# Faster Block Tree Construction

**31st European Symposium on Algorithms (ESA 2023)**

Dominik Köppl, *Florian Kurpicz*, and Daniel Meyer

# Motivation

- Git repositories 
- DNA 
- proteins 
- user generated content $W$
- XML data 
- book collections 
- . . .

highly repetitive
&
huge input

Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023          Institute of Theoretical Informatics, Algorithm Engineering

- Git repositories
- DNA
- proteins
- user generated content $\mathbb{W}$
- XML data
- book collections
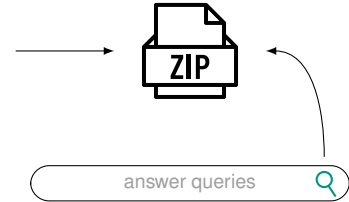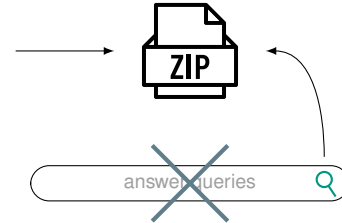- ...

highly repetitive
&
huge input

$\longrightarrow$

ZIP

# Motivation



- Git repositories
- DNA
- proteins
- user generated content $\mathbb{W}$
- XML data
- book collections
- . . .

highly repetitive
&
huge input

→ ZIP ←

answer queries

Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023          Institute of Theoretical Informatics, Algorithm Engineering

# Motivation



- Git repositories
- DNA
- proteins
- user generated content $\mathbb{W}$
- XML data
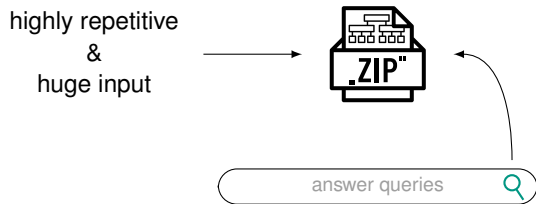- book collections
- . . .

highly repetitive
&
huge input

ZIP

answer queries

Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023    Institute of Theoretical Informatics, Algorithm Engineering

# Compressed Self-Indices

highly repetitive
&
huge input

Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023    Institute of Theoretical Informatics, Algorithm Engineering

# Compressed Self-Indices



highly repetitive
&
huge input

# Compressed Self-Indices

# Compressed Self-Indices



Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023                Institute of Theoretical Informatics, Algorithm Engineering

# Compressed Self-Indices
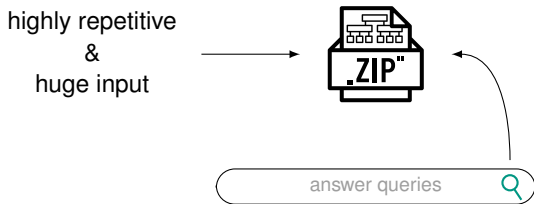


highly repetitive
&
huge input

answer queries 🔍

- access
- rank
- select
- . . .
- pattern matching

**Wavelet Tree** (de-facto standard in practice)

- $T \in [1, \sigma]^n$
- access, rank, select: $O(\log \sigma)$ time
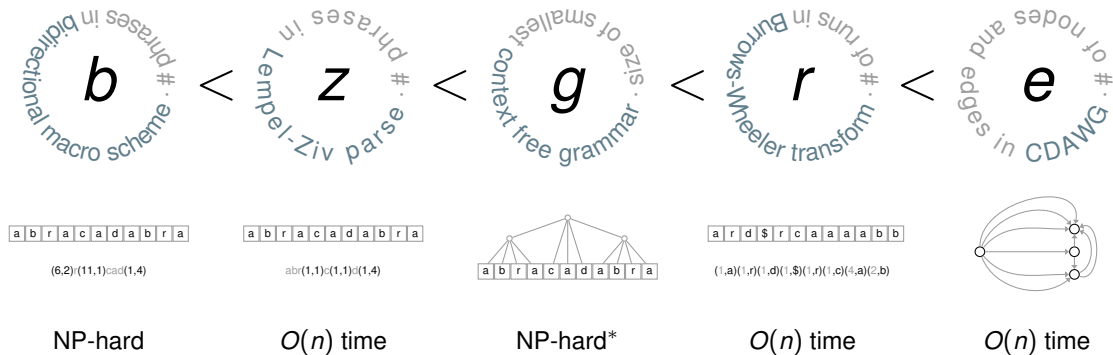- $nH_0(T) + o(n)$ bits space

