



Compressed indexation structure for analysing collections of similar genomes

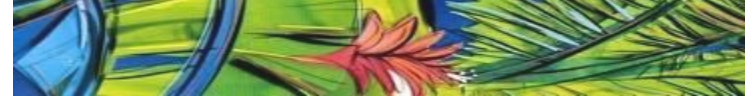
Annie CHATEAU & Alban MANCHERON (LIRMM, équipe MAB),
Gautier SARAH, Gaetan DROC & Manuel RUIZ (CIRAD, UMR AGAP, équipe ID/CIAT)

Clément AGRET

Équipes : ID : Intégration Données

MAB : Méthodes et algorithmes pour la bioinformatique





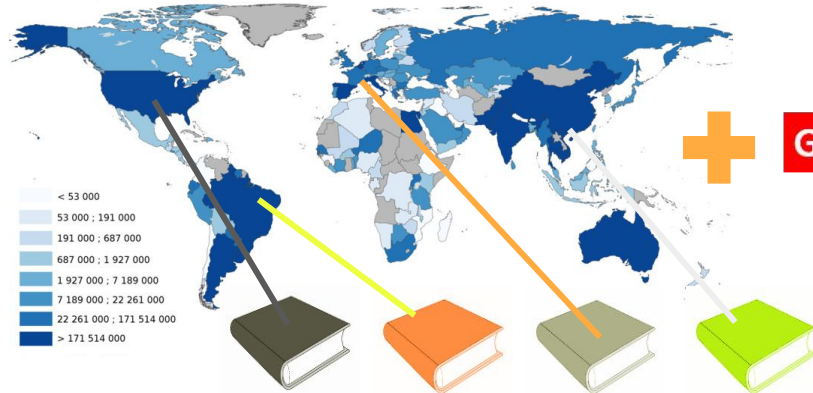
Data & Computation

Let's simplify, what's happen if we see the **GENOME** as a simple book ?

If



=



+

GenomeHarvest

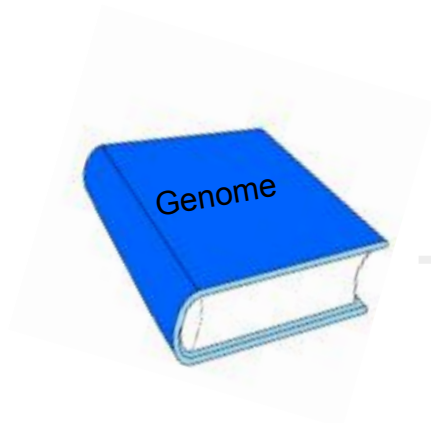
=



A lot of books to read !!



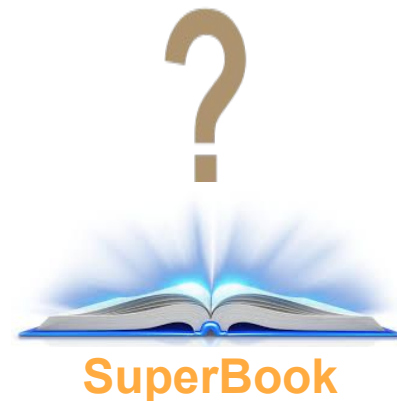
Data & Computation



1
BOOK



+1000
BOOKS



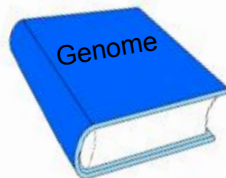
If all books tell more or less the same story, what about writing a **SuperBook** ?



Data & Computation

When you read a book you should be able to answer to questions like:

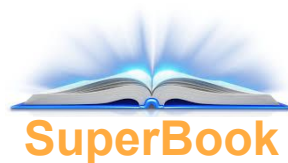
- How many letters are there in the third chapter of the book volume 3 of "Around the World in 80 Days"?
- Is the word "Lustful" appear in the first chapter?
- What is the sentence beginning at 275th paragraph of Chapter 4?



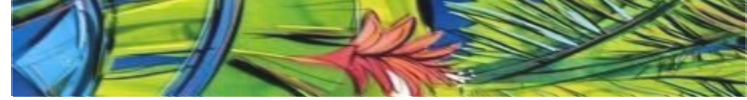
Read the book and answer questions.



Read all books and answer questions !
Can take more than a life.



Read the superBook and answer questions !!



Vocabulary

Alphabet, prefix, factor, suffix

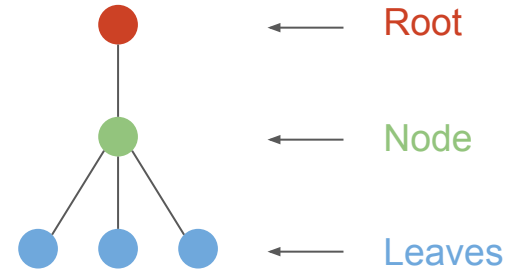
Alphabet $\rightarrow \Sigma = \{A, C, G, T\}$

Prefix Factor Suffix

↓ ↓ ↓

CAGCTGACTAGCACGAACT

Root, node, leaves



K-mer

A fragment of k consecutive nucleotides of a word (a sequence from a reference genome as appropriate)

→ A k-mer is a k size factor of a word



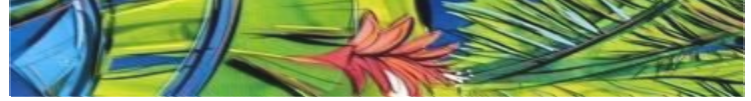
Our hypothesis

In rice genome, the number of distinct k-mers (which appear at least once) tends to stabilize from a x number of genomes.

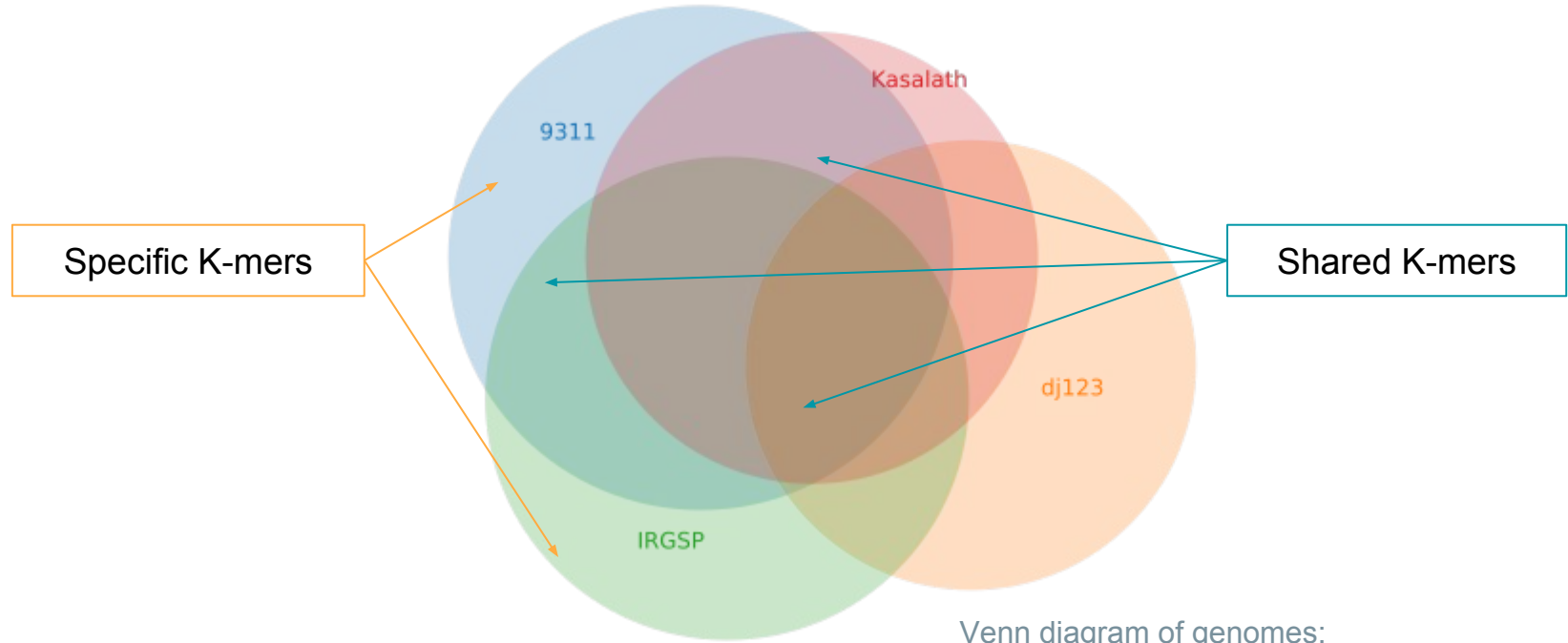
Adding a new genome to an existing index created on 1000 genomes will be equal to adding a reduct set of positions.

To validate this hypothesis:

- 4 genomes + meta-chart ⁽¹⁾ → Venn chart
- 8 genomes + k-mers counter (JellyFish) ⁽²⁾

