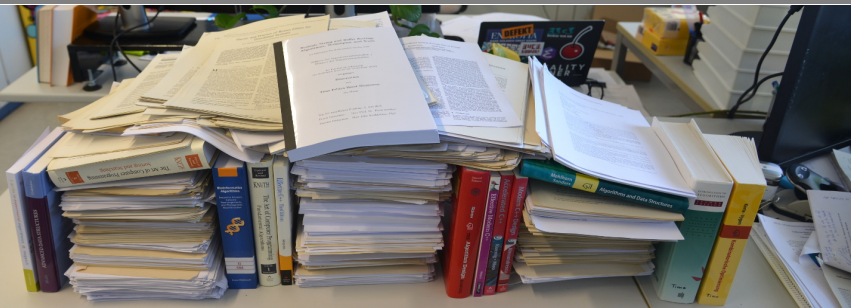# Scalable String and Suffix Sorting: Algorithms, Techniques, and Tools
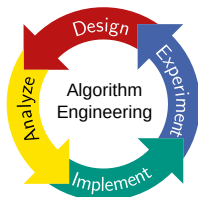
Timo Bingmann · Dissertation Defense · July 3rd, 2018

# Overview



**Multi-Core
Scalable String Sorting**

**External and Distributed
Scalable Suffix Sorting**

# Sorting Strings



Input: *n* strings containing *N* characters in total.

# String Sorting Algorithms

**Theoretical Parallel Algorithms**

- "Optimal Parallel String Algorithms: ..."                                      [Hagerup '94]
  $\mathcal{O}(\log N / \log \log N)$ time and $\mathcal{O}(N \log \log N)$ work on CRCW PRAM

**Existing Basic Sequential Algorithms**

- Radix Sort                          $\mathcal{O}(D + n \log \sigma)$              [McIlroy et al. '95]
- Multikey Quicksort          $\mathcal{O}(D + n \log n)$ exp.         [Bentley, Sedgewick '97]
- Burstsort                          $\mathcal{O}(D + n \log \sigma)$ exp.              [Sinha, Zobel '04]
- Binary LCP-Mergesort     $\mathcal{O}(D + n \log n)$                      [Ng, Kakehi '08]

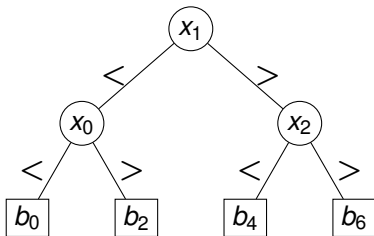**Existing Algorithm Library**

- in C/C++ by Rantala   (for Engineering Radix Sort [Kärkkäinen, Rantala '09])

**Our Contributions: New Basic and Practical Parallel Algorithms**

- Parallel Super Scalar String Sample Sort (pS$^5$)         [B, Sanders, ESA'13]
- Parallel $K$-way LCP-aware Mergesort (and Merge)   [B, et al. Algorithmica'17]

# Super Scalar String Sample Sort (S$^5$)



array0
kit0
arrange0
kayak0
kernel0
kitchen0
kitten0
arcade0
kite0
abacus0
krypton0
alpha0
arcane0

based on Super Scalar Sample Sort
[Sanders, Winkel '04]

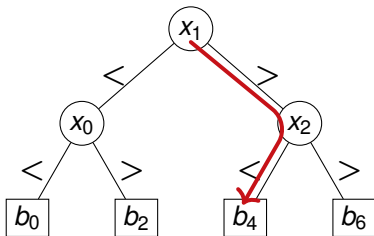# Super Scalar String Sample Sort (S$^5$)



```
a r r a y 0
k i t 0
a r r a n g e 0
k a y a k 0
k e r n e l 0
k i t c h e n 0
k i t t e n 0
a r c a d e 0
k i t e 0
a b a c u s 0
k r y p t o n 0
a l p h a 0
a r c a n e 0
```

based on Super Scalar Sample Sort
[Sanders, Winkel '04]

# Super Scalar String Sample Sort (S$^5$)



| a | r | r | a | y | 0 |
| k | i | t | 0 |
| a | r | r | a | n | g | e | 0 |
| k | a | y | a | k | 0 |
| k | e | r | n | e | l | 0 |
| k | i | t | c | h | e | n | 0 |
| k | i | t | t | e | n | 0 |
| a | r | c | a | d | e | 0 |
| k | i | t | e | 0 |
| a | b | a | c | u | s | 0 |
| k | r | y | p | t | o | n | 0 |
| a | l | p | h | a | 0 |
| a | r | c | a | n | e | 0 |

