# ERA revisited: Theoretical and Experimental evaluation

Matevž Jekovec, Andrej Brodnik

University of Ljubljana
Faculty of Computer and Information Science

KAUST, March 1-5, 2014

# Text indexing problem

## Text indexing problem

### Problem statement

Given unstructured input string $S$ consisting of $N$ characters from alphabet $\Sigma$ of size $\sigma$ build an index such that for the pattern $P$ we:

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

## Text indexing problem

### Problem statement

Given unstructured input string $S$ consisting of $N$ characters from alphabet $\Sigma$ of size $\sigma$ build an index such that for the pattern $P$ we:
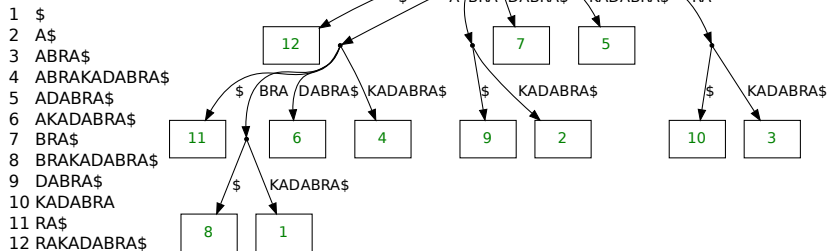
- determine whether $P$ occurs in $S$ in time $O(P)$,
- find all occurrences of $P$ in $S$ in time $O(P + occ)$,
- find the longest common prefix (LCP) of $P$ and any suffix of $S$ in time $O(LCP(P, S))$.

### Solution

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

## Suffix tree — Example

$$T = \overset{1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9\ 10\ 11\ 12}{ABRAKADABRA\$}$$

1  $
2  A$
3  ABRA$
4  ABRAKADABRA$
5  ADABRA$
6  AKADABRA$
7  BRA$
8  BRAKADABRA$
9  DABRA$
10 KADABRA
11 RA$
12 RAKADABRA$

## 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 |

---

[1]Bedathur 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 $\Sigma$ | No | No | Yes |
| Parallel | No | No | PDAM |

- Practical:

|  | Semi-disk-based | | | Out-of-core | | |
|---|---|---|---|---|---|---|
|  | TDD | TRLS. | $B^2ST$ | WF | ERA | PCF |
|  | ('04) | ('07) | ('09) | ('09) | ('11) | ('13) |
| Work w.c. | $O(N^2)$ | $O(N^2)$ | $O(N^2)$ | $O(N^2)$ | $O(N^2)$ | $O(\sqrt{p}N)$ |
| I/O eff. | R. | R. | R.+S. | R.+S. | R.+S. | R.+S. |
| Unbnd. $\Sigma$ | No | No | No | No | No | No |
| Parallel | No | No | No | Yes | Yes | Yes |

[1]Bedathur 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 $\Sigma$ | No | No | Yes |
| Parallel | | | PDAM |

> Huge gap between the theoretical and practical results!

- Practical:

| | Semi-disk-based | | | Out-of-core | | |
|---|---|---|---|---|---|---|
| | TDD ('04) | TRLS. ('07) | B$^2$ST ('09) | WF ('09) | ERA ('11) | PCF ('13) |
| Work w.c. | $O(N^2)$ | $O(N^2)$ | $O(N^2)$ | $O(N^2)$ | $O(N^2)$ | $O(\sqrt{p}N)$ |
| I/O eff. | R. | R. | R.+S. | R.+S. | R.+S. | R.+S. |
| Unbnd. $\Sigma$ | No | No | No | No | No | No |
| Parallel | No | No | No | Yes | Yes | Yes |

[1]Bedathur and Haritsa (2004)

## Suffix tree construction lower bounds

- Sequential:

|  | bounded $\Sigma$ | unbounded $\Sigma$ |
|---|---|---|
| Time | $\Omega(Sort(N))$ | $\Omega(Sort(N))$ |
| I/Os[2] | $\Omega(Sort(N))$ | $\Omega(Sort(N))$ |
| Space[3] | $\Omega(N \lg \sigma)$ bits | $\Omega(N \lg \sigma)$ bits |

---

[2]EM model

[3]Uncompressed index in word RAM

[4]PEM model

## Suffix tree construction lower bounds

- Sequential:

|  | bounded $\Sigma$ | unbounded $\Sigma$ |
|---|---|---|
| Time | $\Omega(Sort(N))$ | $\Omega(Sort(N))$ |
| I/Os[2] | $\Omega(Sort(N))$ | $\Omega(Sort(N))$ |
| Space[3] | $\Omega(N \lg \sigma)$ bits | $\Omega(N \lg \sigma)$ bits |

- Parallel on $p$ processing units:

|  | bounded $\Sigma$ | unbounded $\Sigma$ |
|---|---|---|
| Parallel time | $\Omega(Sort_p(N))$ | $\Omega(Sort_p(N))$ |
| Parallel I/Os[4] | $\Omega(Sort_p(N))$ | $\Omega(Sort_p(N))$ |
| Space[3] | $\Omega(N \lg \sigma)$ bits | $\Omega(N \lg \sigma)$ bits |

---

[2]EM model

[3]Uncompressed index in word RAM

[4]PEM model

## Suffix tree construction lower bounds

- Sequential:

|           | bounded $\Sigma$       | unbounded $\Sigma$                         |
|-----------|------------------------|--------------------------------------------|
| 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 $\Sigma$                    | unbounded $\Sigma$                            |
|-------------------|-------------------------------------|-----------------------------------------------|
| 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                   |

---

[2]EM model

[3]Uncompressed index in word RAM

[4]PEM model

Introduction
0000

Design goals
●○

ERA
000000

Formal analysis
000000000000

Empirical evaluation
000000000000000000

Conclusion

# Challenges

## Challenges

Theoretical algorithms: Lack of locality of reference.

## Challenges

Theoretical algorithms: Lack of locality of reference.
Goal: Use scans only both for input text and the resulting suffix tree!