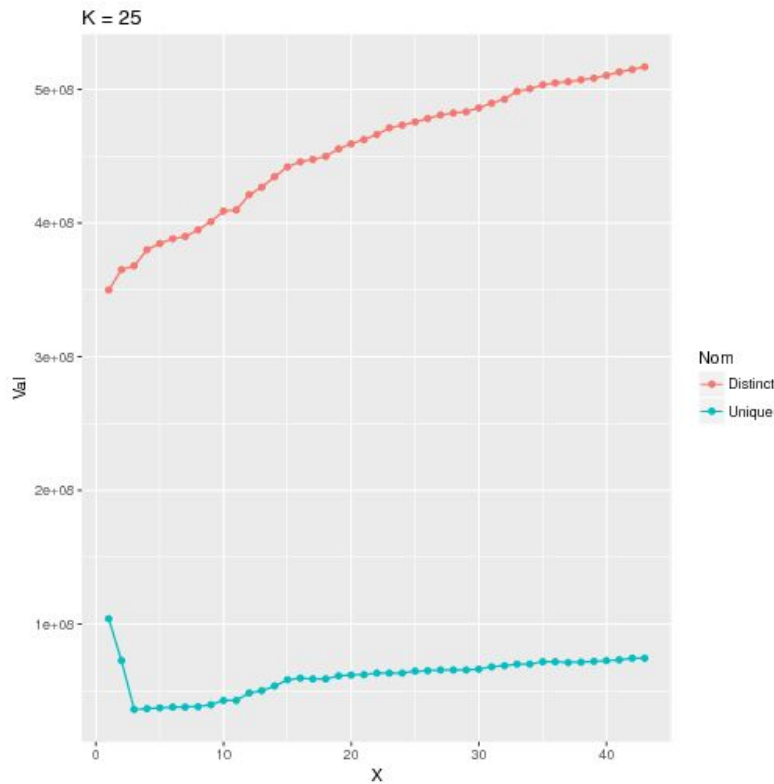


Validation of hypothesis

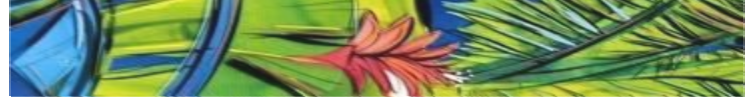




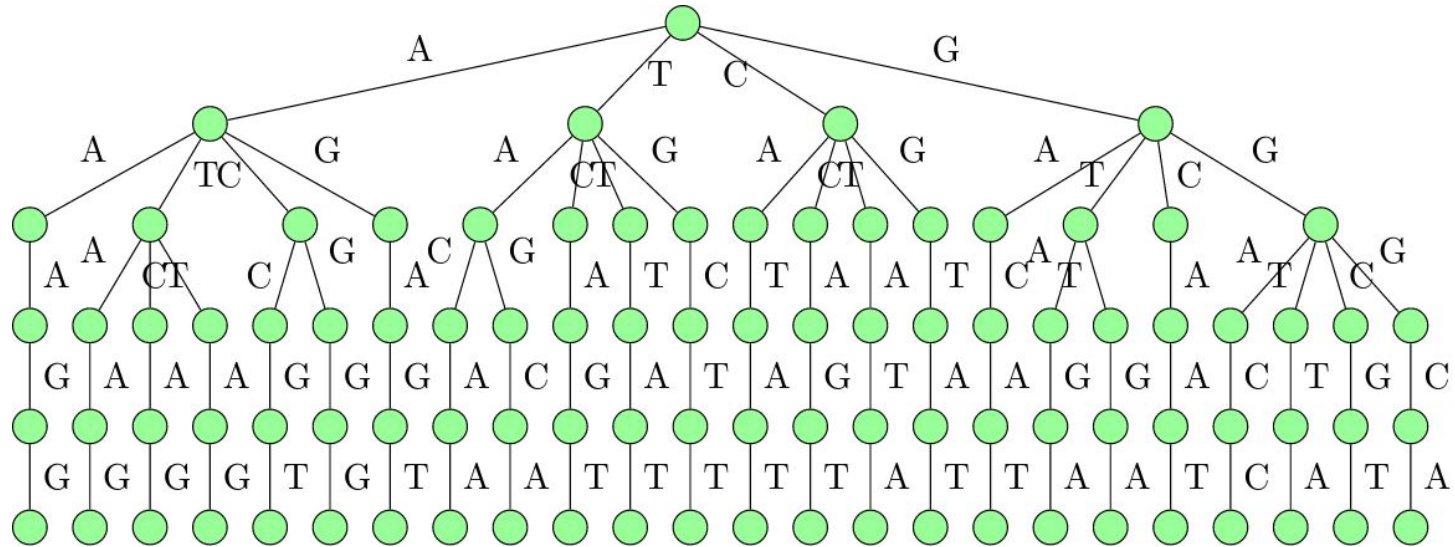
Our study



Distinct: Appears at least once
Unique: Appears exactly once

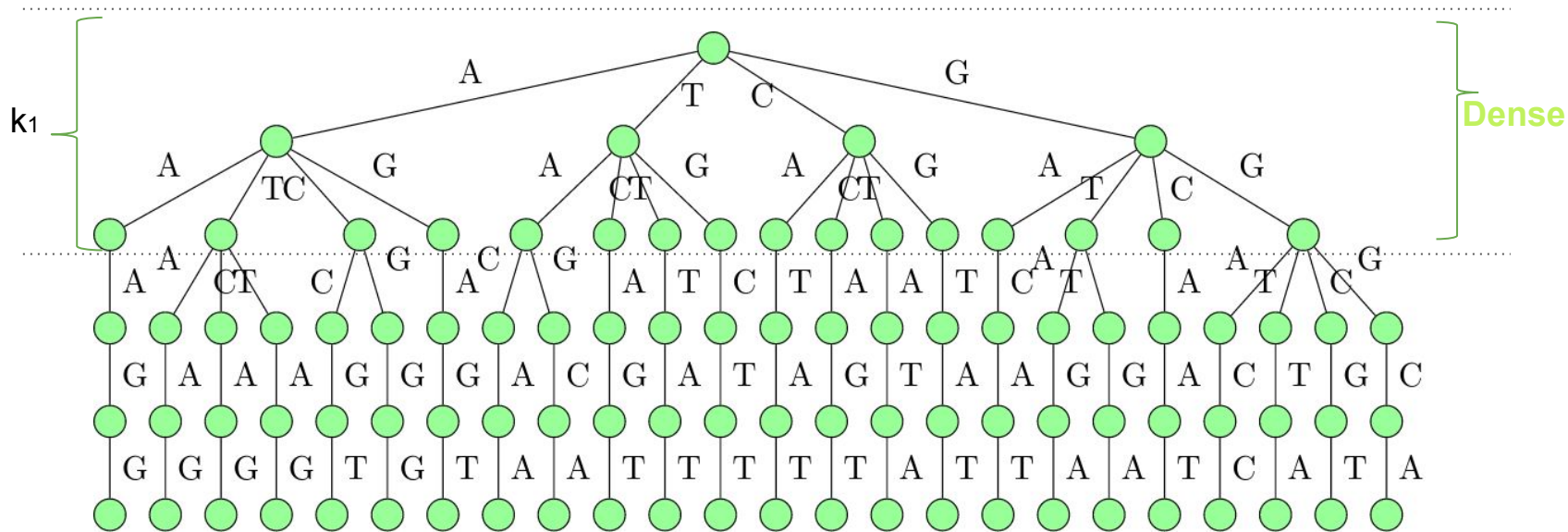


Our approach



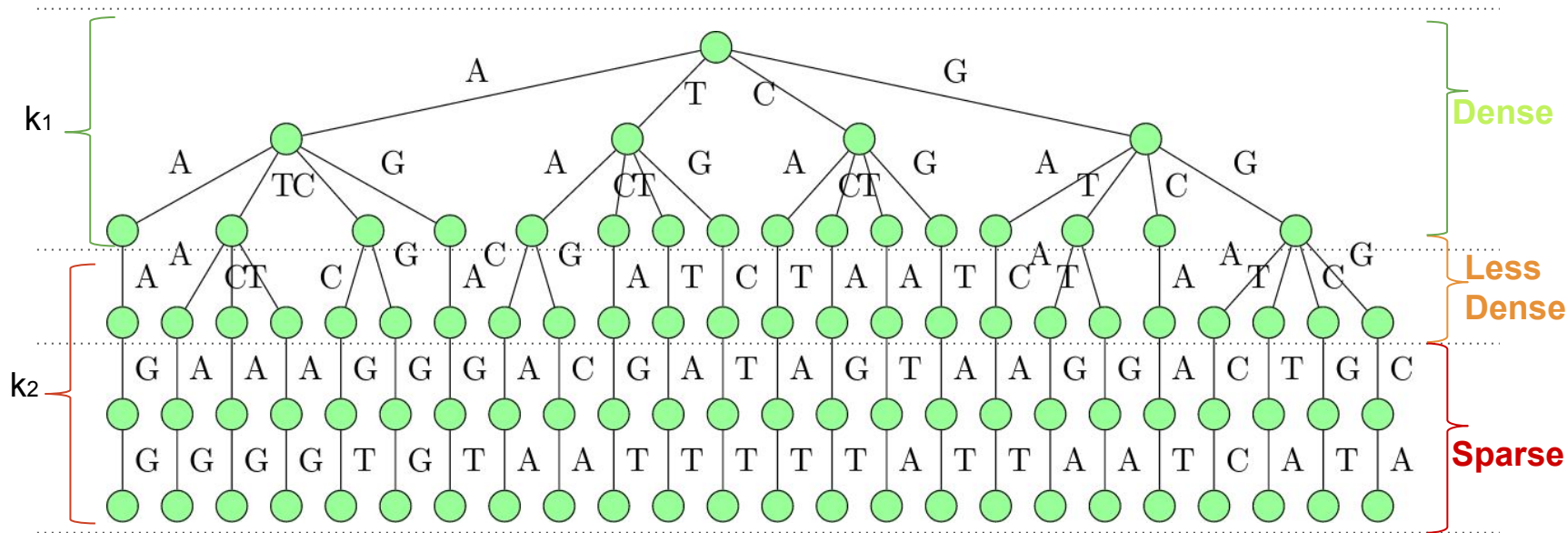


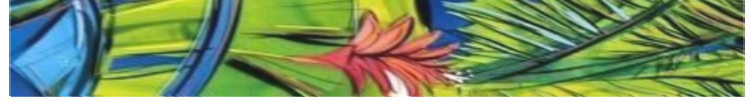
Our approach



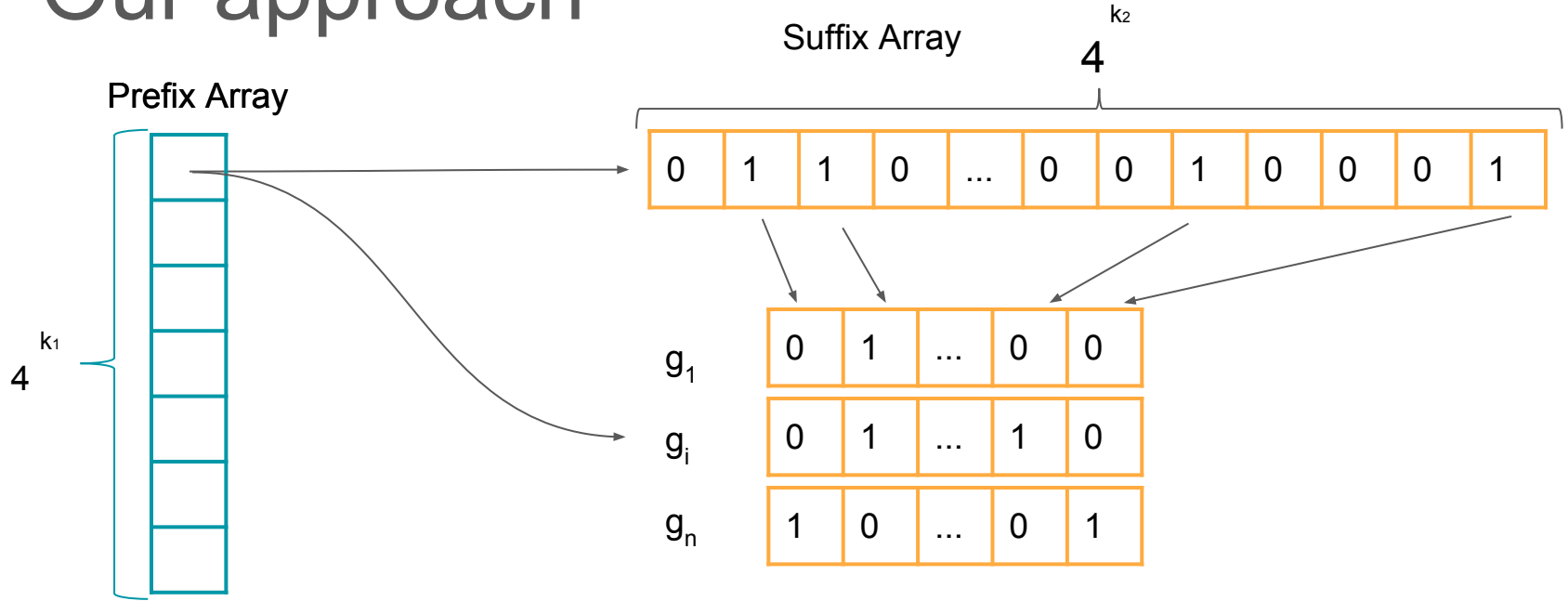


Our approach





Our approach



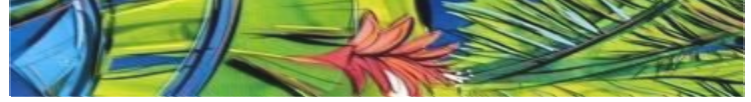


Examples

Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

g1 = ABBABAABAB
g2 = ABBBAABABB
g3 = AABBBABABA

Let's Create the **superBook** !

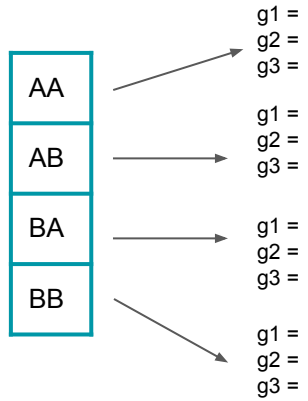
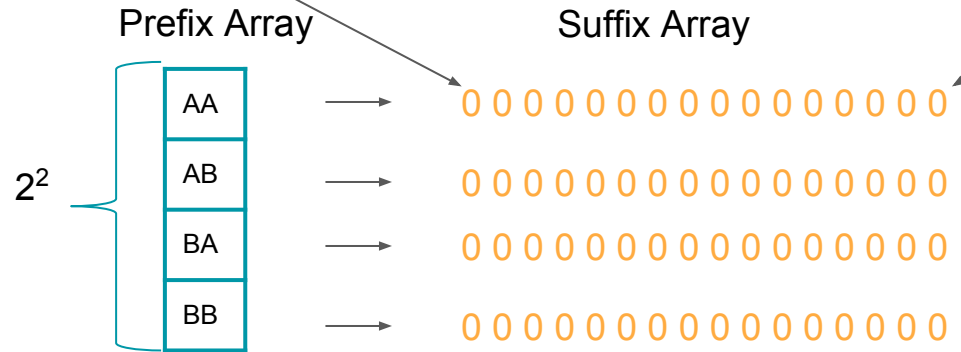


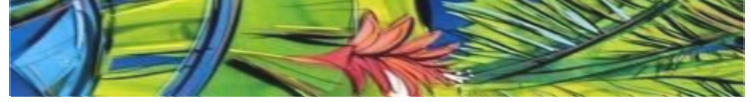
Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

$g_1 = \text{ABBABAABAB}$
 $g_2 = \text{ABBBAABABB}$
 $g_3 = \text{AABBBABABA}$





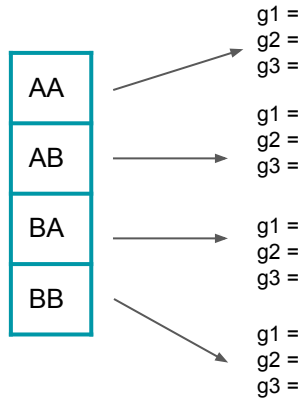
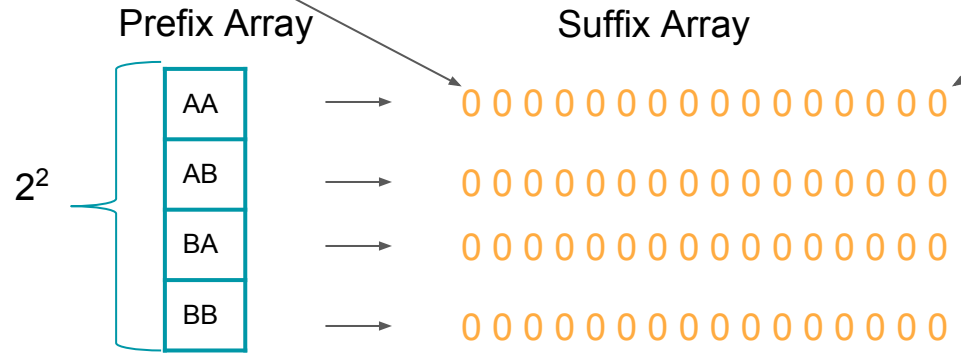
Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

