



# 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 Random Access in BWT
- 7.8 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:

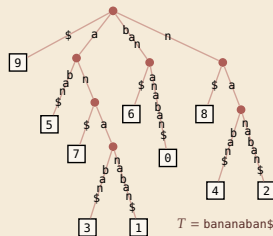
- ▶ **Given:** String  $T[0..n)$
- ▶ **Goal:** Answer LCE queries, i.e.,  
 given positions  $i, j$  in  $T$ ,  
 how far can we read the same text from there?  
 formally:  $\text{LCE}(i, j) = \max\{\ell : T[i..i + \ell) = T[j..j + \ell)\}$

↪ 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 ...



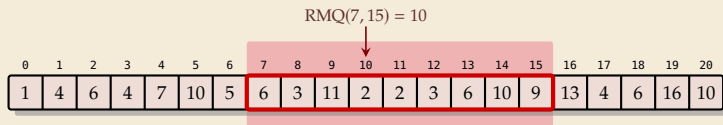
$\rightsquigarrow$  for now, use  $O(1)$  LCA as black box.

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

## 7.1 Range-Minimum Queries

# Range-minimum queries (RMQ)

- array / numbers don't change
- ▶ **Given:** Static array  $A[0..n)$  of numbers
  - ▶ **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

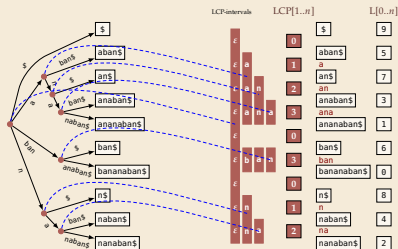
- to understand the fastest algorithm, it is helpful to define the *inverse suffix array*:

- $R[i] = r \iff L[r] = i$   $L = \text{leaf array}$
- $\iff$  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\$	right $R[0] = 6$	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\$	left $L[8] = 4$	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

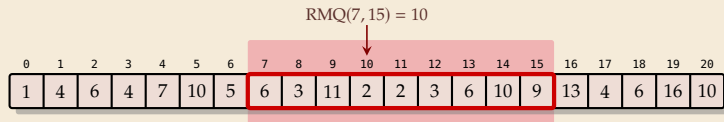
- ▶ For the following, consider RMQ on arbitrary arrays
- ▶ comparison-based  $\rightsquigarrow$  values don't matter, only relative order
- ▶ Two main quantities of interest:
  1. **Preprocessing time:** Running time  $P(n)$  of the preprocessing step  $\rightsquigarrow$  space usage  $\leq P(n)$
  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

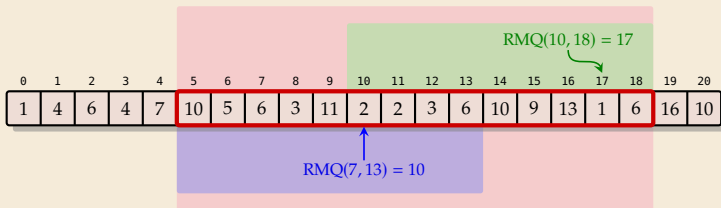
► **Idea:** Like “precompute-all”, but keep only *some* entries

► 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.  $\text{RMQ}(i, j) =$   
 $\arg \min\{A[\text{rmq}_1], A[\text{rmq}_2]\}$

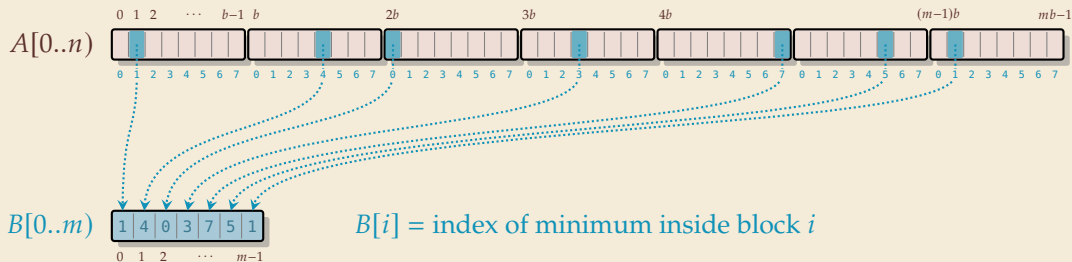
with  $\text{rmq}_1 = \text{RMQ}(i, i + 2^k - 1)$   
 $\text{rmq}_2 = \text{RMQ}(j - 2^k + 1, j)$

