

# Optimal-Time Mapping in Run-Length Compressed PBWT

Paola Bonizzoni, Davide Cozzi, Younan Gao

University of Milano-Bicocca, Italy

DSB 2026, Venice, Italy

# Biological Context: Haplotypes and Bi-allelic Panels

- **Haplotypes:** A haplotype refers to a set of genetic variants (alleles) that are inherited together from a single parent.

$S_1$	0	1	0	1	1

# Biological Context: Haplotypes and Bi-allelic Panels

- **Haplotypes:** A haplotype refers to a set of genetic variants (alleles) that are inherited together from a single parent.
- **Variant Sites ( $w$ ):** Genomic locations where sequences vary across a population are called variant sites.

Haplotype	Site 1	Site 2	Site 3	Site 4	Site 5
$S_1$	0	1	0	1	1

# Biological Context: Haplotypes and Bi-allelic Panels

- **Haplotypes:** A haplotype refers to a set of genetic variants (alleles) that are inherited together from a single parent.
- **Variant Sites ( $w$ ):** Genomic locations where sequences vary across a population are called variant sites.
- **Bi-allelic Panels:** In most genomic studies, these sites are “bi-allelic,” meaning they only exhibit two possible states, typically represented as 0 or 1.

Haplotype	Site 1	Site 2	Site 3	Site 4	Site 5
$S_1$	0	1	0	1	1

# Biological Context: Haplotypes and Bi-allelic Panels

- **Haplotypes:** A haplotype refers to a set of genetic variants (alleles) that are inherited together from a single parent.
- **Variant Sites ( $w$ ):** Genomic locations where sequences vary across a population are called variant sites.
- **Bi-allelic Panels:** In most genomic studies, these sites are “bi-allelic,” meaning they only exhibit two possible states, typically represented as 0 or 1.
- **The Matrix Representation:** Large-scale genomic data is often stored as an  $h \times w$  binary matrix, where:

Haplotype	Site 1	Site 2	Site 3	Site 4	Site 5
$S_1$	0	1	0	1	1
$S_2$	0	0	1	1	0
$S_3$	0	0	1	1	0
$S_4$	0	1	1	1	1
$S_5$	1	0	0	0	1

# A Typical Application: Positional Search

- **Positional Search:** Given a query pattern  $P$  and a specific variant site  $j$ , find the **longest prefix**  $P[1..m']$  that appears in the population starting at site  $j$ .

$P[1..m]$	0	1	1	0	
Haplotype	Site 1	Site 2	Site 3	Site 4	Site 5
$S_1$	0	1	0	1	1
$S_2$	0	0	1	1	0
$S_3$	0	0	1	1	0
$S_4$	0	1	1	1	1
$S_5$	1	0	0	0	1

**Table:** Alignment of Pattern  $P$  starting at Site 2 of Matrix  $M$ .

# A Typical Application: Positional Search

- **Positional Search:** Given a query pattern  $P$  and a specific variant site  $j$ , find the **longest prefix**  $P[1..m']$  that appears in the population starting at site  $j$ .
- **Identification:** Upon finding the match, either identify a representative haplotype or enumerate all haplotypes in the panel that contain the sequence.

$P[1..m]$		0	1	1	0
Haplotype	Site 1	Site 2	Site 3	Site 4	Site 5
$S_1$	0	1	0	1	1
$S_2$	0	0	1	1	0
$S_3$	0	0	1	1	0
$S_4$	0	1	1	1	1
$S_5$	1	0	0	0	1

Table: Alignment of Pattern  $P$  starting at Site 2 of Matrix  $M$ .

# A Typical Application: Positional Search

## Motivation

This query serves as the core building block for complex genomic tasks, such as computing **matching statistics** and identifying **Set-Maximal Exact Matches (SMEs)**.

# A Typical Application: Positional Search

## Motivation

This query serves as the core building block for complex genomic tasks, such as computing **matching statistics** and identifying **Set-Maximal Exact Matches (SMEs)**.

To **facilitate** these queries at scale, the **Positional Burrows–Wheeler Transform (PBWT)** is employed, enabling high-speed prefix matching across massive populations.

# Positional Burrows–Wheeler Transform

- **Clustering Identical Alleles:** At any variant site  $j$ , the PBWT sorts haplotypes based on the co-lexicographic order of their sequences from site 1 through site  $j - 1$ .

M					
S1	1	0	1	1	0
S2	0	0	1	1	0
S3	0	1	1	0	1
S4	1	0	1	1	1
S5	1	0	0	0	1

co-lex

0	0	1
1	0	1
1	0	1
1	0	0
0	1	1

PBWT<sub>3</sub>

PA	
	2
	1
...	4
	5
	3

PA<sub>3</sub>

# Positional Burrows–Wheeler Transform

- **Clustering Identical Alleles:** At any variant site  $j$ , the PBWT sorts haplotypes based on the co-lexicographic order of their sequences from site 1 through site  $j - 1$ .
- **The Prefix Array (PA):** In each column  $j$ , the PA is a permutation of  $\{1, \dots, h\}$  that records the co-lexicographical order of the haplotypes.

