

Strutture dati per indicizzazione testi (genomici)

1

Cosa abbiamo visto

Confronto tra sequenze

- Calcolo della LCS tra 2 sequenze (PD) per studiare la similarità, globale, locale
- Allineamento multiplo

sottosequenza...
score...

Sequenziamento

- Shotgun sequencing
- Fragment assembly
 - approccio graph based

2

Cosa abbiamo visto

Confronto tra sequenze

- Calcolo della LCS tra 2 sequenze (PD) per studiare la similarità, globale, locale
- Allineamento multiplo

sottosequenza...
score...

ma potremmo volere
determinare dove una read è
maggiormente simile ad un
genoma di riferimento

Sequenziamento

- Shotgun sequencing
- Fragment assembly
 - approccio graph based

calcolo overlap
gestione k-mer
(bloom filter)



3

2 problemi «nascosti»

calcolo overlap

(**BWT** e altre strutture dati)

k-mer counting

(conteggio, rappresentazione con
Bloom Filter, Hash, minhash, **BWT**)

4

Argomenti da trattare e «raccontare»

Pan-genome graph

- ✓ Survey

Strutture di indicizzazione testi

- ✓ ST, SA, BWT, LCP [*pattern matching e overlap*]
 - Algoritmo lineare per la costruzione di OG
- ✓ BWA e BOWTIE [aligners]  [*seminario-esame*]
- ✓ Lyndon & overlap & sequence embedding

Tecniche alignment-free

- ✓ minsketch

5

Milestones

1959: Suffix Tree

1990: Suffix Array (Myers e Manber)

1993: Longest Common Prefix (LCP) Array
(Myers e Manber)

1994: Burrow-Wheeler Transform (BWT)

6

How Do We Locate Disease- Causing Mutations?

Combinatorial Pattern Matching

Phillip Compeau and Pavel Pevzner.

Bioinformatics Algorithms: an Active Learning Approach

©2018 by Compeau and Pevzner. All rights reserved

7

Outline

- **Why do we map reads? [reference-based]**
 - Using the Trie
 - From a Trie to a Suffix Tree
 - String Compression and the Burrows-Wheeler Transform
 - Inverting Burrows-Wheeler
 - Using Burrows-Wheeler for Pattern Matching
 - Finding the Matched Patterns
 - Setting Up Checkpoints
 - Inexact Matching
 - Further Applications of Read Mapping

8

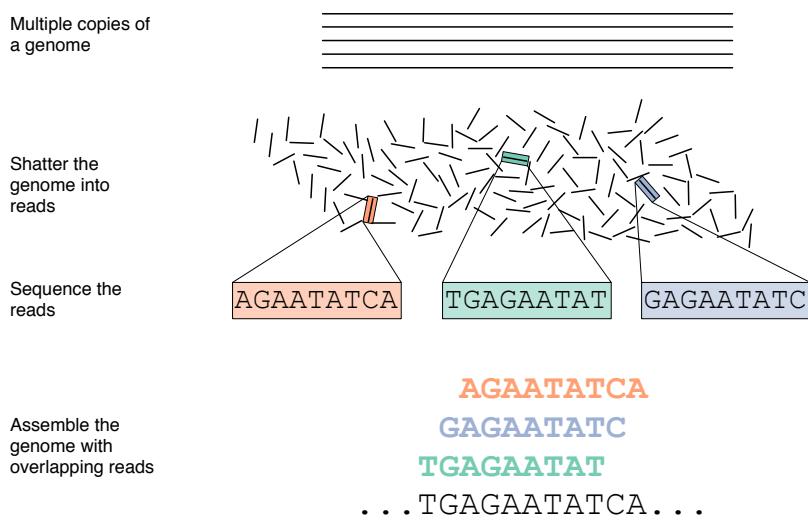
Toward a Computational Problem

- **Reference genome:** database genome used for comparison.
- **Question:** How can we assemble individual genomes (or represented by reads) efficiently using the reference?

CTGATGATGGACTACGCTACTACTGCTAGCTGTAT Individual
 CTGAGGATGGACTACGCTACTACTGATAGCTGTTT Reference

9

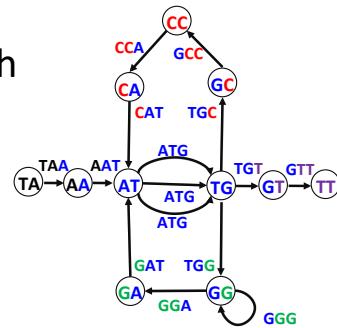
Why Not Use Assembly?



10

Answer1: Why Not Use Assembly?

- Constructing a de Bruijn graph takes a lot of memory.
- Hope: a machine in a clinic that would collect and map reads in 10 minutes.
- Idea: use existing structure of reference genome to help us sequence a patient's genome.



11

Answer2:

<https://training.galaxyproject.org/training-material/>

... In principle, we could do a BLAST analysis to figure out where the sequenced pieces fit best in the known genome.

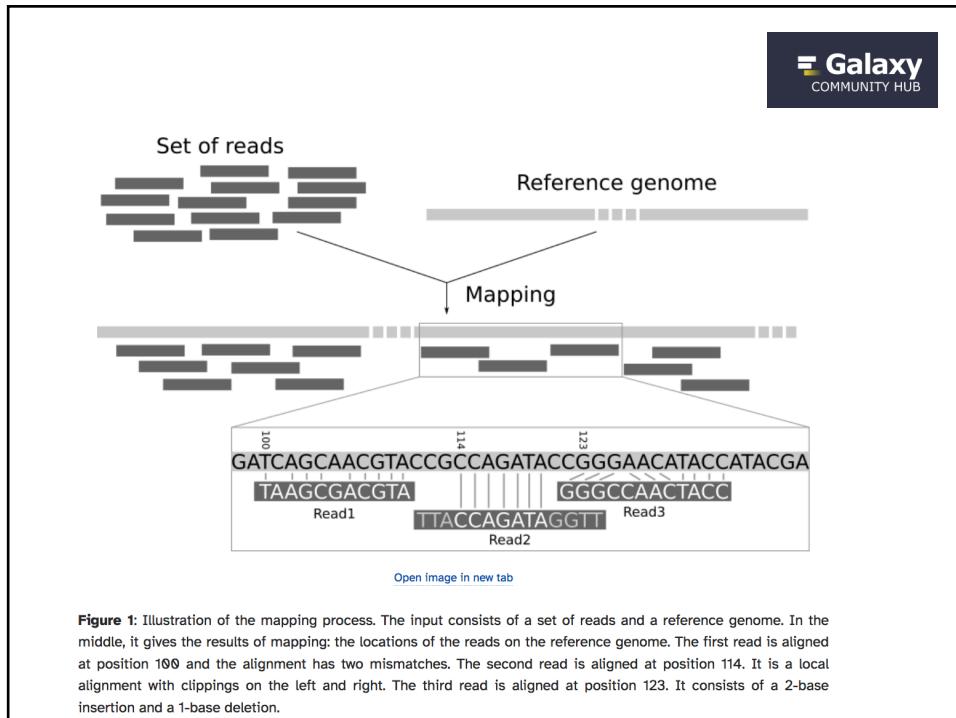


We would need to do that for each of the millions of reads in our sequencing data. Aligning millions of short sequences this way may, however, take a couple of weeks.

And we do not care about the exact base to base correspondence (alignment). What we are interested in is "**where these reads came from**".

This approach is called **mapping**.

12



13

- Input: Set of reads + a reference genome
- Output: BAM file ([Binary Alignment Map](#)) file is a compressed binary file storing the read sequences, whether they have been aligned to a reference sequence (e.g. a chromosome), and if so, the position on the reference sequence at which they have been aligned

14

Read Mapping

- **Read mapping:** determine where each read has high similarity to the reference genome.

without assembling

CTGAGGATGGACTACGCTACTACTGATAGCTGTTT	Reference
GAGGA	Reads
CCACG	TGA-A
exact	mismatch
	insdel

15

Exact Pattern Matching

- Focus on a simple question: where do the reads match the reference genome *exactly*?
- **Single Pattern Matching Problem:**
 - **Input:** A string *Pattern* (read) and a string *Genome*.
 - **Output:** All positions in *Genome* where *Pattern* appears as a substring.

16

Exact Pattern Matching

- Focus on a simple question: where do the reads match the reference genome *exactly*?
- **Multiple Pattern Matching Problem:**
 - **Input:** A collection of strings *Patterns* (reads) and a string *Genome*.
 - **Output:** All positions in *Genome* where a string from *Patterns* appears as a substring.

17

A Brute Force Approach

- We can simply iterate a brute force approach method, *sliding* each *Pattern* down *Genome*.

pan am aban anas
Genome
Pattern

- **Note:** we use words instead of DNA strings for convenience.

18

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

panamabananas	<i>Genome</i>
nana	<i>Pattern</i>

- **Note:** we use words instead of DNA strings for convenience.

19

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

p anamabananas	<i>Genome</i>
n ana	<i>Pattern</i>

- **Note:** we use words instead of DNA strings for convenience.

20

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

p a n a m a b a n a n a s *Genome*
n a n a *Pattern*

- **Note:** we use words instead of DNA strings for convenience.

21

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

p a n a m a b a n a n a s *Genome*
n a n a *Pattern*

- **Note:** we use words instead of DNA strings for convenience.

22

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

p a n a m a b a n a n a s	<i>Genome</i>
n a n a	<i>Pattern</i>

- **Note:** we use words instead of DNA strings for convenience.

23

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

p a n a m a b a n a n a s	<i>Genome</i>
n a n a	<i>Pattern</i>

- **Note:** we use words instead of DNA strings for convenience.

24

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

pan a m a b a n a n a s	Genome
n a n a	Pattern

- **Note:** we use words instead of DNA strings for convenience.

25

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

panamabanan as *Genome*
n a n a *Pattern*

- **Note:** we use words instead of DNA strings for convenience.

26

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

panam**a**bananas *Genome*
 nana *Pattern*

- **Note:** we use words instead of DNA strings for convenience.

27

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

panam**b**ananas *Genome*
 nana *Pattern*

- **Note:** we use words instead of DNA strings for convenience.

28

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

panamab**a**nanas *Genome*
 nana *Pattern*

- **Note:** we use words instead of DNA strings for convenience.

29

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

panamab**a**nanas *Genome*
 nana *Pattern*

- **Note:** we use words instead of DNA strings for convenience.

30

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

panamabanananas *Genome*
 nana *Pattern*

- **Note:** we use words instead of DNA strings for convenience.

31

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

panamabanananas *Genome*
 nana *Pattern*

- **Note:** we use words instead of DNA strings for convenience.

32

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

panamabananas *Genome*
nana *Pattern*

- **Note:** we use words instead of DNA strings for convenience.

33

A Brute Force Approach

- We can simply iterate a brute force approach method, sliding each *Pattern* down *Genome*.

panamabananas *Genome*
nana *Pattern*

- **Note:** we use words instead of DNA strings for convenience.

34

That Was Easy!



35

Brute Force Is Too Slow

- The runtime of the brute force approach is too high!
 - Single Pattern: $O(|Genome| * |Pattern|)$
 - Multiple Patterns: $O(|Genome| * |Patterns|)$
 - $|Patterns|$ = combined length of *Patterns*

36

Outline

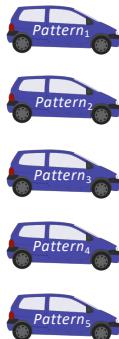
- Why do we map reads?
- **Using the Trie**
- From a Trie to a Suffix Tree
- String Compression and the Burrows-Wheeler Transform
- Inverting Burrows-Wheeler
- Using Burrows-Wheeler for Pattern Matching
- Finding the Matched Patterns
- Setting Up Checkpoints
- Inexact Matching
- Further Applications of Read Mapping

37

Patterns Travel One at a Time

Genome

```
CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAGCTACCACATCGTAGCTAC
```



38

Packing Patterns onto a Bus

Genome

```
CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAGCTACCACATCGTAGCTAC
```

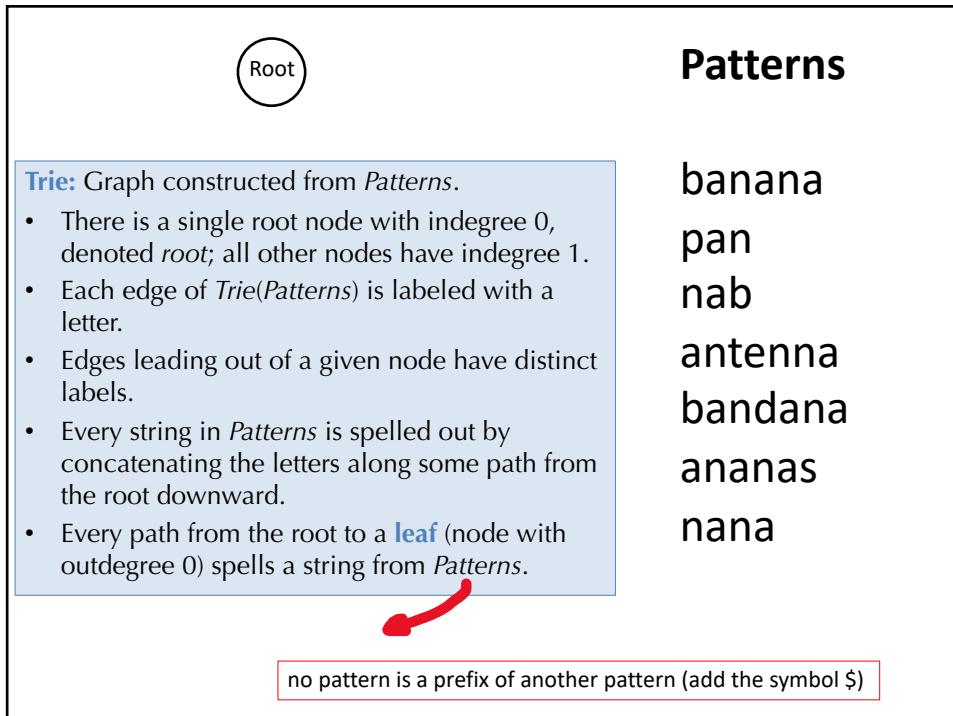


39

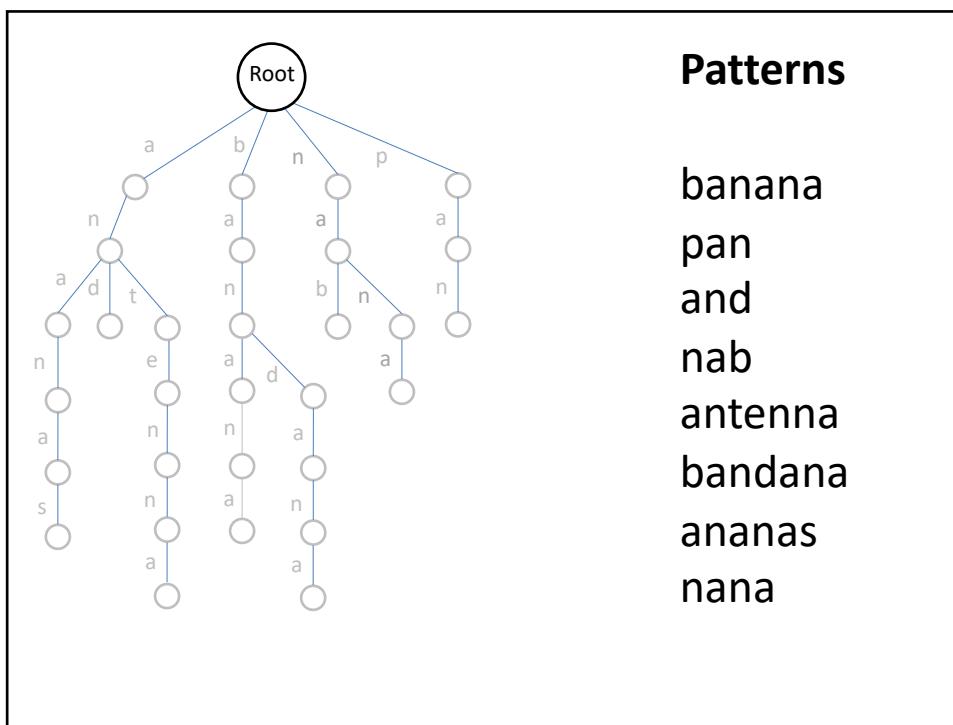
Processing Patterns into a Trie

- Idea: combine reads into a graph. *Each substring of the genome can match at most one read (no read is a prefix of another read)*. So each read will correspond to a unique path through this graph.
- The resulting graph [represents a collection of strings] is called a **trie**.

40



41



42

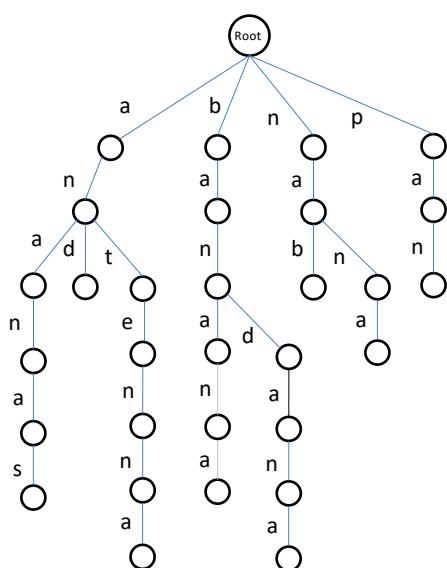
Using the Trie for Pattern Matching

- **TrieMatching:** *Slide* the trie down the genome.
- At each position, walk down the trie and see if we can reach a leaf by matching symbols.
- Analogy: bus stops



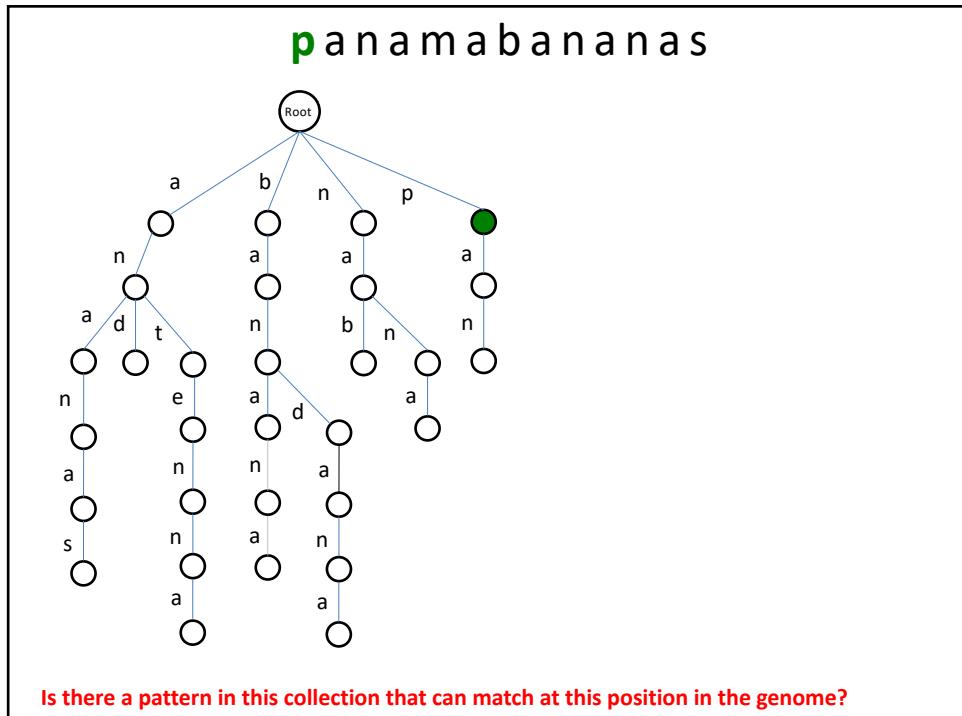
43

panamabananas

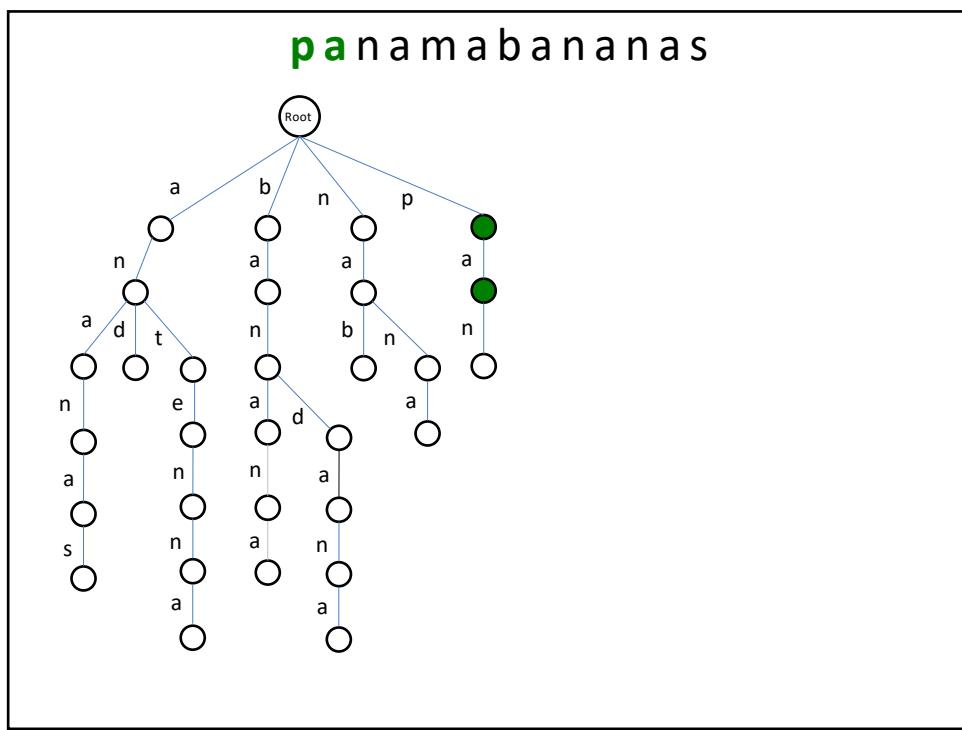


Given the genome and the trie constructed from the set of pattern we want to locate

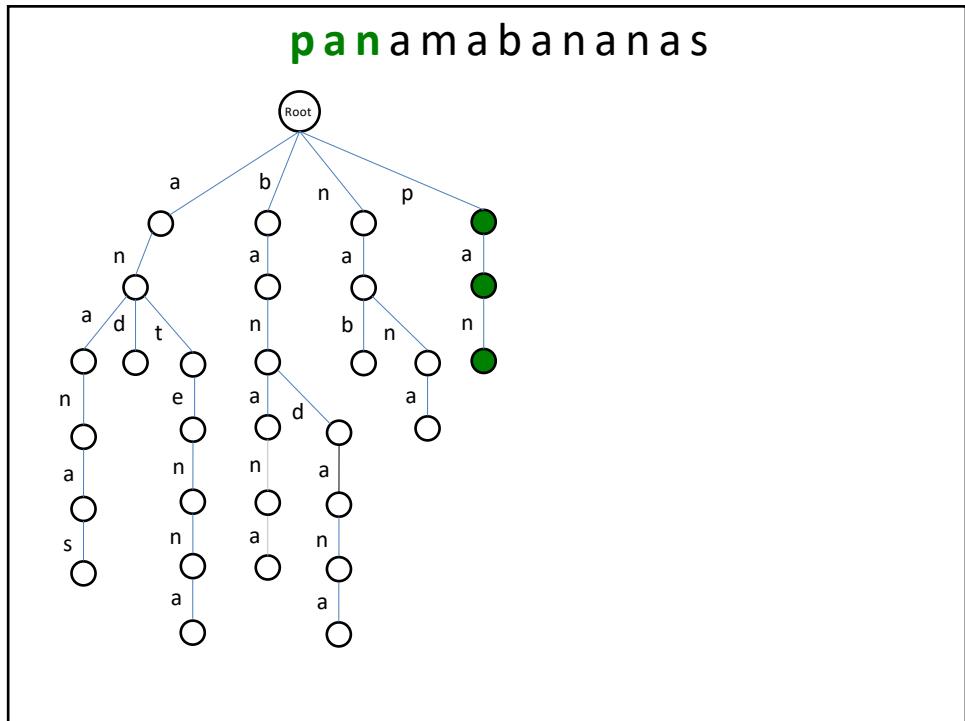
44



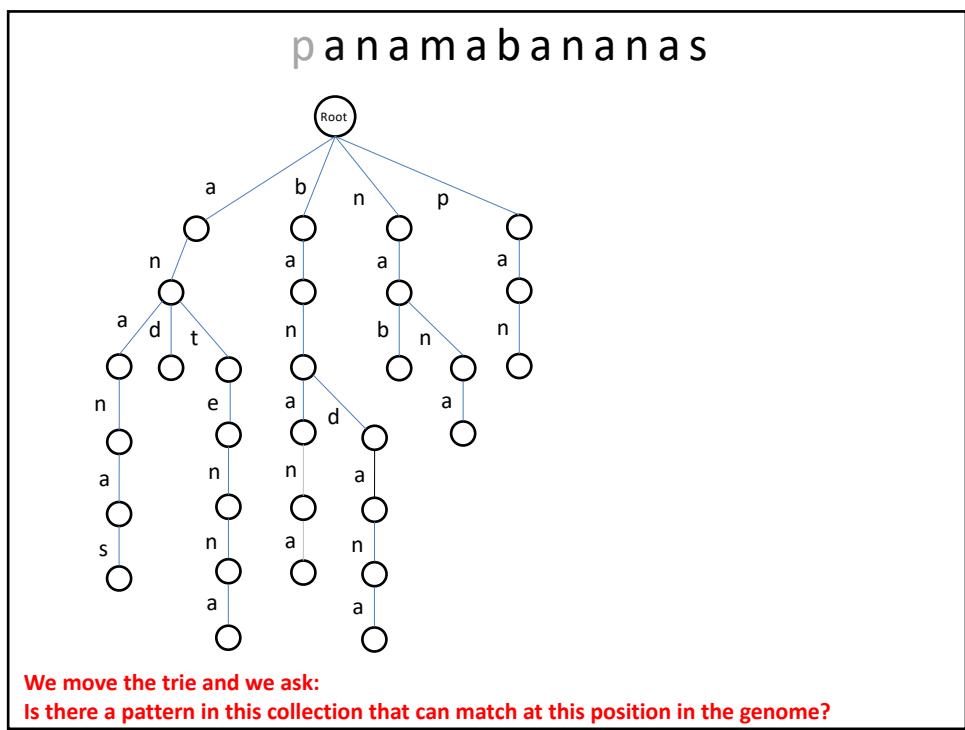
45



46



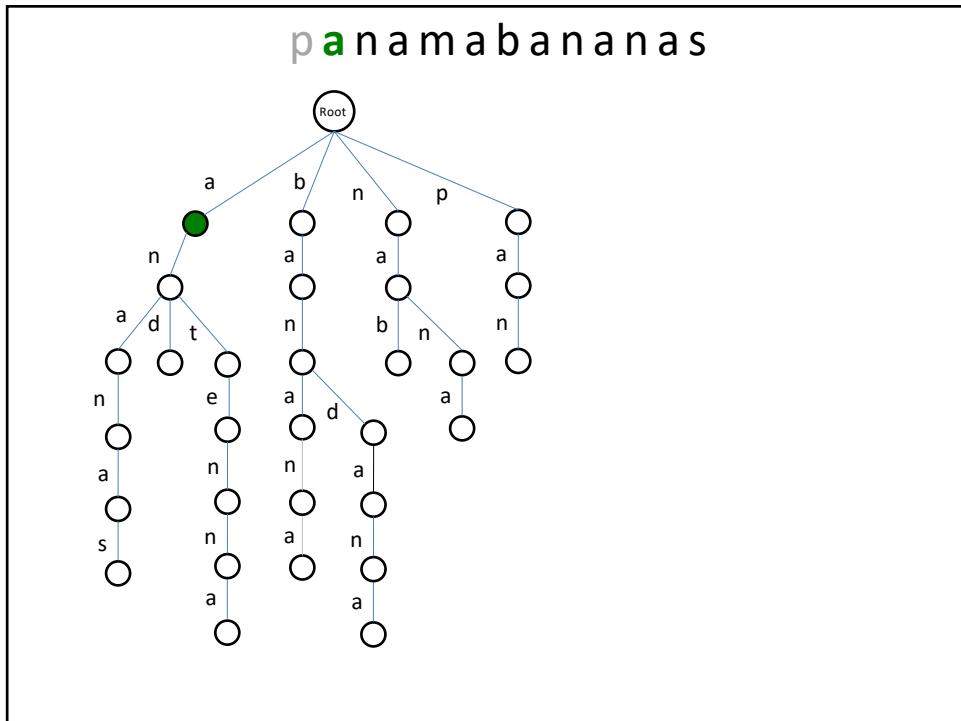
47



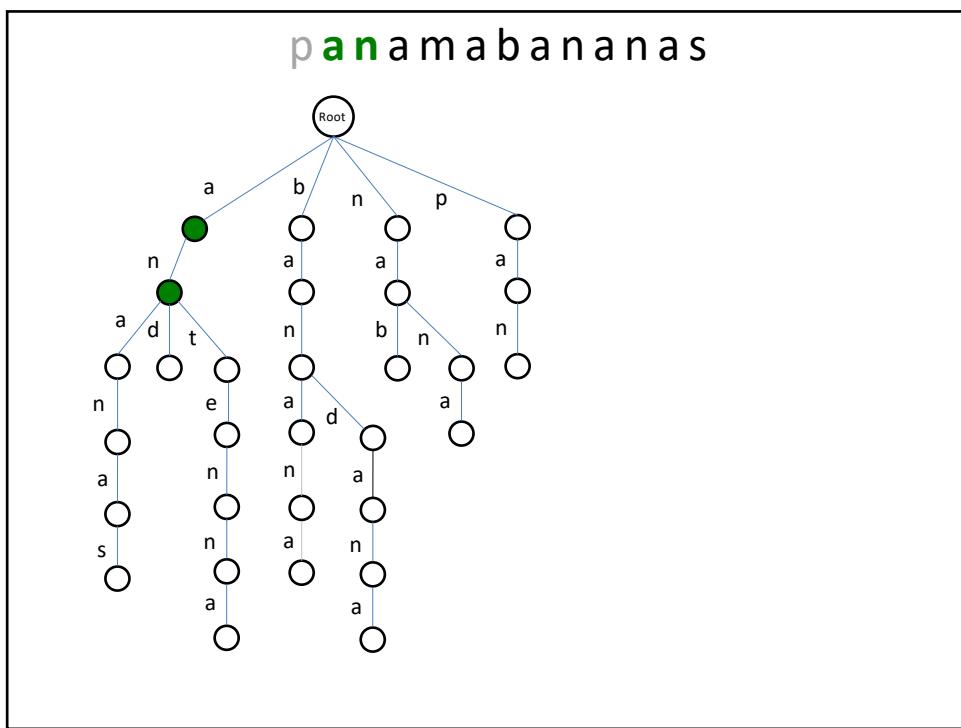
We move the trie and we ask:

Is there a pattern in this collection that can match at this position in the genome?