- partition by $w$ chars
- store in level-order and
  use predicated instructions

|   | 1 | 2 | 3 |
|---|---|---|---|
|   | ar | ab | ki |

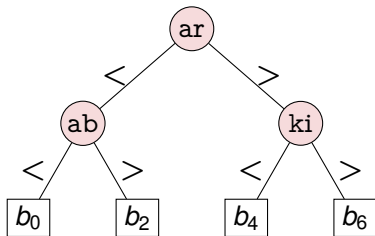$i := 2i + 0/1$

# Super Scalar String Sample Sort ($S^5$)

```
a r r a y 0
k i t 0
a r r a n g e 0
k a y a k 0
k e r n e l 0
k i t c h e n 0
k i t t e n 0
a r c a d e 0
k i t e 0
a b a c u s 0
k r y p t o n 0
a l p h a 0
a r c a n e 0
```
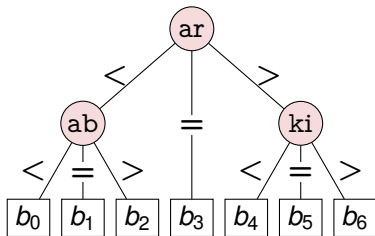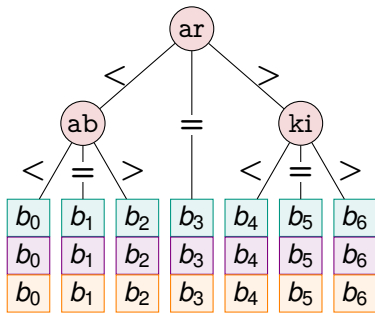
- equality checking:
  1. at each splitter
  2. after full descent

- interleave tree descents:
  classify four strings at once
  $\Rightarrow$ super scalar parallelism

# Super Scalar String Sample Sort (S$^5$)



- easy parallelization
- classification tree in
  L2 caches of processors

# Super Scalar String Sample Sort (S⁵)



prefix

| | |
|---|---|
| a b a c u s 0 | 2 |
| a l p h a 0 | 1 |
| a r r a y 0 | |
| a r r a n g e 0 | 2 |
| a r c a d e 0 | |
| a r c a n e 0 | |
| k a y a k 0 | |
| k e r n e l 0 | 0 |
| k i t 0 | |
| k i t c h e n 0 | |
| k i t t e n 0 | 2 |
| k i t e 0 | |
| k r y p t o n 0 | 0 |

increase prefix by
LCP of splitters
or key size.

- reorder out-of-place, in-place, and/or in parallel
- top-level algorithm in parallel S⁵

# LCP Loser Tree – *K*-way LCP-Merge



needs at most
$\Delta L + n \log_2 K + K$
character comparisons

$(2, \text{aab})$

$(1, \text{acb})$

$(2, \text{aac})$   LCP-Merge   $(0, \text{bca})$

$(2, \text{aab})$   $(2, \text{aac})$   $(0, \text{bca})$   $(1, \text{acb})$

# LCP Loser Tree – *K*-way LCP-Merge



needs at most
$\Delta L + n \log_2 K + K$
character comparisons

# Contributed String Sorting Algorithms

- Parallel Super Scalar String Sample Sort (pS$^5$)
    - fully parallel S$^5$, sequential S$^5$, and fast base case sorters
    - sequential running time of S$^5$:
      $\mathcal{O}(\frac{D}{w} + n \log n)$ expected time with equality checks, and
      $\mathcal{O}((\frac{D}{w} + n) \log v + n \log n)$ expected time with unrolled descents.
    - parallel running time of a single step of fully parallel S$^5$:
      $\mathcal{O}(\frac{n}{p} \log v + \log p)$ time and $\mathcal{O}(n \log v + pv)$ work.

- Hybrid NUMA-aware pS$^5$ + $K$-way LCP-Merge

- Parallel Multikey Quicksort

- Parallel Radix Sort (Adaptive 16-bit and 8-bit)

## Additional Algorithms:

- (Parallel) Multiway LCP-aware Mergesort $\qquad \mathcal{O}(D + n \log n + \frac{n}{K})$
- Sequential LCP-aware Insertion Sort $\qquad\qquad\qquad \mathcal{O}(D + n^2)$

# 128 GiB GOV2 – Speedup on 32-Core Intel



Input characteristics: $n = 3.1$ G, $N = 128$ Gi, $\frac{D}{N} = 82.7\,\%$.