# Compressed Self-Indices



highly repetitive
&
huge input

answer queries

- access
- rank
- select
- …
- pattern matching

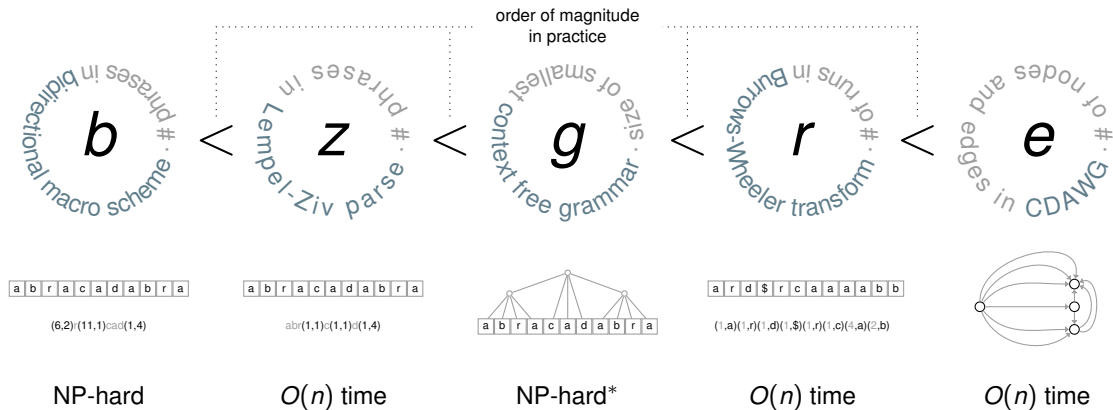**Wavelet Tree** (de-facto standard in practice)

- $T \in [1, \sigma]^n$
- access, rank, select: $O(\log \sigma)$ time
- $nH_0(T) + o(n)$ bits space
- **blind for repetitions**
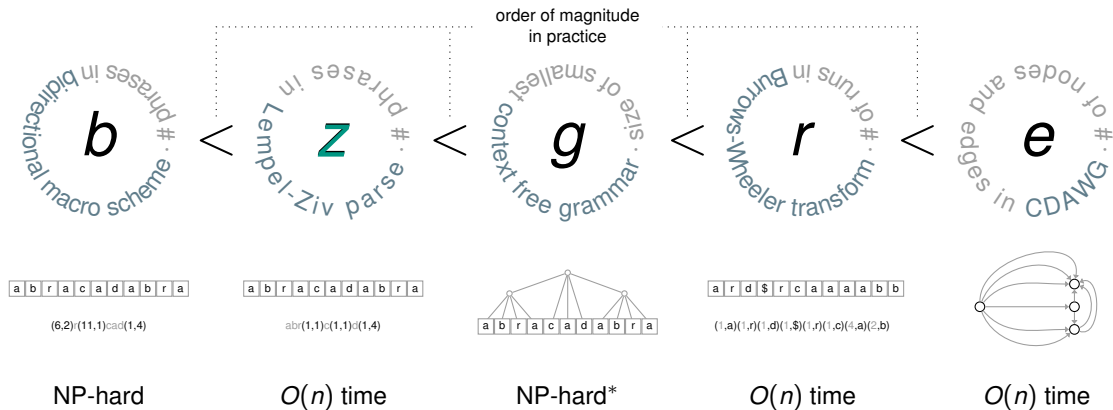
# Measures of Repetitiveness (Excerpt)

$b$ · #·phrases in bidirectional macro scheme

$<$

$z$ · #·phrases in Lempel-Ziv parse

$<$

$g$ · size of smallest context free grammar

$<$

$r$ · #·of runs in Burrows-Wheeler transform

$<$

$e$ · #·of nodes and edges in CDAWG

| a | b | r | a | c | a | d | a | b | r | a |

(6,2)·(11,1)cad(1,4)

| a | b | r | a | c | a | d | a | b | r | a |

abr(1,1)c(1,1)d(1,4)

| a | b | r | a | c | a | d | a | b | r | a |

| a | r | d | $ | r | c | a | a | a | a | b | b |

(1,a)(1,r)(1,d)(1,$)(1,r)(1,c)(4,a)(2,b)

NP-hard

$O(n)$ time

NP-hard*

$O(n)$ time

$O(n)$ time

* there are good heuristics

order of magnitude in practice

$b$ — phrases in bidirectional macro scheme

$z$ — phrases in Lempel-Ziv parse

$g$ — size of smallest context free grammar

$r$ — # of runs in Burrows-Wheeler transform

$e$ — # of nodes and edges in CDAWG

$b < z < g < r < e$

| a | b | r | a | c | a | d | a | b | r | a |

(6,2)·(11,1)cad(1,4)

