

Introduce to Bioinformatics: Distributed Suffix Trees

Instructor: Dr. M. Sakalli
Student: Piotr Jan Dendek

Marmara University, 13.12.2010

Agenda

- ◆ A Suffix Tree – what is it?
- ◆ What is bottleneck of a suffix tree
- ◆ What are workarounds
- ◆ The Distributed Suffix Tree (DST)
- ◆ Conceptual construction demonstration
- ◆ Suffix Tree vs DST
- ◆ Drawbacks
- ◆ Advantages
- ◆ Conclusions

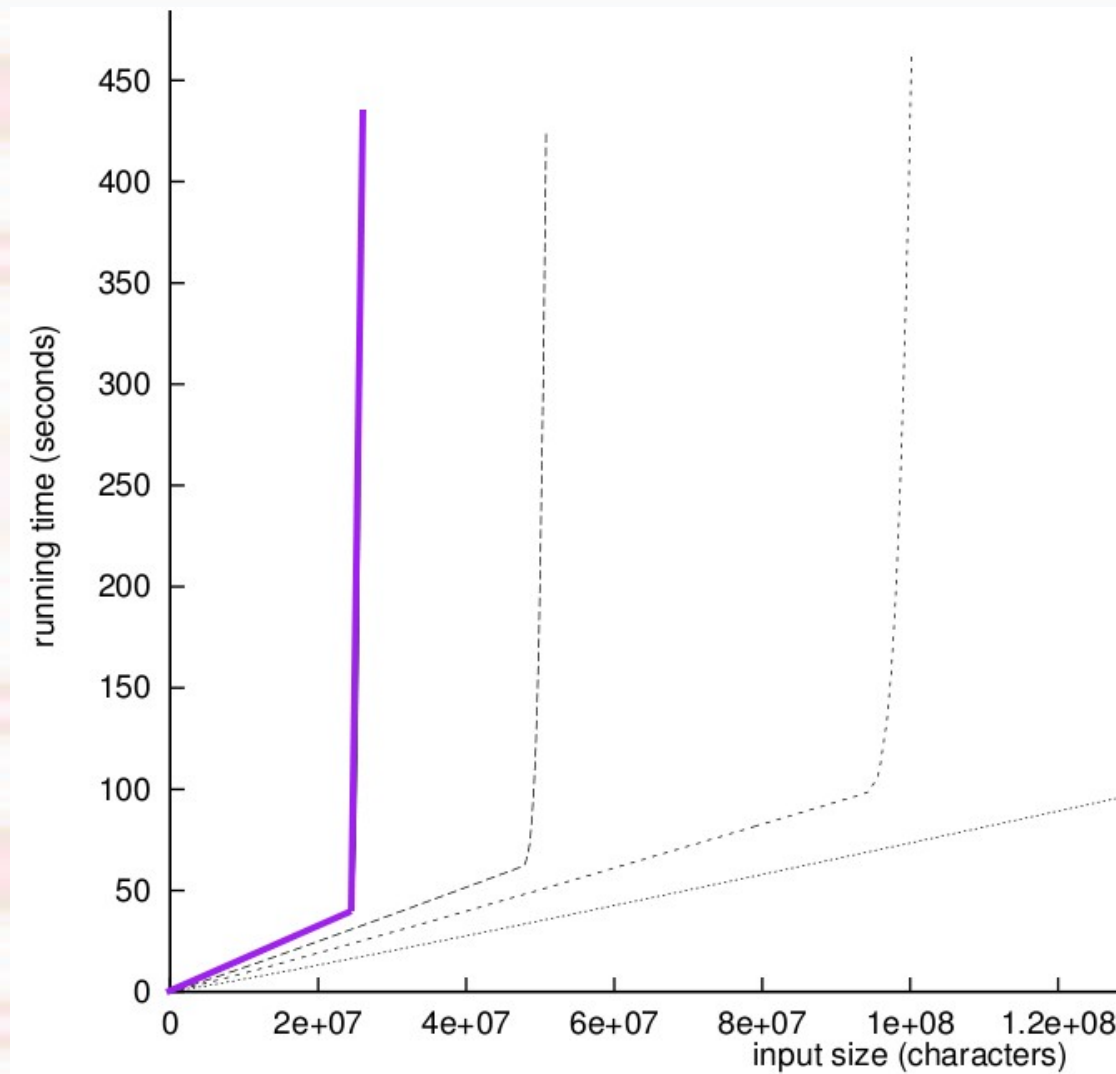
Suffix tree

- ◆ Evolution of Trie
- ◆ Construction time $O(n)$ (where n is length of input string)
- ◆ Good in many application
 - ◆ Check if given word is in a text
 - ◆ Check # of occurrences
 - ◆ Check position of strings occurrences

Bootleneck of a Suffix Tree

In fact construction time is $O(n)$
up to RAM exhaustage.

