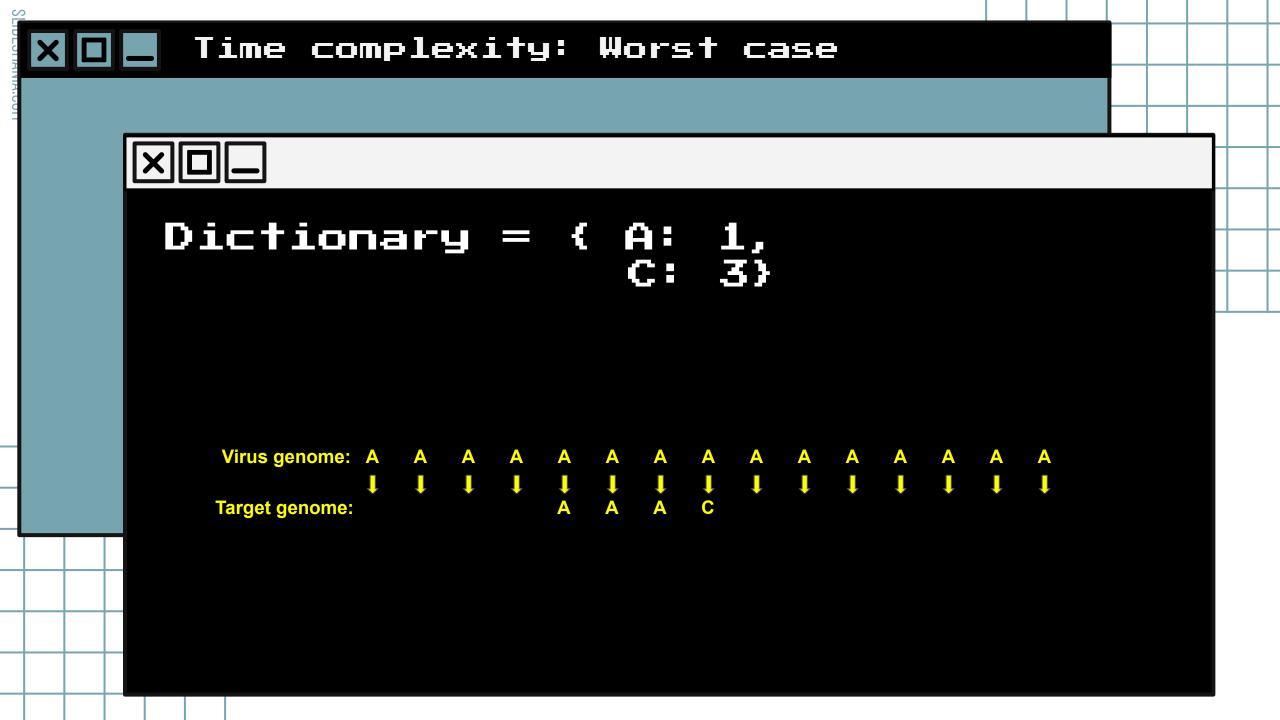


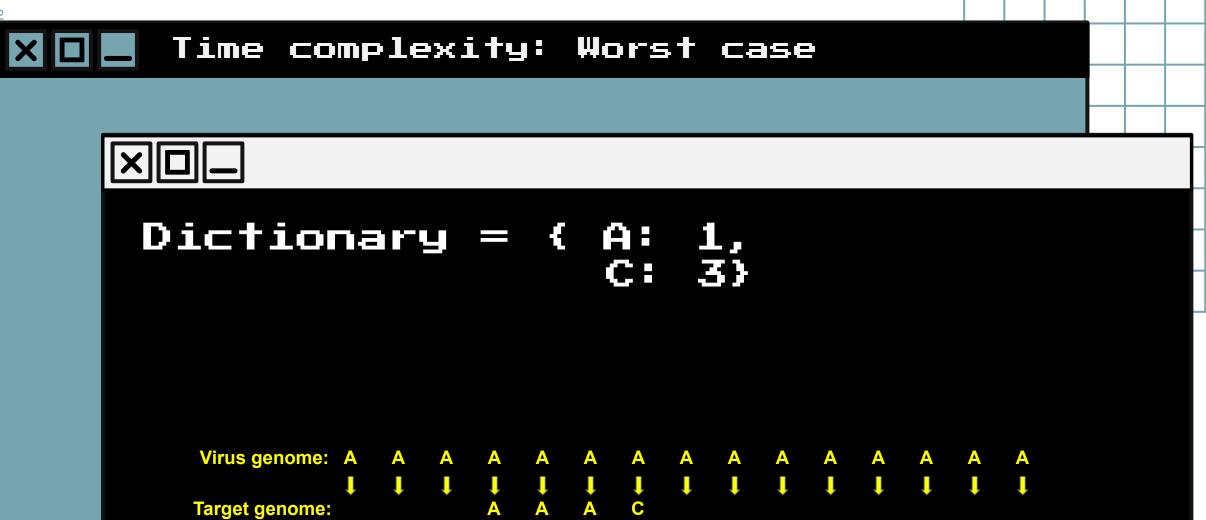












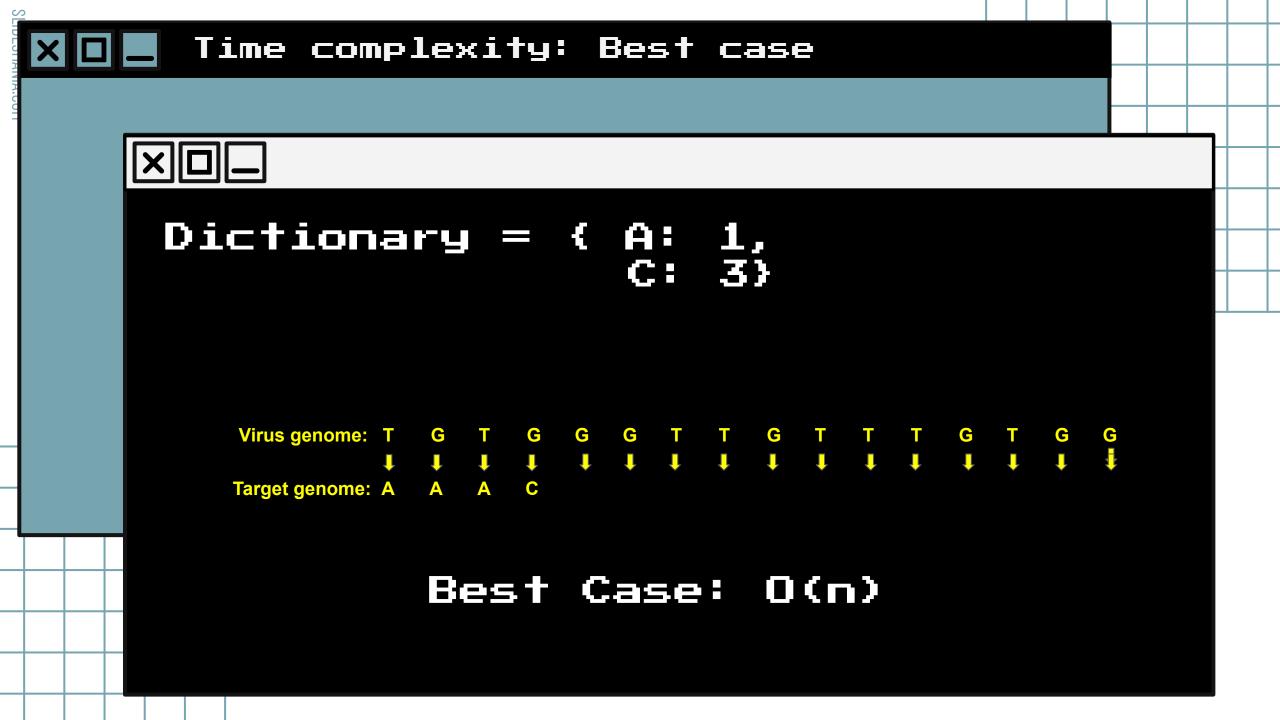
**Target genome:** 

Worst Case: O(nm)













# Conclusion

Algorithm	Preprocessing Time	Matching Time
Naive	0	O(nm)
KMP	O(m)	O(n)
вмн	O(m)	O(nm)



# BEST ALGORITHM?

Ideas from Knuth\_Morris\_Pratt <u>Algorithm</u>







# However, Comparing the three

### Naive Algorithm

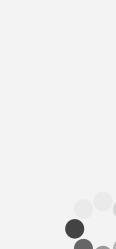
- Easy to think of
- Easy for beginners to implement
- Inefficient

- LPS is not applicable for every case
- Code and concept is Basically brute not easy to think of

### KMP Algorithm BMH Algorithm

- Simple to code, concept based on index of character
- force in worst case





# 