► 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  $m = \lceil n/b \rceil = O(n/\log n)$

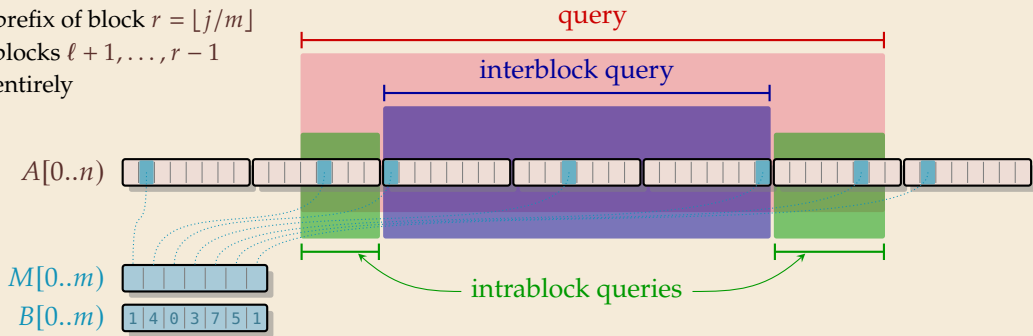


↪ Use sparse table solution for  $B$ .

↪ 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

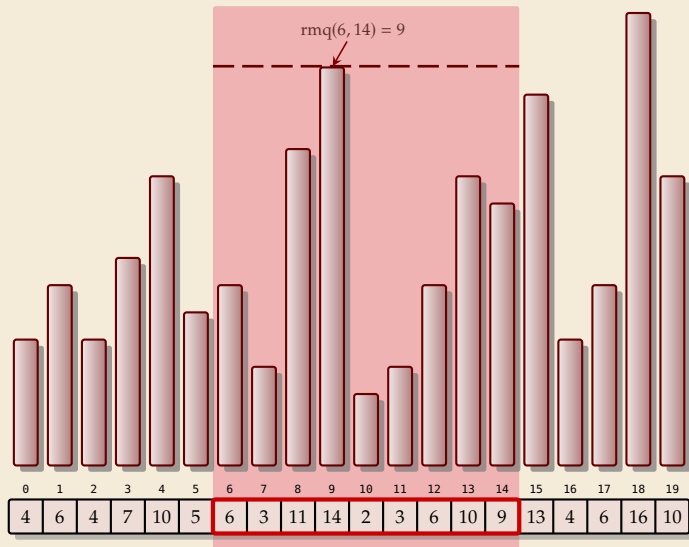


▶  $\text{RMQ}_A(i, j) = \arg \min_{k \in K} A[k]$  with  $K = \left\{ \begin{array}{l} \text{RMQ}_{\text{block } \ell}(i - \ell b, (\ell + 1)b - 1), \\ b \cdot \text{RMQ}_M(\ell + 1, r - 1) + \\ \quad B[\text{RMQ}_M(\ell + 1, r - 1)], \\ \text{RMQ}_{\text{block } r}(rb, j - rb) \end{array} \right\}$

