# *Whole genome alignment in High Performance Computing environments*

Julio César García Vizcaíno     Directores: Antonio Espinosa, Juan Carlos Moure

Computer Architecture & Operating Systems Department
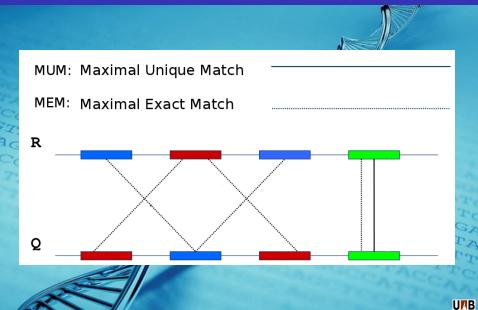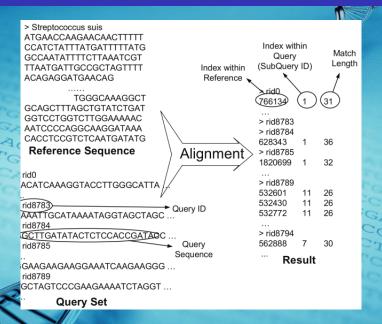Universitat Autónoma de Barcelona

25 de mayo de 2012

# Contents

# Ways of finding exact matches



| Brute Force (3 GB) | Suffix Array (>15 GB) | Suffix Tree (>51 GB) | Hash Table (>15 GB) |
|---|---|---|---|
| BANANA | 6 $ | | |
| BAN | 5 A$ | | |
| ANA | 3 ANA$ | | |
| NAN | 1 ANANA$ | | |
| ANA | 0 BANANA$ | | |
| | 4 NA$ | | |
| | 2 NANA$ | | |
| Naive | Vmatch, PacBio Aligner | MUMmer, MUMmerGPU | BLAST, MAQ, ZOOM, RMAP, CloudBurst |
| Slow & Easy | Binary Search | Tree Searching | Seed-and-extend |

Suffix tree of reference genome

### General objective

Speed up the search of exact matches (distributed) and adapt it to application MUMmer for its execution in HPC environments.
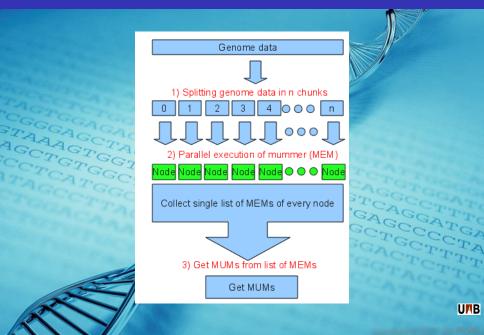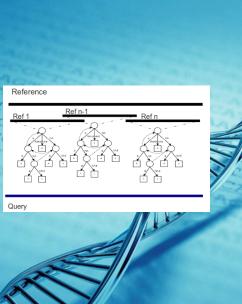
- To have a data structure, efficient usage of memory and processor, that allows a quick search of maximal exact matches.
  - Save relevant information for the search of matches.
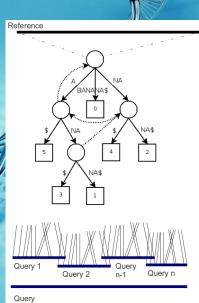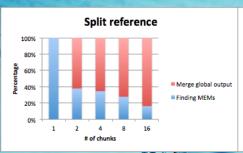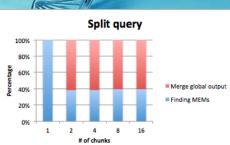  - Be able to nimbly check the data structure.