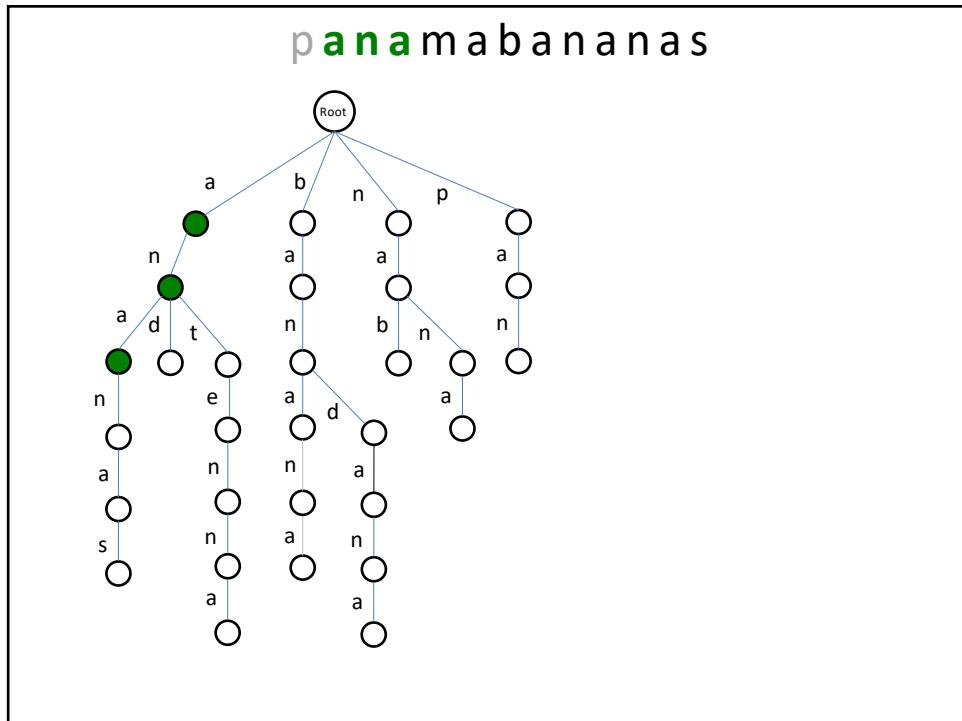
48



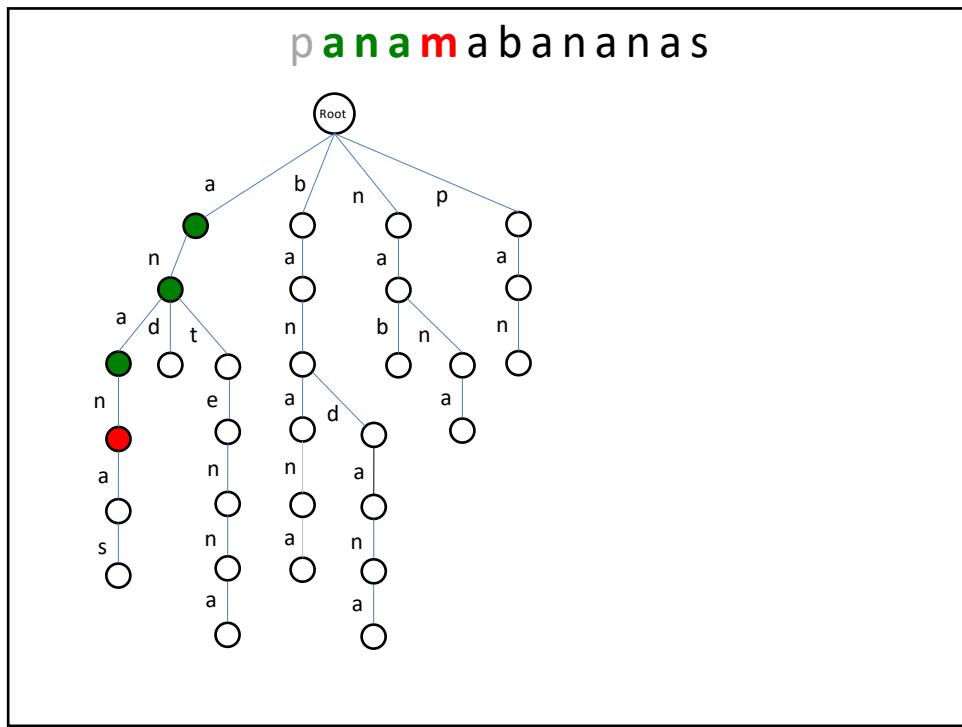
49



50

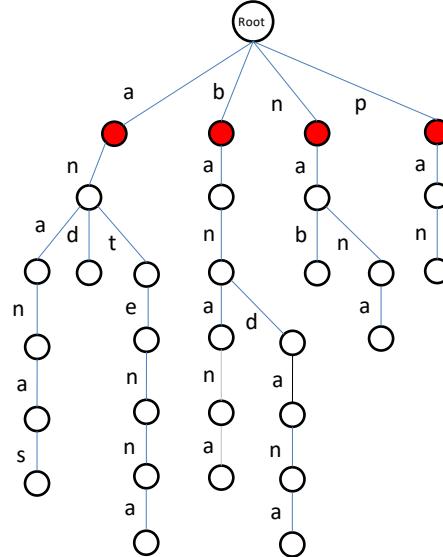


51



52

panamabananas **s**



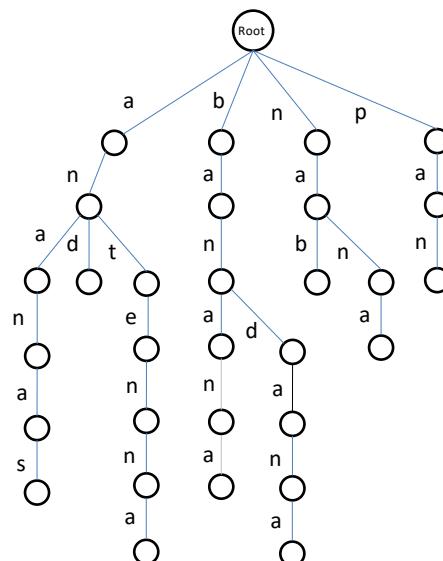
53

panamabananas

STOP: If we add the pattern “band” to our dataset, why would it cause problems?

Answer: “band” is a prefix of “bandana” and would not end at a leaf node ...

Adding an “End” Symbol Fixes the Issue



54

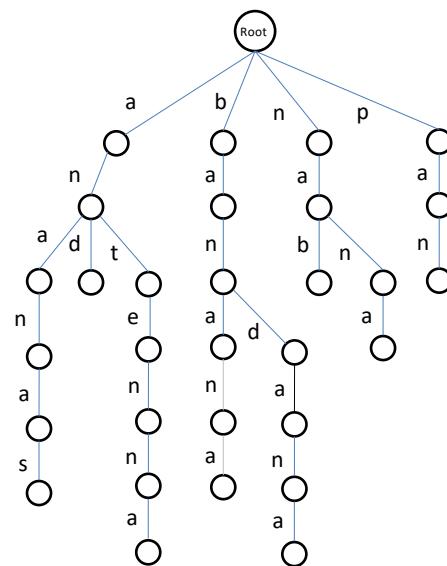
Success!

- Runtime of Brute Force (before):
 - Total: $O(|Genome| * |Patterns|)$
- Runtime of Trie Matching:
 - Trie Construction: $O(|Patterns|)$
 - Pattern Matching: $O(|Genome| * |LongestPattern|)$

55

Memory Analysis of TrieMatching

- We completely forgot about memory!
- Our trie: 30 edges, $|Patterns| = 39$
- Worst case: # edges = $O(|Patterns|)$

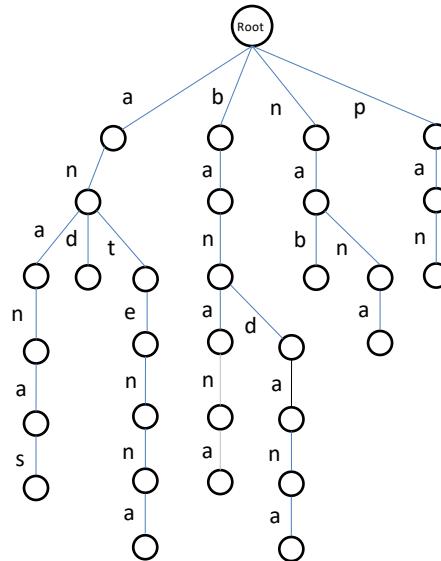


56

Memory Analysis of TrieMatching

- We completely forgot about memory!
- Our trie: 30 edges, $|Patterns| = 39$

Worst case: Trie takes up $O(|Patterns|)$, which for a human sequencing project may be ~ 1 TB!



57

Outline

- Why do we map reads?
- Using the Trie
- **From a Trie to a Suffix Tree**
- String Compression and the Burrows-Wheeler Transform
- Inverting Burrows-Wheeler
- Using Burrows-Wheeler for Pattern Matching
- Finding the Matched Patterns
- Setting Up Checkpoints
- Inexact Matching
- Further Applications of Read Mapping

58

PREPROCESSING THE GENOME INSTEAD

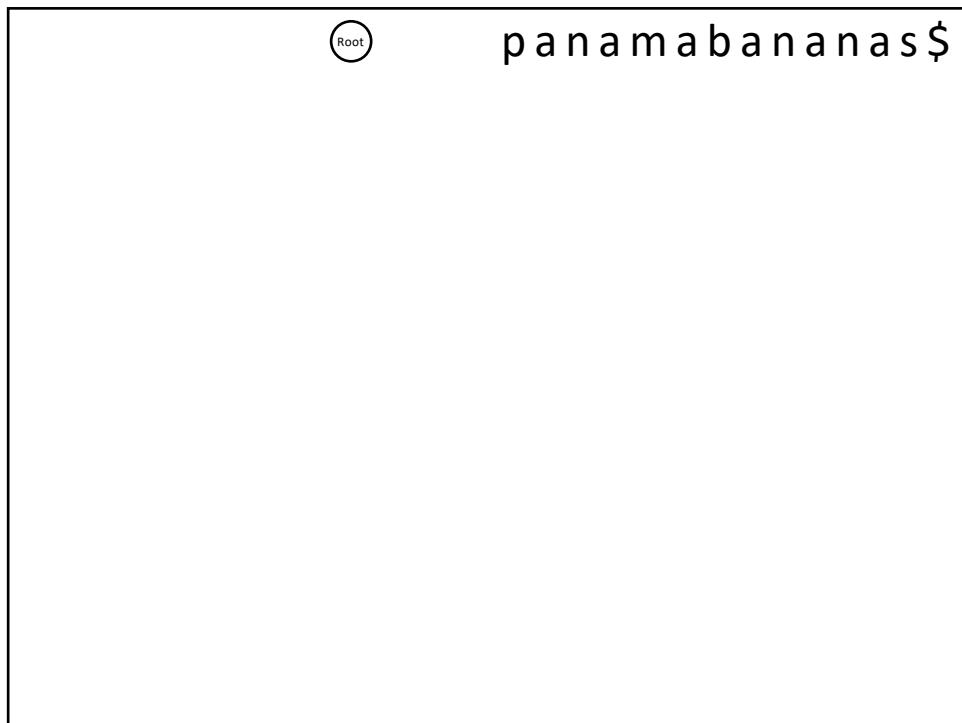
59

Preprocessing the Genome

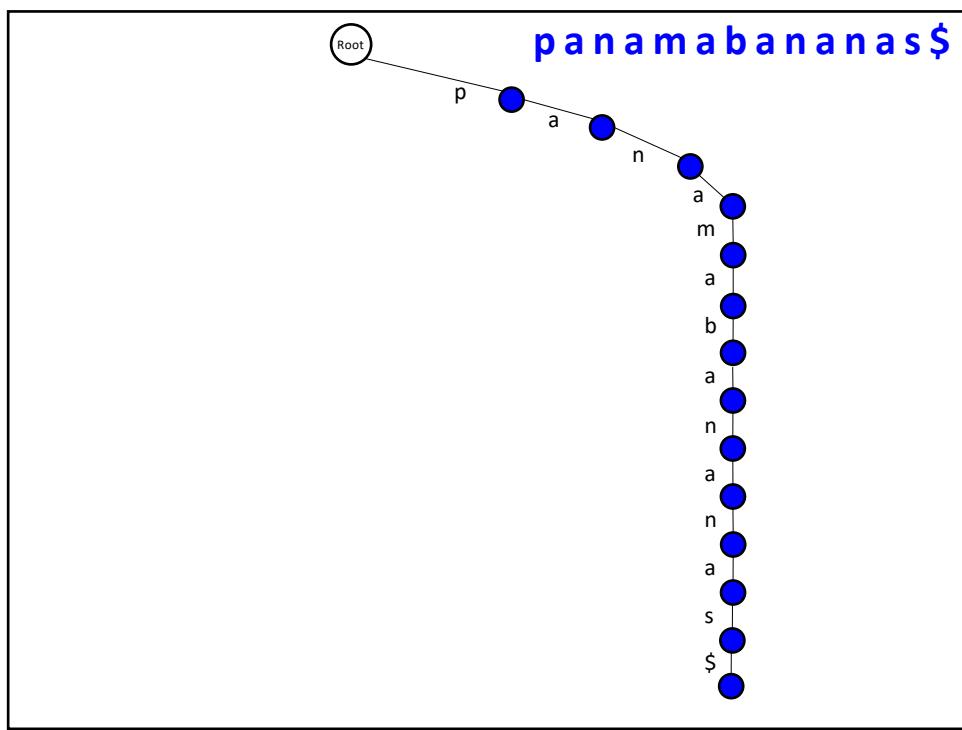
- Split *Genome* into all its suffixes.
- How can we combine these suffixes into a data structure?
- Let's use a trie!

...each matched read is a prefix of a suffix of the Genome... i.e., we find it in a path from the root...

