

CSE 566 Spring 2023

Searching & Constructing Suffix Array

Instructor: Mingfu Shao

Suffix Array

$S = \underline{\text{banana}}\$$

$SA(s) = (7, 6, 4, 2, 1, 5, 3)$

$\$ < a < b < n$

1 banana\$

2 anana\$

3 nana\$

4 ana\$

5 na\$

6 a\$

7 \$

7 \$

6 a\$

4 ana\$

2 anana\$

1 banana\$

5 na\$

3 nana\$

Longest Common Prefix (LCP) Array

$S = \text{banana\$}$

$SA(s) = (7, 6, 4, 2, 1, 5, 3)$

$LCP(s) = (0, 1, 3, 0, 0, 2)$

1 banana\$

2 anana\$

3 nana\$

4 ana\$

5 na\$

6 a\$

7 \$

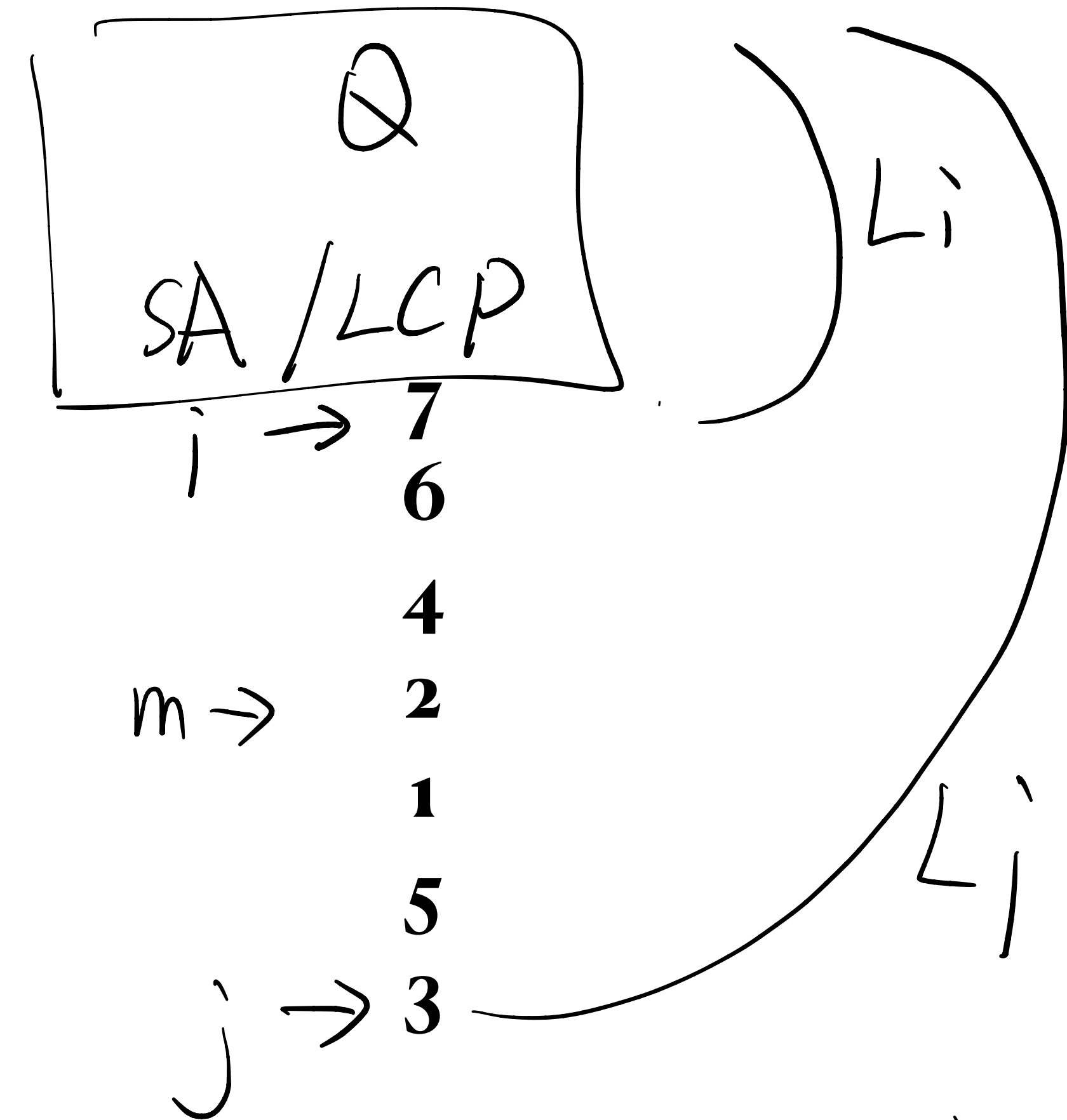
| | |
|----------|-------------------|
| 0 | 7 \$ |
| 1 | 6 a\$ |
| 3 | 4 ana\$ |
| 0 | 2 anana\$ |
| 0 | 1 banana\$ |
| 2 | 5 na\$ |
| | 3 nana\$ |

Searching for a query

- Searching questions:
 - Decide if query Q is the substring of S .
 - Find the longest common substring (LCS) of S and Q such that the LCS starts from $Q[1]$.
- Essentially: find the position of query Q in the sorted list of suffixes
- Key idea: binary search!
- Naive implementation: $O(|q| \log |S|)$.

