



UNIVERSITÀ DI PISA

# Metagenomic analysis through the eBWT

Veronica Guerrini    Giovanna Rosone

University of Pisa, Italy

Supported by the Project MIUR-SIR CMACBioSeq

"Combinatorial Methods for Analysis and Compression of Biological Sequences"

BITS 2019

Bioinformatics Italian Society Annual Meeting

Palermo, June 26 - 28, 2019

# Introduction

**Metagenomics** is the study of genetic material collected from the environment



[Illustration: Spencer Phillips, EMBL-EBI]

Aim to explore the relations between the microbes and their habitats

**Applications.** Clinical microbiology, plant-microbe interactions, monitoring pollution, sustainability, ecology, ...

**Goal:** Identify the taxon of each short read

# Introduction

**Metagenomics** is the study of genetic material collected from the environment



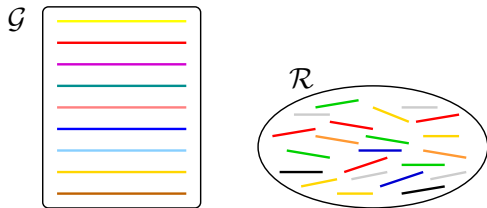
[Illustration: Spencer Phillips, EMBL-EBI]

Aim to explore the relations between the microbes and their habitats

**Applications.** Clinical microbiology, plant-microbe interactions, monitoring pollution, sustainability, ecology, ...

**Goal:** Identify the taxon of each short read

# Our approach based on (e)BWT



We introduce an alignment-free and assembly-free strategy...

... by using the properties of

- the **Burrows-Wheeler Transform** (BWT) of a string,
- an **extension** of BWT to a multiset of strings (eBWT).

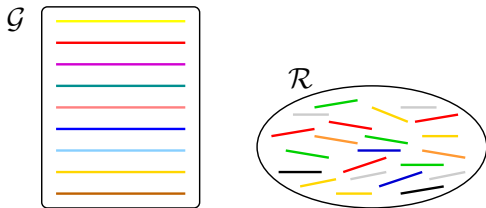
## Definition

Given a string  $v$  (resp. a string collection  $S$ ), the **Burrows-Wheeler Transform** (resp. **extended BWT**)<sup>a</sup> is a **reversible** transformation that produces a permutation of the symbols of  $v$  (resp.  $S$ ), defined over an ordered alphabet.

<sup>a</sup>Burrows and D. Wheeler, A block sorting lossless data compression algorithm, Technical Report 124, Digital Equipment Corporation, 1994  
S. Mantaci, A. Restivo, G. R., M. Sciortino: An extension of the Burrows-Wheeler Transform. TCS 2007

M.Bauer, A. Cox, G. R.: Lightweight algorithms for constructing and inverting the BWT of string collections. TCS 2013

# Our approach based on (e)BWT



We introduce an alignment-free and assembly-free strategy...

... by using the properties of

- the **Burrows-Wheeler Transform** (BWT) of a string,
- an **extension** of BWT to a multiset of strings (eBWT).

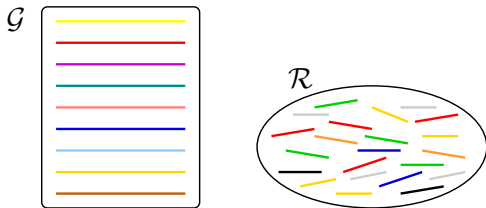
## Definition

Given a string  $v$  (resp. a string collection  $S$ ), the **Burrows-Wheeler Transform** (resp. **extended BWT**)<sup>a</sup> is a **reversible** transformation that produces a permutation of the symbols of  $v$  (resp.  $S$ ), defined over an ordered alphabet.

<sup>a</sup>Burrows and D. Wheeler, A block sorting lossless data compression algorithm, Technical Report 124, Digital Equipment Corporation, 1994  
S. Mantaci, A. Restivo, G. R., M. Sciortino: An extension of the Burrows-Wheeler Transform. TCS 2007

M.Bauer, A. Cox, G. R.: Lightweight algorithms for constructing and inverting the BWT of string collections. TCS 2013

# Our approach based on (e)BWT



We introduce an alignment-free and assembly-free strategy...

... by using the properties of

- the **Burrows-Wheeler Transform** (BWT) of a string,
- an **extension** of BWT to a multiset of strings (eBWT).

## Definition

Given a string  $v$  (resp. a string collection  $S$ ), the **Burrows-Wheeler Transform** (resp. **extended BWT**)<sup>a</sup> is a **reversible** transformation that produces a permutation of the symbols of  $v$  (resp.  $S$ ), defined over an ordered alphabet.

<sup>a</sup>Burrows and D. Wheeler, A block sorting lossless data compression algorithm, Technical Report 124, Digital Equipment Corporation, 1994  
S. Mantaci, A. Restivo, G. R., M. Sciortino: An extension of the Burrows-Wheeler Transform. TCS 2007

M.Bauer, A. Cox, G. R.: Lightweight algorithms for constructing and inverting the BWT of string collections. TCS 2013

# How does EBWT work?

Let  $S = \{S_1, S_2, \dots, S_m\}$  be a collection of strings.

- Append an end-marker  $\$$ , where  $\$_i < A < C < G < T$  (we add subscripts  $\$_i$  only for illustrative purposes), to each string in  $S$  by obtaining a new collection  $S'$ .

String collection S						
	0	1	2	3	4	5
$S_1$	G	C	C	A	A	C
$S_2$	G	A	G	C	T	C
$S_3$	T	C	G	C	T	T

- Sort all the suffixes of the strings in  $S'$ ;
- Take the string  $eBWT(S')$  obtained by concatenating the symbols that (circularly) precede the first symbol of each suffix in the list of (lexicographically) sorted suffixes of  $S'$ .

$eBWT(S') = CCTCAGATCGTGG\$_2\$_1ACTC\$_3C$

Note that the colors and suffixes are only for illustrative purposes.

Sorted Suffixes of $S'$
$\$_1$
$\$_2$
$\$_3$
AACS $_1$
ACS $_1$
AGCTC $\$_2$
C $\$_1$
C $\$_2$
CAAC $\$_1$
CCAAC $\$_1$
CGCTT $\$_3$
CTC $\$_2$
CTT $\$_3$
GAGCTC $\$_2$
GCCAAC $\$_1$
GTC $\$_2$
GCTT $\$_3$
T $\$_3$
TC $\$_2$
TCGCTT $\$_3$

# How does EBWT work?

Let  $S = \{S_1, S_2, \dots, S_m\}$  be a collection of strings.

- Append an end-marker  $\$$ , where  $\$_i < A < C < G < T$  (we add subscripts  $\$_i$  only for illustrative purposes), to each string in  $S$  by obtaining a new collection  $S'$ .

String collection $S'$							
	0	1	2	3	4	5	6
$S_1$	G	C	C	A	A	C	$\$_1$
$S_2$	G	A	G	C	T	C	$\$_2$
$S_3$	T	C	G	C	T	T	$\$_3$

- Sort all the suffixes of the strings in  $S'$ ;
- Take the string  $eBWT(S')$  obtained by concatenating the symbols that (circularly) precede the first symbol of each suffix in the list of (lexicographically) sorted suffixes of  $S'$ .

$eBWT(S') = CCTCAGATCGTGG\$_2\$_1ACTC\$_3C$

Note that the colors and suffixes are only for illustrative purposes.

Sorted Suffixes of  $S'$

$\$_1$   
 $\$_2$   
 $\$_3$   
 AAC $\$_1$   
 AC $\$_1$   
 AGCTC $\$_2$   
 C $\$_1$   
 C $\$_2$   
 CAAC $\$_1$   
 CCAAC $\$_1$   
 CGCTT $\$_3$   
 CTC $\$_2$   
 CTT $\$_3$   
 GAGCTC $\$_2$   
 GCCAAC $\$_1$   
 GCTC $\$_2$   
 GCTT $\$_3$   
 T $\$_3$   
 TC $\$_2$   
 TCGCTT $\$_3$



# How does EBWT work?

Let  $S = \{S_1, S_2, \dots, S_m\}$  be a collection of strings.

- Append an end-marker  $\$$ , where  $\$_i < A < C < G < T$  (we add subscripts  $\$_i$  only for illustrative purposes), to each string in  $S$  by obtaining a new collection  $S'$ .