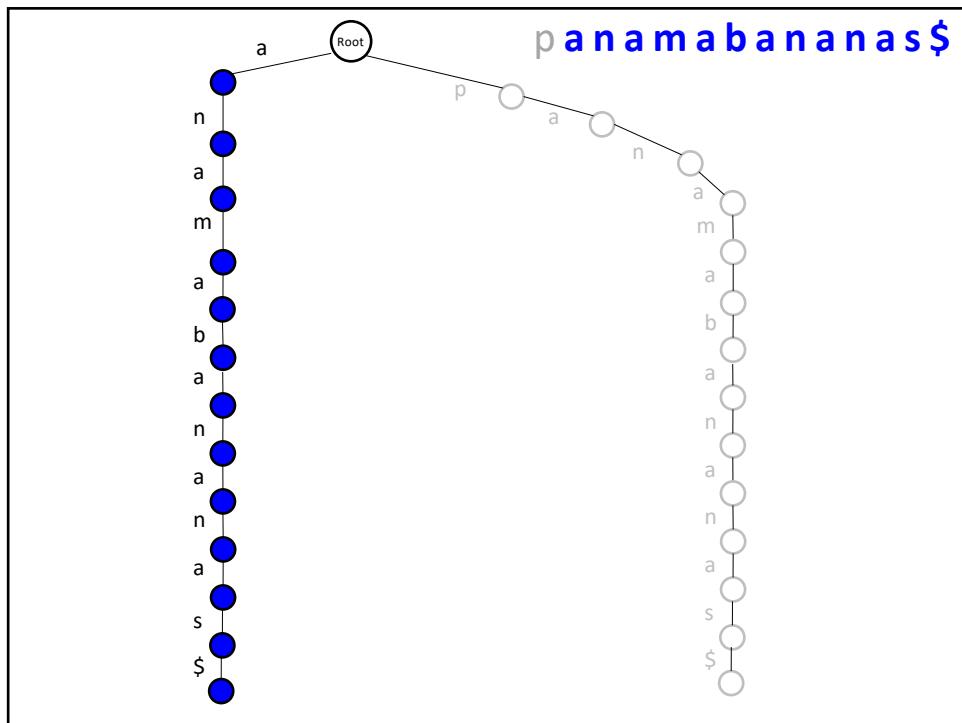
60



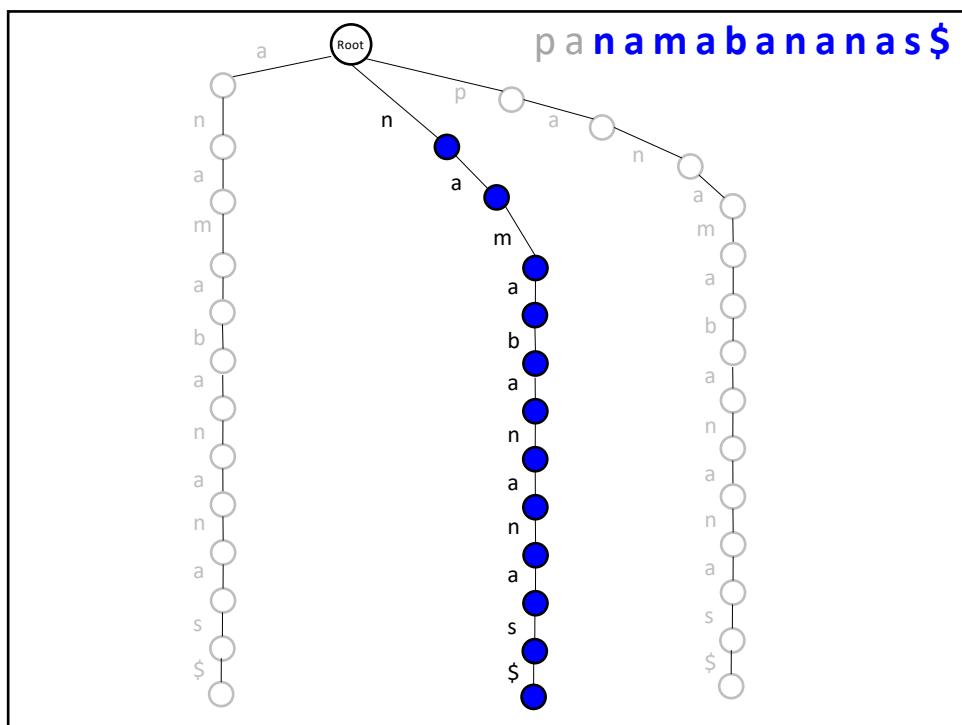
61



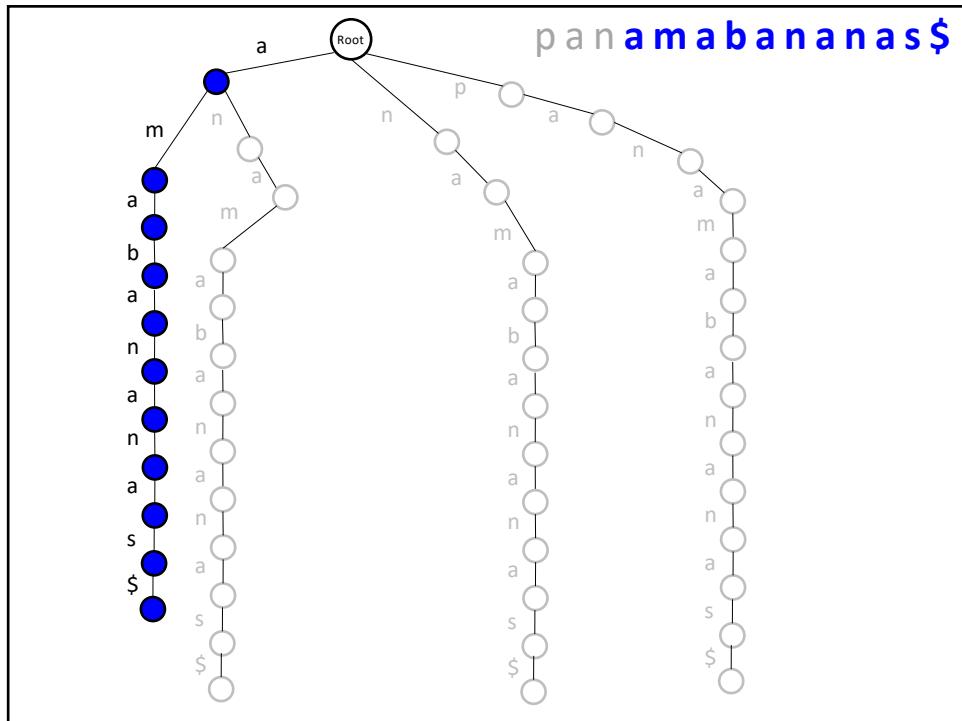
62



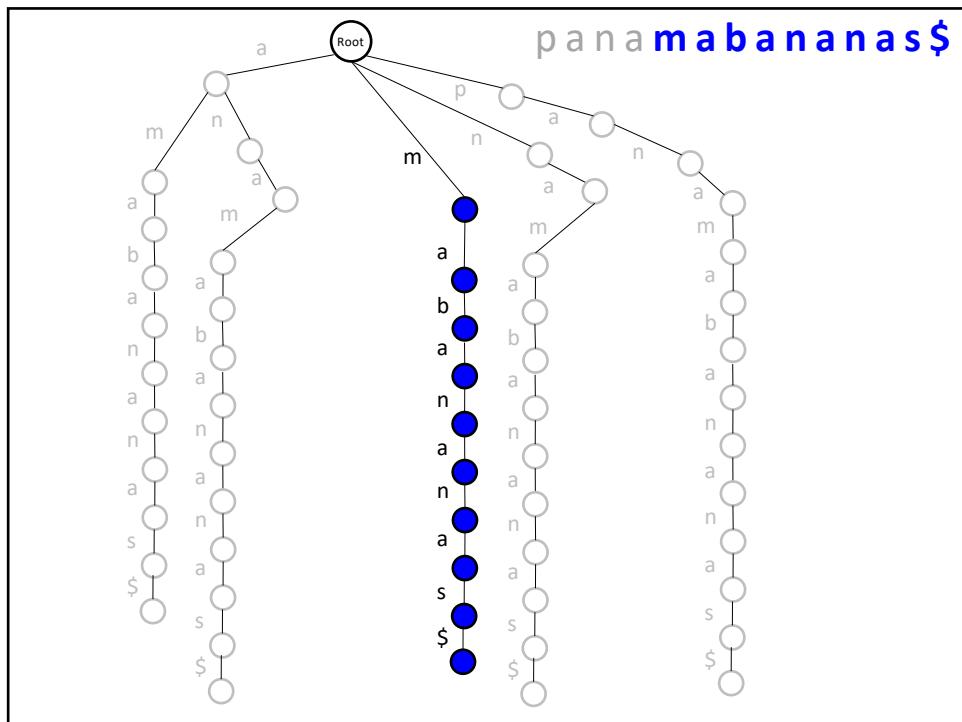
63



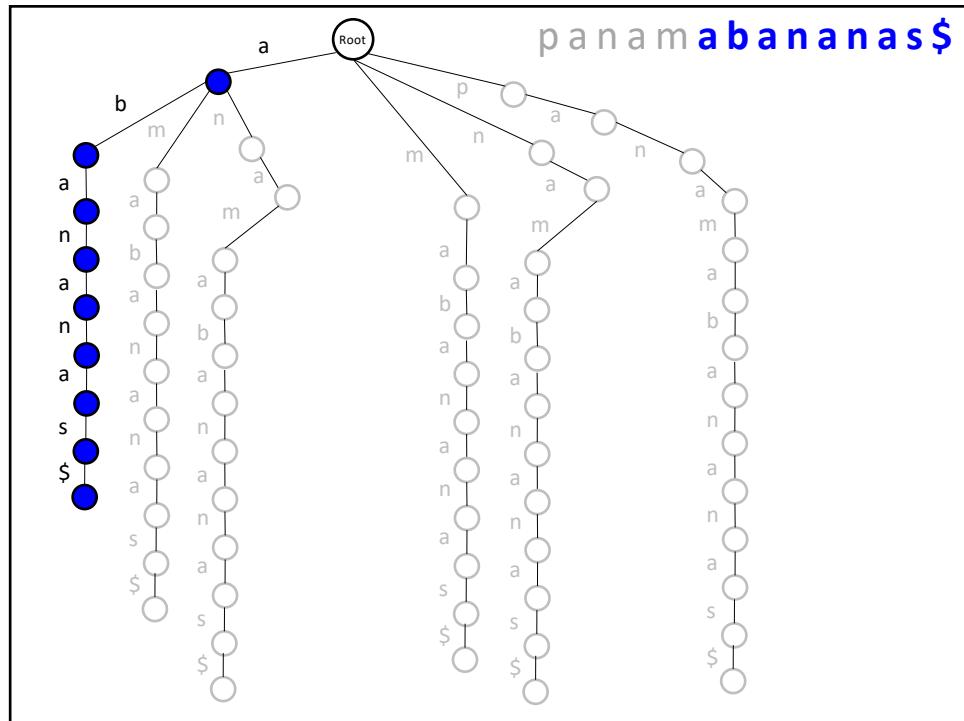
64



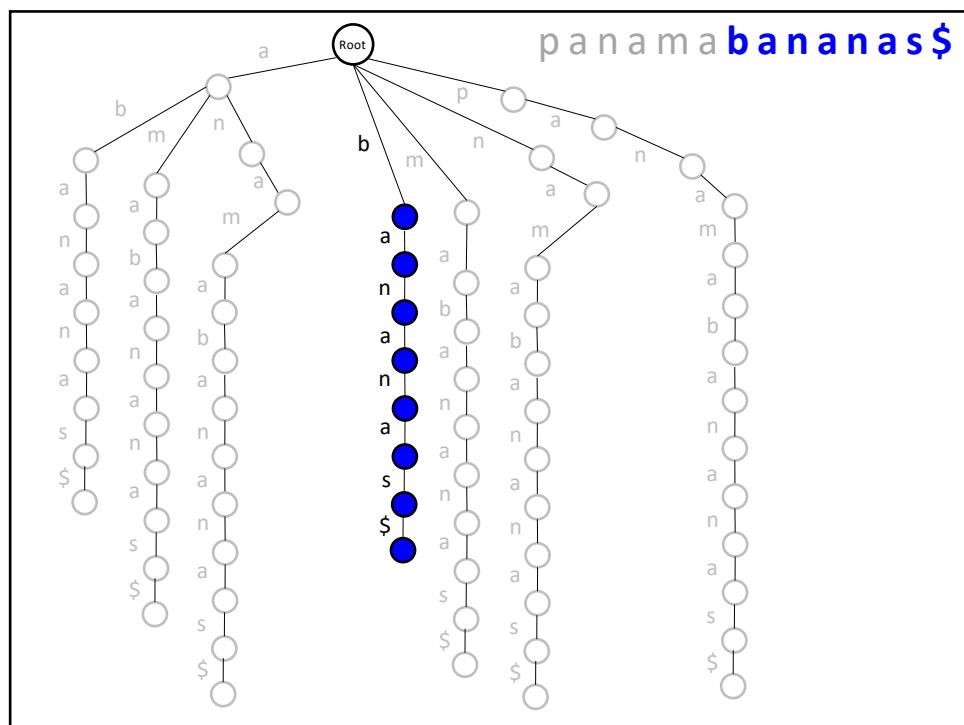
65



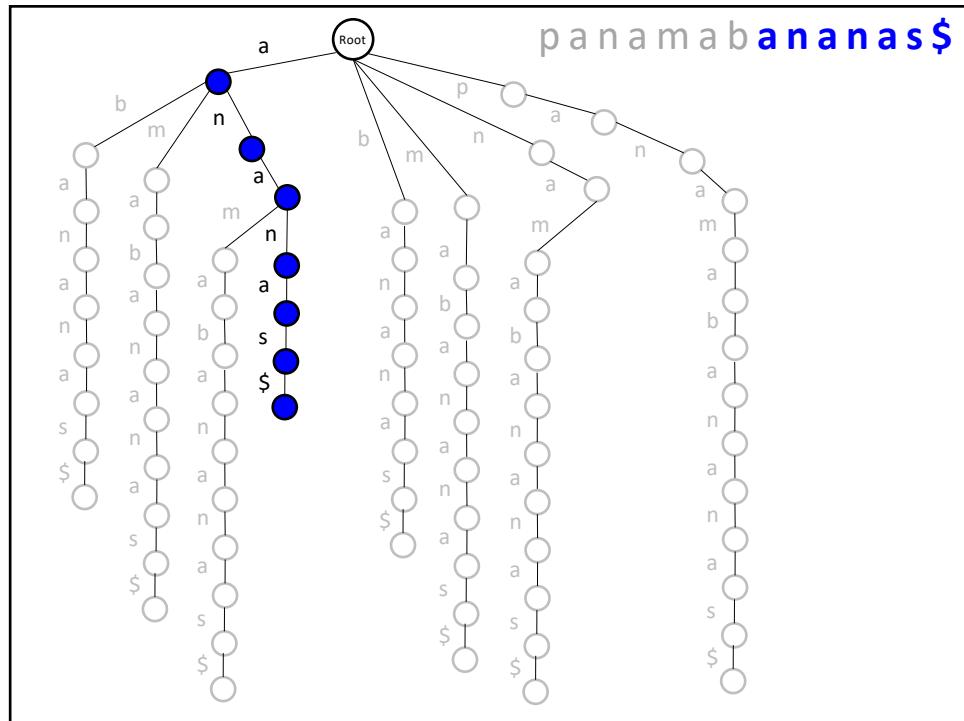
66



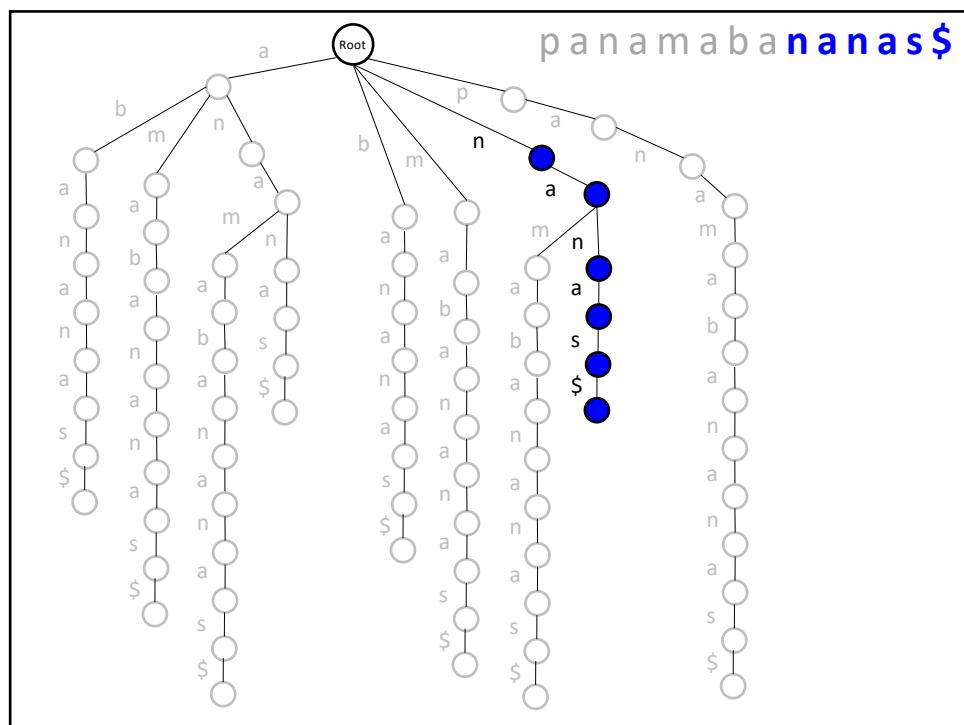
67



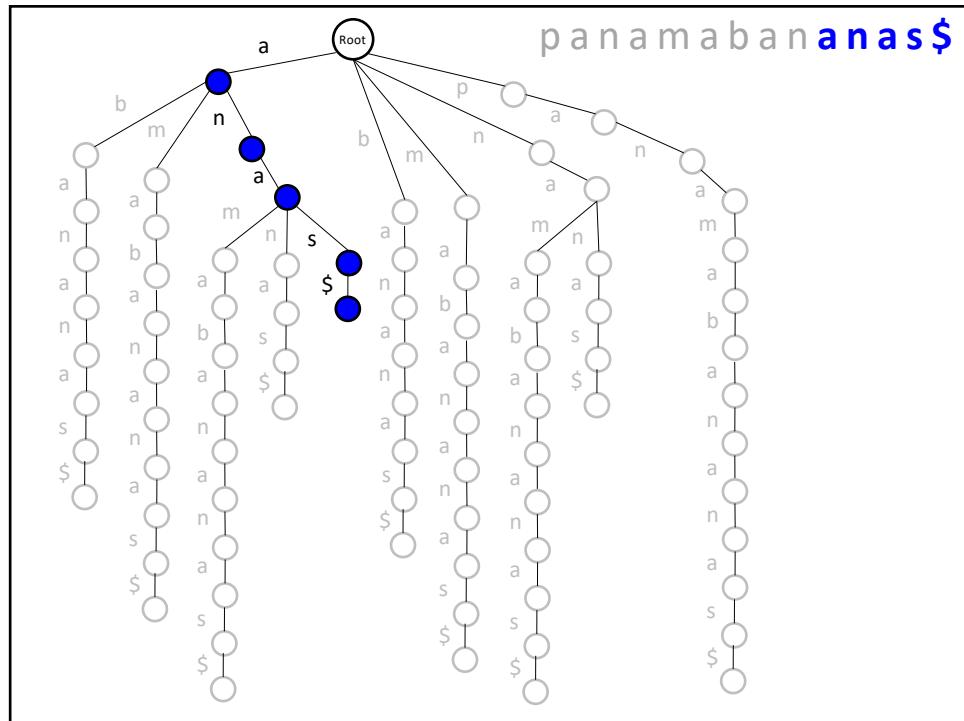
68



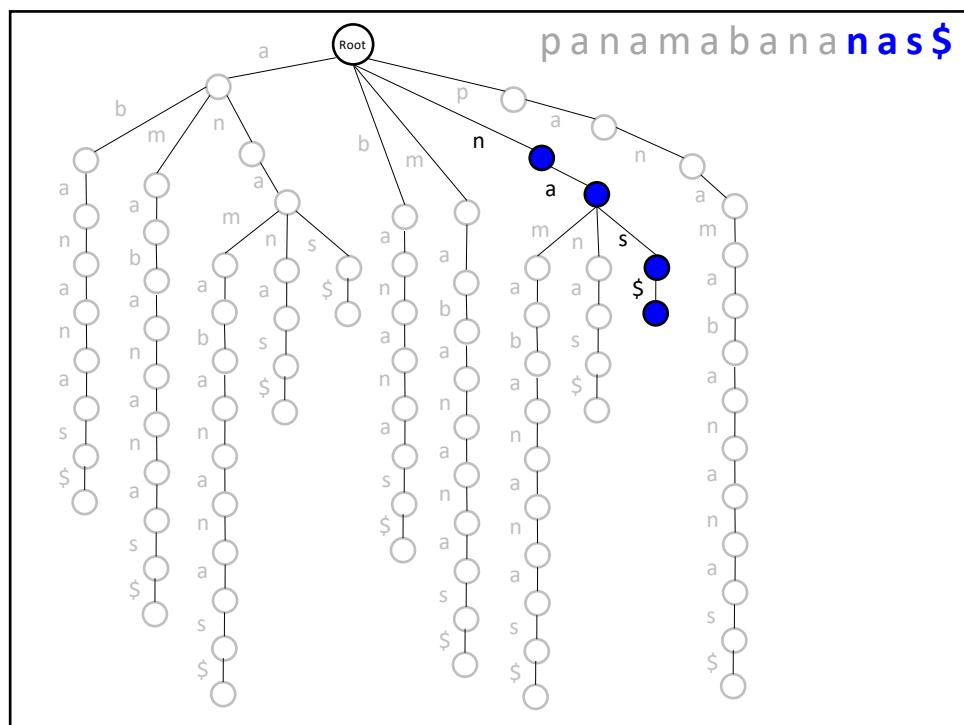
69



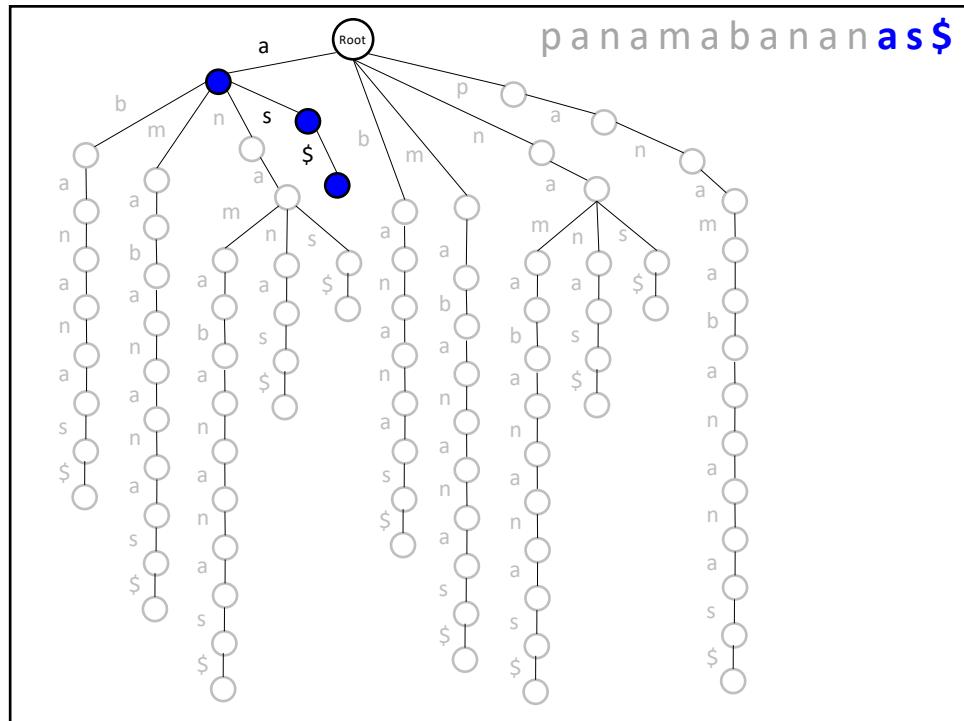
70



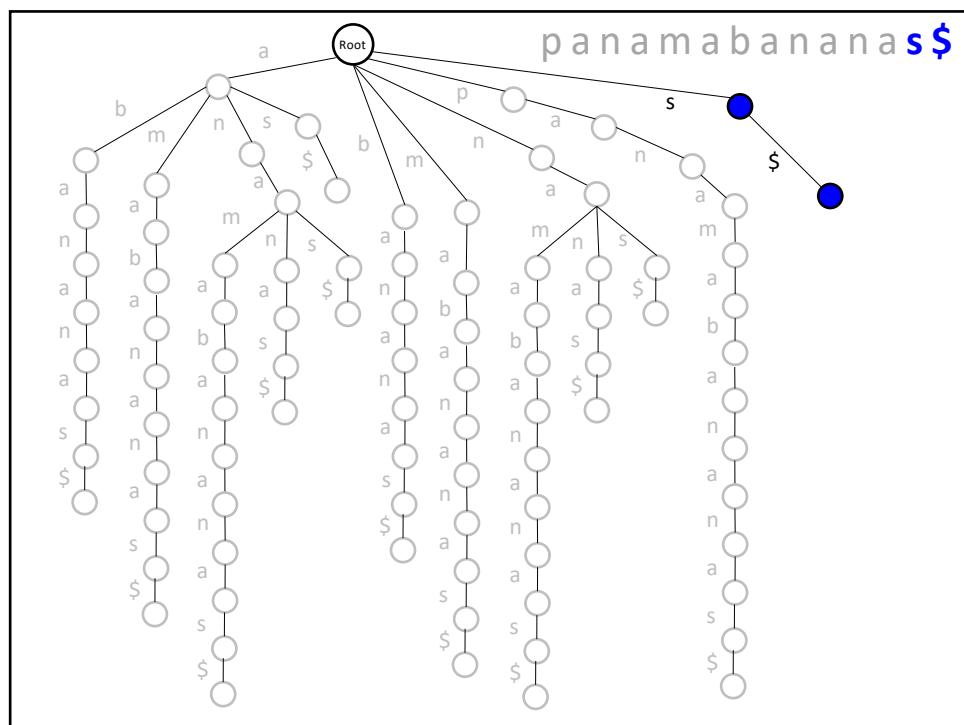
71



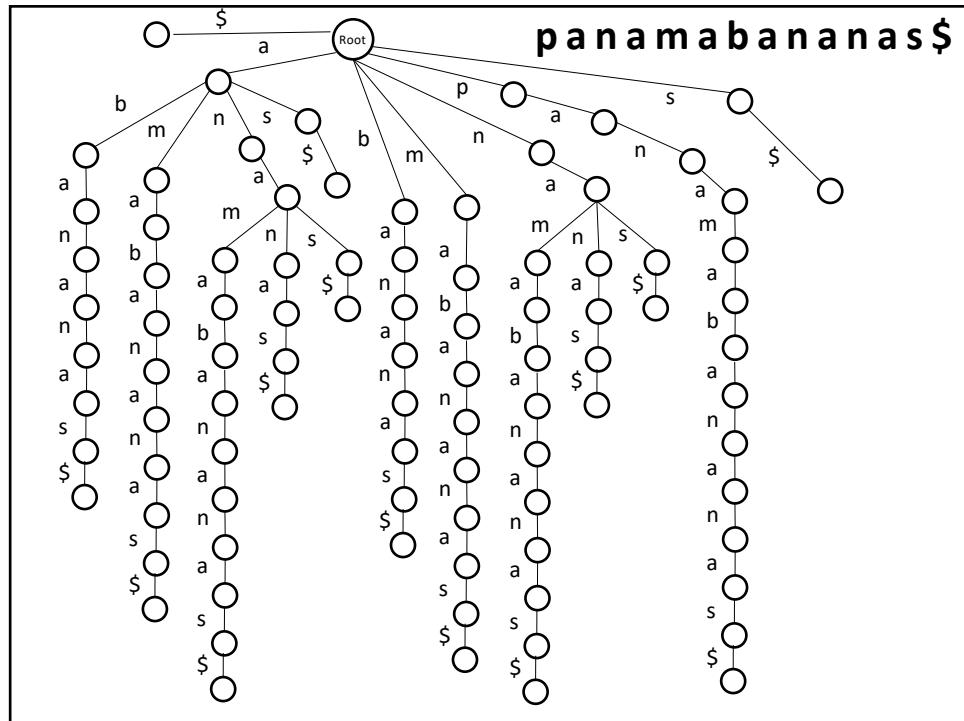
72



73



74



75

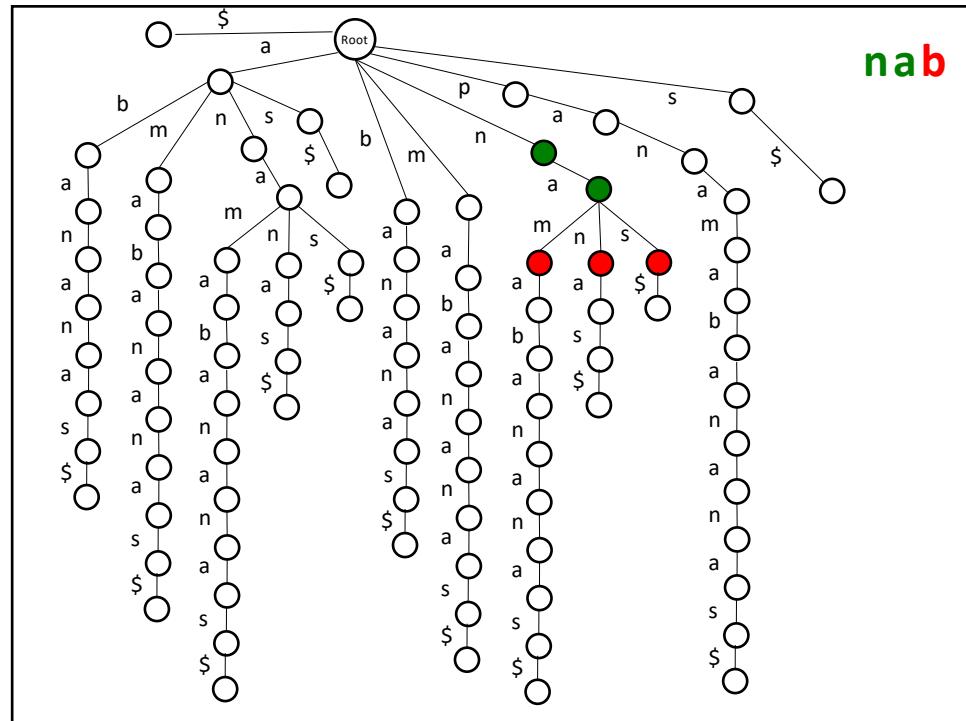
The Suffix Trie and Pattern Matching

- For each *Pattern*, see if *Pattern* can be spelled out from the root downward in the suffix trie.

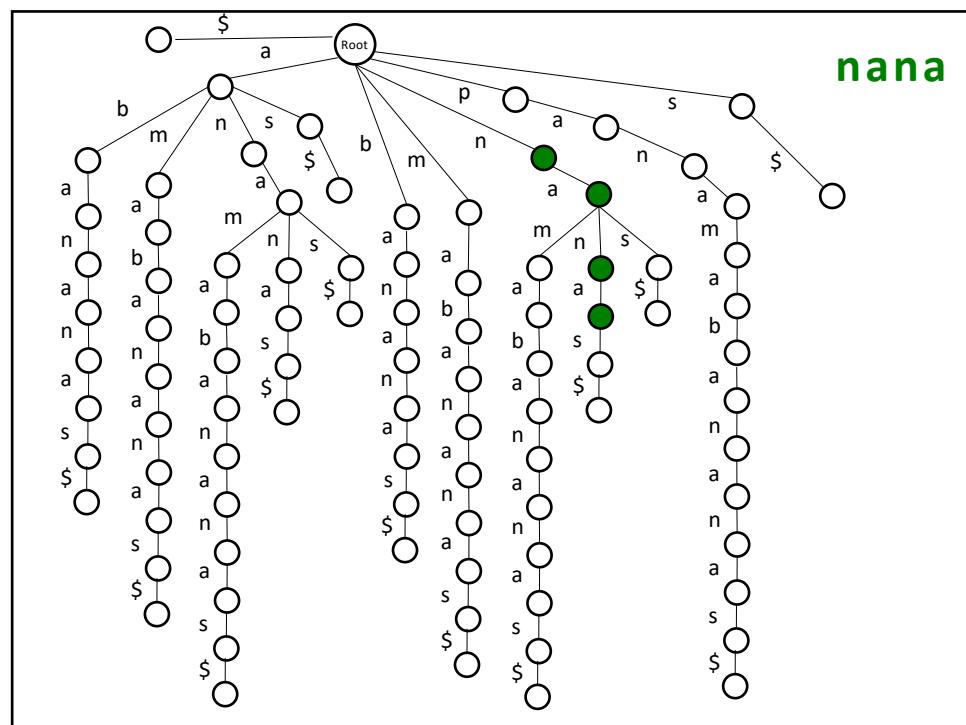


*Suffix
automaton*

76



77

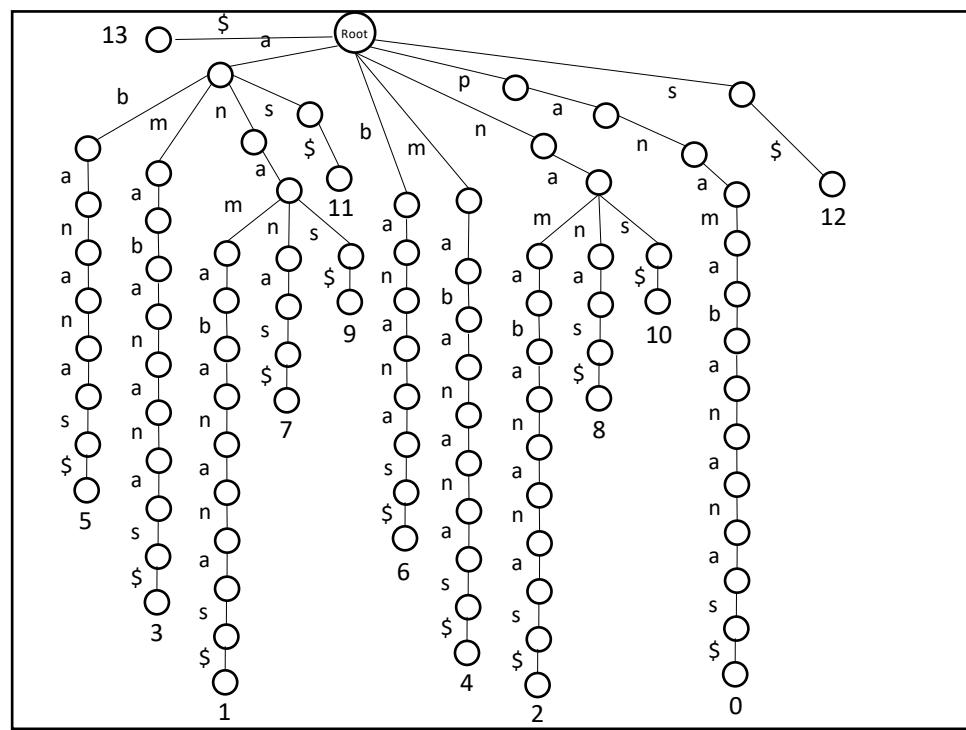


78

Where Are the Matches?

- To find where the pattern matches are, we need to add a little more information to the suffix trie.
- At each leaf (\$), we add the starting position in *Genome* of the suffix ending at that leaf.

79

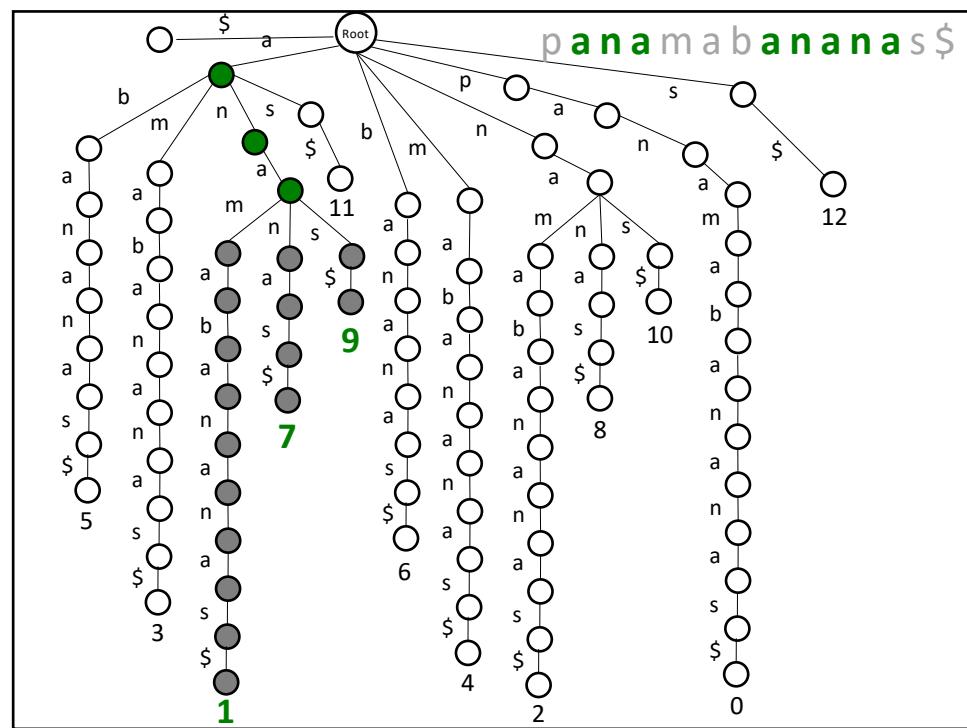


80

Where Are the Matches?

- Once we find a match, we “walk down” to the leaf in order to find the starting position of the match.

81



82

Memory Trouble Once Again

- Worst case: the suffix trie holds $O(|\text{Suffixes}|)$ nodes.
- For a *Genome* of length n ,
 $|\text{Suffixes}| = n(n - 1)/2 = O(n^2)$

Suffixes
 panamabananas\$
 anamabananas\$
 namabananas\$
 amabananas\$
 mabananas\$
 abananas\$
 bananas\$
 ananas\$
 nanas\$
 anas\$
 nas\$
 as\$
 s\$
 \$
 \$

83

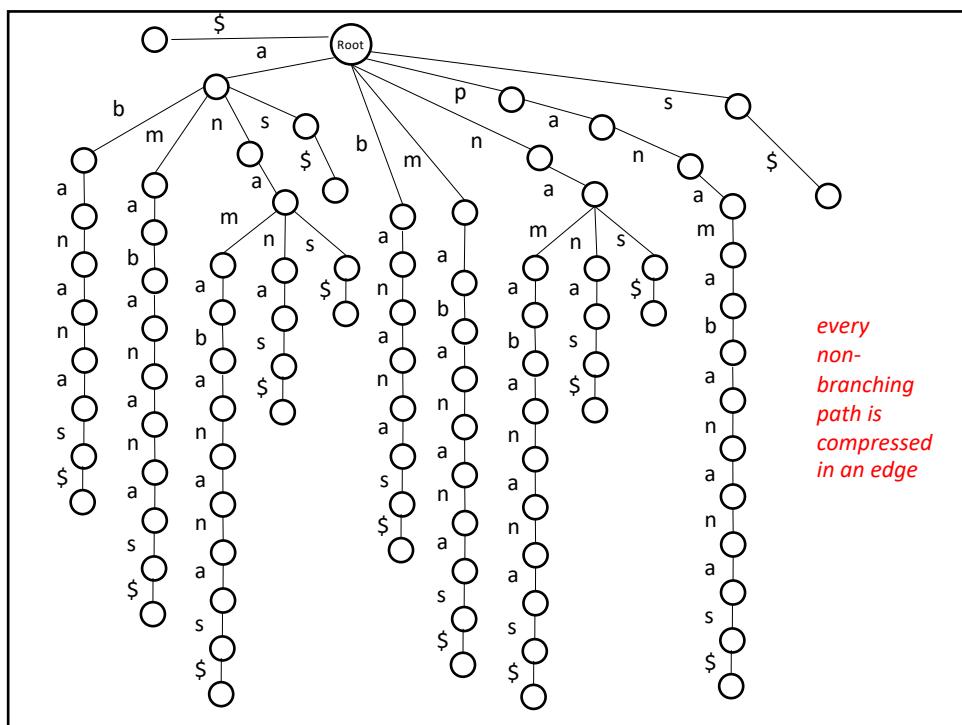
Compressing the Trie

- This doesn't mean that our idea was bad!
- To reduce memory, we can compress each "nonbranching path" of the tree into an edge.

