

# ALGORITHMS OF BIOINFORMATICS

## Googling Genomes

*15 January 2026*

Prof. Dr. Sebastian Wild

# Outline

## 7 Googling Genomes

- 7.1 Range-Minimum Queries
- 7.2 RMQ – Sparse Table Solution
- 7.3 RMQ – Cartesian Trees
- 7.4 String Matching in Enhanced Suffix Array
- 7.5 The Burrows-Wheeler Transform
- 7.6 Inverting the BWT
- 7.7 Searching in the BWT

# Recall Unit 6

## Application 4: Longest Common Extensions

- We implicitly used a special case of a more general, versatile idea:

The *longest common extension (LCE)* data structure:

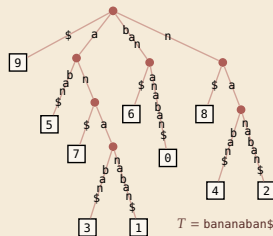
- **Given:** String  $T[0..n)$
- **Goal:** Answer LCE queries, i.e.,  
 given positions  $i, j$  in  $T$ ,  
 how far can we read the same text from there?  
 formally:  $\text{LCE}(i, j) = \max\{\ell : T[i..i + \ell) = T[j..j + \ell)\}$

⇒ use suffix tree of  $T$ !

(length of) longest common prefix  
of  $i$ th and  $j$ th suffix

- In  $\mathcal{T}$ :  $\text{LCE}(i, j) = \text{LCP}(T_i, T_j) \rightsquigarrow$  same thing, different name!  
 $=$  string depth of  
*lowest common ancestor (LCA)* of  
 leaves  $\boxed{i}$  and  $\boxed{j}$

- in short:  $\text{LCE}(i, j) = \text{LCP}(T_i, T_j) = \text{stringDepth}(\text{LCA}(\boxed{i}, \boxed{j}))$



# Recall Unit 6

## Efficient LCA

How to find lowest common ancestors?

- ▶ Could walk up the tree to find LCA  $\rightsquigarrow \Theta(n)$  worst case 🗑️
- ▶ Could store all LCAs in big table  $\rightsquigarrow \Theta(n^2)$  space and preprocessing 🗑️



**Amazing result:** Can compute data structure in  $\Theta(n)$  time and space that finds any LCA in **constant(!) time**.

- ▶ a bit tricky to understand
- ▶ but a theoretical breakthrough
- ▶ and useful in practice

and suffix tree construction inside ...



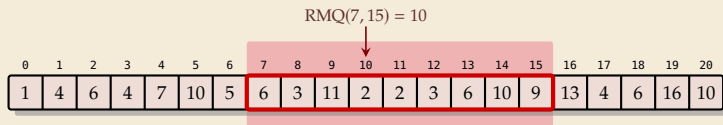
$\rightsquigarrow$  for now, use  $O(1)$  LCA as black box.

$\rightsquigarrow$  After linear preprocessing (time & space), we can find LCEs in  $O(1)$  time.

## 7.1 Range-Minimum Queries

# Range-minimum queries (RMQ)

- ▶ **Given:** Static array  $A[0..n)$  of numbers  
*array/numbers don't change*
- ▶ **Goal:** Find minimum in a range;  
 $A$  known in advance and can be preprocessed



- ▶ **Nitpicks:**
  - ▶ Report *index* of minimum, not its value
  - ▶ Report *leftmost* position in case of ties

# Finally: Longest common extensions

- In Unit 6: Left question open how to compute LCA in suffix trees
- But: Enhanced Suffix Array makes life easier!

$$\text{LCE}(i, j) = \text{LCP}[\text{RMQ}_{\text{LCP}}(\min\{R[i], R[j]\} + 1, \max\{R[i], R[j]\})]$$

## Inverse suffix array: going left & right

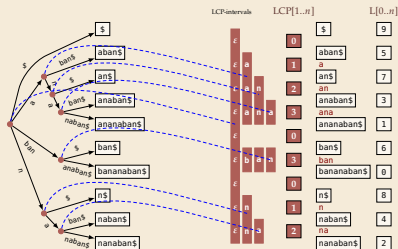
- to understand the fastest algorithm, it is helpful to define the *inverse suffix array*:

- $R[i] = r \iff L[r] = i$   $L = \text{leaf array}$
- $\iff$  there are  $r$  suffixes that come before  $T_i$  in sorted order
- $\iff T_i$  has (0-based) *rank*  $r \rightsquigarrow$  call  $R[0..n]$  the *rank array*