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

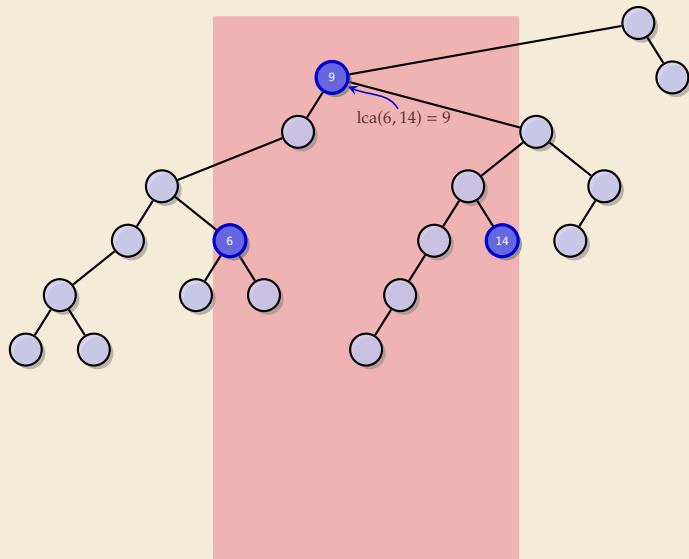
## 7.3 RMQ – Cartesian Trees

# 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!

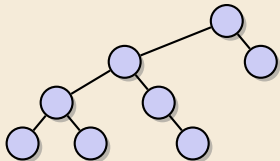
# RMQ & LCA



- **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!
- **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:  $\lg n + 2$  *string* comparisons of string of length  $m$

$\rightsquigarrow O(\log(n) \cdot m)$  character comparisons

- ▶ suffix tree could do  $O(m)$  total time (assuming constant  $\sigma$  or hashing for child links)
- ▶ surely, enhanced suffix arrays can do better than  $O(m \log n)$  🤔

**Idea:** use LCP information to save character comparisons

- ▶ concretely: maintain LCP between lower/upper bound suffixes and  $P$

$$T[a..n] \leq_{\text{lex}} P \leq_{\text{lex}} T[b..n]$$

$$\ell_a = \text{LCP}(T[a..n], P) \text{ and } \ell_b = \text{LCP}(T[b..n], P)$$

- ▶ avoid comparing same characters again
- ▶ Note: with RMQ on LCP array can determine  $\text{LCP}(T_i, T_j)$  for any  $i, j \in [0..n)$

# LCP Binary Search

- ▶ Input:  $\ell_a = \text{LCP}(T_a, P)$   
 $\ell_b = \text{LCP}(T_b, P)$   
 $\rightsquigarrow \ell_m = \text{LCP}(T_m, P) \geq \min\{\ell_a, \ell_b\}$
- ▶ **Case 1:**  $\ell_a = \ell_b$   
 Compare  $P$  and  $T_m$  starting at  $\ell_a$
- ▶ **Case 2:**  $\ell_a \neq \ell_b$ ; w.l.o.g.  $\ell_a > \ell_b$ 
  - ▶ **Case 2a:**  $\text{LCP}(T_a, T_m) > \ell_a$   
 $P >_{\text{lex}} T_m$  w/o cmps!
  - ▶ **Case 2b:**  $\text{LCP}(T_a, T_m) < \ell_a$   
 $P <_{\text{lex}} T_m$  w/o cmps!
  - ▶ **Case 2c:**  $\text{LCP}(T_a, T_m) = \ell_a$   
 Compare  $P$  and  $T_m$  from  $\ell_a$
- ▶ in each case, learn  $\ell_m \rightsquigarrow$  invariant
- ▶ no redundant '='-comparisons

0:	$T_{20}$		\$	
1:	$T_4$		ahbansbananasman\$	
2:	$T_{18}$		an\$	
3:	$T_{11}$	$a \rightarrow$	ananasman\$	Case 2c
4:	$T_{13}$	$m \rightarrow$	anasman\$	$P = \text{anarchy}$
5:	$T_1$		annahbansbananasman\$	
6:	$T_7$	$b \rightarrow$	ansbananasman\$	
7:	$T_{15}$		asman\$	
8:	$T_{10}$		bananasman\$	
9:	$T_6$		bansbananasman\$	
10:	$T_0$		hannahbansbananasman\$	
11:	$T_5$		hbansbananasman\$	
12:	$T_{17}$		man\$	
13:	$T_{19}$		n\$	
14:	$T_3$		nahbansbananasman\$	
15:	$T_{12}$		nanasman\$	
16:	$T_{14}$		nasman\$	
17:	$T_2$		nnahbansbananasman\$	
18:	$T_8$		nsbananasman\$	
19:	$T_9$		sbananasman\$	
20:	$T_{16}$		sman\$	

