

String Algorithms and Data Structures

Z-values and the Z-algorithm

CS 199-225

February 2, 2026

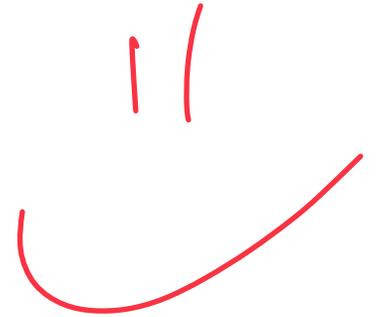
Brad Solomon

This is a
two-part assignment
(this week
+ next week)

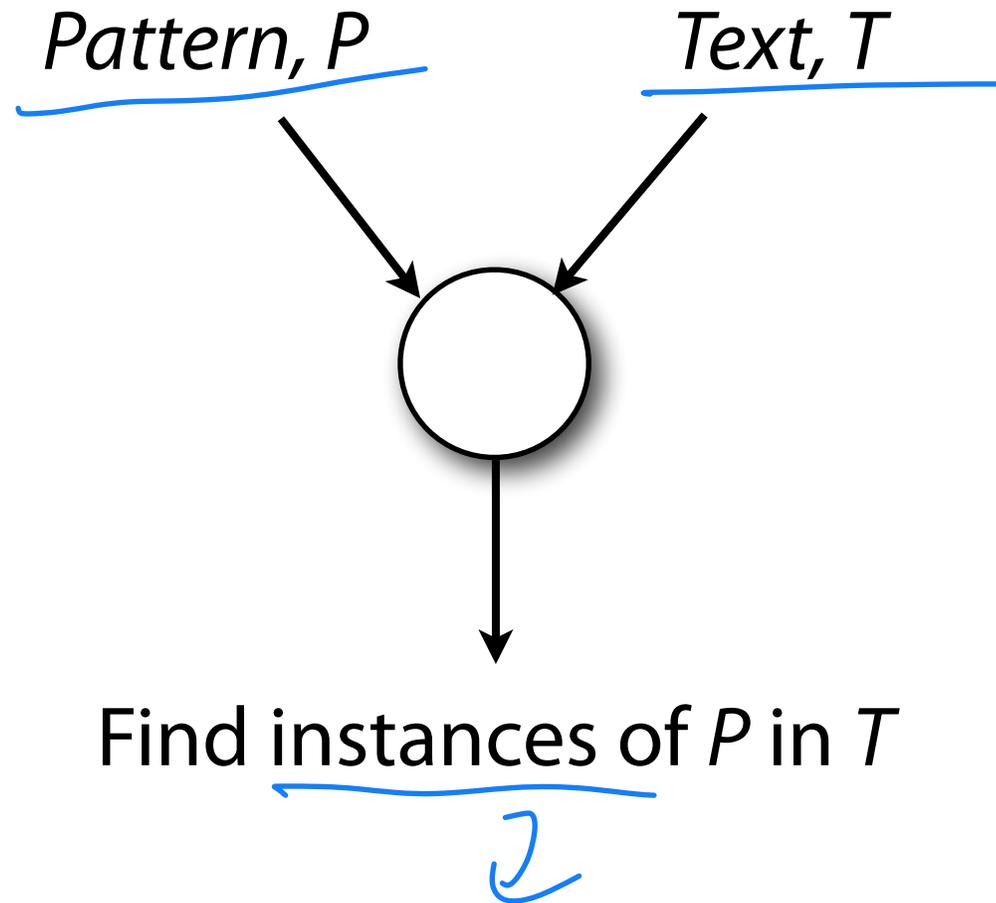


UNIVERSITY OF
ILLINOIS
URBANA - CHAMPAIGN

Department of Computer Science



Exact Pattern Matching



'instances': An exact, full length copy

Exact Pattern Matching

What's a simple algorithm for exact matching?

P: word

T: There would have been a time for such a word

word word word word word word word word word word

word word word word word word word word

word word word word word word word word

word word word word word word word word

word word word word word word word word

← One occurrence

Try all possible alignments. For each, check if it matches. This is the *naïve algorithm*.

Exact Pattern Matching

What is good about the naive solution?

↳ Very easy to code! / Is follow / to debug

↳ Its correct

↳ No preprocessing required

What is bad?

↳ Pretty slow

Exact Pattern Matching

What is our time complexity?



$$(n = |P|, \quad m = |T|)$$

Exact Pattern Matching

What is our time complexity?

$(n = |P|, m = |T|)$

(# of alignments) x (cost of an alignment)



$n - m$ proposed

n / m proposed

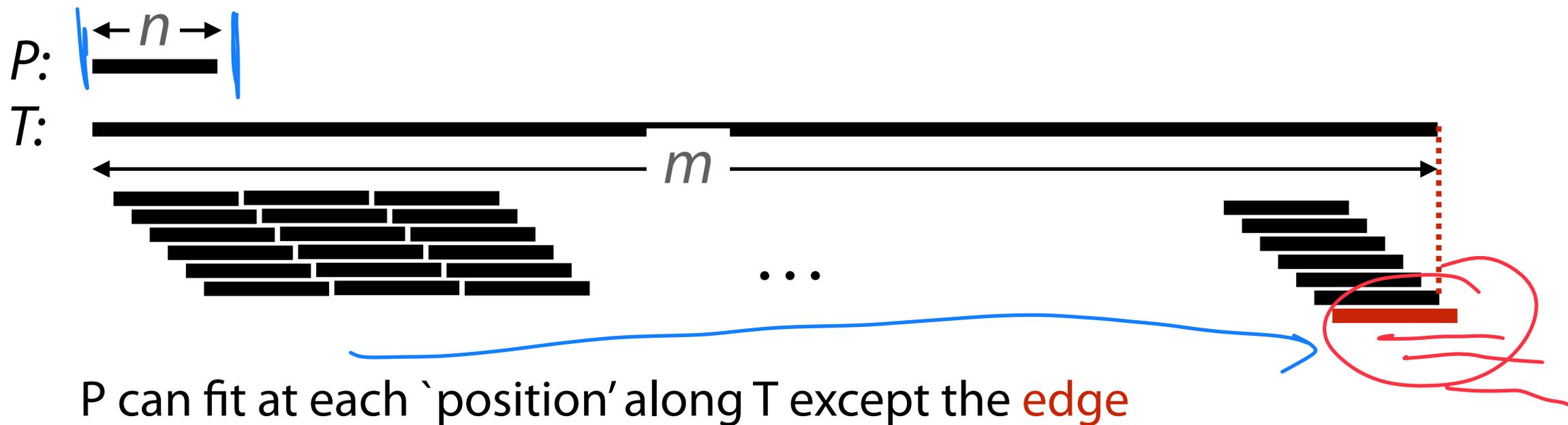
$n - m + 1$ proposed

Exact Pattern Matching

What is our time complexity?