$g_1 = \text{ABBABAABAB}$
 $g_2 = \text{ABBBAABABB}$
 $g_3 = \text{AABBBABABA}$

$g_1 = \text{ABBABAABAB}$





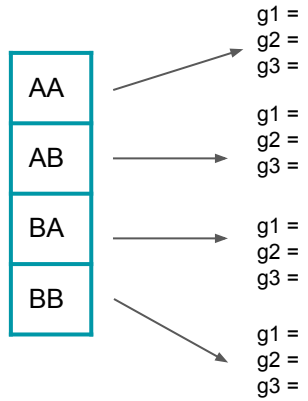
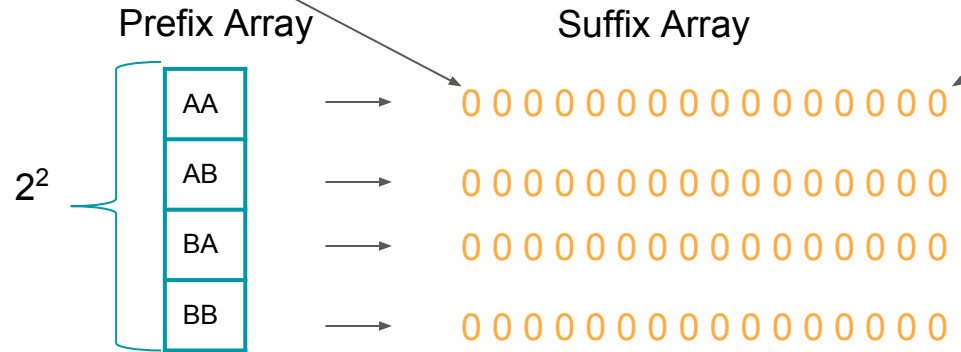
Examples

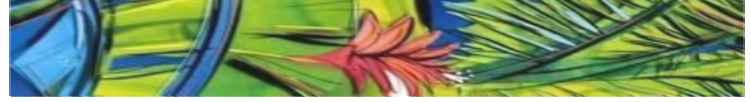
AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

$g_1 = \text{ABBABAABAB}$
 $g_2 = \text{ABBBAABABB}$
 $g_3 = \text{AABBBABABA}$

$g_1 = \text{ABBABAABAB}$
 ABBABA



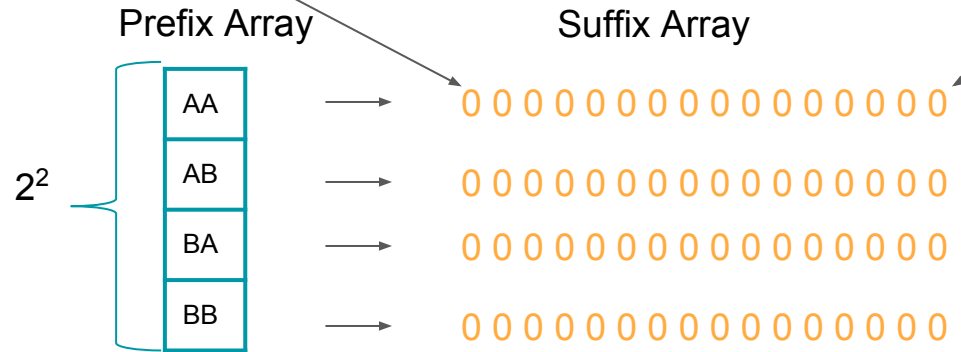


Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

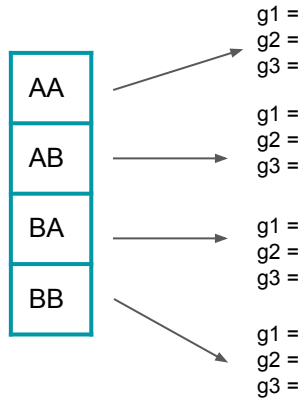
Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

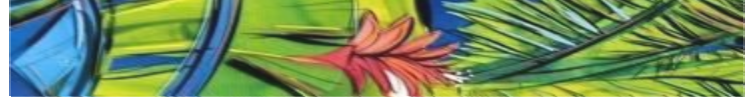
$g_1 = \text{ABBABAABAB}$
 $g_2 = \text{ABBBAABABB}$
 $g_3 = \text{AABBBABABA}$



$g_1 = \text{ABBABAABAB}$
 ABBABA

ABBABA



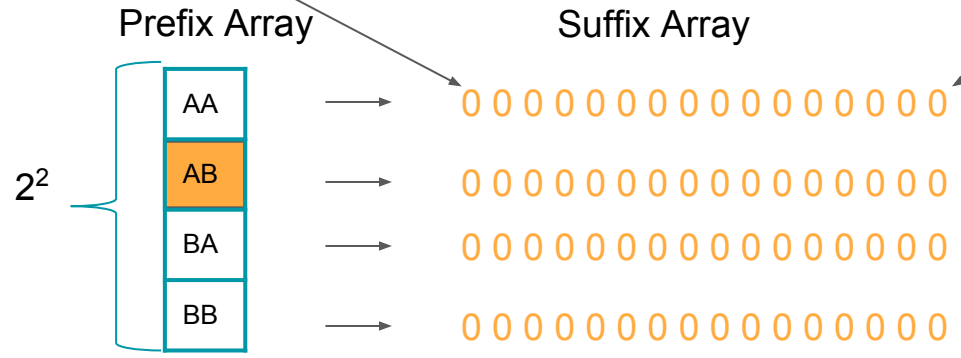


Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

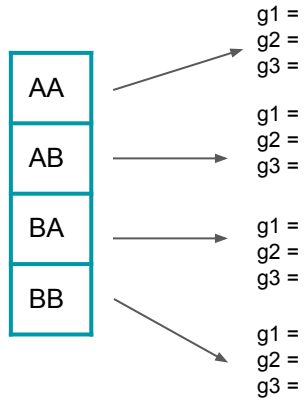
Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

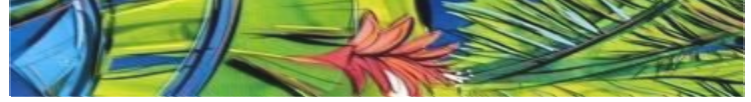
$g_1 = \text{ABBABAABAB}$
 $g_2 = \text{ABBBAABABB}$
 $g_3 = \text{AABBBABABA}$



$g_1 = \text{ABBABAABAB}$
 ABBABA

ABBABA



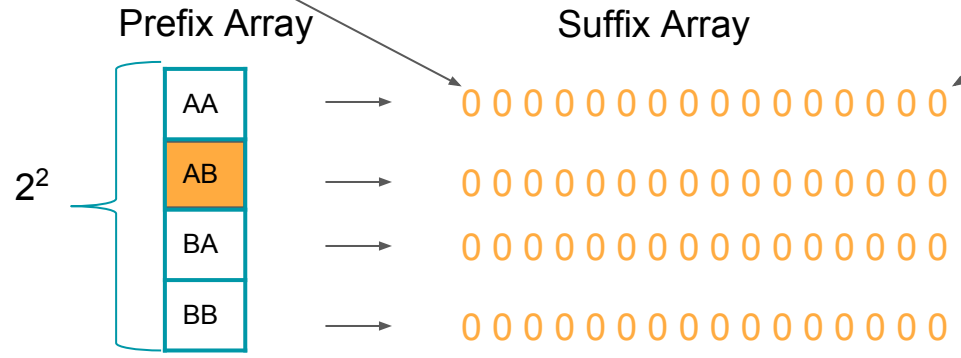


Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

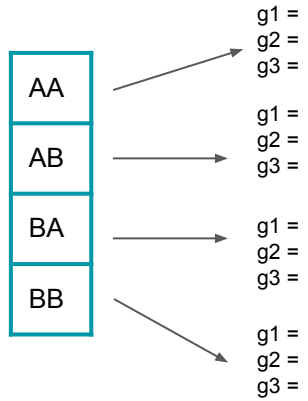
Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

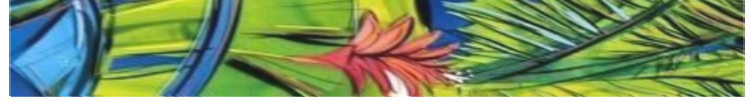
g1 = ABBABAABAB
 g2 = ABBBAABABB
 g3 = AABBBABABA



g1 = ABBABAABAB
 ABBABA

ABBABA



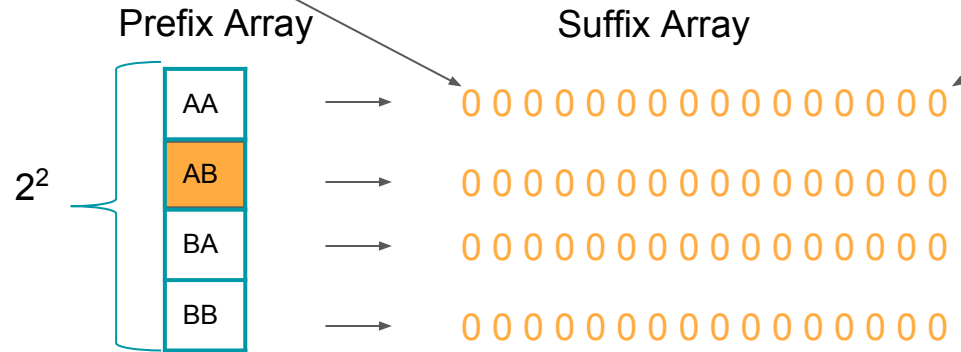


Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, **BABA**, BABB, BBAA, BBAB, BBBA, BBBB

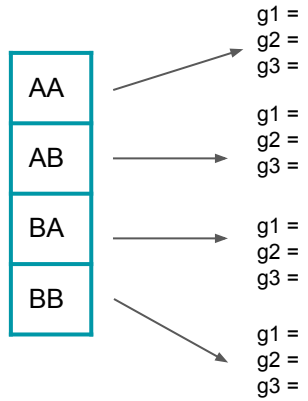
Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

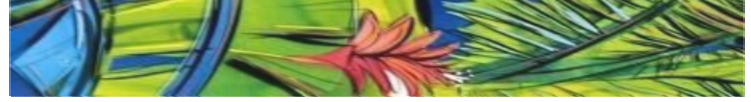
$g_1 = \text{ABBABAABAB}$
 $g_2 = \text{ABBBAABABB}$
 $g_3 = \text{AABBBABABA}$



$g_1 = \text{ABBABAABAB}$
 ABBABA

ABBABA





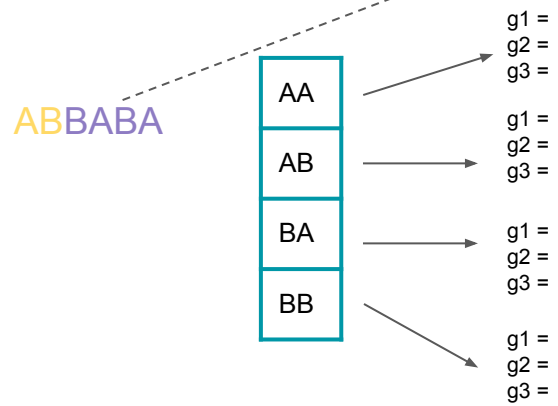
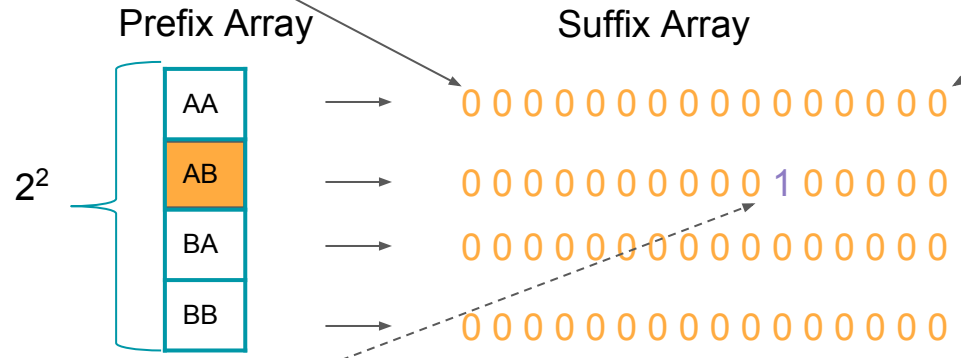
Examples

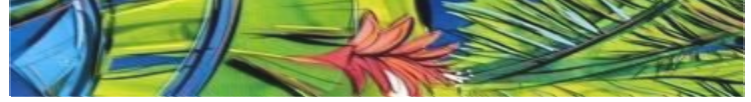
AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, **BABA**, BABB, BBAA, BBAB, BBBA, BBBB

Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

g1 = ABBABAABAB
 g2 = ABBBAABABB
 g3 = AABBBABABA

g1 = ABBABAABAB
 ABBABA





Examples

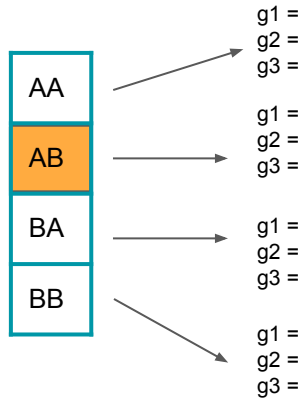
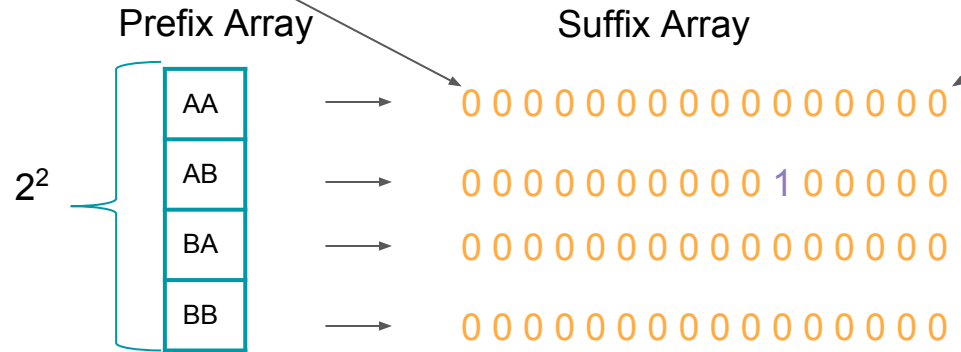
AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