String collection $S'$							
	0	1	2	3	4	5	6
$S_1$	G	C	C	A	A	C	$\$_1$
$S_2$	G	A	G	C	T	C	$\$_2$
$S_3$	T	C	G	C	T	T	$\$_3$

Sorted Suffixes of $S'$
$\$_1$
$\$_2$
$\$_3$
AAC $\$_1$
AC $\$_1$
AGCTC $\$_2$
C $\$_1$
C $\$_2$
CAAC $\$_1$
CCAAC $\$_1$
CGCTT $\$_3$
CTC $\$_2$
CTT $\$_3$
GAGCTC $\$_2$
GCCAAC $\$_1$
GCTC $\$_2$
GCTT $\$_3$
T $\$_3$
TC $\$_2$
TCGCTT $\$_3$
TT $\$_3$

- Sort all the suffixes of the strings in  $S'$ ;
- Take the string  $eBWT(S')$  obtained by concatenating the symbols that (circularly) precede the first symbol of each suffix in the list of (lexicographically) sorted suffixes of  $S'$ .

$eBWT(S') = CCTCAGATCGTGG\$_2\$_1ACTC\$_3C$

Note that the colors and suffixes are only for illustrative purposes.

# How does EBWT work?

Let  $S = \{S_1, S_2, \dots, S_m\}$  be a collection of strings.

- Append an end-marker  $\$$ , where  $\$_i < A < C < G < T$  (we add subscripts  $\$_i$  only for illustrative purposes), to each string in  $S$  by obtaining a new collection  $S'$ .

String collection $S'$							
	0	1	2	3	4	5	6
$S_1$	G	C	C	A	A	C	$\$_1$
$S_2$	G	A	G	C	T	C	$\$_2$
$S_3$	T	C	G	C	T	T	$\$_3$

- Sort all the suffixes of the strings in  $S'$ ;
- Take the string  $eBWT(S')$  obtained by concatenating the symbols that (circularly) precede the first symbol of each suffix in the list of (lexicographically) sorted suffixes of  $S'$ .

$eBWT(S') = CCTCAGATCGTGG\$_2\$_1ACTC\$_3C$

Note that the colors and suffixes are only for illustrative purposes.

$eBWT(S')$	Sorted Suffixes of $S'$
C	$\$_1$
C	$\$_2$
T	$\$_3$
C	AAC $\$_1$
A	AC $\$_1$
G	AGCTC $\$_2$
A	C $\$_1$
T	C $\$_2$
C	CAAC $\$_1$
G	CCAAC $\$_1$
T	CGCTT $\$_3$
G	CTC $\$_2$
G	CTT $\$_3$
$\$_2$	GAGCTC $\$_2$
$\$_1$	GCCAAC $\$_1$
A	GCTC $\$_2$
C	GCTT $\$_3$
T	T $\$_3$
C	TC $\$_2$
$\$_3$	TCGCTT $\$_3$
C	TT $\$_3$

# How does EBWT work?

Let  $S = \{S_1, S_2, \dots, S_m\}$  be a collection of strings.

- Append an end-marker  $\$$ , where  $\$_i < A < C < G < T$  (we add subscripts  $\$_i$  only for illustrative purposes), to each string in  $S$  by obtaining a new collection  $S'$ .

String collection $S'$							
	0	1	2	3	4	5	6
$S_1$	G	C	C	A	A	C	$\$_1$
$S_2$	G	A	G	C	T	C	$\$_2$
$S_3$	T	C	G	C	T	T	$\$_3$

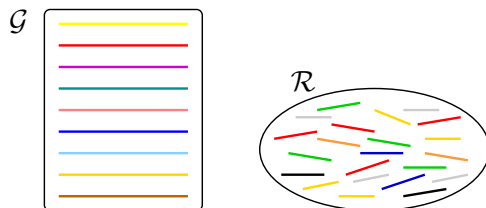
- Sort all the suffixes of the strings in  $S'$ ;
- Take the string  $eBWT(S')$  obtained by concatenating the symbols that (circularly) precede the first symbol of each suffix in the list of (lexicographically) sorted suffixes of  $S'$ .

$$eBWT(S') = CCTCAGATCGTGG\$_2\$_1ACTC\$_3C$$

Note that the colors and suffixes are only for illustrative purposes.

$eBWT(S')$	Sorted Suffixes of $S'$
C	$\$_1$
C	$\$_2$
T	$\$_3$
C	AAC $\$_1$
A	AC $\$_1$
G	AGCTC $\$_2$
A	C $\$_1$
T	C $\$_2$
C	CAAC $\$_1$
G	CCAAC $\$_1$
T	CGCTT $\$_3$
G	CTC $\$_2$
G	CTT $\$_3$
$\$_2$	GAGCTC $\$_2$
$\$_1$	GCCAAC $\$_1$
A	GCTC $\$_2$
C	GCTT $\$_3$
T	T $\$_3$
C	TC $\$_2$
$\$_3$	TCGCTT $\$_3$
C	TT $\$_3$

# Formalization: Metagenomic classification problem



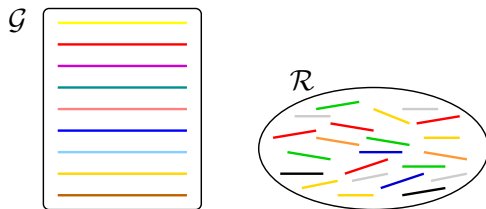
- $\mathcal{R} = \{r_1, \dots, r_{|\mathcal{R}|}\}$  metagenome (collection of short reads)
- $\mathcal{G} = \{g_1, \dots, g_{|\mathcal{G}|}\}$  reference genomes (collection of long sequences)

**Pre-processing step:** build the data structures  $eBWT(\mathcal{S})$ ,  $DA(\mathcal{S})$ ,  $LCP(\mathcal{S})$ , where  $\mathcal{S} = \mathcal{R} \cup \mathcal{G}$ .

The data structures for the reference database can be built only once, and updated with the data structures for the reads collection in order to obtain the data structures for  $\mathcal{S}$ .

**Goal:** to assign each read  $r_i$  in  $\mathcal{R}$  to a unique genome  $g_j$  in  $\mathcal{G}$

# Formalization: Metagenomic classification problem



- $\mathcal{R} = \{r_1, \dots, r_{|\mathcal{R}|}\}$  metagenome (collection of short reads)
- $\mathcal{G} = \{g_1, \dots, g_{|\mathcal{G}|}\}$  reference genomes (collection of long sequences)

**Pre-processing step:** build the data structures  $eBWT(\mathcal{S})$ ,  $DA(\mathcal{S})$ ,  $LCP(\mathcal{S})$ , where  $\mathcal{S} = \mathcal{R} \cup \mathcal{G}$ .

The data structures for the reference database can be built only once, and updated with the data structures for the reads collection in order to obtain the data structures for  $\mathcal{S}$ .

**Goal:** to assign each read  $r_i$  in  $\mathcal{R}$  to a unique genome  $g_j$  in  $\mathcal{G}$

# Preprocessing: eBWT, DA and LCP array

$DA[i]$ : index of the string to which the  $i$ -th suffix belongs.

$LCP[i]$ : length of the Longest Common Prefix between the  $i$ -th and the  $(i - 1)$ -th suffix.

$S = \{GGCGTACCA\$1, GGGGCGTAT\$2, ACGATTAGC\$3\}$

$DA$	$LCP$	$eBWT$	Sorted suffixes	$DA$	$LCP$	$eBWT$	Sorted suffixes
1	0	A	$\$1$	3	0	C	$GATTAGC\$3$
2	0	T	$\$2$	3	1	A	$GC\$3$
3	0	C	$\$3$	1	2	G	$GCGTACCA\$1$
1	0	C	$A\$1$	2	5	G	$GCGTAT\$2$
1	1	T	$ACCA\$1$	1	1	$\$1$	$GGCGTACCA\$1$
3	2	$\$3$	$ACGATTAGC\$3$	2	6	G	$GGCGTAT\$2$
3	1	T	$AGC\$3$	2	2	G	$GGGCGTAT\$2$
2	1	T	$AT\$2$	2	3	$\$2$	$GGGGCGTAT\$2$
3	2	G	$ATTAGC\$3$	1	1	C	$GTACCA\$1$
3	0	G	$C\$3$	2	3	C	$GTAT\$2$
1	1	C	$CA\$1$	2	0	A	$T\$2$
1	1	A	$CCA\$1$	1	1	G	$TACCA\$1$
3	1	A	$CGATTAGC\$3$	3	2	T	$TAGC\$3$
1	2	G	$CGTACCA\$1$	2	2	G	$TAT\$2$
2	4	G	$CGTAT\$2$	3	1	A	$TTAGC\$3$

## Preprocessing: eBWT, DA and LCP array

$DA[i]$ : index of the string to which the  $i$ -th suffix belongs.

$LCP[i]$ : length of the Longest Common Prefix between the  $i$ -th and the  $(i - 1)$ -th suffix.

