

ERa — A Practical Approach to Parallel Construction of Suffix Trees

Andrej (Andy) Brodnik^{1,2}, Matevž Jekovec¹

¹ University of Ljubljana, Faculty of Computer and Information Science

² University of Primorska, Department of Information Sciences and Technologies

Dagstuhl, *Data Structures and Advanced Models of Computation on Big Data*, February 24-28, 2014

Text indexing problem

Text indexing problem

Problem statement

Given unstructured input text T consisting of N characters from alphabet Σ of size σ build an index such that for query pattern P we:

- determine whether P occurs in T in time $O(P)$,
- find all occurrences of P in T in time $O(P + occ)$,
- find the longest common prefix (LCP) of P and any suffix of T in time $O(LCP(P, T))$.

Text indexing problem

Problem statement

Given unstructured input text T consisting of N characters from alphabet Σ of size σ build an index such that for query pattern P we:

- determine whether P occurs in T in time $O(P)$,
- find all occurrences of P in T in time $O(P + occ)$,
- find the longest common prefix (LCP) of P and any suffix of T in time $O(LCP(P, T))$.

Solution

Suffix tree and *suffix array* (SA) with LCP information are fundamental data structures for indexing unstructured text.

Suffix tree construction algorithms

- Theoretical:

	W ('73), McC ('78)	U ('95)	F-C et al. ('00)
Work w. c.	$O(N)$	$O(N)$	$O(N \lg N)$
Online	No	Yes	Yes ¹
I/O efficiency	String	String	Result+String
Unbounded Σ	No	No	Yes
Parallel	No	No	PDAM

¹Bedathur and Haritsa (2004)

Suffix tree construction algorithms

- Theoretical:

	W ('73), McC ('78)	U ('95)	F-C et al. ('00)
Work w. c.	$O(N)$	$O(N)$	$O(N \lg N)$
Online	No	Yes	Yes ¹
I/O efficiency	String	String	Result+String
Unbounded Σ	No	No	Yes
Parallel	No	No	PDAM

- Practical:

	Semi-disk-based			Out-of-core		
	TDD ('04)	TRLS. ('07)	B ² ST ('09)	WF ('09)	ERa ('11)	PCF ('13)
Work w. c.	$O(N^2)$	$O(N^2)$	$O(N^2)$	$O(N^2)$	$O(N^2)$	$O(pN \lg N)$
I/O eff.	R.	R.	R.+S.	R.+S.	R.+S.	R.+S.
Unbnd. Σ	No	No	No	No	No	No
Parallel	No	No	No	Yes	Yes	Yes

¹Bedathur and Haritsa (2004)

Suffix tree construction algorithms

- Theoretical:

	W ('73), McC ('78)	U ('95)	F-C et al. ('00)
Work w. c.	$O(N)$	$O(N)$	$O(N \lg N)$
Online	No	Yes	Yes ¹
I/O efficiency	String	String	Result+String
Unbounded Σ	No	No	Yes
Parallel	No	No	PDAM

- Practical:

	Semi-disk-based			Out-of-core		
	TDD ('04)	TRLS. ('07)	B ² ST ('09)	WF ('09)	ERa ('11)	PCF ('13)
Work w. c.	$O(N^2)$	$O(N^2)$	$O(N^2)$	$O(N^2)$	$O(N^2)$	$O(pN \lg N)$
I/O eff.	R.	R.	R.+S.	R.+S.	R.+S.	R.+S.
Unbnd. Σ	No	No	No	No	No	No
Parallel	No	No	No	Yes	Yes	Yes

¹Bedathur and Haritsa (2004)

Suffix tree construction lower bounds

- Sequential:

	bounded Σ	unbounded Σ
Time	$\Omega(\text{Sort}(N))$	$\Omega(\text{Sort}(N))$
I/Os ²	$\Omega(\text{Sort}(N))$	$\Omega(\text{Sort}(N))$
Space ³	$\Omega(N \lg \sigma)$ bits	$\Omega(N \lg \sigma)$ bits

²EM model

³Uncompressed index in word RAM

⁴PEM model

Suffix tree construction lower bounds

- Sequential:

	bounded Σ	unbounded Σ
Time	$\Omega(\text{Sort}(N))$	$\Omega(\text{Sort}(N))$
I/Os ²	$\Omega(\text{Sort}(N))$	$\Omega(\text{Sort}(N))$
Space ³	$\Omega(N \lg \sigma)$ bits	$\Omega(N \lg \sigma)$ bits