M					
S1	1	0	1	1	0
S2	0	0	1	1	0
S3	0	1	1	0	1
S4	1	0	1	1	1
S5	1	0	0	0	1

co-lex

0	0	1
1	0	1
1	0	1
1	0	0
0	1	1

PBWT<sub>3</sub>

PA	
	2
	1
...	4
	5
	3

PA<sub>3</sub>

# Positional Burrows–Wheeler Transform

- **Creating “Runs”**: This sorting places haplotypes with similar prefixes in adjacent rows, causing identical alleles to cluster together.

M		PBWT		PA
S1	1 0 1 1 0	1 0 1 0 1		1 2 2 5 5
S2	0 0 1 1 0	0 1 1 1 1		2 3 1 2 3
S3	0 1 1 0 1	0 0 1 1 0		3 1 4 1 2
S4	1 0 1 1 1	1 0 0 1 0		4 4 5 4 1
S5	1 0 0 0 1	1 0 1 0 1		5 5 3 3 4

# Positional Burrows–Wheeler Transform

- **Creating “Runs”**: This sorting places haplotypes with similar prefixes in adjacent rows, causing identical alleles to cluster together.
- **Defining a Run**: A maximal contiguous block of identical symbols in a PBWT column is called a “run”.

M		PBWT		PA
S1	1 0 1 1 0	1 0 1 0 1		1 2 2 5 5
S2	0 0 1 1 0	0 1 1 1 1		2 3 1 2 3
S3	0 1 1 0 1	0 0 1 1 0		3 1 4 1 2
S4	1 0 1 1 1	1 0 0 1 0		4 4 5 4 1
S5	1 0 0 0 1	1 0 1 0 1		5 5 3 3 4

# Positional Burrows–Wheeler Transform

- $\mu$ -PBWT Compression:** By storing the starting positions, the lengths and values of these runs rather than every individual allele, the  $\mu$ -PBWT variant compresses massive populations into  $O(\tilde{r})$  space, where  $\tilde{r}$  is the total number of runs.

M		PBWT	1	2	3	4	5	$\mu$ -PBWT	1	2	3	4	5
S1	1 0 1 1 0		1	0	1	0	1		(1,1)	(1,0)	(1,1)	(1,0)	(1,1)
S2	0 0 1 1 0		0	1	1	1	1		(2,0)	(2,1)	(4,0)	(2,1)	(3,0)
S3	0 1 1 0 1		0	0	1	1	0		(4,1)	(3,0)	(5,1)	(5,0)	(5,1)
S4	1 0 1 1 1		1	0	0	1	0						
S5	1 0 0 0 1		1	0	1	0	1						

# Navigating the Panel: Fore and Back Stepping

- **Fore Query (Forward Stepping)**

- **Question:** If a haplotype is in row  $i$  at site  $j$ , what is its row index at site  $j + 1$ ?

M						PBWT						PA					
S1	1	0	1	1	0	1	0	1	0	1		1	2	2	5	5	
S2	0	0	1	1	0	0	1	1	1	1		2	3	1	2	3	
S3	0	1	1	0	1	0	0	1	1	0		3	1	4	1	2	
S4	1	0	1	1	1	1	0	0	1	0		4	4	5	4	1	
S5	1	0	0	0	1	1	0	1	0	1		5	5	3	3	4	

# Navigating the Panel: Fore and Back Stepping

- **Fore Query (Forward Stepping)**

- **Question:** If a haplotype is in row  $i$  at site  $j$ , what is its row index at site  $j + 1$ ?

M		PBWT		PA
S1	1 0 1 1 0	1 0 1 0 1		1 2 2 5 5
S2	0 0 1 1 0	0 1 1 1 1		2 3 1 2 3
S3	0 1 1 0 1	0 0 1 1 0		3 1 4 1 2
S4	1 0 1 1 1	1 0 0 1 0		4 4 5 4 1
S5	1 0 0 0 1	1 0 1 0 1		5 5 3 3 4

- **Back Query (Backward Stepping)**

- **Question:** Given a row index at site  $j + 1$ , what was the rank of that same haplotype at the previous site  $j$ ?

# Navigating the Panel: Fore and Back Stepping

## • Fore Query (Forward Stepping)

- **Question:** If a haplotype is in row  $i$  at site  $j$ , what is its row index at site  $j + 1$ ?

M		PBWT		PA
S1	1 0 1 1 0	1 0 1 0 1		1 2 2 5 5
S2	0 0 1 1 0	0 1 1 1 1		2 3 1 2 3
S3	0 1 1 0 1	0 0 1 1 0		3 1 4 1 2
S4	1 0 1 1 1	1 0 0 1 0		4 4 5 4 1
S5	1 0 0 0 1	1 0 1 0 1		5 5 3 3 4

## • Back Query (Backward Stepping)

- **Question:** Given a row index at site  $j + 1$ , what was the rank of that same haplotype at the previous site  $j$ ?