# Enhanced Suffix Arrays – Update

- ▶ *Enhanced suffix array*:  $L$ ,  $R$  and LCP array with RMQ support
- ▶ **Goal**: simulate any suffix tree operations
  - ▶ string matching in  $O(m + \log n)$  time ✓
  - ▶ string depth of internal nodes = LCP values ✓
  - ▶ internal suffix tree node = LCP interval ✓
    - ↪ storing information per node ✓
      - ▶ bottom-up traversal via enclosing LCP intervals ✓
  - ▶ longest common extension queries ✓
  - ▶ suffix links ✓

## 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

# Towards Self-Indexes

► For large genomes or multiple-genome datasets, can't hold  $T[0..n)$  in fast memory.

► An enhanced suffix array needs additional  $\Theta(n)$  words of space.

↪ When reference genomes first became available, a major show stopper!

► But since string matching can reconstruct  $T$ , can't avoid storing  $T$  somehow!

► A *self-index* is a data structure that answers operations without access to  $T$  at query time

► We get to decide *how* to store  $T$  ↪ might *compress*  $T$  (if compressible)

► Known as "*encoding model*" in space-efficient data structures ↗ genomes highly repetitive!

↪ **Key question:** *How to compress  $T$  while supporting random access and read mapping?*

*"Computing over compressed data"*

# BWT – Definitions

► *cyclic shift* of a string:

► with end-of-word character \$

~> can **recover** original string

$T = \text{time\_flies\_quickly\_}$

$\text{flies\_quickly\_time\_}$



~> cyclic shift



► The Burrows-Wheeler Transform proceeds in three steps:

0. Append end-of-word character \$ to  $S$ .

1. Consider *all cyclic shifts* of  $S$

2. Sort these strings lexicographically

3.  $B$  is the *list of trailing characters* (last column, top-down) of each string

# BWT – Example

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

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

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

alf\_eats\_alfalfa\$  
lf\_eats\_alfalfa\$  
f\_eats\_alfalfa\$  
\_eats\_alfalfa\$  
a\$alf\_eats\_alfal  
alf\_eats\_alfalfa\$  
ats\_alfalfa\$alf\_e  
ts\_alfalfa\$alf\_ea  
s\_alfalfa\$alf\_eat  
\_alfalfa\$alf\_eats  
alfalfa\$alf\_eats\_  
lfalfa\$alf\_eats\_  
falfa\$alf\_eats\_  
alfa\$alf\_eats\_  
lfa\$alf\_eats\_  
fa\$alf\_eats\_  
a\$alf\_eats\_  
s\_alfalfa\$alf\_eat  
ts\_alfalfa\$alf\_ea

sort

\$alf\_eats\_alfalf  
\_alfalfa\$alf\_eat  
\_eats\_alfalfa\$alf  
a\$alf\_eats\_alfal  
alf\_eats\_alfalfa\$  
alfa\$alf\_eats\_  
alfalfa\$alf\_eats\_  
ats\_alfalfa\$alf\_e  
eats\_alfalfa\$alf\_  
f\_eats\_alfalfa\$al  
fa\$alf\_eats\_alfal  
falfa\$alf\_eats\_  
lf\_eats\_alfalfa\$  
lfa\$alf\_eats\_  
lfalfa\$alf\_eats\_  
s\_alfalfa\$alf\_eat  
ts\_alfalfa\$alf\_ea