Introduction
0000

Design goals
●○

ERA
000000

Formal analysis
000000000000

Empirical evaluation
000000000000000000

Conclusion

## Challenges

Theoretical algorithms: Lack of locality of reference.

Goal: Use scans only both for input text and the resulting suffix tree!

Counterintuitive: Input text is arbitrary, suffix tree is lexicographically ordered.

Challenges contd.

I/O efficient solution (eg. WF-ERA, $B^2ST$-PCF):

1. Scan part of the input text from the disk,

Introduction
0000

Design goals
○●

ERA
000000

Formal analysis
000000000000

Empirical evaluation
000000000000000000

Conclusion

# Challenges contd.

I/O efficient solution (eg. WF-ERA, $B^2ST$-PCF):

1. Scan part of the input text from the disk,

2. do in-memory sorting (=random accesses in fast memory only!),

## Challenges contd.

I/O efficient solution (eg. WF-ERA, B$^2$ST-PCF):

1. Scan part of the input text from the disk,
2. do in-memory sorting (=random accesses in fast memory only!),
3. construct the corresponding part of the suffix tree,

Introduction
oooo

Design goals
o●

ERA
oooooo

Formal analysis
oooooooooooo

Empirical evaluation
ooooooooooooooooo

Conclusion

# Challenges contd.

I/O efficient solution (eg. WF-ERA, B$^2$ST-PCF):

1. Scan part of the input text from the disk,
2. do in-memory sorting (=random accesses in fast memory only!),
3. construct the corresponding part of the suffix tree,
4. glue parts together,

## Challenges contd.

I/O efficient solution (eg. WF-ERA, $B^2$ST-PCF):

1. Scan part of the input text from the disk,
2. do in-memory sorting (=random accesses in fast memory only!),
3. construct the corresponding part of the suffix tree,
4. glue parts together,
5. and contiguously write it to disk.

# ERA — Elastic Range[5]

Introduction    Design goals    ERA    Formal analysis    Empirical evaluation    Conclusion
0000            00              ●00000  000000000000      00000000000000000

ERA — Elastic Range[5]

- The fastest practical, parallel suffix tree construction algorithm to date.

---

[5]Mansour, Allam, Skiadopoulos, Kalnis (2011)

# ERA — Elastic Range[5]

- The fastest practical, parallel suffix tree construction algorithm to date.
- Time complexity: $O(N^2)$ w.c. — for extremely skewed text!

---

[5]Mansour, Allam, Skiadopoulos, Kalnis (2011)

# ERA — Elastic Range[5]

- The fastest practical, parallel suffix tree construction algorithm to date.
- 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!

[5]Mansour, Allam, Skiadopoulos, Kalnis (2011)

# ERA contd.



In-memory trie

$Tp_2$     $Tp_n$

$Tp_1$     ...

Horizontal partitioning

Grouping     Elastic range

Vertical partitioning

ERA constructs the suffix tree in two steps:

# ERA contd.



ERA constructs the suffix tree in two steps:

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

## ERA contd.



ERA constructs the suffix tree in two steps:

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

2. The **horizontal partitioning** step builds the actual suffix subtrees.

## ERA contd.

---

**Algorithm 1:** ERA

**Input**: String $S$, Alphabet $\Sigma$, Processors $P$, Private cache size $M$

**Output**: Suffix tree $\mathcal{T}$

**1** $\mathcal{T}_{top}, G \leftarrow VerticalPartitioning(S, \Sigma, M)$

**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**              $\pi \leftarrow G.pop()$

**7**              $\mathcal{T}_\pi \leftarrow HorizontalPartitioning(S, \Sigma, \pi)$

**8**              $Link(\mathcal{T}, \mathcal{T}_\pi)$

**9** return $\mathcal{T}$

---

## Vertical partitioning

Define **S-prefix** $\pi$ as the prefix of the suffixes in the text.

---

[6]Yue (1991)

## Vertical partitioning

Define **S-prefix** $\pi$ as the prefix of the suffixes in the text.
Idea: The S-prefix frequency $f_\pi$ equals $\#$ of leaves in the suffix
subtree corresponding to $\pi$. Assume $2f_\pi$ is the w. c. subtree size.

---

[6]Yue (1991)

## Vertical partitioning

Define **S-prefix** $\pi$ as the prefix of the suffixes in the text.
Idea: The S-prefix frequency $f_\pi$ equals # of leaves in the suffix
subtree corresponding to $\pi$. Assume $2f_\pi$ is the w. c. subtree size.
Vertical partitioning steps:

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

---

[6]Yue (1991)

## Vertical partitioning

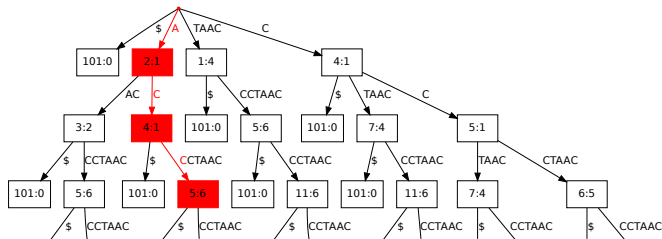Define **S-prefix** $\pi$ as the prefix of the suffixes in the text.
Idea: The S-prefix frequency $f_\pi$ equals # of leaves in the suffix
subtree corresponding to $\pi$. Assume $2f_\pi$ is the w. c. subtree size.
Vertical partitioning steps:

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

2. For each $\pi : f_\pi > M$, expand S-prefix with the right character
   and count the frequency of obtained S-prefixes (now length 2).

---

[6]Yue (1991)

## Vertical partitioning

Define **S-prefix** $\pi$ as the prefix of the suffixes in the text.
Idea: The S-prefix frequency $f_\pi$ equals # of leaves in the suffix
subtree corresponding to $\pi$. Assume $2f_\pi$ is the w. c. subtree size.
Vertical partitioning steps:

1. Scan the text and obtain the characters frequency $f_\pi : \pi \in \Sigma$.
   - We counted all S-prefixes of length 1.
2. For each $\pi : f_\pi > M$, expand S-prefix with the right character
   and count the frequency of obtained S-prefixes (now length 2).
3. Repeat step two for S-prefixes of length $3, 4...$, until all $\mathcal{T}_\pi$ just
   fit into the memory $M$.