## Connection to Text Indexing

Both operations are conceptually similar to the **backward query** (LF-mapping) over the classic **Burrows–Wheeler Transform (BWT)**.

- **Uncompressed PBWT**: Traditional structures achieve **constant-time** fore and back stepping but require  $O(h \cdot w)$  words, which is prohibitive for massive datasets.

- **Uncompressed PBWT:** Traditional structures achieve **constant-time** fore and back stepping but require  $O(\mathfrak{h} \cdot w)$  words, which is prohibitive for massive datasets.
- **Standard Compressed PBWT:** Existing run-length encoded structures ( $\mu$ -PBWT) reduce space to  $O(\tilde{r})$  words, but prior to this work, each mapping step required  $O(\log \log_w \mathfrak{h})$  time.

- **Uncompressed PBWT:** Traditional structures achieve **constant-time** fore and back stepping but require  $O(h \cdot w)$  words, which is prohibitive for massive datasets.
- **Standard Compressed PBWT:** Existing run-length encoded structures ( $\mu$ -PBWT) reduce space to  $O(\tilde{r})$  words, but prior to this work, each mapping step required  $O(\log \log_w h)$  time.
- **Move Data Structures:** Recent advancements in text indexing introduced the "move structure," which supports constant-time backward queries over the BWT using  $O(r)$  words.

# Limitation of the Standard Move Structure

- **Iterative Stepping:** In PBWT applications, fore and back queries are performed sequentially (e.g., stepping from site  $j \rightarrow j + 1 \rightarrow j + 2$ ).

# Limitation of the Standard Move Structure

- **Iterative Stepping:** In PBWT applications, fore and back queries are performed sequentially (e.g., stepping from site  $j \rightarrow j + 1 \rightarrow j + 2$ ).
- Simply applying move structures independently at each site fails because they do not “pass the torch”—the index needed for the next iteration is lost.

# Limitation of the Standard Move Structure

- **Iterative Stepping:** In PBWT applications, fore and back queries are performed sequentially (e.g., stepping from site  $j \rightarrow j+1 \rightarrow j+2$ ).
- Simply applying move structures independently at each site fails because they do not “pass the torch”—the index needed for the next iteration is lost.
  - The classic move structure returns the sub-index that contains  $\text{fore}[j][j]$  at site  $j$ .

PBWT	1	2	3	4	5
	1	0	1	0	1
	0	1	1	1	1
	0	0	1	1	0
	1	0	0	1	0
	1	0	1	0	1

Query:  $\text{Fore}[2][2]$ ,  $x_2 = 2$ ?

Answer:  $\text{Fore}[2][2] = 5$ ,  $x_3 = 5$ .

# Limitation of the Standard Move Structure

- **Iterative Stepping:** In PBWT applications, fore and back queries are performed sequentially (e.g., stepping from site  $j \rightarrow j+1 \rightarrow j+2$ ).
- Simply applying move structures independently at each site fails because they do not “pass the torch”—the index needed for the next iteration is lost.
  - The classic move structure returns the sub-index that contains  $\text{fore}[j][j]$  at site  $j$ .
  - However, to perform the *next* step from  $j+1$  to  $j+2$ , the algorithm requires the **sub-run index** at site  $j+1$ .

PBWT	1	2	3	4	5
	1	0	1	0	1
	0	1	1	1	1
	0	0	1	1	0
	1	0	0	1	0
	1	0	1	0	1

Query:  $\text{Fore}[2][2]$ ,  $x_2 = 2?$

Answer:  $\text{Fore}[2][2] = 5$ ,  $x_3 = 5$ .

## Theoretical Bounds on $\tilde{r}$

Let  $h''$  be the number of adjacent pairs  $(S_i, S_{i+1})$  such that  $S_i \neq S_{i+1}$ .

- **Bounds:**  $h'' + 1 \leq \tilde{r} \leq w(h'' + 1)$ .
- **Corollary:**  $\tilde{r}$  is at least the number of distinct haplotypes.

# Summary of Our Results

## Theoretical Bounds on $\tilde{r}$

Let  $h''$  be the number of adjacent pairs  $(S_i, S_{i+1})$  such that  $S_i \neq S_{i+1}$ .

- **Bounds:**  $h'' + 1 \leq \tilde{r} \leq w(h'' + 1)$ .
- **Corollary:**  $\tilde{r}$  is at least the number of distinct haplotypes.

## Our Contribution: Constant-Time Navigation

Inspired by the move structure, we achieve **constant-time** fore and back stepping while maintaining  $O(\tilde{r})$  word space.

# Summary of Our Results

## Theoretical Bounds on $\tilde{r}$

Let  $h''$  be the number of adjacent pairs  $(S_i, S_{i+1})$  such that  $S_i \neq S_{i+1}$ .