$S = \{GGCGTACCA\$1, GGGGCGTAT\$2, ACGATTAGC\$3\}$

$DA$	$LCP$	$eBWT$	Sorted suffixes	$DA$	$LCP$	$eBWT$	Sorted suffixes
1	0	A	$\$1$	3	0	C	$GATTAGC\$3$
2	0	T	$\$2$	3	1	A	$GC\$3$
3	0	C	$\$3$	1	2	G	$GCGTACCA\$1$
1	0	C	$A\$1$	2	5	G	$GCGTAT\$2$
1	1	T	$ACCA\$1$	1	1	$\$1$	$GGCGTACCA\$1$
3	2	$\$3$	$ACGATTAGC\$3$	2	6	G	$GGCGTAT\$2$
3	1	T	$AGC\$3$	2	2	G	$GGGCGTAT\$2$
2	1	T	$AT\$2$	2	3	$\$2$	$GGGGCGTAT\$2$
3	2	G	$ATTAGC\$3$	1	1	C	$GTACCA\$1$
3	0	G	$C\$3$	2	3	C	$GTAT\$2$
1	1	C	$CA\$1$	2	0	A	$T\$2$
1	1	A	$CCA\$1$	1	1	G	$TACCA\$1$
3	1	A	$CGATTAGC\$3$	3	2	T	$TAGC\$3$
1	2	G	$CGTACCA\$1$	2	2	G	$TAT\$2$
2	4	G	$CGTAT\$2$	3	1	A	$TTAGC\$3$

## Preprocessing: eBWT, DA and LCP array

$DA[i]$ : index of the string to which the  $i$ -th suffix belongs.

$LCP[i]$ : length of the Longest Common Prefix between the  $i$ -th and the  $(i - 1)$ -th suffix.

$S = \{GGCGTACCA\$1, GGGGCGTAT\$2, ACGATTAGC\$3\}$

$DA$	$LCP$	$eBWT$	Sorted suffixes	$DA$	$LCP$	$eBWT$	Sorted suffixes
1	0	A	$\$1$	3	0	C	$GATTAGC\$3$
2	0	T	$\$2$	3	1	A	$GC\$3$
3	0	C	$\$3$	1	2	G	$GCGTACCA\$1$
1	0	C	$A\$1$	2	5	G	$GCGTAT\$2$
1	1	T	$ACCA\$1$	1	1	$\$1$	$GGCGTACCA\$1$
3	2	$\$3$	$ACGATTAGC\$3$	2	6	G	$GGCGTAT\$2$
3	1	T	$AGC\$3$	2	2	G	$GGGCGTAT\$2$
2	1	T	$AT\$2$	2	3	$\$2$	$GGGCGTAT\$2$
3	2	G	$ATTAGC\$3$	1	1	C	$GTACCA\$1$
3	0	G	$C\$3$	2	3	C	$GTAT\$2$
1	1	C	$CA\$1$	2	0	A	$T\$2$
1	1	A	$CCA\$1$	1	1	G	$TACCA\$1$
3	1	A	$CGATTAGC\$3$	3	2	T	$TAGC\$3$
1	2	G	$CGTACCA\$1$	2	2	G	$TAT\$2$
2	4	G	$CGTAT\$2$	3	1	A	$TTAGC\$3$



# Preprocessing: eBWT, DA and LCP array

$DA[i]$ : index of the string to which the  $i$ -th suffix belongs.

$LCP[i]$ : length of the Longest Common Prefix between the  $i$ -th and the  $(i - 1)$ -th suffix.

$S = \{GGCGTACCA\$_1, GGGGCGTAT\$_2, ACGATTAGC\$_3\}$

$DA$	$LCP$	$eBWT$	$DA$	$LCP$	$eBWT$
1	0	A	3	0	C
2	0	T	3	1	A
3	0	C	1	2	G
1	0	C	2	5	G
1	1	T	1	1	$\$1$
3	2	$\$3$	2	6	G
3	1	T	2	2	G
2	1	T	2	3	$\$2$
3	2	G	1	1	C
3	0	G	2	3	C
1	1	C	2	0	A
1	1	A	1	1	G
3	1	A	3	2	T
1	2	G	2	2	G
2	4	G	3	1	A

# Property of the eBWT: Mixing and clustering effect

## Intuitive idea

The greater is the number of substrings share by two strings, the smaller is their “distance”

## Key property of the eBWT

The greater is the number of substrings shared by  $u$  and  $v$ , the greater is the mixing of the suffixes of  $u$  and  $v$  in the sorted list and the greater are the runs (clusters) of the same symbol in the eBWT.

$u = \underline{GGCGTACCA}\$1.$      $v = \underline{GGGGCGTAT}\$2$

<i>eBWT</i>	Sorted suffixes
A	$\$1$
T	$\$2$
C	$A\$1$
T	$ACCA\$1$
T	$AT\$2$
C	$CA\$1$
A	$CCA\$1$
G	$CGTACCA\$1$
G	$CGTAT\$2$
G	$GCGTACCA\$1$
G	$GCGTAT\$2$
$\$1$	$GGCGTACCA\$1$
G	$GGCGTAT\$2$
G	$GGGGCGTAT\$2$
$\$2$	$GGGGCGTAT\$2$
C	$GTACCA\$1$
C	$GTAT\$2$
A	$T\$2$
G	$TACCA\$1$
G	$TAT\$2$

# Property of the eBWT: Mixing and clustering effect

## Intuitive idea

The greater is the number of substrings share by two strings, the smaller is their “distance”

## Key property of the eBWT

The greater is the number of substrings shared by  $u$  and  $v$ , the greater is the mixing of the suffixes of  $u$  and  $v$  in the sorted list and the greater are the runs (clusters) of the same symbol in the eBWT.

$u = \underline{GGCGTACCA}\$1$ .     $v = \underline{GGGGCGTAT}\$2$

<i>eBWT</i>	Sorted suffixes
A	$\$1$
T	$\$2$
C	$A\$1$
T	$ACCA\$1$
T	$AT\$2$
C	$CA\$1$
A	$CCA\$1$
G	$CGTACCA\$1$
G	$CGTAT\$2$
G	$GCGTACCA\$1$
G	$GCGTAT\$2$
$\$1$	$GGCGTACCA\$1$
G	$GGCGTAT\$2$
G	$GGGCGTAT\$2$
$\$2$	$GGGGCGTAT\$2$
C	$GTACCA\$1$
C	$GTAT\$2$
A	$T\$2$
G	$TACCA\$1$
G	$TAT\$2$

# Property of the eBWT: Mixing and clustering effect

## Intuitive idea

The greater is the number of substrings share by two strings, the smaller is their “distance”

## Key property of the eBWT

The greater is the number of substrings shared by  $u$  and  $v$ , the greater is the mixing of the suffixes of  $u$  and  $v$  in the sorted list and the greater are the runs (clusters) of the same symbol in the eBWT.

$u = \underline{GGCGTACCA}\$1.$      $v = \underline{GGGGCGTAT}\$2$

<i>eBWT</i>	Sorted suffixes
A	$\$1$
T	$\$2$
C	$A\$1$
T	$ACCA\$1$
T	$AT\$2$
C	$CA\$1$
A	$CCA\$1$
G	$CGTACCA\$1$
G	$CGTAT\$2$
G	$GCGTACCA\$1$
G	$GCGTAT\$2$
$\$1$	$GGCGTACCA\$1$
G	$GGCGTAT\$2$
G	$GGGCGTAT\$2$
$\$2$	$GGGGCGTAT\$2$
C	$GTACCA\$1$
C	$GTAT\$2$
A	$T\$2$
G	$TACCA\$1$
G	$TAT\$2$

# Property of the eBWT: Mixing and clustering effect

## Intuitive idea

The greater is the number of substrings share by two strings, the smaller is their “distance”

## Key property of the eBWT

The greater is the number of substrings shared by  $u$  and  $v$ , the greater is the mixing of the suffixes of  $u$  and  $v$  in the sorted list and the greater are the runs (clusters) of the same symbol in the eBWT.

$u = \underline{GGCGTACCA}\$1$ .     $v = \underline{GGGGCGTAT}\$2$

<i>eBWT</i>	Sorted suffixes
A	$\$1$
T	$\$2$
C	$A\$1$
T	$ACCA\$1$
T	$AT\$2$
C	$CA\$1$
A	$CCA\$1$
G	$CGTACCA\$1$
G	$CGTAT\$2$
G	$GCGTACCA\$1$
G	$GCGTAT\$2$
$\$1$	$GGCGTACCA\$1$
G	$GGCGTAT\$2$
G	$GGGCGTAT\$2$
$\$2$	$GGGGCGTAT\$2$
C	$GTACCA\$1$
C	$GTAT\$2$
A	$T\$2$
G	$TACCA\$1$
G	$TAT\$2$

# Property of the eBWT: Mixing and clustering effect

## Intuitive idea

The greater is the number of substrings share by two strings, the smaller is their “distance”

## Key property of the eBWT

The greater is the number of substrings shared by  $u$  and  $v$ , the greater is the mixing of the suffixes of  $u$  and  $v$  in the sorted list and the greater are the runs (clusters) of the same symbol in the eBWT.