$i$	$R[i]$	$T_i$		$r$	$L[r]$	$T_{L[r]}$
0	6 <sup>(1)</sup>	bananaban\$		0	9	\$
1	4 <sup>(2)</sup>	ananaban\$		1	5	abans\$
2	9 <sup>(3)</sup>	nanaban\$		2	7	ans\$
3	3 <sup>(4)</sup>	anaban\$		3	3	anaban\$
4	8 <sup>(5)</sup>	naban\$		4	1	ananaban\$
5	1 <sup>(6)</sup>	aban\$		5	6	ban\$
6	5 <sup>(7)</sup>	ban\$		6	0	bananaban\$
7	2 <sup>(8)</sup>	an\$		7	8	n\$
8	7 <sup>(9)</sup>	n\$		8	4	naban\$
9	0 <sup>(10)</sup>	\$		9	2	nanaban\$

sort suffixes

## LCP array and internal nodes



$\rightsquigarrow$  Leaf array  $L[0..n]$  plus LCP array  $LCP[1..n]$  encode full tree!

# Rules of the Game

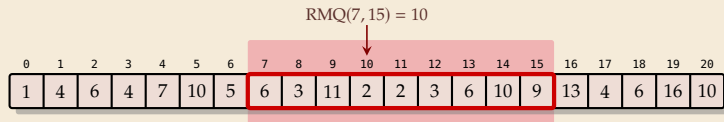
- ▶ For the following, consider RMQ on arbitrary arrays
- ▶ comparison-based  $\rightsquigarrow$  values don't matter, only relative order
- ▶ Two main quantities of interest:
  1. **Preprocessing time:** Running time  $P(n)$  of the preprocessing step  $\rightsquigarrow$  space usage  $\leq P(n)$
  2. **Query time:** Running time  $Q(n)$  of one query (using precomputed data)
- ▶ Write  $\langle P(n), Q(n) \rangle$  **time solution** for short

## RMQ Implications for LCE

- ▶ Recall: Can compute (inverse) suffix array and LCP array in  $O(n)$  time
- $\rightsquigarrow \langle P(n), Q(n) \rangle$  time RMQ data structure implies  
 $\langle P(n) + O(n), Q(n) \rangle$  time LCE data structure



# Trivial Solutions



- ▶ Two easy solutions show extreme ends of scale:

## 1. Scan on demand

- ▶ no preprocessing at all
  - ▶ answer  $\text{RMQ}(i, j)$  by scanning through  $A[i..j]$ , keeping track of min
- $\rightsquigarrow \langle O(1), O(n) \rangle$

## 2. Precompute all

- ▶ Precompute all answers in a big 2D array  $M[0..n][0..n]$
- ▶ queries simple:  $\text{RMQ}(i, j) = M[i][j]$

$\rightsquigarrow \langle O(n^3), O(1) \rangle$

- ▶ Preprocessing can reuse partial results  $\rightsquigarrow \langle O(n^2), O(1) \rangle$

## 7.2 RMQ – Sparse Table Solution

# Sparse Table

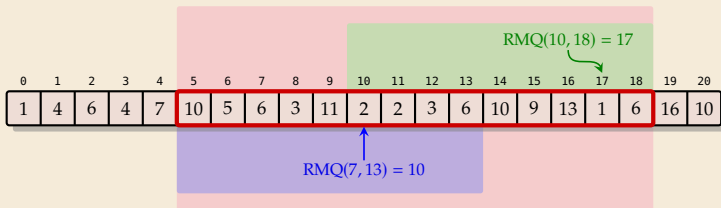
► **Idea:** Like “precompute-all”, but keep only *some* entries

► store  $M[i][j]$  iff  $\ell = j - i + 1$  is  $2^k$ .

↪  $\leq n \cdot \lg n$  entries

↪ Can be stored as  $M'[i][k]$

► How to answer queries?