$$(n = |P|, \quad m = |T|)$$

(# of alignments) x (cost of an alignment)



Exact Pattern Matching

What is our time complexity?

$(n = |P|, m = |T|)$

() x (cost of an alignment)

P: aaaa

T: aa

aaaa aaaa aaaa aaaa aaaa aaaa aaaa aaaa aaaa

aaaa aaaa aaaa aaaa aaaa aaaa aaaa aaaa **aaaa**

aaaa aaaa aaaa aaaa aaaa aaaa aaaa aaaa **aaaa**

aaaa aaaa aaaa aaaa aaaa aaaa aaaa aaaa **aaaa**

aaaa aaaa aaaa aaaa aaaa aaaa aaaa aaaa

*(n-1) alignments
which go
out of
bounds*

There are n-1 positions which extend past the edge of T

Exact Pattern Matching

What is our time complexity?

$$(n = |P|, \quad m = |T|)$$

$$\underline{(m-n+1)} \times (\text{cost of an alignment})$$

P : aaaa



T : aaa

aaaa aaaa aaaa aaaa aaaa aaaa aaaa aaaa aaaa

aaaa aaaa aaaa aaaa aaaa aaaa aaaa aaaa

aaaa aaaa aaaa aaaa aaaa aaaa aaaa aaaa

aaaa aaaa aaaa aaaa aaaa aaaa aaaa aaaa

aaaa aaaa aaaa aaaa aaaa aaaa aaaa aaaa

Each alignment compares characters.

Exact Pattern Matching

What is our time complexity?

$(n = |P|, \quad m = |T|)$

$$\theta((m - n + 1) \times n)$$

$$m \cdot n - n^2 + n$$

$$m \gg n$$



$$O(m \cdot n)$$

$$O(|P| |T|)$$

String Algorithms in Genomics

P: Read ($n = \sim 50-150$)

CTCAAACCTCTGACCTTTGGTGATCCACCCGCCTAGGCCTTC

42

T: Reference ($m = \sim 3$ billion)