- Parallel on p processing units:

	bounded Σ	unbounded Σ
Parallel time	$\Omega(\text{Sort}_p(N))$	$\Omega(\text{Sort}_p(N))$
Parallel I/Os ⁴	$\Omega(\text{Sort}_p(N))$	$\Omega(\text{Sort}_p(N))$
Space ³	$\Omega(N \lg \sigma)$ bits	$\Omega(N \lg \sigma)$ bits

²EM model

³Uncompressed index in word RAM

⁴PEM model

Suffix tree construction lower bounds

- Sequential:

	bounded Σ	unbounded Σ
Time	$\Omega(N)$	$\Omega(N \log N)$
I/Os ²	$\Omega\left(\frac{N}{B}\right)$	$\Omega\left(\frac{N}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$
Space ³	$\Omega(N \lg \sigma)$ bits	$\Omega(N \lg \sigma)$ bits

- Parallel on p processing units:

	bounded Σ	unbounded Σ
Parallel time	$\Omega\left(\frac{N}{p}\right)$	$\Omega\left(\frac{N}{p} \log N\right)$
Parallel I/Os ⁴	$\Omega\left(\frac{N}{pB}\right)$	$\Omega\left(\frac{N}{pB} \log_{\frac{M}{B}} \frac{N}{B}\right)$
Space ³	$\Omega(N \lg \sigma)$ bits	$\Omega(N \lg \sigma)$ bits

²EM model

³Uncompressed index in word RAM

⁴PEM model

Suffix tree construction lower bounds

- Sequential:

	bounded Σ	unbounded Σ
Time	$\Omega(N)$	$\Omega(N \log N)$
I/Os ²	$\Omega\left(\frac{N}{B}\right)$	$\Omega\left(\frac{N}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$
Space ³	$\Omega(N \lg \sigma)$ bits	$\Omega(N \lg \sigma)$ bits

- Parallel on p processing units:

	bounded Σ	unbounded Σ
Parallel time	$\Omega\left(\frac{N}{p}\right)$	$\Omega\left(\frac{N}{p} \log N\right)$
Parallel I/Os ⁴	$\Omega\left(\frac{N}{pB}\right)$	$\Omega\left(\frac{N}{pB} \log_{\frac{M}{B}} \frac{N}{B}\right)$
Space ³	$\Omega(N \lg \sigma)$ bits	$\Omega(N \lg \sigma)$ bits

²EM model

³Uncompressed index in word RAM

⁴PEM model

Theory and Practice

- Substantial gap between the theoretical and practical results.

Theory and Practice

- Substantial gap between the theoretical and practical results.
- Practitioners (often) do not use theoretically the best results.

Theory and Practice

- Substantial gap between the theoretical and practical results.
- Practitioners (often) do not use theoretically the best results.
- Perhaps we should look at practical solutions more carefully.

ERa — Elastic Range (Mansour et al. (2011))

- Currently the fastest practical, parallel suffix tree construction algorithm.

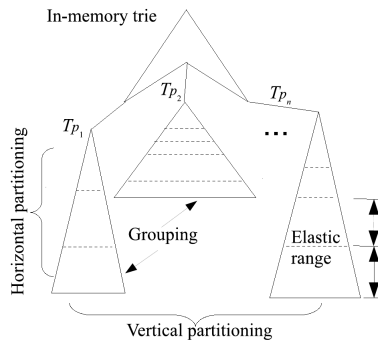
ERa — Elastic Range (Mansour et al. (2011))

- Currently the fastest practical, parallel suffix tree construction algorithm.
- Time complexity: $O(N^2)$ w.c. – for (extremely) skewed text!

ERa — Elastic Range (Mansour et al. (2011))

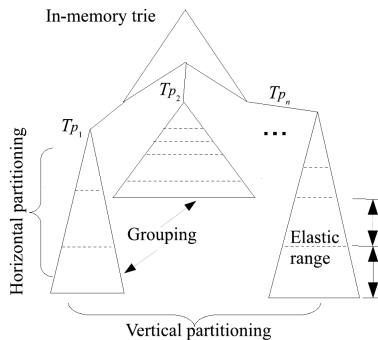
- Currently the fastest practical, parallel suffix tree construction algorithm.
- Time complexity: $O(N^2)$ w.c. – for (extremely) skewed text!
- Yet, it's **fast** in practice: Constructs and stores the human genome's suffix tree in 20 minutes on 16-core desktop PC with HDD or 13 minutes with SSD!