- **Bounds:**  $h'' + 1 \leq \tilde{r} \leq w(h'' + 1)$ .
- **Corollary:**  $\tilde{r}$  is at least the number of distinct haplotypes.

## Our Contribution: Constant-Time Navigation

Inspired by the move structure, we achieve **constant-time** fore and back stepping while maintaining  $O(\tilde{r})$  word space.

- **Applications:** *Prefix Searches and Haplotype Retrieval*

# Summary of Our Results

## Theoretical Bounds on $\tilde{r}$

Let  $h''$  be the number of adjacent pairs  $(S_i, S_{i+1})$  such that  $S_i \neq S_{i+1}$ .

- **Bounds:**  $h'' + 1 \leq \tilde{r} \leq w(h'' + 1)$ .
- **Corollary:**  $\tilde{r}$  is at least the number of distinct haplotypes.

## Our Contribution: Constant-Time Navigation

Inspired by the move structure, we achieve **constant-time** fore and back stepping while maintaining  $O(\tilde{r})$  word space.

- **Applications:** *Prefix Searches* and *Haplotype Retrieval*
- **Multi-Allelic Support:** Our solution naturally generalizes to **multi-allelic panels** without additional overhead.

# Our Results: High-Level Idea

- **Sub-run Partitioning:** We divide the original runs of the PBWT into smaller **sub-runs**.

# Our Results: High-Level Idea

- **Sub-run Partitioning:** We divide the original runs of the PBWT into smaller **sub-runs**.
- **The Three-Overlap Constraint:** Each sub-run is constructed such that it “involves” (overlaps with) **at most three** sub-runs in the adjacent columns.

# Our Results: High-Level Idea

- **Sub-run Partitioning:** We divide the original runs of the PBWT into smaller **sub-runs**.
- **The Three-Overlap Constraint:** Each sub-run is constructed such that it “involves” (overlaps with) **at most three** sub-runs in the adjacent columns.
- **Tabulated Mapping:** We use a lookup table to store this involvement information.

# Our Results: High-Level Idea

- **Sub-run Partitioning:** We divide the original runs of the PBWT into smaller **sub-runs**.
- **The Three-Overlap Constraint:** Each sub-run is constructed such that it “involves” (overlaps with) **at most three** sub-runs in the adjacent columns.
- **Tabulated Mapping:** We use a lookup table to store this involvement information.
- **Constant-Time Navigation:**
  - Given  $(i, j)$  and the index  $x$  of the sub-run containing entry  $(i, j)$ , we compute  $\text{fore}(i, j)$  and the sub-run index at site  $j + 1$  in **constant time**.
  - Given  $(i, j)$  and the same index  $x$ , we compute  $\text{back}(i, j)$  and the sub-run index at site  $j - 1$  in **constant time**.

# Our Results: High-Level Idea

PBWT	1	2	3	4	5
	1	0	1	0	1
	0	1	1	1	1
	0	0	1	1	0
	1	0	0	1	0
	1	0	1	0	1

Query:  $\text{Fore}[2][2]$ ,  $x_2 = 2$ ?

Old Answer:  $\text{Fore}[2][2] = 5$ ,  $x_3 = 5$ .

New Answer:  $\text{Fore}[2][2] = 5$ ,  $x_3 = 4$ .

# The Three-Overlap Constraint

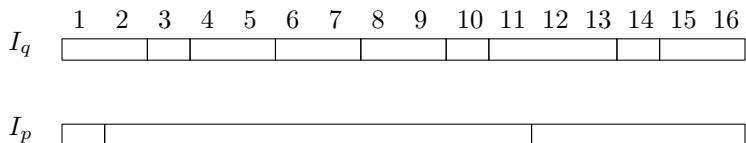
## Definition

A partition  $I_p$  satisfies the **3-overlap constraint** with respect to  $I_q$  if every interval in  $I_p$  overlaps at most three intervals in  $I_q$ .

# The Three-Overlap Constraint

## Definition

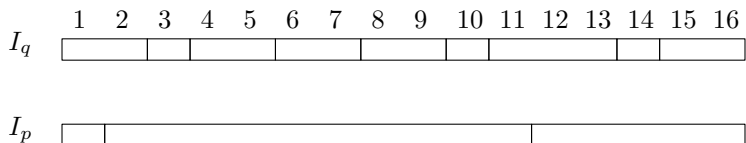
A partition  $I_p$  satisfies the **3-overlap constraint** with respect to  $I_q$  if every interval in  $I_p$  overlaps at most three intervals in  $I_q$ .