---

[6]Yue (1991)

## Vertical partitioning

Define **S-prefix** $\pi$ as the prefix of the suffixes in the text.
Idea: The S-prefix frequency $f_\pi$ equals # of leaves in the suffix
subtree corresponding to $\pi$. Assume $2f_\pi$ is the w. c. subtree size.
Vertical partitioning steps:

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

2. For each $\pi : f_\pi > M$, expand S-prefix with the right character
   and count the frequency of obtained S-prefixes (now length 2).

3. Repeat step two for S-prefixes of length 3, 4..., until all $\mathcal{T}_\pi$ just
   fit into the memory $M$.

4. 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[6].

---

[6]Yue (1991)

Introduction
oooo

Design goals
oo

ERA
oooooo●o

Formal analysis
ooooooooooooo

Empirical evaluation
oooooooooooooooooo

Conclusion

## Vertical partitioning — Example

$\pi = ACC$

Frequency $f_{ACC} = 12$

TAACCCTA
ACCCTAAC
CCTAACCC
TAACCCTA
ACCCTAAC
CCTAACCC
TAACCCTA
ACCCTAAC
CCTAACCC
TAACCCTA
ACCCTAAC
CCTAACCC
TAAC

## Horizontal partitioning

For each virtual group, construct the corresponding suffix subtrees in parallel:

## Horizontal partitioning

For each virtual group, construct the corresponding suffix subtrees in parallel:

1. Locate and store all positions of S-prefixes for the virtual group. Each of located S-prefixes is to become a leaf in the working suffix subtree.

## Horizontal partitioning

For each virtual group, construct the corresponding suffix subtrees in parallel:

1. Locate and store all positions of S-prefixes for the virtual group. Each of located S-prefixes is to become a leaf in the working suffix subtree.

2. Calculate the optimal buffer length *range*
   - Note: The name *Elastic Range*.

## Horizontal partitioning

For each virtual group, construct the corresponding suffix subtrees
in parallel:

1. Locate and store all positions of S-prefixes for the virtual
   group. Each of located S-prefixes is to become a leaf in the
   working suffix subtree.
2. Calculate the optimal buffer length *range*
   - Note: The name *Elastic Range*.
3. Read the next *range* characters for each S-prefix occurrence.

## Horizontal partitioning

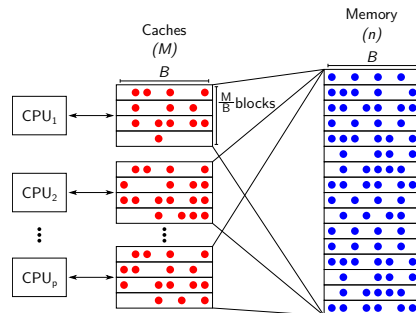For each virtual group, construct the corresponding suffix subtrees in parallel:

1. Locate and store all positions of S-prefixes for the virtual group. Each of located S-prefixes is to become a leaf in the working suffix subtree.

2. Calculate the optimal buffer length *range*
   - Note: The name *Elastic Range*.

3. Read the next *range* characters for each S-prefix occurrence.

4. Do in-memory sorting of read text, remember branching information (=LCP) and the original position (=SA).

## Horizontal partitioning

For each virtual group, construct the corresponding suffix subtrees in parallel:

1. Locate and store all positions of S-prefixes for the virtual group. Each of located S-prefixes is to become a leaf in the working suffix subtree.

2. Calculate the optimal buffer length *range*
   - Note: The name *Elastic Range*.

3. Read the next *range* characters for each S-prefix occurrence.

4. Do in-memory sorting of read text, remember branching information (=LCP) and the original position (=SA).

5. Until all the read buffers are unique, goto step 2.
   - In the next step: While less leafs are orphans, *range* increases, frequency drops.

## Horizontal partitioning

For each virtual group, construct the corresponding suffix subtrees in parallel:

1. Locate and store all positions of S-prefixes for the virtual group. Each of located S-prefixes is to become a leaf in the working suffix subtree.

2. Calculate the optimal buffer length *range*
   - Note: The name *Elastic Range*.

3. Read the next *range* characters for each S-prefix occurrence.

4. Do in-memory sorting of read text, remember branching information ($=$LCP) and the original position ($=$SA).

5. Until all the read buffers are unique, goto step 2.
   - In the next step: While less leafs are orphans, *range* increases, frequency drops.

6. Construct suffix subtree in D-F manner using SA and LCP.

## Model of computation

---
[7]Arge, Goodrich, Nelson, Sitchinava 2008

## Model of computation

Parallel External Memory model (PEM):[7]

- 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.



Caches
(M)

Memory
(n)

$\frac{M}{B}$-blocks

CPU$_1$

CPU$_2$

CPU$_p$

---

[7]Arge, Goodrich, Nelson, Sitchinava 2008

## Assumption

If worst-case input text for ERA is skewed:

$$T = AAA...$$

## Assumption

If worst-case input text for ERA is skewed:

$$T = AAA...$$

then

- vertical partitioning requires $N$ scans $= O(N^2)$ comparisons,
- cache complexity $O\left(\frac{N^2}{B}\right)$.

## Assumption

If worst-case input text for ERA is skewed:

$$T = AAA...$$

then

- vertical partitioning requires $N$ scans $= O(N^2)$ comparisons,
- cache complexity $O\left(\frac{N^2}{B}\right)$.

Our assumption:

- Input text is random (viable for a single human genome, proteins).
- At any place the probability of each character to occur is $\frac{1}{\sigma}$.
- Goal: Calculate expected time and cache complexity.