NP-hard

| a | b | r | a | c | a | d | a | b | r | a |

abr(1,1)c(1,1)d(1,4)

$O(n)$ time

a b r a c a d a b r a

NP-hard*

| a | r | d | $ | r | c | a | a | a | a | b | b |

(1,a)(1,r)(1,d)(1,$)(1,r)(1,c)(4,a)(2,b)

$O(n)$ time

$O(n)$ time

\* there are good heuristics

# Measures of Repetitiveness (Excerpt)



order of magnitude in practice

$b$ — # · phrases in bidirectional macro scheme

$z$ — # · phrases in Lempel-Ziv parse

$g$ — size of smallest context free grammar

$r$ — # · of runs in Burrows-Wheeler transform

$e$ — # · of nodes and edges in CDAWG

$b < z < g < r < e$

| a | b | r | a | c | a | d | a | b | r | a |

(6,2)·(11,1)cad(1,4)

| a | b | r | a | c | a | d | a | b | r | a |

abr(1,1)c(1,1)d(1,4)

| a | b | r | a | c | a | d | a | b | r | a |

a r d $ r c a a a a b b

(1,a)(1,r)(1,d)(1,$)(1,r)(1,c)(4,a)(2,b)

NP-hard          $O(n)$ time          NP-hard*          $O(n)$ time          $O(n)$ time

* there are good heuristics

# The Block Tree



- root has degree $s = 4$
- other nodes have degree $\tau = 2$
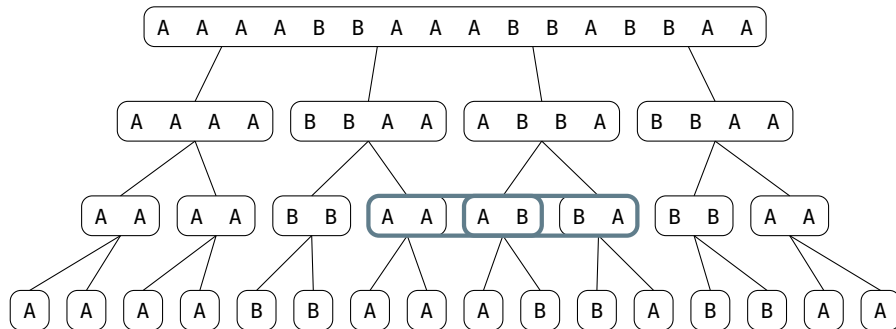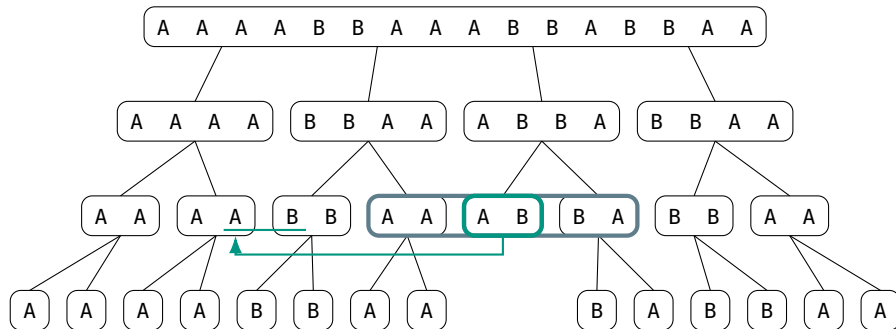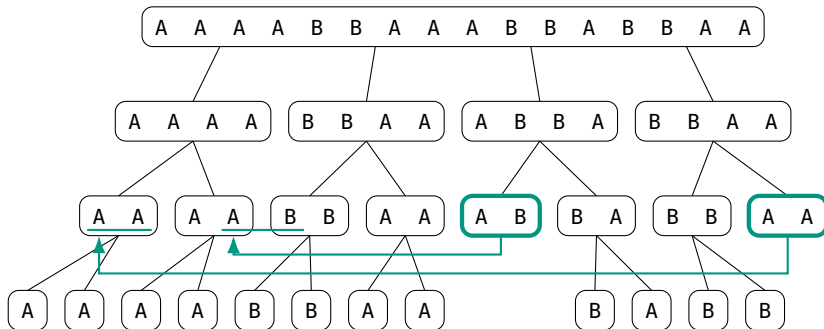- all levels (except the first) have $\leq 3z\tau$ blocks