Timo Bingmann – Scalable String and Suffix Sorting: Algorithms, Techniques, and Tools
Institute of Theoretical Informatics – Algorithmics

July 3rd, 2018        14 / 29

# Overview

**Multi-Core
Scalable String Sorting**

**External and Distributed
Scalable Suffix Sorting**



| ⊥ | a l p h a 0 |
| 1 | a r c a d e 0 |
| 2 | a r r a y 0 |
| 0 | k a y a k 0 |
| 1 | k e r n e l 0 |
| 1 | k i t 0 |
| 3 | k i t c h e n 0 |
| 3 | k i t t e n 0 |
| 1 | k r y p t o n 0 |

Algorithm
Engineering

Design
Experiment
Implement
Analyze

| ⊥ | $ |
| 0 | a $ |
| 1 | a c b a $ |
| 4 | a c b a c b a $ |
| 0 | b a $ |
| 2 | b a c b a $ |
| 5 | b a c b a c b a $ |
| 0 | c b a $ |
| 3 | c b a c b a $ |

# Example $T = \begin{bmatrix} 0\ 1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9\ 10\ 11\ 12\ 13 \\ \texttt{t o b e o r n o t t o b e \$} \end{bmatrix}$

| i | $T_i$ |
|---|---|
| 0 | t o b e o r n o t t o b e \$ |
| 1 | o b e o r n o t t o b e \$ |
| 2 | b e o r n o t t o b e \$ |
| 3 | e o r n o t t o b e \$ |
| 4 | o r n o t t o b e \$ |
| 5 | r n o t t o b e \$ |
| 6 | n o t t o b e \$ |
| 7 | o t t o b e \$ |
| 8 | t t o b e \$ |
| 9 | t o b e \$ |
| 10 | o b e \$ |
| 11 | b e \$ |
| 12 | e \$ |
| 13 | \$ |

# Example $T = \big[$ t o b e o r n o t t o b e \$ $\big]$

$\quad$ 0 1 2 3 4 5 6 7 8 9 10 11 12 13

| $SA_i$ | $LCP_i$ | $T_{SA_i \ldots n}$ |
|---|---|---|
| 13 | - | \$ |
| 11 | 0 | b e \$ |
| 2 | 2 | b e o r n o t t o b e \$ |
| 12 | 0 | e \$ |
| 3 | 1 | e o r n o t t o b e \$ |
| 6 | 0 | n o t t o b e \$ |
| 10 | 0 | o b e \$ |
| 1 | 3 | o b e o r n o t t o b e \$ |
| 4 | 1 | o r n o t t o b e \$ |
| 7 | 1 | o t t o b e \$ |
| 5 | 0 | r n o t t o b e \$ |
| 9 | 1 | t o b e \$ |
| 0 | 4 | t o b e o r n o t t o b e \$ |
| 8 | 1 | t t o b e \$ |

Prefix Doubling

Induced Copying

Recursion

| | | |
|---|---|---|
| MM90 original | BW94 BWT | Far97 $\mathcal{O}(n)$ tree |
| LS99 qsufsort | | |
| | Sew00 1/2 copy / IT99 A/B copy | |
| | | BK03 diffcover / KS03 DC / KSPP03 mod2 split |
| | MF02 deep-shallow | KA03 L/S split / HSS03 mod2 |
| SS05 bpr | | KJP04 fixed $\Sigma$ |
| | Mor06 divsufsort / MP06 ISA chains | NZ07 $\mathcal{O}(n \log |\Sigma|)$ |
| | MP08 cache aware | |
| | NZC09a SA-IS / NZC09b SA-DS | AN10 SFE-coding |
| | Non13 SACA-K | |
| | Bai16 prev ptr | |
| | LLH16 $\mathcal{O}(1)$ space / Got17 $\mathcal{O}(1)$ space | |

1999
2000

2003

2004

2005
2006
2008

2009

2010

2013

2016

2017

**External Memory Algorithms**

MM90
original

LS99
qsufsort

prefix doubling – [CF02]
$\mathcal{O}(\text{sort}(n) \cdot \log(\text{maxlcp}))$ I/Os

KS03
DC

prefix doubling – [DKMS05]
$\text{sort}(5n) \cdot \log_2(\text{maxlcp}) + \mathcal{O}(\text{sort}(n))$ I/Os