ERa



ERa constructs the suffix tree in two steps:

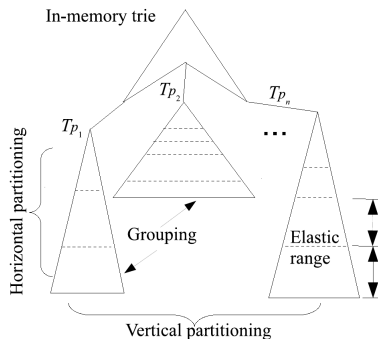
ERa



ERa constructs the suffix tree in two steps:

- 1 The **vertical partitioning** step determines 1) the suffix subtrees just fitting into M and 2) constructs the suffix tree top.

ERa



ERa constructs the suffix tree in two steps:

- 1 The **vertical partitioning** step determines 1) the suffix subtrees just fitting into M and 2) constructs the suffix tree top.
- 2 The **horizontal partitioning** step builds the actual suffix subtrees.

ERa

Algorithm 1: ERa

Input: String S , Alphabet Σ , Processors P , Private cache size M **Output:** Suffix tree \mathcal{T}

```
1  $\mathcal{T}_{top}, G \leftarrow \text{VerticalPartitioning}(S, \Sigma, M)$ 
2  $\mathcal{T} \leftarrow \mathcal{T}_{top}$ 
3 while  $|G| > 0$  do
4   for  $p \in P$  do in parallel
5     if  $|G| > 0$  then
6        $\pi \leftarrow G.\text{pop}()$ 
7        $\mathcal{T}_\pi \leftarrow \text{HorizontalPartitioning}(S, \Sigma, \pi)$ 
8        $\text{Link}(\mathcal{T}, \mathcal{T}_\pi)$ 
9     end
10  end
11 end
12 return  $\mathcal{T}$ 
```

Vertical partitioning

Define **S-prefix** π as the prefix of the suffix in the text;
 $f_\pi = \#$ of suffixes starting with π ; and \mathcal{T}_π is a subtree with a root
corresponding π .

Vertical partitioning

Define **S-prefix** π as the prefix of the suffix in the text;
 $f_\pi = \#$ of suffixes starting with π ; and \mathcal{T}_π is a subtree with a root
corresponding π .

Idea: Approximate the size of \mathcal{T}_π as cf_π for some constant c and
expand the π so much, that $cf_\pi \leq M$.

Vertical partitioning

Define **S-prefix** π as the prefix of the suffix in the text;
 $f_\pi = \#$ of suffixes starting with π ; and \mathcal{T}_π is a subtree with a root
corresponding π .

Idea: Approximate the size of \mathcal{T}_π as cf_π for some constant c and
expand the π so much, that $cf_\pi \leq M$.

- 1 Scan the text and obtain the characters frequency $f_\pi : \pi \in \Sigma$
(counted all S-prefixes of length 1)

Vertical partitioning

Define **S-prefix** π as the prefix of the suffix in the text;
 $f_\pi = \#$ of suffixes starting with π ; and \mathcal{T}_π is a subtree with a root
corresponding π .

Idea: Approximate the size of \mathcal{T}_π as cf_π for some constant c and
expand the π so much, that $cf_\pi \leq M$.

- 1 Scan the text and obtain the characters frequency $f_\pi : \pi \in \Sigma$
(counted all S-prefixes of length 1)
- 2 For each $\pi : cf_\pi > M$, $\pi' = \pi\Sigma$ and count $f_{\pi'}$.

Vertical partitioning

Define **S-prefix** π as the prefix of the suffix in the text;
 $f_\pi = \#$ of suffixes starting with π ; and \mathcal{T}_π is a subtree with a root
corresponding π .

Idea: Approximate the size of \mathcal{T}_π as cf_π for some constant c and
expand the π so much, that $cf_\pi \leq M$.