Faster Search Algorithm

- Check if Q is less than $SA[1]$, or Q is larger than $SA[n]$
- Init: $i = 1$ and $j = n$
- Compute $L_i := \text{LCP}(SA[i], Q)$ and $L_j := \text{LCP}(SA[j], Q)$
- FUNCTION **BS**(i, j, L_i, L_j)
 - Let $m = (i + j) / 2$
 - Case 1: $L_i = L_j$:
 - Case 2: $L_i > L_j$:
 - Case 3: $L_i < L_j$; //symmetric to Case 2
- END FUNCTION

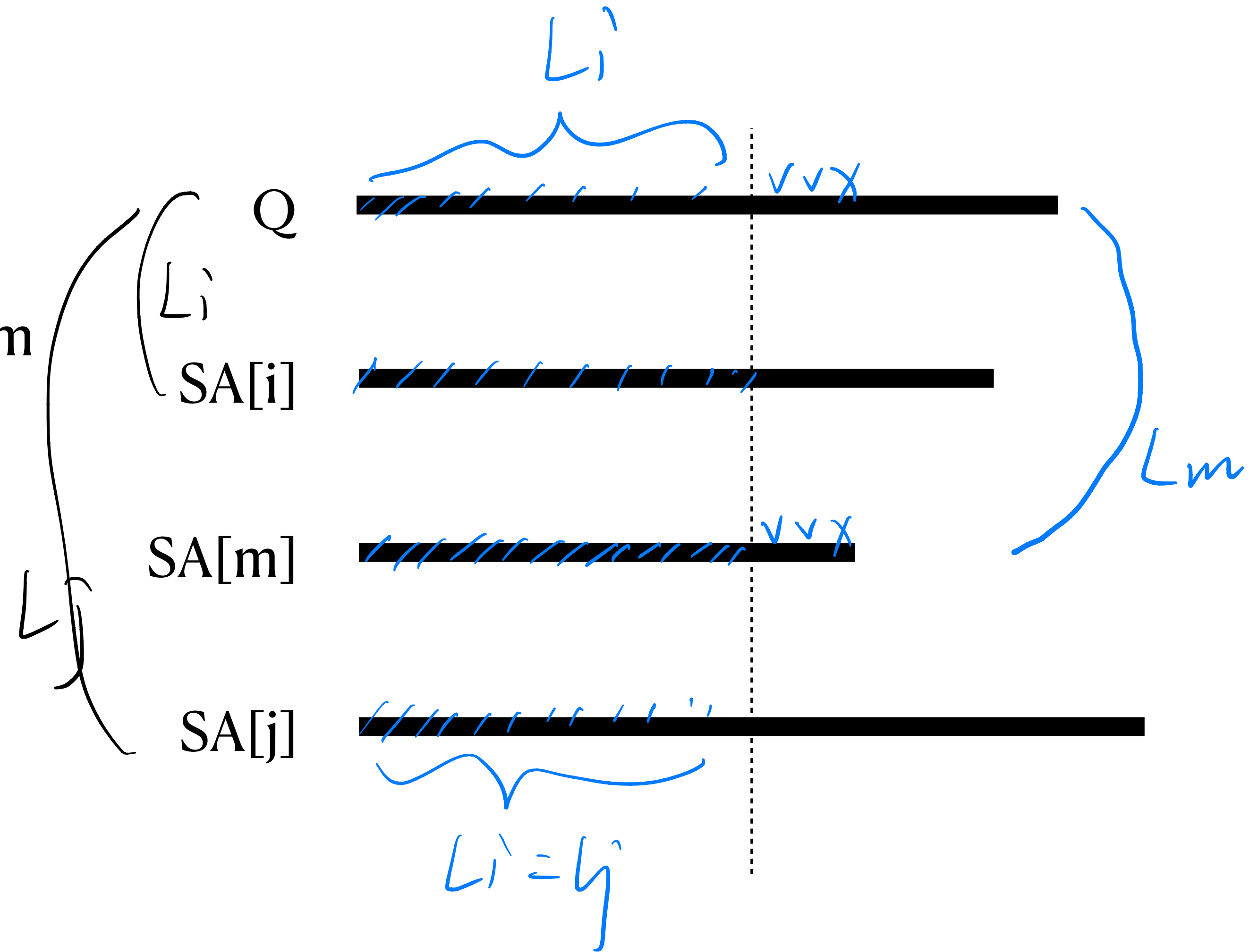


Running Time:

$$\underline{O(\log |S| + O(|Q|))}$$

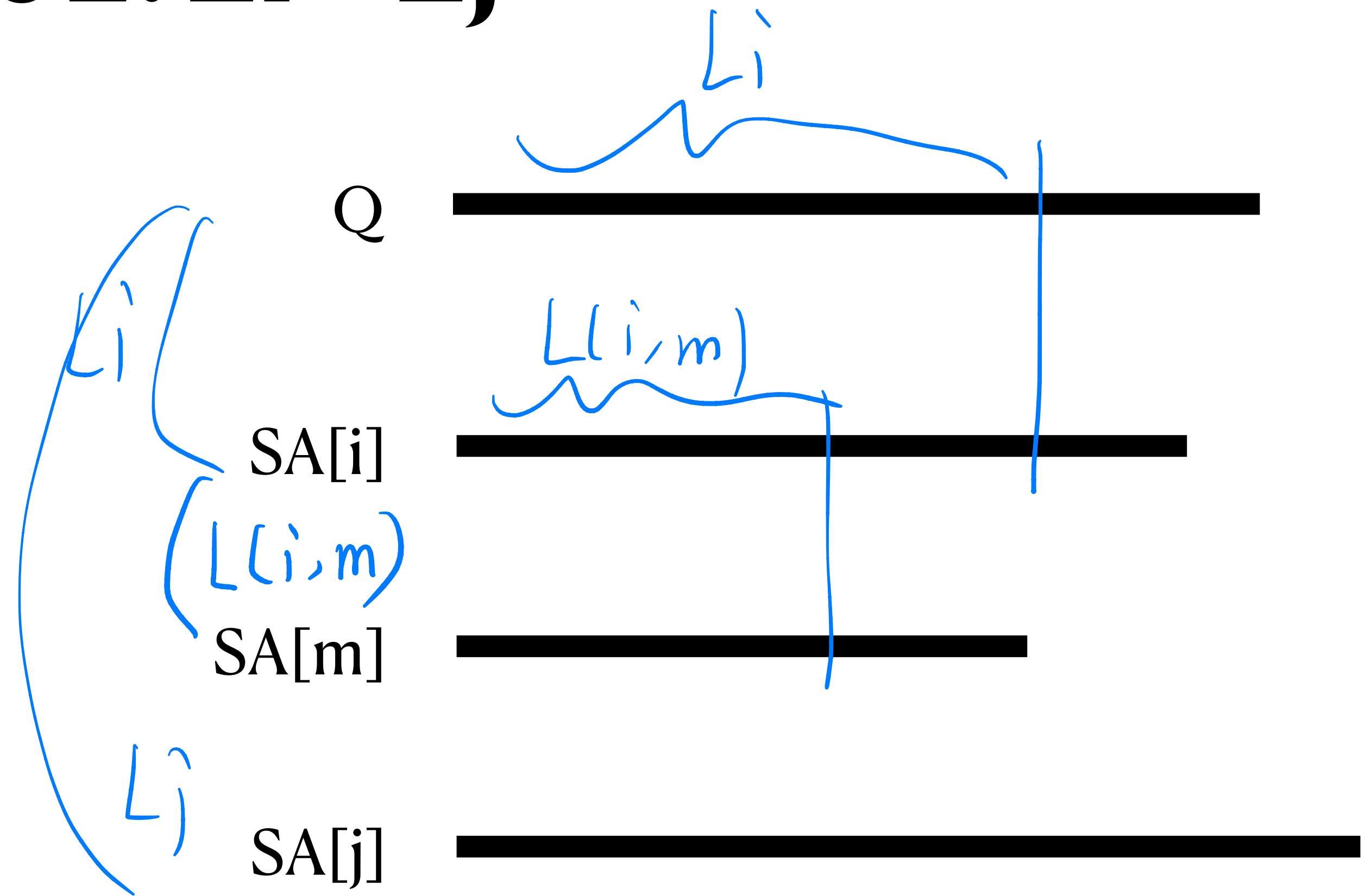
Case 1: $L_i = L_j$

- **Fact:** $L_m := \text{LCP}(\text{SA}[m], Q) \geq L_i = L_j$
- PROCEDURE:
 - Compare Q and $\text{SA}[m]$ starting from position $L_i + 1 \Rightarrow \underline{L_m}$
 - If Q gets exhausted: return m
 - If $Q < \text{SA}[m]$: $\text{BS}(i, m, L_i, L_m)$
 - If $Q > \text{SA}[m]$: $\text{BS}(m, j, L_m, L_j)$
- END PROCEDURE