BWT  
↓

# 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  $L[i] = 0, B[i] = \$$

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

- ▶ more direct methods now also available

	$r$		$\downarrow L[r]$
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	alf\$aalf_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
alf\$aalf_eats_alf	12	lf_eats_alfalfa\$ <b>a</b>	1
lfa\$alf_eats_alf	13	lfa\$alf_eats_alf <b>a</b>	13
fa\$alf_eats_alfalf	14	lfalfa\$alf_eats_ <b>a</b>	10
a\$alf_eats_alfalf	15	s_alfalfa\$alf_eat <b>t</b>	7
\$alf_eats_alfalfa	16	ts_alfalfa\$alf_ <b>e</b> a	6

# BWT – Properties

$r$		$\downarrow L[r]$
0	\$alf_eats_alfalfa	16
1	_alfalfa\$alf_eats	8
2	_eats_alfalfa\$alf	3
3	a\$alf_eats_alfalf	15
4	alf_eats_alfalfa\$	0
5	alfa\$alf_eats_alf	12
6	alfalfa\$alf_eats_	9
7	ats_alfalfa\$alf_e	5
8	eats_alfalfa\$alf_	4
9	f_eats_alfalfa\$a	2
10	fa\$alf_eats_alfalf	14
11	falfa\$alf_eats_alf	11
12	lf_eats_alfalfa\$a	1
13	lfa\$alf_eats_alfalfa	13
14	lfalfa\$alf_eats_alf	10
15	s_alfalfa\$alf_eat	7
16	ts_alfalfa\$alf_ea	6

## Why does BWT help for compression?

- ▶ sorting *groups* characters by *what follows*
  - ▶ Example: lf always preceded by a
  - ▶ more generally: BWT can be partitioned into letters following a given context

(formally: low higher-order empirical entropy)

↪ If  $S$  allows predicting symbols from context,  
 $B$  has locally low entropy of characters.

- ▶ that makes MTF (move-to-front) transformation effective!
- ↪ use in compression pipeline for bzip2:  
BTW → MTF → RLE → Huffman

## A Bigger Example

[illegible][illegible]

$T =$  have, had, hadn't, hasn't, haven't, has, what,

$B =$    t e d t t t s h h h h h h a a v v   u u u u w \$   u e d s a a a n n n a a ,

# Run-length BWT Compression

- ▶ amazingly, just run-length compressing the BWT is already powerful!
- ▶  $r$  = number of runs in BWT

## Example:

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

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

$RL(B) = \begin{bmatrix} a \\ 1 \end{bmatrix} \begin{bmatrix} s \\ 1 \end{bmatrix} \begin{bmatrix} f \\ 2 \end{bmatrix} \begin{bmatrix} \$ \\ 1 \end{bmatrix} \begin{bmatrix} f \\ 1 \end{bmatrix} \begin{bmatrix} \_ \\ 1 \end{bmatrix} \begin{bmatrix} e \\ 1 \end{bmatrix} \begin{bmatrix} \_ \\ 1 \end{bmatrix} \begin{bmatrix} l \\ 3 \end{bmatrix} \begin{bmatrix} a \\ 3 \end{bmatrix} \begin{bmatrix} t \\ 1 \end{bmatrix} \begin{bmatrix} a \\ 1 \end{bmatrix}$

$\rightsquigarrow r = |RL(B)| = 12; \quad n = 17$

## Larger Example:

$S = \text{have\_had\_hadnt\_hasnt\_havent\_has\_what\$}$

$B = \text{tedtttshhhhhhhaavv\_w\$\_edsaaannnaa\_}$

$\rightsquigarrow r = 19; \quad n = 36$

- ▶ Indeed:  $r = O(z \log^2(n))$ ,  $z$  number of LZ77 phrases  
proven in 2019 (!)