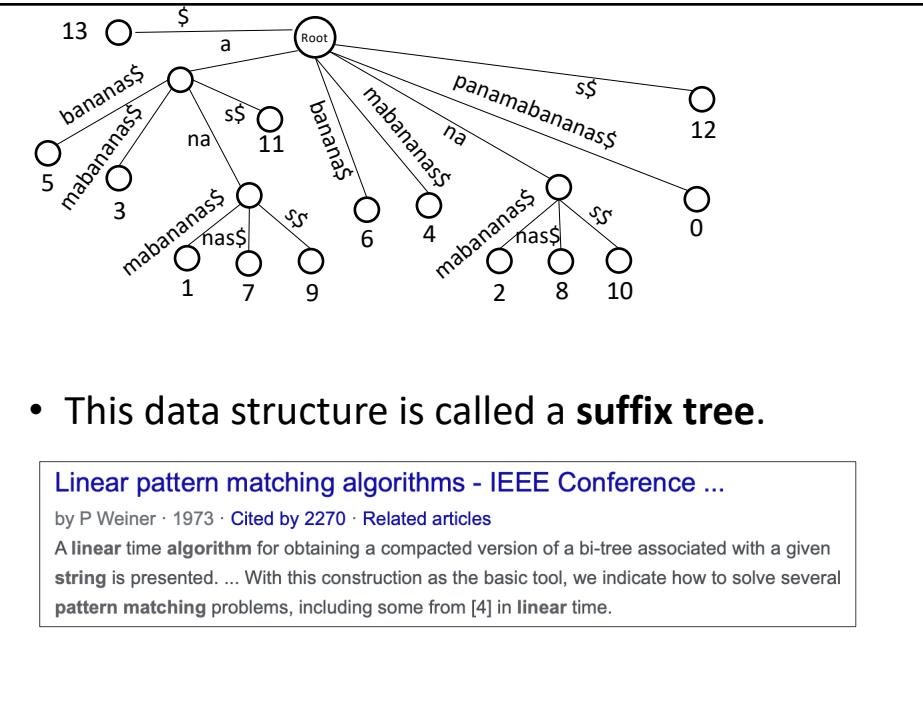
84

SUFFIX TREES

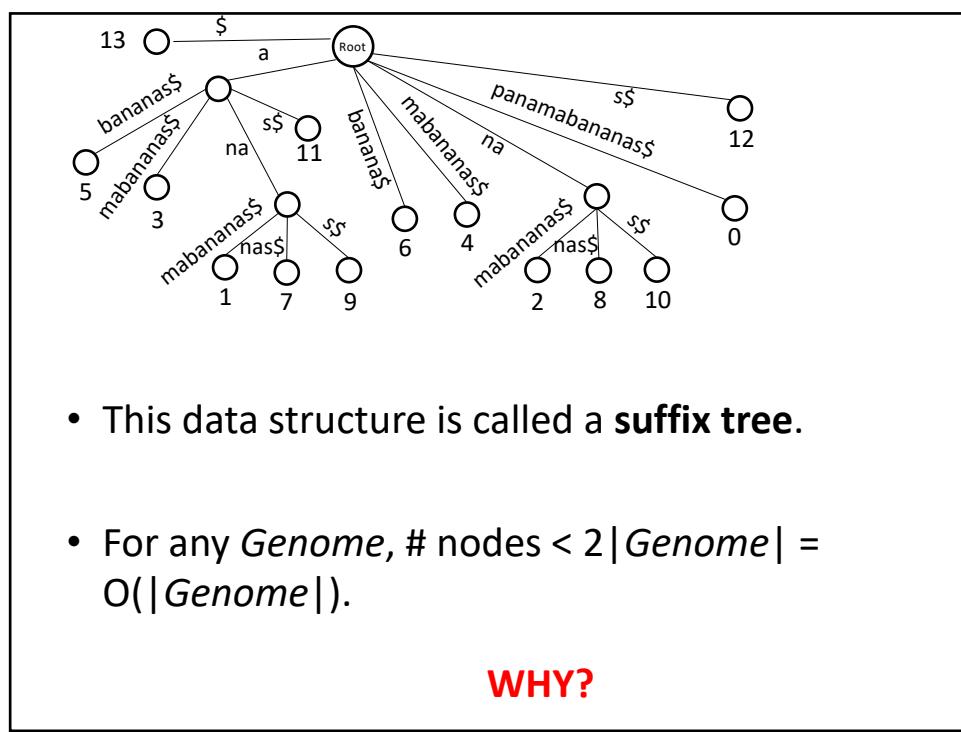
85



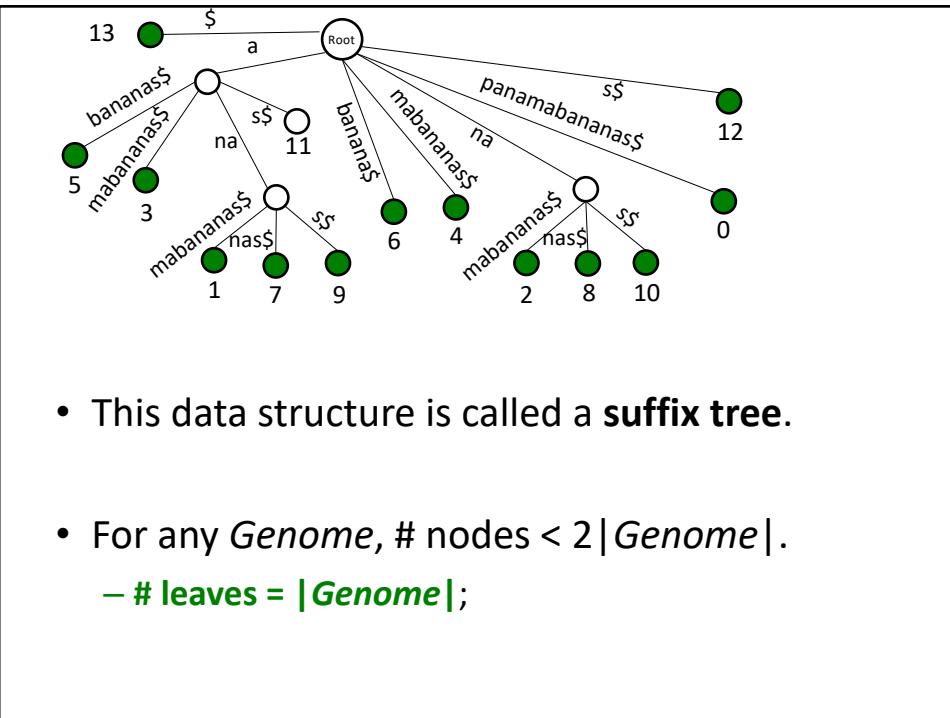
86



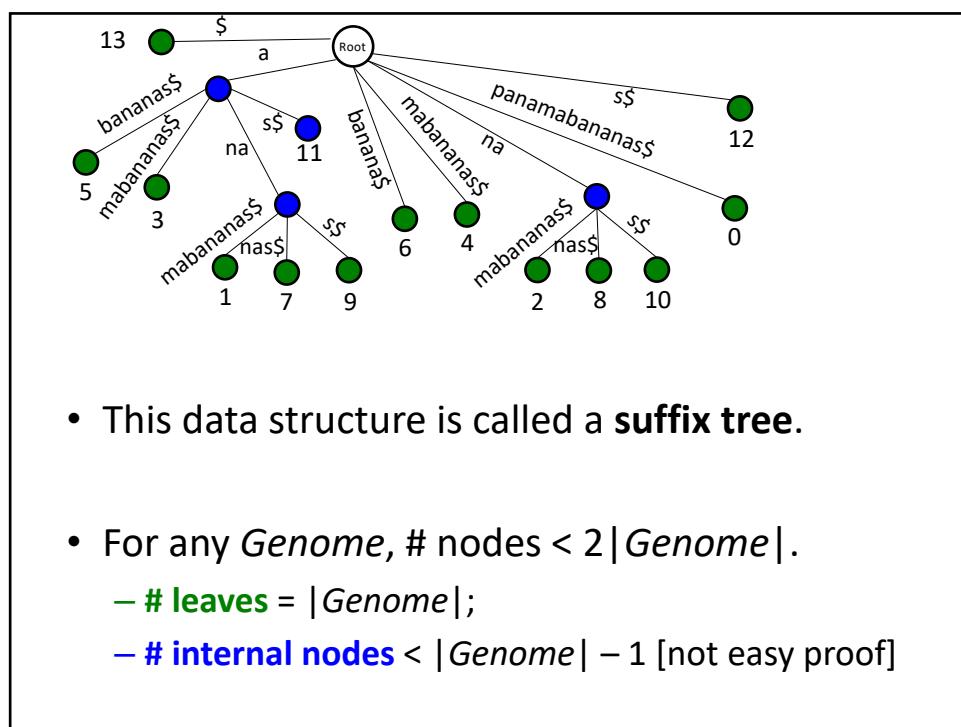
87



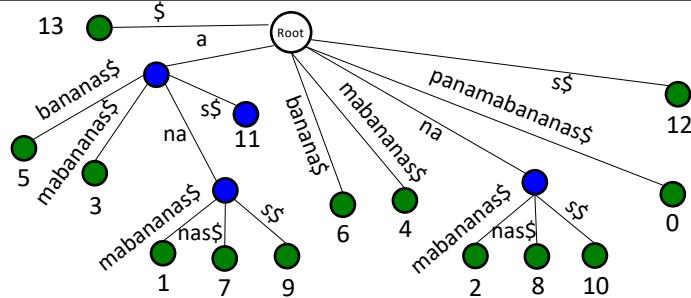
88



89



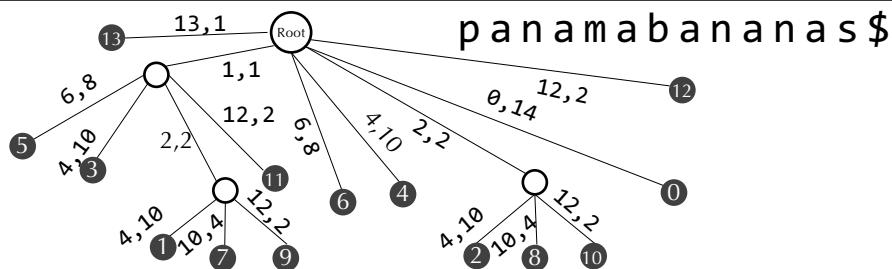
90



STOP: But wait ... we are still storing all the suffixes! So why are suffix trees more memory efficient than suffix tries?

Answer: We can replace each edge with two integers: the starting position of the string, and its length.

91



STOP: But wait ... we are still storing all the suffixes! So why are suffix trees more memory efficient than suffix tries?

Answer: We can replace each edge with two integers: the starting position of the string, and its length.

92

Runtime and Memory Analysis

- Runtime:
 - $O(|Genome|^2)$ to construct the suffix tree.
 - $O(|Genome| + |Patterns|)$ to find pattern matches.
- Memory:
 - $O(|Genome|^2)$ to construct the suffix tree.
 - $O(|Genome|)$ to store the suffix tree.

93

Runtime and Memory Analysis

- Runtime:
 - $O(|Genome|)$ to construct the suffix tree **directly**.
- Memory:
 - $O(|Genome|)$ to construct the suffix tree **directly**.

94

Time complexity

- ✓ Weiner's algorithm (1973)
 - linear time → $O(n)$
- ✓ McCreight's algorithm (1976)
 - linear time → $O(n)$
 - memory efficient
- ✓ Ukkonen's algorithm (1995)
 - linear time → $O(n)$
 - online

95

Runtime and Memory Analysis

- Runtime:
 - $O(|Genome|)$ to construct the suffix tree ***directly***.
 - $O(|Genome| + |Patterns|)$ to find pattern matches.
 - **Total: $O(|Genome| + |Patterns|)$**
- Memory:
 - $O(|Genome|)$ to construct the suffix tree ***directly***.
 - $O(|Genome|)$ to store the suffix tree.
 - **Total: $O(|Genome| + |Patterns|)$**

96

We are Not Finished Yet

- I am happy with the suffix tree, but I am not completely satisfied.
 - Runtime: $O(|Genome| + |Patterns|)$
 - Memory: $O(|Genome|)$
- **However, big-O notation ignores constants!**
 - The best known suffix tree implementations require ~ 20 times the length of $|Genome|$.
 - Can we reduce this constant factor?

97

From Theory to Practice

We are trained to see an algorithm with $O(|Text|)$ runtime and memory and celebrate.

But a data structure with $2*|Text|$ memory is very different from one with $1000*|Text|$ memory.

It takes ~ 60 GB to store a suffix tree for a 3 GB human genome (4 bytes for each edge of tree). Can we do better?

98

SUFFIX ARRAYS

99

Suffix Array of “panamabananas\$”

20 years go by...

Suffix array: Starting positions of sorted suffixes of *Text*.

Sorted Suffixes	suffix array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

100

Suffix Array of “panamabananas\$”

20 years go by...

Suffix array: Starting positions of sorted suffixes of *Text*.

Can be constructed in linear time, and takes only about 12 GB for human genome.

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

101

Suffix Array of “panamabananas\$”



fattorizzazione di
Lyndon per costruire
il suffix array

Can be constructed in linear time, and takes only about 12 GB for human genome.

102

Suffix Array of “panamabananas\$”

STOP: How can we use the suffix array for pattern matching?

...each matched read is a prefix of a suffix of the Genome... i.e., we find it as a prefix in a “segment” of consecutive suffixes ...

SORTED...

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

103

Suffix Array of “panamabananas\$”

STOP: How can we use the suffix array for pattern matching?

Answer: we can pick a value, obtain “higher” or “lower”, and move the next value to the middle of the unsearched suffixes.

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

104

Suffix Array of “panamabananas\$”

This approach, shown for Pattern = “ana”, is the **binary search**.

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

105

Suffix Array of “panamabananas\$”

This approach, shown for Pattern = “ana”, is the **binary search**.

First, we find the *first* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

106

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is the **binary search**.

First, we find the *first* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

107

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is the **binary search**.

First, we find the *first* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

108

Suffix Array of “panamabananas\$”

This approach, shown for Pattern = “ana”, is the **binary search**.

First, we find the *first* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

109

Suffix Array of “panamabananas\$”

This approach, shown for Pattern = “ana”, is the **binary search**.

First, we find the *first* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

110

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is the **binary search**.

First, we find the *first* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

111

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is the **binary search**.

First, we find the *first* occurrence of “ana”.

Next, we find the *last* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

112

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is the **binary search**.

First, we find the *first* occurrence of “ana”.

Next, we find the *last* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
<i>anas\$</i>	<i>9</i>
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

113

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is the **binary search**.

First, we find the *first* occurrence of “ana”.

Next, we find the *last* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
<i>as\$</i>	<i>11</i>
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

114

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is the **binary search**.

First, we find the *first* occurrence of “ana”.

Next, we find the *last* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

115

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is the **binary search**.

STOP: How do we know the starting locations of these three occurrences?

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
ana mabananas\$	1
ana nas\$	7
ana s\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

116

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is the **binary search**.

Answer: The suffix array tells us!

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
ana mabananas\$	1
ana nas\$	7
ana s\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

117

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is the **binary search**.

STOP: What is the runtime of the binary search approach using the suffix array?

Sorted Suffixes	Suffix Array
\$	13
abanananas\$	5
amabananas\$	3
ana mabananas\$	1
ana nas\$	7
ana s\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

118

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is the **binary search**.

Answer: Each pattern takes $O(\log(|Text|))$.

Sorted Suffixes	suffix array
\$	13
abanananas\$	5
amabananas\$	3
ana mabananas\$	1
ana nas\$	7
ana s\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

119

Binary Search Pseudocode

```

PatternMatchingSuffixArray(Text, Pattern)
    s ← SuffixArray(Text)
    minIndex ← 0
    maxIndex ← |Text| - 1
    while minIndex <= maxIndex
        midIndex ← floor((minIndex + maxIndex)/2)
        if Pattern > suffix of Text starting at s(midIndex)
            minIndex ← midIndex + 1
        else
            maxIndex ← midIndex - 1
        if Pattern matches suffix of Text starting at s(midIndex)
            first ← minIndex
        else
            return no occurrences
        minIndex ← first
        maxIndex ← |Text| - 1
        while minIndex <= maxIndex
            midIndex ← floor((minIndex + maxIndex)/2)
            if Pattern matches suffix of Text starting at s(midIndex)
                minIndex ← midIndex + 1
            else
                maxIndex ← midIndex - 1
        last ← maxIndex
        return (first, last)
    
```

120

Suffix Arrays are Another Example of Computational Biologists Influencing CS

Suffix arrays: a new method for on-line string searches

U Manber, G Myers - siam Journal on Computing, 1993 - SIAM

A new and conceptually simple data structure, called a suffix array, for on-line string searches is introduced in this paper. Constructing and querying suffix arrays is reduced to a sort and search paradigm that employs novel algorithms. The main advantage of suffix arrays over suffix trees is that, in practice, they use three to five times less space. From a complexity standpoint, suffix arrays permit on-line string searches of the type, "Is W a substring of A?" to be answered in time $O(P+\log N)$, where P is the length of W and N is the ...

[☆](#) [99](#) Cited by 2870 Related articles All 43 versions [»»](#)

Note: This is a paper produced by a computational biologist (and future head of search for Google) in a non-biological context seven years before the publication of the first draft human genome.

121

Suffix Arrays are Another Example of Computational Biologists Influencing CS

Suffix arrays: a new method for on-line string searches

U Manber, G Myers - siam Journal on Computing, 1993 - SIAM

A new and conceptually simple data structure, called a suffix array, for on-line string searches is introduced in this paper. Constructing and querying suffix arrays is reduced to a sort and search paradigm that employs novel algorithms. The main advantage of suffix arrays over suffix trees is that, in practice, they use three to five times less space. From a complexity standpoint, suffix arrays permit on-line string searches of the type, "Is W a substring of A?" to be answered in time $O(P+\log N)$, where P is the length of W and N is the ...

[☆](#) [99](#)

Kärkkäinen & Sanders (2003)

Suffix Array del testo T in tempo $O(n)$

Note: This is a paper produced by a computational biologist (and future head of search for Google) in a non-biological context seven years before the publication of the first draft human genome.

122

FATTORIZZAZIONI DI LYNDON E SA...

123

Outline

- Why do we map reads?
- Using the Trie
- From a Trie to a Suffix Tree
- **String Compression and the Burrows-Wheeler Transform**
- Inverting Burrows-Wheeler
- Using Burrows-Wheeler for Pattern Matching
- Finding the Matched Patterns
- Setting Up Checkpoints
- Inexact Matching
- Further Applications of Read Mapping

124

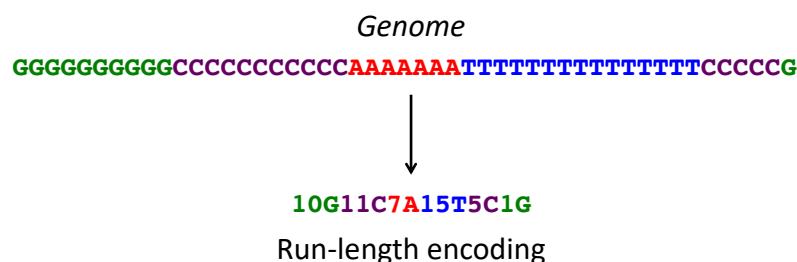
O(|G|) ... Genome Compression

- Idea: decrease the amount of memory required to hold *Genome*.
- This indicates that we need methods of **compressing** a large genome, which is seemingly a separate problem.

125

Idea #1: Run-Length Encoding

- **Run-length encoding:** compresses a run of n identical symbols.

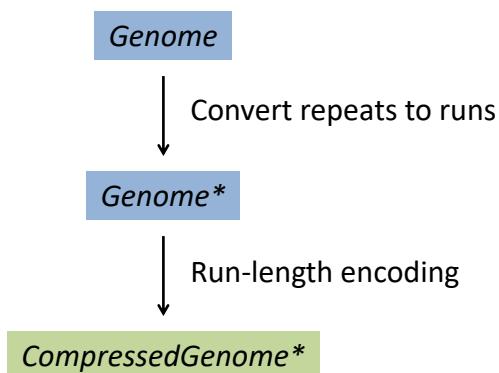


- Problem: Genomes don't have lots of runs...

126

Converting Repeats to Runs

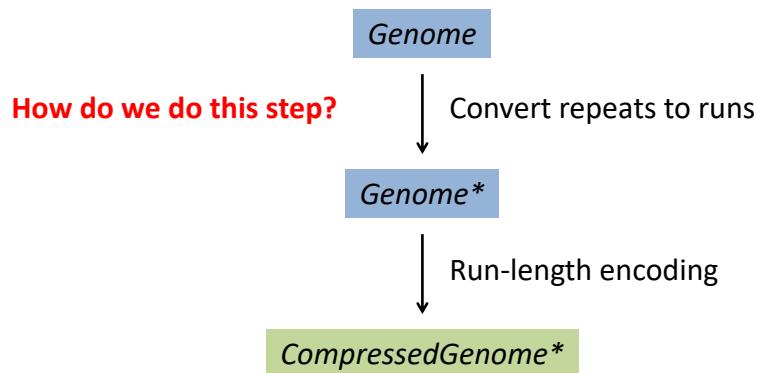
- ...but they do have lots of repeats!



127

Converting Repeats to Runs

- ...but they do have lots of repeats!



128

The Burrows-Wheeler Transform

spiegata nella prima parte del corso,
dove sono precisamente riportate:

- definizioni
 - algoritmi
 - correttezza degli algoritmi

ne vedremo solo l'uso per
il mapping

129

The Burrows-Wheeler Transform

The diagram illustrates the string matching step for the input string "panamabananas\$" against the pattern "\$panana". A red arrow points from the character 'n' in "panamabananas" to the character 'p' in "\$panana".

Below, a state transition diagram shows a circle with an incoming arrow labeled 'a' and an outgoing arrow labeled 'a', with states 'n', 'a', 'b', and 'm' around it.

Form all cyclic rotations of
“panamabananas\$”

130

The Burrows-Wheeler Transform

```

panamabananas$  

$panamabananas  

s$panamabana  

as$panamaban  

nas$panamaba  

anas$panamab  

nanas$panamab  

ananas$panamab  

bananas$panama  

abananas$panam  

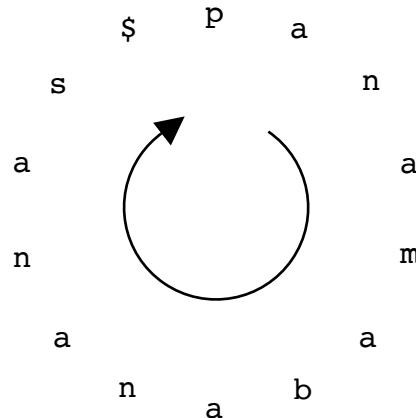
mabananas$pana  

amabananas$pan  

namabananas$pa  

anamabananas$p

```



Form all cyclic rotations of
“panamabananas\$”

131

The Burrows-Wheeler Transform

```

panamabananas$  

$panamabananas  

s$panamabana  

as$panamaban  

nas$panamaba  

anas$panamab  

nanas$panamab  

ananas$panamab  

bananas$panama  

abananas$panam  

mabananas$pana  

amabananas$pan  

namabananas$pa  

anamabananas$p

```



```

$panamabananas  

abanas$panam  

amabananas$pan  

anamabananas$p  

ananas$panamab  

anas$panamab  

as$panamab  

bananas$panama  

mabananas$pana  

namabananas$pa  

nas$panamab  

panamabananas$  

s$panamabana

```

Form all cyclic rotations of
“panamabananas\$”

Sort the strings lexicographically
(\$ comes first)

the color shows that the strings are sorted

132

The Burrows-Wheeler Transform

panamabananas\$
 \$panamabananas
 s\$panamabana
 as\$panamaban
 nas\$panamaba
 anas\$panamab
 nanas\$panamab
 ananas\$panamab
 bananas\$panama
 abananas\$panam
 mabananas\$pana
 amabananas\$pan
 namabananas\$pa
 anamabananas\$p

\$panamabananas
 abananas\$panam
 amabananas\$pan
 anamabananas\$p
 ananas\$panama**b**
 anas\$panamaba**n**
 as\$panamabana**n**
 bananas\$panama**a**
 mabananas\$pana**a**
 namabananas\$p**a**
 nanas\$panamab**a**
 nas\$panamaba**a**
 panamabananas\$
 s\$panamabana**a**

Form all cyclic rotations of
 "panamabananas\$"

Burrows-Wheeler Transform:
 Last column = smnpbnnaaaa\$a

133

the LAST column... Why Should We Care?

Watson-Crick paper on double structure of DNA

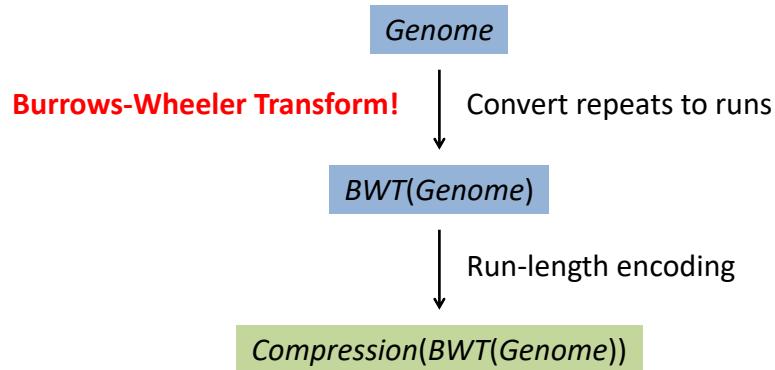
nd Corey (1). They kindly made their manuscript availa a
 nd criticism, especially on interatomic distances. We a
 nd cytosine. The sequence of bases on a single chain d a
 nd experimentally (3,4) that the ratio of the amounts o u
 nd for this reason we shall not comment on it. We wish a
 nd guanine (purine) with cytosine (pyrimidine). In oth a
 nd ideas of Dr. M. H. F. Wilkins, Dr. R. E. Franklin a
 nd its water content is rather high. At lower water co a
 nd pyrimidine bases. The planes of the bases are perpe a
 nd stereochemical arguments. It has not escaped our no a
 nd that only specific pairs of bases can bond together u
 nd the atoms near it is close to Furberg's 'standard co a
 nd the bases on the inside, linked together by hydrogen a
 nd the bases on the outside. In our opinion, this stru a
 nd the other a pyrimidine for bonding to occur. The hy a
 nd the phosphates on the outside. The configuration of a
 nd the ration of guanine to cytosine, are always very c a
 nd the same axis (see diagram). We have made the usual u
 nd their co-workers at King's College, London. One of a

qui no...

the text has a lot of repetitive words... (the word «and»)

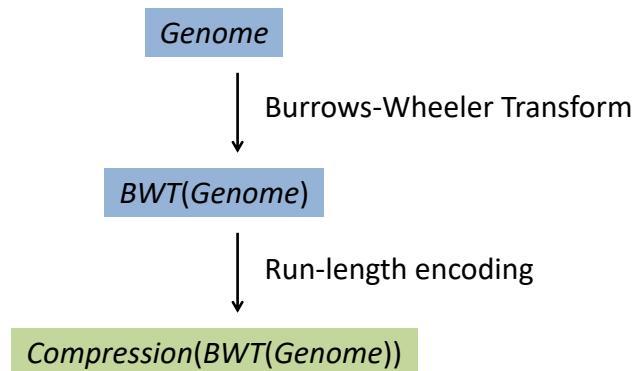
134

BWT: Converting Repeats to Runs



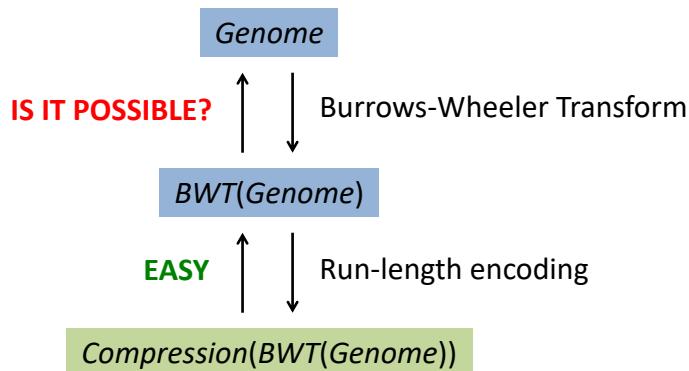
135

How Can We Decompress?



136

How Can We Decompress?



137

Outline

- Why do we map reads?
- Using the Trie
- From a Trie to a Suffix Tree
- String Compression and the Burrows-Wheeler Transform
- **Inverting Burrows-Wheeler**
- Using Burrows-Wheeler for Pattern Matching
- Finding the Matched Patterns
- Setting Up Checkpoints
- Inexact Matching
- Further Applications of Read Mapping

138

Reconstructing banana

```
$banana
a$banan
ana$ban
anana$b
banana$
na$banan
nana$b
```

We have the last column, we obtain the string

139

Reconstructing banana

```
$banana
a$banan
ana$ban
anana$b
banana$
na$banan
nana$b
```

- Sorting all elements of “annb\$aa” gives first column of matrix [since the strings are lexicographically sorted]

140

Reconstructing banana

\$banana	a\$
a\$banan	na
ana\$ban	na
anana\$b	ba
banana\$	2-mers
na\$banan	\$b
nana\$b	an

- We now know 2-mer composition of the circular string banana\$

141

Reconstructing banana

\$banana	a\$	\$b
a\$banan	na	a\$
ana\$ban	na	an
anana\$b	ba	an
banana\$	2-mers	Sort
na\$banan	\$b	ba
nana\$b	an	na

- We now know 2-mer composition of the circular string banana\$
- Sorting gives us the first 2 columns of the matrix.

142

Reconstructing banana

\$banana	a\$	\$b
a\$banan	na	a\$
ana\$ban	na	an
anana\$b	ba	an
banana\$	→ 2-mers \$b	Sort ba
na\$ban	an	na
nana\$ba	an	na

- We now know 2-mer composition of the circular string banana\$
- Sorting gives us the first 2 columns of the matrix.

143

Reconstructing banana

```
$banana
a$banan
ana$ban
anana$b
banana$
na$ban
nana$ba
```

We use the first two columns + the last column,
and we obtain the 3-mer composition of the circular
string «banana»

144

Reconstructing banana

\$banana	a\$b
a\$banan	na\$
ana\$ban	nan
anana\$b	ban
banana\$	3-mers
na\$bana	\$ba
nana\$ba	ana

- We now know 3-mer composition of the circular string banana\$

145

Reconstructing banana

\$banana	a\$b	\$ba
a\$banan	na\$	a\$b
ana\$ban	nan	ana
anana\$b	ban	ana
banana\$	3-mers	Sort
na\$bana	\$ba	ban
nana\$ba	ana	na\$

- We now know 3-mer composition of the circular string banana\$
- Sorting gives us the first 3 columns of the matrix.

146

Reconstructing banana

\$banana	a\$b	\$ba
a\$banan	na\$	a\$b
ana\$ban	nan	ana
anana\$b	ban	ana
banana\$	\$ba	ban
na\$bana	ana	na\$
nana\$ba	ana	nan

→ 3-mers → Sort

- We now know 3-mer composition of the circular string banana\$
- Sorting gives us the first 3 columns of the matrix.

147

Reconstructing banana

```
$banana
a$banan
ana$ban
anana$b
banana$
na$bana
nana$ba
```

148

Reconstructing banana

\$banana	a\$ba
a\$banan	na\$b
ana\$ban	nana
anana\$b	→ bana
banana\$	4-mers \$ban
na\$bana	ana\$
nana\$ba	anana

- We now know 4-mer composition of the circular string banana\$

149

Reconstructing banana

\$banana	a\$ba	\$ban
a\$banan	na\$b	a\$bb
ana\$ban	nana	anaa
anana\$b	→ bana	→ anaa
banana\$	4-mers \$ban	Sort bann
na\$bana	ana\$	na\$b
nana\$ba	anana	nana

- We now know 4-mer composition of the circular string banana\$
- Sorting gives us the first 4 columns of the matrix.

150

Reconstructing banana

\$banana	a\$ba	\$ban
a\$banan	na\$b	a\$bb
ana\$ban	nana	anaa
anana\$b	bana	anaa
banana\$	\$ban	bann
na\$bana	ana\$	na\$b
nana\$ba	an	nana

4-mers Sort

- We now know 4-mer composition of the circular string banana\$
- Sorting gives us the first 4 columns of the matrix.

151

Reconstructing banana

```
$banana
a$banan
ana$ban
anana$b
banana$
na$bana
nana$ba
```

152

Reconstructing banana

\$banana	a\$ban
a\$banan	na\$ba
ana\$ban	nana\$
anana\$b	banan
banana\$	5-mers \$bana
na\$bana	ana\$b
nana\$ba	anana

- We now know 5-mer composition of the circular string banana\$

153

Reconstructing banana

\$banana	a\$ban	\$bana
a\$banan	na\$ba	a\$bbn
ana\$ban	nana\$	anaab
anana\$b	banan	—> anaaa
banana\$	5-mers \$bana	Sort bannn
na\$bana	ana\$b	na\$ba
nana\$ba	anana	nana\$

- We now know 5-mer composition of the circular string banana\$
- Sorting gives us the first 5 columns of the matrix.