$u = \underline{GGCGTACCA}\$1.$      $v = \underline{GGGGCGTAT}\$2$

<i>eBWT</i>	Sorted suffixes
A	$\$1$
T	$\$2$
C	$A\$1$
T	$ACCA\$1$
T	$AT\$2$
C	$CA\$1$
A	$CCA\$1$
G	$CGTACCA\$1$
G	$CGTAT\$2$
G	$GCGTACCA\$1$
G	$GCGTAT\$2$
$\$1$	$GGCGTACCA\$1$
G	$GGCGTAT\$2$
G	$GGGCGTAT\$2$
$\$2$	$GGGGCGTAT\$2$
C	$GTACCA\$1$
C	$GTAT\$2$
A	$T\$2$
G	$TACCA\$1$
G	$TAT\$2$

# Property of the eBWT: Mixing and clustering effect

## Intuitive idea

The greater is the number of substrings share by two strings, the smaller is their “distance”

## Key property of the eBWT

The greater is the number of substrings shared by  $u$  and  $v$ , the greater is the mixing of the suffixes of  $u$  and  $v$  in the sorted list and the greater are the runs (clusters) of the same symbol in the eBWT.

$u = \underline{GGCGTACCA}\$1$ .     $v = \underline{GGG}\underline{GCGTAT}\$2$

<i>eBWT</i>	Sorted suffixes
A	$\$1$
T	$\$2$
C	$A\$1$
T	$ACCA\$1$
T	$AT\$2$
C	$CA\$1$
A	$CCA\$1$
G	$CGTACCA\$1$
G	$CGTAT\$2$
G	$GCGTACCA\$1$
G	$GCGTAT\$2$
$\$1$	$GGCGTACCA\$1$
G	$GGCGTAT\$2$
G	$GGGCGTAT\$2$
$\$2$	$GGGGCGTAT\$2$
C	$GTACCA\$1$
C	$GTAT\$2$
A	$T\$2$
G	$TACCA\$1$
G	$TAT\$2$

# Property of the eBWT: Mixing and clustering effect

## Intuitive idea

The greater is the number of substrings share by two strings, the smaller is their “distance”

## Key property of the eBWT

The greater is the number of substrings shared by  $u$  and  $v$ , the greater is the mixing of the suffixes of  $u$  and  $v$  in the sorted list and the greater are the runs (clusters) of the same symbol in the eBWT.

$u = \underline{GGCGTACCA}\$1$ .     $v = \underline{GGGGCGTAT}\$2$

<i>eBWT</i>	Sorted suffixes
A	$\$1$
T	$\$2$
C	$A\$1$
T	$ACCA\$1$
T	$AT\$2$
C	$CA\$1$
A	$CCA\$1$
G	$CGTACCA\$1$
G	$CGTAT\$2$
G	$GCGTACCA\$1$
G	$GCGTAT\$2$
$\$1$	$GGCGTACCA\$1$
G	$GGCGTAT\$2$
G	$GGGCGTAT\$2$
$\$2$	$GGGGCGTAT\$2$
C	$GTACCA\$1$
C	$GTAT\$2$
A	$T\$2$
G	$TACCA\$1$
G	$TAT\$2$



# Step 1: Build $\alpha$ -clusters and Similarity Arrays (Part I)

Minimum  $LCP$  value  $\alpha = 3$

$i$	$LCP$	$eBWT$	Sorted suffixes
1	0	A	$\$i$
2	0	T	$\$j$
3	0	C	$A\$i$
4	1	T	$ACCA\$i$
5	1	T	$AGTTT\$j$
6	1	T	$ATGTATTAGTTT\$j$
7	2	T	$ATTAGTTT\$j$
8	0	C	$CA\$i$
9	1	A	$CCA\$i$
10	1	G	$CGGGGCGTA\dots\$j$
11	2	G	$CGTACCA\$i$
12	4	G	$CGTATGAT\dots\$j$
13	1	T	$CTTTTGGCG\dots\$j$
14	0	G	$GCGGGGCGT\dots\$j$
15	3	G	$GCGTACCA\$i$
16	5	G	$GCGTATGTAA\dots\$j$
$\vdots$	$\vdots$	$\vdots$	$\vdots$

$$r_i = \text{KGGCGTACCA}\$i$$

$$g_j = \text{TTATTTGGCGGGGCGTATGTATTAGTTT}\$j$$

**Detect blocks ( $\alpha$ -clusters) of  $eBWT(S)$ :** An  $\alpha$ -cluster  $\mathcal{C}_\alpha$  of  $eBWT(S)$  is any pair of indices  $(pS, pE)$  such that

- $LCP[pS] < \alpha$  and  $LCP[pE + 1] < \alpha$ ,
- $LCP[k] \geq \alpha$ ,  $pS < k \leq pE$ ,
- $DA[s] \in \mathcal{R}$  and  $DA[t] \in \mathcal{G}$ ,  $pS \leq s, t \leq pE$ .

$$\mathcal{C}_\alpha(r_i, g_j) = \{$$

**Compute the similarity between any read and any genome:**

$$Sim_r[g] = \sum_{x \in \mathcal{C}_\alpha} \sum_{a \in \Sigma} \min(n_r, n_g)$$

$n_r$  = number of indices  $s$  in  $x$  such that  $eBWT[s] = a$  and  $DA[s] = r$ ,

$n_g$  = number of indices  $t$  in  $x$  such that  $eBWT[t] = a$  and  $DA[t] = g$ .

i.e. the total number of bases of  $r$  and  $g$  that match (by using IUPAC list) in each  $\alpha$ -cluster.

$$Sim_{r_i}[g_j] =$$

# Step 1: Build $\alpha$ -clusters and Similarity Arrays (Part I)

Minimum *LCP* value  $\alpha = 3$

$i$	<i>LCP</i>	<i>eBWT</i>	Sorted suffixes
1	0	A	$\$i$
2	0	T	$\$j$
3	0	C	$A\$i$
4	1	T	$ACCA\$i$
5	1	T	$AGTTT\$j$
6	1	T	$ATGTATTAGTTT\$j$
7	2	T	$ATTAGTTT\$j$
8	0	C	$CAS\$i$
9	1	A	$CCA\$i$
10	1	G	$CGGGGCGTA\dots\$j$
11	2	G	$CGTACCA\$i$
12	4	G	$CGTATGAT\dots\$j$
13	1	T	$CTTTTGGCG\dots\$j$
14	0	G	$GCGGGGCGT\dots\$j$
15	3	G	$GCGTACCA\$i$
16	5	G	$GCGTATGTAA\dots\$j$
$\vdots$	$\vdots$	$\vdots$	$\vdots$

$$r_i = \text{KGGCGTACCA}\$i$$

$$g_j = \text{TTATTTTGGCGGGGCGTATGTATTAGTTT}\$j$$

**Detect blocks ( $\alpha$ -clusters) of *eBWT*(*S*):** An  $\alpha$ -cluster  $\mathcal{C}_\alpha$  of *eBWT*(*S*) is any pair of indices ( $pS, pE$ ) such that

- $LCP[pS] < \alpha$  and  $LCP[pE + 1] < \alpha$ ,
- $LCP[k] \geq \alpha$ ,  $pS < k \leq pE$ ,
- $DA[s] \in \mathcal{R}$  and  $DA[t] \in \mathcal{G}$ ,  $pS \leq s, t \leq pE$ .

$$\mathcal{C}_\alpha(r_i, g_j) = \{(11, 12),$$

**Compute the similarity between any read and any genome:**

$$Sim_r[g] = \sum_{x \in \mathcal{C}_\alpha} \sum_{a \in \Sigma} \min(n_r, n_g)$$

$n_r$  = number of indices  $s$  in  $x$  such that  $eBWT[s] = a$  and  $DA[s] = r$ ,  
 $n_g$  = number of indices  $t$  in  $x$  such that  $eBWT[t] = a$  and  $DA[t] = g$ .

i.e. the total number of bases of  $r$  and  $g$  that match (by using IUPAC list) in each  $\alpha$ -cluster.