1. Find  $k$  with  $\ell/2 \leq 2^k \leq \ell$

2. Cover range  $[i..j]$  by  
 $2^k$  positions right from  $i$  and  
 $2^k$  positions left from  $j$

3.  $RMQ(i, j) =$   
 $\arg \min\{A[rmq_1], A[rmq_2]\}$

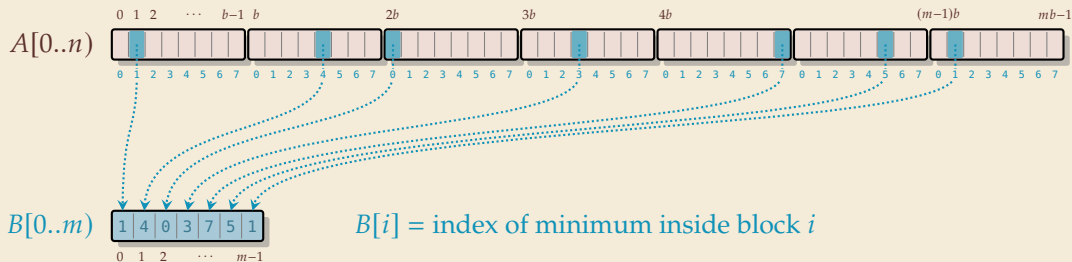
with  $rmq_1 = RMQ(i, i + 2^k - 1)$   
 $rmq_2 = RMQ(j - 2^k + 1, j)$

► Preprocessing can be done in  $O(n \log n)$  times

↪  $\langle O(n \log n), O(1) \rangle$  time solution!

# Bootstrapping

- ▶ We know a  $\langle O(n \log n), O(1) \rangle$  time solution
- ▶ If we use that for  $m = \Theta(n/\log n)$  elements,  $O(m \log m) = O(n)$ !
- ▶ Break  $A$  into blocks of  $b = O(\log n)$  numbers
- ▶ Create array of block minima  $B[0..m]$  for  $m = \lceil n/b \rceil = O(n/\log n)$

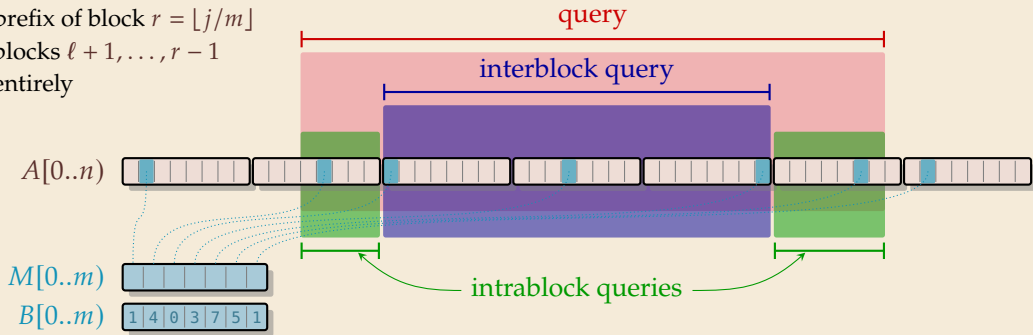


$\rightsquigarrow$  Use sparse table solution for  $B$ .

$\rightsquigarrow$  Can solve RMQs in  $B[0..m]$  in  $\langle O(n), O(1) \rangle$  time

# Query decomposition

- ▶ Query  $\text{RMQ}_A(i, j)$  covers
  - ▶ suffix of block  $\ell = \lfloor i/m \rfloor$
  - ▶ prefix of block  $r = \lfloor j/m \rfloor$
  - ▶ blocks  $\ell + 1, \dots, r - 1$  entirely

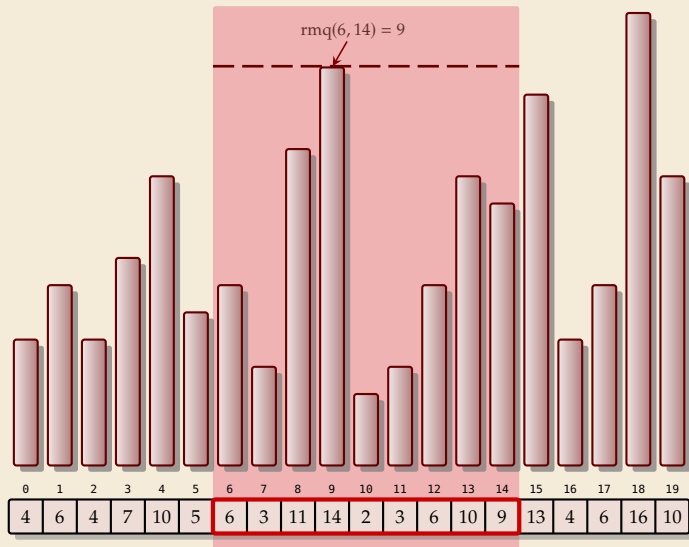


▶  $\text{RMQ}_A(i, j) = \arg \min_{k \in K} A[k]$  with  $K = \left\{ \begin{array}{l} \text{RMQ}_{\text{block } \ell}(i - \ell b, (\ell + 1)b - 1), \\ b \cdot \text{RMQ}_M(\ell + 1, r - 1) + \\ \quad B[\text{RMQ}_M(\ell + 1, r - 1)], \\ \text{RMQ}_{\text{block } r}(rb, j - rb) \end{array} \right\}$