154

Reconstructing banana

\$banana	a\$ban	\$bana
a\$banan	na\$ba	a\$bbn
ana\$b a	nana\$	anaab
anana\$b	banan	aaaaa
banana\$	\$bana	bannn
na\$bana	ana\$b	na\$ba
nana\$ba	anana	nana\$

5-mers Sort

- We now know 5-mer composition of the circular string banana\$
- Sorting gives us the first 5 columns of the matrix.

155

Reconstructing banana

```
$banana
a$banan
ana$ba
anana$b
banana$
na$bana
nana$ba
```

156

Reconstructing banana

\$banana	a\$bana
a\$banan	na\$ban
ana\$ban	nana\$b
anana\$b	banana
banana\$	6-mers
na\$bana	\$banan
nana\$ba	ana\$ba

→

anana\$

- We now know 6-mer composition of the circular string banana\$

157

Reconstructing banana

\$banana	a\$bana	\$banan
a\$banan	na\$ban	a\$bbna
ana\$ban	nana\$b	anaaba
anana\$b	banana	anaaa\$
banana\$	6-mers	Sort
na\$bana	\$banan	bannna
nana\$ba	ana\$ba	na\$ban

→

anana\$

- We now know 6-mer composition of the circular string banana\$
- Sorting gives us the first 6 columns of the matrix.

158

Reconstructing banana

\$banana	a\$bana	\$banan
a\$banan	na\$ban	a\$bbna
ana\$ban	nana\$b	anaaba
anana\$b	banana	anaaa\$
banana\$	→ 6-mers	Sort →
na\$bana	\$banan	bannna
nana\$ba	ana\$ba	na\$ban

- We now know 6-mer composition of the circular string banana\$
- Sorting gives us the first 6 columns of the matrix.

159

Reconstructing banana

\$banana
 a\$banan
 ana\$ban
 anana\$b
 banana\$
 na\$bana
 nana\$ba

- We now know the entire matrix!

160

Reconstructing banana

```
$banana
a$banan
ana$ban
anana$b
banana$
na$bana
nana$ba
```

- We now know the entire matrix!
- Taking all elements in the first row (after \$) produces banana.

161

More Memory Issues

- Reconstructing *Genome* from $BWT(\text{Genome})$ required us to store $|\text{Genome}|$ copies of $|\text{Genome}|$.

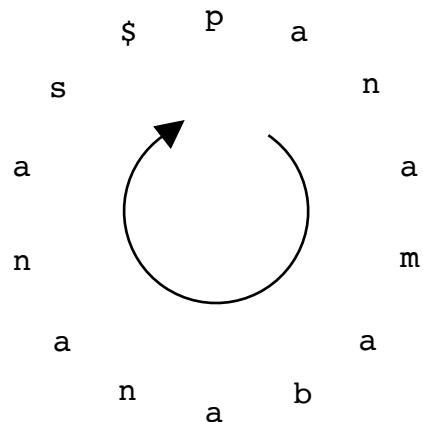
```
$banana
a$banan
ana$ban
anana$b
banana$
na$bana
nana$ba
```

- **Can we invert BWT with less space?**

162

A “Strange” Observation

```
$panamabananas
abanas$panam
amabananas$pan
anamabananas$p
ananas$panamab
anas$panamaban
as$panamaban
bananas$panama
mabananas$pana
namabananas$p
anas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```

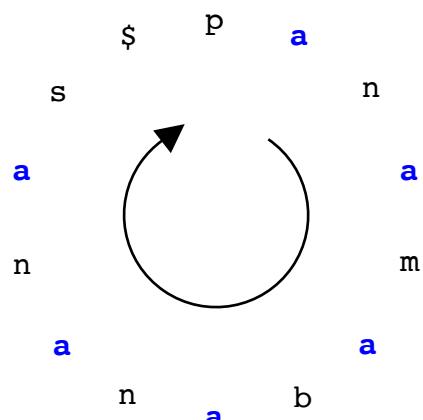


We have the last column, and the first column for free.

163

A Strange Observation

```
$panamabananas
abanas$panam
amabananas$pan
anamabananas$p
ananas$panamab
anas$panamaban
as$panamaban
bananas$panama
mabananas$pana
namabananas$p
anas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```

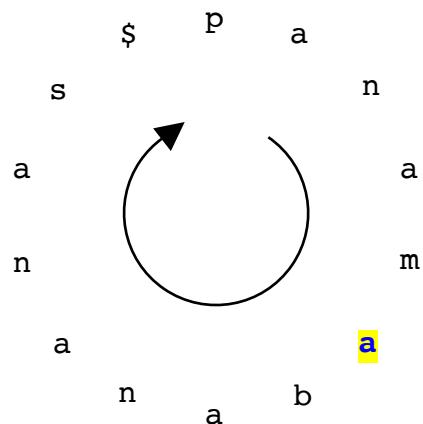


See the six occurrence of a in the first column, in the last column and in the circle representation

164

A Strange Observation

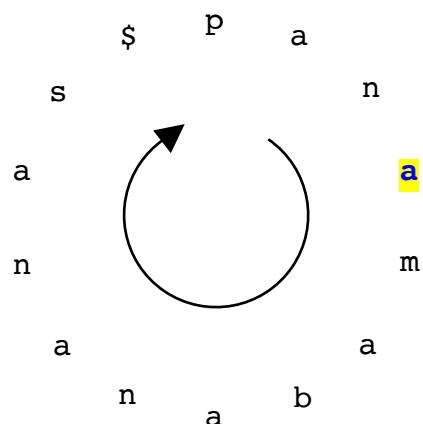
```
$panamabananas
abanas$panam
amabananas$pan
anamabananas$p
ananas$panamab
anas$panamaban
as$panamabanan
bananas$panama
mabananas$pana
namabananas$p
anas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```



165

A Strange Observation

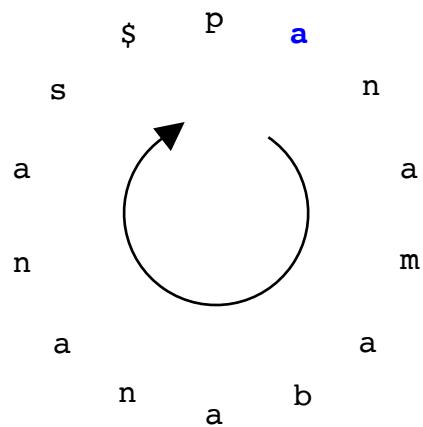
```
$panamabananas
abanas$panam
amabananas$pan
anamabananas$p
ananas$panamab
anas$panamaban
as$panamabanan
bananas$panama
mabananas$pana
namabananas$p
anas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```



166

A Strange Observation

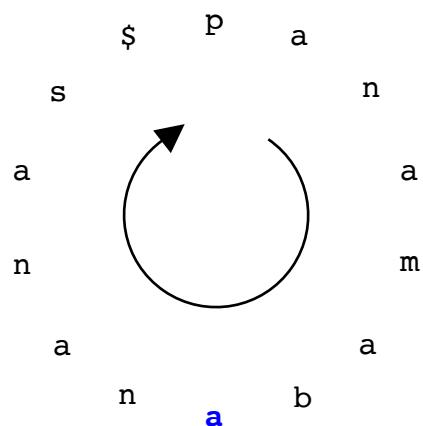
```
$panamabananas
abanas$panam
amabananas$pan
anamabananas$p
ananas$panamab
anas$panamaban
as$panamabanan
bananas$panama
mabananas$pana
namabananas$p
nanas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```



167

A Strange Observation

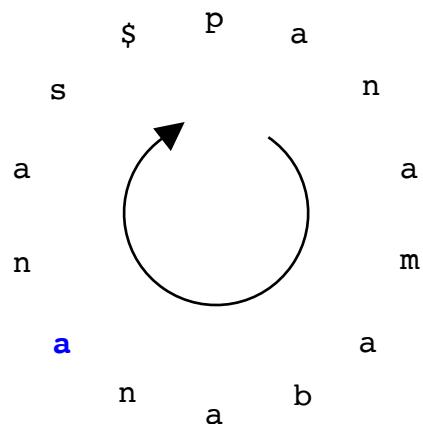
```
$panamabananas
abanas$panam
amabananas$pan
anamabananas$p
ananas$panamab
anas$panamaban
as$panamabanan
bananas$panama
mabananas$pana
namabananas$p
nanas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```



168

A Strange Observation

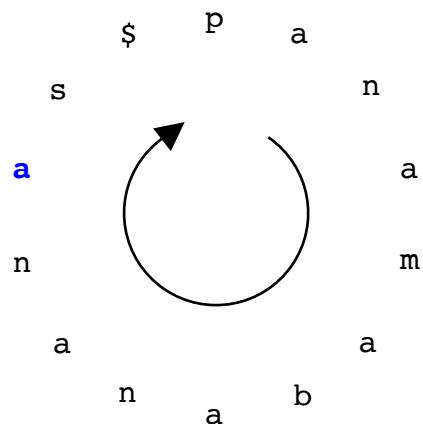
```
$panamabananas
abanas$panam
amabananas$pan
anamabananas$p
ananas$panamab
nas$panamaban
as$panamaban
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```



169

A Strange Observation

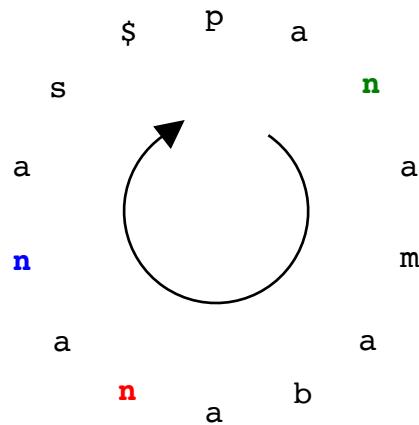
```
$panamabananas
abanas$panam
amabananas$pan
anamabananas$p
ananas$panamab
nas$panamaban
as$panamaban
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```



170

A Strange Observation

```
$panamabananas
abanas$panam
amabananas$pan
anamabananas$p
ananas$panamab
anas$panamaba
as$panamaban
bananas$panama
mabananas$pana
namabananas$pa
anas$panamaba
nas$panamabana
panamabananas$
spanamabanana
```



171

Is It True in General?

We want to find a faster way of inverting BWT

172

Is It True in General?

```
$panamabananas
1 abananas$panam
2 amabananas$pan
3 anamabananas$p
4 ananas$panamab
5 anas$panamaban
6 as$panamaban
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```

These strings are sorted

173

Is It True in General?

```
$panamabananas
1 abananas$panam
2 amabananas$pan
3 anamabananas$p
4 ananas$panamab
5 anas$panamaban
6 as$panamaban
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```

bananas\$panam
mabananas\$pan
namabananas\$p
nanas\$panamab
nas\$panamaban
s\$panamabanana

Chop off a →

These strings are sorted

174

Is It True in General?

```
$panamabananas
1 abananas$panam
2 amabananas$pan
3 anamabananas$p
4 ananas$panamab
5 anas$panamaban
6 as$panamabanan
```

bananas\$panam
mabananas\$pan
namabananas\$p
nanas\$panamab
nas\$panamaban
s\$panamabanan

Chop off a

Still sorted

These strings are sorted

175

Is It True in General?

```
$panamabananas
1 abananas$panam
2 amabananas$pan
3 anamabananas$p
4 ananas$panamab
5 anas$panamaban
6 as$panamabanan
```

bananas\$panam
mabananas\$pan
namabananas\$p
nanas\$panamab
nas\$panamaban
s\$panamabanan

Chop off a

Still sorted

Add a
to end

bananas\$panam a
mabananas\$pan a
namabananas\$p a
nanas\$panamab a
nas\$panamaban a
s\$panamabanana a

These strings are sorted

176

Is It True in General?

```
$panamabananas
1 abananas$panam
2 amabananas$pan
3 anamabananas$p
4 ananas$panamab
5 anas$panamaban
6 as$panamabanan
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
panamabananas$
s$panamabanana
```

These strings are sorted

Chop off a

```
bananas$panam
mabananas$pan
namabananas$p
nanas$panamab
nas$panamaban
s$panamabanan
```

Still sorted

Add a
to end

```
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
s$panamabanana
```

Still sorted

177

Is It True in General?

```
$panamabananas
1 abananas$panam
2 amabananas$pan
3 anamabananas$p
4 ananas$panamab
5 anas$panamaban
6 as$panamabanan
bananas$panama 1
mabananas$pana 2
namabananas$pa 3
nanas$panamaba 4
nas$panamabana 5
panamabananas$
s$panamabanana 6
```

These strings are sorted

Chop off a

```
bananas$panam
mabananas$pan
namabananas$p
nanas$panamab
nas$panamaban
s$panamabanan
```

Still sorted

Add a
to end

```
bananas$panama
mabananas$pana
namabananas$pa
nanas$panamaba
nas$panamabana
s$panamabanana
```

Still sorted

178

Is It True in General?

- **First-Last Property:** The k -th occurrence of *symbol* in *FirstColumn* and the k -th occurrence of *symbol* in *LastColumn* correspond to the same position of *symbol* in *Genome*.

spiegata nella prima parte del corso,
dove sono precisamente riportate:

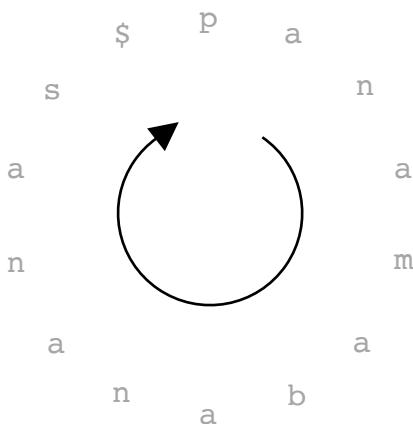
- definizioni
- algoritmi
- correttezza degli algoritmi

\$₁panamabananas**_{s₁}**
a₁bananas**\$panam₁**
a₂mabananas**\$pan₁**
a₃namabananas**\$p₁**
a₄nanas**\$panama_{b₁}**
a₅nas**\$panamaba_{n₂}**
a₆**\$panamabana_{n₃}**
b₁ananas**\$panama_{a₁}**
m₁abananas**\$pana_{a₂}**
n₁amabananas**\$pa_{a₃}**
n₂anas**\$panamaba_{a₄}**
n₃as**\$panamaba_{a₅}**
p₁anamabananas**\$₁**
s₁**\$panamabanan_{a₆}**

179

More Efficient BWT Decompression

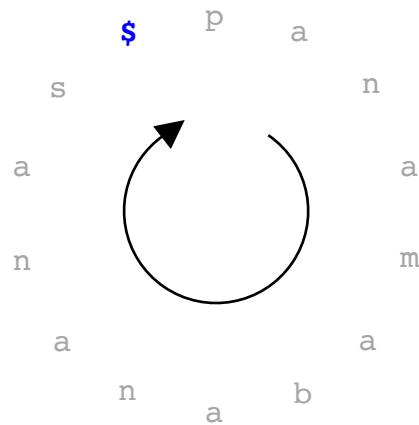
```
$1panamabananas$1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananass$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamabanans3  
b1ananas$panama1  
m1abananass$pana2  
n1amabananass$p3  
n2anas$panamaba4  
n3as$panamabanas5  
p1anamabananas$1  
s1$panamabanana6
```



180

More Efficient BWT Decompression

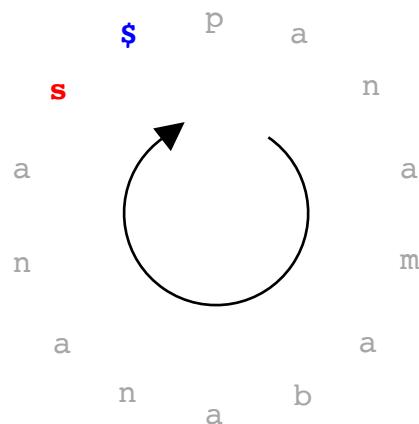
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



181

More Efficient BWT Decompression

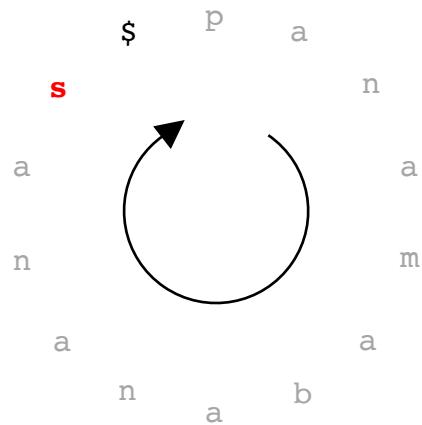
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



182

More Efficient BWT Decompression

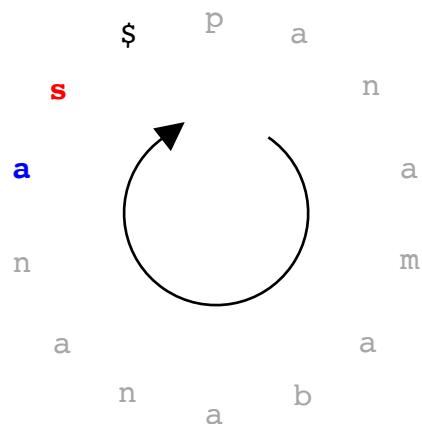
```
$1panamabananass1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamaban3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$\sub{1}
s1$panamabananaa6
```



183

More Efficient BWT Decompression

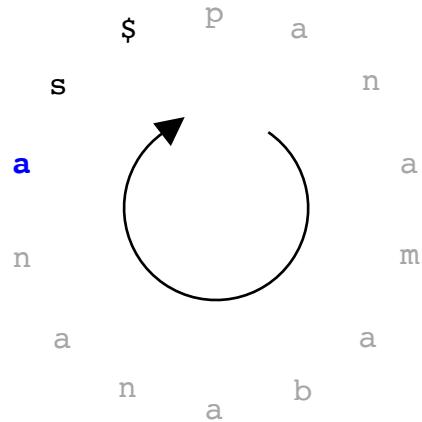
```
$1panamabananass1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamaban3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$\sub{1}
s1$panamabananaa6
```



184

More Efficient BWT Decompression

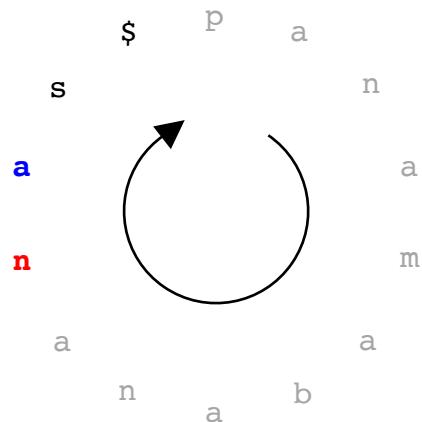
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabanaa5
p1anamabananas$1
s1$panamabananaa6
```



185

More Efficient BWT Decompression

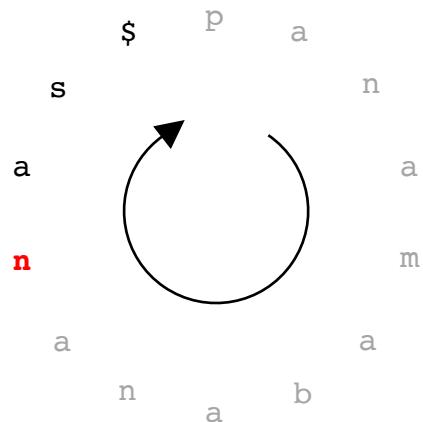
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabanaa5
p1anamabananas$1
s1$panamabananaa6
```



186

More Efficient BWT Decompression

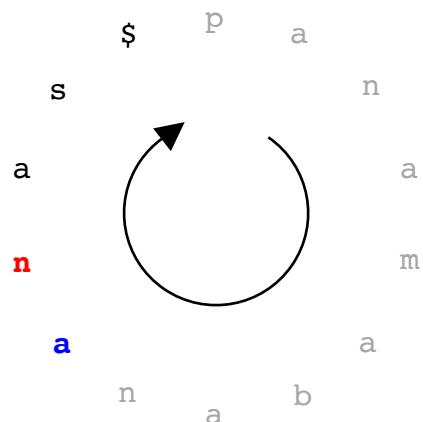
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



187

More Efficient BWT Decompression

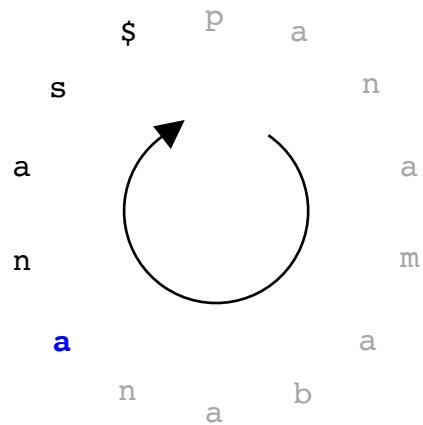
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



188

More Efficient BWT Decompression

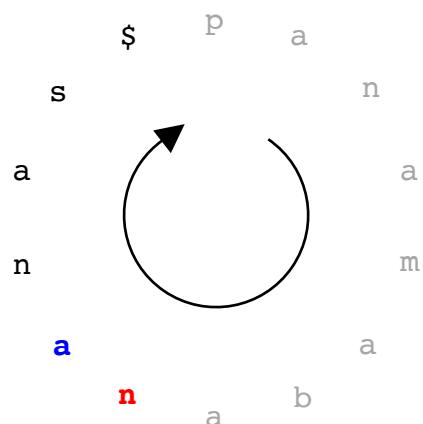
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaba2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabaa5
p1anamabananas$1
s1$panamabananaa6
```



189

More Efficient BWT Decompression

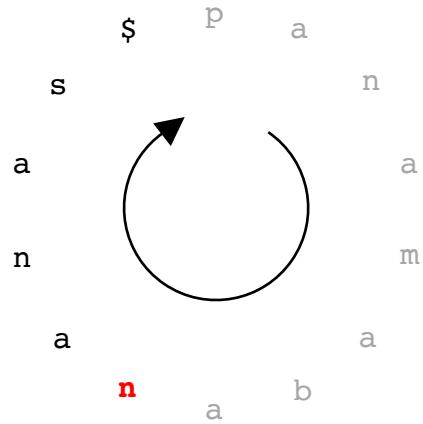
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamabanan3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabaa5
p1anamabananas$1
s1$panamabananaa6
```



190

More Efficient BWT Decompression

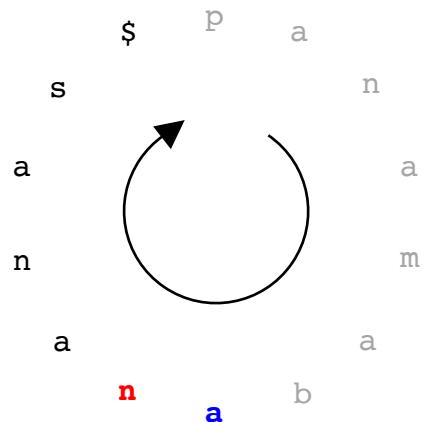
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaba2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabanaa5
p1anamabananas$1
s1$panamabananaa6
```



191

More Efficient BWT Decompression

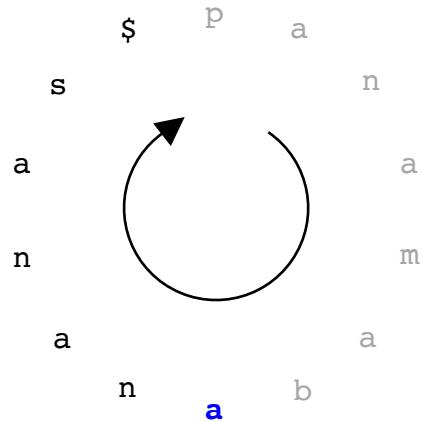
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaba2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabanaa5
p1anamabananas$1
s1$panamabananaa6
```



192

More Efficient BWT Decompression

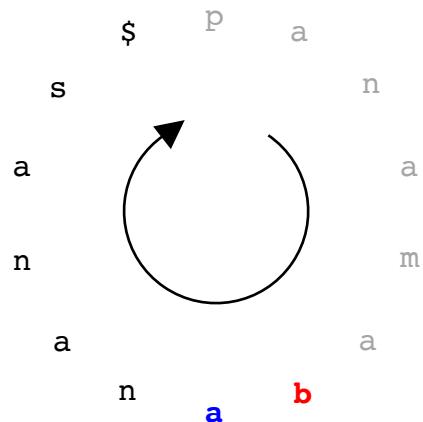
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabanaa5
p1anamabananas$1
s1$panamabananaa6
```



193

More Efficient BWT Decompression

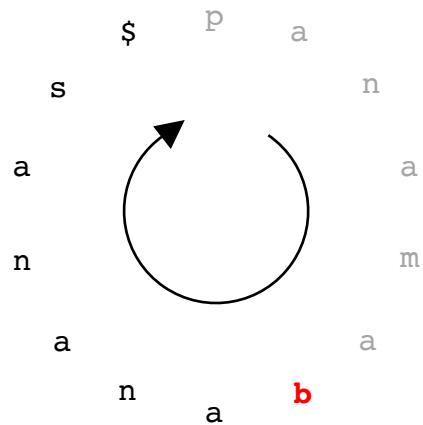
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabanaa5
p1anamabananas$1
s1$panamabananaa6
```



194

More Efficient BWT Decompression

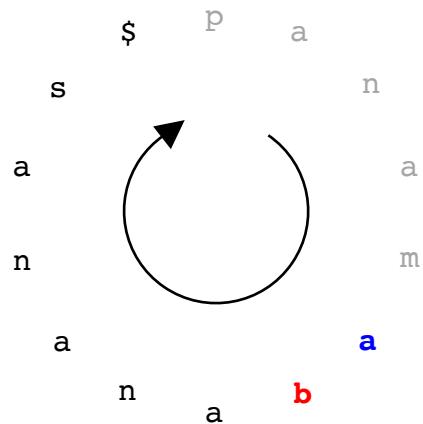
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
1ananas$panama1
m1abananas$pana2
n1amabananas$p3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



195

More Efficient BWT Decompression

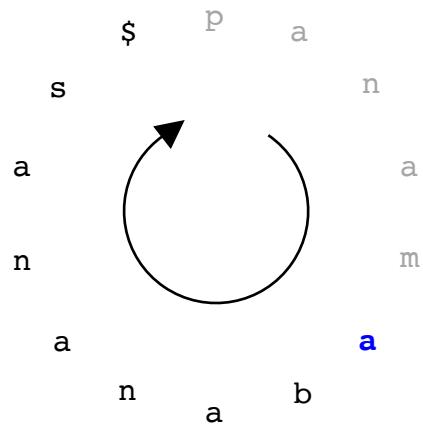
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
1ananas$panama1
m1abananas$pana2
n1amabananas$p3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



196

More Efficient BWT Decompression

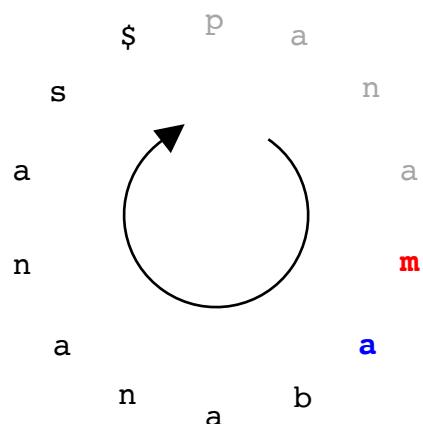
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



197

More Efficient BWT Decompression

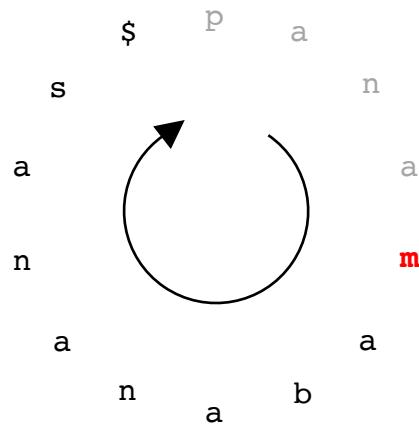
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



198

More Efficient BWT Decompression

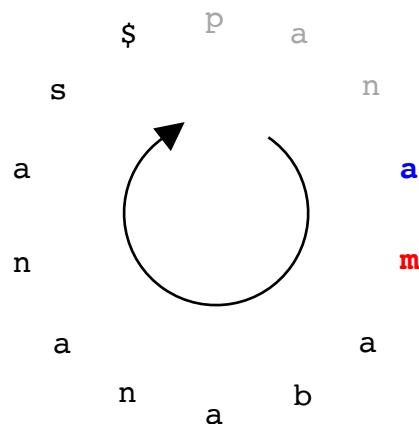
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