- 1 Scan the text and obtain the characters frequency $f_\pi : \pi \in \Sigma$
(counted all S-prefixes of length 1)
- 2 For each $\pi : cf_\pi > M$, $\pi' = \pi\Sigma$ and count $f_{\pi'}$.
- 3 Repeat step two for S-prefixes of length 3, 4, ..., until all \mathcal{T}_π
just fits into the memory M .

Vertical partitioning

Define **S-prefix** π as the prefix of the suffix in the text;
 $f_\pi = \#$ of suffixes starting with π ; and \mathcal{T}_π is a subtree with a root
corresponding π .

Idea: Approximate the size of \mathcal{T}_π as cf_π for some constant c and
expand the π so much, that $cf_\pi \leq M$.

- 1 Scan the text and obtain the characters frequency $f_\pi : \pi \in \Sigma$
(counted all S-prefixes of length 1)
- 2 For each $\pi : cf_\pi > M$, $\pi' = \pi\Sigma$ and count $f_{\pi'}$.
- 3 Repeat step two for S-prefixes of length 3, 4, ..., until all \mathcal{T}_π
just fits into the memory M .

Extra : To optimally fill the main memory, combine the S-prefixes
into *virtual groups* G , fitting into the main memory as tight as
possible (use First-Fit Decreasing heuristic for bin packing
problem)

$$\pi = ACC, f_{ACC} = 12$$

Horizontal partitioning

For a prefix π (virtual group G) construct the suffix subtree:

Horizontal partitioning

For a prefix π (virtual group G) construct the suffix subtree:

- 1 Locate and store all $n = f_\pi$ positions of π – we will sort n strings in-memory; i.e. all strings are ambiguous.

Horizontal partitioning

For a prefix π (virtual group G) construct the suffix subtree:

- 1 Locate and store all $n = f_\pi$ positions of π – we will sort n strings in-memory; i.e. all strings are ambiguous.
- 2 Optimal prefix length range of chars following π $r = \rho \frac{M}{n}$ – *Elastic Range*.

Horizontal partitioning

For a prefix π (virtual group G) construct the suffix subtree:

- 1 Locate and store all $n = f_\pi$ positions of π – we will sort n strings in-memory; i.e. all strings are ambiguous.
- 2 Optimal prefix length range of chars following π $r = \rho \frac{M}{n}$ – *Elastic Range*.
- 3 Read the next r characters for each string.

Horizontal partitioning

For a prefix π (virtual group G) construct the suffix subtree:

- 1 Locate and store all $n = f_\pi$ positions of π – we will sort n strings in-memory; i.e. all strings are ambiguous.
- 2 Optimal prefix length range of chars following π $r = \rho \frac{M}{n}$ – *Elastic Range*.
- 3 Read the next r characters for each string.
- 4 Sort strings in-memory and record branching information (=LCP) and the original position (=SA).

Horizontal partitioning

For a prefix π (virtual group G) construct the suffix subtree:

- 1 Locate and store all $n = f_\pi$ positions of π – we will sort n strings in-memory; i.e. all strings are ambiguous.
- 2 Optimal prefix length range of chars following π $r = \rho \frac{M}{n}$ – *Elastic Range*.
- 3 Read the next r characters for each string.
- 4 Sort strings in-memory and record branching information (=LCP) and the original position (=SA).
- 5 Let n be the number of strings that are still ambiguous; go to step 2 until $n = 0$ (as r increases n decreases).

Horizontal partitioning

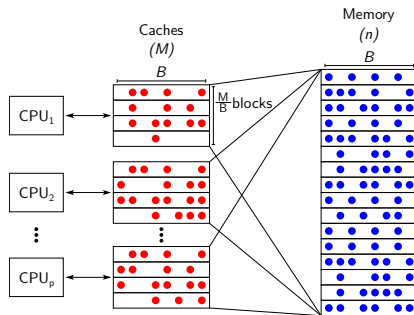
For a prefix π (virtual group G) construct the suffix subtree:

- 1 Locate and store all $n = f_\pi$ positions of π – we will sort n strings in-memory; i.e. all strings are ambiguous.
- 2 Optimal prefix length range of chars following π $r = \rho \frac{M}{n}$ – *Elastic Range*.
- 3 Read the next r characters for each string.
- 4 Sort strings in-memory and record branching information (=LCP) and the original position (=SA).
- 5 Let n be the number of strings that are still ambiguous; go to step 2 until $n = 0$ (as r increases n decreases).
- 6 Construct suffix subtree in DF manner using SA and LCP.