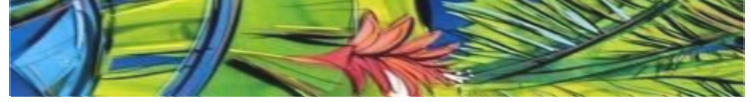
Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

g1 = ABBABAABAB
 g2 = ABBBAABABB
 g3 = AABBBABABA

g1 = ABBABAABAB
 ABBABA

ABBABA



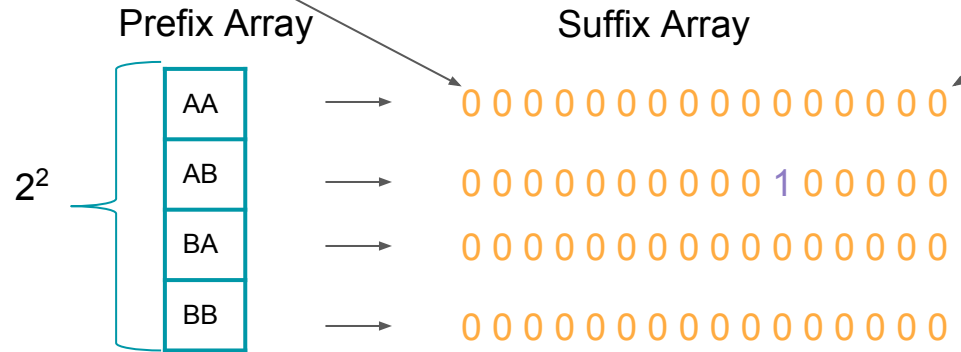


Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

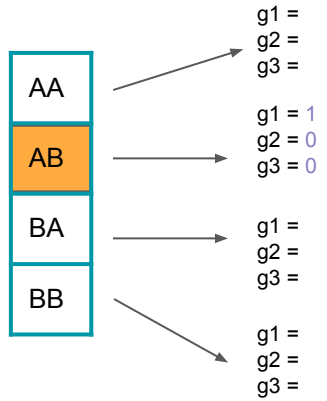
Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

g1 = ABBABAABAB
 g2 = ABBBAABABB
 g3 = AABBBABABA



g1 = ABBABAABAB
 ABBABA

ABBABA





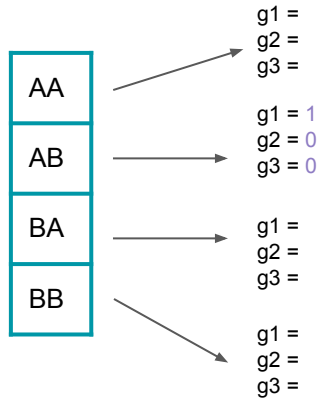
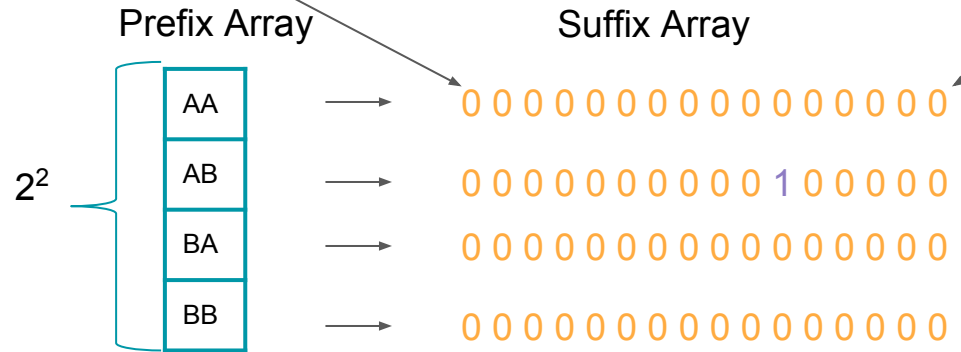
Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

$g_1 = \text{ABBABAABAB}$
 $g_2 = \text{ABBBAABABB}$
 $g_3 = \text{AABBBABABA}$

$g_1 = \text{ABBABAABAB}$
 ABBABA



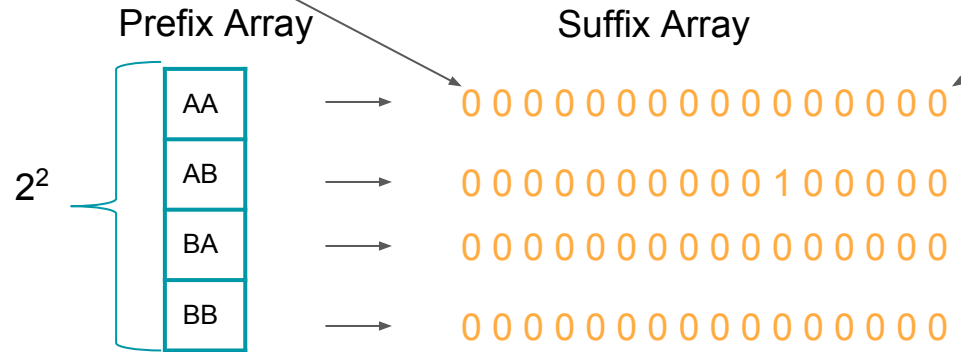


Examples

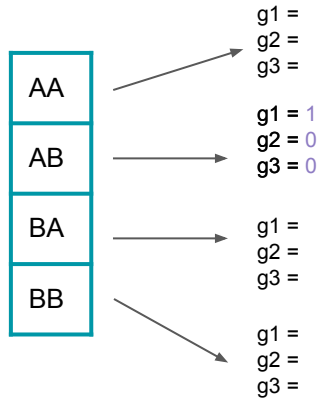
AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

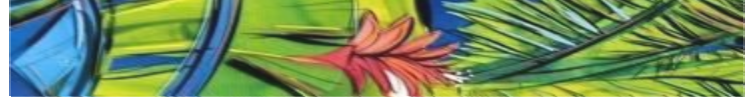
Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

g1 = ABBABAABAB
 g2 = ABBBAABABB
 g3 = AABBBABABA



g1 = ABBABAABAB
 ABBABA
 BBABAA
 BABAAB
 ABAABA
 BAABAB



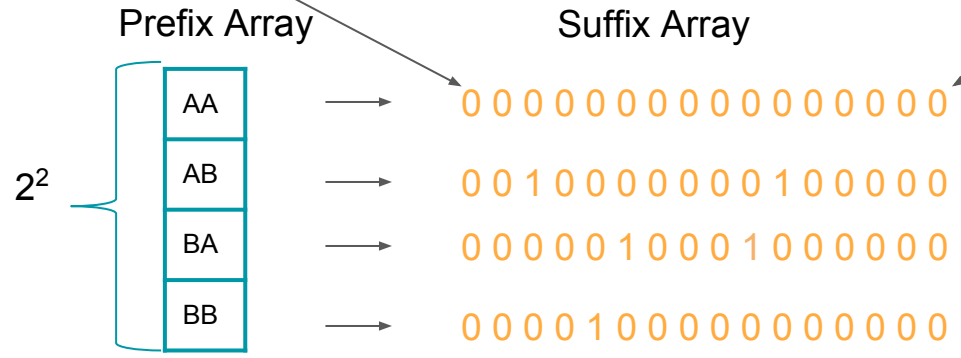


Examples

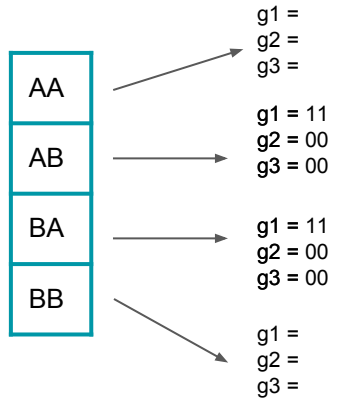
AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBA, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

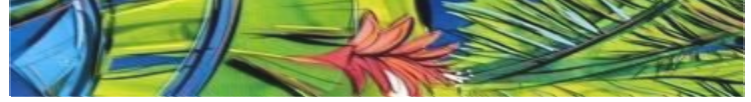
Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

g1 = ABBABAABAB
 g2 = ABBBAABABB
 g3 = AABBBABABA



g1 = ABBABAABAB
 ABBABA
 BBABAA
 BABAAB
 ABAABA
 BAABAB



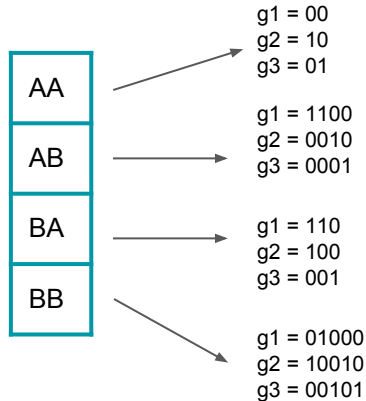
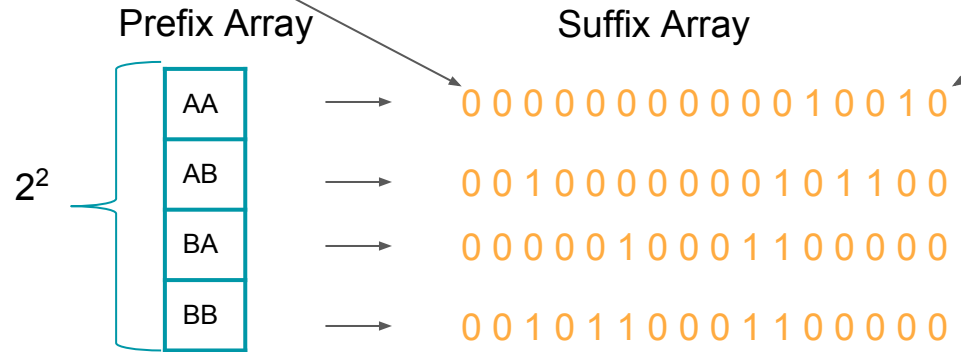


Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

$g_1 = \text{ABBABAABAB}$
 $g_2 = \text{ABBBAABABB}$
 $g_3 = \text{AABBBABABA}$



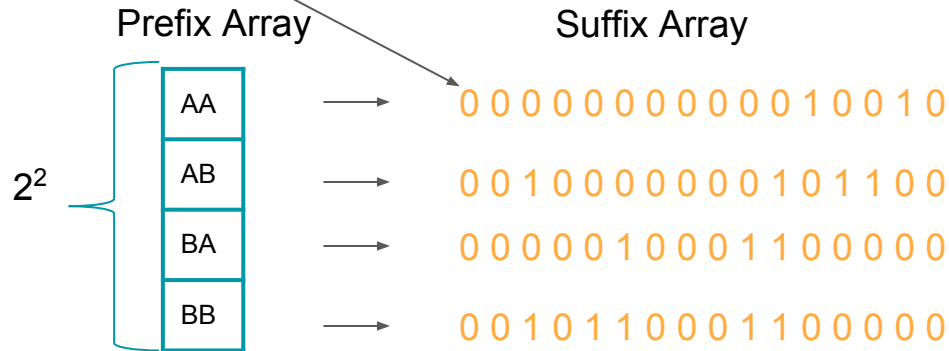


Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

Alphabet $\rightarrow \Sigma = \{A, B\}$
 $K = 6 \rightarrow K_1 = 2 \text{ \& } K_2 = 4$

$g_1 = \text{ABBABAABAB}$
 $g_2 = \text{ABBBAABABB}$
 $g_3 = \text{AABBBABABA}$

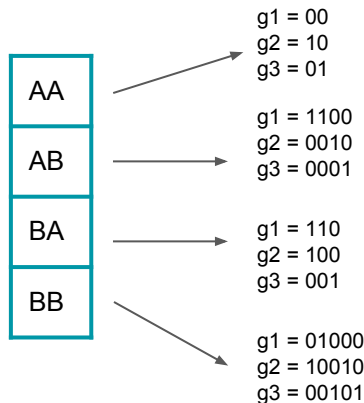


$\text{Rank}_x(i)$:

Rank return the number of elements x in the range $[0, i]$.

$\text{Select}_x(i)$:

Select is the inverse operation to rank; it answers the question “at which position is the i^{th} set bit?”