DC7 – [DKMS05]
$\text{sort}(24.75n) + \text{scan}(6n)$ I/Os

NZC09a
SA-IS

eSAIS – [B, Fischer, Osipov, ALENEX'13, JEA'16]
$\text{sort}(17n) + \text{scan}(9n)$ I/Os

fSAIS – [KKPZ17]
$\mathcal{O}(\text{scan}(n) \cdot \log^2_{M/B}(n/B))$ I/Os

**ST xXL**

1999
2000

2003

2004

2005
2006
2008

2009

2010

2013
2016

2017

amazon
web services

Microsoft Azure

Google Cloud Platform

bwUniCluster
© KIT (SCC)

# Big Data Batch Processing



Fast

Performance

Slow

New Framework:

Thrill

MPI

Our Requirements:
- compound primitives into complex algorithms
- efficient simple data types,
- overlap computation and communication,
- automatic disk usage,
- C++, and much more...

Apache Spark / Apache Flink

MapReduce Hadoop

Low Level
Difficult

Interface
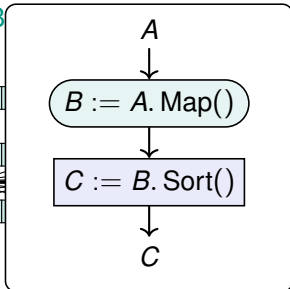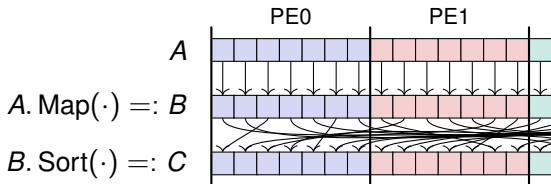
High Level
Simple

# Distributed Immutable Array (DIA)

**User Programmer's View:**

- DIA<T> = distributed array of items T on the cluster
- Cannot access items directly, instead use small set of scalable primitives, for example: Map, Sort, ReduceByKey, Zip, Window, etc.

# Distributed Immutable Array (DIA)

**User Programmer's View:**

- $DIA<T>$ = distributed array of items T on the cluster
- Cannot access items directly, instead use small set of scalable primitives, for example: Map, Sort, ReduceB

# Distributed Immutable Array (DIA)

**User Programmer's View:**

- DIA<T> = distributed array of items T on the cluster
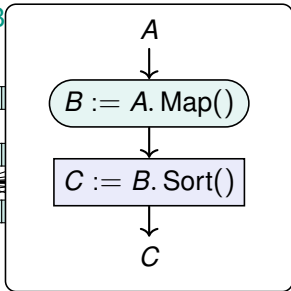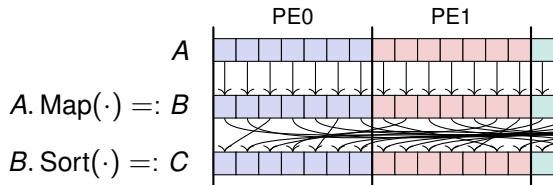- Cannot access items directly, instead use small set of scalable primitives, for example: Map, Sort, ReduceB...



**Framework Designer's View:**

- Goals: distribute work, optimize execution on cluster, add redundancy where applicable. $\Longrightarrow$ build data-flow graph.
- DIA<T> = pipelined chain of computations

# Thrill's Goal and Current Status

An easy way to program fast distributed algorithms in C++.