# The Three-Overlap Constraint

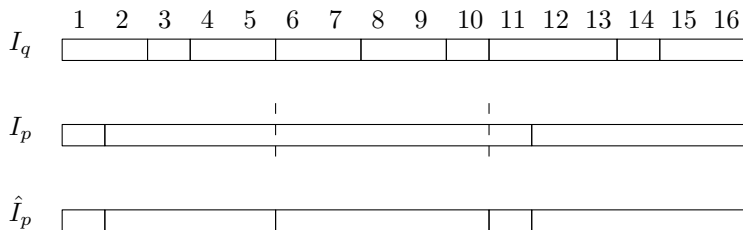
## Definition

A partition  $I_p$  satisfies the **3-overlap constraint** with respect to  $I_q$  if every interval in  $I_p$  overlaps at most three intervals in  $I_q$ .



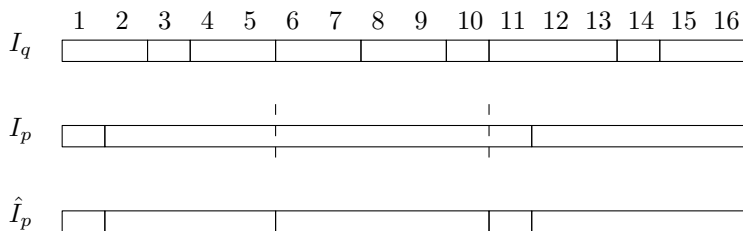
The 3-overlap constraint for  $I_p$  is violated.

# The Three-Overlap Constraint



- **Normalization Algorithm:** We can transform any partition  $I_p$  into a new partition  $\hat{I}_p$  that satisfies the constraint in  $O(|I_p| + |I_q|)$  time.

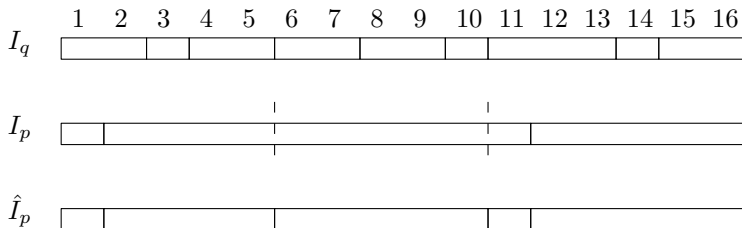
# The Three-Overlap Constraint



- **Normalization Algorithm:** We can transform any partition  $I_p$  into a new partition  $\hat{I}_p$  that satisfies the constraint in  $O(|I_p| + |I_q|)$  time.
- **Size Bound:** The number of the sub-intervals in  $\hat{I}_p$  is bounded by:

$$|\hat{I}_p| \leq |I_p| + \left\lceil \frac{|I_q|}{2} \right\rceil$$

# The Three-Overlap Constraint



- **Normalization Algorithm:** We can transform any partition  $I_p$  into a new partition  $\hat{I}_p$  that satisfies the constraint in  $O(|I_p| + |I_q|)$  time.
- **Size Bound:** The number of the sub-intervals in  $\hat{I}_p$  is bounded by:

$$|\hat{I}_p| \leq |I_p| + \left\lfloor \frac{|I_q|}{2} \right\rfloor$$

- **Implication:** This bound ensures that while we increase the number of intervals to satisfy the constraint, the total space remains linear in terms of the number of runs  $O(\tilde{r})$ .

# Constructing Sub-Runs SubIB's: Defining the Bijection

For each column  $j$ , we define a bijection  $\text{foreL}_j$  to construct the list of sub-runs  $\text{SubIB}_{j+1}$  required for constant-time back queries.

# Constructing Sub-Runs SubIB's: Defining the Bijection

- The Mapping Function:** Given a set of intervals  $L = \{[b_\tau, e_\tau]\}_{1 \leq \tau \leq |L|}$  at site  $j$ , the bijection  $\text{foreL}_j(L)$  maps them to site  $j+1$ :  $\text{foreL}_j(L) = \{[\text{fore}(b_\tau, j), \text{fore}(e_\tau, j)] \mid 1 \leq \tau \leq |L|\}$

PBWT						PA					
	1	0	1	0	1		1	2	2	5	5
	0	1	1	1	1		2	3	1	2	3
	0	0	1	1	0		3	1	4	1	2
	1	0	0	1	0		4	4	5	4	1
	1	0	1	0	1		5	5	3	3	4

$$\text{Fore}_3(\{[1,3], [4, 4], [5, 5]\}) = \{[1,1], [2,4], [5, 5]\}$$

# Constructing Sub-Runs SubIB's: Defining the Bijection

- **The Mapping Function:** Given a set of intervals  $L = \{[b_\tau, e_\tau]\}_{1 \leq \tau \leq |L|}$  at site  $j$ , the bijection  $\text{foreL}_j(L)$  maps them to site  $j+1$ :  $\text{foreL}_j(L) = \{[\text{fore}(b_\tau, j), \text{fore}(e_\tau, j)] \mid 1 \leq \tau \leq |L|\}$
- **Ordered Output:** sorted in increasingly order by left endpoints.