Kempa, Kociumaka: Resolution of the Burrows-Wheeler Transform Conjecture, CACM 2022



## 7.6 Inverting the BWT

# Inverse BWT

- Great, can compute BWT efficiently and it helps compression. *But can we get  $T$  back?*

► **“Magic” solution:**

1. Create array  $D[0..n]$  of pairs:  
 $D[r] = (B[r], r)$ .
2. Sort  $D$  *stably* with respect to *first entry*.
3. Use  $D$  as linked list with (char, next entry)

**Example:**

$B = \text{ard\$rcaaaabb}$

$S = \text{abracadabra\$}$

	$D$	sorted $D$
		char next
0	(a, 0)	0 (\$, 3)
1	(r, 1)	1 (a, 0)
2	(d, 2)	2 (a, 6)
3	(\$, 3)	3 (a, 7)
4	(r, 4)	4 (a, 8)
5	(c, 5)	5 (a, 9)
6	(a, 6)	6 (b, 10)
7	(a, 7)	7 (b, 11)
8	(a, 8)	8 (c, 5)
9	(a, 9)	9 (d, 2)
10	(b, 10)	10 (r, 1)
11	(b, 11)	11 (r, 4)

# Inverse BWT – The magic revealed

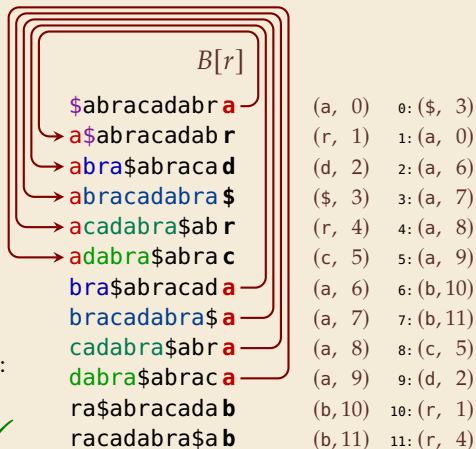
- ▶ Inverse BWT very easy to compute:
  - ▶ only sort individual characters in  $B$  (not suffixes)
  - ↪  $O(n)$  with counting sort

▶ *but why does this work!?*

- ▶ decode char by char
  - ▶ can find unique \$ ↪ starting row

- ▶ to get next char, we need
  - char in *first* column of *current row*
  - find row with that char's copy in BWT
 ↪ then we can walk through and decode

- ▶ for (i): first col = chars of  $B$  in sorted order ✓
- ▶ for (ii): relative order of same character stays same:  
 $i$ th a in first column =  $i$ th a in BWT  
 ↪ stably sorting  $(B[r], r)$  by first entry enough ✓



# Random Access Decoding

Can similarly output **any substring**  $T[i..i + \ell)$  if we know inverse suffix array:

Simply do  $\ell$  steps of the inverse BWT starting at  $r = R[i - 1]!$

$i$	$R[i]$	$T_i$	$r$	$L[r]$	$T_{L[r]}$	$B[r]$	$D$	$\text{sort}(D)$
0	6 <sup>th</sup>	bananabana\$	0	9	\$bananabana	n	(n, 0)	0: (\$, 6)
1	4 <sup>th</sup>	ananabana\$b	1	5	aban\$banana	n	(n, 1)	1: (a, 5)
2	9 <sup>th</sup>	nanabana\$b	2	7	an\$bananab	a	(b, 2)	2: (a, 7)
3	3 <sup>th</sup>	anabana\$ban	3	3	anabana\$ban	a	(n, 3)	3: (a, 8)
4	8 <sup>th</sup>	nabana\$ban	4	1	ananabana\$b	a	(b, 4)	4: (a, 9)
5	1 <sup>th</sup>	aban\$banana	5	6	ban\$banana	a	(a, 5)	5: (b, 2)
6	5 <sup>th</sup>	ban\$banana	6	0	bananabana\$	a	(\$, 6)	6: (b, 4)
7	2 <sup>th</sup>	an\$bananab	7	8	n\$bananab	a	(a, 7)	7: (n, 0)
8	7 <sup>th</sup>	n\$bananaba	8	4	nabana\$ban	a	(a, 8)	8: (n, 1)
9	0 <sup>th</sup>	\$bananabana	9	2	nanabana\$b	a	(a, 9)	9: (n, 3)

*sort suffixes*

Decoding only needs access to

1.  $i$ th char  $c$  of  $\text{sort}(T) = \text{sort}(B)$
2. *position* of (that copy of)  $c$  in  $B$

$\rightsquigarrow$  If we have that, can skip sorting / storing all of  $D$ !