$$Sim_{r_i}[g_j] = 1 +$$

# Step 1: Build $\alpha$ -clusters and Similarity Arrays (Part I)

Minimum *LCP* value  $\alpha = 3$

$i$	<i>LCP</i>	<i>eBWT</i>	Sorted suffixes
1	0	A	$\$i$
2	0	T	$\$j$
3	0	C	$A\$i$
4	1	T	$ACCA\$i$
5	1	T	$AGTTT\$j$
6	1	T	$ATGTATTAGTTT\$j$
7	2	T	$ATTAGTTT\$j$
8	0	C	$CAS\$i$
9	1	A	$CCA\$i$
10	1	G	$CGGGGCGTA\dots\$j$
11	2	G	$CGTACCA\$i$
12	4	G	$CGTATGAT\dots\$j$
13	1	T	$CTTTTGGCG\dots\$j$
14	0	G	$GCGGGGCGT\dots\$j$
15	3	G	$GCGTACCA\$i$
16	5	G	$GCGTATGTAA\dots\$j$
$\vdots$	$\vdots$	$\vdots$	$\vdots$

$$r_i = \text{KGGCGTACCA}\$i$$

$$g_j = \text{TTATTTTGGCGGGGCGTATGTATTAGTTT}\$j$$

**Detect blocks ( $\alpha$ -clusters) of *eBWT*(*S*):** An  $\alpha$ -cluster  $C_\alpha$  of *eBWT*(*S*) is any pair of indices ( $pS, pE$ ) such that

- $LCP[pS] < \alpha$  and  $LCP[pE + 1] < \alpha$ ,
- $LCP[k] \geq \alpha$ ,  $pS < k \leq pE$ ,
- $DA[s] \in \mathcal{R}$  and  $DA[t] \in \mathcal{G}$ ,  $pS \leq s, t \leq pE$ .

$$C_\alpha(r_i, g_j) = \{(11, 12), (14, 16), \dots\}$$

**Compute the similarity between any read and any genome:**

$$Sim_r[g] = \sum_{x \in C_\alpha} \sum_{a \in \Sigma} \min(n_r, n_g)$$

$n_r$  = number of indices  $s$  in  $x$  such that  $eBWT[s] = a$  and  $DA[s] = r$ ,

$n_g$  = number of indices  $t$  in  $x$  such that  $eBWT[t] = a$  and  $DA[t] = g$ .

i.e. the total number of bases of  $r$  and  $g$  that match (by using IUPAC list) in each  $\alpha$ -cluster.

$$Sim_{r_i}[g_j] = 1+1+$$

# Step 1: Build $\alpha$ -clusters and Similarity Arrays (Part I)

Minimum *LCP* value  $\alpha = 3$

$i$	<i>LCP</i>	<i>eBWT</i>	Sorted suffixes
1	0	A	$\$i$
2	0	T	$\$j$
3	0	C	$A\$i$
4	1	T	$ACCA\$i$
5	1	T	$AGTTT\$j$
6	1	T	$ATGTATTAGTTT\$j$
7	2	T	$ATTAGTTT\$j$
8	0	C	$CA\$i$
9	1	A	$CCA\$i$
10	1	G	$CGGGGCGTA\dots\$j$
11	2	G	$CGTACCA\$i$
12	4	G	$CGTATGAT\dots\$j$
13	1	T	$CTTTTGGCG\dots\$j$
14	0	G	$GCGGGGCGT\dots\$j$
15	3	G	$GCGTACCA\$i$
16	5	G	$GCGTATGTAA\dots\$j$
$\vdots$	$\vdots$	$\vdots$	$\vdots$

$$r_i = \text{KGGCGTACCA}\$i$$

$$g_j = \text{TTATTTTGGCGGGGCGTATGTATTAGTTT}\$j$$

**Detect blocks ( $\alpha$ -clusters) of *eBWT*(*S*):** An  $\alpha$ -cluster  $C_\alpha$  of *eBWT*(*S*) is any pair of indices ( $pS, pE$ ) such that

- $LCP[pS] < \alpha$  and  $LCP[pE + 1] < \alpha$ ,
- $LCP[k] \geq \alpha$ ,  $pS < k \leq pE$ ,
- $DA[s] \in \mathcal{R}$  and  $DA[t] \in \mathcal{G}$ ,  $pS \leq s, t \leq pE$ .

$$C_\alpha(r_i, g_j) = \{(11, 12), (14, 16), \dots\}$$

**Compute the similarity between any read and any genome:**

$$Sim_r[g] = \sum_{x \in C_\alpha} \sum_{a \in \Sigma} \min(n_r, n_g)$$

$n_r$  = number of indices  $s$  in  $x$  such that  $eBWT[s] = a$  and  $DA[s] = r$ ,

$n_g$  = number of indices  $t$  in  $x$  such that  $eBWT[t] = a$  and  $DA[t] = g$ .

i.e. the total number of bases of  $r$  and  $g$  that match (by using IUPAC list) in each  $\alpha$ -cluster.

$$Sim_{r_i}[g_j] = 1+1+\dots$$

# Step 1: Build $\alpha$ -clusters and Similarity Arrays (Part II)

Minimum  $LCP$  value  $\alpha = 3$

$i$	$LCP$	$eBWT$	Sorted suffixes
⋮	⋮	⋮	⋮
17	1	T	GGCGGGGCG...\$ <sub>j</sub>
18	4	K	GGCGTACCA\$ <sub>i</sub>
19	6	G	GGCGTATGTAT...\$ <sub>j</sub>
20	2	G	GGGCGTAT...\$ <sub>j</sub>
21	3	C	GGGGCGTAT...\$ <sub>j</sub>
22	1	C	GTACCA\$ <sub>i</sub>
23	3	C	GTATGTA...\$ <sub>j</sub>
24	4	C	GTATTA...\$ <sub>j</sub>
25	2	A	GTTT\$ <sub>j</sub>
26	0	\$ <sub>i</sub>	KGGCGTACCA\$ <sub>i</sub>
27	0	T	T\$ <sub>j</sub>
28	1	G	TACCA\$ <sub>i</sub>
29	2	T	TAGTTT\$ <sub>j</sub>
⋮	⋮	⋮	⋮

$r_i = \text{KGGCGTACCA$}_i$   
 $g_j = \text{TTATTTTGGCGGGGCGTATGTATTAGTTT$}_j$

**Detect blocks ( $\alpha$ -clusters) of  $eBWT(S)$ :** An  $\alpha$ -cluster  $C_\alpha$  of  $eBWT(S)$  is any pair of indices  $(pS, pE)$  such that

- $LCP[pS] < \alpha$  and  $LCP[pE + 1] < \alpha$ ,
- $LCP[k] \geq \alpha$ ,  $pS < k \leq pE$ ,
- $DA[i] \in \mathcal{R}$  and  $DA[j] \in \mathcal{G}$ ,  $pS \leq i, j \leq pE$ .

$C_\alpha = \{(11, 12), (14, 16),$

**Compute the similarity between any read and any genome:**

$$Sim_r[g] = \sum_{x \in C_\alpha} \sum_{a \in \Sigma} \min(n_r, n_g)$$

$n_r$  = number of indices  $s$  in  $x$  such that  $eBWT[s] = a$  and  $DA[s] = r$ ,

$n_g$  = number of indices  $t$  in  $x$  such that  $eBWT[t] = a$  and  $DA[t] = g$ .

i.e. the total number of bases of  $r$  and  $g$  that match (by using IUPAC list) in each  $\alpha$ -cluster.

$$Sim_r[g] = 1+1+$$

# Step 1: Build $\alpha$ -clusters and Similarity Arrays (Part II)

Minimum  $LCP$  value  $\alpha = 3$

$i$	$LCP$	$eBWT$	Sorted suffixes
⋮	⋮	⋮	⋮
17	1	T	GGCGGGGCG...\$ <sub>j</sub>
18	4	K	GGCGTACCA\$ <sub>i</sub>
19	6	G	GGCGTATGTAT...\$ <sub>j</sub>
20	2	G	GGGCGTAT...\$ <sub>j</sub>
21	3	C	GGGGCGTAT...\$ <sub>j</sub>
22	1	C	GTACCA\$ <sub>i</sub>
23	3	C	GTATGTA...\$ <sub>j</sub>
24	4	C	GTATTA...\$ <sub>j</sub>
25	2	A	GTTT\$ <sub>j</sub>
26	0	\$ <sub>i</sub>	KGGCGTACCA\$ <sub>i</sub>
27	0	T	T\$ <sub>j</sub>
28	1	G	TACCA\$ <sub>i</sub>
29	2	T	TAGTTT\$ <sub>j</sub>
⋮	⋮	⋮	⋮

$K \rightarrow \{G \text{ or } T\}$  in the IUPAC list.