Model of computation

Parallel External Memory model (PEM):⁵

- Shared memory model,
- 2-level memory hierarchy:
 - p processors, each with private cache of size M bytes.
 - parallel memory transfers in blocks of size B bytes.
- Performance metrics:
 - parallel time,
 - parallel block transfers (cache complexity).
- Concurrent reads assumed.



⁵Arge, Goodrich, Nelson, Sitchinava 2008

Assumptions and Definitions

- Alphabet Σ of size σ ; the number of processors p ; the size of block B ; the size of memory M ; and the length of input string N .
- The worst-case the input text is skewed: $T = AAA\dots$

Assumptions and Definitions

- Alphabet Σ of size σ ; the number of processors p ; the size of block B ; the size of memory M ; and the length of input string N .
- The worst-case the input text is skewed: $T = AAA\dots$
- Vertical partitioning expands S-prefixes by one character at the time.
- Requires N scans, that equals $\frac{N^2}{2}$ comparisons.
- Cache complexity $O(\frac{N^2}{B})$.

Assumptions and Definitions

- Alphabet Σ of size σ ; the number of processors p ; the size of block B ; the size of memory M ; and the length of input string N .
- The worst-case the input text is skewed: $T = AAA\dots$
- Vertical partitioning expands S-prefixes by one character at the time.
- Requires N scans, that equals $\frac{N^2}{2}$ comparisons.
- Cache complexity $O(\frac{N^2}{B})$.

Our assumption:

- Input text is random (viable for a single genome, proteins).
- At any place the probability of each character to occur is $\frac{1}{\sigma}$.
- The suffix tree build from a random string is shallowest (Szpankowski 1993).

Analysis: Vertical partitioning

Assume $M < \sqrt{N}$.

- Extending π by one from 1 till $\log_{\sigma} \frac{N}{M}$ and hence this many scans of the text.
- Sorting $P = N/M$ prefixes.
- Packing prefixes into virtual groups.

Analysis: Vertical partitioning

Assume $M < \sqrt{N}$.

- Extending π by one from 1 till $\log_{\sigma} \frac{N}{M}$ and hence this many scans of the text.
- Sorting $P = N/M$ prefixes.
- Packing prefixes into virtual groups.
- Consequently the number of I/Os

$$O\left(\log_{\sigma} \frac{N}{M} \cdot \left(\frac{N}{B} + M^2\right) + \frac{N}{M \cdot B} \log_{\frac{M}{B}} \frac{N}{M \cdot B}\right)$$

Analysis: Horizontal partitioning

For building of P subtrees by p processors, where each subtree:

- Designed that data structure fits into memory.

Analysis: Horizontal partitioning

For building of P subtrees by p processors, where each subtree:

- Designed that data structure fits into memory.
- Consequently the number of I/Os

$$O\left(\frac{N}{B} \cdot \log_{\sigma} \frac{N}{M}\right)$$

Evaluation: environment

Computer and environment:

- 2× 16-core AMD Opteron 6272 @2,100 MHz
- 128 GiB RAM
- Seagate Baracuda 250 GB, 7,200 RPM, 32 MiB cache, SATA
- Ubuntu server 12.04, Linux kernel 3.11.0
- ext4 file system, deadline I/O scheduler
- MPI programming

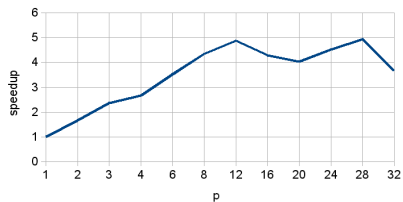
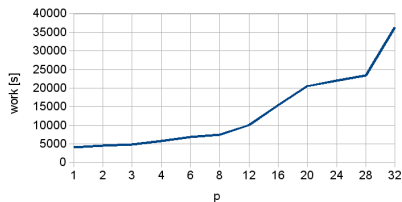
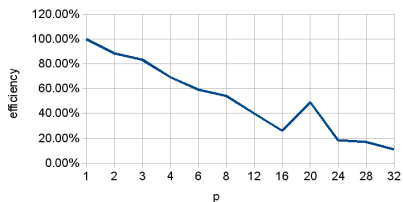
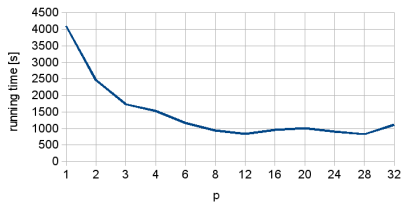
ERa parameters:

- Memory size per core: 2 GiB
- Input text: Human genome HG18.txt, 2.8 Gbp

ERa output:

- Total suffix tree size: 77.3 GB stored in 187 files
- \mathcal{T}_{top} size: 10.2 KB

Results – 1



The time increases as we increase the number of cores.

Results – 2

So what is the machine doing?

Results – 2

So what is the machine doing?

init Initialization including broadcasting to MPI clients.

cnt*, **cnt1** Vertical partitioning: counting the frequency and locating the S-prefixes in the input text.

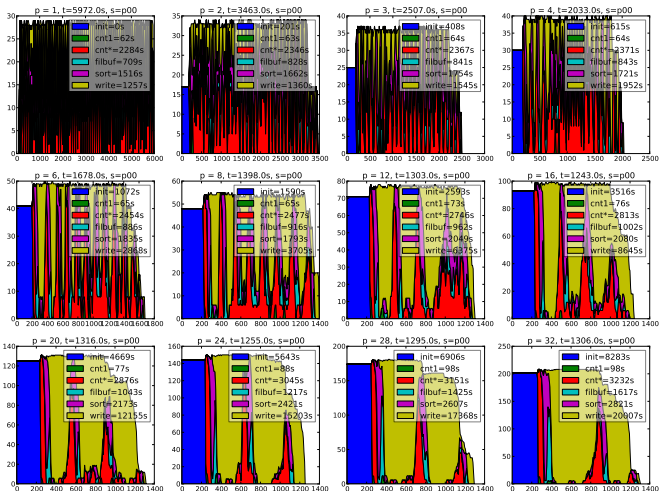
filbuf Horizontal partitioning: reading the input text.

sort Horizontal partitioning: in-memory string sorting.

write Horizontal partitioning: writing the final result to disk.

Results – 2 contd.

parallel10_devnullprobability CPU times p00



Hypothesis 1

Observation 1: The majority of time is spent writing the final result to the disk.

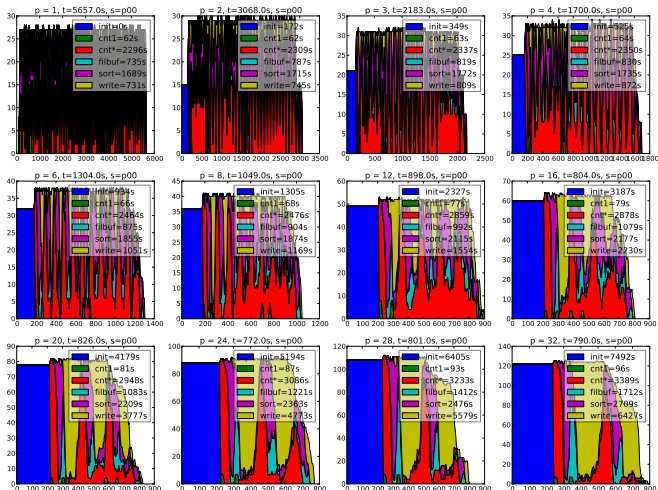
Hypothesis 1

Observation 1: The majority of time is spent writing the final result to the disk.

Hypothesis 1: Problem is the disk performance, so replace HDD with SSD.

Results – 3

parallel10_devnullprobability_ssd CPU times p00



Hypothesis 2

Observation 2: The amount of time for writting decreased, but as the number of cores grows, it is still substantial.

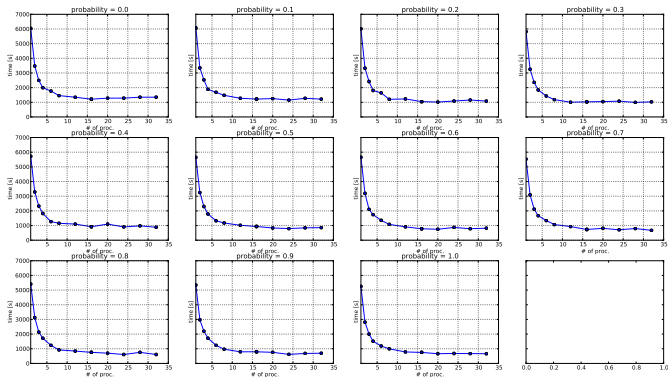
Hypothesis 2

Observation 2: The amount of time for writting decreased, but as the number of cores grows, it is still substantial.

