## 9.1.2. STRING MATCHING ON A 2D GRID

↳ GIVEN A 2D GRID FIND THE OCURRENCE(S) OF PATTERN P IN THAT GRID.

→ SEARCH CAN BE MADE IN 4 OR 8 DIRECTION
→ PATTERN CAN BE FOUND IN A STRAIGHT LINE OR DIAGONAL

```
A B C D E F G H I G G
H E B K W A L D O R K
F T Y A W A L D O R M
E T S I M R L Q S R C
B Y O A R B R D E Y V
K L L B Q W I K O M K
S T R E B G A D H R B
Y U I Q L X C N B J F
```

* SOLUTION CAN BE THOUGHT AS GRAPH SEARCHING
* ALSO DEPTH LIMITED SEARCH WHERE THE DEPTH IS THE LENGTH OF THE PATTERN

## EXTRA: Z - FUNCTION

↳ FOR A GIVEN STRING S, THE Z-FUNCTION IS AN ARRAY OF $n$ LENGTH WHERE THE $i$-TH ELEMENT IS EQUAL TO THE GREATEST NUMBER OF CHARS. STARTING FROM POSITION $i$ THAT COINCIDE WITH THE FIRST CHARS. OF S

- $z[i]$ : LENGTH OF THE LONGEST STRING THAT IS A PREFIX OF S AND A PREFIX OF THE SUFFIX OF S STARTING AT $i$

EXAMPLE :

"AAAAA"          " AACBACCA"              " ABACABA"              " ABCDE"
Z = [0 4 3 2 1]    Z = [0 1 0 0 1 0 0 1]    Z = [0 0 1 0 3 0 1]    Z = [0 0 0 0 0]

* GO TO TRIVIAL IMPLEMENTATION

        ↳ $O(n^2)$

## EFFICIENT ALGORITHM : $O(n)$

(IDEA) $\rightarrow$ COMPUTE $z[i]$ FROM $1, ..., n-1$ (ZERO BASED) TRYING TO USE PREVIOUSLY COMPUTED VALUES

$\hookrightarrow$ WE'LL KEEP THE $[L, R]$ INDICES OF THE RIGHTMOST SEGMENT MATCH

SUBSTRINGS THAT COINCIDE WITH $\leftarrow$
A PREFIX OF S

- $L$ : STARTING INDEX OF S.B.
- $R$ : ENDING " " " $\rightarrow$ BOUNDARY OF WHATEVER WE HAVE SCANNED STRING S

### S.M. :

$$S_0 S_1 S_2 ... S_i S_{i+1} ... S_{n-1}$$

STEGMENT MATCH

$\hookrightarrow$ LENGTH OF S.M. : $\boxed{i + z[i] - 1}$

$$S_0 S_1 = S_i S_{i+1}$$

### ALGORITHM :

TWO CASES :

① $i > R$ : CURRENT POSITION IS OUTSIDE OF SEGMENT MATCH

- COMPUTE $z[i]$ USING TRIVIAL ALGORITHM
- IF $z[i] > 0$ WE UPDATE VALUES $L, R$ SINCE LENGTH OF S.M. IS BETTER THAN PREVIOUS $R$

② $i \leq R$ : CURRENT POSITION IS INSIDE S.M. $[L, R]$

- INITIALIZE $z[i]$ TO SOME VALUE BETTER THAN ZERO

OBSERVATION :

- $S[L...R] = S[0...R-L]$ THESE SUBSTRINGS MATCH

$$S_0 S_1 S_2 ... S_L S_{L+1} ... S_R ... S_{n-1}$$
$\underbrace{\qquad}_{R-L} \qquad \underbrace{\qquad}_{R-L}$

- FOR INITIAL VALUE $z[i] = \dfrac{z[i-L]}{\uparrow}$

VALUE FOR SEGMENT MATCH $S[0...R-L]$

- $z[i-L]$ MIGHT BE TOO LARGE FOR CURRENT $i$ INDEX
  THUS :

$$z[i] = MIN \underbrace{(R-i+1}_{\downarrow}, z[i-L])$$

IF WE ARE AT INDEX
$i = R = n-1$ THEN MAX
POSSIBLE VALUE WITHOUT
EXCEEDING INDEX $R$ IS $\underline{1}$

- AFTER INITIALIZING $z[i]$ WE TRY TO INCREMENT $z[i]$ USING
  THE TRIVIAL ALGORITHM

\* GO TO Z_FUNCTION.CPP


APPLICATIONS :

- [ SEARCH THE SUBSTRING ]

  • FIND OCCURRENCES OF STRING $s$ IN TEXT $t$
  • CONCATENATE $s \# t$
  • BUILD Z- FUNCTION
  • IF $k = |s| \longrightarrow$ OCCURRENCE OF $s$ IN $t$ AT INDEX $i$ : $i \in [0,..., |t|-1]$

    $\hookrightarrow k = z[i + |s| + 1]$


- [ NUMBER OF DISTINCT SUBSTRINGS IN A STRING ]    $O(n^2)$

  • COUNT # OF DIFFERENT SUBSTRINGS OF STRING $s$       $n=|s|$
  • APPEND NEW CHAR $c$ TO $s$
  • DEFINE STRING $T = c + s$   (CONCATENATION)
  • REVERSE T
  • COMPUTE Z FUNCTION OF T
  • GET MAX ELEMENT IN $z \rightarrow z_{max}$
  • ADD $|T| - z_{max}$ TO $k$

    $\hookrightarrow k$ INITIALLY AS ZERO

      $\hookrightarrow$ STORES # OF DISTINCT SUBSTRINGS


  [EXAMPLE ]
  $\downarrow$

$$ABC \quad \longrightarrow \left| \{ A, B, C, AB, AC, ABC \} \right| = k = 6$$

| ABCA | ABCAB | ABCABC |
|---|---|---|
| ACBA | BACBA | CBACBA |
| $z = [0001]$ | $z = [00021]$ | $z = [000321]$ |
| $z_{MAX} = 1$ | $z_{MAX} = 2$ | $z_{MAX} = 3$ |

$$4 - 1 = 3 \qquad 4 - 2 = 2 \qquad 4 - 3 = 1$$

$$\underbrace{\qquad\qquad\qquad\qquad}$$

$$\boxed{k = 6}$$