$r_i = \text{KGGCGTACCA$}_i$   
 $g_j = \text{TTATTTTGGCGGGGCGTATGTATTAGTTT$}_j$

**Detect blocks ( $\alpha$ -clusters) of  $eBWT(S)$ :** An  $\alpha$ -cluster  $C_\alpha$  of  $eBWT(S)$  is any pair of indices  $(pS, pE)$  such that

- $LCP[pS] < \alpha$  and  $LCP[pE + 1] < \alpha$ ,
- $LCP[k] \geq \alpha$ ,  $pS < k \leq pE$ ,
- $DA[i] \in \mathcal{R}$  and  $DA[j] \in \mathcal{G}$ ,  $pS \leq i, j \leq pE$ .

$C_\alpha = \{(11, 12), (14, 16), (17, 19)\}$ ,

**Compute the similarity between any read and any genome:**

$$Sim_r[g] = \sum_{x \in C_\alpha} \sum_{a \in \Sigma} \min(n_r, n_g)$$

$n_r$  = number of indices  $s$  in  $x$  such that  $eBWT[s] = a$  and  $DA[s] = r$ ,

$n_g$  = number of indices  $t$  in  $x$  such that  $eBWT[t] = a$  and  $DA[t] = g$ .

i.e. the total number of bases of  $r$  and  $g$  that match (by using IUPAC list) in each  $\alpha$ -cluster.

$$Sim_r[g] = 1+1+1+$$

# Step 1: Build $\alpha$ -clusters and Similarity Arrays (Part II)

Minimum *LCP* value  $\alpha = 3$

$i$	<i>LCP</i>	<i>eBWT</i>	Sorted suffixes
⋮	⋮	⋮	⋮
17	1	T	GGCGGGGCG...\$ <sub>j</sub>
18	4	K	GGCGTACCA\$ <sub>i</sub>
19	6	G	GGCGTATGTAT...\$ <sub>j</sub>
20	2	G	GGGCGTAT...\$ <sub>j</sub>
21	3	C	GGGGCGTAT...\$ <sub>j</sub>
22	1	C	GTACCA\$ <sub>i</sub>
23	3	C	GTATGTA...\$ <sub>j</sub>
24	4	C	GTATTA...\$ <sub>j</sub>
25	2	A	GTTT\$ <sub>j</sub>
26	0	\$ <sub>i</sub>	KGGCGTACCA\$ <sub>i</sub>
27	0	T	T\$ <sub>j</sub>
28	1	G	TACCA\$ <sub>i</sub>
29	2	T	TAGTTT\$ <sub>j</sub>
⋮	⋮	⋮	⋮

it is not cluster (only blue symbols)

$r_i = \text{KGGCGTACCA}\$_i$   
 $g_j = \text{TTATTTTGGCGGGGCGTATGTATTAGTTT}\$_j$

**Detect blocks ( $\alpha$ -clusters) of *eBWT*(*S*):** An  $\alpha$ -cluster  $C_\alpha$  of *eBWT*(*S*) is any pair of indices ( $pS, pE$ ) such that

- $LCP[pS] < \alpha$  and  $LCP[pE + 1] < \alpha$ ,
- $LCP[k] \geq \alpha$ ,  $pS < k \leq pE$ ,
- $DA[i] \in \mathcal{R}$  and  $DA[j] \in \mathcal{G}$ ,  $pS \leq i, j \leq pE$ .

$C_\alpha = \{(11, 12), (14, 16), (17, 19)\}$ ,

**Compute the similarity between any read and any genome:**

$$Sim_r[g] = \sum_{x \in C_\alpha} \sum_{a \in \Sigma} \min(n_r, n_g)$$

$n_r$  = number of indices  $s$  in  $x$  such that  $eBWT[s] = a$  and  $DA[s] = r$ ,

$n_g$  = number of indices  $t$  in  $x$  such that  $eBWT[t] = a$  and  $DA[t] = g$ .

i.e. the total number of bases of  $r$  and  $g$  that match (by using IUPAC list) in each  $\alpha$ -cluster.

$$Sim_r[g] = 1+1+1+$$

# Step 1: Build $\alpha$ -clusters and Similarity Arrays (Part II)

Minimum  $LCP$  value  $\alpha = 3$

$i$	$LCP$	$eBWT$	Sorted suffixes
⋮	⋮	⋮	⋮
17	1	T	GGCGGGGCG...\$ <sub>j</sub>
18	4	K	GGCGTACCA\$ <sub>i</sub>
19	6	G	GGCGTATGTAT...\$ <sub>j</sub>
20	2	G	GGGCGTAT...\$ <sub>j</sub>
21	3	C	GGGGCGTAT...\$ <sub>j</sub>
22	1	C	GTACCA\$ <sub>i</sub>
23	3	C	GTATGTA...\$ <sub>j</sub>
24	4	C	GTATTA...\$ <sub>j</sub>
25	2	A	GTTT\$ <sub>j</sub>
26	0	\$ <sub>i</sub>	KGGCGTACCA\$ <sub>i</sub>
27	0	T	T\$ <sub>j</sub>
28	1	G	TACCA\$ <sub>i</sub>
29	2	T	TAGTTT\$ <sub>j</sub>
⋮	⋮	⋮	⋮

$$r_i = \text{KGGCGTACCA}\$i$$

$$g_j = \text{TTATTTTGGCGGGGCGTATGTATTAGTTT}\$j$$

**Detect blocks ( $\alpha$ -clusters) of  $eBWT(S)$ :** An  $\alpha$ -cluster  $C_\alpha$  of  $eBWT(S)$  is any pair of indices  $(pS, pE)$  such that

- $LCP[pS] < \alpha$  and  $LCP[pE + 1] < \alpha$ ,
- $LCP[k] \geq \alpha$ ,  $pS < k \leq pE$ ,
- $DA[i] \in \mathcal{R}$  and  $DA[j] \in \mathcal{G}$ ,  $pS \leq i, j \leq pE$ .

$$C_\alpha = \{(11, 12), (14, 16), (17, 19), (22, 24)\}$$

**Compute the similarity between any read and any genome:**

$$Sim_r[g] = \sum_{x \in C_\alpha} \sum_{a \in \Sigma} \min(n_r, n_g)$$

$n_r$  = number of indices  $s$  in  $x$  such that  $eBWT[s] = a$  and  $DA[s] = r$ ,

$n_g$  = number of indices  $t$  in  $x$  such that  $eBWT[t] = a$  and  $DA[t] = g$ .

i.e. the total number of bases of  $r$  and  $g$  that match (by using IUPAC list) in each  $\alpha$ -cluster.

$$Sim_r[g] = 1+1+1+1 = 4$$



## Step 2: Classification

Given a threshold value  $\beta$ , the read  $r_i$  is

- **assigned** to  $g_j$  if  $g_j$  is the only genome such that  $Sim_{r_i}[g_j] \sim \max_g Sim_{r_i}[g]$  and  $Sim_{r_i}[g_j] > \beta$ .
- **not classified** if  $\max_g Sim_{r_i}[g] \leq \beta$ .
- **ambiguous** if  $\max_g Sim_{r_i}[g] > \beta$ , but there exist at least two genomes  $g_p$  and  $g_q$  s.t.  $Sim_{r_i}[g_p] \sim Sim_{r_i}[g_q] \sim \max_g Sim_{r_i}[g]$

### Example

Let  $\alpha = 3$  and  $\beta = 0.4$ .

Suppose the  $\alpha$ -similarity between  $r_i$  and  $g_1, g_2, g_3, g_4, g_5$  is

$Sim_{r_i}[g_1] = 0.5, Sim_{r_i}[g_2] = 0, Sim_{r_i}[g_3] = , Sim_{r_i}[g_4] = 0.2, Sim_{r_i}[g_5] = 0.$

## Step 2: Classification

Given a threshold value  $\beta$ , the read  $r_i$  is

- **assigned** to  $g_j$  if  $g_j$  is the only genome such that  $Sim_{r_i}[g_j] \sim \max_g Sim_{r_i}[g]$  and  $Sim_{r_i}[g_j] > \beta$ .
- **not classified** if  $\max_g Sim_{r_i}[g] \leq \beta$ .
- **ambiguous** if  $\max_g Sim_{r_i}[g] > \beta$ , but there exist at least two genomes  $g_p$  and  $g_q$  s.t.  $Sim_{r_i}[g_p] \sim Sim_{r_i}[g_q] \sim \max_g Sim_{r_i}[g]$

### Example

Let  $\alpha = 3$  and  $\beta = 0.4$ .