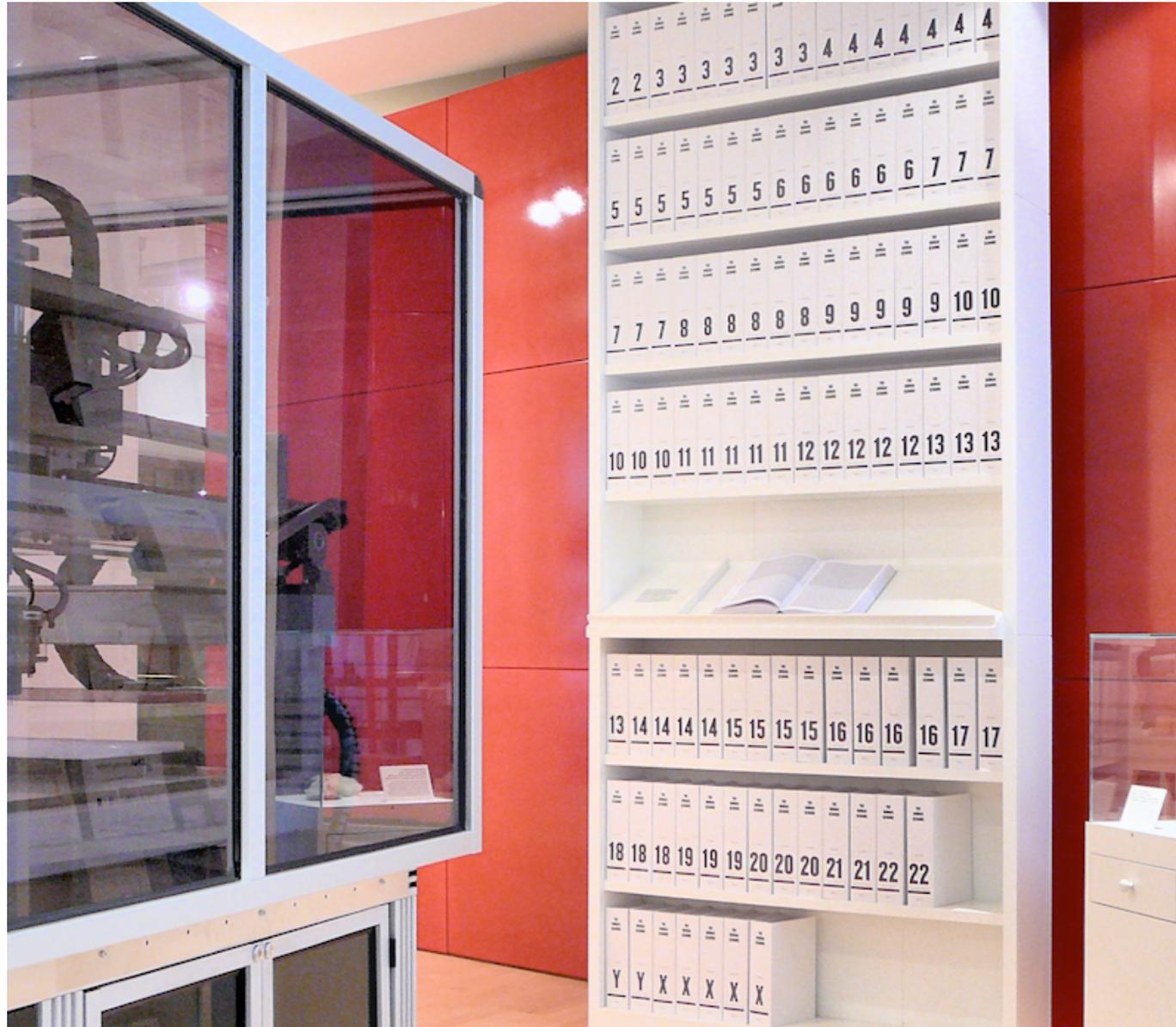
GATCACAGGTCTATCACCTATTAACCACTCACGGGAGCTCTCCATGCATTTGGTATTTT
CGTCTGGGGGATGCACGCGATAGCATTGCGAGACGCTGGAGCCGGAGCACCTATGTC
GCAGTATCTGTCTTTGATTCTGCCTCATCTATTATTTATCGCACCTACGTTCAATATT
ACAGGCGAACATACTTACTAAAGTGTGTTAATTAATTAATGCTTGTAGGACATAATA
ACAATTGAATGTCTGCACAGCCACTTCCACACAGACATCATAACAAAAATTTCCACCA
AACCCCCCTCCCCGCTTCTGGCCACAGCACTTAAACACATCTCTGCCAAACCCAAAA
ACAAAGAACCCTAACACCAGCCTAACAGATTTCAAATTTTATCTTTTGGCGGTATGCAC
TTTTAACAGTCACCCCCCACTAACACATTATTTTCCCCTCCCCTCCACTCCATACTAAT
CTCATCAATACAACCCCCGCCATCCTACCCAGCACACACACACCCGCTGCTAACCCATA
CCCCGAACCAACCAACCCCAAGACACCCCCACAGTTTATGTAGCTTACCTCCTCAA
GCAATACACTGACCCGCTCAAACCTCTGGATTTTGGATCCACCCAGCGCCTTGGCCTAAA
CTAGCCTTTCTATTAGCTCTTAGTAAGATTACACATGCAAGCATCCCCGTTCCAGTGAGT
TCACCTCTAAATCACCACGATCAAAGGAACAAGCATCAAGCACGCAGCAATGCAGCTC
AAAACGCTTAGCCTAGCCACACCCCCACGGGAAACAGCAGTGATTAACCTTAGCTATAA
ACGAAAGTTAACTAAGCTATACTAACCCAGGGTTGGTCAATTTGGTGCCAGCCACCT
GGTCACACGATTAACCAAGTCAATAGAAGCCGGCGTAAAGAGGTTTGTAGATCACCC
TCCCCAATAAAGCTAAAACCTCACCTGAGTTGTA AAAA ACTCCGTTGACACAAAATAGAC
TACGAAAGTGGCTTTAACATATCTGAACACACAATAGCTAAGCCCAA ACTGGGAT TAGA
TACCCACTATGCTTAGCCCTAAACCTCAACAGTTAAATCAA AAAACTGCTCGCCAGAA
CACTACGAGCCACAGCTTAAAACCTCAAAGGACCTGGCGGTGCTCATATCCCTCTAGAGG
AGCCTGTTCTGTAATCGATAAACCCCGATCAACCTCACCACCTCTGCTCAGCCTATAT
CCGCCATCTTCAGCAAACCCTGATGAAGGCTACAAAGTAAGCGCAAATACCCACGTA
ACGTTAGGTCAAGGTGTAGCCCATGAGGTGGCAAGAAATGGGCTACATTTCTACCCCA
AAAAC TACGATAGCCCTTATGAAACTTAAAGGGTGAAGGTGGATTTAGCAGTAAACTAAG
AGTAGAGTGCTTAGTTGAACAGGGCCCTGAAGCGGTACACACCCGCCGTCACCTCCTC
AAGTATACTTCAAAGGACATTTAACTAAAACCCCTACGCATTTATATAGAGGAGACAAGT
CGTAACCTCAAACCTCTGCCTTTGGTGATCCACCCGCCTTGGCCTACCTGCATAATGAAG
AAGCACCCTAACCTTACACTTAGGAGATTTCAACTTAACTTACCCCTCTGAGCTAAAGCTA



String Algorithms in Genomics



String Algorithms in Genomics



Improving exact pattern matching



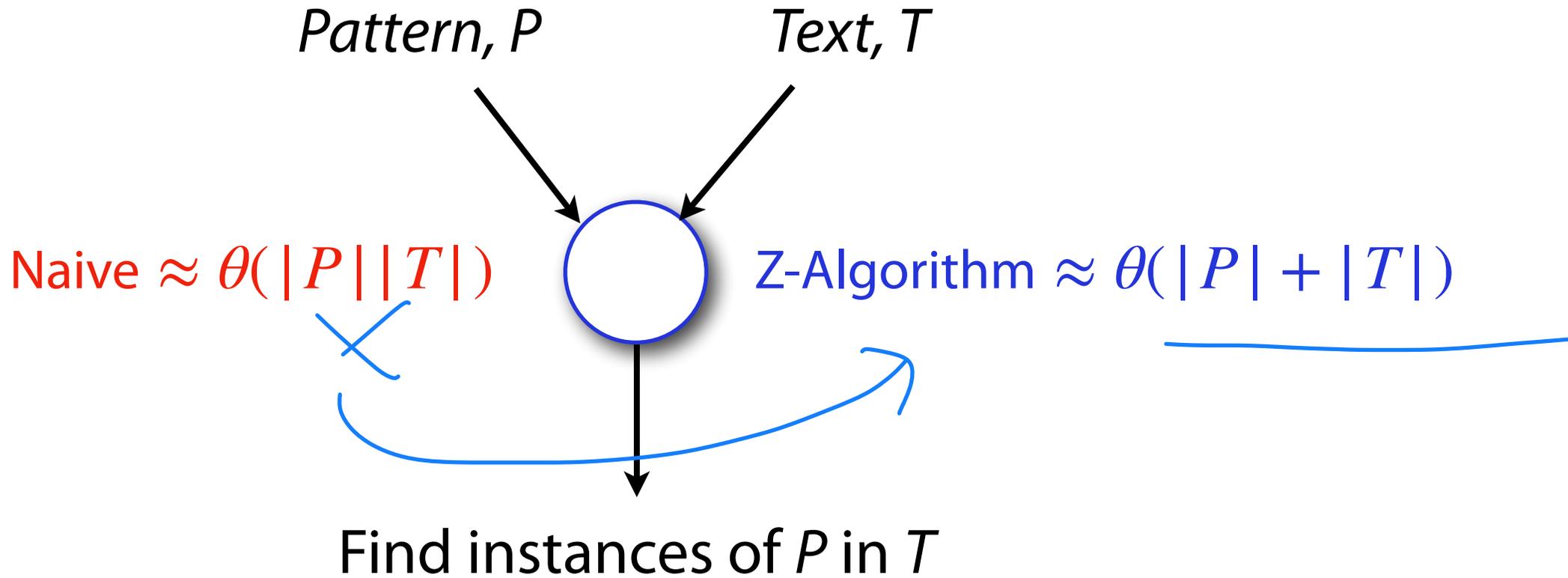
How can we do better than the naïve algorithm?