- root has degree $s = 4$
- other nodes have degree $\tau = 2$
- all levels (except the first) have $\leq 3z\tau$ blocks

- keep $B_i$ if consecutive $B_{i-1} \cdot B_i$ or $B_i \cdot B_{i+1}$ are leftmost occurrences
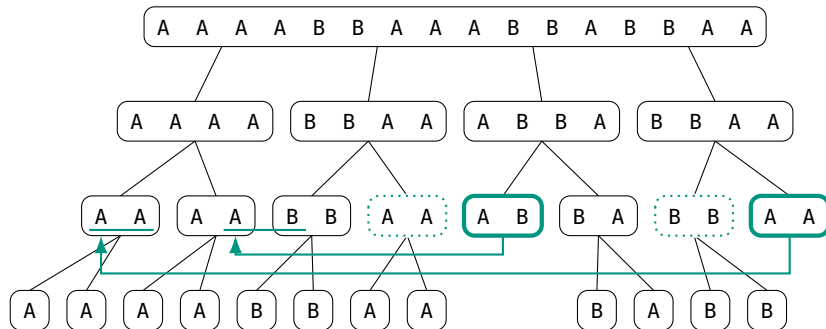
# The Block Tree



- root has degree $s = 4$
- other nodes have degree $\tau = 2$
- all levels (except the first) have $\leq 3z\tau$ blocks

- keep $B_i$ if consecutive $B_{i-1} \cdot B_i$ or $B_i \cdot B_{i+1}$ are leftmost occurrences

- root has degree $s = 4$
- other nodes have degree $\tau = 2$
- all levels (except the first) have $\leq 3z\tau$ blocks

- keep $B_i$ if consecutive $B_{i-1} \cdot B_i$ or $B_i \cdot B_{i+1}$ are leftmost occurrences

Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023

Institute of Theoretical Informatics, Algorithm Engineering

# The Block Tree



- root has degree $s = 4$
- other nodes have degree $\tau = 2$
- all levels (except the first) have $\leq 3z\tau$ blocks

- keep $B_i$ if consecutive $B_{i-1} \cdot B_i$ or $B_i \cdot B_{i+1}$ are leftmost occurrences

- root has degree $s = 4$
- other nodes have degree $\tau = 2$
- all levels (except the first) have $\leq 3z\tau$ blocks

- keep $B_i$ if consecutive $B_{i-1} \cdot B_i$ or $B_i \cdot B_{i+1}$ are leftmost occurrences

# The Block Tree



- root has degree $s = 4$
- other nodes have degree $\tau = 2$
- all levels (except the first) have $\leq 3z\tau$ blocks

- keep $B_i$ if consecutive $B_{i-1} \cdot B_i$ or $B_i \cdot B_{i+1}$ are leftmost occurrences

- root has degree $s = 4$
- other nodes have degree $\tau = 2$
- all levels (except the first) have $\leq 3z\tau$ blocks

- keep $B_i$ if consecutive $B_{i-1} \cdot B_i$ or $B_i \cdot B_{i+1}$ are leftmost occurrences

# The Block Tree



- root has degree $s = 4$
- other nodes have degree $\tau = 2$
- all levels (except the first) have $\leq 3z\tau$ blocks

- keep $B_i$ if consecutive $B_{i-1} \cdot B_i$ or $B_i \cdot B_{i+1}$ are leftmost occurrences
- remove more blocks with pruning

# State-of-Block-Tree-Construction

| Method | Reference | Working Space | Time | Implementation |
|--------|-----------|---------------|------|----------------|
| Aho-Corasic | [Bel+21] | $O(n)$ | $O(n(1 + \log_\tau(z\tau/s)))$ | no |
| Fingerprints | [Bel+21] | $O(s + z\tau \log_\tau(\frac{n \log \sigma}{s \log n}))$ | $O(n(1 + \log_\tau(z\tau/s)))$ expected | yes (slow) |
| LPF Array | [here] | $O(n)$ | $O(n(1 + \log_\tau(z\tau/s)))$ | yes (fast) |

Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023    Institute of Theoretical Informatics, Algorithm Engineering

A   A   A   A   B   B   A   A   A   B   B   A   B   B   A   A

Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023        Institute of Theoretical Informatics, Algorithm Engineering