Bootleneck of a Suffix Tree



Workarounds

- ◆ Do we need workaround?
 - ◆ Yes, e.g. Human Genom Project produced great amount of data on which pople would like to navigate. Those data can not be stored in singular RAM unit
- ◆ Name some workarounds
 - ◆ On-disk Suffix Trees (rather slow)
 - ◆ Distributed Suffix Trees

Tips of a day

- ◆ Suffix Tree gives us ...
- ◆ paths from its root to leafs ...
- ◆ which provides in t (or less then t) hops ...
- ◆ leaf representing t -long word.
- ◆ Especially for
 - ◆ input text,
 - ◆ input text without first letter (suffix with index 1),
 - ◆ input text without first two letter (suffix with index 2),
 - ◆ ... and so on
- ◆ So suffixes of all indexes from 0 to t , where t is input text length, are covered

Distributed Suffix Tree

- ◆ Let's go back to Distributed Suffix Tree!
- ◆ Having a set $\langle V \rangle_{\langle z \rangle}$ and input text $\langle t \rangle$
 - ◆ Where $\langle z \rangle$ is substring of original input string
 - ◆ Where items in $\langle V \rangle_{\langle z \rangle}$ are indexes of concatenations of $\langle z \rangle$ and text after $\langle z \rangle$ in text $\langle t \rangle$
 - ◆ e.g. $\langle t \rangle = \text{"aaabbaabab"}$, $\langle z \rangle = \text{"aa"}$
 - ◆ $\langle t \rangle_{\langle z \rangle} = \text{"12abb6abab"}$
 - ◆ $\langle V \rangle_{\langle z \rangle} = \{1, 2, 6\}$

Distributed Suffix Tree

- ◆ e.g. $\langle t \rangle = \text{"aaabbaabab"}$, $\langle z \rangle = \text{"aa"}$
- ◆ $\langle t \rangle _ \langle z \rangle = \text{"12abb6abab"}$
- ◆ $\langle V \rangle _ \{aa\} = \{1,2,6\}$
- ◆ There are still unused suffixes: 3,4,5,7,8,9,10,11
- ◆ 11 is for empty string better known as „\$”

- ◆ Let's follow the exhaustive example from the beginning
- ◆ Starting letters of used suffixes are uppercase

Distributed Suffix Tree

- ◆ There are still unused suffixes: 1,2,3,4,5,6,7,8,9,10,11
- ◆ $\langle t \rangle = \text{"aaabbaabab"}$, $\langle z \rangle = \text{"aa"}$, $\langle V \rangle_{\{\text{aa}\}} = \{1,2,6\}$
- ◆ There are still unused suffixes: 3,4,5,7,8,9,10,11
- ◆ $\langle t \rangle = \text{"AAabbAabab"}$, $\langle z \rangle = \text{"ab"}$, $\langle V \rangle_{\{\text{ab}\}} = \{3,7,9\}$
- ◆ There are still unused suffixes: 4,5,8,10,11
- ◆ $\langle t \rangle = \text{"AAAbbAAbAb"}$, $\langle z \rangle = \text{"ba"}$, $\langle V \rangle_{\{\text{ba}\}} = \{5,8\}$

Distributed Suffix Tree

- ◆ There are still unused suffixes: 4,5,8,10,11
- ◆ $\langle t \rangle = \text{"AAAbbAAbAb"}$, $\langle z \rangle = \text{"ba"}$, $\langle V \rangle_{\{ba\}} = \{5,8\}$
- ◆ There are still unused suffixes: 4,10,11
- ◆ $\langle t \rangle = \text{"AAAbBAABAb"}$, $\langle z \rangle = \text{"bb"}$, $\langle V \rangle_{\{bb\}} = \{4\}$
- ◆ There are still unused suffixes: 10,11
- ◆ $\langle t \rangle = \text{"AAABBAABAB"}$, $\langle z \rangle = \text{"b\$"}$, $\langle V \rangle_{\{b\$\}} = \{10\}$
- ◆ There are still unused suffixes: 11
- ◆ $\langle t \rangle = \text{"AAABBAABAB"}$, $\langle z \rangle = \text{"\$"}$, $\langle V \rangle_{\{\$\}} = \{11\}$

