

# Finding tandem repeats in genomic data

Projects in Bioinformatics, 10 ECTS  
Astrid Christiansen, 201404423  
Supervised by Thomas Mailund  
Aarhus University

June 2022

# Table of Contents

Introduction

Suffix tree algorithm

Suffix array algorithm

Implementation

Experiments

# Table of Contents

Introduction

Suffix tree algorithm

Suffix array algorithm

Implementation

Experiments

# Introduction

Tandem repeat: *ATCTG ATCTG*.

# Introduction

Tandem repeat: *ATCTG ATCTG*.

- ▶ Find all TRs in string  $x$  of length  $n$ .

# Introduction

Tandem repeat: *ATCTG ATCTG*.

- ▶ Find all TRs in string  $x$  of length  $n$ .
- ▶ Suffix tree- vs. suffix array algorithm.

# Introduction

Tandem repeat: *ATCTG ATCTG*.

- ▶ Find all TRs in string  $x$  of length  $n$ .
- ▶ Suffix tree- vs. suffix array algorithm.
- ▶ Both time  $O(n \log n + z)$ .

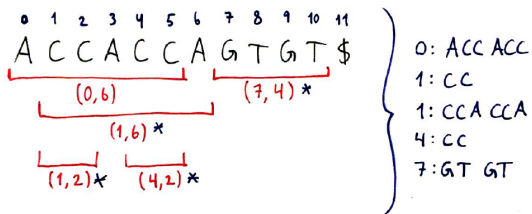
# Tandem repeats and branching tandem repeats

- ▶ Non-branching vs. branching TRs.



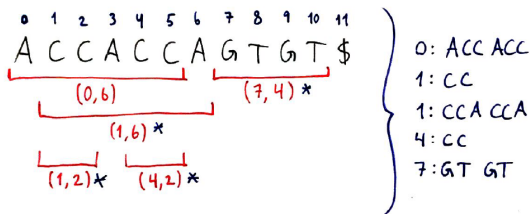
# Tandem repeats and branching tandem repeats

- Non-branching vs. branching TRs.



# Tandem repeats and branching tandem repeats

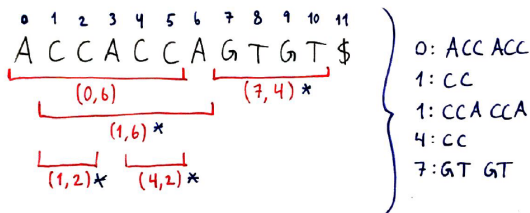
- Non-branching vs. branching TRs.



- Every *non-branching* TR is a *left-rotation* of another TR that starts one place to its right.

# Tandem repeats and branching tandem repeats

- Non-branching vs. branching TRs.



- Every *non-branching* TR is a *left-rotation* of another TR that starts one place to its right.
- Find all TRs from branching ones in time  $O(z)$ .

# Table of Contents

Introduction

Suffix tree algorithm

Suffix array algorithm

Implementation

Experiments

# Suffix tree

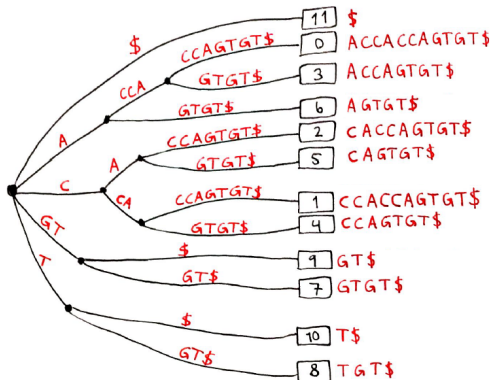
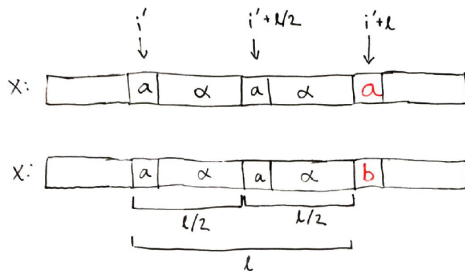
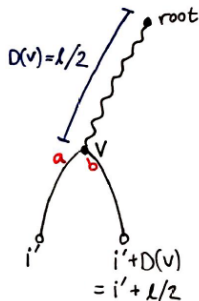
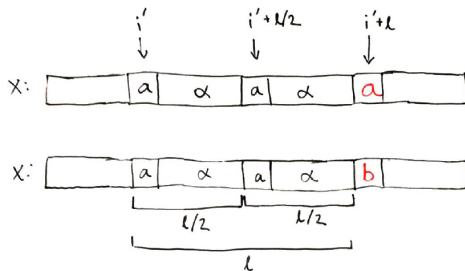


Figure: Suffix tree for string  $x = \text{ACCACCAGTGT}\$$ .