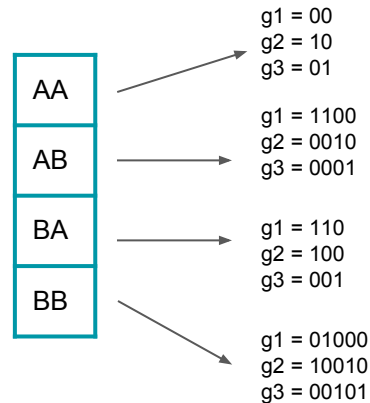
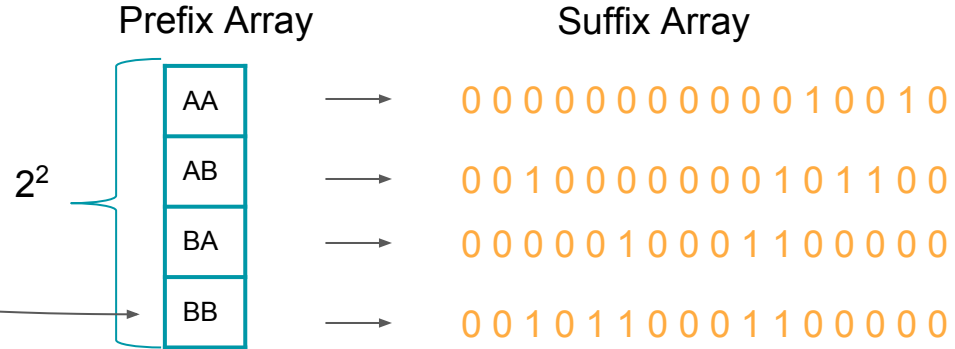


Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

Is the word BBABAA exist ?

BBABAA



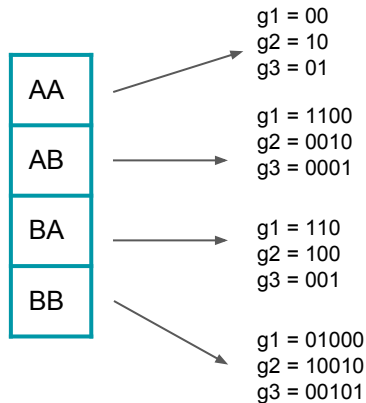
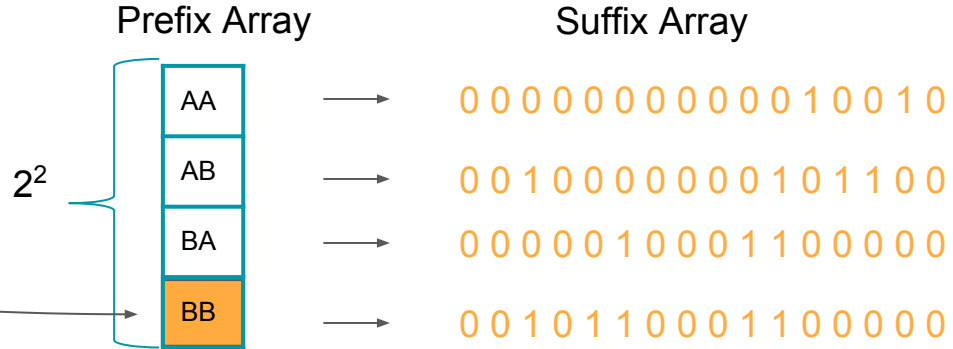


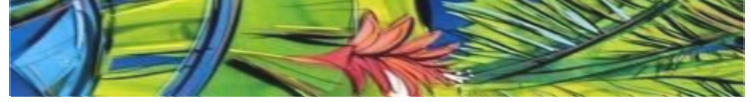
Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

Is the word BBABAA exist ?

BBABAA





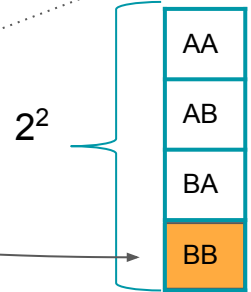
Examples

AAAA, AAAB, AABA, AABB, **ABAA**, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

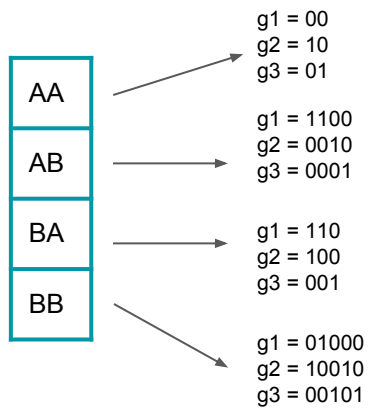
Is the word BBABAA exist ?

BBABAA

Prefix Array



Suffix Array





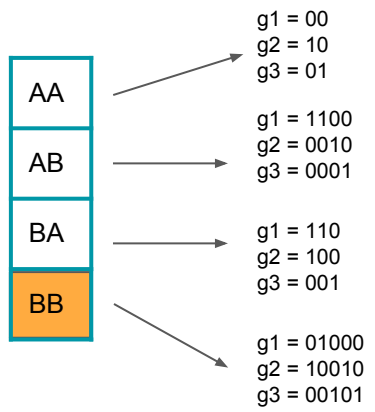
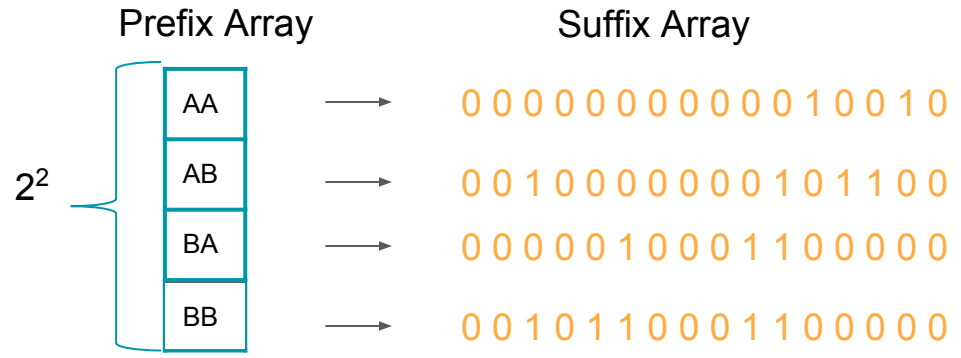
Examples

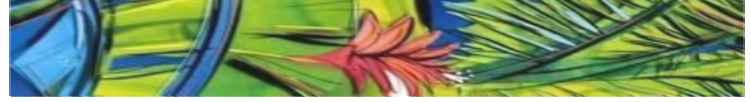
AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

Is the word BBABAA exist ?

BBABAA

In which genomes ?
 $\text{Rank}_1(\text{ABAA}) = 2$





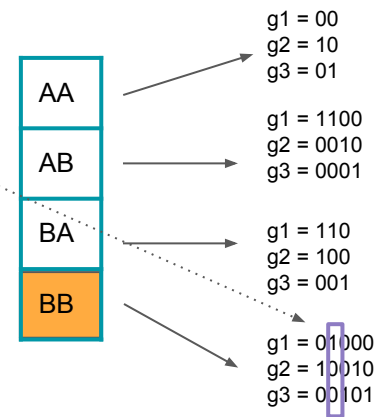
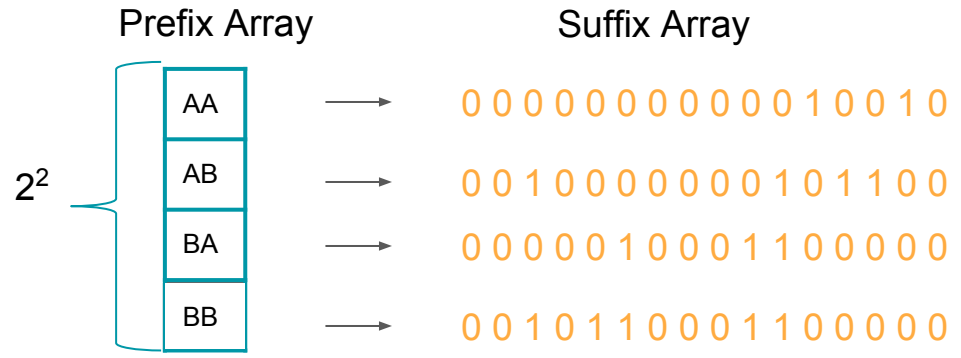
Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBA, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

Is the word BBABAA exist ?

BBABAA

In which genomes ?
 $\text{Rank}_1(\text{ABAA}) = 2$





Examples

AAAA, AAAB, AABA, AABB, ABAA, ABAB, ABBA, ABBB, BAAA, BAAB, BABA, BABB, BBAA, BBAB, BBBA, BBBB

Is the word BBABAA exist ?

BBABAA

In which genomes ?

Rank₁(ABAA) = 2

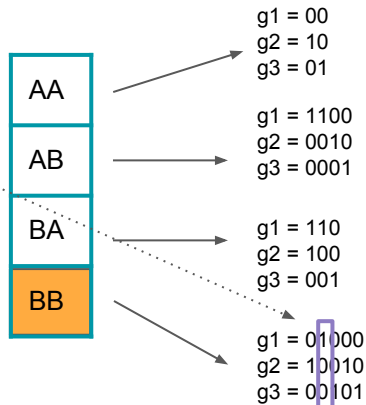
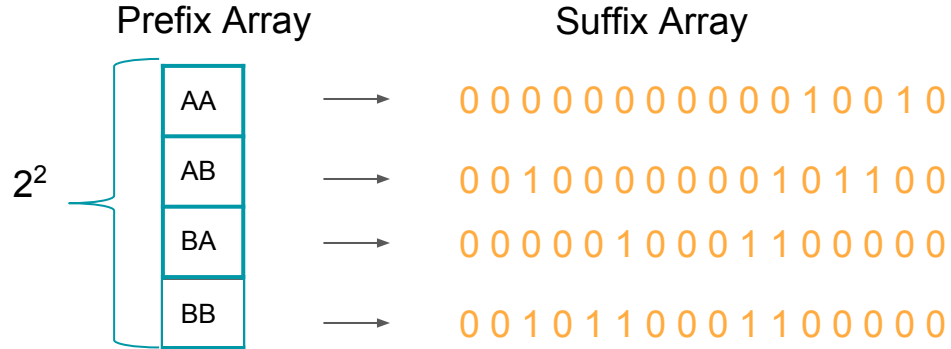
→ g1 only!

Reminder :

g1 = ABBABAABAB

g2 = ABBBAABABB

g3 = AABBBABABA



RRR

RRR was first proposed by Raman et al [1]

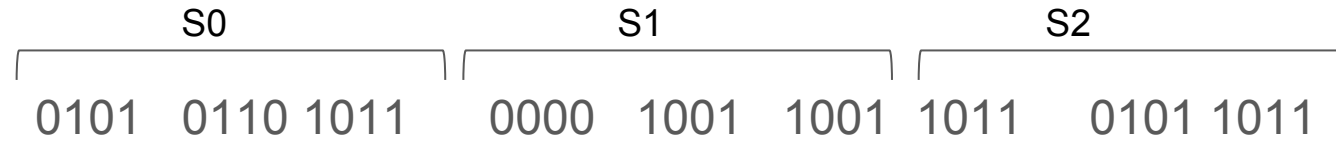
→ $O(1)$ time binary rank queries

→ $N H_0(S) + o(N)$ $H_0(S)$ is the zeroth-order empirical entropy of S

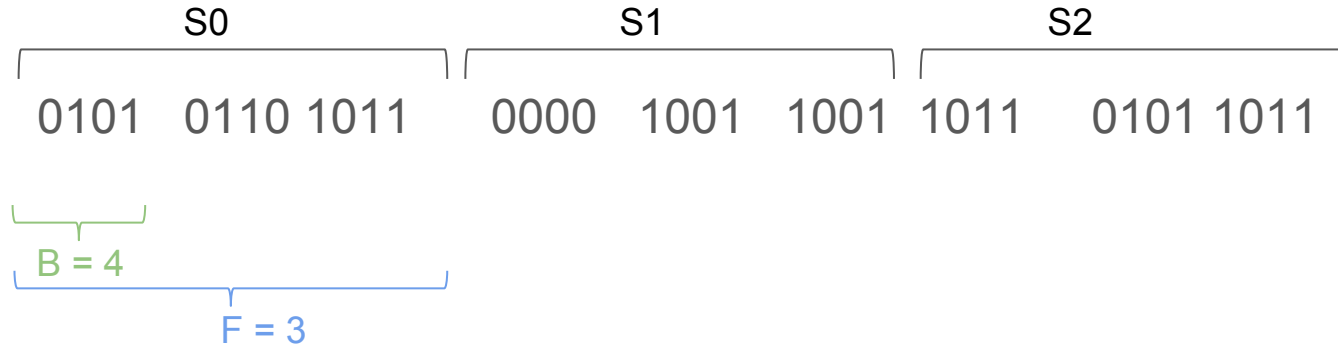
RRR

0101 0110 1011 0000 1001 1001 1011 0101 1011

RRR



RRR



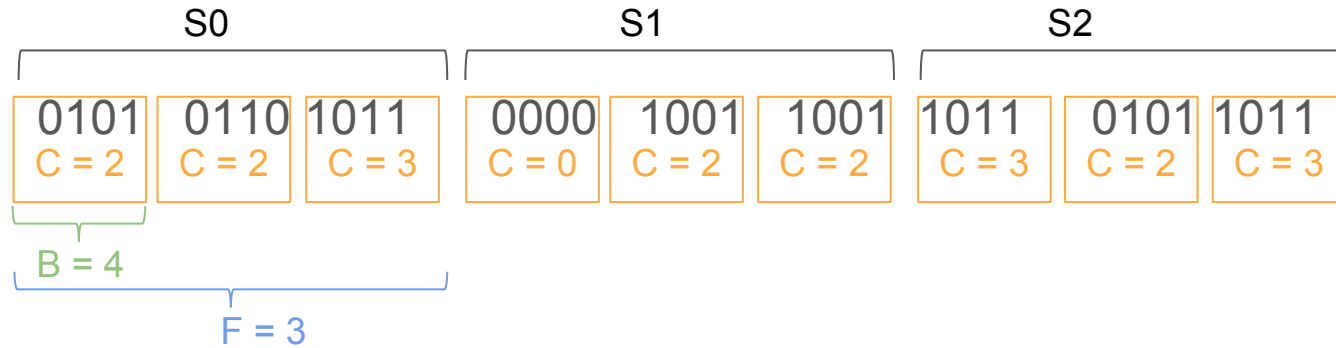
B : Size of the blocks = Numbers of 1 and 0 in the block

