

CSE 494/598: ALGORITHMS IN COMPUTATIONAL BIOLOGY

Fall 2022

Instructor: Heewook Lee **Date:** 09/12/2022

String Algorithms: Exact Pattern Matching

Note: This note may contain errors.

4 Exact Pattern Matching

In the following few lectures, we will explore algorithms for performing exact pattern matching

4.1 Motivation

Exact pattern matching is the most basic string operation and useful.

1. Application is endless. (Finding a word within a document)
2. Similarly, its application in biology is also very strong. Finding restriction enzyme sites (proteins that cut the DNA at specific location in the genome) within a genome. Restriction enzyme sites are short k -mers. Finding all of these sites is an instance of exact matching.
3. We have already seen an application of exact string matching in ori-finding problem.

4.2 Problem Formulation

ExactPatternMatching: Given a pattern (query) P , does it exist in the text (reference) T or find all occurrences of the pattern in the text?

Example:

Text position	1	i	n = 17
Text	ATTCATATTTCGGCTAT		
Pattern	GCAT		
Pattern Pos	1..m = 4		

4.3 naïve Approach

We already came up with a naïve approach where the algorithm scans the text from left to right and perform string comparisons.

For each pattern, the algorithm performs up to m comparisons (pattern length) while scanning the text at $(n - m + 1)$ positions, resulting in $m(n-m+1)$ operations, giving us $O(nm)$.

Q: Does the naïve algorithm need to perform $m(n-m+1)$ operations?

A: It depends on whether we stop as soon as we find a match or we continue until we have found all the matches of the pattern within the text. In the former case, we maybe able to perform less number of steps.

Q: Assume we just need to find the first match, does it ever need to perform $m(n-m+1)$ operations?

A: Yes, the worst case, where the pattern is located at the tail of the string.

Q: Why is the algorithm not so applicable?

A: Imagine we are trying to match a sequence (say approximately 1 kb sequence) in human genome (3 billion bp). The total run-time would be on the order of 3×10^{12} operations. On a 3GHz processor (and assuming each comparison takes exactly one CPU cycle), the calculation would take 1,000 seconds (16 mins). Does this sound not too bad?? Modern sequencing experiments generate tens to hundreds of millions of sequences. This is probably too slow.

Key Question: How can we improve on this algorithm? The key idea is to use information we've already figured out when making earlier comparisons. We will first look at Z algorithm.

4.4 Z Algorithm

Let's switch gear and first talk about a simple function Z that captures the internal structure of the text. For every position, i in $(0, n]$ in text T , we define Z array $Z[i]$ to be the length of the longest prefix of $T[i...n]$ that matches exactly a prefix of $T[0...n]$. An example is shown here:

```
Text:  A T T C A C T A T T C G G C T A T
Z[i]:  0 0 0 0 1 0 0 4 0 0 0 0 0 0 2 0
```

Q: Assume we have an efficient way to generate Z -array for a given text, can Z values in aligning a pattern to a text? How would we use the Z -array to find occurrences of a pattern in text?

Possible Answer: Create a string S by concatenating P and T , $S = PT$. Find all $Z[i]$ values in S and any $Z[i] \geq |P|$ identifies an occurrence of P in T .

Better approach: We can construct the following string $S = P\$T$, where P and T are the pattern and the text, respectively, and $\$$ is a character that does not match any other character in either pattern or text.

Q: How do Z values computed for the string S defined above help us find matches between P and T ? It suffices to look at the Z values corresponding to characters in T . Any Z -value equal to the length of P indicates the pattern matches.

A: Find all $Z[i]$ values in S , and any $Z[i] = |P|$ identifies an occurrence of P in T .

Q: Can we use the reverse text $S' = TP$ or $T\$P$. NO! It's not symmetric.

Q: What is the runtime and memory usage of this algorithm?

A: First the memory: Need to store the $Z[i]$'s in addition to the pattern and text, for a total of $2(n+m)+1$ memory locations. More memory but it's still linear in the size of the inputs. Once $Z[i]$'s are computed, we simply need to examine each position in the Z array to find all matches between P and T , or to decide that a match doesn't exist, resulting in $O(n)$ if $Z[i]$'s can be computed efficiently.

4.4.1 Can Z -values also be computed in linear time and space?

Let's first discuss how a naïve algorithm computes.

For each position from 2 to $|S|$, compare and write down $Z[i]$, resulting in quadratic time.

Q: can we somehow reuse the work done earlier?

A: The intuition is that when computing $Z[i]$ we can use the $Z[j]$ values for $j < i$.

Here is the algorithm. The algorithm computes $Z[i]$, $i = 2$ to $i = |S|$

First, let's define a few things we need to further explain. Suppose we have $Z[i] > 0$, meaning we have a prefix of $S[i..|S|]$ that matches a prefix of S and this is the longest prefix at position i . Any prefix found at any position i is a Z -box. At position i , among all Z -boxes, let r be the rightmost position among all Z -boxes found so far and l be the leftmost position in Z -box whose right-end is r (Fig 1). Since the each prefix represented by a box in the figure is the longest prefix at each position, $x \neq y$.

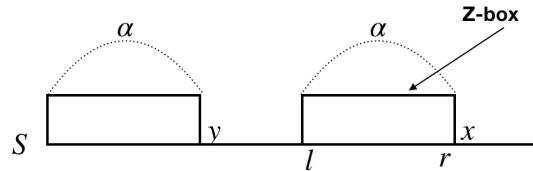


Figure 1: Z -box: r is the rightmost position in any Z -box found so far. l is the leftmost position in Z -box whose right end is r . x and y are the characters after the last position of the Z -box bounded by l and r and the matching prefix of S respectively.