# Finding branching TRs



# Finding branching TRs



# Finding branching TRs

---

## Algorithm BRANCHING-REPEATS( $T$ )

---

```
1: for each inner node,  $v$ , in  $T$  do
2:   for each leaf,  $i'$ , in sub-tree of  $v$  do
3:     if leaf  $i' + D(v)$  is also in sub-tree of  $v$  and  $x[i'] \neq x[i' + 2D(v)]$  then
4:       report ( $i', 2D(v)$ )
5:     end if
6:   end for
7: end for
```

---



# Finding branching TRs

---

## Algorithm BRANCHING-REPEATS( $T$ )

---

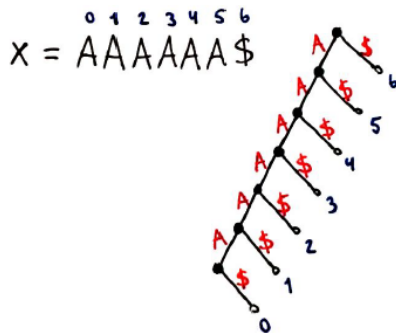
```
1: for each inner node,  $v$ , in  $T$  do
2:   for each leaf,  $i'$ , in sub-tree of  $v$  do
3:     if leaf  $i' + D(v)$  is also in sub-tree of  $v$  and  $x[i'] \neq x[i' + 2D(v)]$  then
4:       report ( $i', 2D(v)$ )
5:     end if
6:   end for
7: end for
```

---

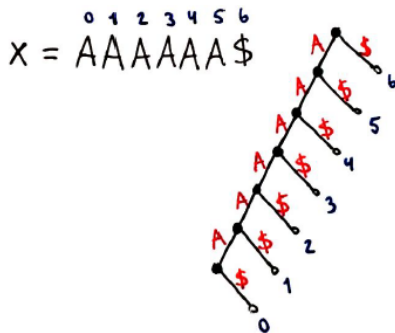
How to determine if leaf is in sub-tree of node  $v$ ?

► DFS numbering  $\Rightarrow$  constant query time.

# Running time



# Running time



Worst case:  $\Theta(n^2)$ .

# The smaller half trick

---

## Algorithm BRANCHING-REPEATS-SMALLER-HALF( $T$ )

---

```
1: for each inner node,  $v$ , in  $T$  do
2:   for each leaf,  $i'$ , in "sub-tree of  $v$  except  $\text{widest}(v)$ " do
3:     if leaf  $i' + D(v)$  is also in sub-tree of  $v$  and  $x[i'] \neq x[i' + 2D(v)]$  then
4:       report ( $i', 2D(v)$ )
5:     end if
6:      $j' \leftarrow i' - D(v)$ 
7:     if leaf  $j'$  is in sub-tree of  $\text{widest}(v)$  and  $x[j'] \neq x[j' + 2D(v)]$  then
8:       report ( $j', 2D(v)$ )
9:     end if
10:   end for
11: end for
```

---

# The smaller half trick

---

## Algorithm BRANCHING-REPEATS-SMALLER-HALF( $T$ )

---

```
1: for each inner node,  $v$ , in  $T$  do
2:   for each leaf,  $i'$ , in "sub-tree of  $v$  except  $\text{widest}(v)$ " do
3:     if leaf  $i' + D(v)$  is also in sub-tree of  $v$  and  $x[i'] \neq x[i' + 2D(v)]$  then
4:       report ( $i', 2D(v)$ )
5:     end if
6:      $j' \leftarrow i' - D(v)$ 
7:     if leaf  $j'$  is in sub-tree of  $\text{widest}(v)$  and  $x[j'] \neq x[j' + 2D(v)]$  then
8:       report ( $j', 2D(v)$ )
9:     end if
10:   end for
11: end for
```

---

Smaller half trick:  $O(n \log n)$ .

# The smaller half trick

---

## Algorithm BRANCHING-REPEATS-SMALLER-HALF( $T$ )

---

```
1: for each inner node,  $v$ , in  $T$  do
2:   for each leaf,  $i'$ , in "sub-tree of  $v$  except widest( $v$ )" do
3:     if leaf  $i' + D(v)$  is also in sub-tree of  $v$  and  $x[i'] \neq x[i' + 2D(v)]$  then
4:       report ( $i', 2D(v)$ )
5:     end if
6:      $j' \leftarrow i' - D(v)$ 
7:     if leaf  $j'$  is in sub-tree of widest( $v$ ) and  $x[j'] \neq x[j' + 2D(v)]$  then
8:       report ( $j', 2D(v)$ )
9:     end if
10:   end for
11: end for
```

---

Smaller half trick:  $O(n \log n)$ .