199

More Efficient BWT Decompression

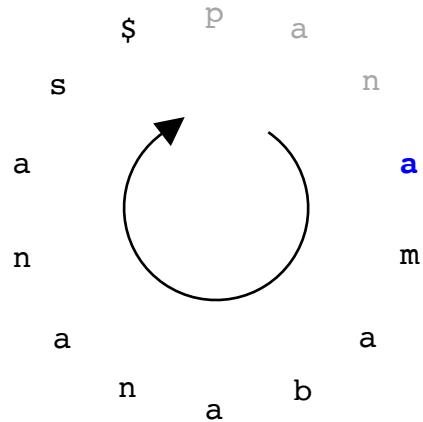
```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



200

More Efficient BWT Decompression

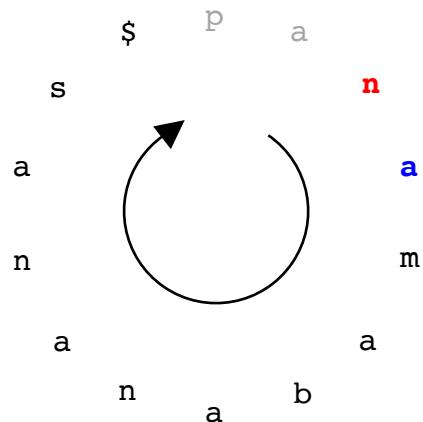
```
$1panamabananass1
a1bananas$panam1
a2mabananas$pann1
a3namabananas$p1
a4anas$panamab1
a5nas$panamabann2
a6s$panamabann3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabananaa6
```



201

More Efficient BWT Decompression

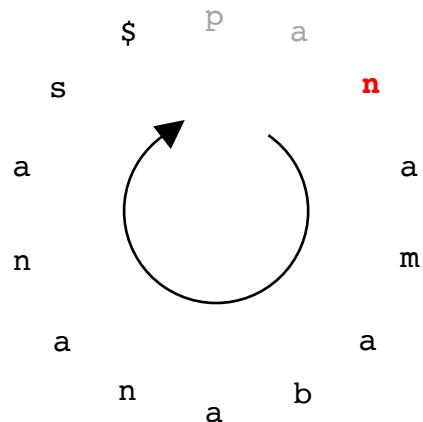
```
$1panamabananass1
a1bananas$panam1
a2mabananas$pann1
a3namabananas$p1
a4anas$panamab1
a5nas$panamabann2
a6s$panamabann3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabananaa6
```



202

More Efficient BWT Decompression

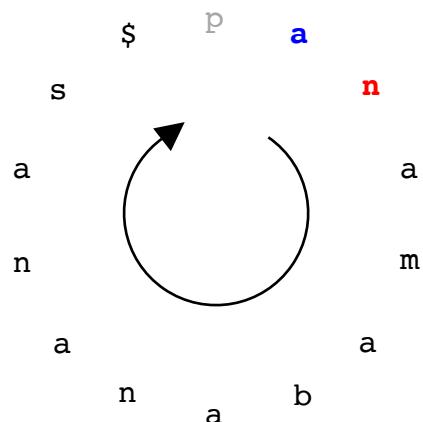
```
$1panamabananass1  
a1bananas$panam1  
a2mabananas$pann1  
a3namabananas$p1  
a4anas$panamab1  
a5nas$panamabann2  
a6s$panamabann3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabanana6
```



203

More Efficient BWT Decompression

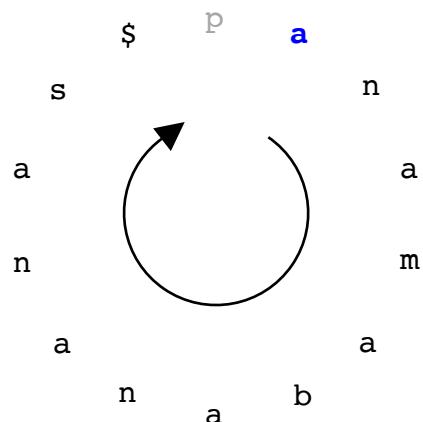
```
$1panamabananass1  
a1bananas$panam1  
a2mabananas$pann1  
a3namabananas$p1  
a4anas$panamab1  
a5nas$panamabann2  
a6s$panamabann3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabanana6
```



204

More Efficient BWT Decompression

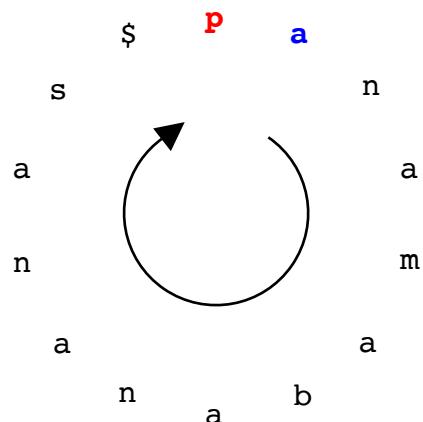
\$1panamabananas\$1
a1bananas\$panam1
a2mabananas\$pan1
a3namabananasp1
a4ananas\$panamab1
a5nas\$panamaban2
a6\$panamabanans3
b1ananas\$panama1
m1abanananas\$pana2
n1amabanananas\$pa3
n2anass\$panamaba4
n3as\$panamabanas5
p1anamabanananas\$1
s1\$panamabananana6



205

More Efficient BWT Decompression

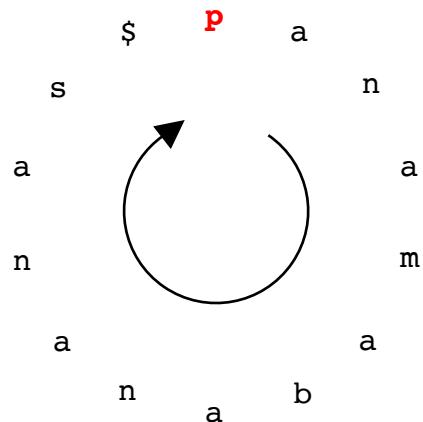
\$1panamabananas\$1
a1bananas\$panam1
a2mabananas\$pan1
a3namabananass\$pri
a4ananas\$panamab1
a5nas\$panamaban2
a6s\$panamabanans3
b1ananas\$panama1
m1abanananas\$pana2
n1amabanananas\$p3
n2anas\$panamaba4
n3as\$panamabana5
p1anamabananass1
s1\$panamabananana6



206

More Efficient BWT Decompression

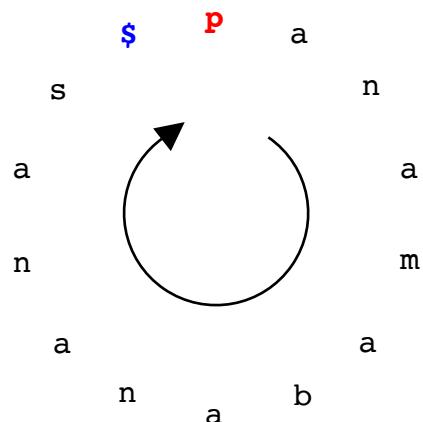
```
$1panamabananass1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamabann2  
a6s$panamabanan3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabanaa5  
p1anamabananas$1  
s1$panamabananaa6
```



207

More Efficient BWT Decompression

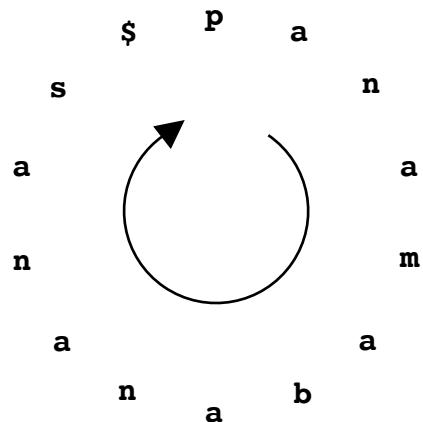
```
$1panamabananass1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamabann2  
a6s$panamabanan3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabanaa5  
p1anamabananas$1  
s1$panamabananaa6
```



208

More Efficient BWT Decompression

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamaban3
b1ananas$panama1
m1abananas$pana2
n1amabananas$p3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



- Memory: $2|Genome| = O(|Genome|)$.

209

What happened to pattern matching?



210

Outline

- Why do we map reads?
- Using the Trie
- From a Trie to a Suffix Tree
- String Compression and the Burrows-Wheeler Transform
- Inverting Burrows-Wheeler
- **Using Burrows-Wheeler for Pattern Matching**
- Finding the Matched Patterns
- Setting Up Checkpoints
- Inexact Matching
- Further Applications of Read Mapping

211

Recalling Our Goal

- Suffix Tree/array Pattern Matching:
 - Runtime: $O(|Genome| + |Patterns|)$
 - Memory: $O(|Genome|)$
 - Problem: suffix tree takes $20 \times |Genome|$ space,
but suffix array takes $\sim 4 \times |Genome|$ space.
- Can we use $BWT(Genome)$ as our data structure instead? BWT is a column...

212

Finding Pattern Matches Using BWT

- Searching for **ana** in **panamabananas**

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3na1mabananas$p1
a4na1nas$panamab1
a5na1s$panamab2
a6s$panamab1an3
b1ananas$panama1
m1abananas$pana2
n1amabananas$p3a3
n2anas$panamaba4
n3as$panamab5a5
p1anamabananas$1
s1$panamabanana6
```

213

Finding Pattern Matches Using BWT

- Searching for **ana** in **panamabananas**

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3na1mabananas$p1
a4na1nas$panamab1
a5na1s$panamab2
a6s$panamab1an3
b1ananas$panama1
m1abananas$pana2
n1amabananas$p3a3
n2anas$panamaba4
n3as$panamab5a5
p1anamabananas$1
s1$panamabanana6
```

214

Finding Pattern Matches Using BWT

- Searching for **ana** in **panamabananas**

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamabanan3
b1ananas$panama1
m1abananas$pana2
n1amabananas$p3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$\sub{1}
s1$panamabanan6
```

215

Finding Pattern Matches Using BWT

- Searching for **ana** in **panamabananas**

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamabanan3
b1ananas$panama1
m1abananas$pana2
n1amabananas$p3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$\sub{1}
s1$panamabanan6
```

216

Finding Pattern Matches Using BWT

- Searching for **ana** in **panamabananas**

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaba2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$p3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabana6
```

217

Finding Pattern Matches Using BWT

- Searching for **ana** in **panamabananas**

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaba2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$p3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabana6
```

218

Finding Pattern Matches Using BWT

- Searching for **ana** in **panamabananas**

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamabanan3
b1ananas$panama1
m1abananas$pana2
n1a1mabananas$p3
n2a2nas$panamab4
n3a3s$panamaban5
p1anamabananas$1
s1$panamabanan6
```

219

Finding Pattern Matches Using BWT

- Searching for **ana** in **panamabananas**

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamabanan3
b1ananas$panama1
m1abananas$pana2
n1a1mabananas$p3
n2a2nas$panamab4
n3a3s$panamaban5
p1anamabananas$1
s1$panamabanan6
```

220

Finding Pattern Matches Using BWT

- Searching for **ana** in **panamabananas**

consecutive occurrences!

 \$₁panamabananas₁
 a₁bananas\$panam₁
 a₂mabananas\$pan₁

a₃n₄amabananas\$p₁
a₄n₅anas\$panamab₁
a₅n₆as\$panamaban₂
 a₆s\$panamaban₃
 b₁ananas\$panama₁
 m₁abananas\$pana₂
 n₁amabananas\$p₃a₃
 n₂anas\$panamaba₄
 n₃as\$panamabana₅
 p₁anamabananas\$\sub{1}
 s₁\$panamabanana₆

221

Outline

- Why do we map reads?
- Using the Trie
- From a Trie to a Suffix Tree
- String Compression and the Burrows-Wheeler Transform
- Inverting Burrows-Wheeler
- Using Burrows-Wheeler for Pattern Matching
- Finding the Matched Patterns**
- Setting Up Checkpoints
- Inexact Matching
- Further Applications of Read Mapping

222

Where Are the Matches?

- **Multiple Pattern Matching Problem:**
 - **Input:** A collection of strings *Patterns* and a string *Genome*.
 - **Output:** All positions in *Genome* where one of *Patterns* appears as a substring.

223

Where Are the Matches?

- **Multiple Pattern Matching Problem:**
 - **Input:** A collection of strings *Patterns* and a string *Genome*.
 - **Output:** All **positions** in *Genome* where one of *Patterns* appears as a substring.
- Where are the **positions**? BWT has not revealed them.

224

Where Are the Matches?

- Example: We know that **ana** occurs 3 times, but where?

```
$1panamabananass1  
a1bananas$panamm1  
a2mabananas$pann1  
a3namabananas$pp1  
a4nanas$panamabb1  
a5nas$panamaban2  
a6s$panamaban3  
b1ananas$panamaa1  
m1abananas$panaa2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabaa5  
p1anamabananas$\sub{s1  
s1$panamabananaa6
```

225

Using the Suffix Array to Find Matches

- Suffix array:** holds starting position of each suffix beginning a row.

From BWT to SA
(mentre costruisco la BWT,
creo il SA)

```
$1panamabananass1  
a1bananas$panamm1  
a2mabananas$pann1  
a3namabananas$pp1  
a4nanas$panamabb1  
a5nas$panamaban2  
a6s$panamaban3  
b1ananas$panamaa1  
m1abananas$panaa2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabaa5  
p1anamabananas$\sub{s1  
s1$panamabananaa6
```

226

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panamabananass\$

1 3

\$₁panamabananas₁
a₁bananas\$panam₁
a₂mabananas\$pan₁
a₃namabanananas\$p₁
a₄nanas\$panamab₁
a₅nas\$panamaban₂
a₆s\$panamaban₃
b₁ananas\$panama₁
m₁abananas\$pana₂
n₁amabanananas\$pa₃
n₂anas\$panamaba₄
n₃as\$panamabana₅
p₁anamabanananas\$₁
s₁\$panamabana₆

227

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panam**a**bananas\$

1 3
5

\$₁panamabananas₁
a₁**bananas\$panam**₁
a₂mabananas\$pan₁
a₃namabanananas\$p₁
a₄nanas\$panamab₁
a₅nas\$panamaban₂
a₆s\$panamaban₃
b₁ananas\$panama₁
m₁abananas\$pana₂
n₁amabanananas\$pa₃
n₂anas\$panamaba₄
n₃as\$panamabana₅
p₁anamabanananas\$₁
s₁\$panamabana₆

228

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panamabananass\$

1	3	\$ ₁ panamabananas ₁
5		a ₁ bananas\$panam ₁
3		a ₂ mabananas\$pan ₁
		a ₃ namabananas\$p ₁
		a ₄ anas\$panamab ₁
		a ₅ nas\$panamaban ₂
		a ₆ s\$panamaban ₃
		b ₁ ananas\$panama ₁
		m ₁ abananas\$pana ₂
		n ₁ amabananas\$pa ₃
		n ₂ anas\$panamaba ₄
		n ₃ as\$panamabana ₅
		p ₁ anamabananas\$ ₁
		s ₁ \$panamabana ₆

229

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panamabananass\$

1	3	\$ ₁ panamabananas ₁
5		a ₁ bananas\$panam ₁
3		a ₂ mabananas\$pan ₁
1		a ₃ namabananas\$p ₁
		a ₄ anas\$panamab ₁
		a ₅ nas\$panamaban ₂
		a ₆ s\$panamaban ₃
		b ₁ ananas\$panama ₁
		m ₁ abananas\$pana ₂
		n ₁ amabananas\$pa ₃
		n ₂ anas\$panamaba ₄
		n ₃ as\$panamabana ₅
		p ₁ anamabananas\$ ₁
		s ₁ \$panamabana ₆

230

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panamab**a**nanas\$

1	3	\$ ₁ panamabananas ₁
5		a₁bananas\$ panam ₁
3		a₂mabananas\$ pan ₁
1		a₃namabananas\$ p ₁
7		a₄nanas\$ panamab ₁
		a ₅ nas\$panamaban ₂
		a ₆ s\$panamaban ₃
		b ₁ ananas\$panama ₁
		m ₁ abananas\$pana ₂
		n ₁ amabananas\$pa ₃
		n ₂ anas\$panamaba ₄
		n ₃ as\$panamabana ₅
		p ₁ anamabananas\$ ₁
		s₁\$panamabana₆

231

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panamab**a**nanas\$

1	3	\$ ₁ panamabananas ₁
5		a₁bananas\$ panam ₁
3		a₂mabananas\$ pan ₁
1		a₃namabananas\$ p ₁
7		a₄nanas\$ panamab ₁
9		a₅nas\$ panamaban ₂
		a ₆ s\$panamaban ₃
		b ₁ ananas\$panama ₁
		m ₁ abananas\$pana ₂
		n ₁ amabananas\$pa ₃
		n ₂ anas\$panamaba ₄
		n ₃ as\$panamabana ₅
		p ₁ anamabananas\$ ₁
		s₁\$panamabana₆

232

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panamabananas**a**s\$

1 3	\$ ₁ panamabananas ₁
5	a ₁ bananas \$panam ₁
3	a ₂ mabananas \$pan ₁
1	a ₃ namabananas \$p ₁
7	a ₄ nanas \$panamab ₁
9	a ₅ nas \$panamaban ₂
1 1	a ₆ s \$panamaban ₃
	b ₁ ananas \$panama ₁
	m ₁ abananas \$pana ₂
	n ₁ amabananas \$pa ₃
	n ₂ anas \$panamaba ₄
	n ₃ as \$panamabana ₅
	p ₁ anamabananas \$ ₁
	s ₁ \$panamabana ₆

233

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panama**bananas**\$

1 3	\$ ₁ panamabananas ₁
5	a ₁ bananas \$panam ₁
3	a ₂ mabananas \$pan ₁
1	a ₃ namabananas \$p ₁
7	a ₄ nanas \$panamab ₁
9	a ₅ nas \$panamaban ₂
1 1	a ₆ s \$panamaban ₃
6	b ₁ ananas \$panama ₁
	m ₁ abananas \$pana ₂
	n ₁ amabananas \$pa ₃
	n ₂ anas \$panamaba ₄
	n ₃ as \$panamabana ₅
	p ₁ anamabananas \$ ₁
	s ₁ \$panamabana ₆

234

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panamabananass\$

1 3	\$ ₁ panamabananas ₁
5	a ₁ bananas\$panam ₁
3	a ₂ mabananas\$pan ₁
1	a ₃ namabananas\$p ₁
7	a ₄ nanas\$panamab ₁
9	a ₅ nas\$panamaban ₂
1 1	a ₆ s\$panamaban ₃
6	b ₁ ananas\$panama ₁
4	m ₁ abananas\$p _{ana} ₂
	n ₁ amabananas\$p _a ₃
	n ₂ anas\$panamaba ₄
	n ₃ as\$panamabana ₅
	p ₁ anamabananas\$p ₁
	s ₁ \$panamabana ₆

235

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panamabananass\$

1 3	\$ ₁ panamabananas ₁
5	a ₁ bananas\$panam ₁
3	a ₂ mabananas\$pan ₁
1	a ₃ namabananas\$p ₁
7	a ₄ nanas\$panamab ₁
9	a ₅ nas\$panamaban ₂
1 1	a ₆ s\$panamaban ₃
6	b ₁ ananas\$panama ₁
4	m ₁ abananas\$p _{ana} ₂
2	n ₁ amabananas\$p _a ₃
	n ₂ anas\$panamaba ₄
	n ₃ as\$panamabana ₅
	p ₁ anamabananas\$p ₁
	s ₁ \$panamabana ₆

236

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panamabana**nas\$**

13	\$ ₁ panamabananas ₁
5	a ₁ bananas\$panam ₁
3	a ₂ mabananas\$pan ₁
1	a ₃ namabananas\$p ₁
7	a ₄ nanas\$panamab ₁
9	a ₅ nas\$panamaban ₂
11	a ₆ s\$panamaban ₃
6	b ₁ ananas\$panama ₁
4	m ₁ abananas\$p _{ana} ₂
2	n ₁ amabananas\$p _a ₃
8	n ₂ anas\$panamaba ₄
	n ₃ as\$panamabana ₅
	p ₁ anamabananas\$p ₁
	s ₁ \$panamabana ₆

237

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panamabana**nas\$**

13	\$ ₁ panamabananas ₁
5	a ₁ bananas\$panam ₁
3	a ₂ mabananas\$pan ₁
1	a ₃ namabananas\$p ₁
7	a ₄ nanas\$panamab ₁
9	a ₅ nas\$panamaban ₂
11	a ₆ s\$panamaban ₃
6	b ₁ ananas\$panama ₁
4	m ₁ abananas\$p _{ana} ₂
2	n ₁ amabananas\$p _a ₃
8	n ₂ anas\$panamaba ₄
10	n ₃ as\$panamabana ₅
	p ₁ anamabananas\$p ₁
	s ₁ \$panamabana ₆

238

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panamabananas\$

13	\$ ₁ panamabananas ₁
5	a ₁ bananas\$panam ₁
3	a ₂ mabananas\$pan ₁
1	a ₃ namabananas\$p ₁
7	a ₄ nanas\$panamab ₁
9	a ₅ nas\$panamaban ₂
11	a ₆ s\$panamaban ₃
6	b ₁ ananas\$panama ₁
4	m ₁ abananas\$p _{ana} ₂
2	n ₁ amabananas\$p _a ₃
8	n ₂ anas\$panamaba ₄
10	n ₃ as\$panamabana ₅
0	p ₁ anamabananas\$\sub{1}
	s ₁ \$panamabana ₆

239

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panamabananas\$ \$

13	\$ ₁ panamabananas ₁
5	a ₁ bananas\$panam ₁
3	a ₂ mabananas\$pan ₁
1	a ₃ namabananas\$p ₁
7	a ₄ nanas\$panamab ₁
9	a ₅ nas\$panamaban ₂
11	a ₆ s\$panamaban ₃
6	b ₁ ananas\$panama ₁
4	m ₁ abananas\$p _{ana} ₂
2	n ₁ amabananas\$p _a ₃
8	n ₂ anas\$panamaba ₄
10	n ₃ as\$panamabana ₅
0	p ₁ anamabananas\$\sub{1}
12	s ₁ \$panamabana ₆

240

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

panamabananass\$

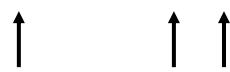
1 3	\$ ₁ panamabananass ₁
5	a ₁ bananas\$panam ₁
3	a ₂ mabananas\$pan ₁
1	a ₃ namabanananas\$p ₁
7	a ₄ nanas\$panamab ₁
9	a ₅ nas\$panamaban ₂
1 1	a ₆ s\$panamabanan ₃
6	b ₁ ananas\$panama ₁
4	m ₁ abanananas\$pana ₂
2	n ₁ amabanananas\$p ₃
8	n ₂ anas\$panamaba ₄
1 0	n ₃ as\$panamaban ₅
0	p ₁ anamabananas\$s ₁
1 2	s ₁ \$panamabanan ₆

241

Using the Suffix Array to Find Matches

- **Suffix array:** holds starting position of each suffix beginning a row.

- Thus, **ana** occurs at positions **1, 7, 9** of **panamabananass\$**.



1 3	\$ ₁ panamabananass ₁
5	a ₁ bananas\$panam ₁
3	a ₂ mabananas\$pan ₁
1	a₃na mabanananas\$p ₁
7	a₄na nas\$panamab ₁
9	a₅na s\$panamaban ₂
1 1	a ₆ s\$panamabanan ₃
6	b ₁ ananas\$panama ₁
4	m ₁ abanananas\$pana ₂
2	n ₁ amabanananas\$p ₃
8	n ₂ anas\$panamaba ₄
1 0	n ₃ as\$panamaban ₅
0	p ₁ anamabananas\$s ₁
1 2	s ₁ \$panamabanan ₆

242

The Suffix Array: Memory Once Again

- Memory: $\sim 4 \times |Genome|$.

time complexity later...

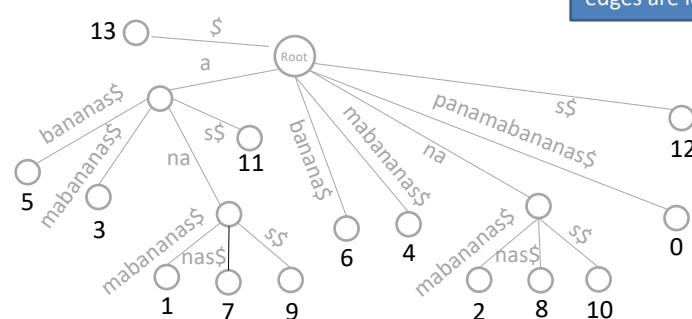
243

The Suffix Array: Memory Once Again

- Memory: $\sim 4 \times |Genome|$.

From ST to SA

ST hides the SA
(in-order visit if the edges are lex. ordered)



[13 5 3 1 7 9 11 6 4 2 8 10 0 12]

244

From SA to BWT

$s = \text{appellee\$}$
123456789

9
1
8
7
4
6
5
3
2

Suffix array
(start position
for the suffixes)

245

From SA to BWT

$s = \text{appellee\$}$
123456789

9
1
8
7
4
6
5
3
2

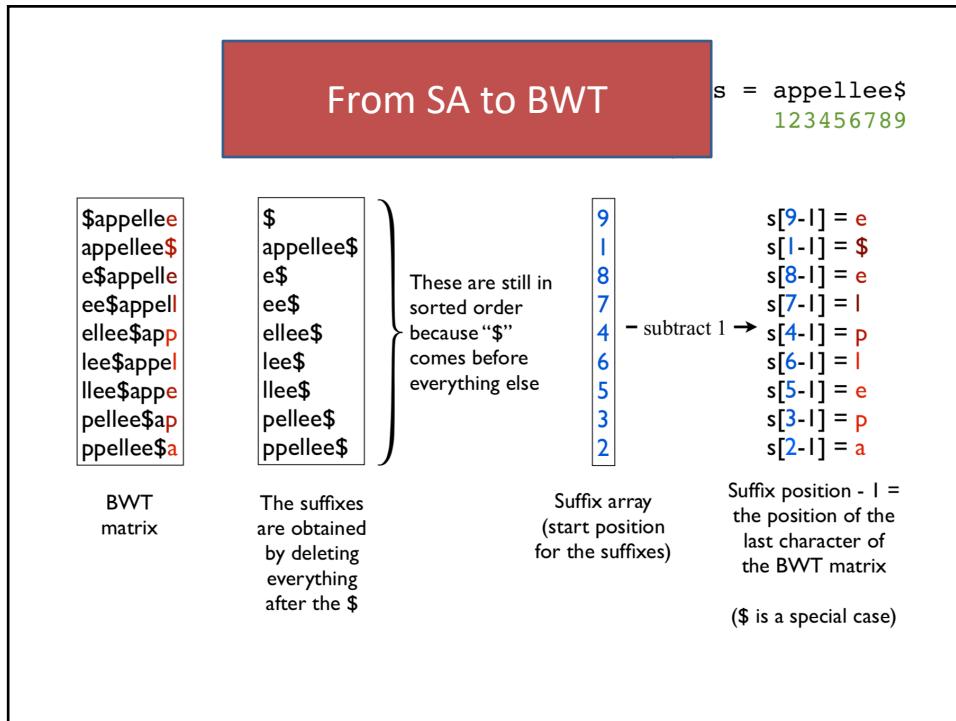
$s[9-1] = e$
 $s[1-1] = \$$
 $s[8-1] = e$
 $s[7-1] = l$
 - subtract 1 → $s[4-1] = p$
 $s[6-1] = l$
 $s[5-1] = e$
 $s[3-1] = p$
 $s[2-1] = a$

Suffix array
(start position
for the suffixes)

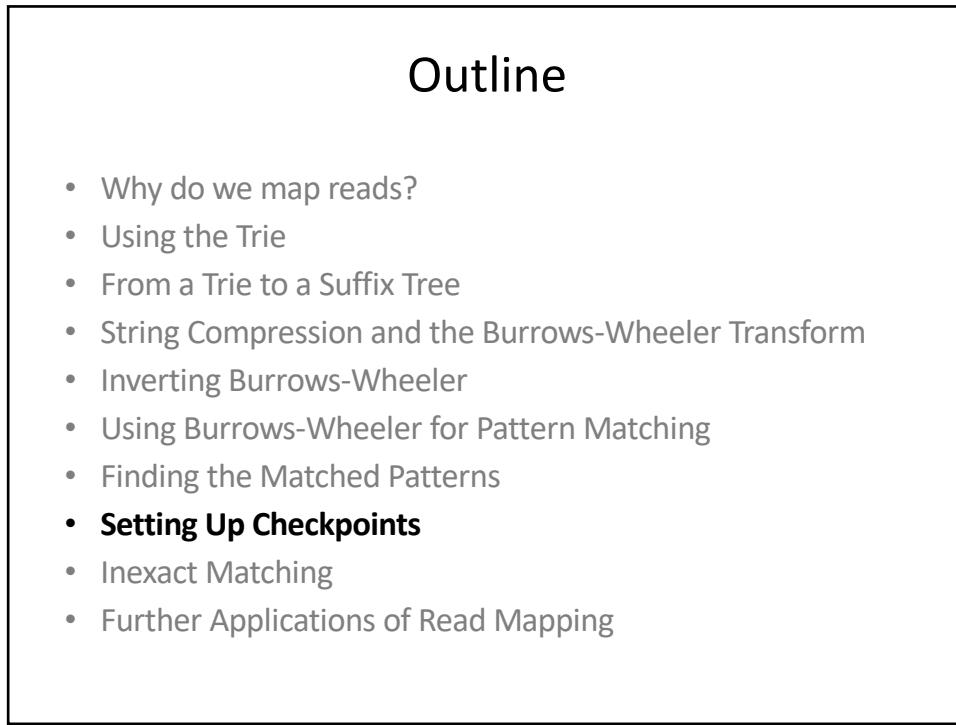
Suffix position - 1 =
the position of the
last character of
the BWT matrix

(\$ is a special case)

246



247



248

Reducing Suffix Array Size

- We don't want to have to store all of the suffix array; can we store only part of it? Show how checkpointing can be used to store 1/100 the suffix array.

chi vuole approfondire, può seguire
dal testo/youtube
(progetto)

249

Outline

- Why do we map reads?
- Using the Trie
- From a Trie to a Suffix Tree
- String Compression and the Burrows-Wheeler Transform
- Inverting Burrows-Wheeler
- Using Burrows-Wheeler for Pattern Matching
- Finding the Matched Patterns
- Setting Up Checkpoints
- **Inexact Matching**
- Further Applications of Read Mapping

250

Read Mapping

- **Read mapping:** determine where each read has high similarity to the reference genome.

without assembling

CTGAGGATGGACTACGCTACTACTGATAGCTGTTT	Reference
GAGGA	Reads
CCACG	TGA-A
exact	mismatch
	insdel

251

Returning to Our Original Problem

- We need to look at INEXACT matching in order to find variants.
- **Approx. Pattern Matching Problem:**
 - **Input:** A string *Pattern*, a string *Genome*, and an integer *d*.
 - **Output:** All positions in *Genome* where *Pattern* appears as a substring with at most *d* mismatches.

252

Returning to Our Original Problem

- We need to look at INEXACT matching in order to find variants.
- **Multiple Approx. Pattern Matching Problem:**
 - **Input:** A **collection** of strings *Patterns*, a string *Genome*, and an integer d .
 - **Output:** All positions in *Genome* where a string from *Patterns* appears as a substring with at most d mismatches.

253

Method 1: Seeding

- Say we are looking for at most d mismatches.
- Divide each of our strings into $d + 1$ smaller pieces, called **seeds**.
- Check if each *Pattern* has a seed that matches *Genome* exactly.
- If so, check the entire *Pattern* against *Genome*.

254

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamaban3
b1anas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```

255

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaban2
a6s$panamaban3
b1anas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```

256

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaba2
a6s$panamabana3
b1anas$panama1
m1abananas$pana2
n1amabananas$p3
n2anas$panamab4
n3as$panamaba5
p1anamabananas$1
s1$panamabana6
```

257

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

If we allow 1 mismatch,
then we need to
keep the
red letters around.