Distributed Suffix Tree

- ◆ All suffixes of input text are covered!
- ◆ We can divide in similar way Suffix Tree using each $\langle z \rangle$ as subroot
- ◆ Sadly – we must drop suffix tree links between nodes in separate subtrees
- ◆ Luckily – this is not equal with losing of $O(n)$ construction time

Distributed Suffix Tree

- ◆ >> Sadly – we must drop suffix tree links between nodes in separate subtrees<<...
- ◆ ... but how we'll construct a tree then?!
- ◆ With „proper” suffix
- ◆ E.g. „acacc”
 - ◆ Normally there would be suffix tree link between „acac” and „cac” and from „acacc” to „cacc”
 - ◆ But because strings under „ac”-root starts with „ac”, the proper suffix of „acac” is „ac” (instead of using „cac” we are using next suffix)
 - ◆ Similarly, proper suffix of „acacc” is „acc”

Time for a demonstration

- ◆ Let's take input string:
 - ◆ aacacccacacaccacaaa\$
- ◆ Let's point its lower mers (letters and 2-mers)
 - ◆ 1-mer(aacacccacacaccacaaa\$)={a,c,\$}
 - ◆ 2-mer(aacacccacacaccacaaa\$)={aa,ac,ca,cc,a\$}
- ◆ Now point out mers which are not covered by others
 - ◆ a_{1-mer} IN aa_{2-mer} or ac_{2-mer} or a\$_{2-mer}
 - ◆ c_{1-mer} IN cc_{2-mer} or ca_{2-mer} [[[or c\$_{2-mer}]]]
 - ◆ \$_{1-mer} NOT IN 2-mer

Time for a demonstration

- ◆ Now point out mers which are not covered by others
 - ◆ $aa_{\{2\text{-mer}\}}$ $ac_{\{2\text{-mer}\}}$ $a\$_{\{2\text{-mer}\}}$
 - ◆ $cc_{\{2\text{-mer}\}}$ $ca_{\{2\text{-mer}\}}$ $c\$_{\{2\text{-mer}\}}$
 - ◆ $\$_{\{1\text{-mer}\}}$
- ◆ $c\$_{\{2\text{-mer}\}}$ have not occurred but should be mentioned to make division general for alphabet $\{a,c,\$\}$

