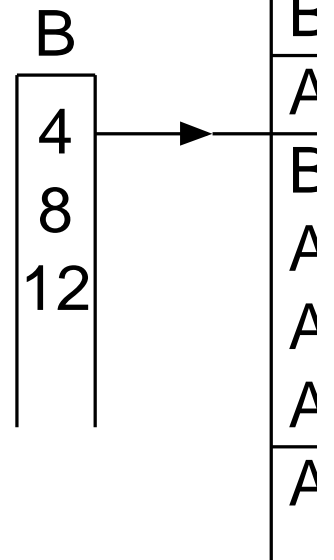


L RANK		
A	1	
B	1	A B
A	0	4 4
B	0	
A	1	
A	2	A B
A	0	7 5
A	1	



A new improvement is to use variable length encoding, where a frequent symbol uses less bits for symbol and more bits for rank.

Finally, we can trade time for space RANK with scanning from the next point [1].

L RANK

A	3	
B	3	
A	0	<div style="display: flex; flex-direction: column; align-items: center;"> <div style="margin-bottom: 5px;">A</div> <div style="margin-bottom: 5px;">4</div> <div style="margin-bottom: 5px;">8</div> <div style="margin-bottom: 5px;">12</div> </div>
B	0	
A	1	
A	2	
A	3	
A	0	

L RANK

A	0
B	0
A	1

re by replacing
arest reference

WAVELET TREES

Wavelet tree is a text rep both compressed and pre ries with little additiona with compressed text inc *eral rank queries*:

$$\text{RANK}_c(j) = |\{i \mid i \cdot$$

We need wavelet trees fo

$$\text{RANK}(j) = I$$

We use our own wavele optimized for special ran

We combine wavelet tre ranks, obtaining the *most*

presentation that can be reprocessed for rank queries in space. They are used as indexes [2] to answer general

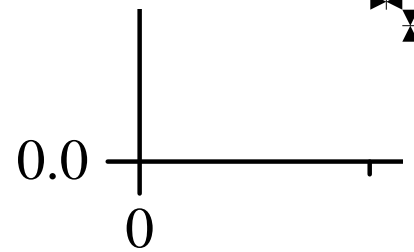
$$\{i < j \text{ and } L[i] = c\}.$$

for special rank queries:

$$\text{RANK}_{L[j]}(j).$$

Let tree implementations support rank queries.

Let us see with reference point a *space-efficient* algorithm.



REFERENCE

- [1] U. Lauther and A. R. Meyer. Algorithms for the rank of a matrix. In *Proc. 13th Annual Symposium on Combinatorics and Probability*, volume 1, pages 1-10. Springer, 2000.
- [2] G. Navarro and A. R. Meyer. Rank indexes. *ACM Computing Surveys*, 40(2):1-33, 2007.
- [3] J. Seward. Rank transform. In *Proceedings of the 1994 ACM SIGMOD International Conference on Management of Data*, pages 4-13, 1994.



Memory (bytes/symbol)

CES

and T. Lukovszki. Space efficient algorithm for the Burrows-Wheeler backtransformation. In *10th Annual European Symposium on Algorithms*, volume 3669 of *LNCS*, pages 293–304. Springer, 2005.

and V. Mäkinen. Compressed full-text indexing. *ACM Computing Surveys*, 39(1):Article 2, 2006.

Space-time tradeoffs in the inverse Burrows-Wheeler transform. In *Proc. IEEE Data Compression Conference*, pages 439–448. IEEE, 2001.

The basic inversion algorithm described has linear time and space complexity. It dominates the time and space required for decompression in programs like bzip2.

It is slow because each memory access during the permutation traversal is essential, causing many cache misses.

It needs a lot of space for the RANK array.

$$|\text{RANK}| = n \log n \text{ bits} = 4n \log n$$

$$|\text{text}| = n \log \sigma \text{ bits} = n \log \sigma$$

where n = text length and σ = alphabet size.

