

Today we will discuss in detail the exact ^{string} ~~pattern~~ matching problem we presented in the last lecture:

[Exact String Matching Problem: Given a (long) string T and a shorter string P find all occurrences of P in T . Occurrences of P can overlap.]

Applications:

- ① Mapping second-generation sequencing short reads (~ 150 bp) to reference genome (~ 3 Bbp)
- ② Find word occurrences in text documents. Ctrl-S in emacs.
- ③ Keyword search in books, webpages. ~~Construct an inverted index.~~

Example

$T = \text{ABCABCDABABCDABDABDE}$

$P = \text{ABCDABD}$

We will discuss a few algorithms to solve this problem. In their analysis we will treat 'character comparison' as the primitive operation.

Naive algorithm:

for $i = 1, \dots, |T|$

$j = 1$

while $j \leq |P|$ and $T[i+j-1] == P[j]$:

$j++$

if $j == |P| + 1$: report occurrence at position i

How many character comparisons: $|T| \cdot |P| \Rightarrow O(|T| \cdot |P|)$

Q: what is the worst case?

To improve this we should use information in previous loops when comparing P to $T[k\dots]$. $j = 1, \dots, k-i$

What can we infer about $T[k\dots]$ from previous loops:

ABC ABC DAB ...

ABC ~~D~~AB D

ABC ~~D~~ABC

ABC DAB ...

The similarities of substrings in P allow us to make inferences about unnecessary comparisons.

First approach is to preprocess P to find similarities

Review

S : ABCDABD

Prefix of S : substring of S starting at position 1:

A, AB, ABC, ABCDAB

Suffix of S : substring of S ending at position $|S|$:

D, BD, ABD, BCDABD

Z-algorithm (Gusfield Ch. 1): $O(|T| + |P|)$

Fundamental Preprocessing

Def. Given string S and position $i > 1$, $z_i(S)$ is then length of the longest substring of S starting at position i that matches a prefix of S .



$$S[i, \dots, i+z_i] = S[1, \dots, z_i]$$

$$|x| = z_i$$

ABCDABD

\uparrow
 $z_5 = 2$

How do we use z_i to solve the exact string matching problem?

$$S = P\$T$$

positions i and j

$$z_j \neq |P| \ \& \ i > |P| + 1 \Rightarrow P \text{ occurs at } T[i - |P| + 1]$$

$$P = S[1, \dots, |P|] = S[i, \dots, i + |P|] = T[i - |P| + 1, \dots, i]$$

Linear time alg.

to compute Z-scores
gives us linear time
alg. for exact
matching

compute Z-values gives us

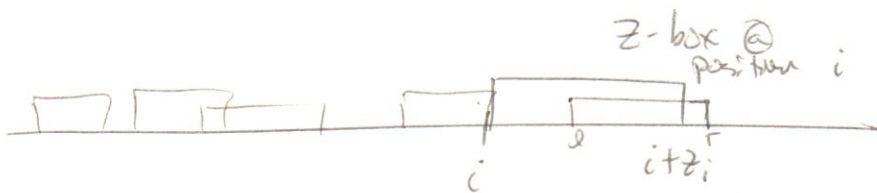
1950-07-04

→ Directly compare $S[2, \dots]$ and $S[1, \dots]$ until first mismatching position
i.e., compute z_2 directly

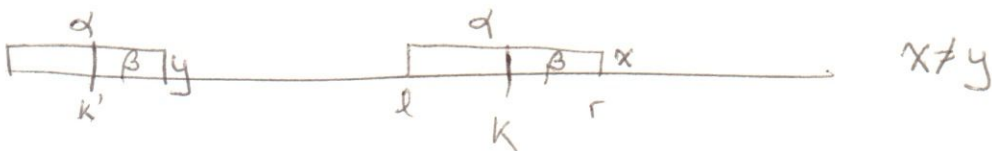
~~Step~~ $k = 2, \dots, \text{~~max}~~ \text{~~iter}~~}$