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4anas$panamab1
a5nas$panamaba2
a6s$panamabana3
b1anas$panama1
m1abananas$pana2
n1amabananas$p3
n2anas$panamab4
n3as$panamaba5
p1anamabananas$1
s1$panamabana6
```

258

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

If we allow 1 mismatch,
then we need to
keep the
red letters around.

Mismatches
\$ ₁ panamabananas _{s₁}
a₁ bananas\$pana m₁
a₂ mabananas\$pan n₁
a₃ namabananas\$p _{p₁}
a₄ nanas\$panamab _{b₁}
a₅ nass\$panamaba _{n₂}
a₆ s\$panamaba _{n₃}
b₁ ananas\$panama _{a₁}
m₁ abananas\$pana _{a₂}
n₁ amabananas\$pa _{a₃}
n₂ anas\$panamaba _{a₄}
n₃ as\$panamaba _{a₅}
p₁ anamabananas\$ _{s₁}
s₁ \$panamabanana _{a₆}

259

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

If we allow 1 mismatch,
then we need to
keep the
red letters around.

Mismatches
\$ ₁ panamabananas _{s₁}
a₁ bananas\$pana m₁ 1
a₂ mabananas\$pan n₁ 0
a₃ namabananas\$p _{p₁} 1
a₄ nanas\$panamab _{b₁} 1
a₅ nass\$panamaba _{n₂} 0
a₆ s\$panamaba _{n₃} 0
b₁ ananas\$panama _{a₁}
m₁ abananas\$pana _{a₂}
n₁ amabananas\$pa _{a₃}
n₂ anas\$panamaba _{a₄}
n₃ as\$panamaba _{a₅}
p₁ anamabananas\$ _{s₁}
s₁ \$panamabanana _{a₆}

260

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

Now we extend all strings with at most 1 mismatch.

	# Mismatches
\$ ₁ panamabananas ₁	
a ₁ bananas\$panam ₁	1
a ₂ mabananas\$pan ₁	0
a ₃ namabananas\$p ₁	1
a ₄ nanas\$panamab ₁	1
a ₅ nass\$panamaban ₂	0
a ₆ s\$panamabanan ₃	0
b ₁ ananas\$panama ₁	
m ₁ abananas\$pana ₂	
n ₁ amabananas\$pa ₃	
n ₂ anas\$panamaba ₄	
n ₃ as\$panamaban ₅	
p ₁ anamabananas\$ ₁	
s ₁ \$panamabanana ₆	

261

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

Now we extend all strings with at most 1 mismatch.

	# Mismatches
\$ ₁ panamabananas ₁	
a ₁ bananas\$panam ₁	1
a ₂ mabananas\$pan ₁	0
a ₃ namabananas\$p ₁	1
a ₄ nanas\$panamab ₁	1
a ₅ nass\$panamaban ₂	0
a ₆ s\$panamabanan ₃	0
b ₁ a ₁ nanas\$panama ₁	
m ₁ a ₁ bananas\$pana ₂	
n ₁ a ₁ mabananas\$pa ₃	
n ₂ a ₁ nas\$panamaba ₄	
n ₃ a ₁ s\$panamaban ₅	
p ₁ a ₁ namabananas\$ ₁	
s ₁ \$panamabanana ₆	

262

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

Now we extend all strings with at most 1 mismatch.

	# Mismatches
\$ ₁ panamabananas ₁	
a ₁ bananas\$panam ₁	1
a ₂ mabananas\$pan ₁	0
a ₃ namabananas\$p ₁	1
a ₄ nanas\$panamab ₁	1
a ₅ nass\$panamaban ₂	0
a ₆ s\$panamabanan ₃	0
b ₁ ananas\$panama ₁	
m ₁ a ₁ bananas\$pana ₂	
n ₁ a ₂ mabananas\$p ₃	
n ₂ a ₃ nas\$panamaba ₄	
n ₃ a ₄ s\$panamabana ₅	
p ₁ anamabananas\$ ₁	
s ₁ \$panamabanana ₆	

263

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

Now we extend all strings with at most 1 mismatch.

	# Mismatches
\$ ₁ panamabananas ₁	
a ₁ bananas\$panam ₁	
a ₂ mabananas\$pan ₁	
a ₃ namabananas\$p ₁	
a ₄ nanas\$panamab ₁	
a ₅ nass\$panamaban ₂	
a ₆ s\$panamabanan ₃	
b ₁ ananas\$panama ₁	1
m ₁ a ₁ bananas\$pana ₂	1
n ₁ a ₂ mabananas\$p ₃	0
n ₂ a ₃ nas\$panamaba ₄	0
n ₃ a ₄ s\$panamabana ₅	0
p ₁ anamabananas\$ ₁	1
s ₁ \$panamabanana ₆	

264

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

One string
produces a second
mismatch (the \$),
so we discard it.

	# Mismatches
\$ ₁ panamabananas ₁	
a ₁ bananas\$panam ₁	
a ₂ mabananas\$pan ₁	
a ₃ namabananas\$p ₁	
a ₄ anas\$panamab ₁	
a ₅ nass\$panamaban ₂	
a ₆ s\$panamaban ₃	
b₁a anas\$panama ₁	1
m₁a bananas\$pana ₂	1
n₁a mabananas\$p ₃	0
n₂a nas\$panamaba ₄	0
n₃a s\$panamaban ₅	0
p₁a namabananas\$ ₁	2
s ₁ \$panamaban ₆	

265

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

In the end, we
have five 3-mers
with at most 1
mismatch.

	# Mismatches
\$ ₁ panamabananas ₁	
a₁b ananas\$panam ₁	1
a₂m abananas\$pan ₁	1
a₃n amabananas\$p ₁	0
a₄n as\$panamab ₁	0
a₅n as\$panamaban ₂	0
a ₆ s\$panamaban ₃	
b₁a nas\$panama ₁	
m₁a bananas\$pana ₂	
n₁a mabananas\$p ₃	
n₂a ns\$panamaba ₄	
n₃a s\$panamaban ₅	
p₁a namabananas\$ ₁	
s ₁ \$panamaban ₆	

266

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

In the end, we have five 3-mers with at most 1 mismatch.

	# Mismatches
\$ ₁ panamabananas ₁	
a ₁ b ₁ ananas\$panam ₁	1
a ₂ m ₁ abananas\$pan ₁	1
a ₃ n ₁ amabananas\$p ₁	0
a ₄ n ₁ anas\$panamab ₁	0
a ₅ n ₁ a\$spanamaban ₂	0
a ₆ s\$panamaban ₃	
b ₁ ananas\$panama ₁	
m ₁ abananas\$pana ₂	
n ₁ amabananas\$pa ₃	
n ₂ anas\$panamaba ₄	
n ₃ as\$panamabana ₅	
p ₁ anamabananas\$ ₁	
s ₁ \$panamabana ₆	

267

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

In the end, we have five 3-mers with at most 1 mismatch.

	Suffix Array
\$ ₁ panamabananas ₁	
a ₁ b ₁ ananas\$panam ₁	5
a ₂ m ₁ abananas\$pan ₁	3
a ₃ n ₁ amabananas\$p ₁	1
a ₄ n ₁ anas\$panamab ₁	7
a ₅ n ₁ a\$spanamaban ₂	9
a ₆ s\$panamaban ₃	
b ₁ ananas\$panama ₁	
m ₁ abananas\$pana ₂	
n ₁ amabananas\$pa ₃	
n ₂ anas\$panamaba ₄	
n ₃ as\$panamabana ₅	
p ₁ anamabananas\$ ₁	
s ₁ \$panamabana ₆	

268

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamab****a**nanas

In the end, we have five 3-mers with at most 1 mismatch.

Suffix Array		
\$ ₁	panamab	a nanas
a ₁	b	ananas\$panam
a ₂	m	abananas\$pana
a ₃	n	namabananas\$p
a ₄	a	anas\$panamab
a ₅	s	nas\$panamaba
a ₆		s\$panamabana
b ₁	a	ananas\$panama
m ₁	a	abananas\$pana
n ₁	a	namabananas\$p
n ₂	a	anas\$panamab
n ₃	s	s\$panamaba
p ₁	a	namabananas\$
s ₁		s\$panamabana

269

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamab****a**nanas

In the end, we have five 3-mers with at most 1 mismatch.

Suffix Array		
\$ ₁	panamab	a nanas
a ₁	bananas	\$panam
a ₂	m	abananas\$pana
a ₃	a	namabananas\$p
a ₄	n	anas\$panamab
a ₅	s	nas\$panamaba
a ₆	s	s\$panamabana
b ₁	a	ananas\$panama
m ₁	a	abananas\$pana
n ₁	a	namabananas\$p
n ₂	a	anas\$panamab
n ₃	s	s\$panamaba
p ₁	a	namabananas\$
s ₁		s\$panamabana

270

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

In the end, we have five 3-mers with at most 1 mismatch.

Suffix Array		
\$ ₁	panamabananas	s ₁
a ₁	bananas\$panam	1
a ₂	mabananas\$pan	3
a₃na	mabananas\$p ₁	1
a ₄	nanas\$panamab	7
a ₅	nass\$panamaba	9
a ₆ s	panamabanan	3
b ₁ ana	anas\$panama	1
m ₁ aba	abananas\$pana	2
n ₁ ama	amabananas\$pa	3
n ₂ ana	anas\$panamaba	4
n ₃ as	s\$panamaba	5
p ₁ anamaba	nanamas\$	1
s ₁ \$panamaba	anana	6

271

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

In the end, we have five 3-mers with at most 1 mismatch.

Suffix Array		
\$ ₁	panamabananas	s ₁
a ₁	bananas\$panam	1
a ₂	mabananas\$pan	3
a₃nam	abananas\$p ₁	1
a₄na s	nas\$panamab	7
a ₅ n	ass\$panamaba	9
a ₆ s\$panamaba	nanana	3
b ₁ ana	anas\$panama	1
m ₁ aba	abananas\$pana	2
n ₁ ama	amabananas\$pa	3
n ₂ ana	anas\$panamaba	4
n ₃ as	s\$panamaba	5
p ₁ anamaba	nanamas\$	1
s ₁ \$panamaba	anana	6

272

Method 2: BWT Saves the Day Again

- Recall: searching for **ana** in **panamabananas**

In the end, we have five 3-mers with at most 1 mismatch.

Suffix Array		
\$ ₁	panamabananas ₁	
a ₁	bananas\$panam ₁	5
a ₂	mabananas\$pan ₁	3
a ₃	namabananas\$p ₁	1
a ₄	anas\$panamab ₁	7
a₅	nas\$panamaban₂	9
a ₆	s\$panamaban ₃	
b ₁	ananas\$panama ₁	
m ₁	abananas\$pana ₂	
n ₁	amabananas\$pa ₃	
n ₂	anas\$panamaba ₄	
n ₃	as\$panamaban ₅	
p ₁	anamabananas\$ ₁	
s ₁	\$panamabanana ₆	

273

Citations

Compression

[PDF] A block-sorting lossless data compression algorithm

www.hpl.hp.com/techreports/Compaq-DEC/SRC-RR-124.pdf ▾

by AB Lossless - 1994 - Related articles

May 10, 1994 - A Block-sorting Lossless. Data Compression Algorithm. M. Burrows and D.J. Wheeler digital. Systems Research Center. 130 Lytton Avenue.

Pattern Matching

Opportunistic data structures with applications - IEEE Conference ...

<https://ieeexplore.ieee.org/document/892127>

by P Ferragina - 2000 - Cited by 1016 - Related articles

Opportunistic data structures with applications. Abstract: We address the issue of compressing and indexing data. We devise a data structure whose space occupancy is a function of the entropy of the underlying data set.

Read Mapping

Fast and accurate short read alignment with Burrows–Wheeler transform

<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2705234/> ▾

by H Li - 2009 - Cited by 17996 - Related articles

May 18, 2009 - Results: We implemented Burrows-Wheeler Alignment tool (**BWA**), a new read alignment package that is based on backward search with ...

Ultrafast and memory-efficient alignment of short DNA sequences to ...

<https://genomebiology.biomedcentral.com/articles/10.1186/gb-2009-10-3-r25> ▾

by B Langmead - 2009 - Cited by 13797 - Related articles

Mar 4, 2009 - **Bowtie** is an ultrafast, memory-efficient alignment program for aligning short DNA sequence reads to large genomes. For the human genome, ...

274

From BWT + FM to SA

The FM-index

- It is a 2D array whose size is $[|text| + 1, |\Sigma|]$, where $|\Sigma|$ is the alphabet size
- It keeps track of how many of each symbol have been seen in the BWT prior to its i^{th} symbol
- The last m row is the totals for each symbol. By accumulating these totals you can determine the BWT index corresponding to the first of each symbol in the suffix array (Offset).
- Can be generated by a single scan through the BWT
- Memory overhead $O(m|\Sigma|)$

Counts(x)=numero di occorrenze di x
 Occ(x, i)=numero di caratteri x in $\text{BWT}[0, i-1]$, $i > 0$.
 Offset: accumulo i conteggi per scoprire la prima occorrenza nel SA
 \$ è in 0
 a è in 0+1=1
 b è in 1+3=4
 n è in 4+1=5

		FM-index				
Index	Suffix Array	BWT	\$	a	b	n
0	\$banana	a	0	0	0	0
1	a\$banan	n	0	1	0	0
2	ana\$ban	n	0	1	0	1
3	anana\$b	b	0	1	0	2
4	banana\$	\$	0	1	1	2
5	na\$bana	a	1	1	1	2
6	nana\$ba	a	1	2	1	2
7	Counts		1	3	1	2
	Offset		0	1	4	5

275

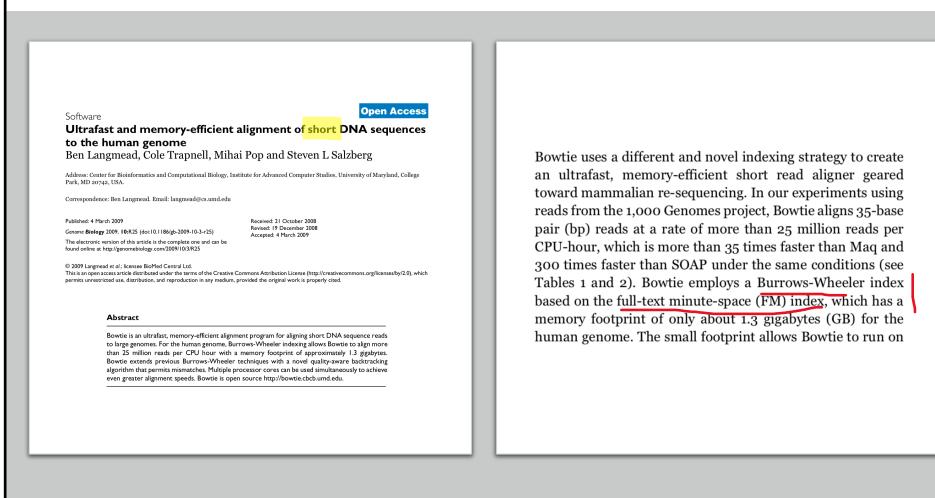
BWT+FM: migliora la ricerca binaria già vista...

Exact Matching with FM Index

- In progressive rounds, **top** & **bot** delimit the range of rows beginning with progressively longer suffixes of Q

276

BWT Aligner: *Bowtie* (BWT+FM-index)



277

Bowtie2: Extending Burrows-Wheeler-based read alignment to longer reads and gapped alignments

Ben Langmead^{1,2}, Mihai Pop¹, Rafael A. Irizarry² and Steven L. Salzberg¹

¹ Center for Bioinformatics and Computational Biology, University of Maryland, College Park, MD, USA, ² Johns Hopkins Bloomberg School of Public Health, Department of Biostatistics, Baltimore, MD, 21205

Website: <http://bowtie.cbcb.umd.edu>, mailing list: <https://lists.sourceforge.net/lists/listinfo/bowtie-bio-announce>

Since its release in 2009, the Bowtie [1] short read aligner has been widely used (50,000 downloads) and cited (hundreds of citations, over 5000 papers). While it was designed for aligning short sequencing reads to 35 bp genomes, the Burrows-Wheeler index allows Bowtie to align more than 25 million reads per CPU hour with a memory footprint of approximately 1.3 gigabytes. Bowtie also provides Burrows-Wheeler transform (BWT)-based quality-aware bowtie2 and bowtie2r programs that can process multiple parallel bowtie processes. Multiple processor cores can be used to align reads in parallel to achieve greater alignment speeds. Bowtie is open source (<http://bowtie.cbcb.umd.edu>).

Abstract

Bowtie is an ultrafast, memory-efficient aligner program for aligning short DNA sequence reads to large genomes. For the human genome, Burrows-Wheeler indexing allows Bowtie to align more than 25 million reads per CPU hour with a memory footprint of approximately 1.3 gigabytes. Bowtie also provides Burrows-Wheeler transform (BWT)-based quality-aware bowtie2 and bowtie2r programs that can process multiple parallel bowtie processes. Multiple processor cores can be used to align reads in parallel to achieve greater alignment speeds. Bowtie is open source (<http://bowtie.cbcb.umd.edu>).

Figure 2 illustrates the hybrid approach taken by Bowtie 2. In Bowtie 1, the entire alignment problem is solved "in Burrows-Wheeler space," using queries to the Burrows-Wheeler index. In Bowtie 2, alignment labor is divided between a Burrows-Wheeler alignment component (which handles all aligned reads) and a dynamic programming component (which handles unaligned reads). A key point is that the Burrow-Wheeler approach is playing to its strengths: whereas dynamic programming is flexible about gaps and affine gap penalties, dynamic programming is much less flexible about short unaligned reads.

Seeds are extracted from various points along the read and its reverse complement according to a configurable policy; a typical policy is to use seeds of length 16 (e.g., 28) every N positions (e.g., 14), where the user-specified parameter N is typically 1000. Once seeds are aligned by the Burrows-Wheeler aligner, alignments are passed to a dynamic programming step. This step samples from among the seed alignments to find anchors for dynamic programming problems. These anchor problems are solved by dynamic programming in the surrounding region of the reference, with padding included to allow for gaps. The dynamic programming problem can be forced to align the entire read end-to-end, or can align it locally.

Figure 1 shows the Burrows-Wheeler matrix of T and the Burrows-Wheeler matrix of reverse(T). The two matrices are identical except for the first column, which is the reverse of the last column. The Burrows-Wheeler transform (BWT) saves time and space by rapidly converting between backward moves in the forward index and forward moves in the backward index, or vice versa.

Figure 2 shows the architecture of Bowtie 2. It consists of three main components: a BWT search module, a dynamic programming module, and a BWT index. The BWT search module takes a query sequence and finds its position in the BWT index. The dynamic programming module then uses this information to find the best alignment. The BWT index is a compressed representation of the genome, allowing for fast and efficient search.

Figure 3 shows how Bowtie 2 decides when to look for discordant and unpaired mate alignments given paired-end reads. It shows four cases: (a) both reads have found concordant pairs, (b) one read has found a discordant pair and the other has found a concordant pair, (c) one read has found a discordant pair and the other has found no pairs, and (d) neither read has found any pairs.

Figure 4 shows the relative performance of Bowtie 2, BWA, and SOAP. The x-axis is the number of reads (from 10^4 to 10^7) and the y-axis is the time taken in seconds. BWT aligns 100M reads in approximately 10 seconds, while BWA and SOAP take significantly longer, up to 300 seconds.

278

139

BWT Aligner: BWA (BWT+FM-index...)

> *Bioinformatics*. 2009 Jul 15;25(14):1754-60. doi: 10.1093/bioinformatics/btp324.
Epub 2009 May 18.

Fast and accurate short read alignment with Burrows-Wheeler transform

Heng Li¹, Richard Durbin
Affiliations + expand
PMID: 19451168 PMCID: PMC2705234 DOI: 10.1093/bioinformatics/btp324
[Free PMC article](#)

Abstract
Motivation: The enormous amount of short reads generated by the new DNA sequencing technologies call for the development of fast and accurate read alignment programs. A first generation of hash table-based methods has been developed, including MAQ, which is accurate, feature rich and fast enough to align short reads from a single individual. However, MAQ does not support gapped alignment for single-end reads, which makes it unsuitable for alignment of longer reads where indels may occur frequently. The speed of MAQ is also a concern when the alignment is scaled up to the resequencing of hundreds of individuals.
Results: We implemented Burrows-Wheeler Alignment tool (BWA), a new read alignment package that is based on backward search with Burrows-Wheeler Transform (BWT), to efficiently align short sequencing reads against a large reference sequence such as the human genome, allowing mismatches and gaps. BWA supports both base space reads, e.g. from Illumina sequencing machines, and color space reads from AB SOLiD machines. Evaluations on both simulated and real data suggest that BWA is approximately 10-20x faster than MAQ, while achieving similar accuracy. In addition, BWA outputs alignment in the new standard SAM (Sequence Alignment/Map) format. Variant calling and other downstream analyses after the alignment can be achieved with the open source SAMtools software package.
Availability: <http://maq.sourceforge.net>.

279

BWA VS. Bowtie

Hatem et al. *BMC Bioinformatics* 2013, **14**:184
<http://www.biomedcentral.com/1471-2105/14/184>

RESEARCH ARTICLE **Open Access**

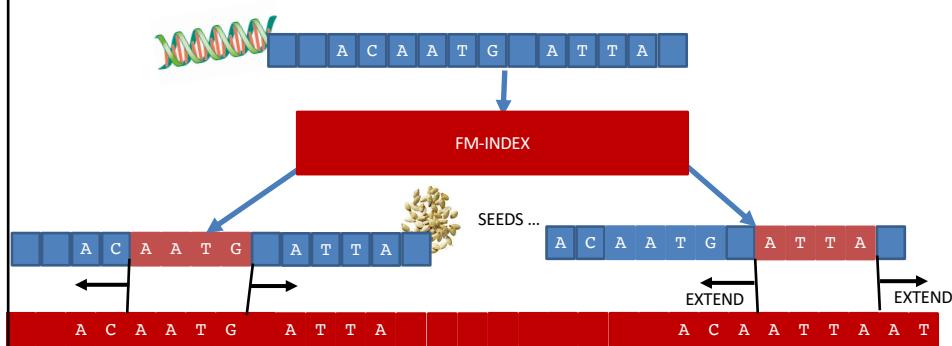
Benchmarking short sequence mapping tools

Ayat Hatem^{1,2}, Doruk Bozdağ², Amanda E Toland³ and Ümit V Çatalyürek^{1,2*}

Abstract
Background: The development of next-generation sequencing instruments has led to the generation of millions of short sequences in a single run. The process of aligning these reads to a reference genome is time consuming and demands the development of fast and accurate alignment tools. However, the current proposed tools make different compromises between the accuracy and the speed of mapping. Moreover, many important aspects are overlooked while comparing the performance of a newly developed tool to the state of the art. Therefore, there is a need for an objective evaluation method that covers all the aspects. In this work, we introduce a benchmarking suite to extensively analyze tools with respect to various aspects and provide an objective comparison.
Results: We applied our benchmarking tests on 9 well known mapping tools, namely, Bowtie, Bowtie2, BWA, SOAP2, MAQ, RMAP, GSNAF, Novoalign, and msrFAST (msrFAST) using synthetic data and real RNA-Seq data. MAQ and RMAP are based on building hash tables for the reads, whereas the remaining tools are based on indexing the reference genome. The benchmarking tests reveal the strengths and weaknesses of each tool. The results show that no single tool outperforms all others in all metrics. However, Bowtie maintained the best throughput for most of the tests while BWA performed better for longer read lengths. The benchmarking tests are not restricted to the mentioned tools and can be further applied to others.
Conclusion: The mapping process is still a hard problem that is affected by many factors. In this work, we provided a benchmarking suite that reveals and evaluates the different factors affecting the mapping process. Still, there is no tool that outperforms all of the others in all the tests. Therefore, the end user should clearly specify his needs in order to choose the tool that provides the best results.
Keywords: Short sequence mapping, Next-generation sequencing, Benchmark, Sequence analysis

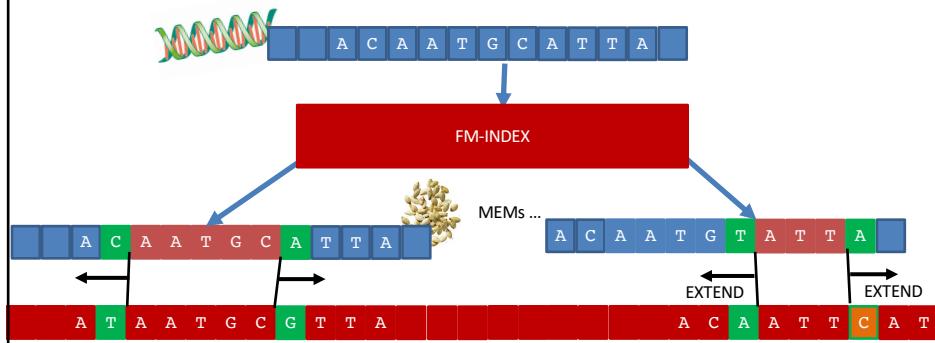
280

Use of FM-index in Read Alignment



281

MEMs Make Good Seeds



Given a text $S[1..n]$ and a pattern $P[1..m]$. The substring $P[i .. i + \ell - 1]$ of length ℓ is a maximal exact match (MEM) of P in S if: $P[i .. i + \ell - 1]$ occurs in S and $P[i - 1..i + \ell - 1]$ and $P[i .. i + \ell]$ does not occur in S

282

Received: 20 December 2021 | Revised: 24 January 2022 | Accepted: 24 January 2022
DOI: 10.1002/im2.4

iMeta WILEY

Applications of de Bruijn graphs in microbiome research

Keith Dufault-Thompson | Xiaofang Jiang

Intramural Research Program, National Library of Medicine, National Institutes of Health, Bethesda, Maryland, USA

Correspondence
Xiaofang Jiang, Intramural Research Program, National Library of Medicine, National Institutes of Health, Building 38A, Room 6N607, 8600 Rockville Pike, Bethesda, MD 20894, USA.
Email: xiaofang.jiang@nih.gov

Funding information
Intramural Research Program of the NIH, National Library of Medicine

Abstract
High-throughput sequencing has become an increasingly central component of microbiome research. The development of de Bruijn graph-based methods for assembling high-throughput sequencing data has been an important part of the broader adoption of sequencing as part of biological studies. Recent advances in the construction and representation of de Bruijn graphs have led to new approaches that utilize the de Bruijn graph data structure to aid in different biological analyses. One type of application of these methods has been in alternative approaches to the assembly of sequencing data like targeted assembly, where only gene sequences are assembled out of larger metagenomes, and differential assembly, where sequences that are differentially present between two samples are assembled. de Bruijn graphs have also been applied for comparative genomics where they can be used to represent large sets of multiple genomes or metagenomes where structural features in the graphs can be used to identify variants, indels, and homologous regions in sequences. These de Bruijn graph-based representations of sequencing data have even begun to be applied to whole sequencing databases for large-scale searches and experiment discovery. de Bruijn graphs have played a central role in how high-throughput sequencing data is worked with, and the rapid development of new tools that rely on these data structures suggests that they will continue to play an important role in biology in the future.

KEY WORDS
de Bruijn graphs, microbiome, Omics

Highlights

- de Bruijn graph-based sequence assembly approaches have been an essential part of the broad application of sequencing methods, especially in microbiome research.
- de Bruijn graphs can be used to efficiently represent sequencing data in a format that is highly scalable and can be extended and modified to address different research questions.

283

Uwe Baier*, Thomas Büchler*, Enno Ohlebusch*, Pascal Weber*

*Institute of Theoretical Computer Science
Ulm University, D-89069 Ulm, Germany

{uwe.baier, thomas.buechler, enno.ohlebusch, pascal-1.weber}@uni-ulm.de

Abstract

This paper introduces the de Bruijn graph edge minimization problem, which is related to the compression of de Bruijn graphs: find the order- k de Bruijn graph with minimum edge count among all orders. We describe an efficient algorithm that solves this problem. Since the edge minimization problem is connected to the BWT compression technique called “tunneling”, the paper also describes a way to minimize the length of a tunneled BWT in such a way that useful properties for sequence analysis are preserved. Although being a restriction, this is significant progress towards a solution to the open problem of finding optimal disjoint blocks that minimize space, as stated in Alanko et al. (DCC 2019).

Introduction

284

Space-Efficient Merging of Succinct de Bruijn Graphs



Lavinia Egidi^{1[0000-0002-9745-0942]}, Felipe A. Louza^{2[0000-0003-2931-1470]}, and
Giovanni Manzini^{1,3[0000-0002-5047-0196]}

¹ University of Eastern Piedmont, Alessandria, Italy
{lavinia.egidi, giovanni.manzini}@uniupo.it

² Department of Computing and Mathematics, University of São Paulo, Brazil
louza@usp.br

³ IIT CNR, Pisa Italy

Abstract. We propose a new algorithm for merging succinct representations of *de Bruijn* graphs introduced in [Bowe *et al.* WABI 2012]. Our algorithm is based on the lightweight BWT merging approach by Holt and McMillan [Bioinformatics 2014, ACM-BCB 2014]. Our algorithm has the same asymptotic cost of the state of the art tool for the same problem presented by Muggli *et al.* [bioRxiv 2017, Bioinformatics 2019], but it uses less than half of its working space. A novel important feature of our algorithm, not found in any of the existing tools, is that it can compute the *Variable Order* succinct representation of the union graph within the same asymptotic time/space bounds.

285

CALCOLO DEGLI OVERLAP

286

LCP