**Algorithm 2:** VerticalPartitioning

**Input**: Input string $S$, alphabet $\Sigma$, $1^{st}$ level memory size $M$

**Output**: Set of *VirtualTrees*

1 *VirtualTrees* $\leftarrow \emptyset$

2 $P \leftarrow \emptyset$

3 $P' \leftarrow \{\forall$ symbol $s \in \Sigma$ generate a S-prefix $\pi_i \in P'\}$

4 **repeat**

5     scan input string $S$

6     count in $S$ the frequency $f_{\pi_i}$ of every S-prefix $\pi_i \in P'$

7     **forall the** $\pi_i \in P'$ **do**

8         **if** $0 < f_{\pi_i} \leq M$ **then** add $\pi_i$ to $P$

9         **else forall the** *symbol* $s \in \Sigma$ **do** add $\pi_i s$ to $P'$

10         remove $\pi_i$ from $P'$

11 **until** $P' = \emptyset$

12 sort $P$ in descending $f_{\pi_i}$ order

13 **repeat**

14     $G \leftarrow \emptyset$

15     add $P.head$ to $G$ and remove the item from $P$

16     $curr \leftarrow$ next item in $P$

17     **while** *NOT end of P* **do**

18         **if** $f_{curr} + SUM_{\gamma_i \in G}(f_{\gamma_i}) \leq M$ **then**

19             add *curr* to $G$ and remove the item from $P$

20         $curr \leftarrow$ next item in $P$

21     add $G$ to *VirtualTrees*

22 **until** $P = \emptyset$

23 **return** *VirtualTrees*

Introduction
0000

Design goals
00

ERA
000000

Formal analysis
000●00000000

Empirical evaluation
000000000000000

Conclusion

# Analysis: Vertical partitioning

## Analysis: Vertical partitioning

Expected behavior:

1. Extension of S-prefixes:
   - Initially $\sigma$ S-prefixes of frequency $f_\pi = \frac{N}{\sigma}$ each.
   - $f_\pi$ divided by $\sigma$ each iteration until $f_\pi < M$.
   - Total $\log_\sigma N - \log_\sigma M = \log_\sigma \frac{N}{M}$ iterations.
   - Finally $\frac{N}{M}$ unique S-prefixes with frequency $\frac{M}{\sigma} < f_\pi \leq M$.

Introduction
oooo

Design goals
oo

ERA
oooooo

Formal analysis
ooo●ooooooooo

Empirical evaluation
oooooooooooooooo

Conclusion

## Analysis: Vertical partitioning

Expected behavior:

1. Extension of S-prefixes:
   - Initially $\sigma$ S-prefixes of frequency $f_\pi = \frac{N}{\sigma}$ each.
   - $f_\pi$ divided by $\sigma$ each iteration until $f_\pi < M$.
   - Total $\log_\sigma N - \log_\sigma M = \log_\sigma \frac{N}{M}$ iterations.
   - Finally $\frac{N}{M}$ unique S-prefixes with frequency $\frac{M}{\sigma} < f_\pi \leq M$.

2. Atomic sorting the frequencies using one of the comparison-based sorting algorithms.

## Analysis: Vertical partitioning

Expected behavior:

1. Extension of S-prefixes:
   - Initially $\sigma$ S-prefixes of frequency $f_\pi = \frac{N}{\sigma}$ each.
   - $f_\pi$ divided by $\sigma$ each iteration until $f_\pi < M$.
   - Total $\log_\sigma N - \log_\sigma M = \log_\sigma \frac{N}{M}$ iterations.
   - Finally $\frac{N}{M}$ unique S-prefixes with frequency $\frac{M}{\sigma} < f_\pi \leq M$.

2. Atomic sorting the frequencies using one of the comparison-based sorting algorithms.

3. Virtual trees construction (bin packing problem):
   - At least 1 and at most $\sigma$ S-prefixes are packed each iteration.
   - External loop iterated between $\frac{N}{\sigma M}$ and $\frac{N}{M}$ times.

Introduction
oooo

Design goals
oo

ERA
oooooo

Formal analysis
ooooo●ooooooo

Empirical evaluation
oooooooooooooooo

Conclusion

# Analysis: Vertical partitioning — Time

Introduction
oooo

Design goals
oo

ERA
oooooo

Formal analysis
ooooo●oooooo

Empirical evaluation
ooooooooooooooooo

Conclusion

## Analysis: Vertical partitioning — Time

1. Extension of S-prefixes

$$\sum_{i=1}^{\log_\sigma \frac{N}{M}} \left( scan(n) + \sigma^{i+1} \right) = \log_\sigma \frac{N}{M} \cdot scan(n) + \frac{\sigma^2(N-M)}{M \cdot \sigma - M} =$$

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

Introduction
0000

Design goals
00

ERA
000000

Formal analysis
0000●0000000

Empirical evaluation
000000000000000

Conclusion

## Analysis: Vertical partitioning — Time

1. Extension of S-prefixes
$$\sum_{i=1}^{\log_\sigma \frac{N}{M}} \left(scan(n) + \sigma^{i+1}\right) = \log_\sigma \frac{N}{M} \cdot scan(n) + \frac{\sigma^2(N-M)}{M \cdot \sigma - M} =$$
$$O\left(N \log_\sigma \frac{N}{M} + \frac{\sigma N}{M}\right)$$

2. Sorting
$$O\left(\frac{N}{M} \lg \frac{N}{M}\right)$$

# Analysis: Vertical partitioning — Time

1. Extension of S-prefixes
$$\sum_{i=1}^{\log_\sigma \frac{N}{M}} \left( scan(n) + \sigma^{i+1} \right) = \log_\sigma \frac{N}{M} \cdot scan(n) + \frac{\sigma^2(N-M)}{M \cdot \sigma - M} =$$
$$O\left( N \log_\sigma \frac{N}{M} + \frac{\sigma N}{M} \right)$$