... If we have infinite space?

... If I tell you the pattern ahead of time?

... If I tell you the text ahead of time?

Exact Pattern Matching w/ Z-algorithm



'instances': An exact, full length copy

The Z-value [$Z_i(S)$]

Given a string S , $Z_i(S)$ is the length of the longest substring in S , starting at position i , that matches a prefix of S .



S: 0 1 2 3 4 5 6 7 8 9
 T T C G T T A G C G
 T T C G T T A G C G

~~$Z_0(S) = 10$~~

$Z_1(S) = 1$

$Z_2(S) = 0$

$Z_3(S) =$

$Z_4(S) =$

$Z_5(S) =$

The Z-value [$Z_i(S)$]

Given a string S , $Z_i(S)$ is the length of the longest substring in S , starting at position i , that matches a prefix of S .

0 1 2 3 4 5 6 7 8 9
S: T T C G T T A G C G

$$Z_0(S) = 10$$

$$Z_1(S) = 1$$

$$Z_2(S) = 0$$

$$Z_3(S) = 0$$

$$Z_4(S) = 2$$

$$Z_5(S) = 1$$

The Z-value [$Z_i(S)$]

Given a string S , $Z_i(S)$ is the length of the longest substring in S , starting at position $i > 0$, that matches a prefix of S .

 0 1 2 3 4 5 6 7 8 9
S: **T T C G T T A G C G**

$$\del{Z_0(S) = 10}$$

$$Z_3(S) = 0$$

$$Z_1(S) = 1$$

$$Z_4(S) = 2$$

$$Z_2(S) = 0$$

$$Z_5(S) = 1$$

Calculating the Z-values

Naive: Compute the Z-values by *explicitly* comparing characters (left-to-right scan):

$$Z_1 = 3$$

A A A A B A A C A A B A A ...
↑ ↑ ↑ ↑ X

A A A A B A A C A A B A A ...
↑ ↑ ↑ ↑ X

$$Z_5 = 2$$

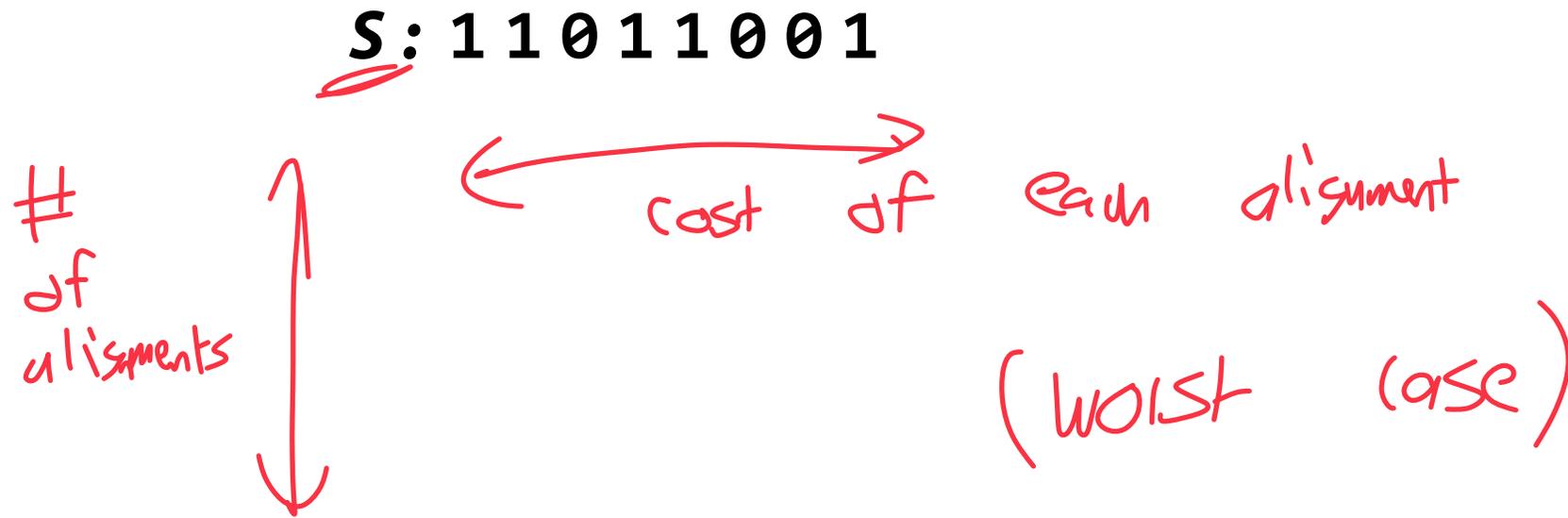
A A A A B A A C A A B A A ...
X

A A A A B A A C A A B A A ...
X

What is our time complexity?