LCP (Longest Common Prefix) Array è stato introdotto da Myers e Manber nel 1993 come struttura ausiliaria per il Suffix Array.

DEF: dato un testo T (\$-terminato) di n caratteri, LCP Array LCP è un array di n interi da 1 a n , tale che $LCP[i] = \ell$ se e solo se ℓ è la lunghezza del più lungo prefisso comune tra $T[S[i]:]$ e $T[S[i-1]:]$.

$LCP[1] = -1$ by default.

287

LCP

LCP (Longest Common Prefix) Array è stato introdotto da Myers e Manber nel 1993 come struttura ausiliaria per il Suffix Array.

DEF

Kasai et al. (2001)

Array LCP da Suffix Array in tempo lineare

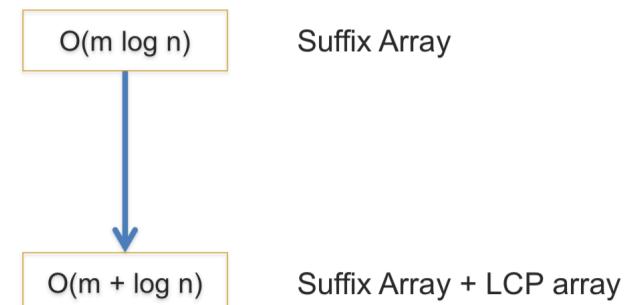
LCP

prefisso comune tra $T[S[i]:]$ e $T[S[i-1]:]$.

$LCP[1] = -1$ by default.

288

Pattern matching



289

Indicizzazione per una collezione di stringhe

- ✓ Suffix Tree, Suffix Array, BWT e LCP Array generalizzati a una collezione di testi

290

Louza et al. *Algorithms Mol Biol.* (2020) 15:18
<https://doi.org/10.1186/s13015-020-00177-y>



Algorithms for
Molecular Biology

SOFTWARE ARTICLE **Open Access**

gsufsort: constructing suffix arrays, LCP arrays and BWTs for string collections

Felipe A. Louza^{1*} , Guilherme P. Telles², Simon Gog³, Nicola Prezza⁴ and Giovanna Rosone^{5*}

Abstract

Background: The construction of a suffix array for a collection of strings is a fundamental task in Bioinformatics and in many other applications that process strings. Related data structures, as the Longest Common Prefix array, the Burrows–Wheeler transform, and the document array, are often needed to accompany the suffix array to efficiently solve a wide variety of problems. While several algorithms have been proposed to construct the suffix array for a single string, less emphasis has been put on algorithms to construct suffix arrays for string collections.

Result: In this paper we introduce **gsufsort**, an open source software for constructing the suffix array and related data indexing structures for a string collection with N symbols in $O(N)$ time. Our tool is written in ANSI/C and is based on the algorithm gSACA-K (Louza et al. in *Theor Comput Sci* 678:22–39, 2017), the fastest algorithm to construct suffix arrays for string collections. The tool supports large fasta, fastq and text files with multiple strings as input. Experiments have shown very good performance on different types of strings.

Conclusions: **gsufsort** is a fast, portable, and lightweight tool for constructing the suffix array and additional data structures for string collections.

Keywords: Suffix array, LCP array, Burrows–Wheeler transform, Document array, String collections

291

Indicizzazione per una collezione di stringhe

- ✓ Suffix Tree, Suffix Array, BWT e LCP Array generalizzati a una collezione di testi
- ✓ Riconoscimento di overlaps all'interno di una collezione di testi (*reads*)



SA+LCP+BWT → overlap in tempo
 lineare

292

Anche la BWT consente di gestire i k-mer (dBG)

Space-efficient construction of succinct de Bruijn graphs

Felipe A. Louza

University of São Paulo, Brazil

Joint work with
Lavinia Egidi and Giovanni Manzini.

293

C. Hohlweg, C. Reutenauer
Lyndon words, permutations and trees
Theor. Comput. Sci., 307 (1) (2003), pp. 173-178

FATTORIZZAZIONI DI LYNDON E SA



294

147

FATTORIZZAZIONI DI LYNDON E BWT



295

FATTORIZZAZIONI DI LYNDON E OVERLAP



296

KFinger: Capturing Overlaps between Long Reads by Using Lyndon Fingerprints*

Paola Bonizzoni^{1[0000-0001-7289-4988]}, Alessia Petescia^{1[0000-0002-9000-3664]},
 Yuri Pirola^{1[0000-0002-8479-7592]}, Raffaella Rizzi^{1[0000-0001-9730-7516]}, Rocco Zaccagnino^{2[0000-0002-9089-5957]}, and Rosalba Zizza^{2[0000-0001-9144-3074]}

¹ Dip. di Informatica, Sistemistica e Comunicazione, University of Milano-Bicocca,
 viale Sarca 336, 20126 Milan, Italy
 paola.bonizzoni@unimib.it, a.petescia@campus.unimib.it,

yuri.pirola@unimib.it, raffaella.rizzi@unimib.it

² Dip. di Informatica, University of Salerno,
 via Giovanni Paolo II 132, 84084 Fisciano, Italy
 {rzaccagnino, rzizza}@unisa.it

Abstract. Detecting common regions and overlaps between DNA sequences is crucial in many Bioinformatics tasks. One of them is genome assembly based on the use of the overlap graph which is constructed by detecting the overlap between genomic reads. When dealing with long reads this task is further complicated by the length of the reads and the high sequencing error rate. This paper proposes a novel alignment-free method for detecting the overlaps in a set of long reads which exploits a signature (called *fingerprint*) of reads built from a factorization of the read based on the notion of Lyndon words. The method has been implemented in the tool **KFinger** and tested over a simulated and a real PacBio HiFi dataset of genomic reads; its results have been compared with the well-known aligner **Minimap2**. **KFinger** is available at <https://github.com/AlgoLab/kfinger>.

297

Can We Replace Reads by Numeric Signatures? Lyndon Fingerprints as Representations of Sequencing Reads for Machine Learning

Paola Bonizzoni^{1(✉)}, Clelia De Felice³, Alessia Petescia¹, Yuri Pirola¹,
 Raffaella Rizzi¹, Jens Stoye², Rocco Zaccagnino³, and Rosalba Zizza³

Abstract. Representations of biological sequences facilitating sequence comparison are crucial in several bioinformatics tasks. Recently, the Lyndon factorization has been proved to preserve common factors in overlapping reads [6], thus leading to the idea of using factorizations of sequences to define measures of similarity between reads. In this paper we propose as a signature of sequencing reads the notion of *fingerprint*, i.e., the sequence of lengths of consecutive factors in Lyndon-based factorizations of the reads. Surprisingly, fingerprints of reads are effective in preserving sequence similarities while providing a compact representation of the read, and so, k -mers extracted from a fingerprint, called k -fingers, can be used to capture sequence similarity between reads.

We first provide a probabilistic framework to estimate the behaviour of fingerprints. Then we experimentally evaluate the effectiveness of this representation for machine learning algorithms for classifying biological sequences. In particular, we considered the problem of assigning RNA-Seq reads to the most likely gene from which they were generated. Our results show that fingerprints can provide an effective machine learning interpretable representation, successfully preserving sequence similarity.

298

149

Numeric Lyndon-based feature embedding of sequencing reads for machine learning approaches

P. Bonizzoni^a, M. Constantini^b, C. De Felice^c, A. Petresci^a, Y. Pirola^a, M. Previtali^a, R. Rizzi^a, J. Stoye^b, R. Zaccagnino^c, R. Zizza^c

Show more ▾

+ Add to Mendeley Share Cite

<https://doi.org/10.1016/j.ins.2022.06.005> Get rights and content

Abstract

Feature embedding methods have been proposed in the literature to represent sequences as numeric vectors to be used in some *bioinformatics* investigations, such as family classification and protein structure prediction.

Recent theoretical results showed that the well-known Lyndon factorization preserves common factors in overlapping strings [1]. Surprisingly, the *fingerprint* of a sequencing read, which is the sequence of lengths of consecutive factors in variants of the Lyndon factorization of the read, is effective in capturing sequence similarities, suggesting it as basis for the definition of novel representations of sequencing reads.

We propose a novel feature embedding method for Next-Generation Sequencing (NGS) data using the notion of *fingerprint*. We provide a theoretical and experimental framework to estimate the behaviour of fingerprints and of the k -mers extracted from it, called k -*fingers*, as possible feature embeddings for sequencing reads. As a case study to assess the effectiveness of such embeddings, we use fingerprints to represent RNA-Seq reads in order to assign them to the most likely gene from which they originated as fragments of transcripts of the gene.

We provide an implementation of the proposed method in the tool *lyn2vec*, which produces Lyndon-based feature embeddings of sequencing reads.

299

Appunti di viaggio (e di progetto)

300

Zielezinski et al. *Genome Biology* (2017) 18:186
DOI 10.1186/s13059-017-1319-7



Genome Biology

REVIEW **Open Access**



Alignment-free sequence comparison:
benefits, applications, and tools

Andrzej Zielezinski¹, Susana Vinga², Jonas Almeida³ and Wojciech M. Karlowski^{1*}

Abstract

Alignment-free sequence analyses have been applied to problems ranging from whole-genome phylogeny to the classification of protein families, identification of horizontally transferred genes, and detection of recombined sequences. The strength of these methods makes them particularly useful for next-generation sequencing data processing and analysis. However, many researchers are unclear about how these methods work, how they compare to alignment-based methods, and what their potential is for use for their research. We address these questions and provide a guide to the currently available alignment-free sequence analysis tools.

301

Read mapping again

We need a *heuristic* for multiple alignment of overlapping reads. (Harder than it seems!)

Read 2 _____

Read 1 _____

Read 4 _____

Read 0 _____

Read 3 _____

Genome _____

302

Question

If we have 30,000 reads (e.g., a bacterial dataset), then we are going to have to perform $(30,000)^2 = 900 \text{ million}$ alignments. This is before we have to work with the overlap network!

Key Point: we can hit the boundary of what a (fast, modern) computer can do very quickly ...

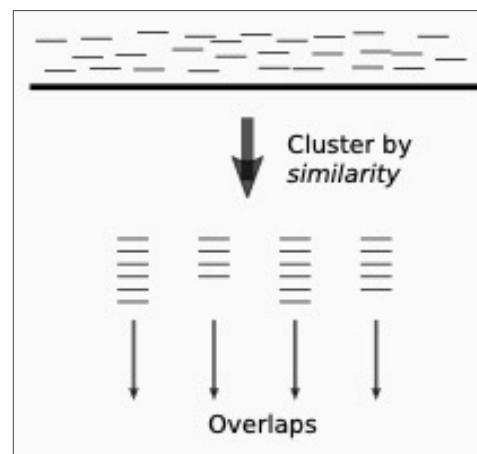
© 2020 by Phillip Compeau

303

Finding similar strings before overlapping

Idea: make a quick decision at the start to find strings that may be similar.

STOP: any ideas on how we could we quickly cluster reads based on similarity?



© 2020 by Phillip Compeau

304

The Brilliant Insight

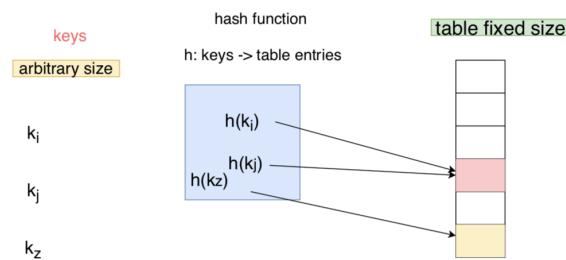
The more two strings *overlap*, the more k -mers that they will share.



© 2020 by Phillip Compeau

305

Hashing: accesso costante al dato



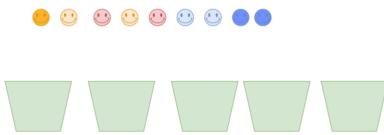
se oggetti simili collidono nella stessa posizione... allora ho un metodo veloce (costante!) per stabilire somiglianza tra oggetti (opportuna collisione diventa strumento positivo)

306

Hashing e alcune misure di somiglianza

- ▶ Idea di processare buckets, contenitori dove elementi finiscono con stesso hash
- ▶ misure di somiglianza: distanza di Jaccard, distanza di Hamming Distance e distanza di edit

Trilioni di buckets e dobbiamo hashare elementi nei buckets



307

Nozioni: Jaccard distance

La **distanza** d or **differenza** tra due stringhe x, y , mentre la somiglianza è data da $1 - d$.

Esempi:

- ▶ **Jaccard distance** tra due insiemi A e B definita come:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$
- ▶ **Hamming distance** per due stringhe di lunghezza n , numero di posizioni in cui differiscono

misura di similarità
tra insiemi (la
applico ai due
insiemi di kmer di
due read per
dedurre l'overlap)

308

Hashing e alcune misure di somiglianza

- ▶ Idea di processare buckets, contenitori dove elementi finiscono con stesso hash
 - ▶ misure di somiglianza: distanza di Jaccard, distanza di Hamming Distance e distanza di edit
- Oggetti che hanno stesso hash hanno alta probabilità di essere simili



309

(localmente sensibile)
Una famiglia di funzioni hash è **locality sensitive** se elementi simili hanno alta probabilità di avere stessa funzione di Hash

Bioinformatics, 35, 2019, i127–i135
doi: 10.1093/bioinformatics/btz354
ISMB/ECCB 2019



Locality-sensitive hashing for the edit distance

Guillaume Marçais*, Dan DeBlasio, Prashant Pandey and Carl Kingsford*

Computational Biology Department, Carnegie Mellon University, Pittsburgh, PA 15213, USA

*To whom correspondence should be addressed.

Abstract

Motivation: Sequence alignment is a central operation in bioinformatics pipelines and, despite many improvements, remains a computationally challenging problem. Locality-sensitive hashing (LSH) is one method used to estimate the likelihood of two sequences to have a proper alignment. Using an LSH, it is possible to separate, with high probability and relatively low computation, the pairs of sequences that do not have high-quality alignment from those that may. Therefore, an LSH reduces the overall computational requirement while not introducing many false negatives (i.e. omitting to report a valid alignment). However, current LSH methods treat sequences as a bag of k -mers and do not take into account the relative ordering of k -mers in sequences. In addition, due to the lack of a practical LSH method for edit distance, in practice, LSH methods for Jaccard similarity or Hamming similarity are used as a proxy.

Results: We present an LSH method, called Order Min Hash (OMH), for the edit distance. This method is a refinement of the minhash LSH used to approximate the Jaccard similarity, in that OMH is sensitive not only to the k -mer contents of the sequences but also to the relative order of the k -mers in the sequences. We present theoretical guarantees of the OMH as a gapped LSH.

Availability and implementation: The code to generate the results is available at <http://github.com/>

310

Local sensitive Hashing

(localmente sensibile)

Una famiglia di funzioni hash è **locality sensitive** se elementi simili hanno alta probabilità di avere stessa funzione di Hash

Definizione

LSH per una relazione $s()$ di somiglianza è una famiglia F di hash functions tale che:

$\Pr_{h \in F}[h(x) = h(y)] = s(x, y)$, per ogni x, y nel dominio di F .

Come stimiamo $s(x, y)$?

311

Local sensitive Hashing: come stimiamo $s(x, y)$?

kmer comuni
→ parte comune

Confronto tra x, y usando lo sketch di x, y :

$\text{sketch}(x) = \langle h_1(x) \dots h_t(x) \rangle$,

$\text{sketch}(y) = \langle h_1(y) \dots h_t(y) \rangle$, usando t hash functions

applicate a x, y , $h_1(\cdot) \dots h_t(\cdot)$.

stima di $s(x, y)$

$s(x, y) = \frac{\#\{h_i(x) = h_i(y)\}}{t}$, per ogni x, y nel dominio di F dove la collisione indica somiglianza

costruzione di F ?

non confronto i k-mer a coppie.
Uso sketch e indicizzazione dei k-mer con tabelle hash

Se le t hash hanno lo stesso valore per x e y , allora x e y sono simili

Obiettivo: definire F in modo tale che il numero di volte due oggetti condividono le collisioni è la somiglianza...

312

Le idee principali di LSH

La procedura LSH è un metodo di riduzione delle dimensioni

Step 1: le sequenze (o parti di esse) vengono riassunte in sketches - cosa è uno sketch? E' una rappresentazione più piccola della sequenza originaria

Step 2: confrontando gli sketches come keys di tabelle hash, il software trova coppie di sequenze che sono simili con alta probabilità

Note: minHash LSH (Broder, 1997) per la Jaccard distance o per Hamming distance come approssimazione della somiglianza.

313

MinHash sketch

definition: LSH for a similarity is a family of hash functions that satisfies the following:

$Pr[h(x) = h(y)] = s(x, y)$

- minHash LSH for the Jaccard distance between two sets A and B is a family of hash functions such that

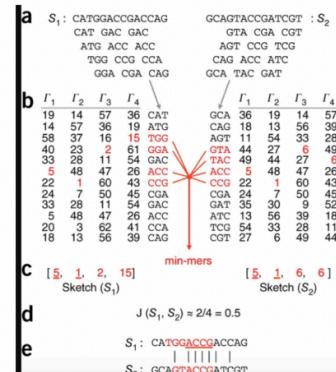
$Pr[h(A) = h(B)] = J(A, B)$, for A and B sets.

314

MinHash sketch: example of applications

creation of a MinHash sketch for a DNA sequence S

- ▶ given $k = 3$, S_1 and S_2 have 12 k-mers
- ▶ k-mers are converted into integer fingerprints
- ▶ 4 hash functions $\Gamma_1, \Gamma_2, \Gamma_3, \Gamma_4$, and for each we choose the minimum, producing the sketch
- ▶ the sketch is smaller in size than the collection of k-mers
- ▶ The fraction of entries shared between the sketches of two sequences S_1 and S_2 (0.5) serves as an estimate of their true Jaccard similarity (0.22)
- ▶ the size of the sketch is used to obtain accurate estimates



questi sketch rappresentano i k-mer
Quindi usiamo minhash per indicizzare i k-mer. Altro metodo sono i Bloom Filter

315

An example of application

Ondov et al. *Genome Biology* (2016) 17:132
DOI 10.1186/s13059-016-0997-x

Genome Biology

SOFTWARE

Open Access



Mash: fast genome and metagenome distance estimation using MinHash

Brian D. Ondov¹, Todd J. Treangen¹, Pál Melsted², Adam B. Mallonee¹, Nicholas H. Bergman¹, Sergey Koren³ and Adam M. Phillippy^{3*}

Abstract

Mash extends the MinHash dimensionality-reduction technique to include a pairwise mutation distance and P value significance test, enabling the efficient clustering and search of massive sequence collections. Mash reduces large sequences and sequence sets to small, representative sketches, from which global mutation distances can be rapidly estimated. We demonstrate several use cases, including the clustering of all 54,118 NCBI RefSeq genomes in 33 CPU h; real-time database search using assembled or unassembled Illumina, Pacific Biosciences, and Oxford Nanopore data; and the scalable clustering of hundreds of metagenomic samples by composition. Mash is freely released under a BSD license (<https://github.com/marbl/mash>).

Keywords: Comparative genomics, Genomic distance, Alignment, Sequencing, Nanopore, Metagenomics

316

Mash: the idea

- ▶ Take genomes that need to be pairwise compared in a fast way.
- ▶ For each genome extracts the k-mers and hash them via single hash function (**canonical k-mers**: the smallest lexicographic one by the forward and the reverse complement)
- ▶ Fix a sketch size s: output the s smallest hashes output over all k-mers
- ▶ Compression: 618 Gbp of genomic sequence and sketches total only 93 MB

317

k-mer counting

(conteggio, rappresentazione con Bloom Filter, Hash, **minhash**, BWT)

318

minimizer

smallest cyclic permutation of a kmer (also w.r.t. reverse and complement)

CGTTGATCAATT TG

K=4;

AATT è il più piccolo?

Ma il suo reverse-complement è AATT.

Scelgo ATTT: ha come reverse-complement è **AAAT**,
smallest cyclic permutation of a kmer (also w.r.t. reverse and complement)

319

minimizer

smallest cyclic permutation of a kmer (also w.r.t. reverse and complement)

CGTTGATCAATT TG

K=4;

AATT è il più piccolo?

Ma il suo reverse-complement è AATT.

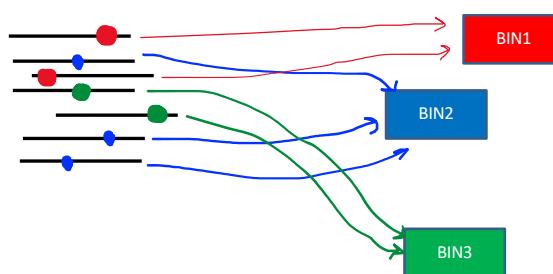
Scelgo ATTT: ha come reverse-complement è **AAAT**,
smallest cyclic permutation of a kmer (also w.r.t. reverse and complement)

Se 2 read sono in
overlap, se il minimizer
cade all'interno, decido
overlap!

320

Uso (in KMC) – kmer counting

Se ho varie read, ho i rispettivi minimizer (etichette di BIN1, BIN2, BIN3), uso i minimizer per tenere traccia dei kmer (andranno nello stesso contenitore)



dato k , posso predeterminare i BIN (considero i kmer possibili in ordine lessicografico) – ottimizzato per farlo on-line...

e poi posso fare il calcolo in parallelo sui BIN

bilanciamento dei BIN; scelta del k ; non va bene per tutti i contesti

321

- Experimental evaluation

Table 4. k -mers counting results for *H.sapiens* 2

Algorithm	$k = 28$			$k = 55$		
	RAM	Disk	Time	RAM	Disk	Time
SSD						
Jellyfish 2	62	0	3212	<i>Out of memory</i>		
KAnalyze	<i>Out of disk (> 650 GB)</i>			<i>Unsupported k</i>		
DSK	6	263	5487	6	256	7732
Turtle	<i>Out of memory</i>			<i>Out of memory</i>		
MSPKC	<i>Out of time (> 10 h)</i>			<i>Out of time (> 10 h)</i>		
KMC 1	17	396	2998	<i>Out of disk (> 650 GB)</i>		
KMC 2 (12 GB)	12	101	1615	13	70	2038
KMC 2 (6 GB)	6	101	1706	13	70	2446
KMC 2 (4 GB)	4	101	1843	4	70	2802
HDD						
Jellyfish 2	62	0	3231	<i>Out of memory</i>		
DSK	6	263	18493	6	256	22432
KMC 1	17	396	4898	<i>Out of disk (> 650 GB)</i>		
KMC 2 (12 GB)	12	101	2259	13	70	2640
KMC 2 (4 GB)	4	101	2707	4	70	3471

Timings in seconds and RAM/disk consumption in GB.

Table 3. k -mers counting results for *M.balsiana*

Algorithm	$k = 28$			$k = 55$		
	RAM	Disk	Time	RAM	Disk	Time
SSD						
Jellyfish 2	17	0	1080	26	0	853
Jellyfish 2-BF	<i>Out of memory</i>			56	30	2865
KAnalyze	9	354	8249	—	—	—
DSK	6	164	2356	6	138	2962
Turtle	46	0	1484	<i>Out of memory</i>		
MSPKC	10	185	8729	<i>Out of time (> 10 h)</i>		
KMC 1	13	165	1229	15	279	2622
KMC 2 (12 GB)	12	41	755	12	29	834
KMC 2 (6 GB)	6	41	685	6	29	895
KMC 2 (4 GB)	4	41	833	4	29	896
KMC 2 (in-mem)	49	0	663	38	0	780
HDD						
Jellyfish 2	17	0	1115	26	0	881
DSK	6	164	6216	6	138	7228
Turtle	46	0	1498	<i>Out of memory</i>		
MSPKC	10	185	12152	<i>Out of time (> 10 h)</i>		
KMC 1	13	165	2194	15	279	3367
KMC 2 (12 GB)	12	41	960	12	29	1041
KMC 2 (4 GB)	4	41	1051	4	29	1207

Timings in seconds and RAM/disk consumption in GB.

322

Using Minimizers for Read Mapping?

We can pre-process the genome by building an index containing the k -mer minimizer of each window for some window length w .



© 2020 by Phillip Compeau

323

Using Minimizers for Read Mapping?

We can pre-process the genome by building an index containing the k -mer minimizer of each window for some window length w .

Then, for each read, we compute its minimizers.



© 2020 by Phillip Compeau

324

Using Minimizers for Read Mapping?

We can pre-process the genome by building an index containing the k -mer minimizer of each window for some window length w .

Then, for each read, we compute its minimizers.



STOP: How will we know if a (long) read aligns well to the genome?

© 2020 by Phillip Compeau

325

Using Minimizers for Read Mapping?

We can pre-process the genome by building an index containing the k -mer minimizer of each window for some window length w .

Then, for each read, we compute its minimizers.



Answer: If it matches multiple minimizers that come from a *small* region of the genome.

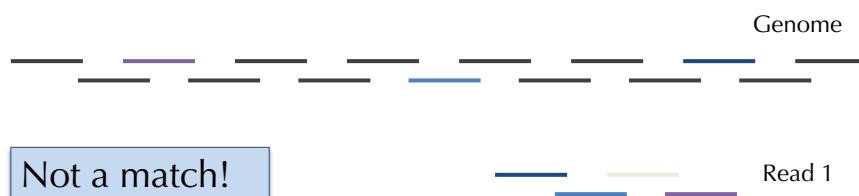
© 2020 by Phillip Compeau

326

Using Minimizers for Read Mapping?

We can pre-process the genome by building an index containing the k -mer minimizer of each window for some window length w .

Then, for each read, we compute its minimizers.



327

Using Minimizers for Read Mapping?

We can pre-process the genome by building an index containing the k -mer minimizer of each window for some window length w .

Then, for each read, we compute its minimizers.



328

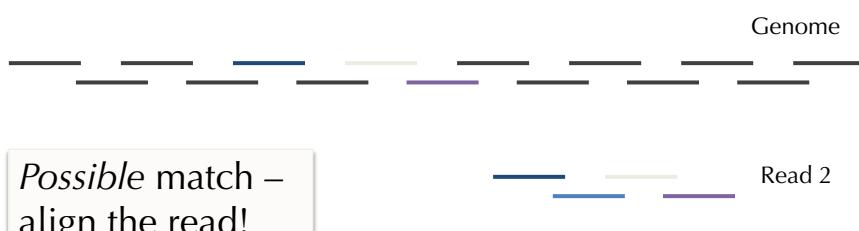
Using Minimizers for Read Mapping?

This is the approach applied by “Minimap”.

www.ncbi.nlm.nih.gov/pmc/articles/PMC4937194

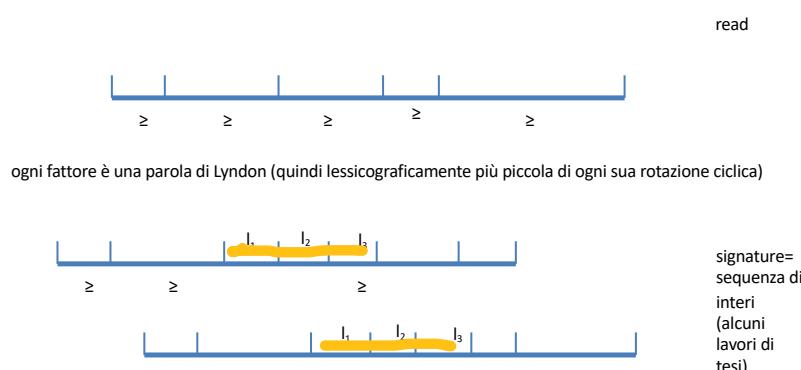
Minimap and miniasm: fast mapping and de novo ... - NCBI - NIH

Mar 19, 2016 — Results: We present a new mapper, **minimap** and a de novo assembler, **miniasm**, ... Notably, our tool **minimap** is a raw read overlapper and **miniasm** is an assembler. We do ... **Minimap** is heavily influenced by all these **works**.
by H Li · 2016 · Cited by 462 · Related articles



329

Signature con Fattorizzazione Lyndon

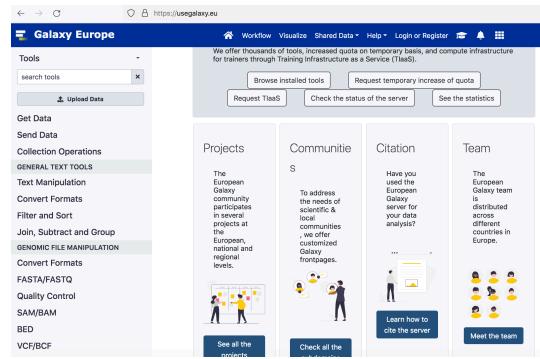


330

Per l'esame

Nei prossimi giorni

- elenco dei progetti con breve descrizione sarà aggiornato
- assignments e training
- scegliete quello che volete sviluppare
- anche in team
- non sparite ☺



331

Per l'esame

Sito

- elenco dei progetti con breve descrizione sarà aggiornato
- assignments e training
- scegliete quello che volete sviluppare
- anche in team
- non sparite ☺

Date d'esame pubblicate:

- ci metteremo d'accordo per fare le presentazioni alla classe, orientativamente in quelle date
- (prenotatevi)
- scegliete qualcosa che vi incuriosisce ...

332



333