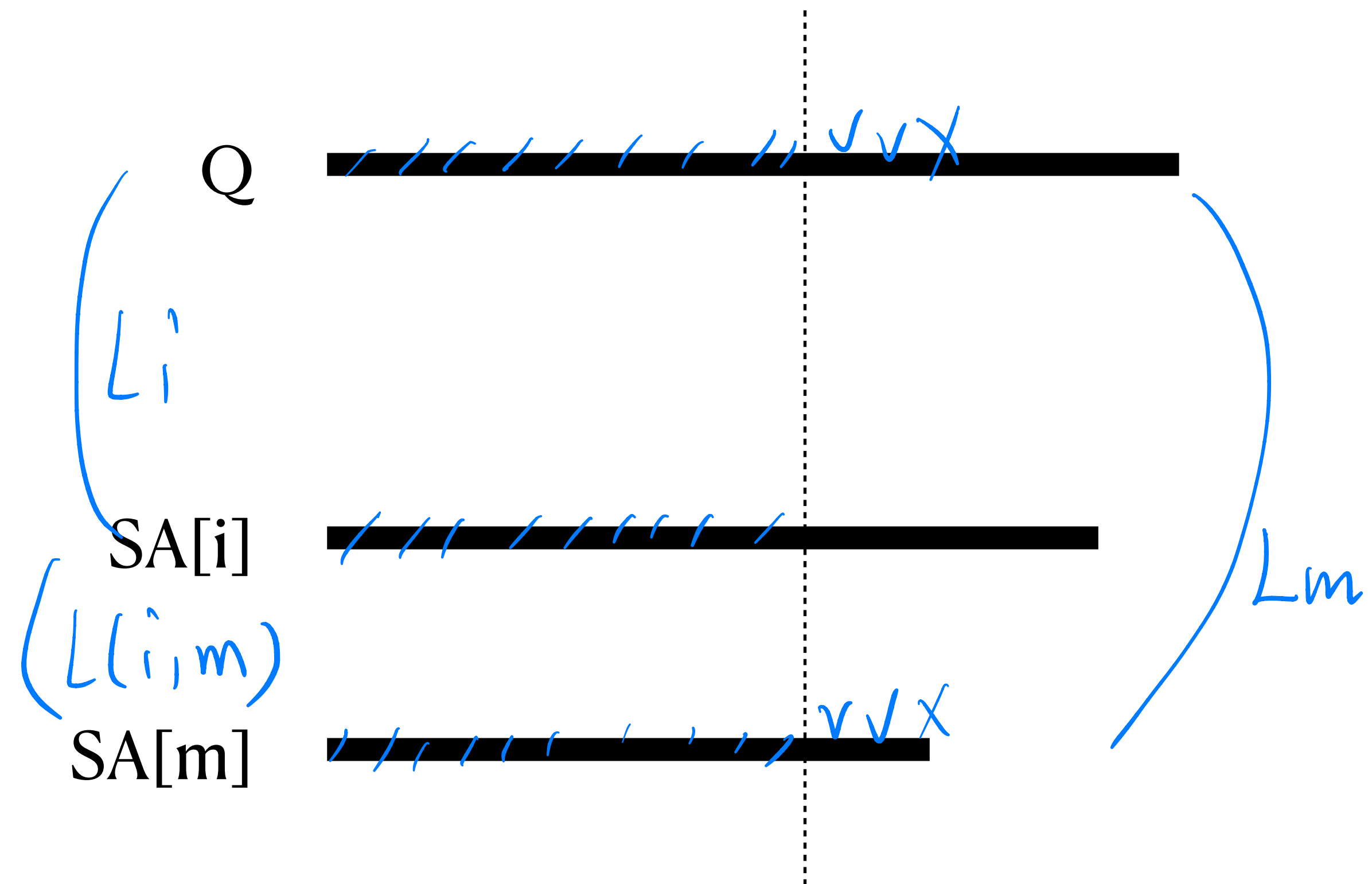
Case 2: $L_i > L_j$

- Query $L(i,m) := \text{LCP}(\text{SA}[i], \text{SA}[m])$
 - Can be done in constant time!
- Case 2a: $L(i, m) = L_i$
- Case 2b: $L(i, m) < L_i$
- Case 2c: $L(i, m) > L_i$