Genome data

1) Splitting genome data in n chunks

0 1 2 3 4 ○ ○ ○ n

2) Parallel execution of mummer (MEM)

Node Node Node Node Node ○ ○ ○ Node

Collect single list of MEMs of every node

3) Get MUMs from list of MEMs

Get MUMs

Split reference

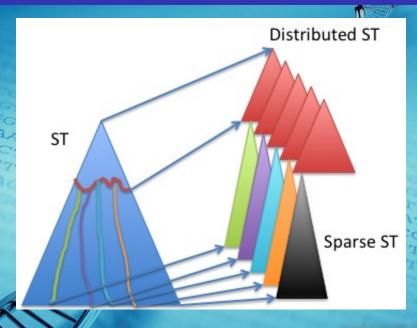Split query

- New variant of the suffix tree.
- Handle of large strings efficiently.
- Based on linear time construction algorithm for subtrees of a suffix tree.
- It tackles the memory bottleneck problem by constructing these subtrees indepently and in parallel.
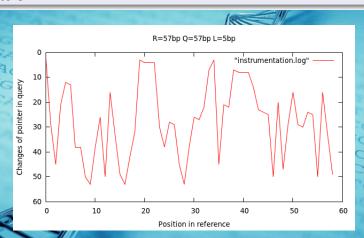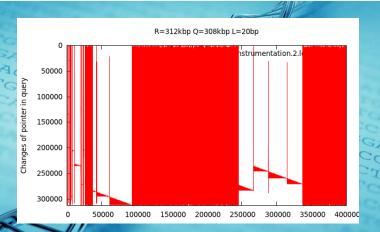
Every suffix of query (pointer) is searched in suffix tree. By using suffix links we jump to other depth of suffix tree and we avoid to check $x$ characters.
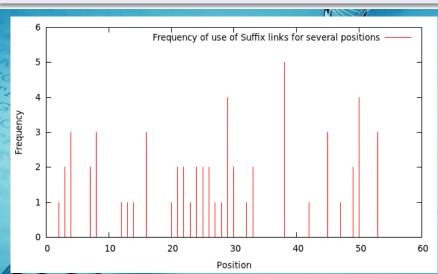


R=57bp Q=57bp L=5bp

The jumps in suffix tree are done while checking the parent node after finishing the last match.



R=312kbp Q=308kbp L=20bp

# Suffix links

The location of suffix links are made during suffix tree construction.



Frequency of use of Suffix links for several positions
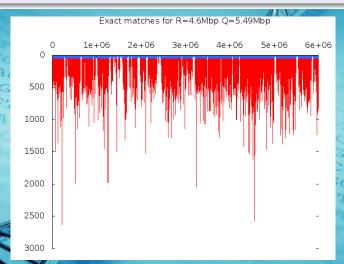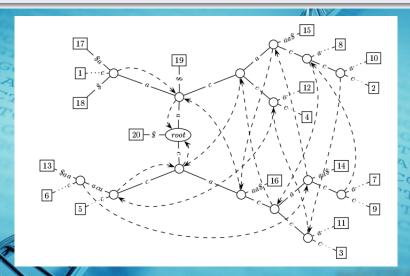
The location of suffix links are more likely to be in deeper regions of suffix tree.

Finding of maximal matches (path from root) are marked in suffix tree. Improved detection of maximal matches.
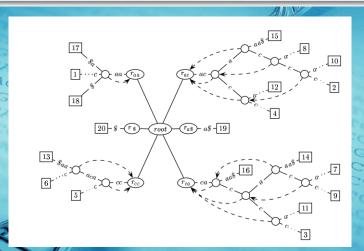


Exact matches for R=4.6Mbp Q=5.49Mbp

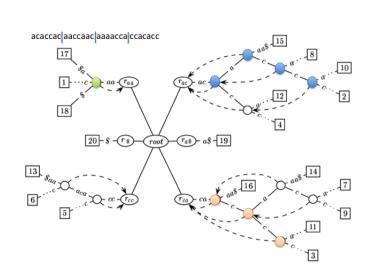Standard suffix tree of aacacccacacaccacaaa$ with standard suffix links.

# Distributed suffix tree

The SSTs for aacacccacacaccacaaa$ with their respective root nodes labelled $r_{aa}$, $r_{ac}$, $r_{ca}$, $r_{cc}$, $r_{a\$}$ and $r_\$$.

## *Conclusions*

### First year

It has been adapted a data structure which can be deployed in HPC environments.

It may be used to implement parallel and distributed techniques for search of maximal matches.

This data structure is able of handling large input sequences to search maximal exact matches.

Perform massive searchs of maximal matches: design and test a parallel and distributed algorithm to perform the search of maximal matches in HPC environments.

Thanks!