⇒ only 3 possible values to check  
if **intrablock** and **interblock** queries known ✓

## 7.3 RMQ – Cartesian Trees

# RMQ & LCA

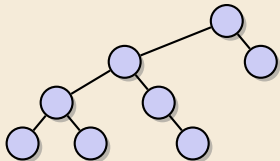


- **Range-max queries** on array  $A$ :  
$$\text{rmq}_A(i, j) = \arg \max_{i \leq k \leq j} A[k]$$
  
$$= \text{index of max}$$
- **Task:** Preprocess  $A$ ,  
then answer RMQs fast  
ideally constant time!

- ▶ **Range-max queries** on array  $A$ :  
$$\text{rmq}_A(i, j) = \arg \max_{i \leq k \leq j} A[k]$$
$$= \text{index of max}$$
- ▶ **Task:** Preprocess  $A$ ,  
then answer RMQs fast  
ideally constant time!
- ▶ **Cartesian tree:** (cf. *treap*)  
construct binary tree by  
sweeping line down
- ▶  $\text{rmq}(i, j) =$  inorder of  
**lowest common ancestor** (LCA)  
of  $i$ th and  $j$ th node in inorder



# Counting binary trees



- ▶ Given the Cartesian tree,  
all RMQ answers are determined  
and vice versa!

- How many different Cartesian trees are there for arrays of length  $n$ ?

- known result: *Catalan numbers*  $\frac{1}{n+1} \binom{2n}{n}$

- easy to see:  $\leq 2^{2n}$

↪ many arrays will give rise to the same Cartesian tree

*Can we exploit that?*

# Intrablock queries

↪ It remains to solve the **intrablock** queries!

► Want  $\langle O(n), O(1) \rangle$  time overall

↖ must include preprocessing for all  $m = \left\lceil \frac{n}{b} \right\rceil = \Theta\left(\frac{n}{\log n}\right)$  blocks!

► Choose  $b = \left\lceil \frac{1}{4} \lg n \right\rceil$

► many blocks, but just  $b$  numbers long

↪ Cartesian tree of  $b$  elements can be encoded using  $2b = \frac{1}{2} \lg n$  bits

↪ # different Cartesian trees is  $\leq 2^{2b} = 2^{\frac{1}{2} \lg n} = \left(2^{\lg n}\right)^{1/2} = \sqrt{n}$

↪ many *equivalent* blocks!

↪ **Recall: Exhaustive-Tabulation Technique:**

1. represent each subproblem by storing its *type* (here: encoding of Cartesian tree)
2. *enumerate* all possible subproblem types and their solutions
3. use type as index in a large *lookup table*

# Exhaustive Tabulation

1. For each block, compute  $2b$  bit representation of Cartesian tree
  - ▶ can be done in linear time
2. Compute large lookup table

Block type	$i$	$j$	RMQ( $i, j$ )
⋮			
⋮			

▶  $\leq \sqrt{n}$  block types

▶  $\leq b^2$  combinations for  $i$  and  $j$

$\rightsquigarrow \Theta(\sqrt{n} \cdot \log^2 n)$  rows

▶ each row can be computed in  $O(\log n)$  time

$\rightsquigarrow$  overall preprocessing:  $O(n)$  time!