PBWT	PA
1 0 1 0 1	1 2 2 5 5
0 1 1 1 1	2 3 1 2 3
0 0 1 1 0	3 1 4 1 2
1 0 0 1 0	4 4 5 4 1
1 0 1 0 1	5 5 3 3 4

$$\text{Fore}_3(\{[1,3], [4, 4], [5, 5]\}) = \{[1,1], [2,4], [5, 5]\}$$

# Constructing Sub-Runs SubIB: The Algorithm

The list of sub-runs  $\text{SubIB}_j$  is constructed inductively for each column  $j = 1 \dots w$ :

# Constructing Sub-Runs SubIB: The Algorithm

The list of sub-runs  $\text{SubIB}_j$  is constructed inductively for each column  $j = 1 \dots w$ :

- **Base Case** ( $j = 1$ ): Set  $\text{SubIB}_1 = \text{intervals}_1$  (the original runs of the first column).

# Constructing Sub-Runs SubIB: The Algorithm

- **Inductive Step** ( $j > 1$ ): To construct  $\text{SubIB}_j$ , we reconcile the list  $\text{intervals}_j$  of the runs at column  $j$  with  $\text{SubIB}_{j-1}$ :

Input:  $\text{subIB}_{j-1}$

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
17	2	14	5	11	1	16	8	4	12	6	3	10	15	9	13

$\text{subIB}_{j-1}$

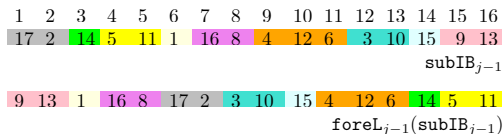
# Constructing Sub-Runs SubIB: The Algorithm

- **Inductive Step** ( $j > 1$ ): To construct  $\text{SubIB}_j$ , we reconcile the list  $\text{intervals}_j$  of the runs at column  $j$  with  $\text{SubIB}_{j-1}$ :
  - ① **Map**: Project the previous sub-runs forward:  $l_q = \text{foreL}_{j-1}(\text{SubIB}_{j-1})$ .

Input:  $\text{subIB}_{j-1}$

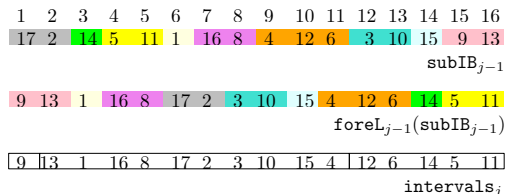
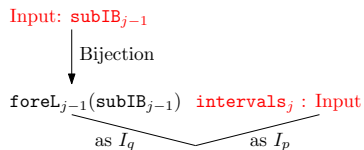
↓  
Bijection

$\text{foreL}_{j-1}(\text{subIB}_{j-1})$



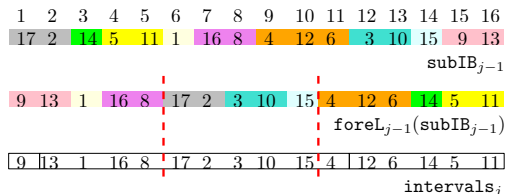
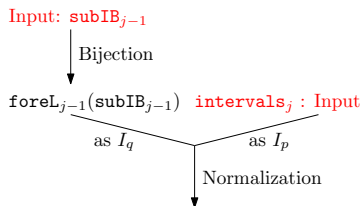
# Constructing Sub-Runs SubIB: The Algorithm

- Inductive Step** ( $j > 1$ ): To construct  $\text{SubIB}_j$ , we reconcile the list  $\text{intervals}_j$  of the runs at column  $j$  with  $\text{SubIB}_{j-1}$ :
  - Map**: Project the previous sub-runs forward:  $I_q = \text{foreL}_{j-1}(\text{SubIB}_{j-1})$ .



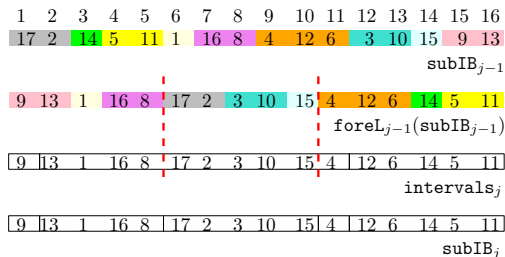
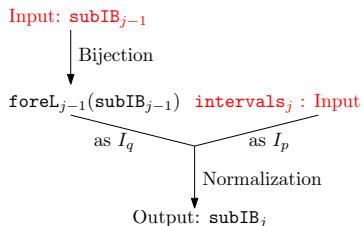
# Constructing Sub-Runs SubIB: The Algorithm

- Inductive Step ( $j > 1$ ):** To construct  $\text{SubIB}_j$ , we reconcile the list  $\text{intervals}_j$  of the runs at column  $j$  with  $\text{SubIB}_{j-1}$ :
  - Map:** Project the previous sub-runs forward:  $I_q = \text{foreL}_{j-1}(\text{SubIB}_{j-1})$ .
  - Normalize:** Apply  $\text{normalization}(I_p, I_q)$  where  $I_p = \text{intervals}_j$ .