2. Sorting
$$O\left( \frac{N}{M} \lg \frac{N}{M} \right)$$

3. Virtual trees construction
$$O\left( \left( \frac{N}{M} \right)^2 \right)$$

## Analysis: Vertical partitioning — Time

1. Extension of S-prefixes
$$\sum_{i=1}^{\log_\sigma \frac{N}{M}} \left( scan(n) + \sigma^{i+1} \right) = \log_\sigma \frac{N}{M} \cdot scan(n) + \frac{\sigma^2(N-M)}{M \cdot \sigma - M} =$$
$$O\left( N \log_\sigma \frac{N}{M} + \frac{\sigma N}{M} \right)$$

2. Sorting
$$O\left( \frac{N}{M} \lg \frac{N}{M} \right)$$

3. Virtual trees construction
$$O\left( \left(\frac{N}{M}\right)^2 \right)$$

Overall:

- If $\sigma < M$: $O\left( N \log_\sigma \frac{N}{M} + \left(\frac{N}{M}\right)^2 \right)$

- If $\sigma \geq M$: $O\left( N \log_\sigma \frac{N}{M} + \frac{\sigma N}{M} + \frac{N}{M} \lg \frac{N}{M} + \left(\frac{N}{M}\right)^2 \right)$

## Analysis: Vertical partitioning — I/O

1. Extension of S-prefixes
   - Line 6: $scan(N)$ for reading

Introduction
0000

Design goals
00

ERA
000000

Formal analysis
000000●000000

Empirical evaluation
000000000000000

Conclusion

# Analysis: Vertical partitioning — I/O

1. Extension of S-prefixes
   - Line 6: $scan(N)$ for reading
     $|P'| = O\left(\frac{N}{M}\right)$
     If $|P'| \leq M$: no I/Os for writing $f_\pi$
     If $|P'| > M$: $\frac{M}{|P'|} = \frac{M^2}{N}$ I/Os for storing $f_\pi$

# Analysis: Vertical partitioning — I/O

1. Extension of S-prefixes
   - Line 6: $scan(N)$ for reading
     $|P'| = O\left(\frac{N}{M}\right)$
     If $|P'| \leq M$: no I/Os for writing $f_\pi$
     If $|P'| > M$: $\frac{M}{|P'|} = \frac{M^2}{N}$ I/Os for storing $f_\pi$
   - Lines 7-10:
     If $|P'| \leq M$: no I/Os
     If $|P'| > M$: $\frac{P'}{B} = \frac{N}{M \cdot B}$ I/Os

# Analysis: Vertical partitioning — I/O

1. Extension of S-prefixes
   - Line 6: $scan(N)$ for reading
     $|P'| = O\left(\frac{N}{M}\right)$
     If $|P'| \le M$: no I/Os for writing $f_\pi$
     If $|P'| > M$: $\frac{M}{|P'|} = \frac{M^2}{N}$ I/Os for storing $f_\pi$
   - Lines 7-10:
     If $|P'| \le M$: no I/Os
     If $|P'| > M$: $\frac{P'}{B} = \frac{N}{M \cdot B}$ I/Os
2. Sorting $|P| = \frac{N}{M}$ elements:
   If $M \ge \sqrt{N}$: no I/Os
   If $M < \sqrt{N}$: $O(\frac{N}{M \cdot B} \log_{\frac{M}{B}} \frac{N}{M \cdot B})$ I/Os

## Analysis: Vertical partitioning — I/O

1. Extension of S-prefixes
   - Line 6: $scan(N)$ for reading
     $|P'| = O\left(\frac{N}{M}\right)$
     If $|P'| \leq M$: no I/Os for writing $f_\pi$
     If $|P'| > M$: $\frac{M}{|P'|} = \frac{M^2}{N}$ I/Os for storing $f_\pi$
   - Lines 7-10:
     If $|P'| \leq M$: no I/Os
     If $|P'| > M$: $\frac{P'}{B} = \frac{N}{M \cdot B}$ I/Os

2. Sorting $|P| = \frac{N}{M}$ elements:
   If $M \geq \sqrt{N}$: no I/Os
   If $M < \sqrt{N}$: $O(\frac{N}{M \cdot B} \log_{\frac{M}{B}} \frac{N}{M \cdot B})$ I/Os

3. Virtual tree $G \leq M$:
   $M \geq \sqrt{N}$: no I/Os
   $M < \sqrt{N}$: $\frac{|P|}{B} = \frac{N}{M \cdot B}$ I/Os

Introduction
oooo

Design goals
oo

ERA
oooooo

Formal analysis
ooooooo●ooooo

Empirical evaluation
oooooooooooooooooo

Conclusion

## Analysis: Vertical partitioning — I/O contd.

Overall:

- If $M \geq \sqrt{N}$:
  $O\left(\frac{N}{B} \log_\sigma \frac{N}{M}\right)$

- If $M < \sqrt{N}$:
  $O\left(\log_\sigma \frac{N}{M} \cdot \left(\frac{N}{B} + M^2\right) + \frac{N}{M \cdot B} \log_{\frac{M}{B}} \frac{N}{M \cdot B} + \left(\frac{N}{M \cdot B}\right)^2\right)$

**Algorithm 3:** HorizontalPartitioning.SubTreePrepare

**Input**: Input string $S$, S-prefix $\pi$

**Output**: Arrays $SA$ and $LCP$ corresponding suffix sub-tree $\mathcal{T}_\pi$

1 $SA$ contains the locations of S-prefix $\pi$ in string $S$

2 $LCP \leftarrow \{\}$

3 $ISA \leftarrow \{0, 1, ..., |SA| - 1\}$

4 $A \leftarrow \{0, 0, ..., 0\}$

5 $Buf \leftarrow \{\}$

6 $P \leftarrow \{0, 1, ..., |L| - 1\}$

7 $start \leftarrow |\pi|$

8 **while** there exists an undefined $Buf[i]$, $1 \leq i \leq |SA| - 1$ **do**

9    $range \leftarrow GetRangeOfSymbols$

10    **for** $i \leftarrow 0$ to $|SA| - 1$ **do**

11      **if** $ISA[i] \neq done$ **then**

12        $Buf[ISA[i]] \leftarrow ReadRange(S, SA[ISA[i]] + start, range)$
       // ReadRange(S,a,b) reads $b$ symbols of $S$ starting at
       position $a$