F : Superblock factor

C : Class number = Numbers of 1 in the block b

O : Offset = Index into the table

RRR



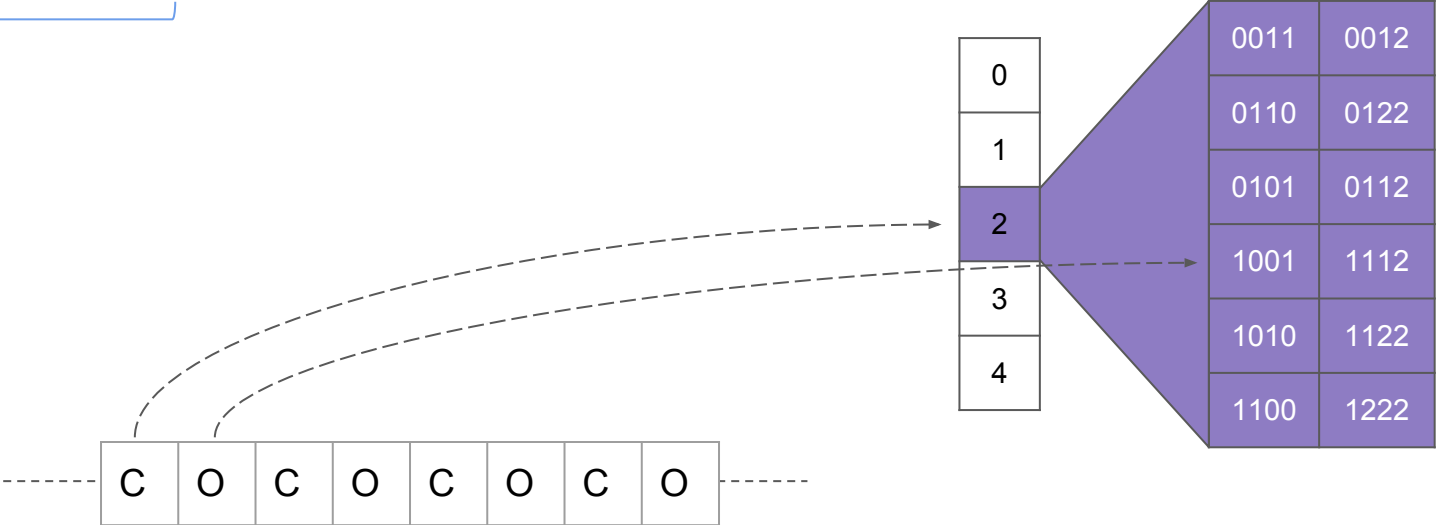
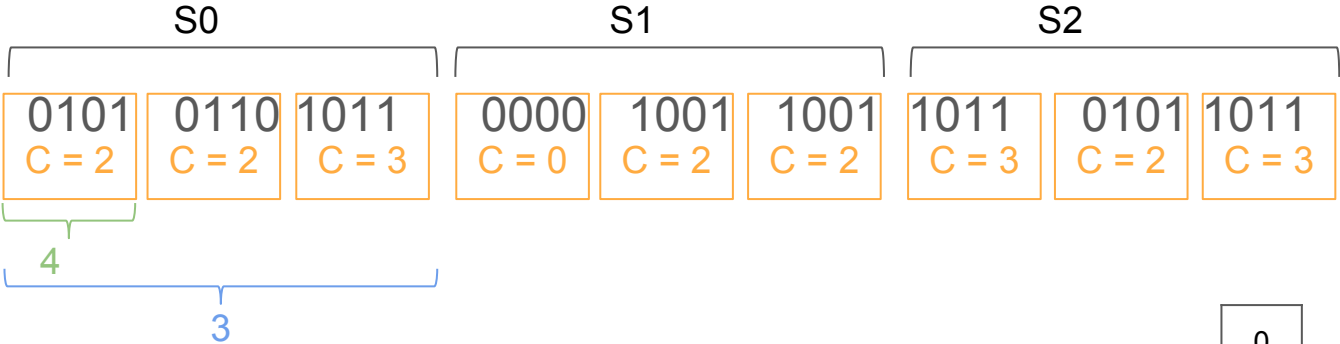
B : Size of the blocks = Numbers of 1 and 0 in the block

F : Superblock factor

C : Class number = Numbers of 1 in the block b

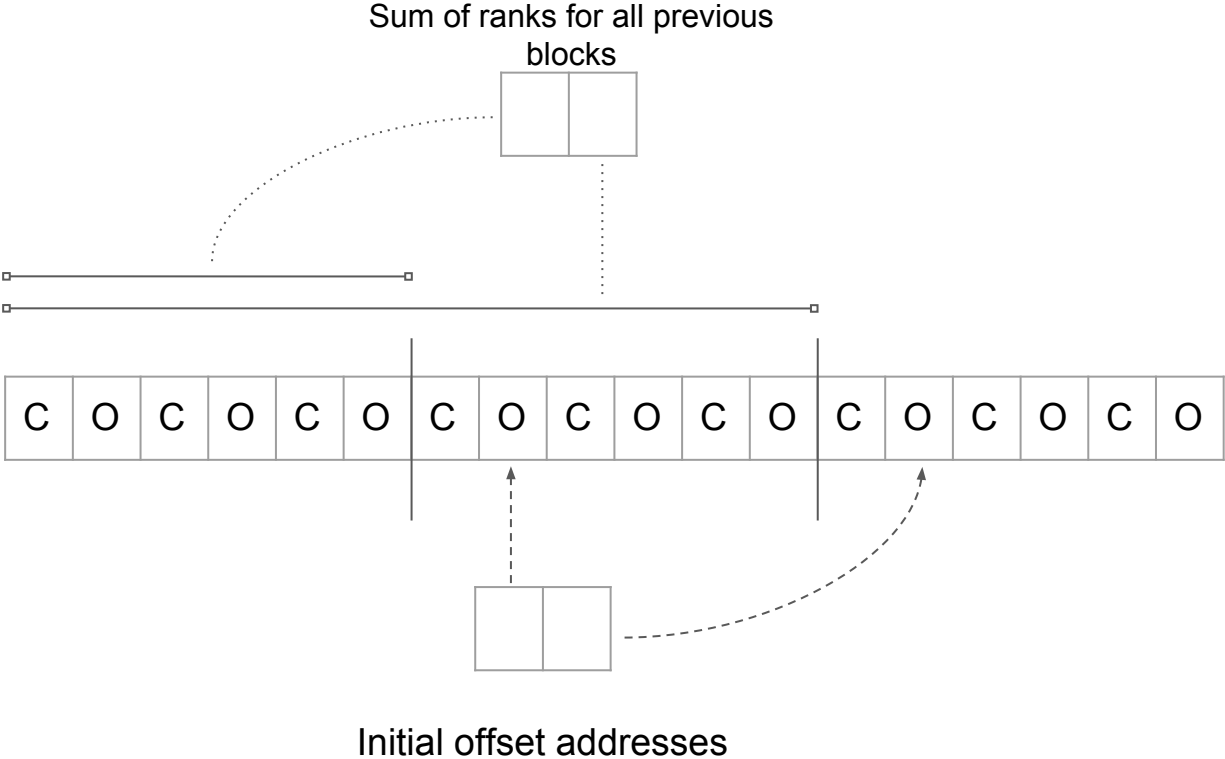
O : Offset = Index into the table

RRR



[1] R. Raman, V. Raman, and S. R. Satti. Succinct indexable dictionaries with applications to encoding k-ary trees, prefix sums and multisets. ACM Transactions on Algorithms, 3(4), 2007.

RRR

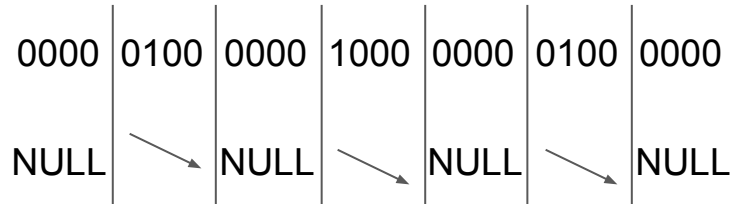


RUBIKS :RRR Update for Bit Indexing in K-mer Structure

0000 0100 0000 1000 0000 0100 0000

- Prefixes : 4^8
- Kmers : 10^9
→ 15258 “1”
- Suffixes: 4^{20}

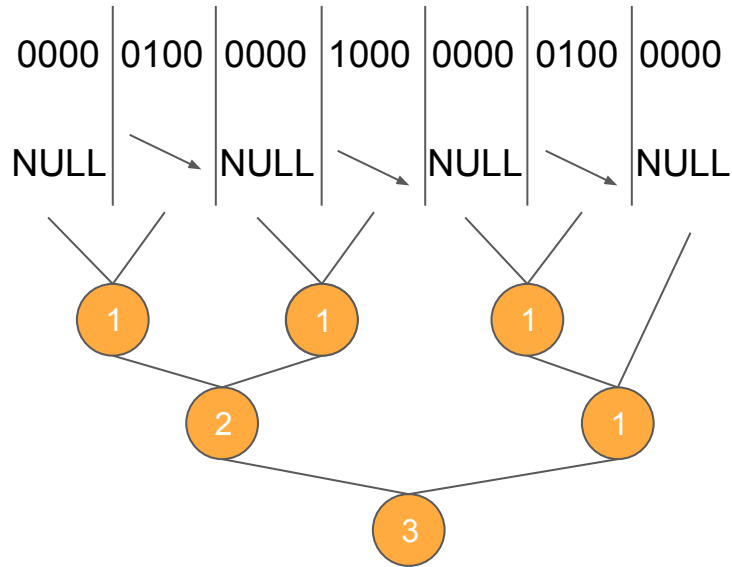
RUBIKS :RRR Update for Bit Indexing in K-mer Structure



- Prefixes : 4^8
- Kmers : 10^9
→ 15258 "1"
- Suffixes : 4^{20}

↘ Pointeur vers une RRR

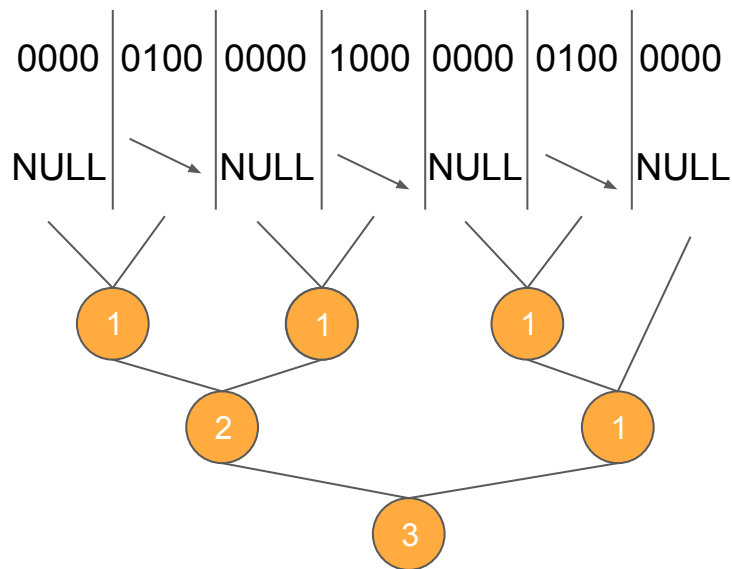
RUBIKS :RRR Update for Bit Indexing in K-mer Structure



- Prefixes : 4^8
- Kmers : 10^9
→ 15258 "1"
- Suffixes : 4^{20}

→ Pointeur vers une RRR

RUBIKS :RRR Update for Bit Indexing in K-mer Structure



- Prefixes : 4^8
- Kmers : 10^9
→ 15258 "1"
- Suffixes : 4^{20}

→ Pointeur vers une RRR

We still can do a Rank and Select !



Our outlook

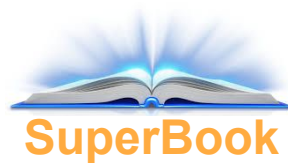
How to answer the following questions:

Is the word "ATATAAGATTACA" present in the first chromosome of genomes 644?

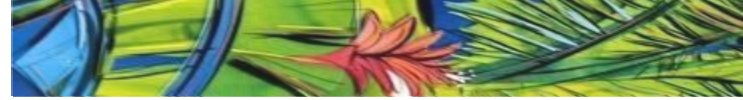
Which sequences are common to the kasalath and 9311 genomes?