Suppose the  $\alpha$ -similarity between  $r_i$  and  $g_1, g_2, g_3, g_4, g_5$  is

$Sim_{r_i}[g_1] = 0.5, Sim_{r_i}[g_2] = 0, Sim_{r_i}[g_3] = 0.8, Sim_{r_i}[g_4] = 0.2, Sim_{r_i}[g_5] = 0.$

## Step 2: Classification

Given a threshold value  $\beta$ , the read  $r_i$  is

- **assigned** to  $g_j$  if  $g_j$  is the only genome such that  $Sim_{r_i}[g_j] \sim \max_g Sim_{r_i}[g]$  and  $Sim_{r_i}[g_j] > \beta$ .
- **not classified** if  $\max_g Sim_{r_i}[g] \leq \beta$ .
- **ambiguous** if  $\max_g Sim_{r_i}[g] > \beta$ , but there exist at least two genomes  $g_p$  and  $g_q$  s.t.  $Sim_{r_i}[g_p] \sim Sim_{r_i}[g_q] \sim \max_g Sim_{r_i}[g]$

### Example

Let  $\alpha = 3$  and  $\beta = 0.4$ .

Suppose the  $\alpha$ -similarity between  $r_i$  and  $g_1, g_2, g_3, g_4, g_5$  is

$Sim_{r_i}[g_1] = 0.5$ ,  $Sim_{r_i}[g_2] = 0$ ,  $Sim_{r_i}[g_3] = 0.8$ ,  $Sim_{r_i}[g_4] = 0.2$ ,  $Sim_{r_i}[g_5] = 0$ .

$\Rightarrow r_i$  is **assigned** to  $g_3$ .

## Step 2: Classification

Given a threshold value  $\beta$ , the read  $r_i$  is

- **assigned** to  $g_j$  if  $g_j$  is the only genome such that  $Sim_{r_i}[g_j] \sim \max_g Sim_{r_i}[g]$  and  $Sim_{r_i}[g_j] > \beta$ .
- **not classified** if  $\max_g Sim_{r_i}[g] \leq \beta$ .
- **ambiguous** if  $\max_g Sim_{r_i}[g] > \beta$ , but there exist at least two genomes  $g_p$  and  $g_q$  s.t.  $Sim_{r_i}[g_p] \sim Sim_{r_i}[g_q] \sim \max_g Sim_{r_i}[g]$

### Example

Let  $\alpha = 3$  and  $\beta = 0.4$ .

Suppose the  $\alpha$ -similarity between  $r_i$  and  $g_1, g_2, g_3, g_4, g_5$  is

$Sim_{r_i}[g_1] = 0.3, Sim_{r_i}[g_2] = 0, Sim_{r_i}[g_3] = 0.34, Sim_{r_i}[g_4] = 0.2, Sim_{r_i}[g_5] = 0.$

## Step 2: Classification

Given a threshold value  $\beta$ , the read  $r_i$  is

- **assigned** to  $g_j$  if  $g_j$  is the only genome such that  $Sim_{r_i}[g_j] \sim \max_g Sim_{r_i}[g]$  and  $Sim_{r_i}[g_j] > \beta$ .
- **not classified** if  $\max_g Sim_{r_i}[g] \leq \beta$ .
- **ambiguous** if  $\max_g Sim_{r_i}[g] > \beta$ , but there exist at least two genomes  $g_p$  and  $g_q$  s.t.  $Sim_{r_i}[g_p] \sim Sim_{r_i}[g_q] \sim \max_g Sim_{r_i}[g]$

### Example

Let  $\alpha = 3$  and  $\beta = 0.4$ .

Suppose the  $\alpha$ -similarity between  $r_i$  and  $g_1, g_2, g_3, g_4, g_5$  is

$Sim_{r_i}[g_1] = 0.5$ ,  $Sim_{r_i}[g_2] = 0$ ,  $Sim_{r_i}[g_3] = 0.5$ ,  $Sim_{r_i}[g_4] = 0.2$ ,  $Sim_{r_i}[g_5] = 0$ .

$\Rightarrow r_i$  is **ambiguous**.

# Preliminary experiments on simulated reads

- Positive and negative control datasets (about 20 millions of paired-reads) designed in [Lindgreen et al., 2016].
- Reference database  $\mathcal{G}$ : 930 genomes from 686 species

<i>setA2</i>	CLARK-S -highconf	LIGHTMETAEBWT $\alpha$ 16 $\beta$ 0.25	LIGHTMETAEBWT $\alpha$ 16 $\beta$ 0.35	Centrifuge -min-hitlen 16	Centrifuge -min-hitlen 22
<b>SEN (%)</b>	93.03	<b>93.13</b>	92.69	<b>95.65</b>	93.01
<b>PREC (%)</b>	99.06	99.73	<b>99.74</b>	97.64	99.66
<b>F1 (%)</b>	95.95	<b>96.32</b>	96.08	<b>96.63</b>	96.22
<i>setB2</i>					
<b>SEN (%)</b>	92.84	<b>93.01</b>	92.51	<b>95.53</b>	92.94
<b>PREC (%)</b>	99.11	99.77	<b>99.78</b>	97.68	99.69
<b>F1 (%)</b>	95.87	<b>96.27</b>	96.01	<b>96.59</b>	96.20
<i>setA2Ran</i>					
<i>TN</i>	5,726,336	5,726,294	5,726,357	150,971	5,712,085
<i>FP</i>	22	64	1	5,575,387	14,273
<b>SPEC (%)</b>	99.99	99.99	<b>100.00</b>	<b>2.64</b>	99.75
<i>setB2Ran</i>					
<i>TN</i>	5,406,642	5,406,601	5,406,658	141,994	5,393,260
<i>FP</i>	17	58	1	5,264,665	13,399
<b>SPEC (%)</b>	99.99	99.99	<b>100.00</b>	<b>2.63</b>	99.75

**SEN** = proportion of the actual positives identified by the method.

**PREC** = proportion of positives that are correctly identified by the method.

**F1** = harmonic mean between **SEN** and **PREC**.

**SPEC** = proportion of actual negatives that are correctly identified as such.

# Preliminary experiments on simulated reads

- Positive and negative control datasets (about 20 millions of paired-reads) designed in [Lindgreen et al., 2016].
- Reference database  $\mathcal{G}$ : 930 genomes from 686 species

<i>setA2</i>	CLARK-S -highconf	LIGHTMETAEBWT $\alpha$ 16 $\beta$ 0.25	LIGHTMETAEBWT $\alpha$ 16 $\beta$ 0.35	Centrifuge -min-hitlen 16	Centrifuge -min-hitlen 22
<b>SEN (%)</b>	93.03	<b>93.13</b>	92.69	<b>95.65</b>	93.01
<b>PREC (%)</b>	99.06	99.73	<b>99.74</b>	97.64	99.66
<b>F1 (%)</b>	95.95	<b>96.32</b>	96.08	<b>96.63</b>	96.22
<i>setB2</i>					
<b>SEN (%)</b>	92.84	<b>93.01</b>	92.51	<b>95.53</b>	92.94
<b>PREC (%)</b>	99.11	99.77	<b>99.78</b>	97.68	99.69
<b>F1 (%)</b>	95.87	<b>96.27</b>	96.01	<b>96.59</b>	96.20
<i>setA2Ran</i>					
<i>TN</i>	5,726,336	5,726,294	5,726,357	150,971	5,712,085
<i>FP</i>	22	64	1	5,575,387	14,273
<b>SPEC (%)</b>	99.99	99.99	<b>100.00</b>	<b>2.64</b>	99.75
<i>setB2Ran</i>					
<i>TN</i>	5,406,642	5,406,601	5,406,658	141,994	5,393,260
<i>FP</i>	17	58	1	5,264,665	13,399
<b>SPEC (%)</b>	99.99	99.99	<b>100.00</b>	<b>2.63</b>	99.75

**SEN** = proportion of the actual positives identified by the method.

**PREC** = proportion of positives that are correctly identified by the method.

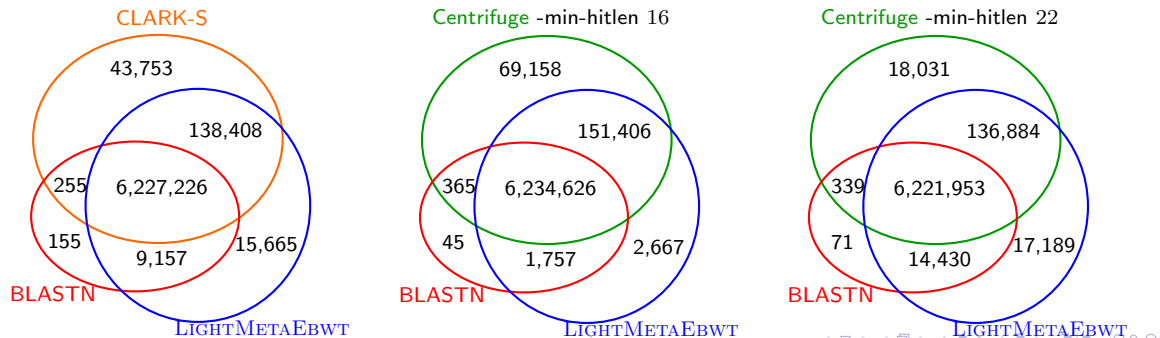
**F1** = harmonic mean between **SEN** and **PREC**.