Total time:  $O(n \log n + z)$ .

# Table of Contents

Introduction

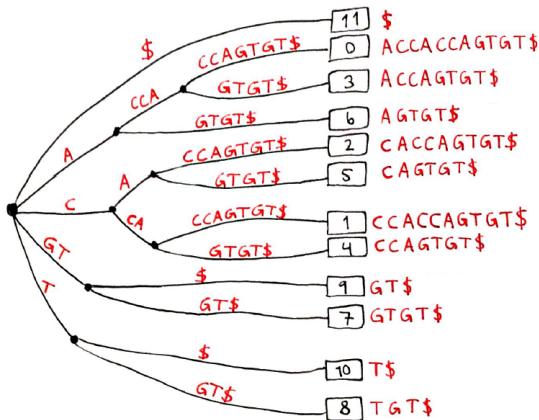
Suffix tree algorithm

Suffix array algorithm

Implementation

Experiments

# Suffix array



SA[0] = 11

SA[1] = 0

SA[2] = 3

SA[3] = 6

SA[4] = 2

SA[5] = 5

SA[6] = 1

SA[7] = 4

SA[8] = 9

SA[9] = 7

SA[10] = 10

SA[11] = 8

ISA[11] = 0

ISA[0] = 1

ISA[3] = 2

ISA[6] = 3

ISA[2] = 4

ISA[5] = 5

ISA[1] = 6

ISA[4] = 7

ISA[9] = 8

ISA[7] = 9

ISA[10] = 10

ISA[8] = 11



# $\ell$ -intervals

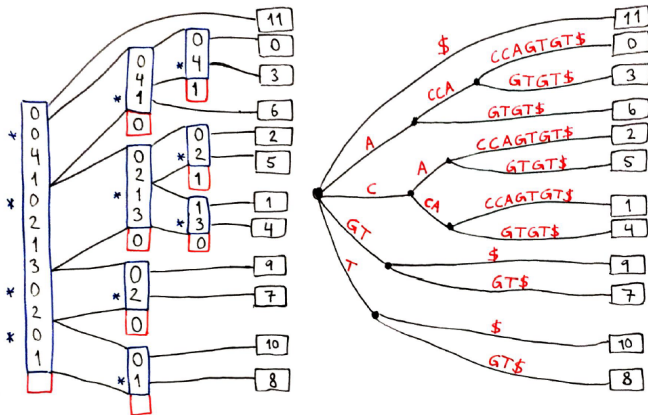
SA:	LCP:
\$	0
ACCACCAGTGT\$	0
ACCAGTGT\$	4
AGTGT\$	1
CACCAGTGT\$	0
CAGTGT\$	2
CCACCAGTGT\$	1
CCAGTGT\$	3
GT\$	0
GTGT\$	2
T\$	0
TGT\$	1

# $\ell$ -intervals

SA:	LCP:
\$	0
ACCACCAGTGT\$	0
ACCAGTGT\$	4
AGTGT\$	1
CACCAGTGT\$	0
CAGTGT\$	2
CCACCAGTGT\$	1
CCAGTGT\$	3
GT\$	0
GTGT\$	2
T\$	0
TGT\$	1

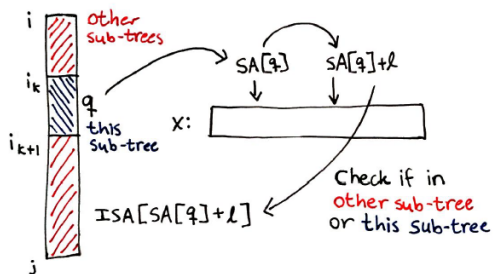
$\ell = \min_{k \in (i,j)} LCP[k]$  and either  
 $i = 0$  or  $LCP[i] < \ell$  and either  
 $j = n$  or  $LCP[j] < \ell$

# $\ell$ -intervals

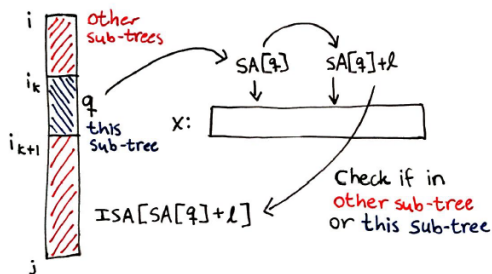


- Use RMQ to find  $\ell$ .