Time for a demonstration

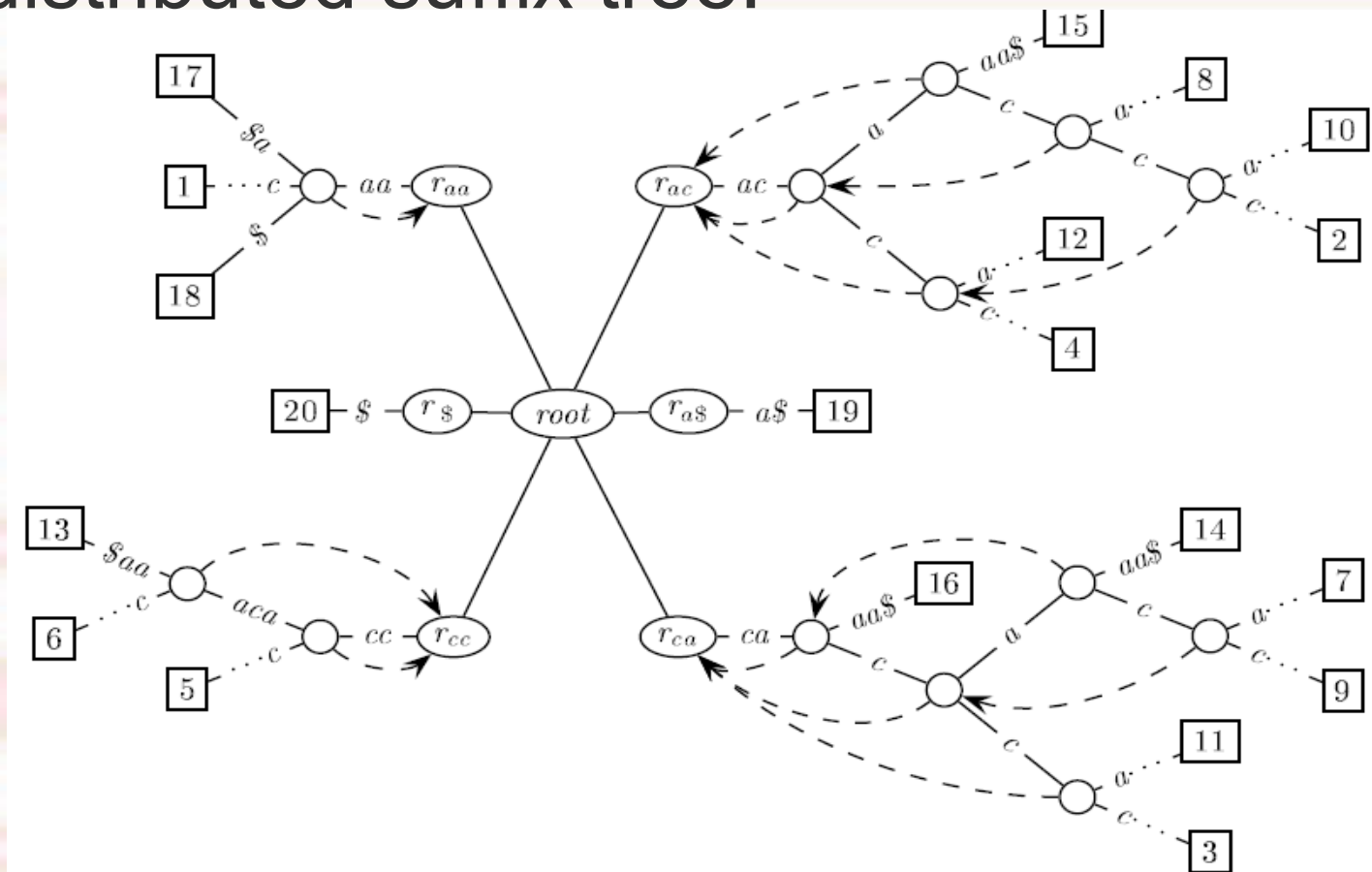
- ◆ Now point out mers which are not covered by others
 - ◆ aa
 - ◆ ac
 - ◆ a\$
 - ◆ cc
 - ◆ ca
 - ◆ c\$
 - ◆ \$

Time for a demonstration

- ◆ Now point out mers which are not covered by others (aa, ac, a\$, cc, ca, c\$, \$)
- ◆ And make them sub-roots in Distributed Suffix Tree
- ◆ If max # of digits in subroot is $\langle N \rangle$ and # of alphabet letters without \$ is $\langle L \rangle$ then # of subroots $\langle S \rangle$ is
- ◆ $\langle S \rangle = \langle N \rangle^{\langle L-1 \rangle} * \langle N+1 \rangle$ | regular roots
+ $\langle N \rangle^{\langle L-2 \rangle} * \langle 1 \rangle + \dots$ | \$-ended roots
+ $\langle N \rangle^{\langle 1 \rangle} * \langle 1 \rangle$ | till $\langle \text{digit} \rangle \$$
+ 1 | and pure \$

