- Given a text
- Consider all suffixes
- Sort them alphabetically

$T =$ mississippi#

$P =$ si

1) Suffixes that share a prefix are __contiguos__

2) Pattern is "positioned" before the occurrencies
$\Rightarrow$ just scan to find the occurrencies

```
#
i#
i ppi#
issippi#
ississippi#
mississippi#
pi#
ppi#
sippi#
sissippi#
ssippi#
ssissippi#
```

$P = \boxed{si}$

$\Rightarrow$ __Suffix array__ : STORES for every suffix the offset of the 1st character

```
   1 2 3 4 5 6 7 8 9 10 11 12
T = M i s s i s s i p  p  i  #
```

length of the text strings

$SA = \Theta (N \log \sigma)$ BITS

$N = $ # characters of the text +

SA
| 12 |
| 11 |
| 8 |
| 5 |
| 2 |
| 1 |
| 10 |
| 9 |
| 7 |
| 4 |
| 6 |
| 3 |

- SEARCH A PATTERN on SA

=> Binary search          $P = Si$

  ① => pick the middle element of SA

```
SA
12
11
8
5
2
1  ----------▷
10
9
7
4·  ----------▷ COMPARE
6
3
```

$1 = \text{MISSISSIPPI}$   ② COMPARE WITH $P = Si$

        => $S > M$

           => P will be in the second part of the array

③ REPEAT

        => ① pick middle element

=> I do $\boxed{\log_2 n \text{ steps}}$

      ∀ step I compare ⑫ characters at most

=> Cost of search $= O(p \log_2 m)$

(OCCURRENCE)

=> when I've found position of Ⓟ

        => SCAN the SA     => BUT is and compare   COSTLY

=> Search for     P #
             and
              P $       where $\# < \Sigma < \$$
                                 as char of alphabet

$\Rightarrow$ 2 BINARY SEARCHES TO FIND BEGINNING and the

end $\Rightarrow$ to know # of occurrences

$\Rightarrow O(\underbrace{P \cdot \log m}_{\substack{\text{only to} \\ \text{count}}} + \underbrace{\text{occurrences}}_{\substack{\text{TO COUNT out} \\ \text{locate the} \\ \text{occurrences}}})$ time

$\Rightarrow$ Counting occurrencies $O(p \cdot \log m)$ time

Retrival $O(p \cdot \log_2 m + occ)$ time

---

Building SA (char *T, int *m, char **SA)

```
for (i = 0; i < m; i++) {
        SA[i] = T + i;
    }
Qsort ( SA, m, sizeof (char *), suffix_cmp)
```

$-\circ-$

suffix_cmp ( char **p, char **q) {

Return  str_cmp (*p, *q) }

$\underset{\nearrow 0 \text{ scan}}{\cup} \left( \frac{M}{B} \cdot m \log m \right) \quad I/0$

$\cup \left( m \cdot m \log m \right)$

# (LCP)-ARRAY

└─ longest common prefix

(total number of characters shared the longest prefix
by adjacent suffixes)

$$T = \overset{1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9\ 10\ 11\ 12}{M\ I\ S\ S\ I\ S\ S\ I\ P\ P\ I\ \#}$$

| LCP | SA | |
|-----|-----|-----|
| 0 | 12 | # |
| 0 | 11 | i # |
| 1 | 8 | i ppi # |
| 4 | 5 | issippi # |
| 0 | 2 | ississippi # |
| 0 | 1 | mississippi # |
| 0 | 10 | pi # |
| 1 | 9 | ppi # |
| 0 | 7 | sippi # |
| 2 | 4 | sissippi # |
| 1 | 6 | ssippi # |
| 3 | 3 | ssissippi # |

P = $S_i$

max (bracket over LCP values 1,4,0)

$\uparrow$ n−1 position

$LCP(5,3) = 0 = min(4, 0, 0) = 0$

```
Build SUFFIX-ARRAY (char *T, char *m, char **SA){
    for(i=0; i<m; i++){
        SA[i]=T; }
        QSORT (SA, m, sizeof(char *), suffix_cmp)
    Suffix_cmp (char **P, char **Q)
        return string_cmp (*P, *P)
    cost    O( M/B  u log u )   ✓
```

```
Buildy SA ( char *T, char *m, char &&st)
        for (i=0; i<m; i++){
                SA[i] = T + i; }
        QSORT ( SA, m, sizeof(char*), suffix-cmp)
        }

    suffix-cmp ( char *p, char *q) {
            return stry-cmp (*p, *q)

    => Complexity O( m/B  m lg m
```

SNOW_PLOW

- Require U array of unsorted items
- H = build a min-heap over U-items
- set M = ∅

while ( M ≠ 0)

        min = extract the minimum item from H
        next =: read next item from U
        if (next < min)

                put next in U
        else
                put next in the heap

        end if
end while

$$O\left(\frac{M}{B} \ \textbf{log}_{2} \ \frac{M}{2M}\right)$$

=> Snow plow produces sorted runs of avgt 2M

        => long with ed if ( ~~min~~ $\overset{next}{}$ < min)
        end assuming RANDOM DISTRIBUTION ~~over~~ fetching items

                $$\left((k-M) + M = k\right)$$

=> M items ends up in U      P = $\frac{1}{2}$ => I red k items

$$M = \frac{K}{2} \qquad \boxed{K = 2M}$$

```
Build SA (char *T, char *m, char **SA) {
    for (i=0; j<i; i++) {
        SA[i] = T+i;
        Qsort (char SA, char m, sizeof (char *), suffix_cmp)
    }
}

suffix_cmp (char **p, char **q)
    return strij_cmp (*p, *q)
```

$$O\left( \frac{M}{B} \circ M \lg_2 n \right)$$