REFERENCE POINT RANKS

We reduce space by storing ranks relative to reference points, which can be placed at regular intervals.

Every k th position [1]

Every k th occurrence

described above
xity, but it still
uirements dur-
ke bzip2.

access during
ntially random

K array:

n bytes
bytes

phabet size.

S

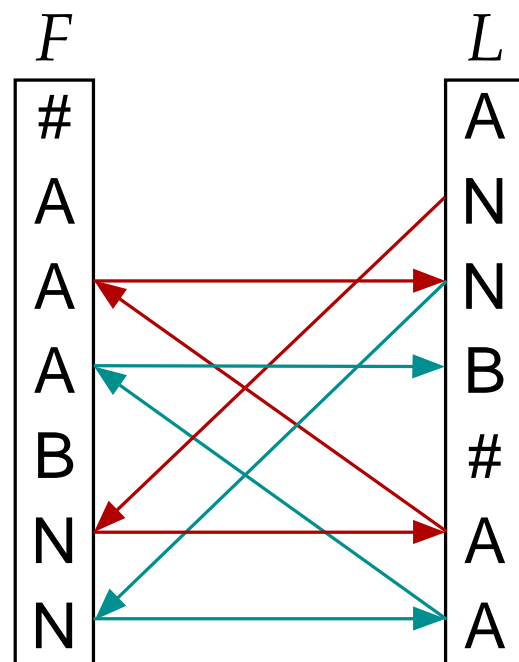
relative to ref-
d in two ways:

ccurrence [new]

REPETITION SHORTCUTS

Repetitions in the text may be used to find *paths* (PPP) in the inverse search. We can use this as follows.

1. On the **first pass**, observe the PPP (due to repeated ANA)



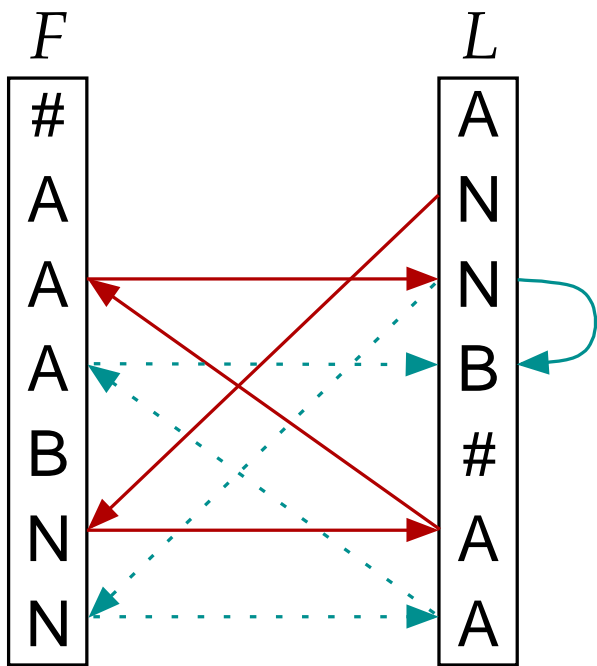
The shortcuts reduce the
This is the *fastest* known