# RMQ Discussion

►  $\langle O(n), O(1) \rangle$  time solution for RMQ

$\rightsquigarrow \langle O(n), O(1) \rangle$  time solution for LCE in strings!

 optimal preprocessing and query time!

 a bit complicated

## 7.4 String Matching in Enhanced Suffix Array

## Binary searching the suffix array

Recall: Can solve the string matching problem by binary searching  $P[0..m)$  in  $L[0..n]$


- ▶ worst-case cost:  $O(\log(n) \cdot m)$  character comparisons


↪ use LCP information to speed up string comparisons


- ▶ with RMQ on LCP array can determine lcp with middle


# Conclusion

- ▶ (Enhanced) Suffix Arrays are the modern version of suffix trees
  - ▶ directly simulate suffix tree operations on  $L$  and LCP arrays

 can be harder to reason about

 can support same algorithms as suffix trees

 but use much less space

 simple(r) linear-time construction

## Outlook:

- ▶ enhanced suffix arrays still need original text  $T$  to work
- ▶ a *self-index* avoids that
  - ▶ can store  $T$  in *compressed* form **and** support operations like string matching

## **7.5 The Burrows-Wheeler Transform**



# Digression: Recall BWT

## Burrows-Wheeler Transform

1. Take all cyclic shifts of  $S$
2. Sort cyclic shifts
3. Extract last column

$S = \text{alf\_eats\_alfalfa\$}$

$B = \text{asff\$f\_e\_lllaaata}$

alf\_eats\_alfalfa\$  
lf\_eats\_alfalfa\$  
f\_eats\_alfalfa\$  
\_eats\_alfalfa\$  
eats\_alfalfa\$  
ats\_alfalfa\$  
ts\_alfalfa\$  
s\_alfalfa\$  
\_alfalfa\$  
alfalfa\$  
lfalfa\$  
falffa\$  
alfa\$  
lfa\$  
fa\$  
a\$  
\$alf\_eats\_alfalfa

sort

\$alf\_eats\_alfalf**a**  
\_alfalfa\$alf\_eat**s**  
\_eats\_alfalfa\$alf**f**  
a\$alf\_eats\_alfalf**f**  
alf\_eats\_alfalfa\$**f**  
alfalfa\$alf\_eats\_alf**f**  
alfalfa\$alf\_eats\_**e**  
ats\_alfalfa\$alf\_**t**  
eats\_alfalfa\$alf\_**t**  
f\_eats\_alfalfa\$a**l**  
fa\$alf\_eats\_alfalf**a**  
falffa\$alf\_eats\_alf**a**  
lf\_eats\_alfalfa\$a**a**  
lfa\$alf\_eats\_alf**a**  
lfalfa\$alf\_eats\_**a**  
s\_alfalfa\$alf\_eat**t**  
ts\_alfalfa\$alf\_ea**a**

BWT  
↓

# Digression: Computing the BWT

*How can we compute the BWT of a text efficiently?*

- ▶ cyclic shifts  $S \hat{=}$  suffixes of  $S$ 
  - ▶ comparing cyclic shifts stops at first \$
  - ▶ for comparisons, anything after \$ irrelevant!
- ▶ BWT is essentially suffix sorting!
  - ▶  $B[i] = S[L[i] - 1]$
  - ▶ where  $L[i] = 0, B[i] = \$$

$\rightsquigarrow$  Can compute  $B$  in  $O(n)$  time from  $L$

	$r$		$\downarrow L[r]$
alf_eats_alfalfa\$	0	\$alf_eats_alfalf <b>a</b>	16
lf_eats_alfalfa\$	1	_alfalfa\$alf_eat <b>s</b>	8
f_eats_alfalfa\$	2	_eats_alfalfa\$alf <b>f</b>	3
_eats_alfalfa\$	3	a\$alf_eats_alfalf <b>f</b>	15
eats_alfalfa\$	4	alf_eats_alfalfa\$ <b>f</b>	0
ats_alfalfa\$	5	alfa\$alf_eats_alf <b>f</b>	12
ts_alfalfa\$	6	alfalfa\$alf_eats_ <b>_</b>	9
s_alfalfa\$	7	ats_alfalfa\$alf_ <b>e</b>	5
_alfalfa\$	8	eats_alfalfa\$alf_ <b>_</b>	4
alfalfa\$	9	f_eats_alfalfa\$alf <b>l</b>	2
lfalfa\$	10	fa\$alf_eats_alfalf <b>l</b>	14
falfa\$	11	falfa\$alf_eats_alf <b>l</b>	11
alfa\$	12	lf_eats_alfalfa\$ <b>a</b>	1
lfa\$	13	lfa\$alf_eats_alf <b>a</b>	13
fa\$	14	lfalfa\$alf_eats_ <b>a</b>	10
a\$	15	s_alfalfa\$alf_eat <b>t</b>	7
\$	16	ts_alfalfa\$alf_ <b>e</b> a	6

## 7.6 Inverting the BWT



## 7.7 Searching in the BWT