Calculating the Z-values

Naive: Compute the Z-values by *explicitly* comparing characters (left-to-right scan):

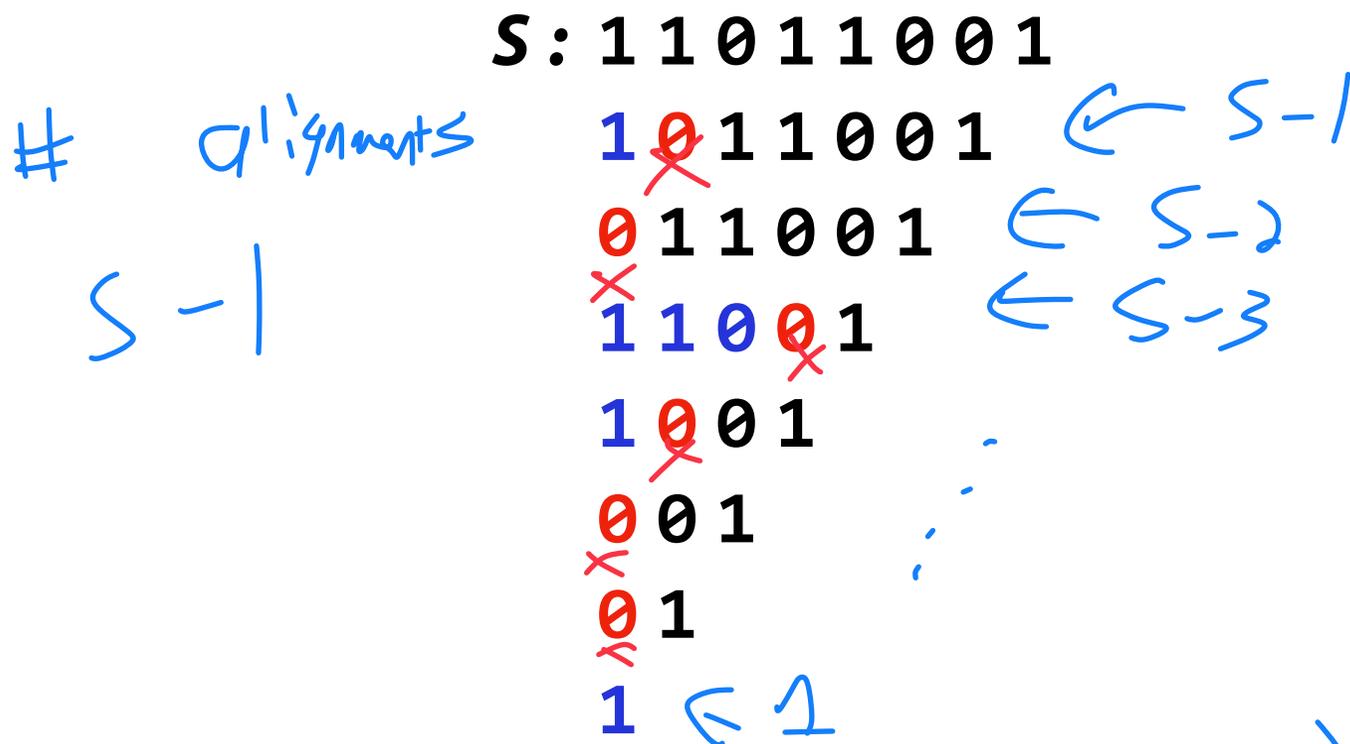


What is our time complexity?

Calculating the Z-values



Naive: Compute the Z-values by *explicitly* comparing characters (left-to-right scan):



$$\sum_{i=1}^{s-1} i = \frac{s(s+1)}{2}$$

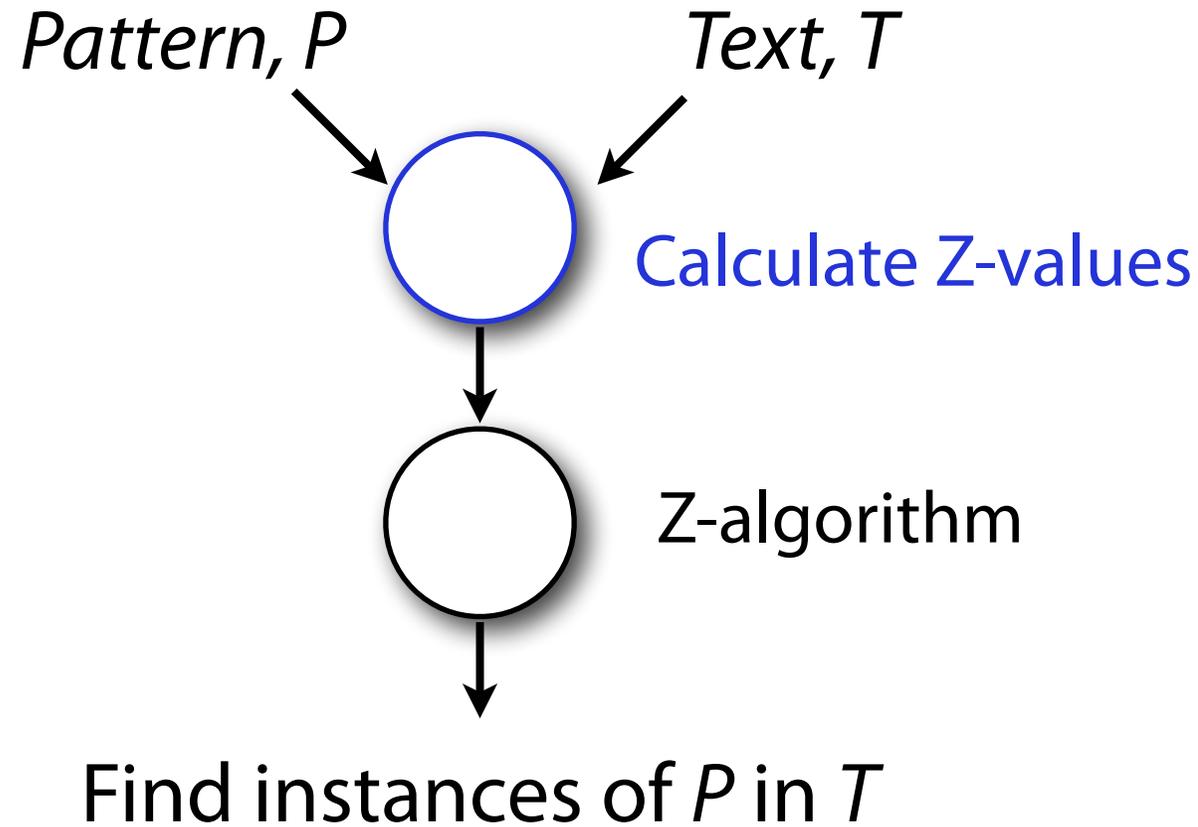
What is our time complexity?

$$O(s^2)$$

||

Pattern matching with the Z-value

Given a Z_i value calculator, how do we solve pattern matching?