WAVFI FT TRFFS

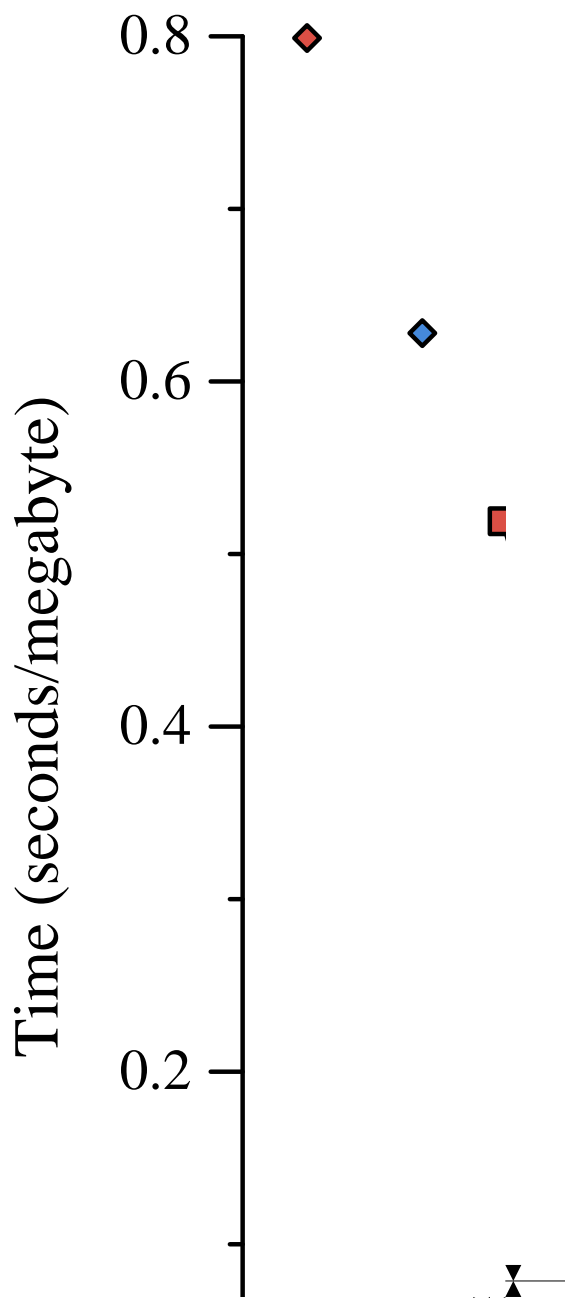
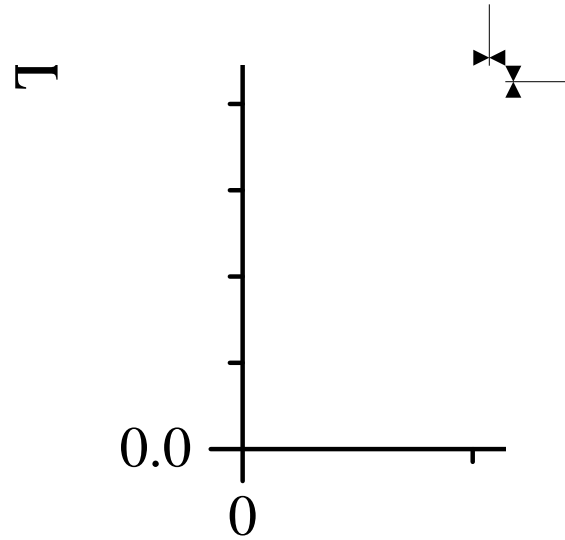
RT CUTS