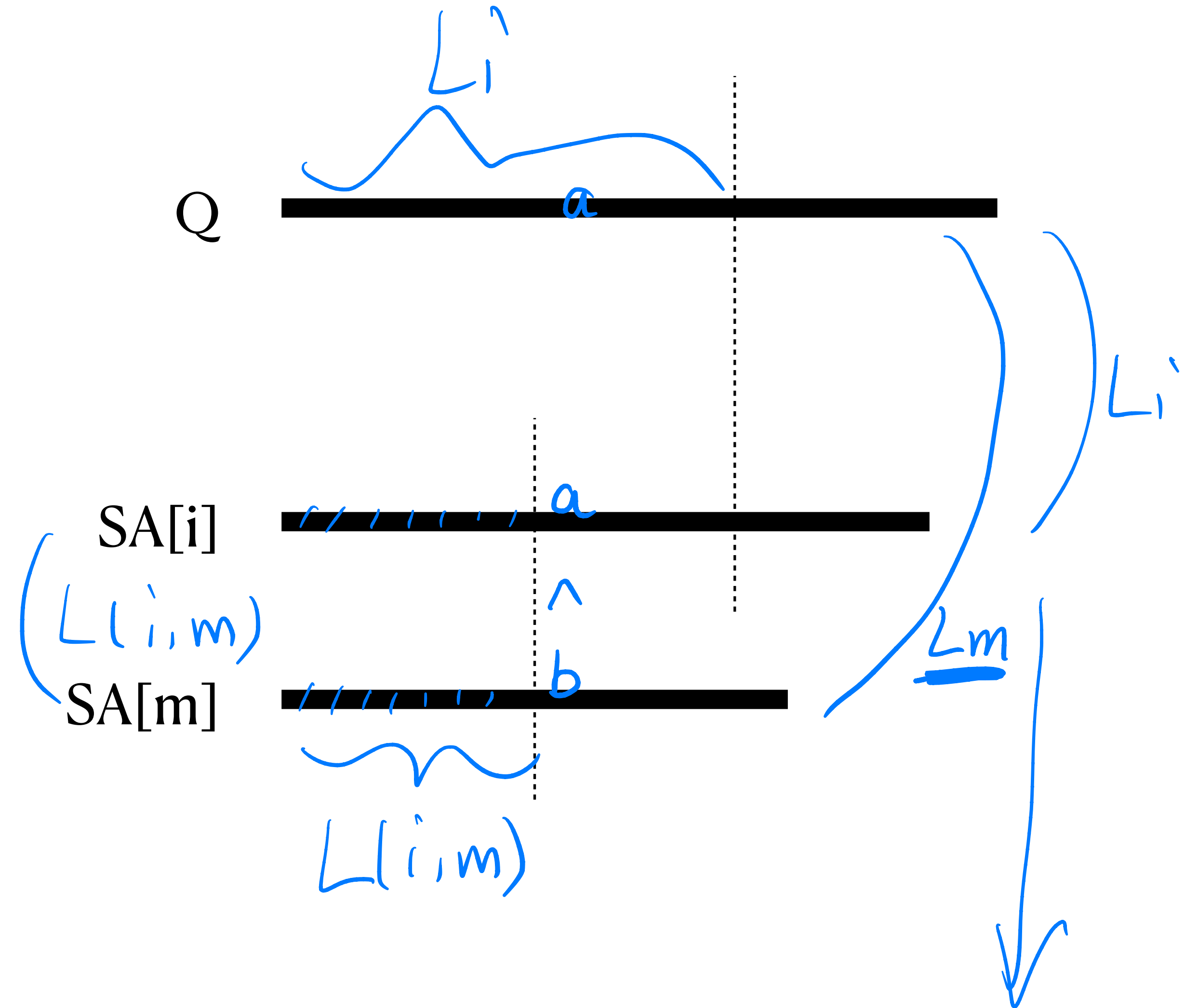
Case 2a: $L_i > L_j$ and $L(i, m) = L_i$

- **Fact:** $L_m \geq L_i$
- **PROCEDURE:**
 - Compare Q and SA[m] starting from position $L_i + 1$; this calculates L_m
 - If Q gets exhausted: return m
 - If $Q < SA[m]$: BS(i, m, L_i , L_m)
 - If $Q > SA[m]$: BS(m, j, L_m , L_j)
- **END PROCEDURE**