13    **for** every active area $AA$ **do**

14      Reorder the elements of $Buf$, $P$ and $SA$ in $AA$ so that $Buf$
     is lexicographically sorted. In the process maintain the
     index $ISA$

15      If two or more elements $\{a_1, ..., a_t\} \in AA$, $2 \leq t$, exist such
     that $Buf[a_1] = ... = Buf[a_i]$ introduce for them a new
     active area

16    **for** all $i$ such that $Buf[i]$ is not defined, $1 \leq i \leq |SA| - 1$ **do**

17      $cp$ is the common prefix of $Buf[i - 1]$ and $Buf[i]$

18      **if** $|cp| < range$ **then**

19        $Buf[i] \leftarrow (Buf[i - 1][|cp|], Buf[i][|cp|], start + |cp|)$

20        **if** $Buf[i - 1]$ is defined or $i = 1$ **then**

21          Mark $ISA[P[i - 1]]$ and $A[i - 1]$ as done

22        **if** $Buf[i + 1]$ is defined or $i = [SA] - 1$ **then**

23          Mark $ISA[P[i]]$ and $A[i]$ as done // last element of
         an active area

24    $start \leftarrow start + range$

25 **return** $(SA, LCP)$

## Analysis: Horizontal partitioning

## Analysis: Horizontal partitioning

Expected behaviour:

## Analysis: Horizontal partitioning

Expected behaviour:

- Define $n$ the number of unfinished branches, then $n \cdot range = O(M)$.

## Analysis: Horizontal partitioning

Expected behaviour:

- Define $n$ the number of unfinished branches, then
  $n \cdot range = O(M)$.
- Intuitively $n$ decreases and $range$ increases during execution of
  lines 8-24.

## Analysis: Horizontal partitioning

Expected behaviour:

- Define $n$ the number of unfinished branches, then $n \cdot range = O(M)$.
- Intuitively $n$ decreases and *range* increases during execution of lines 8-24.
- For length $k$, there can be at most $\sigma^k$ unique strings. For random text and step $1 \le i \le k$, strings are non-unique until $k$ is reached.

## Analysis: Horizontal partitioning

Expected behaviour:

- Define $n$ the number of unfinished branches, then $n \cdot range = O(M)$.
- Intuitively $n$ decreases and *range* increases during execution of lines 8-24.
- For length $k$, there can be at most $\sigma^k$ unique strings. For random text and step $1 \leq i \leq k$, strings are non-unique until $k$ is reached.
- If $O(M)$ random strings need to be processed, then lines 8-24 is iterated $O(\log_\sigma M)$ times. The big-oh constant depends on *range*.

## Analysis: Horizontal partitioning — Time

## Analysis: Horizontal partitioning — Time

Each iteration:

1. Lines 10-12 required $n$ time to fill the buffers (constant time read).

## Analysis: Horizontal partitioning — Time

Each iteration:

1. Lines 10-12 required $n$ time to fill the buffers (constant time read).

2. String sorting requires $O(n \cdot range)$ time since the average *distinguising prefix size* equals $O(range)$.

## Analysis: Horizontal partitioning — Time

Each iteration:

1. Lines 10-12 required $n$ time to fill the buffers (constant time read).

2. String sorting requires $O(n \cdot range)$ time since the average *distinguising prefix size* equals $O(range)$.

3. Lines 16-23 require $O(n \cdot range)$ time in the worst case.

## Analysis: Horizontal partitioning — Time

Each iteration:

1. Lines 10-12 required *n* time to fill the buffers (constant time read).

2. String sorting requires $O(n \cdot range)$ time since the average *distinguising prefix size* equals $O(range)$.

3. Lines 16-23 require $O(n \cdot range)$ time in the worst case.

Overall: Assuming *p* processors equally balanced after processing $O(N/M)$ virtual groups require

$$O\left(\frac{N}{M}\frac{M\log_{\sigma}M}{p}\right) = O\left(\frac{N}{p}\log_{\sigma}M\right)$$

Introduction
oooo

Design goals
oo

ERA
oooooo

Formal analysis
ooooooooooo●o

Empirical evaluation
ooooooooooooooooo

Conclusion

Analysis: Horizontal partitioning — I/O

## Analysis: Horizontal partitioning — I/O

1. Cache misses occur in lines 10-12 only:

## Analysis: Horizontal partitioning — I/O

1. Cache misses occur in lines 10-12 only:
   - If $n \geq \frac{N}{B}$, then $O\left(\frac{N}{B}\right)$ I/Os.

## Analysis: Horizontal partitioning — I/O

1. Cache misses occur in lines 10-12 only:
   - If $n \geq \frac{N}{B}$, then $O\left(\frac{N}{B}\right)$ I/Os.
   - Else: $O(n)$ I/Os

## Analysis: Horizontal partitioning — I/O

1. Cache misses occur in lines 10-12 only:
   - If $n \geq \frac{N}{B}$, then $O\left(\frac{N}{B}\right)$ I/Os.
   - Else: $O(n)$ I/Os

2. When does the change from $n \geq \frac{N}{B}$ to $n < \frac{N}{B}$ occur?

## Analysis: Horizontal partitioning — I/O

1. Cache misses occur in lines 10-12 only:
   - If $n \geq \frac{N}{B}$, then $O\left(\frac{N}{B}\right)$ I/Os.
   - Else: $O(n)$ I/Os
2. When does the change from $n \geq \frac{N}{B}$ to $n < \frac{N}{B}$ occur?
3. Assuming uniformly random text, $n = c \cdot M$ for some constant $c$ **all the time**! (all branches are open until the last iteration)