→ Assume z_2, \dots, z_{k-1} are computed

Def: Z-box at position i is interval $i, \dots, i+z_i$

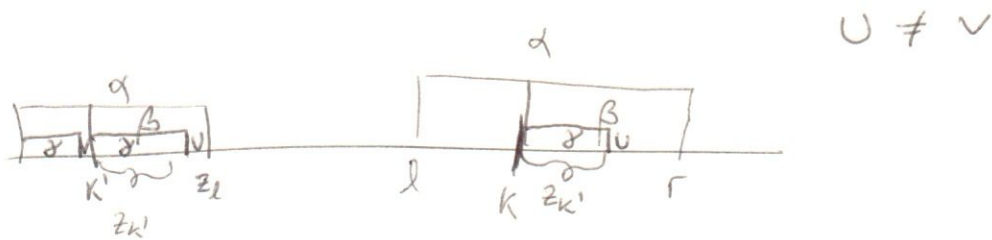


Def: r : end of right-most z -box found through $2, \dots, k-1$
 l : start of right-most z -box found through $2, \dots, k-1$

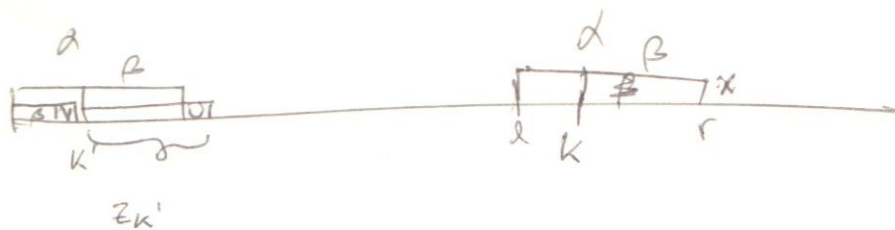


Case 2: $R \leq V$

Case 2a: $z_k' < |\beta| \Rightarrow z_k = \cancel{z_k'}$



Case 2b:

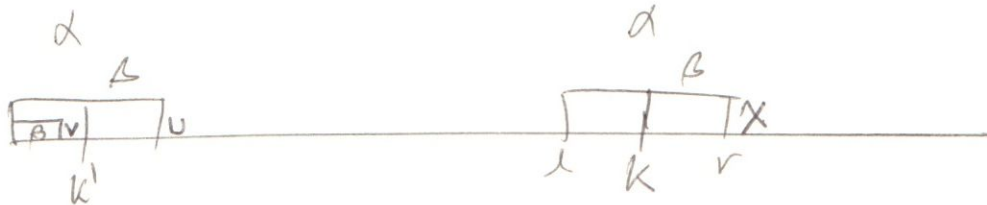


$$U = V$$

$$X \neq U$$

$$\Rightarrow X \neq V$$

Case 2c: $z_{k'} = |\beta| \Rightarrow z_k = |\beta| + q$



$$U \neq V$$

$$X \neq U$$

$$X \neq V$$

compare $S[r+1, \dots]$ to $S[\beta+1, \dots]$ until first mismatch (suppose q matches)

$$\forall \quad r_k = r_{k-1} + q$$

$$l = k$$

Case 1 ~~2b~~ $k > r$

compare $S[k, \dots]$ to $S[1, \dots]$ until first mismatch (suppose $q \geq 0$ matches)

$$r_k = k + q$$

$$l = k$$

Obs: $r_k \geq r_{k-1}$ for all k !!

Analysis:

Correctness: by construction

Runtime: $\underbrace{\# \text{ matching}}_{|S|} + \underbrace{\# \text{ mismatching}}_{|S|}$

$$O(|S|)$$

Step Next class:

KMP

Boyer-Moore

Also $O(|P| + |T|)$. B-M is better in practice, KMP will help you in project 2.

Step

Compression with z-scores

Def

~~if~~ if $z_i = 0$ send $S[i]$

o.w. send z_i instead of $S[i, \dots, i+z_i]$

Def

Case

Case