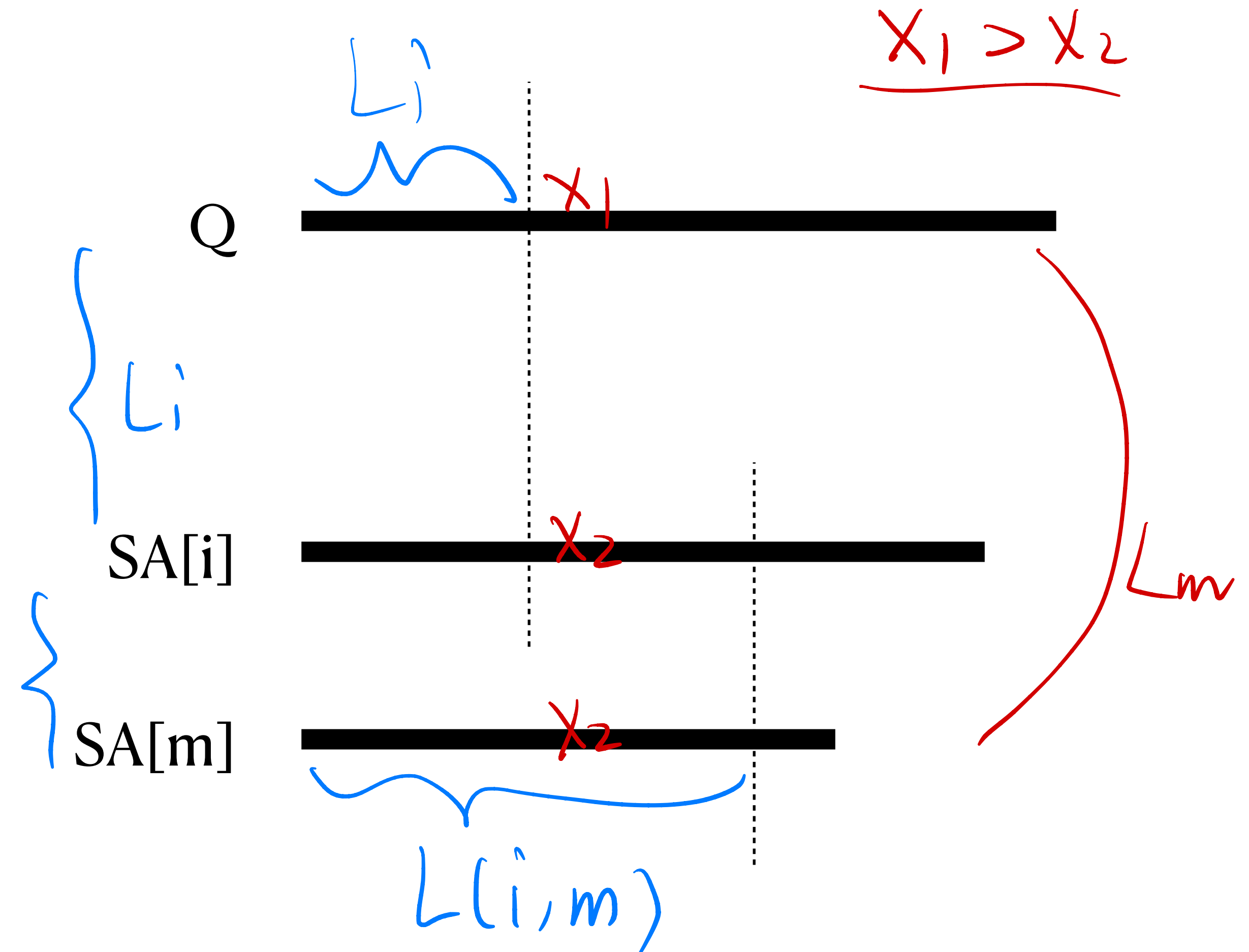
Case 2b: $Li > Lj$ and $L(i, m) < Li$

- **Fact:** $Q < SA[m]$
- **Fact:** $Lm = L(i, m)$
- PROCEDURE:
 - $BS(i, m, Li, \underline{Lm})$
- END PROCEDURE



Case 2c: $L_i > L_j$ and $L(i, m) > L_i$

- **Fact:** $Q > SA[m]$
- **Fact:** $L_m = L_i$
- **PROCEDURE:**
 - $BS(m, j, L_m, L_j)$
- **END PROCEDURE**