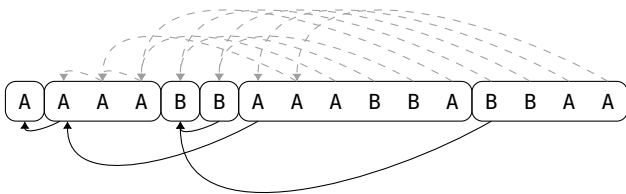# Lempel-Ziv Parse
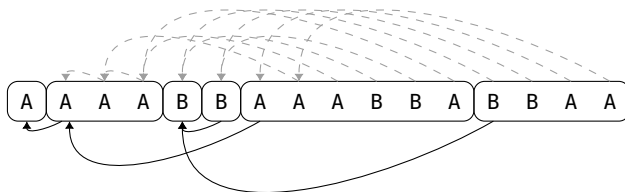
# Lempel-Ziv Parse

# Lempel-Ziv Parse

# Lempel-Ziv Parse

# Lempel-Ziv Parse



Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023          Institute of Theoretical Informatics, Algorithm Engineering

# Lempel-Ziv Parse



Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023    Institute of Theoretical Informatics, Algorithm Engineering

Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023

Institute of Theoretical Informatics, Algorithm Engineering

# Our Algorithm (Marking of Nodes)

# Our Algorithm (Marking of Nodes)

# Our Algorithm (Marking of Nodes)

# Our Algorithm (Marking of Nodes)

# Our Algorithm (Marking of Nodes)

# Our Algorithm (Marking of Nodes)



Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023                    Institute of Theoretical Informatics, Algorithm Engineering

# Our Algorithm (Marking of Nodes)



Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023    Institute of Theoretical Informatics, Algorithm Engineering

# Our Algorithm (Marking of Nodes)

Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023

Institute of Theoretical Informatics, Algorithm Engineering

# Experimental Evaluation

- highly tuned implementation
- tree consists only of bit and compact vectors
- tuning parameter
  - degree root $s = \{1, z\}$ (only we have $s = z$)
  - degree other nodes $\tau = \{2, 4, 8, 16\}$
  - number characters in leaves $b = \{2, 4, 8, 16\}$

Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023     Institute of Theoretical Informatics, Algorithm Engineering
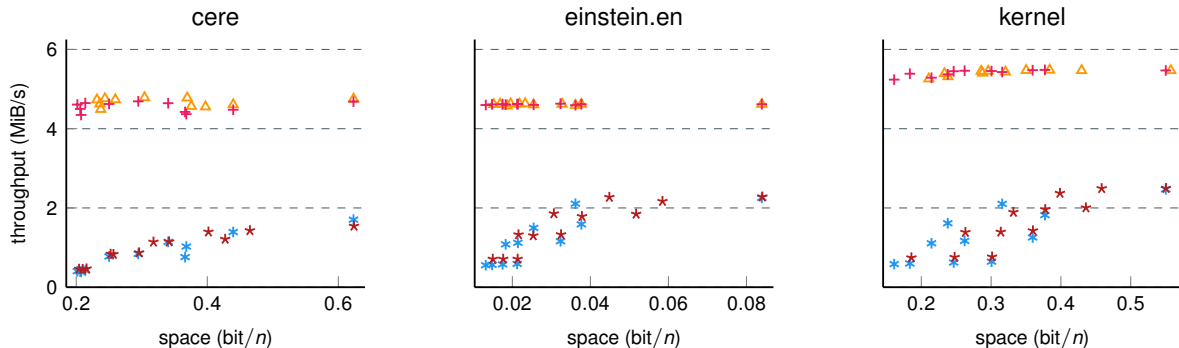
# Experimental Evaluation

- highly tuned implementation
- tree consists only of bit and compact vectors
- tuning parameter
  - degree root $s = \{1, z\}$ (only we have $s = z$)
  - degree other nodes $\tau = \{2, 4, 8, 16\}$
  - number characters in leaves $b = \{2, 4, 8, 16\}$