# Constructing Sub-Runs SubIB: The Algorithm

- **Inductive Step** ( $j > 1$ ): To construct  $\text{SubIB}_j$ , we reconcile the list  $\text{intervals}_j$  of the runs at column  $j$  with  $\text{SubIB}_{j-1}$ :
  - 1 **Map**: Project the previous sub-runs forward:  $I_q = \text{foreL}_{j-1}(\text{SubIB}_{j-1})$ .
  - 2 **Normalize**: Apply  $\text{normalization}(I_p, I_q)$  where  $I_p = \text{intervals}_j$ .
  - 3 **Assign**: The resulting partition becomes  $\text{SubIB}_j$ .



# Constructing Sub-Runs SubIB: The Algorithm

## Key Properties of the Construction

- Every interval in  $\text{SubIB}_j$  is a **sub-interval** of an interval in  $\text{intervals}_j$  (and thus corresponds to a **sub-run** at column  $j$ ).

9	13	1	16	8	17	2	3	10	15	4	12	6	14	5	11
---	----	---	----	---	----	---	---	----	----	---	----	---	----	---	----

$\text{intervals}_j$

9	13	1	16	8	17	2	3	10	15	4	12	6	14	5	11
---	----	---	----	---	----	---	---	----	----	---	----	---	----	---	----

$\text{subIB}_j$

# Constructing Sub-Runs SubIB: The Algorithm

## Key Properties of the Construction

- Every interval in  $\text{SubIB}_j$  is a **sub-interval** of an interval in  $\text{intervals}_j$  (and thus corresponds to a **sub-run** at column  $j$ ).
- Every interval in  $\text{SubIB}_j$  overlaps **at most three** intervals in  $\text{foreL}_{j-1}(\text{SubIB}_{j-1})$ .

9 13 1 16 8 17 2 3 10 15 4 12 6 14 5 11  
 $\text{foreL}_{j-1}(\text{subIB}_{j-1})$

9 13 1 16 8 17 2 3 10 15 4 12 6 14 5 11  
 $\text{subIB}_j$

# Space Complexity: Bounding the Number of Sub-Runs

## Lemma

*The total number of sub-runs is linearly bounded by  $2\tilde{r}$ .*

## Lemma

*The total number of sub-runs is linearly bounded by  $2\tilde{r}$ .*

- **Recursive Step:** From the Normalization Lemma, we have:

$$|\text{SubIB}_j| \leq r_j + \frac{|\text{SubIB}_{j-1}|}{2}$$

# Space Complexity: Bounding the Number of Sub-Runs

## Lemma

*The total number of sub-runs is linearly bounded by  $2\tilde{r}$ .*

- **Recursive Step:** From the Normalization Lemma, we have:

$$|\text{SubIB}_j| \leq r_j + \frac{|\text{SubIB}_{j-1}|}{2}$$

- **Geometric Expansion:** Solving this recursion reveals a geometric sequence:

$$|\text{SubIB}_j| \leq \sum_{1 \leq \tau \leq j} \frac{r_\tau}{2^{j-\tau}} = r_j + \frac{r_{j-1}}{2} + \frac{r_{j-2}}{4} + \dots$$

# Space Complexity: Bounding the Number of Sub-Runs

## Lemma

*The total number of sub-runs is linearly bounded by  $2\tilde{r}$ .*

- **Recursive Step:** From the Normalization Lemma, we have:

$$|\text{SubIB}_j| \leq r_j + \frac{|\text{SubIB}_{j-1}|}{2}$$

- **Geometric Expansion:** Solving this recursion reveals a geometric sequence:

$$|\text{SubIB}_j| \leq \sum_{1 \leq \tau \leq j} \frac{r_\tau}{2^{j-\tau}} = r_j + \frac{r_{j-1}}{2} + \frac{r_{j-2}}{4} + \dots$$

- **Total Summation:** By summing across all columns  $j$ :

$$\sum_{1 \leq j \leq w} |\text{SubIB}_j| \leq \sum_{1 \leq j \leq w} \sum_{1 \leq \tau \leq j} \frac{r_\tau}{2^{j-\tau}} < 2 \sum_{j=1}^w r_j = 2\tilde{r}$$

# Main Result: Constant-Time Navigation

## Theorem

*There exists a data structure of  $O(\tilde{r})$  words that supports *fore*, *back*, and *character retrieval* in  $O(1)$  **time** per step, without accessing the original matrix or PBWT.*

# Main Result: Constant-Time Navigation

`back[7][j], xj = 3`

9	13	1	16	8	17	2	3	10	15	4	12	6	14	5	11
---	----	---	----	---	----	---	---	----	----	---	----	---	----	---	----

SubIB<sub>j</sub>

# Main Result: Constant-Time Navigation

$\text{back}[7][j], x_j = 3$

9	13	1	16	8	17	2	3	10	15	4	12	6	14	5	11
---	----	---	----	---	----	---	---	----	----	---	----	---	----	---	----