Fast Querying LCP

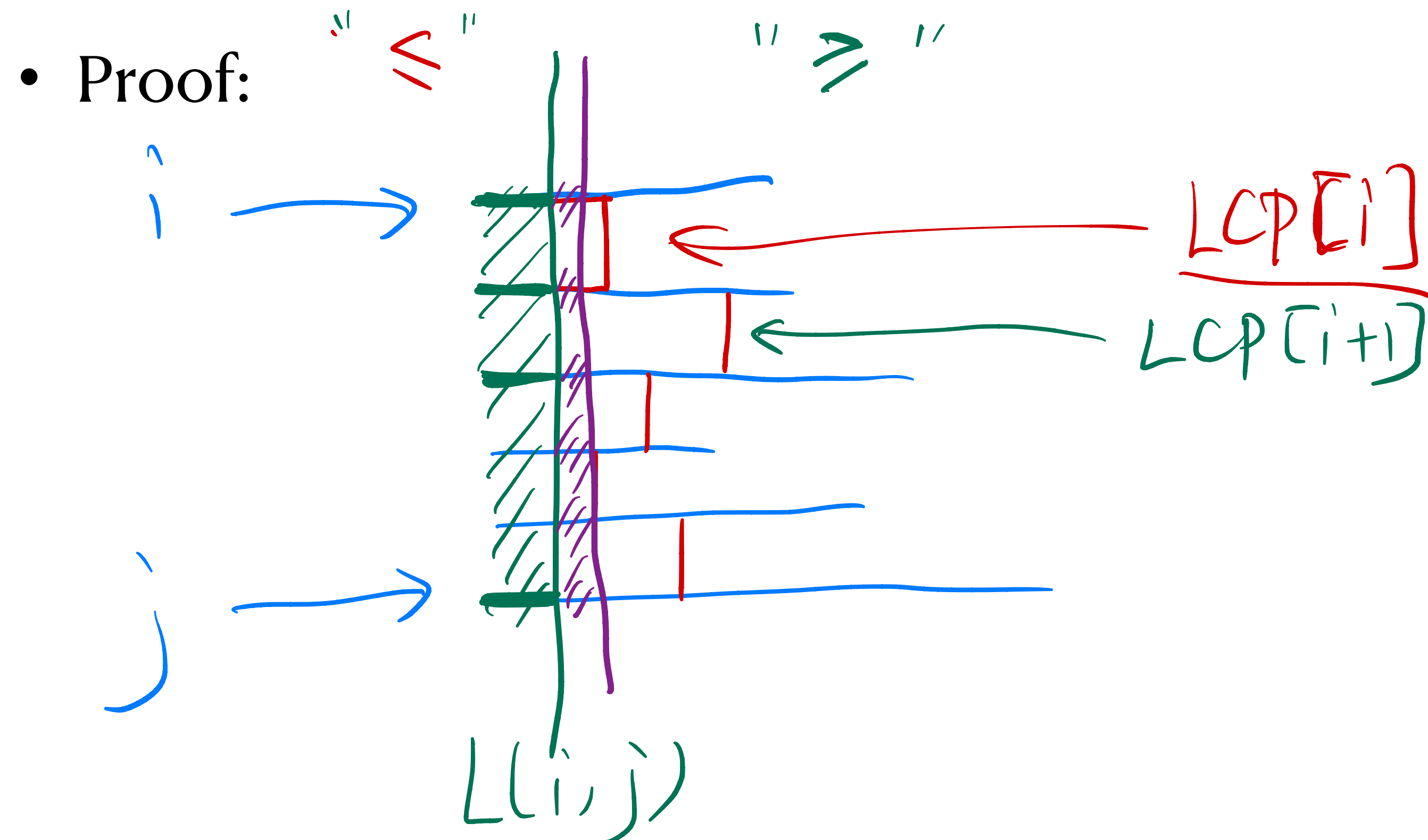
$$LCP(SA[i], SA[j])$$

- **Question:** given LCP array, query $L(i, j)$ in constant time, for arbitrary i and j . $i < j$

$$SA(s) = (7, 6, 4, 2, 1, 5, 3)$$

$$LCP(s) = (0, 1, 3, 0, 0, 2)$$

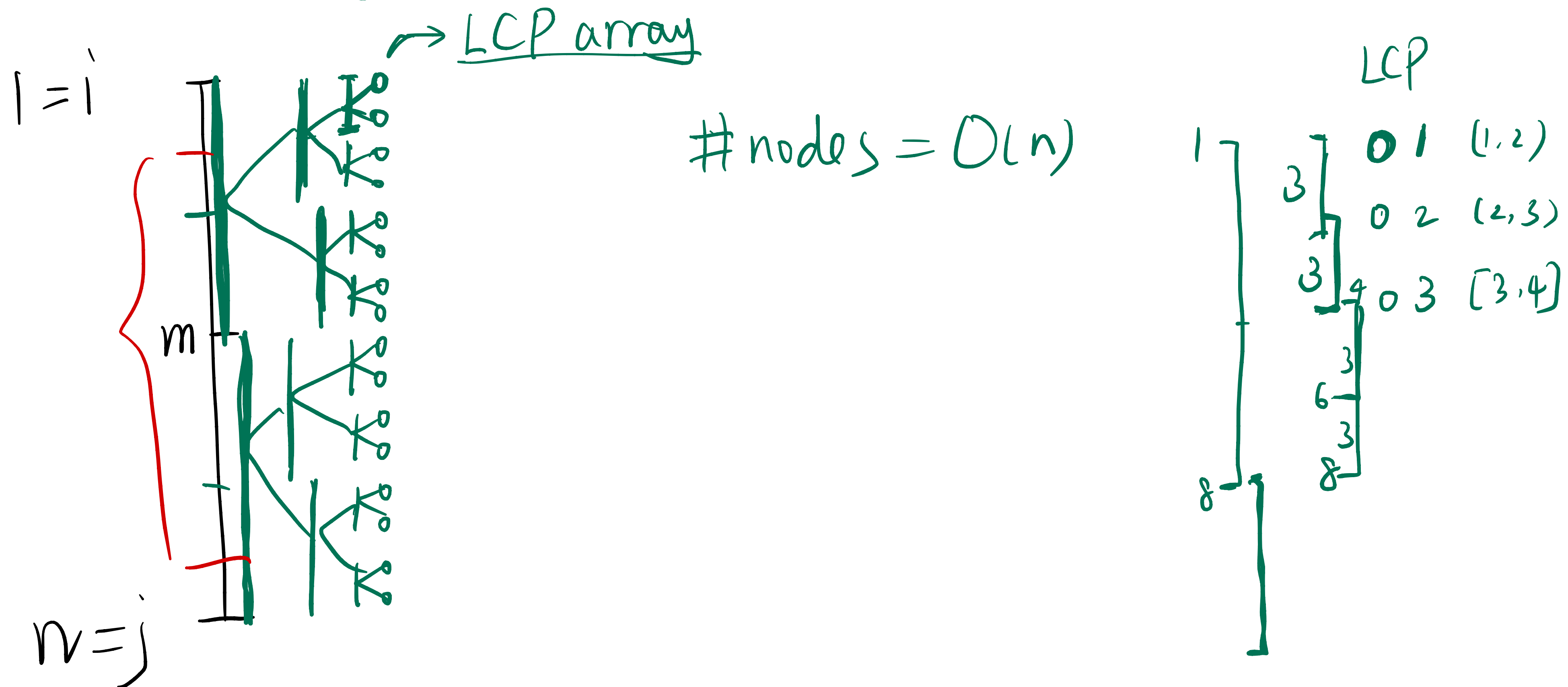
- **Fact:** $L(i, j) = \min_{\{i \leq k < j\}} LCP[k]$.



| | |
|---|---------------|
| 7 | \$ |
| 0 | 6 a\$ ← i |
| 1 | 4 ana\$ |
| 3 | 2 anana\$ ← j |
| 0 | 1 banana\$ |
| 0 | 5 na\$ |
| 2 | 3 nana\$ ← j |