# Analysis: Horizontal partitioning — I/O

1. Cache misses occur in lines 10-12 only:
   - If $n \geq \frac{N}{B}$, then $O\left(\frac{N}{B}\right)$ I/Os.
   - Else: $O(n)$ I/Os

2. When does the change from $n \geq \frac{N}{B}$ to $n < \frac{N}{B}$ occur?

3. Assuming uniformly random text, $n = c \cdot M$ for some constant $c$ **all the time**! (all branches are open until the last iteration)

4. Suffix subtree construction from SA and LCP requires a single $scan(N)$ I/Os only and is omitted.

## Analysis: Horizontal partitioning — I/O

1. Cache misses occur in lines 10-12 only:
   - If $n \geq \frac{N}{B}$, then $O\left(\frac{N}{B}\right)$ I/Os.
   - Else: $O(n)$ I/Os

2. When does the change from $n \geq \frac{N}{B}$ to $n < \frac{N}{B}$ occur?

3. Assuming uniformly random text, $n = c \cdot M$ for some constant $c$ **all the time**! (all branches are open until the last iteration)

4. Suffix subtree construction from SA and LCP requires a single $scan(N)$ I/Os only and is omitted.

5. I/O complexity for horizontal partitioning is thus

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

Wrap-up

## Wrap-up

Parallel time complexity of ERA (assuming $\sigma \leq M$):

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

Parallel cache complexity of ERA (assuming $M \geq \sqrt{N}$):

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

Introduction
oooo

Design goals
oo

ERA
oooooo

Formal analysis
oooooooooooo

Empirical evaluation
●ooooooooooooooo

Conclusion

# Empirical evaluation

## Empirical evaluation

Testing environment:

- 2x 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

ERA parameters:

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

ERA modification: Call fsync() after writing each file.

ERA output:

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

## Results – 1

# 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?

string cpy    Parsing and copying the string.

vertpart    Vertical partitioning.

cnt1, cnt*    Horizontal partitioning: determining locations of S-prefix in virtual trees of size 1 or $> 1$.

filbuf    Horizontal partitioning: reading range characters from S-prefix locations.

sort    Horizontal partitioning: string sorting, implicit LCP, SA construction.

write    Horizontal partitioning: extraction from LCP and SA to suffix tree, write to disk.

# Results – 2 contd.



parallel10_devnullprobability CPU times p00

# Results – 2 contd.



parallel10_devnullprobability iostat p00

Introduction
OOOO

Design goals
OO

ERA
OOOOOO

Formal analysis
OOOOOOOOOOOO

Empirical evaluation
OOOOOO●OOOOOOOOOOO

Conclusion

## Results – 2 contd.

## Hypotheis 1

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

Introduction    Design goals    ERA    Formal analysis    **Empirical evaluation**    Conclusion
0000            00             000000  000000000000        000000●000000000

Hypotheis 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

# Hypotheis 2

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

## Hypotheis 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`.
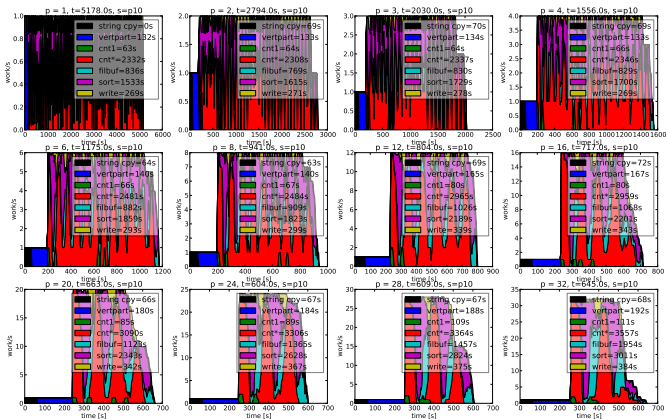
# Results – 4

# Results – 4 contd.



parallel10_devnullprobability CPU times p10

## Hypotheis 3

*Observation 3:* Things are getting better, but there is still an increase in time when the number of cores is increased.
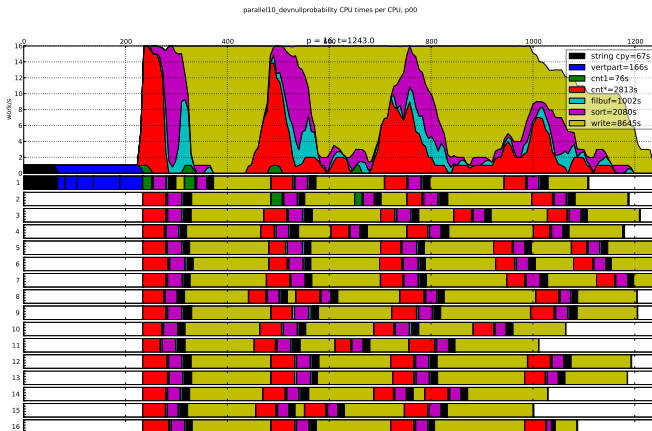
Introduction
0000

Design goals
00

ERA
000000

Formal analysis
000000000000

Empirical evaluation
000000000000000000

Conclusion

## Hypotheis 3

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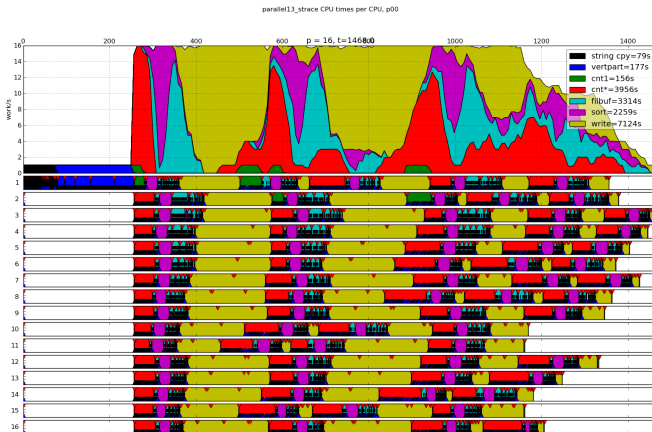
*Hypothesis 3:* ??