Hypothesis 2: There it is still a problem with a disk performance and consequently further speed-up disk by writting to `/dev/null`.

Results – 4

times per # of proc., idealval prob. wise



Hypothesis 3

Observation 3: Things are getting better, but there is still an increase in time when the number of cores is increased.

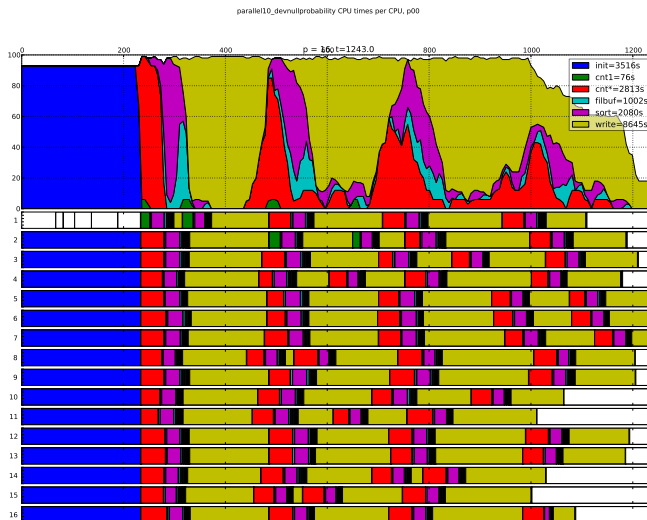
Hypothesis 3

Observation 3: Things are getting better, but there is still an increase in time when the number of cores is increased.

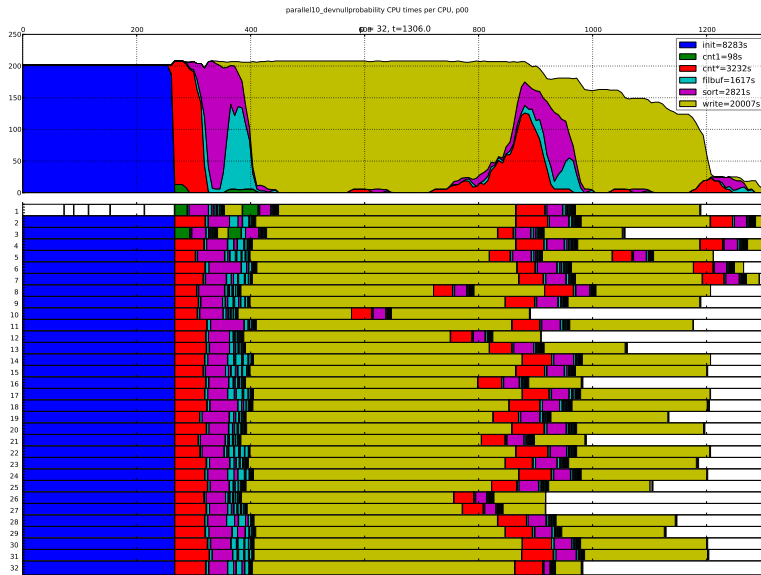
Hypothesis 3: ??

Check in more detail what the processes are doing.

Results – 5 ($p = 16$)



Results – 5 ($p = 32$)



Conclusion

- There is a substantial gap between theoretical results and practically used solutions.
- ERa despite being practically the fastest algorithm is **not theoretically tight** – even for random input strings with uniform substring distribution.

Conclusion

- There is a substantial gap between theoretical results and practically used solutions.
- ERa despite being practically the fastest algorithm is **not theoretically tight** – even for random input strings with uniform substring distribution.

Open challenges:

- Analyse ERa bottlenecks for further improvements (see if they match the critical terms in time and I/O complexities).
- Shall we choose some other basic technique for the implementation of a practical algorithm?
- Design a theoretically tight yet practically competitive parallel algorithm for suffix tree construction.

Thanx for your attention!