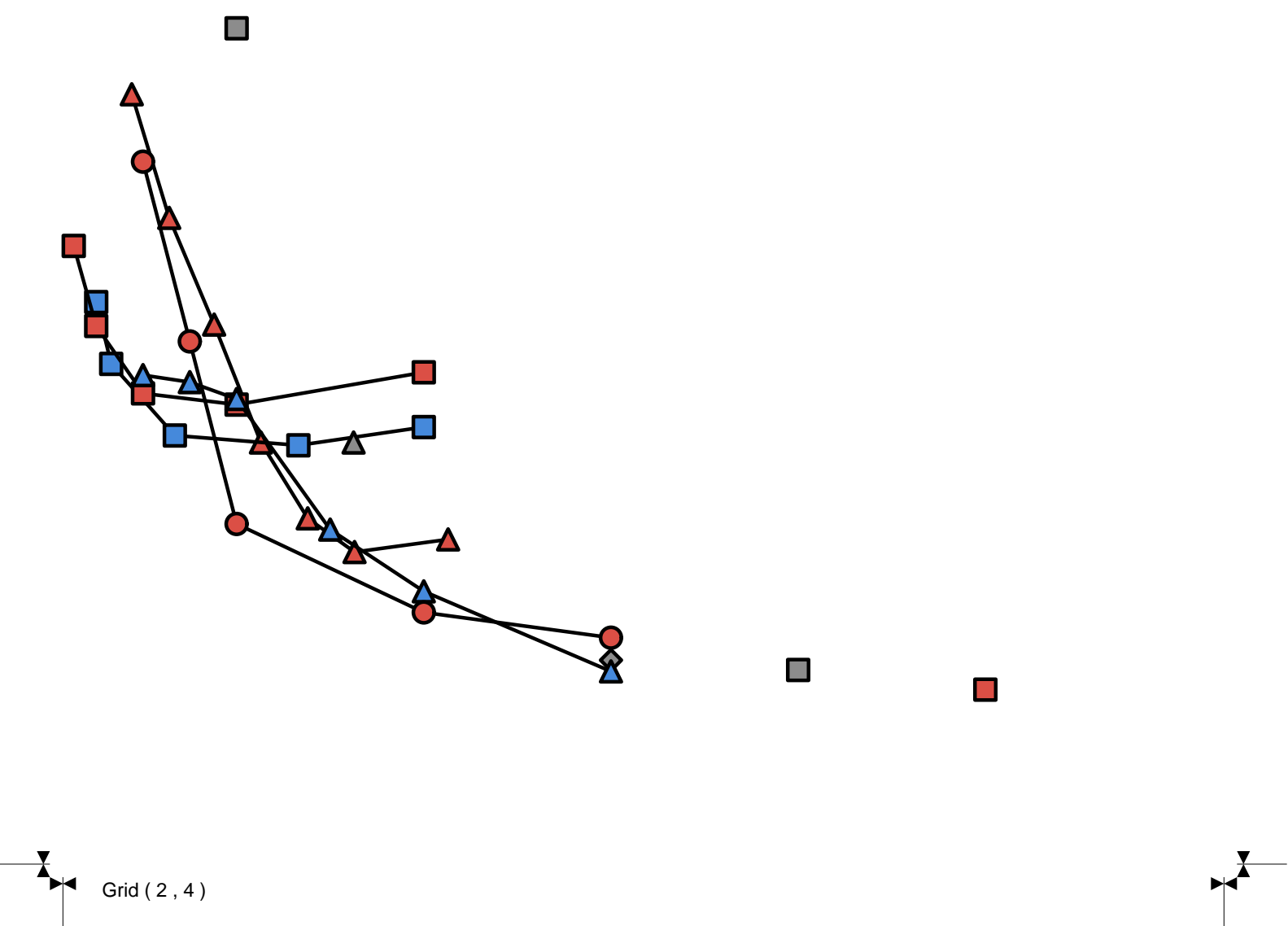
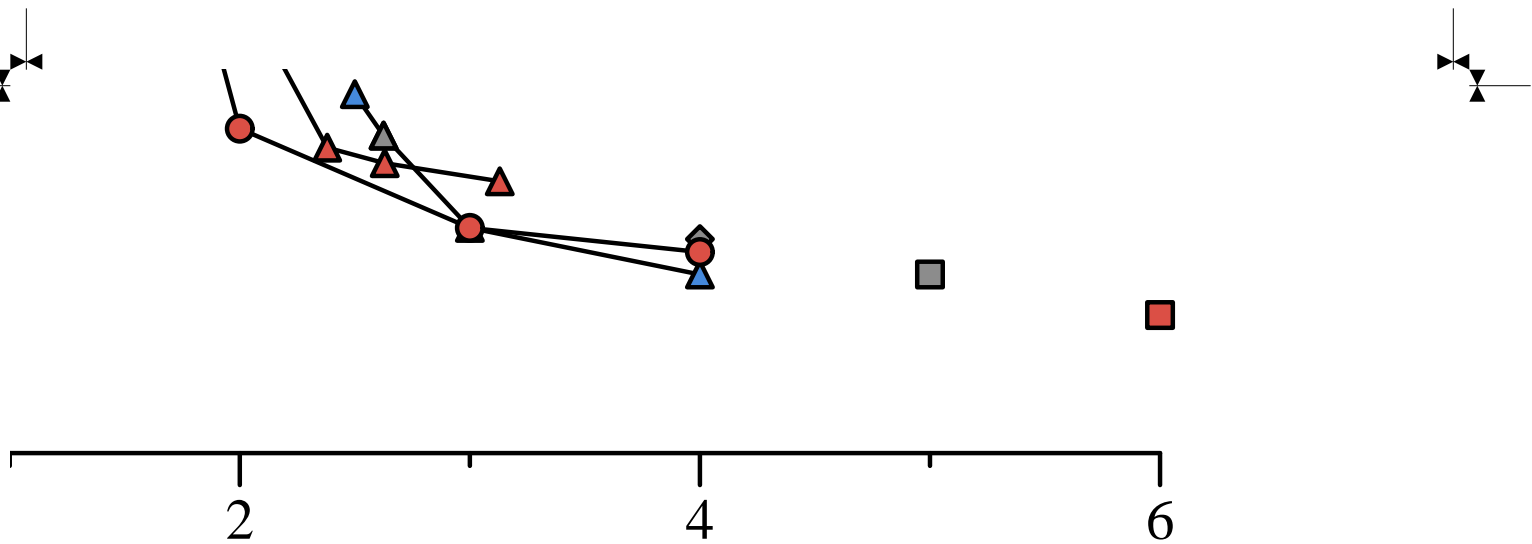
manifest as *pairs of parallel*
the BWT permutation. We