**SPEC** = proportion of actual negatives that are correctly identified as such.

# Preliminary experiments on real data: Mock community

Reference database  $\mathcal{G}$ : 61 genomes from 22 species

SRR172902	BLASTN -eval $10^{-5}$ -id 90	CLARK-S -c 0.75 -g 0.08	LIGHTMETAEBWT $\alpha$ 15 $\beta$ 0.2	Centrifuge -min-hitlen 16	Centrifuge -min-hitlen 22
<b>Classified</b>	6,236,793	6,409,642	6,390,456	6,455,555	6,377,207
<b>Ambiguous</b>	48,310	0	10,570	11,150	8,422
<b>Unclassified</b>	276,962	152,423	161,039	95,360	176,436
<b>SEN (%)</b>	95.04	97.68	97.39	98.38	97.18

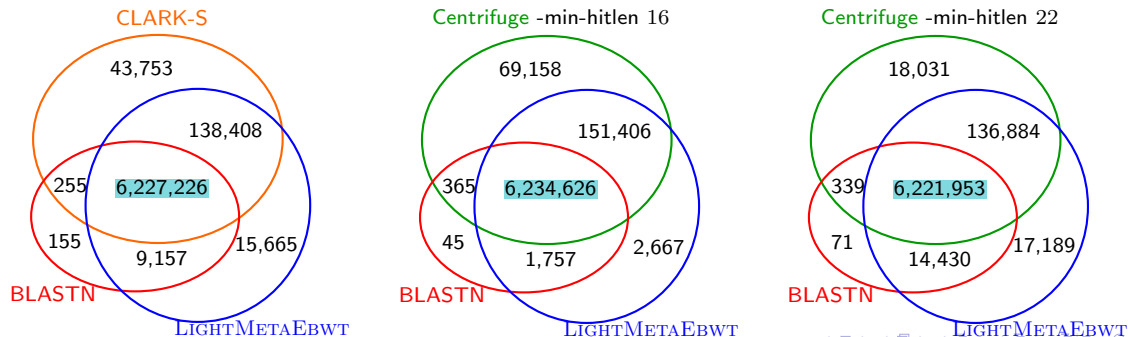




# Preliminary experiments on real data: Mock community

Reference database  $\mathcal{G}$ : 61 genomes from 22 species

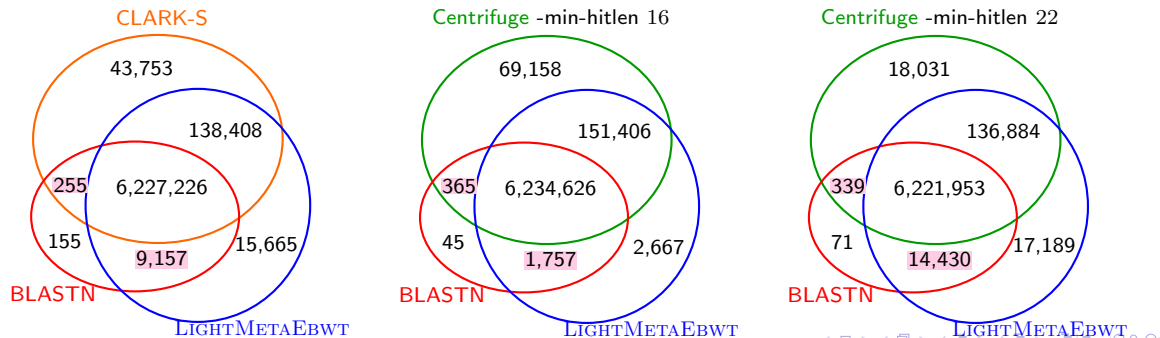
SRR172902	BLASTN -eval $10^{-5}$ -id 90	CLARK-S -c 0.75 -g 0.08	LIGHTMETAEBWT $\alpha$ 15 $\beta$ 0.2	Centrifuge -min-hitlen 16	Centrifuge -min-hitlen 22
<b>Classified</b>	6,236,793	6,409,642	6,390,456	6,455,555	6,377,207
<b>Ambiguous</b>	48,310	0	10,570	11,150	8,422
<b>Unclassified</b>	276,962	152,423	161,039	95,360	176,436
<b>SEN (%)</b>	95.04	97.68	97.39	98.38	97.18



# Preliminary experiments on real data: Mock community

Reference database  $\mathcal{G}$ : 61 genomes from 22 species

SRR172902	BLASTN -eval $10^{-5}$ -id 90	CLARK-S -c 0.75 -g 0.08	LIGHTMETAEBWT $\alpha$ 15 $\beta$ 0.2	Centrifuge -min-hitlen 16	Centrifuge -min-hitlen 22
<b>Classified</b>	6,236,793	6,409,642	6,390,456	6,455,555	6,377,207
<b>Ambiguous</b>	48,310	0	10,570	11,150	8,422
<b>Unclassified</b>	276,962	152,423	161,039	95,360	176,436
<b>SEN (%)</b>	95.04	97.68	97.39	98.38	97.18



## Final Remarks

Notice that a like-for-like comparison on the time-consuming and on memory-consuming between LIGHTMETAEBWT, CLARK-S and Cenfrifuge is not possible

LIGHTMETAEBWT	Centrifuge	Clark-S
One-core	multi-thread	multi-thread
Analyze all paired-reads at the same time, this implies keeping in memory the whole similarity matrix	Analyze one read at a time and Pair-reads analyzed in parallel	Analyze one read at a time and Concatenate paired-reads
no engineered implementation	engineered implementation	engineered implementation

	Classification for setA2 (21,461,160 paired-reads and 930 genomes)			
	LIGHTMETAEBWT	Centrifuge	CLARK-S	
setA2	$\alpha$ 16	-min-hitlen 16	-highconf	last step
RAM	$\sim$ 1GB (+ $\sim$ 18GB for the matrix)	$\sim$ 2GB	$\sim$ 121GB	$\sim$ 78GB
Time	$\sim$ 54 min	$\sim$ 27min	$\sim$ 8h	$\sim$ 28min

Here, we improved our previous version of this method presented to AICoB 2019: for instance, considering the IUPAC list. We are now working on a more engineered implementation

Thank you!

# Preliminary experiments

Dataset	$ \mathcal{R} $	$ \mathcal{G} $	Instance	Efficiency	Wall Clock	RAM
<i>setA2</i>	21,461,160	930	CLARK-S (8cores)	99%	8:25:17	121.58
			LIGHTMETAEBWT	99%	2:29:57	19.45
			Centrifuge	99%	1:55:10+26:38	9.32
<i>setB2</i>	20,249,373	930	CLARK-S (8cores)	69%	14:26:31	121.58
			LIGHTMETAEBWT	96%	3:06:05	18.55
			Centrifuge	99%	1:55:10+21:16	9.32
<i>setA2Ran</i>	5,726,358	930	CLARK-S (8cores)	92%	5:02:22	121.58
			LIGHTMETAEBWT	98%	18:46	9.94
			Centrifuge	99%	1:55:10+4:08	9.32
<i>setB2Ran</i>	5,406,659	930	CLARK-S (8cores)	99%	5:16:50	121.58
			LIGHTMETAEBWT	99%	17:01	9.61
			Centrifuge	99%	1:55:10+3:20	9.32
<i>SRR172902</i>	6,562,065	61	CLARK-S (8cores)	99%	33:40	72.23
			LIGHTMETAEBWT	99%	5:57	0.69
			Centrifuge	99%	1:21+2:29	0.37

- The three data structures for LIGHTMETAEBWT setA2 take 1:23:55 and 4GB of RAM to be constructed.
- For LIGHTMETAEBWT, post-processing time and space usage differ from step to step. For instance, per each set of paired end read in setA2 step 1 takes 1:15 with 1GB of RAM (changeable), step 2 around 36:17 with 19.45GB (18.59GB to keep the matrix of similarities) and step 3 takes 15:55 with 0.06GB.