Z-value Pattern Matching

To solve pattern matching (given P and T), let $S = P\$T$

$\$$ = 'terminal character', outside alphabet

will never match anything

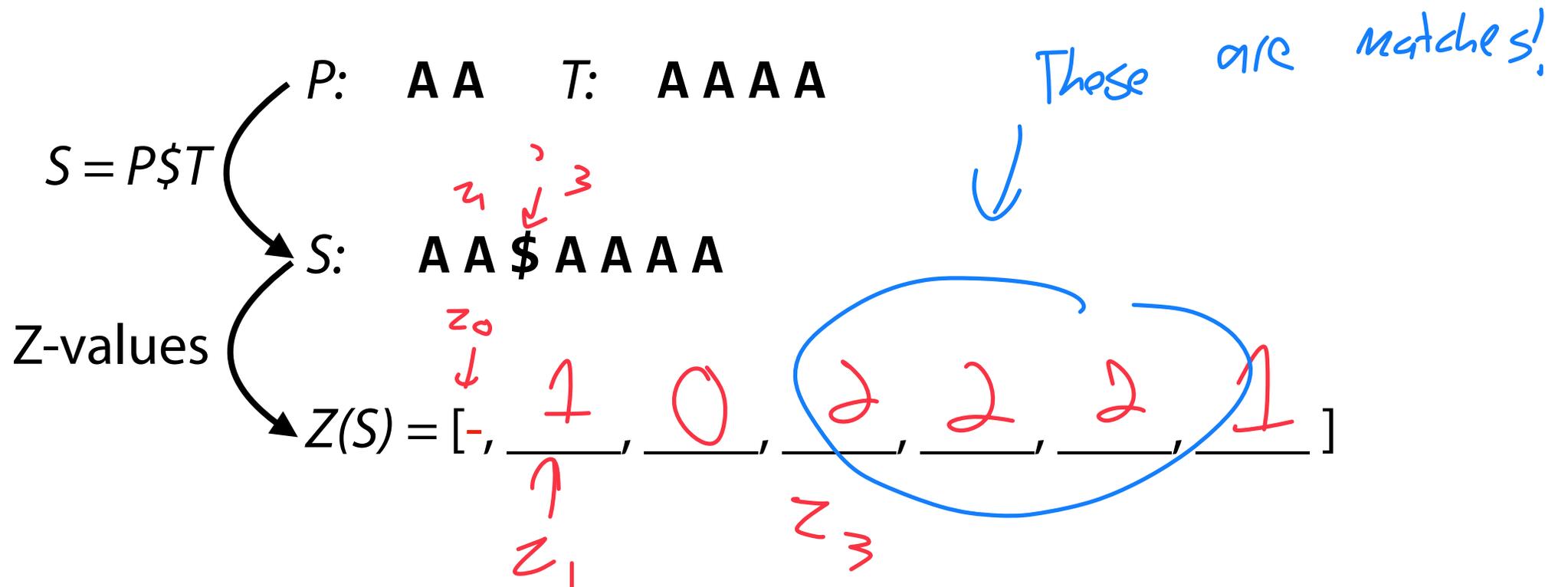
$S = P\$T$

$P:$	AA	$T:$	AAAA
$S:$	AA\$AAAA		

Z-value Pattern Matching

To solve pattern matching (given P and T), let $S = P\$T$

$\$$ = 'terminal character', outside alphabet



Z-value Pattern Matching

To solve pattern matching (given P and T), let $S = P\$T$

$\$$ = 'terminal character', outside alphabet

P : AA T : AAAA

0 1 2 3 4 5 6
 S : AA\$AAAA

$Z(S) = [-, 1, 0, 2, 2, 2, 1]$

What Z_i values are matches?

Z_3, Z_4, Z_5 ← length of P
 $i - |P| - 1$ ← $\$$

What are the matching indices in T ?

Z-value Pattern Matching

To solve pattern matching (given P and T), let $S = P\$T$

$\$$ = 'terminal character', outside alphabet

P : AA T : AAAA

0 1 2 3 4 5 6
 S : AA\$AAAA
 0 1 2 3

$Z(S) = [-, 1, 0, 2, 2, 2, 1]$

What Z_i values are matches? Any $|P|$ value in $Z(S)$

What are the matching indices in T ? 0, 1, 2 which are $(i - |P| - 1)$

Z-value Pattern Matching



P: TT T: CTTA

S: TT\$CTTA

Z(S): (-, 2, 0, 0, 2, 1, 0)

z_i

Z-value search pseudo-code

1. Concatenate (S=P\$T)

2. Calculate Z-values for S

3. For $i > 0$, match if $Z_i = \underline{1P}$

Match is **not** at i , but instead at

$i - |P| - 1$

Assignment 2: a_zval

Learning Objective:

Construct a Z-value calculator and measure its efficiency

Demonstrate use of Z-values in pattern matching

Consider: Our goal is $\theta(|P| + |T|)$. Does Z-value search match this?