**Current Status:**
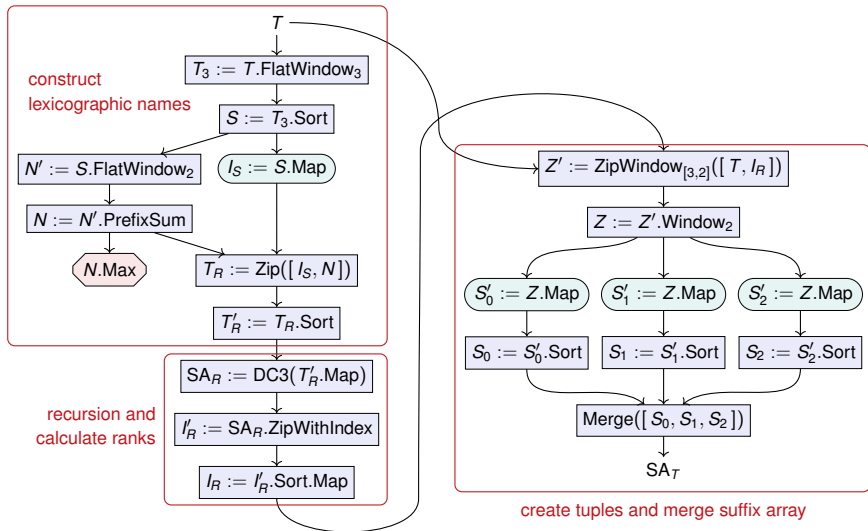
- Open-source prototype at http://github.com/thrill/thrill.
- $\approx 60\,$K lines of C++14 code, 70–80 % written by B, $\geq 12$ contributors
- Published at IEEE Conference on Big Data          [B, et al. '16]
- Faster than Apache Spark and Apache Flink on five micro benchmarks: WordCount1000, WordCountCC, PageRank, TeraSort, and K-Means.
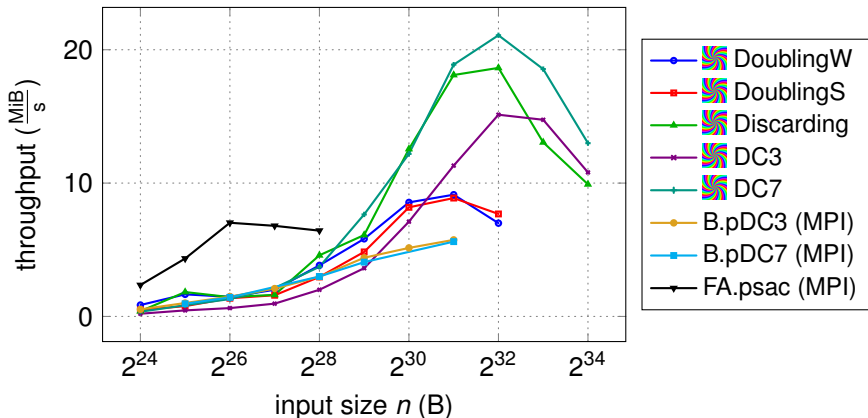
**Case Studies:**

- Five suffix sorting algorithms          [B, Gog, Kurpicz, arXiv'17]
- Louvain graph clustering algorithm          [Hamann et al. arXiv'17]
- More examples: stochastic gradient descent, triangle counting, etc.
- Future: fault tolerance, scalability, and more applications.

# Data-Flow Graph of DC3 with Recursion



construct lexicographic names

$T$

$T_3 := T.\text{FlatWindow}_3$

$S := T_3.\text{Sort}$

$N' := S.\text{FlatWindow}_2$

$I_S := S.\text{Map}$

$N := N'.\text{PrefixSum}$

$N.\text{Max}$

$T_R := \text{Zip}([\,I_S, N\,])$

$T'_R := T_R.\text{Sort}$

recursion and calculate ranks

$SA_R := \text{DC3}(T'_R.\text{Map})$

$I'_R := SA_R.\text{ZipWithIndex}$

$I_R := I'_R.\text{Sort.Map}$

$Z' := \text{ZipWindow}_{[3,2]}([\,T, I_R\,])$

$Z := Z'.\text{Window}_2$

$S'_0 := Z.\text{Map}$

$S'_1 := Z.\text{Map}$

$S'_2 := Z.\text{Map}$

$S_0 := S'_0.\text{Sort}$

$S_1 := S'_1.\text{Sort}$

$S_2 := S'_2.\text{Sort}$

$\text{Merge}([\,S_0, S_1, S_2\,])$

$SA_T$

create tuples and merge suffix array

Timo Bingmann – Scalable String and Suffix Sorting: Algorithms, Techniques, and Tools
Institute of Theoretical Informatics – Algorithmics

July 3rd, 2018          27 / 29

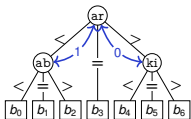# Suffix Sorting Wikipedia with 32 Hosts



Run on $32 \times$ i3.4xlarge AWS EC2 instances containing 16-core Intel Xeon E5-2686 CPUs with 2.30 GHz, 8 GB of RAM, and $2 \times 1.9$ TB NVMe SSDs.
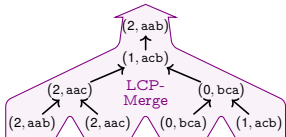
# Overview: Main Contributions

**Multi-Core
Scalable String Sorting**

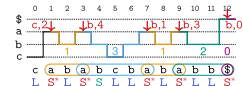- Parallel Super Scalar String Sample Sort (pS⁵) [BS13]



- Parallel Multiway LCP-Merge, Merge Sort, and More [BES17]



**External and Distributed
Scalable Suffix Sorting**

- Induced Sorting in External Memory: eSAIS [BFO13, BFO16]



- New High-Performance Distributed Framework in C++: Thrill [BAJ+16]



- Distributed External Suffix Sorting