- original FP BT [Bel+21]
- our reimplementation of the original FP BT
- our LPF BT construction with $s = 1$ and $s = z$
- dynamic programming variants
- parallelization
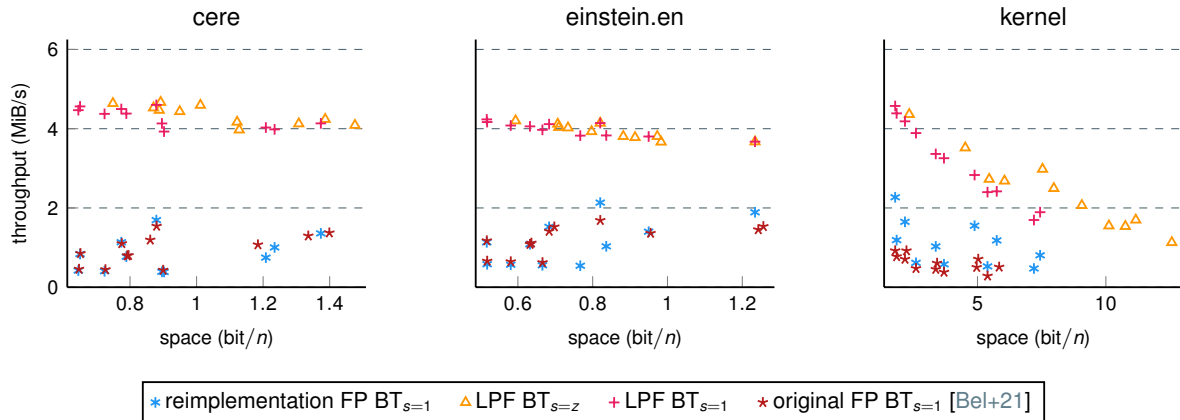- no comparison with wavelet trees (faster)

- repetitive instances from P&C corpus
- non-repetitive instances from P&C corpus

Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023    Institute of Theoretical Informatics, Algorithm Engineering

# Highly Repetitive Inputs (Access Only)



Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023    Institute of Theoretical Informatics, Algorithm Engineering

# Highly Repetitive Inputs (with Rank and Select Support)

Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023    Institute of Theoretical Informatics, Algorithm Engineering
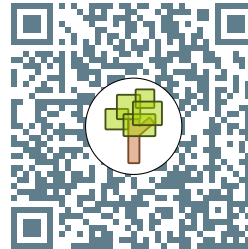
# Conclusion and Future Work



- fastest block tree construction algorithm
- first parallel block tree construction
- works in practice for non-repetitive inputs

- better scaling parallel construction
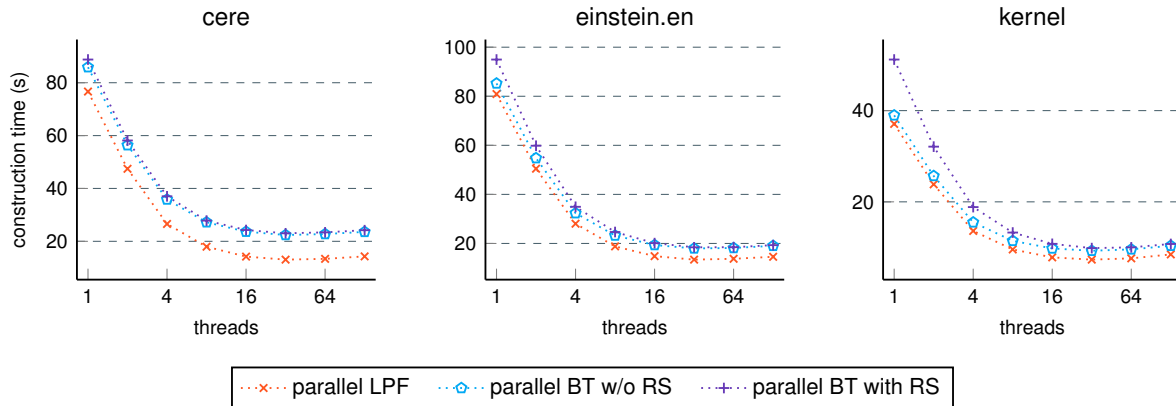- better marking rules (less pruning)



/pasta-toolbox/block_tree

# Parallel Construction (Strong Scaling)



cere     einstein.en     kernel

... parallel LPF    ... parallel BT w/o RS    ... parallel BT with RS

Köppl, *Kurpicz*, and Meyer | Faster Block Tree Construction | ESA 2023     Institute of Theoretical Informatics, Algorithm Engineering

# Bibliography

[Bel+21]  Djamal Belazzougui, Manuel Cáceres, Travis Gagie, Pawel Gawrychowski, Juha Kärkkäinen, Gonzalo Navarro, Alberto Ordóñez Pereira, Simon J. Puglisi, and Yasuo Tabei. "Block Trees". In: *J. Comput. Syst. Sci.* 117 (2021), pages 1–22. DOI: 10.1016/j.jcss.2020.11.002.