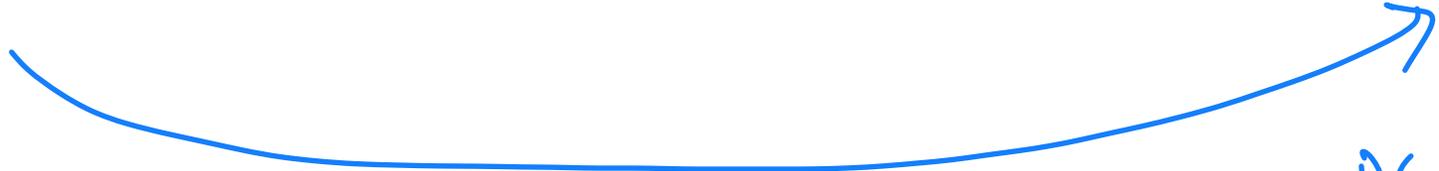
↳ Z value calculation is simply too slow!
 $\Theta(|P|(P+|T|))$

End-of-class brainstorm

What information does a single Z-value tell us?

If I know $Z_{i-1}(S)$, can I use that information to help me compute $Z_i(S)$?


Past calc


present
index

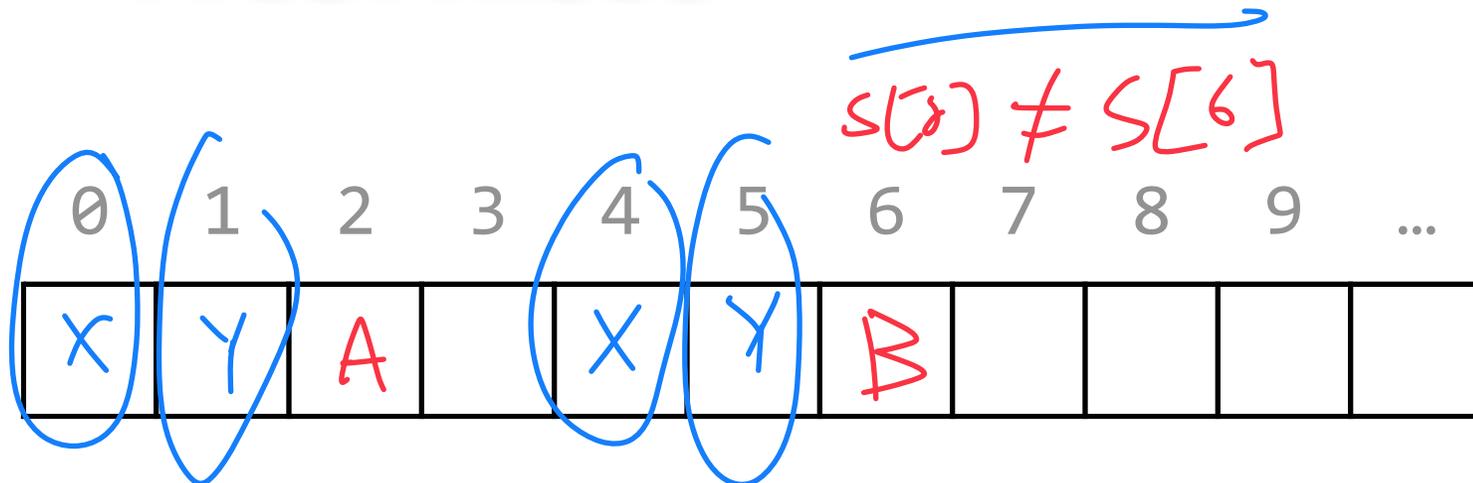
The Z-value (Take 2)

Given a string S , $Z_i(S)$ is the length of the longest substring in S , starting at position i , that matches a prefix of S .

What information does this give us?

S : XYCATAB

$$Z_4 = 2$$



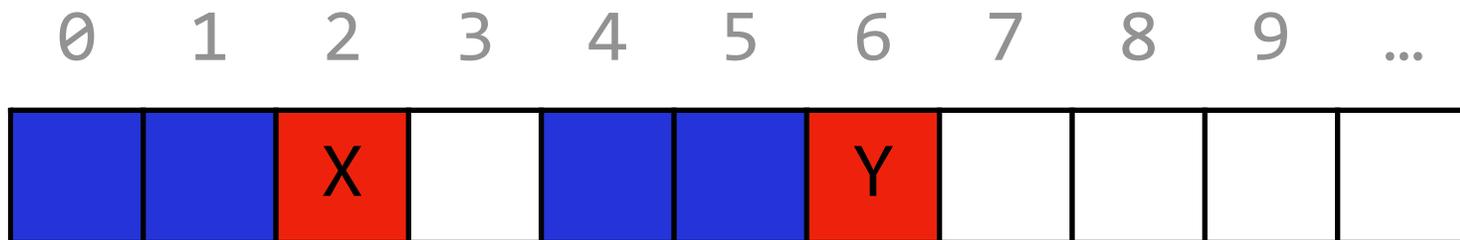
The Z-value (Take 2)

Given a string S , $Z_i(S)$ is the length of the longest substring in S , starting at position i , that matches a prefix of S .

What information does this give us?

S : **XYXYXABCD**

$$Z_4 = 2$$

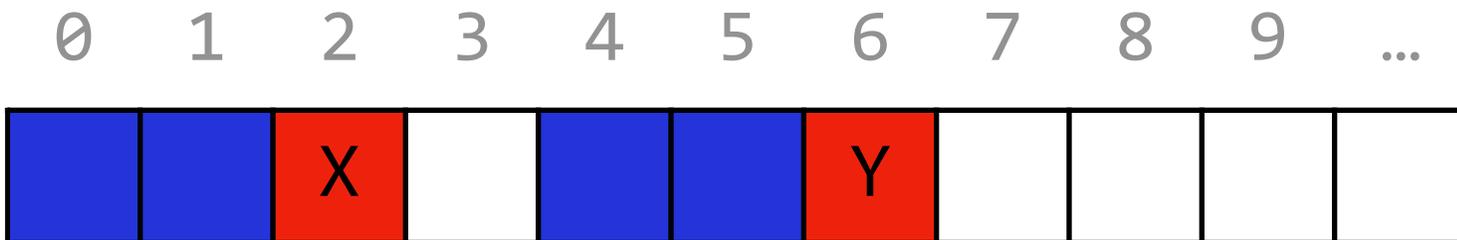


The Z-value (Take 2)

Given a string S , $Z_i(S)$ is the length of the longest substring in S , starting at position i , that matches a prefix of S .

What information does this give us?