Index of 3000 genomes

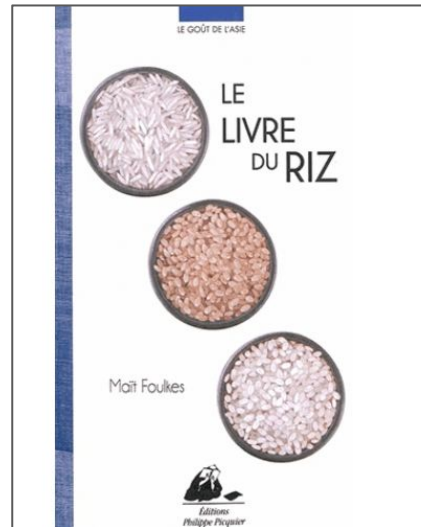
- Tools based on this index
- Integrate these tools into the GenomeHarvest project
- Create tools to make these index structure easy to use



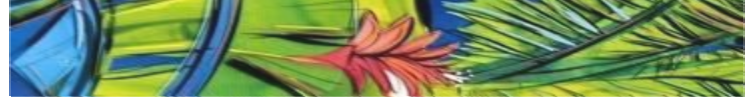
Read the SuperBook and answer questions



Thank you for your attention.
Do you have questions ?







Annexes

If searching for a word takes 1 sec per line:

→ Search for a word in a book of 500 000 000 lines:

500,000,000 sec: 5,787 Days

→ Search for a word in a dictionary of 500 000 000 lines:

Log (500,000,000) sec: 28,897 sec



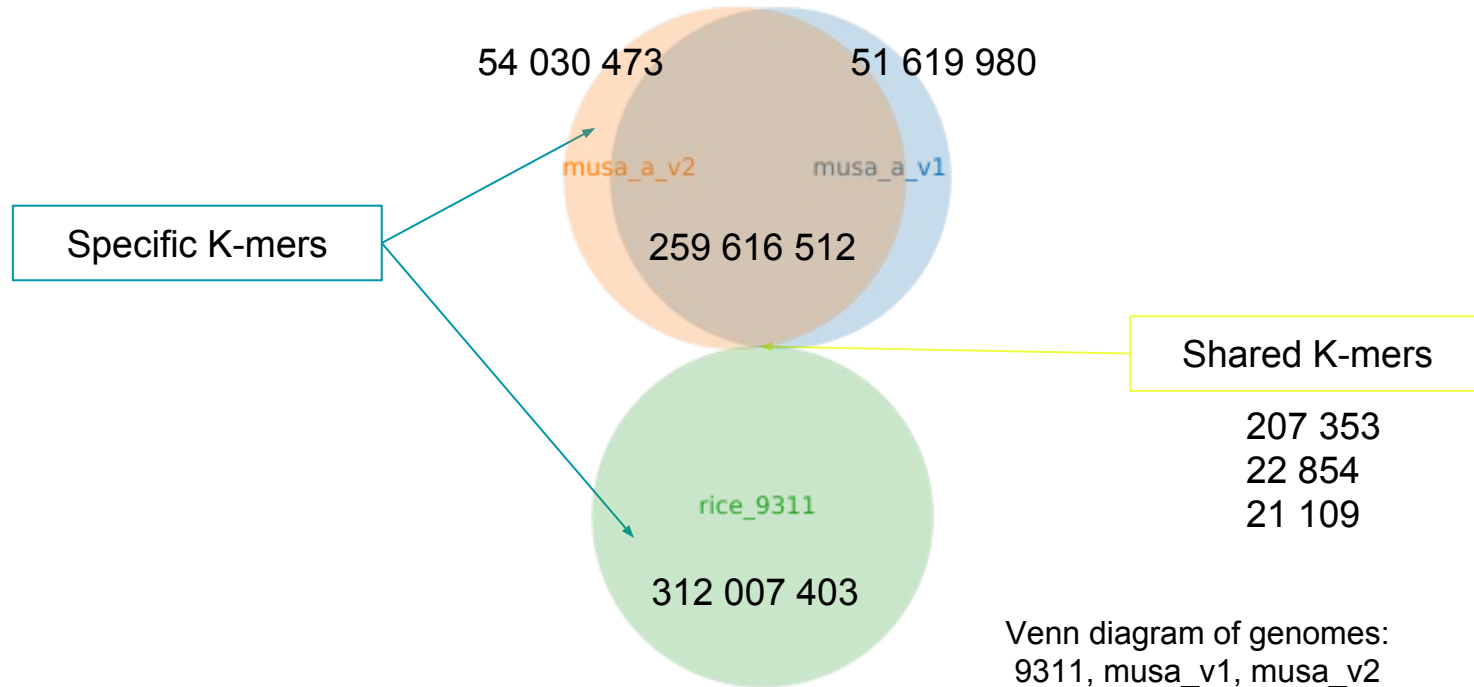
Results

(1) SDSL Lite

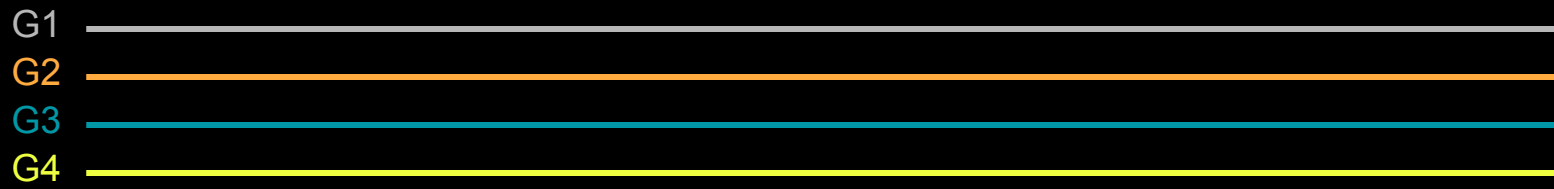
K_2	compression	b	rrrb $\log \frac{n}{m}$
14		32 Mo	2 Mo
15		128 Mo	4 Mo
16		512 Mo	32 Mo
17		2 Go	128 Mo

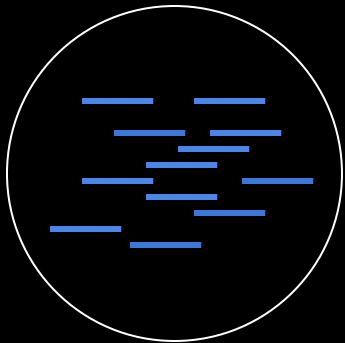
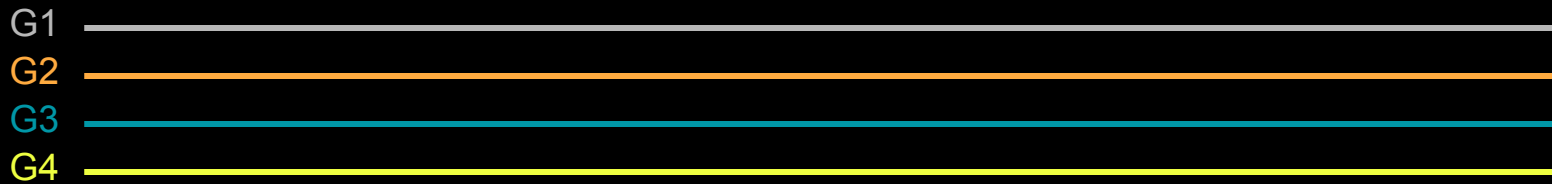


Validation of hypothesis

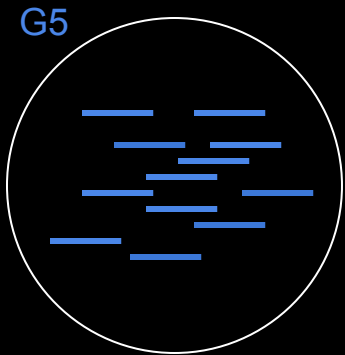
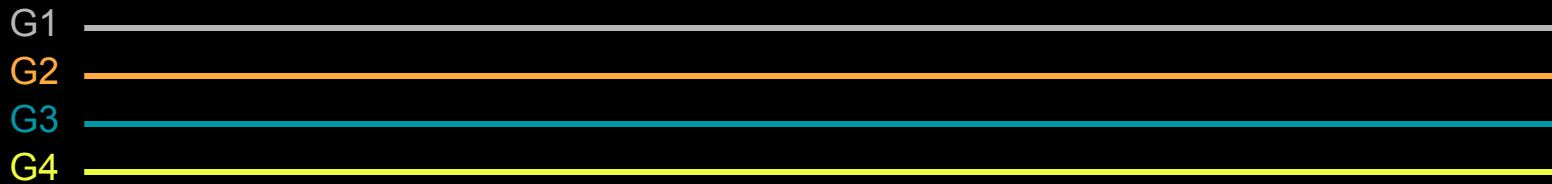


```
#ifdef __aquoicasert
```



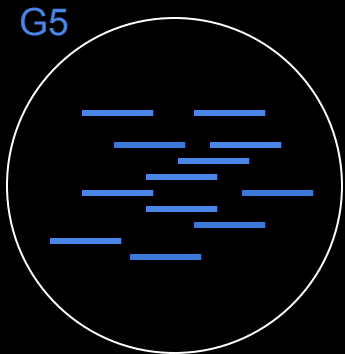
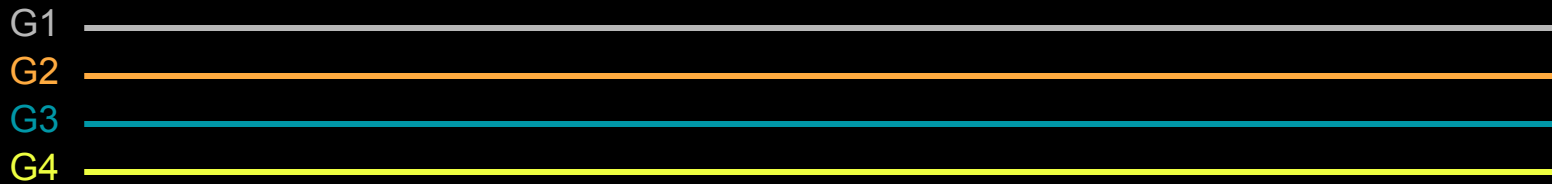


ReadSet1

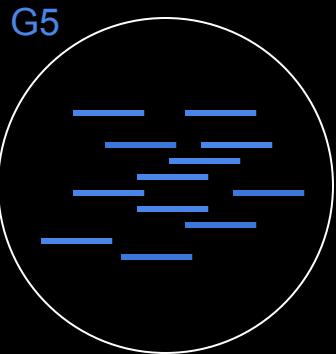
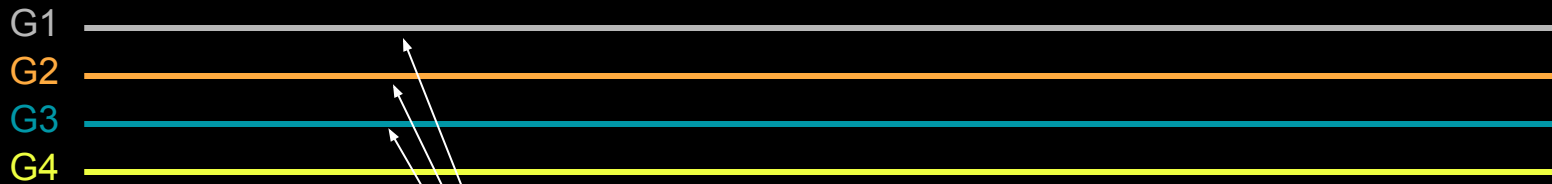


ReadSet1

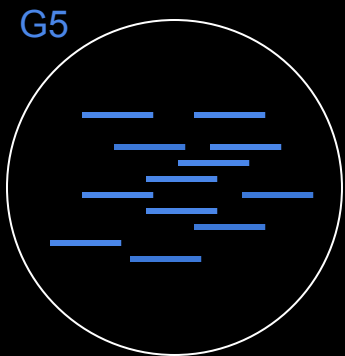
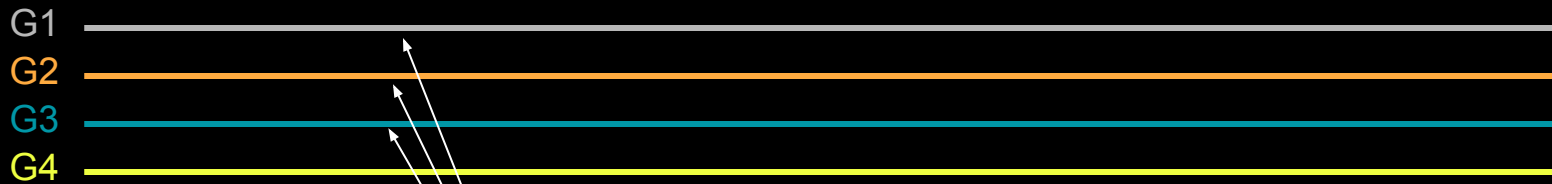