Now, assume that we have already computed $Z[2] \dots Z[k-1]$, and we want to compute $Z[k]$. We will break down into 2 main cases based on relative location of k and r .

Case 1: ($k > r$). In this case, there isn't any information for us to utilize to deduce $Z[k]$. Therefore, we need to compute $Z[k]$ explicitly

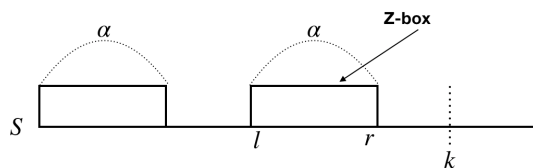


Figure 2: Case 1 ($k > r$)

Case 2: ($k \leq r$)

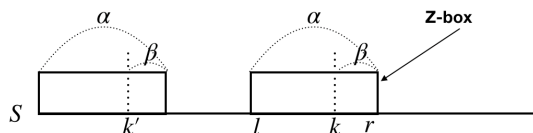
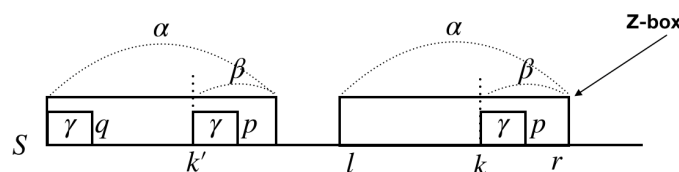


Figure 3: Case 2 ($k \leq r$)

We will break this case into 3 sub-cases.

Case 2a: ($Z[k'] < |\beta|$), meaning $|\gamma| < |\beta|$



Since $p \neq q$, we can conclude that $Z[k] = Z[k']$ without any character comparisons.

On third iteration ($i = 3$), we have the following alignment (mismatch is shown in red, also marked as 'x' denoting the adjacent position of the last matched position)

T Pos	1	i	x	n														
Text	X	Y	A	B	C	X	A	B	C	X	A	D	C	D	A	F	E	A
Pattern	A	B	C	D														

Q: Can we shift the pattern by more than 1 for next iteration?

A: Notice that the pattern is made up of different characters. We can conclude that we can shift the pattern to position $i = 6$ skipping positions 4 and 5.

Let's consider another example where we have a different pattern but same text. At iteration $i = 3$ we have the following alignment:

T Pos	1	i	x	n														
Text	X	Y	A	B	C	X	A	B	C	X	A	D	C	D	A	F	E	A
Pattern	A	B	C	X	A	B	C	D	E									

Q: How far can we shift?

A: Notice that the suffix of $P[1..(x-1)]$ is also the prefix of P . We can simply shift the pattern over until its prefix matches a previously matched section of the text (suffix of $P[1..(x-1)]$).

The intuition is that once we run into a mismatch or run out of characters to match (a complete match), we can shift as much as the length of the **longest** suffix of $P[1..(x-1)]$ matches exactly a prefix of P .

Here is another example where the pattern results is a complete match:

T Pos	1	x	n												
Text	A	C	A	B	A	C	A	C	D	B	A	C	A	C	D
Pattern	A	C	A	B	A	C	A								

You can see that P matches (ran out of character to match, x denotes the position after last match) exactly at position $i = 1$. We need to shift to find more occurrences of P in T if any.

Q: How far can we shift?

A: You can see that 'ACA' is the longest suffix of $P[1..(x-1)]$ (which is P in this case) that matches exactly a prefix of P . Therefore, we can shift the pattern to position 5 and start comparing from $x = 8$.

In order to use the intuition obtained from our examples, we first define lps (longest prefix suffix) array for pattern P . Let $lps[i]$ be the length of the longest non-trivial suffix of $P[1..i]$ that matches exactly a prefix of the pattern (non-trivial meaning the suffix and the prefix are same as $P[1..i]$ itself).

Q: How do we compute lps array efficiently?

Let us assume that we are trying to compute the value $lps[i+1]$ and that we already know the values $lps[j]$ for all $j \leq i$. Also, let c be the character at position $i+1$ in the pattern ($c = P[i+1]$), and x be the character at position $lps[i]+1$ ($x = P[lps[i]+1]$) (i.e. the character that follows the prefix of P that matched the suffix ending at i).

If the two characters match ($x == c$), we simply extend the matching prefix/suffix pair by one character to obtain $lps[i+1] = lps[i] + 1$.

Q: What happens, however, if $x \neq c$?

A: In this case, $lps[i]$ provides us no additional information, but we might be able to use $lps[lps[i]]$, the suffix of $P[1..lps[i]]$ that matches a prefix of P . Let x' be $P[lps[lps[i]]]$, the character that follows the character at position $lps[lps[i]]$. Here we will compare c with x' and if they are match, $lps[i+1] = lps[lps[i]] + 1$. If

they don't, we simply continue trying further $\text{lps}[\text{lps}[\dots[\text{lps}[i]]\dots]]$ values until we either find a match, or we "bottom out" indicating no such prefix exists.

In the "bottomed out" case, we set $\text{lps}[i + 1] = 1$ if $c = P[1]$, or $\text{lps}[i + 1] = 0$, otherwise.