## 7.7 Random Access in BWT

# Rank & Select on Sequences

Recall: Decoding only needs access to

1.  $i$ th char  $c$  of  $\text{sort}(T) = \text{sort}(B)$
2. *position* of (that copy of)  $c$  in  $B$

Both can be supported using  
rank/select on sequences.

- ▶  $\text{rank}_c(T[0..n], i) = |T[0..i]|_c$  #occurrences of  $c$   
= # $c$  in first  $i$  characters of  $T$
- ▶  $\text{select}_c(T[0..n], r)$   
=  $\min\{j : |T[0..j]|_c \geq r\} \cup \{n\}$   
= index of  $r$ th  $c$  in  $T$ , ( $r = 1, 2, \dots$ )

## Random Access in BWT

- ▶ store offsets  $O[c] = \sum_{c'=0}^{c-1} |B|_{c'}$  for  $c \in \Sigma$
- ▶  $i$ th char of  $\text{sort}(B) =$  unique  $c$  for which  $O[c] \leq i < O[c+1]$
- ▶ position of  $r$ th  $c$  in  $B = \text{select}_c(B, r)$

	0	1	2	3	4	5	6	7	8
$T[0..9]$	b	a	n	a	n	a	b	a	n

	0	1	2	3	4	5	6	7	8	9
$\text{rank}_a(T, i)$	0	0	1	1	2	2	3	3	4	4
$\text{rank}_b(T, i)$	0	1	1	1	1	1	1	2	2	2
$\text{rank}_n(T, i)$	0	0	0	1	1	2	2	2	2	3

$\text{select}_a(T, r)$	/	1	3	5	7	9	9	9	9	9
$\text{select}_b(T, r)$	/	0	6	9	9	9	9	9	9	9
$\text{select}_n(T, r)$	/	2	4	8	9	9	9	9	9	9

	0	1	2	3	4	5	6	7	8
$\text{sort}(T)$	a	a	a	a	b	b	n	n	n

$$O[0..\sigma] = [0, 4, 6, 9]$$

# Wavelet Trees



The **Wavelet Trees** for a  $T \in [0..n)$  over  $\Sigma = [0..\sigma)$

- ▶ supports access to  $T[i]$  in  $O(\log \sigma)$  time,
- ▶  $\text{rank}_c(T, i)$  and  $\text{select}_c(T, r)$  in  $O(\log \sigma)$  time, and
- ▶ occupies  $\sim n \lg \sigma$  bits of space. (Further compression possible!)  $\rightsquigarrow$  Advanced Data Structures
- ▶ The generalized  $\sigma$ - $\text{rank}_c(T, i) = \text{rank}_c(T, i) + \sum_{c' < c} |T|_{c'}$  is also supported in  $O(\log \sigma)$  time

Storing  $B[0..n]$  as a wavelet tree  $\rightsquigarrow$  reconstruct  $\ell$  chars from  $T$  in  $O(\ell \log \sigma)$  time  
if starting position known

e.g.,  $t = \lg n$

Storing **every**  $t$ <sup>th</sup> entry of  $R[0..n]$   $\rightsquigarrow$  may need to go back  $t$  characters for access  
 $\rightsquigarrow O((\ell + t) \log \sigma)$  time for decode  
using  $\sim n \lg n / t$  extra bits of space

## Locally decodable BWT

- ▶ no longer need to store  $T[0..n)$ !
- ▶ compressible (e.g., Wavelet trees with compressed bitvectors)

## 7.8 Searching in the BWT



# Backwards Search

Recall how the sorted suffixes in a suffix array  $L[0..n]$  made **string matching** very easy.

- Simply binary search the pattern  $P[0..m]$  in  $L$ !

↪ all occurrences must form interval

With wavelet tree BWT, we can replace binary search by **backwards radix search**!