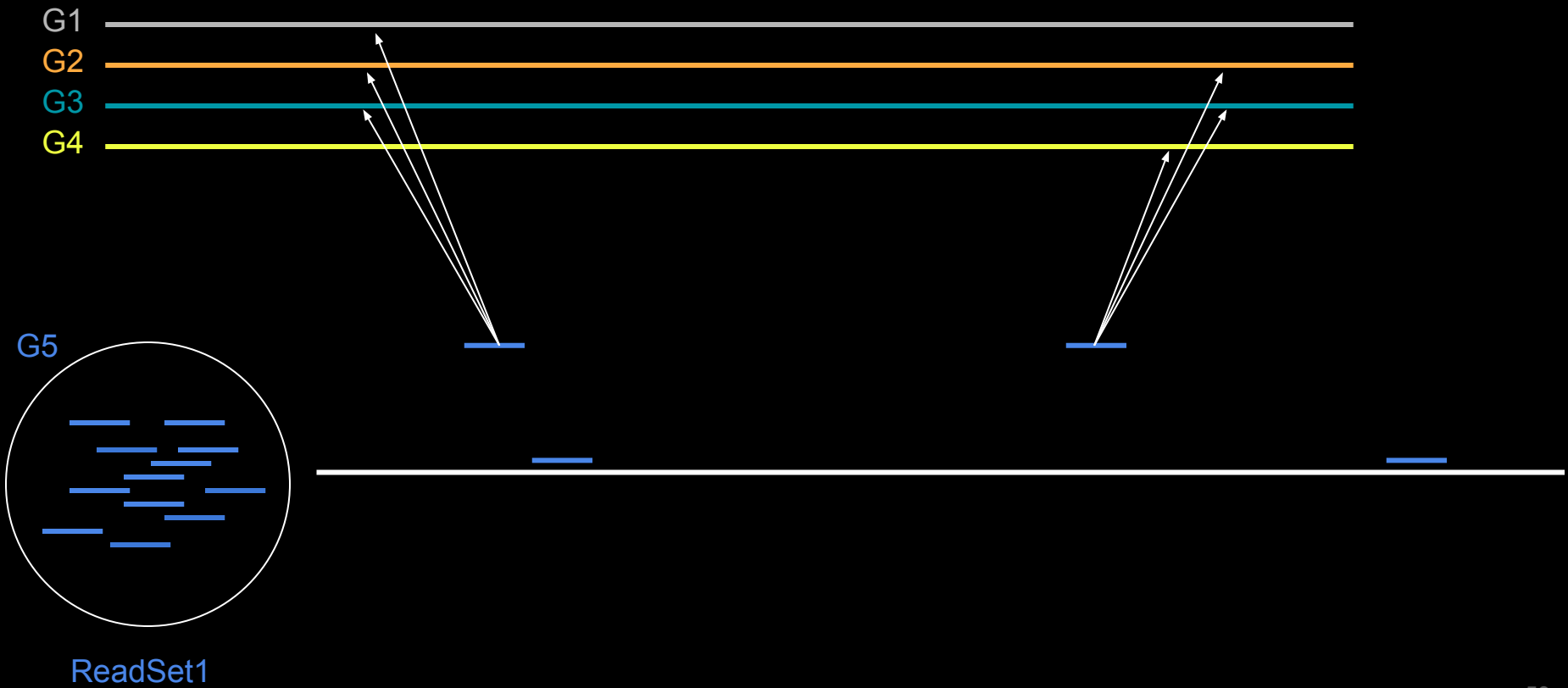
ReadSet1

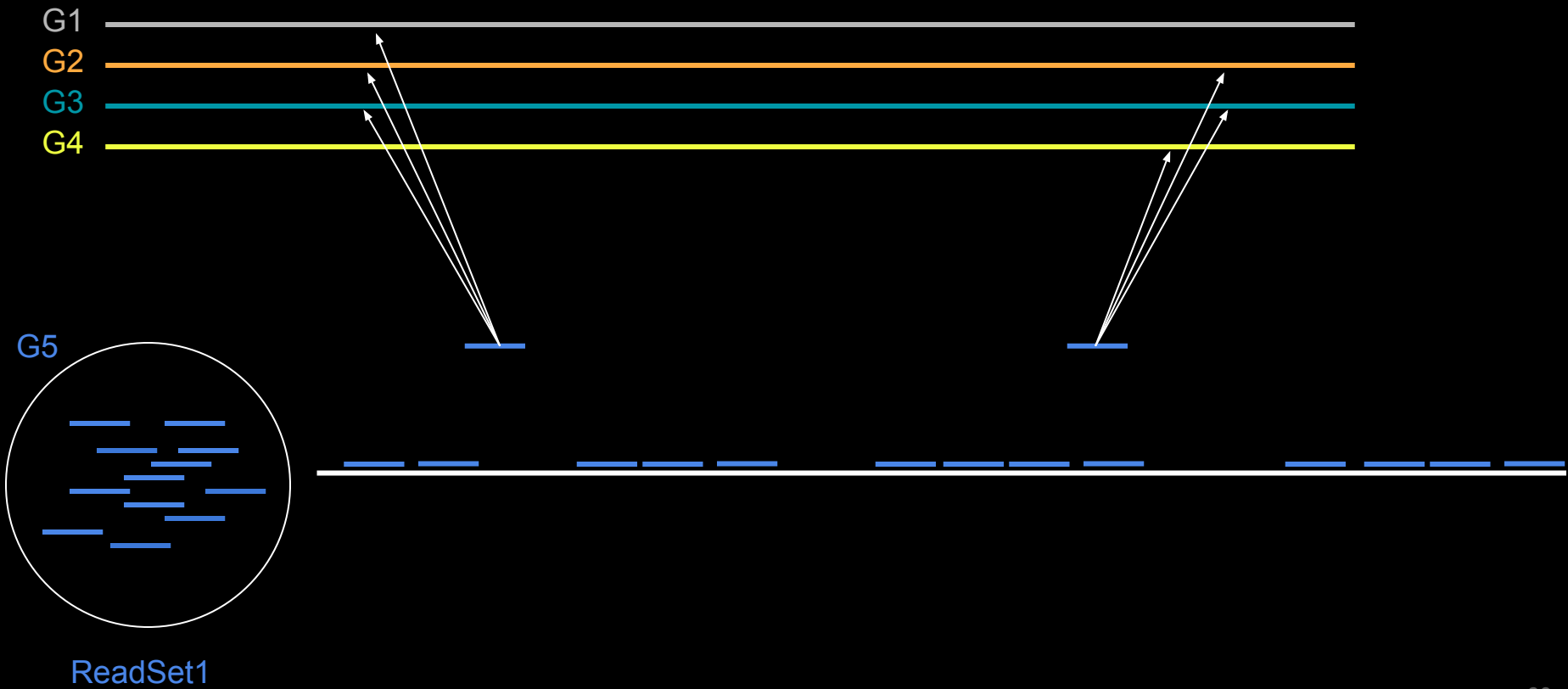


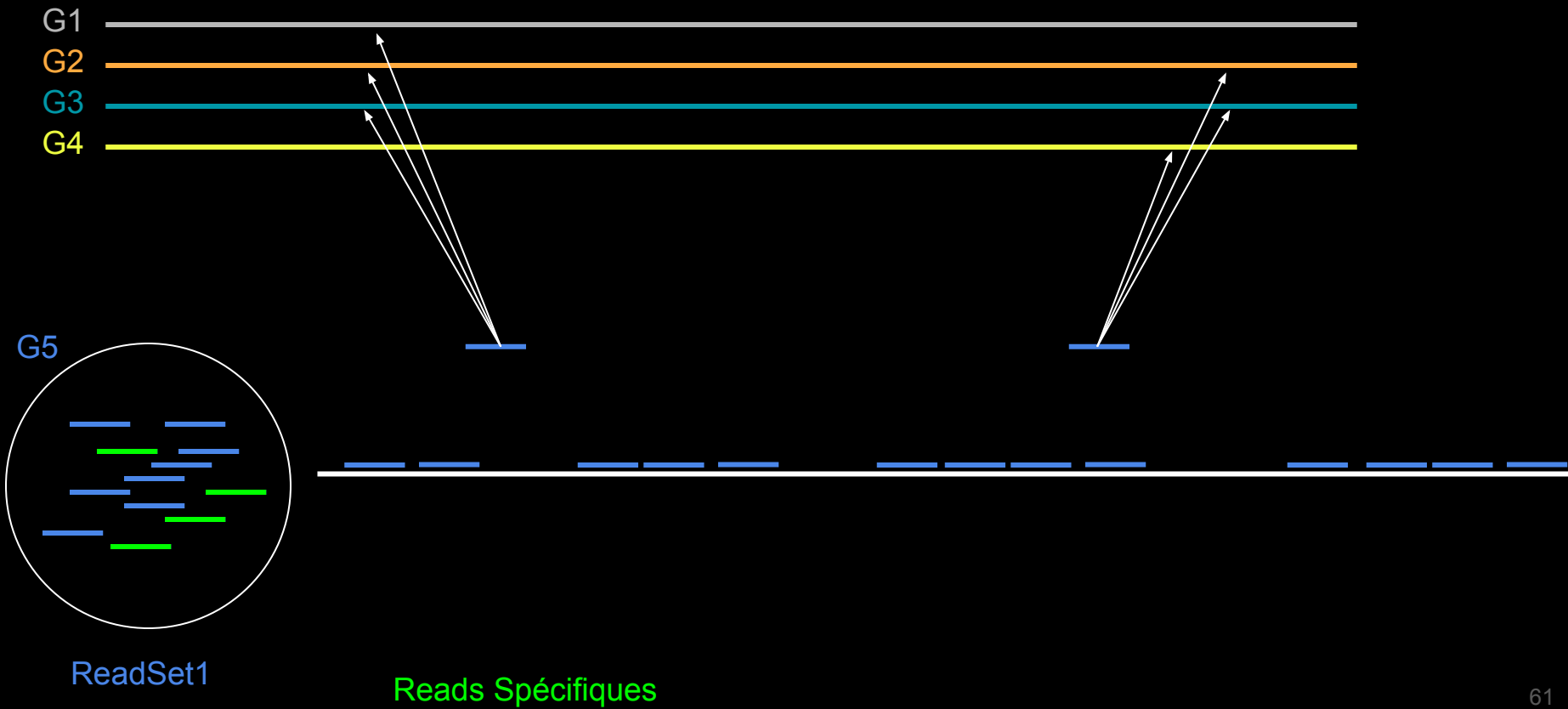
ReadSet1

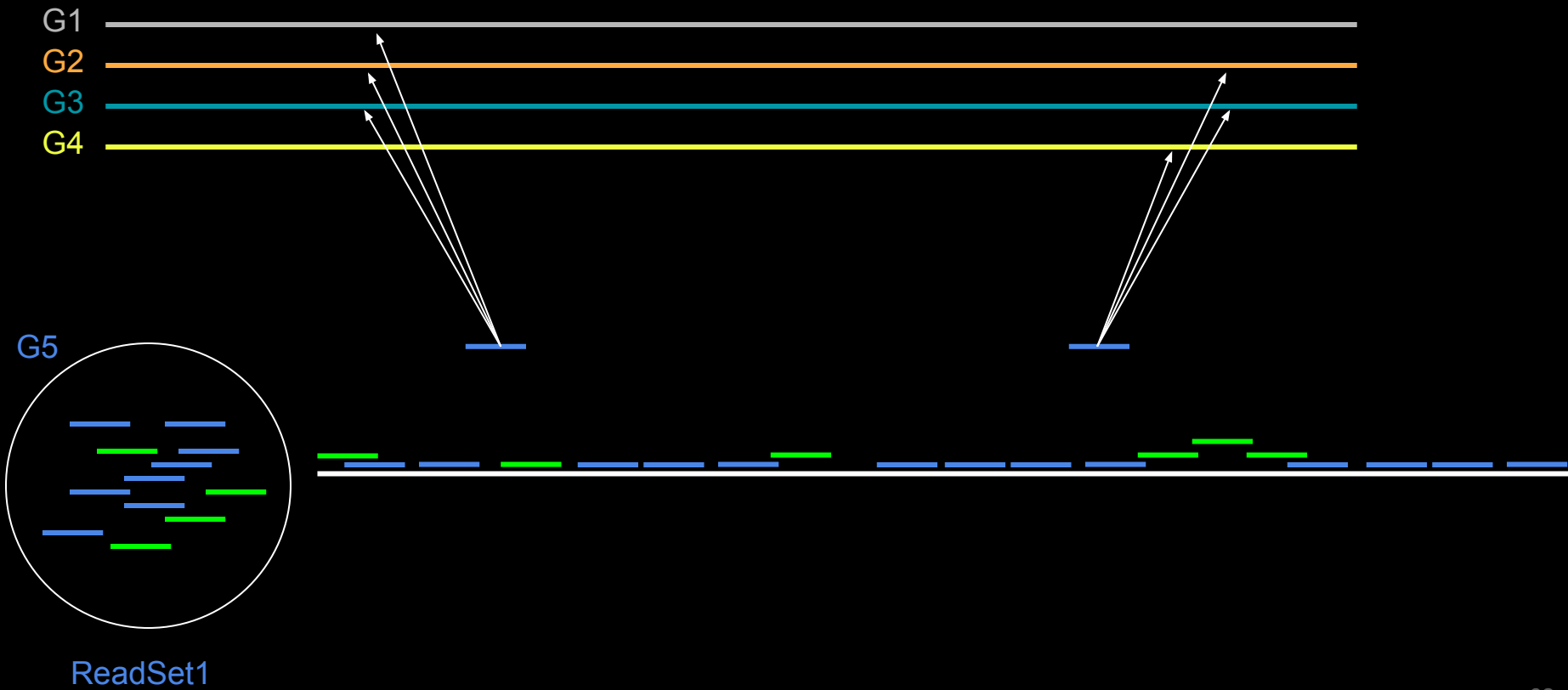


ReadSet1









```
#endif
```