2. Replace the **other path**
by shortcut and follow it
on the second pass



the number of cache misses.
the algorithm.





The Burrows–Wheeler transform (BWT) is a reversible text transform defined as follows:

Input: text $T = \text{BANANA}\#$

1. Build a matrix with the text *rotations* as rows

2. Sort the rows

B	A	N	A	N	A	#
A	N	A	N	A	#	B
N	A	N	A	#	B	A
A	N	A	#	B	A	N
N	A	#	B	A	N	A
A	#	B	A	N	A	N
#	B	A	N	A	N	A

F

#	B	A
A	#	B
A	N	A
A	N	A
B	A	N
N	A	#
N	A	N

Output: BWT $L = \text{ANNB}\#\text{AA}$ (the last column of the sorted matrix)

The properties of the BWT make it a good choice for data compression. It is used as a preprocessing stage in many compression programs, including the widely used bzip2 (thus the b).

NEW ALGORITHMS FOR INDEXING

(BWT) is an in-
follows.

Define $RANK(j) = |\{i \mid i \leq j \text{ and } L[i] = L[j]\}|$
The BWT can be inverted

Input: BWT $L = \text{ANNB\#A}$

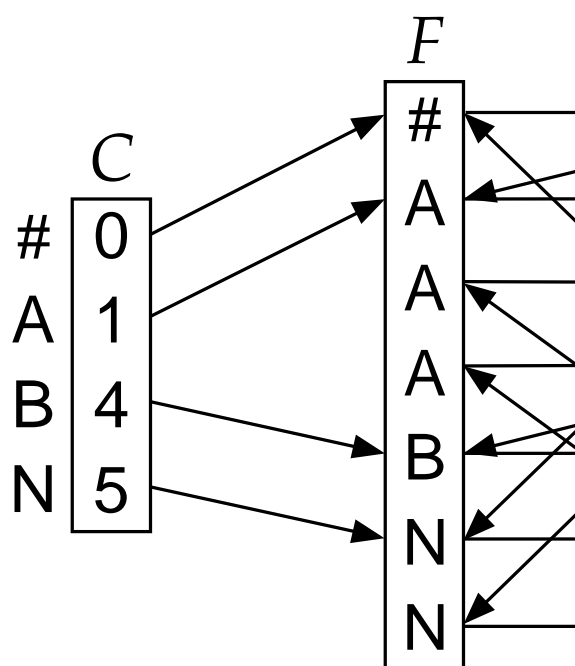
Sort the rows

				L
A	N	A	N	A
B	A	N	A	N
A	#	B	A	N
A	N	A	#	B
N	A	N	A	#
#	B	A	N	A
N	A	#	B	A

1st column)

is easier to com-
puted as the first
columns including
)

1. Compute C and $RANK$ array



2. Starting at $L[i] = \#$, follow

$$i \mapsto C[L[i]]$$

Output $L[i]$ at each step

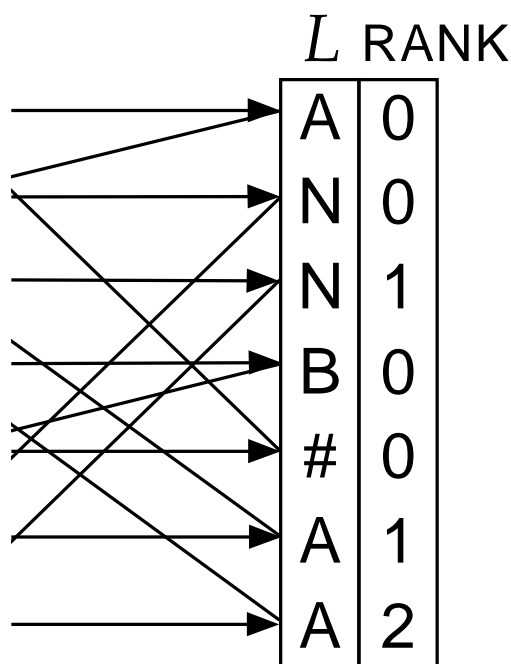
Output: reverse text $T^R =$

INVERSE BWT

$\{i < j \text{ and } L[i] = L[j]\} |$.
 ed as follows.

AA

arrays by scanning L



ow the permutation:

$] + \text{RANK}(i)$

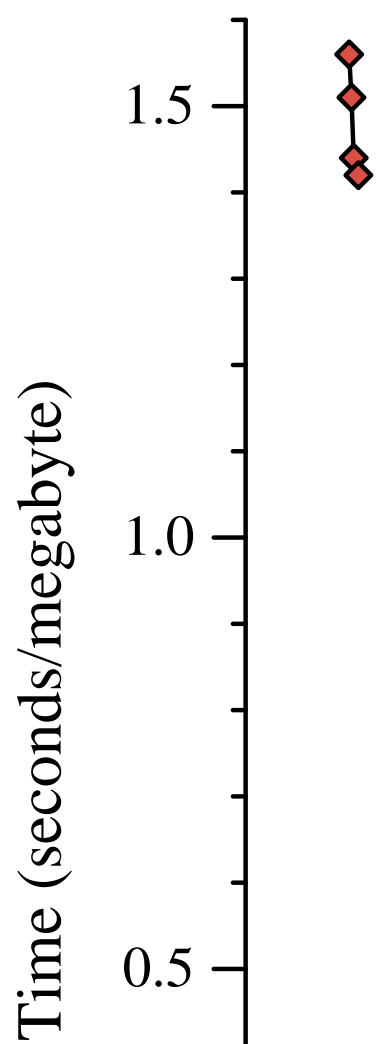
$= \#ANANAB$

The graphs be
 quirements of
 The algorithms

New algorithm
 repetitio

Improved imp
 algorithm

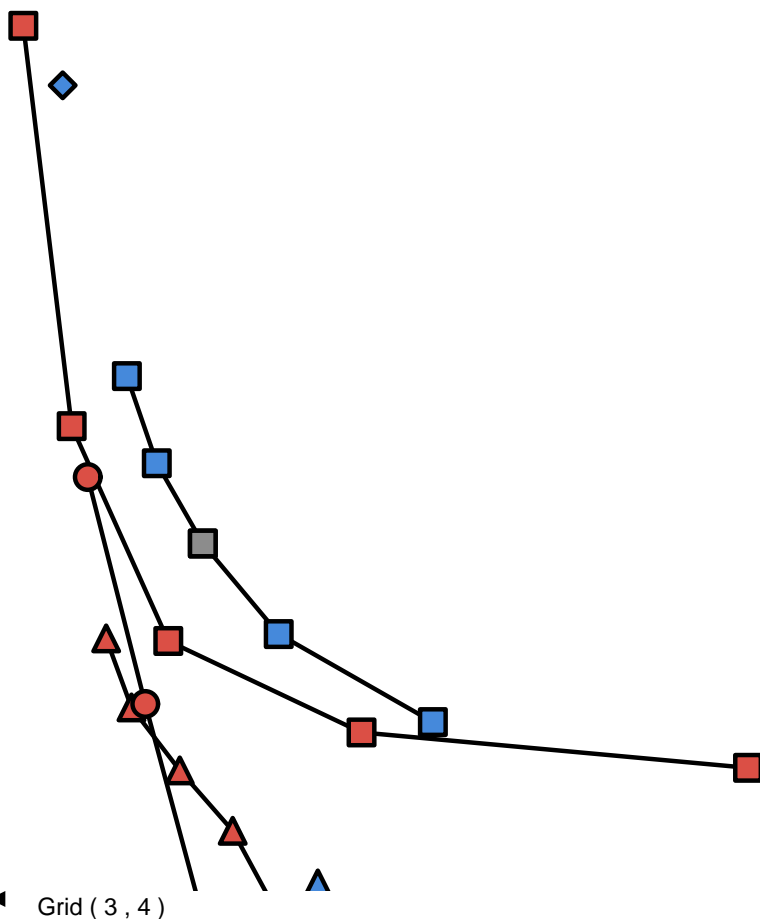
Prior algorithm



Below show the time and space requirements of several algorithms on two texts. The results are divided into three groups:

- Algorithms based on reference point ranks, position shortcuts and wavelet trees
- Implementations of wavelet trees and wavelet matrices from [1]
- Algorithms from [3, 1]

ENGLISH 100MB





IMPROVING INV THROUGHOUT

The Burrows-Wheeler transform
pression used for example in the
The *inverse* BWT is usually the b
with respect to both space and ti

BURROWS–WHEELER TRA

The Burrows–Wheeler transform (B

INVERSE BURROWS- THE SPACE-TIME S

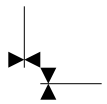
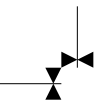
1 (BWT) is a powerful tool for data com-
2 a popular compression program bzip2.
3 bottleneck in the decompression phase
4 time.

TRANSFORM

INVERSE BWT

1 (BWT) is an in-
Grid (4, 2)

Define $DANK(i) = |S_i| i$



MATEMAATTIS-LUON
MATEMATISK-NATU

WHEELER TRANSFO SPECTRUM

Juha I

Our new algorithms improve the perfor
range from the fastest known algorithm
and cover the whole space-time tradeo

EXPERIMENT

i and $T[i] = T[i] \cup$
Grid (4, 3)

The graphs be

HELSINGIN YLIOPISTO
HELSINGFORS UNIVERSITET
UNIVERSITY OF HELSINKI
LUMINENTTIETEELLINEN TIEDEKUNTA
LUMIVETENSKAPLIGA FAKULTETEN
FACULTY OF SCIENCE

FORM

Kärkkäinen and Simon Puglisi

formance of inverse BWT. They
n to the most space-efficient one,
off spectrum in between.

INITIAL RESULTS

plots show the time and space re-