Time for a demonstration

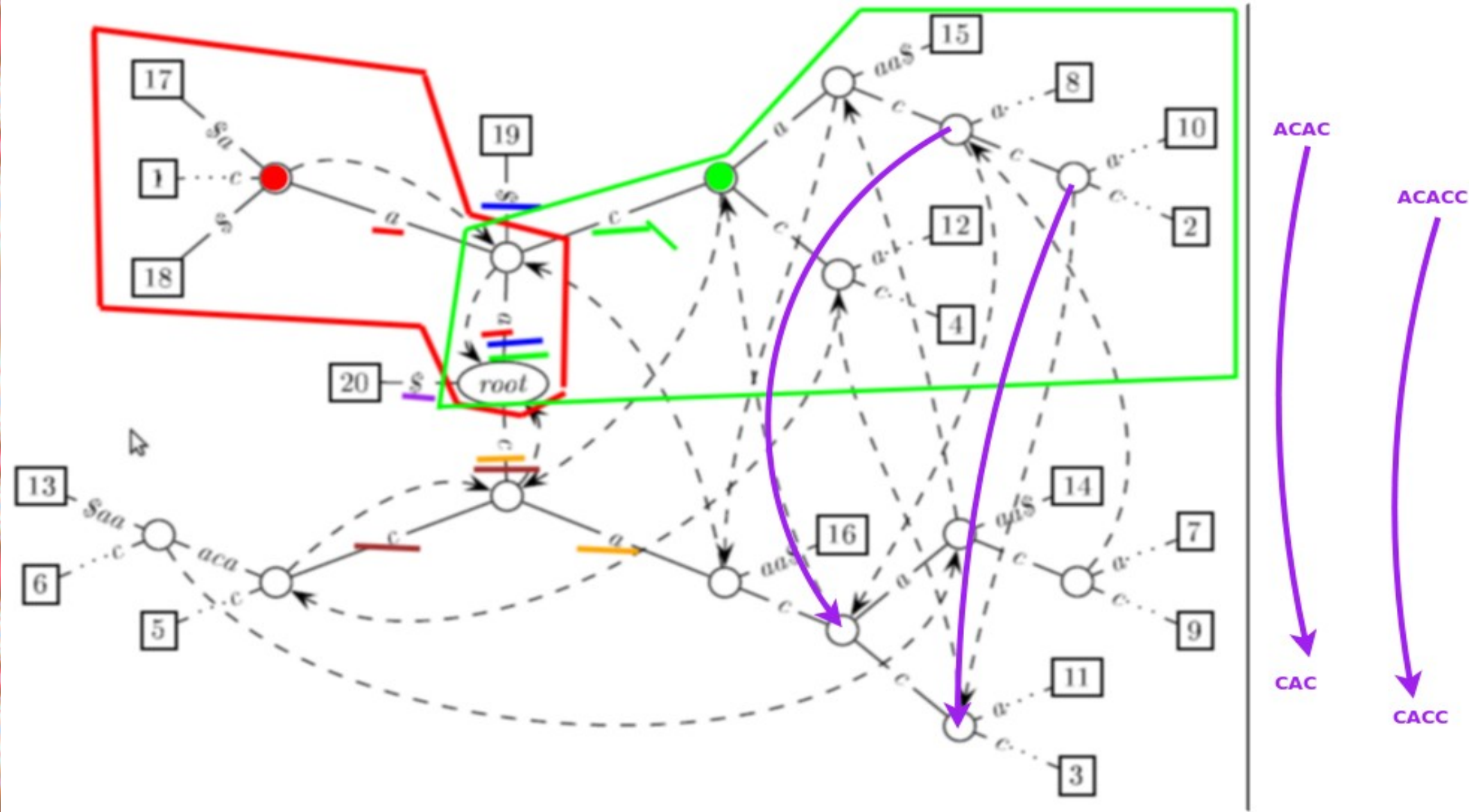
- ◆ the distributed suffix tree:



Normal vs Distributed

- ◆ Clones of nodes in Distributed Suffix Tree (DSF)
- ◆ Deletions of not proper arcs in DSF
- ◆ Look closer!