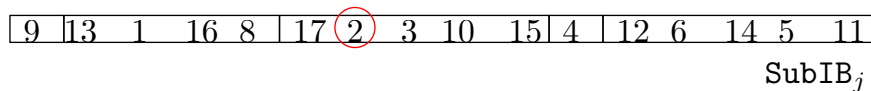
SubIB<sub>j</sub>

9	13	1	16	8	17	2	3	10	15	4	12	6	14	5	11
---	----	---	----	---	----	---	---	----	----	---	----	---	----	---	----

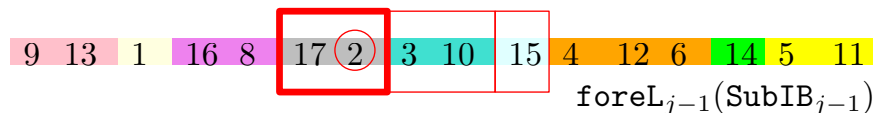
foreL<sub>j-1</sub>(SubIB<sub>j-1</sub>)

# Main Result: Constant-Time Navigation

$\text{back}[7][j], x_j = 3$



offset=2



# Main Result: Constant-Time Navigation

$$\text{back}[7][j] = 2, x_{j-1} = 1$$

9	13	1	16	8	17	2	3	10	15	4	12	6	14	5	11
---	----	---	----	---	----	---	---	----	----	---	----	---	----	---	----

SubIB<sub>j</sub>

offset=2

9	13	1	16	8	17	2	3	10	15	4	12	6	14	5	11
---	----	---	----	---	----	---	---	----	----	---	----	---	----	---	----

foreL<sub>j-1</sub>(SubIB<sub>j-1</sub>)

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
17	2	14	5	11	1	16	8	4	12	6	3	10	15	9	13

SubIB<sub>j-1</sub>

# Application: Haplotype Retrieval

## Theorem

Any  $S_i$  can be retrieved in  $O(w + \log \log_w h)$  time using  $O(\tilde{r})$  words.

M					
S1	1	0	1	1	0
S2	0	0	1	1	0
S3	0	1	1	0	1
S4	1	0	1	1	1
S5	1	0	0	0	1

# Application: Haplotype Retrieval

## Theorem

Any  $S_i$  can be retrieved in  $O(w + \log \log_w h)$  time using  $O(\tilde{r})$  words.

- 1 Predecessor query to find initial sub-run ( $O(\log \log_w h)$ ).

M	PBWT									
S1	1	0	1	1	0	1	0	1	0	1
S2	0	0	1	1	0	0	1	1	1	1
S3	0	1	1	0	1	0	0	1	1	0
S4	1	0	1	1	1	1	0	0	1	0
S5	1	0	0	0	1	1	0	1	0	1

$$x_1 = 2$$



# Application: Prefix Search

- **Goal:** Find longest common prefix  $P[1..m']$  and list occ haplotypes.
- **Complexity:**  $O(m' \log \log_w \sigma + \text{occ})$  time in  $O(h + \tilde{r})$  space.
- **Variable Lengths:** Handled via terminal symbol #.

	1	2	3	4	5	6	7	8	9
$S_1$	a	b	a	a	b	a	#		
$S_2$	a	b	a	a	b	b	#		
$S_3$	a	b	a	c	#				
$S_4$	b	c	b	c	a	b	a	#	
$S_5$	b	c	b	c	a	b	b	#	
$S_6$	b	c	b	c	b	b	a	#	
$S_7$	b	c	b	c	b	b	b	a	#

(a) Arbitrary Lengths

	1	2	3	4	5	6	7	8	9
1	a	b	a	a	b	b	#	#	#
2	a	b	a	a	b	b	a	#	
3	a	b	a	c	#	a	b	#	
4	b	c	b	c	a	b	#	a	
5	b	c	b	c	a	b	a		
6	b	c	b	c	b	b	b		
7	b	c	b	c	b				

(b) Extended PBWT

	1	2	3	4	5	6	7	8	9
1	1	1	1	1	1	4	1	4	7
2	2	2	2	2	2	5	4	6	
3	3	3	3	3	3	1	5	5	
4	4	4	4	4	4	2	2	7	
5	5	5	5	5	5	6	6		
6	6	6	6	6	6	7	7		
7	7	7	7	7	7				

(c) Extended PA

Figure: Handling varying haplotype lengths.

- **Optimal Time:** Constant-time mapping achieved in run-length compressed PBWT space.
- **Lower and Upper Bounds on  $\tilde{r}$ :** Established relationship between  $\tilde{r}$  and non-identical consecutive haplotype pairs.
- **Practical Utility:** Revisited haplotype retrieval and prefix search applications.
- **Future Work:** Accelerating SMEM and MPSC computations.

Thank You for Your Attention!

## Questions?

**Younan Gao**

University of Milano-Bicocca  
younan.gao@unimib.it