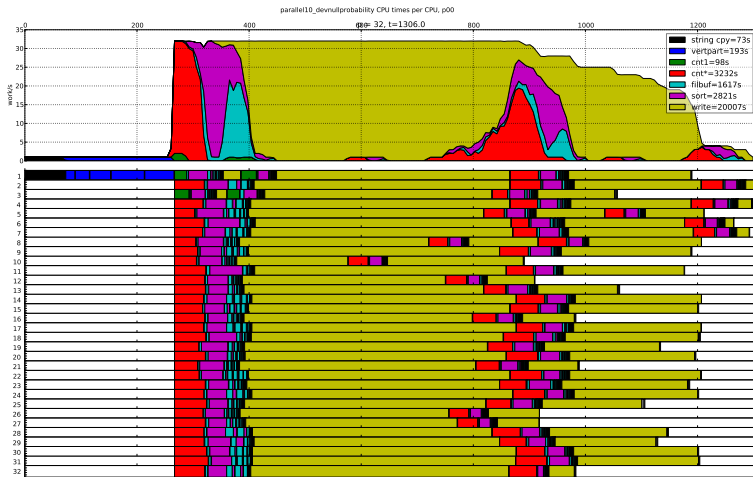Check in more detail what the processes are doing.
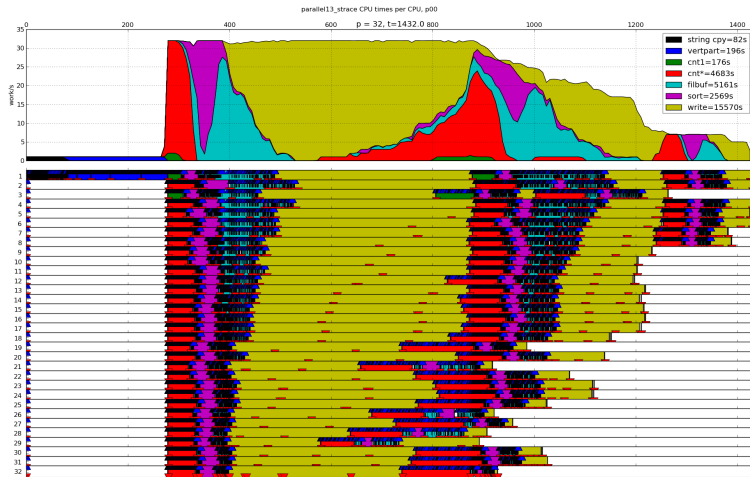
# Results – 5 ($p = 16$)

Introduction
oooo

Design goals
oo

ERA
oooooo

Formal analysis
oooooooooooooo

Empirical evaluation
oooooooooooooooo●oo

Conclusion

## Results – 5 ($p = 16$), strace



parallel13_strace CPU times per CPU, p00

# Results – 5 ($p = 32$)

## Results – 5 ($p = 32$), strace

Introduction
0000

Design goals
00

ERA
000000

Formal analysis
000000000000

Empirical evaluation
000000000000000000

Conclusion

## Conclusion

## Conclusion

- Huge gap between the theoretical time and I/O asymptotically tight algorithms and the practical ones.

## Conclusion

- Huge gap between the theoretical time and I/O asymptotically tight algorithms and the practical ones.
- ERA despite being practically the fastest algorithm is **not theoretically tight** even for random input strings with uniform substring distribution.

## Conclusion

- Huge gap between the theoretical time and I/O asymptotically tight algorithms and the practical ones.
- ERA despite being practically the fastest algorithm is **not theoretically tight** even for random input strings with uniform substring distribution.

Open challenges:

## Conclusion

- Huge gap between the theoretical time and I/O asymptotically tight algorithms and the practical ones.
- 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).

## Conclusion

- Huge gap between the theoretical time and I/O asymptotically tight algorithms and the practical ones.
- 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?

## Conclusion

- Huge gap between the theoretical time and I/O asymptotically tight algorithms and the practical ones.
- 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.

Introduction
oooo
Design goals
oo
ERA
oooooo
Formal analysis
ooooooooooooo
Empirical evaluation
oooooooooooooooooo
Conclusion

## Thank you.



**Matevž Jekovec**
matevz.jekovec@fri.uni-lj.si


University of Ljubljana
Faculty of Computer and
Information Science

**Andrej Brodnik**
andrej.brodnik@fri.uni-lj.si


Laboratorij za vseprisotne sisteme
Laboratory for Ubiquitous SYstems
http://lusy.fri.uni-lj.si

## Extra — List of all experiments

Shown at the presentation:

- Execution time for different $p$, speed-up, efficiency.
- CPU times, `iostat` and `mpstat` per different phases for $p = \{1, 2, 3, 4, 6, 8, 12, 16, 20, 24, 27, 32\}$.
- `iostat` output for various $p$.
- `mpstat` output for various $p$.
- Work per core for single execution for $p = 32$.
- Using `strace`, fetching `read`, `write`, `lseek` syscalls.

Introduction
0000

Design goals
00

ERA
000000

Formal analysis
000000000000

Empirical evaluation
000000000000000

Conclusion

## Extra — List of all experiments contd.

Test scenarios:

- Original code + added fsync(), various # of cores, various mem. size per core.
- Various string buffer sizes BUF_TYPE= $\{8, 16, 32, 64\}$ bit.
- Integration of Multikey cached quicksort (Rantala-Bentley-Sedgewick) instead of GNU qsort.
- Maximum limit of simultaneously opened files for writing $F = \{1, 2, 3, 4, 5, 6, 12, 16\}$.
- Output to /dev/null with probability $Pr = [0, 0.1...1]$.
- Separated disk for writing and reading.
- SSD for reading and/or writing.
- Different file system schedulers: noop, default, cfq.
- Different file system max queue length.
- Output to raw device without file system.
- Execution on 12x Raspberry $\pi$ with shared NFS storage.