Look closer: Suffix Tree



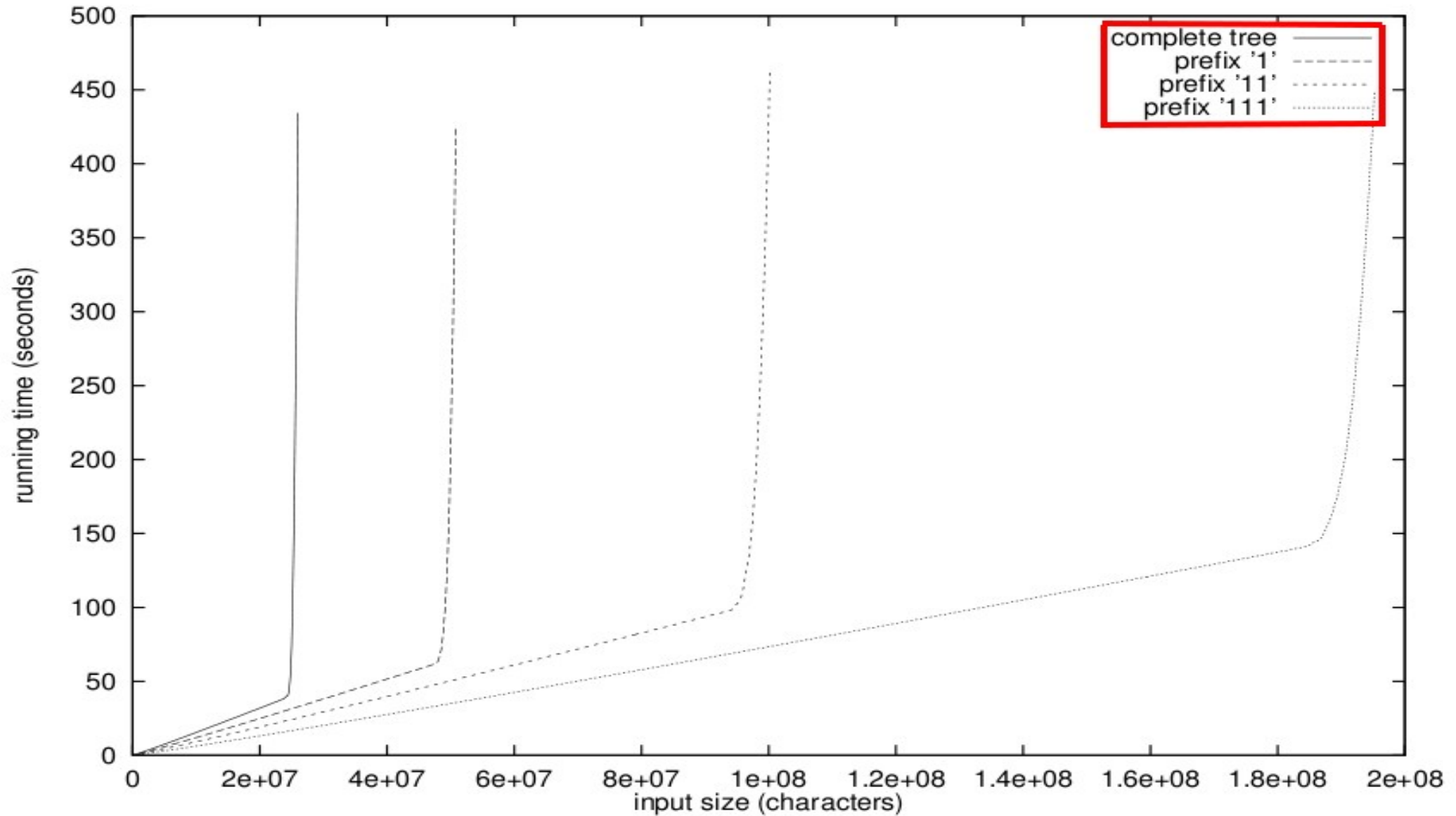
Drawbacks

- ◆ Dependency on hardware connection
- ◆ Wired structure
- ◆ Some statistical operations perform time is much worse

Advantegees

- ◆ Realtime sublinear construction
- ◆ Chance to introduce grid calculations
- ◆ And again the plot from beginning of presentation...

Advantages



Advantegees

- ◆ ... which means that Distributed Suffix Tree allow us to ...
- ◆ ... use more RAM (sum of RAM across all involved computers) ...
- ◆ ... which preserve $O(n)$ for longer time
- ◆ With 3 digit root it is extended to 780% of base input text length
- ◆ You must say „WOW!”

Conclusions

- ◆ Till now topic of Distributed Suffix Trees has not been enough good covered (2 articles by the same authors from Imperial College, UK)
- ◆ An idea of providing multicore implementation of DST may be interesting
- ◆ Using arcs from DST in Suffix Trees concatenated with de Bruijn graph sounds interesting

References

- ◆ „Distributed Suffix Trees”, Raphaël Clifford, King’s College, London, UK
- ◆ „Distributed and Paged Suffix Trees for Large Genetic Databases”, Raphaël Clifford and Marek Sergot, Imperial College London, UK
- ◆ „Distributed Suffix Tree Overlay for Peer-to-Peer Search”, Hai Zhuge and Liang Feng (IEEE Transactions on Knowledge and Data Engineering, vol. 20, no. 2, February 2008)
- ◆ „Distributed Suffix Tree for Peer-to-Peer Search”, Hai Zhuge and Liang Feng, Chinese Academy of Sciences, Beijing, China

Thank you for your attention!