- use  $\text{sort}(B)$  to locate interval for **last** character  $P[m-1]$
- use one step of inverse BWT to narrow down on  $P[m-2..m]$ , repeat.

$i$	$R[i]$	$T_i$		$L[r]$	$r$	$T_{L[r]}$	$B[r]$	$\text{sort}(D)$
0	6 <sup>th</sup>	bananaban\$	<div><math>P = \text{ana}</math></div>	9	0	\$bananaba	n	0: (\$, 6)
1	4 <sup>th</sup>	ananaban\$b		5	1	aban\$bana	n	1: (a, 5)
2	9 <sup>th</sup>	nanaban\$b		7	2	an\$b	anab	2: (a, 7)
3	3 <sup>th</sup>	anaban\$ban		3	3	anab	anab	3: (a, 8)
4	8 <sup>th</sup>	naban\$bana		1	4	anab	anab	4: (a, 9)
5	1 <sup>th</sup>	aban\$b		6	5	ban\$b	anab	5: (b, 2)
6	5 <sup>th</sup>	ban\$b		0	6	bananab	\$	6: (b, 4)
7	2 <sup>th</sup>	an\$b		8	7	n\$b	anab	7: (n, 0)
8	7 <sup>th</sup>	n\$b		4	8	nab	anab	8: (n, 1)
9	0 <sup>th</sup>	\$bananab		2	9	nanab	anab	9: (n, 3)

# Backwards Search – Code

Recall total rank operation supported by wavelet trees

$$\sigma\text{-rank}_c(B, i) = |B[0..i]|_c + \sum_{c' < c} |B|_{c'}$$

---

```
1 procedure backwardSearch( $B[0..n]$ ,  $P[0..m]$ )
2   //  $B[0..n]$  given as wavelet tree
3   // returns range  $[s..e)$  of ranks for suffixes starting with  $P$ 
4    $c := P[m - 1]$ 
5    $s := \sigma\text{-rank}_c(B, 0)$ 
6    $e := \sigma\text{-rank}_c(B, n)$ 
7   for  $j := m - 2, m - 3, \dots, 0$ 
8     if  $s \geq e$  break // no matches
9      $c := P[j]$ 
10     $s := \sigma\text{-rank}_c(B, s)$ 
11     $e := \sigma\text{-rank}_c(B, e)$ 
12  return  $[s..e)$ 
```

---

# Locating Matches

- ▶ Backwards Search finds interval  $[s..e)$  such that

$$P[0..m) = T[L[r] .. L[r]+m) \text{ iff } r \in [s..e)$$

↪ still need suffix array  $L[0..n]$  to locate matches!

- ▶ but can detect and count occurrences even without  $L$

## Sampled Suffix Array

- ▶ As for inverse suffix array, can store  $L[r]$  only for every  $t$ th starting index  $i$  in  $T$ , i. e., only store entries for ranks  $r$  with  $L[r] \equiv 0 \pmod{t}$

↪  $O(n \log n / t)$  bits of extra space

↪ Need to continue backwards search for at most  $t$  extra characters to locate match

↪ String matching in  $O(m \log \sigma + occ \cdot t \log \sigma)$  time

*Wavelet-tree BWT + Sampled Suffix Array = **FM Index***



Ferragina, Manzini: *Indexing compressed text*, JACM 2005

# FM-Index Discussion

- ▶ FM-Index is one of first *compressed self-indexes*
- ▶ can represent text using  $\sim \mathcal{H}_k(T)n$  bits of space  
 $\mathcal{H}_k(T)$  =  $k$ th order empirical entropy
- ▶ still widely used, e. g., as basis of *bowtie2* read alignment tool



Langmead, Salzberg: *Fast gapped-read alignment with Bowtie 2*, Nature Methods 2012

## Ongoing research

- ▶ Reduce space for very repetitive strings (collection of genomes)  
e. g.,  $r$ -index
- ▶ full support of suffix tree functionality with little extra space?



Navarro: *Indexing Highly Repetitive String Collections, Part II: Compressed Indexes*, ACM Comp. Surv. 2021