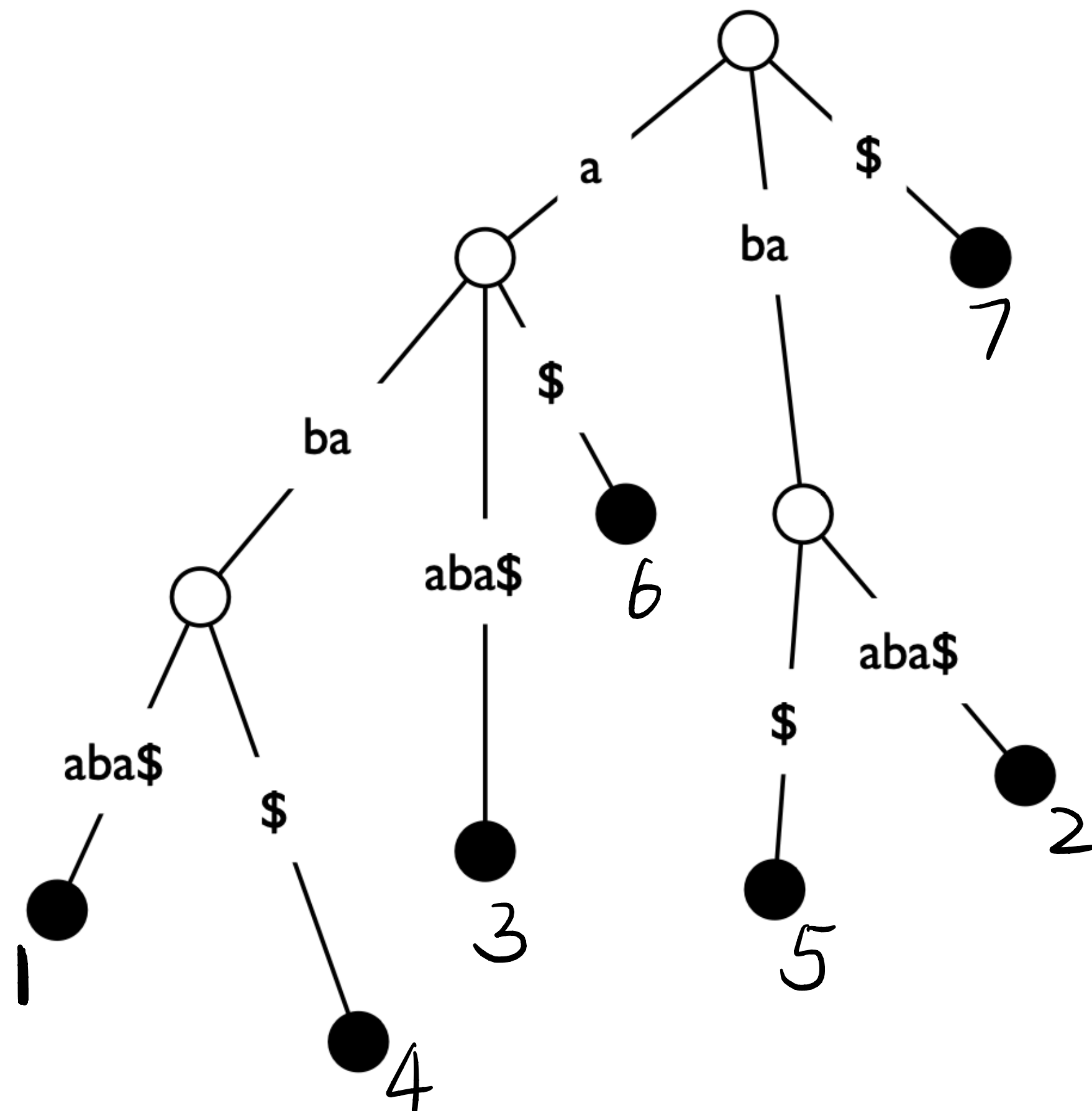
Algorithms

- General algorithm for range-min-query-problem. $O(n)$ space, $O(1)$ time $[i, j]$
for each query
- **Key Observation:** not every pair of query $L(i, j)$ will be used!



Constructing Suffix Array using Suffix Tree

$T = \underline{\text{abaaba\$}}$



$O(|T|)$:
traversal the edges out of
a vertex following the
alphabetical order.

Linear-Time Algorithm

Example: $S = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 \\ a & b & a & b & b & a & b & a & b & b & a & \$ \end{matrix}$

$\begin{matrix} 1 & 2 & 3 & 1 & 2 & 3 & 1 & 2 & 3 & 1 & 2 & 3 \end{matrix}$

$n = |S|$

$\text{DQC} \quad \frac{2}{3}$

$\frac{1}{3}$

Step 1: G_1

G_2

G_3

Step 2: $T = \begin{matrix} \text{token} \\ \downarrow \\ b & a & b & b & a & b & a & \$ & \$ & a & b & b & a & b & b & a & \$ & \$ & \$ \end{matrix}$

$\rightarrow X: \bar{E}$

$|T| = 2n$

$\rightarrow |X| = \frac{2n}{3}$