## Find TRs using suffix array

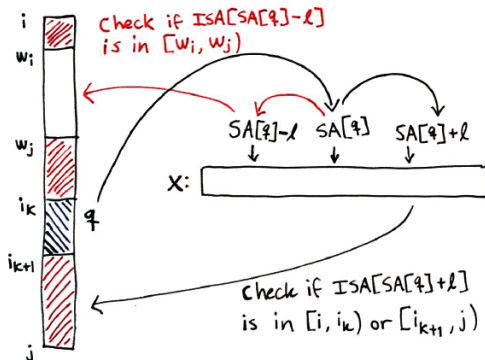


## Find TRs using suffix array

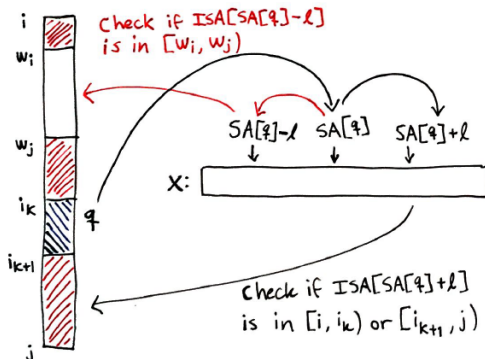


- Again time  $\Theta(n^2) \Rightarrow$  use smaller half trick!

## Find TRs using suffix array - smaller half trick

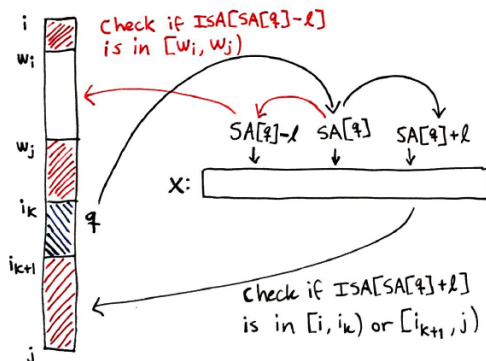


## Find TRs using suffix array - smaller half trick



- Time  $O(n \log n)$ .

# Find TRs using suffix array - smaller half trick



- ▶ Time  $O(n \log n)$ .
- ▶ Total time  $O(n \log n + z)$ .



# Table of Contents

Introduction

Suffix tree algorithm

Suffix array algorithm

**Implementation**

Experiments

# Implementation: Finding branching TRs

```
1 def branching_TR_smaller_half(x, sa, lcp):
2     isa = construct_isa(sa)
3     M = RMQ_preprocess(lcp)
4     for (i, j) in get_inner_nodes(lcp, M, 0, len(x)):
5         child_nodes = list(get_child_nodes(lcp, M, i, j))
6         (w_i, w_j) = widest(child_nodes)
7         (_, L) = RMQ(lcp, M, i + 1, j)
8         for (ii, jj) in child_nodes:
9             if (ii, jj) == (w_i, w_j):
10                 continue
11             for q in valid_isa_index(sa, ii, jj, +L):
12                 r = isa[sa[q] + L]
13                 if (i <= r < j) and not (ii <= r < jj):
14                     yield (sa[q], 2*L)
15             for q in valid_isa_index(sa, ii, jj, -L):
16                 r = isa[sa[q] - L]
17                 if w_i <= r < w_j:
18                     yield (sa[r], 2*L)
```

## Implementation: Finding all TRs

```
1 def find_all_tandem_repeats(x, branching_TRs):
2     for (i, L) in branching_TRs:
3         yield (i, L)
4         while can_rotate(x, i, L):
5             yield (i-1, L)
6             i -= 1
7
8 def can_rotate(x, i, L):
9     return i > 0 and x[i - 1] == x[i + L - 1]
```

# Table of Contents

Introduction

Suffix tree algorithm

Suffix array algorithm

Implementation

Experiments

# Correctness

\*\*\*\*\*

TANDEM REPEATS, (index, length), for string:

x: ACCACCAGTGT\$

(0, 6): ACC ACC

(1, 2): C C

(1, 6): CCA CCA

(4, 2): C C

(7, 4): GT GT

\*\*\*\*\*

\*\*\*\*\*

TANDEM REPEATS, (index, length), for string:

x: aaaaaa\$

(0, 2): a a

(0, 4): aa aa

(0, 6): aaa aaa

(1, 2): a a

(1, 4): aa aa

(2, 2): a a

(2, 4): aa aa

(3, 2): a a

(4, 2): a a

\*\*\*\*\*

# Correctness

\*\*\*\*\*

TANDEM REPEATS, (index, length), for string:

x: ACCACCAGTGT\$

(0, 6): ACC ACC

(1, 2): C C

(1, 6): CCA CCA

(4, 2): C C

(7, 4): GT GT

\*\*\*\*\*

\*\*\*\*\*

TANDEM REPEATS, (index, length), for string:

x: aaaaaa\$

(0, 2): a a

(0, 4): aa aa

(0, 6): aaa aaa

(1, 2): a a

(1, 4): aa aa

(2, 2): a a

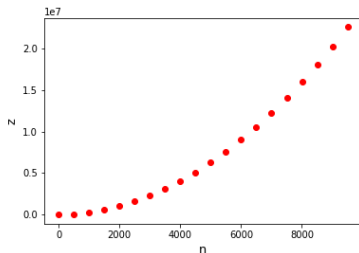
(2, 4): aa aa

(3, 2): a a

(4, 2): a a

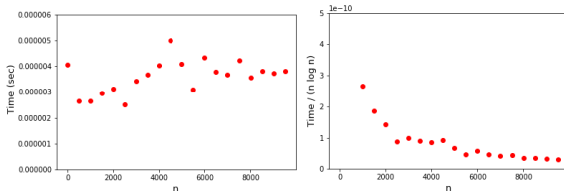
\*\*\*\*\*

Number of found TRs for  $A^n$ :

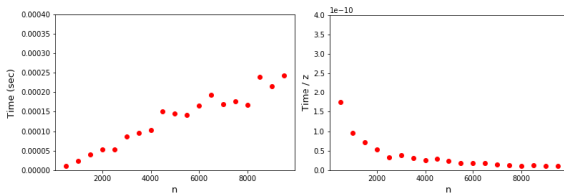


# Running time: Worst case input

Finding branching tandem repeats:

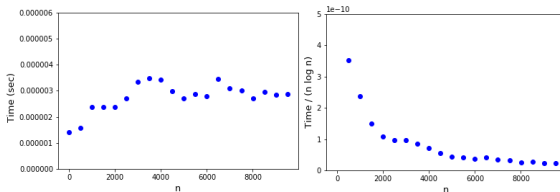


Finding all tandem repeats:

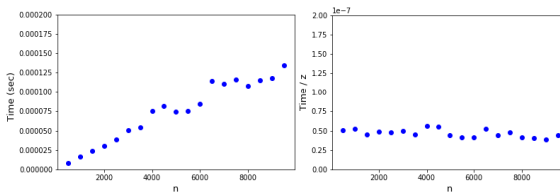


# Running time: Random input

Finding branching tandem repeats:



Finding all tandem repeats:





## Discussion and conclusion

- ▶ TRs can be found using ST or SA.
- ▶ Same time complexity,  $O(n \log n + z)$ .
- ▶ ST better for visualising algorithm.
- ▶ SA more space efficient and simpler data structures.

## Finding all TRs from the branching TRs

Every *non-branching* TR is a *left-rotation* of another TR that starts one place to its right.

# Finding all TRs from the branching TRs

Every *non-branching* TR is a *left-rotation* of another TR that starts one place to its right.

- ▶  $(i', \ell)$  where  $x[i', i' + \ell) = \alpha\alpha = a\beta a\beta$ .

# Finding all TRs from the branching TRs

Every *non-branching* TR is a *left-rotation* of another TR that starts one place to its right.

- ▶  $(i', \ell)$  where  $x[i', i' + \ell) = \alpha\alpha = a\beta a\beta$ .
- ▶ Non-branching  $\Rightarrow x[i', i' + \ell] = a\beta a\beta a$ .

# Finding all TRs from the branching TRs

Every *non-branching* TR is a *left-rotation* of another TR that starts one place to its right.

- ▶  $(i', \ell)$  where  $x[i', i' + \ell) = \alpha\alpha = a\beta a\beta$ .
- ▶ Non-branching  $\Rightarrow x[i', i' + \ell] = a\beta a\beta a$ .
- ▶  $(i', \ell)$  is a left-rotation of  $(i' + 1, \ell)$ .

# Finding all TRs from the branching TRs

Every *non-branching* TR is a *left-rotation* of another TR that starts one place to its right.

- ▶  $(i', \ell)$  where  $x[i', i' + \ell] = \alpha\alpha = a\beta a\beta$ .
- ▶ Non-branching  $\Rightarrow x[i', i' + \ell] = a\beta a\beta a$ .
- ▶  $(i', \ell)$  is a left-rotation of  $(i' + 1, \ell)$ .

We can find all the TRs by repeated left-rotations from the *branching* TRs! Time  $O(z)$ .