$\begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ S: & \boxed{T} & \boxed{T} & C & G & \boxed{T} & \boxed{T} & A & G & C & G \end{matrix} \quad Z_4 = 2$



The Z-Algorithm

Assume we've computed Z_1, \dots, Z_{i-1} and need to calculate Z_i

Case 1: We know nothing about the characters at $S[i]$

$Z_1 = ?$

	0	1	2	3	4	5	6	7
	A	A	A	A	B	B	B	B
	A	A	A	A	B	B	B	B

(Note: In the original image, the character 'A' at index 1 is circled in red. Blue arrows show comparisons between the first 'A' and subsequent 'A's. Red arrows show comparisons between the first 'A' and 'A's in the second row. Red 'X' marks are on the 'A' at index 3 in both rows and the 'B' at index 4 in the second row. A red underline is under indices 5-7.)

Case 2: We know something about the characters at $S[i]$

$Z_2 = ?$

	0	1	2	3	4	5	6	7
	A	A	A	A	B	B	B	B
	A	A	A	A	B	B	B	B

(Note: In the original image, the 'A' at index 2 is circled in red. The cells at (0,0), (1,0), (0,1), (1,1), (0,2), (1,2), (0,3), (1,3) are shaded blue. The cells at (0,3), (1,3), (0,4), (1,4) are shaded red. Red arrows show comparisons between the first 'A' and 'A's in the second row.)

The Z-Algorithm

$$Z_1 = 3$$

$$Z_2 = ?$$

The diagram illustrates the Z-Algorithm on the string "AABAABBA". The string is represented as a 2x8 grid of characters. The first row contains 'A', 'A', 'A', 'A', 'B', 'B', 'B', 'B'. The second row contains 'A', 'A', 'A', 'A', 'B', 'B', 'B', 'B'. The first three columns (indices 1, 2, 3) are shaded blue, representing the first Z-value of 3. The fourth column (index 4) is shaded red, representing the current character being calculated. A dashed line is drawn above the first three columns, indicating the current window. Handwritten annotations include: $l=1$ (black) with a black arrow pointing to index 1; $i=2$ (red) with a red arrow pointing to index 2; and $r=3$ (blue) with a blue arrow pointing to index 3. A blue vertical bar is drawn at index 4, and a black horizontal bar is drawn below the first three columns.

0	1	2	3	4	5	6	7
A	A	A	A	B	B	B	B
A	A	A	A	B	B	B	B

We track our current knowledge of S using three values: i, r, l

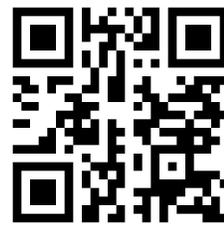
i , the current index position being calculated

r , the index of the rightmost character which has ever been matched

l , the index of Z-value which r belongs too

The Z-Algorithm

Join Code: 225



i , the current index =

Start
2

End
3

r , the furthest matching char =

1

1

l , the furthest reaching Z-value =

1

1

-	1						
0	1	2	3	4	5	6	7
A	A	B	B	A	A	B	A
A	A	B	B	A	A	B	A

only updates if
new characters matches

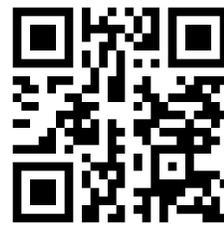
The Z-Algorithm

i , the current index = Start End
 r , the furthest matching char = 1 1
 l , the furthest reaching Z-value = 1 1

	-	1	0	0				
	0	1	2	3	4	5	6	7
A	A	B	B	A	A	B	A	
A	A	B	B	A	A	B	A	

The Z-Algorithm

Join Code: 225



i , the current index =

r , the furthest matching char =

l , the furthest reaching Z-value =

Start

End

4

5

1

6

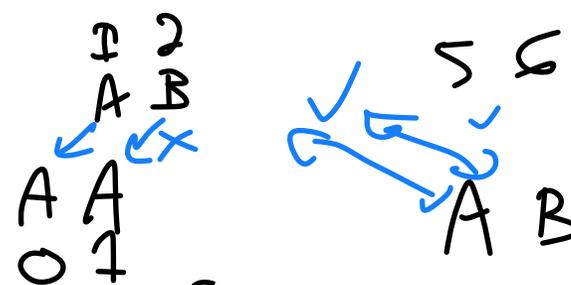
1

4

-	1	0	0	<u>4</u>	<u>5</u>	<u>6</u>	<u>7</u>
0	1	2	3	4	5	6	7
A	A	B	B	A	A	B	A
A	A	B	B	A	A	B	A

The Z-Algorithm

Cool! 😊



i , the current index =

r , the furthest matching char =

l , the furthest reaching Z-value =

Start

End

5

6

6

6

4

4

- 1) B/c $z_4 = 3$ I know 4, 5, 6 is is
- 2) Index 5 matches index 1 1, 2, 3

$z_1 = 1$

-	1	0	0	3	1		
0	1	2	3	4	5	6	7
A	A	B	B	A	A	B	A
A	A	B	B	A	A	B	A

claim: z_5 is is z_1

No work needed

↑ ????

The Z-Algorithm

i , the current index =

Start

7

End

8

r , the furthest matching char =

6

7

l , the furthest reaching Z-value =

4

7

-	1	0	0	3	1	^{z₆} 0	<u>1</u>
0	1	2	3	4	5	6	7
A	A	B	B	A	A	B	A
A	A	B	B	A	A	B	A

The Z-Algorithm

Start

End

i , the current index =

r , the furthest matching char =

l , the furthest reaching Z-value =

-	1	0	0	3	1	0	1
0	1	2	3	4	5	6	7
A	A	B	B	A	A	B	A
A	A	B	B	A	A	B	A

The Z-Algorithm

$\langle P \rangle$ \$ $\langle T \rangle$
⏟



Intuition: We can use the previous Z_1, \dots, Z_i to compute Z_{i+1} !

Track 'what we know' using three integers: i, r, l

Next week: Review how integers are updated to define specific cases.

Claim: If i is monotonically increasing then we never compare same character twice

$$\hookrightarrow O(|P| + |T|)$$