
Fast Lightweight Suffix Array Construction (and Checking)

Stefan Burkhardt Juha Kärkkäinen

MPI für Informatik

CPM 2003, Morelia, June 25–27

Yet Another Suffix Array Construction Algorithm

Fast

- ▶ $\mathcal{O}(n \log n)$ **time** for general alphabet
- ▶ fast in practice

Lightweight

- ▶ **sublinear** ($\mathcal{O}(n/\sqrt{\log n})$) **extra space** in addition to input (string) and output (suffix array)
- ▶ less than n bytes in practice

Tradeoff

- ▶ $\mathcal{O}(vn + n \log n)$ **time** in $\mathcal{O}(n/\sqrt{v})$ **extra space** for $v \in [2, n]$

Related work

time

alphabet

extra space

reference

Fast in the worst case

$\mathcal{O}(n)$

integer

$\geq 4n$ bytes

Larsson & Sadakane '99

$\mathcal{O}(n \log n)$

general

many others

Fast on average (string sorting + heuristics)

$\Omega(n^2)$ worst

general

little

Manzini & Ferragina '02

$\mathcal{O}(n \log n)$ avg.

others

Compressed suffix array

$\mathcal{O}(n \log n)$

constant

$6n$ bits

Lam & al. '02

$\mathcal{O}(n \log \log \sigma)$

integer

$\mathcal{O}(n \log \sigma)$ bits

Hon & al. FOCS '03

Algorithm Outline

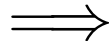
1. Choose a sample of suffixes
2. Sort the sample
3. Sort all suffixes with the help of the sample

suffixes			sample	
0	whowhowho?			
1	howhowho?			
2	owhowho?		2	owhowho?
3	whowho?			
4	howho?	⇒	4	howho?
5	owho?		5	owho?
6	who?			
7	ho?			
8	o?			
9	?		9	?

Step 2: Sort the sample ...

... and store the ranks

1	9	?
2	4	howho?
3	5	owho?
4	2	owhowho?



0	w h 4o w 2h 3o w h o 1?
1	h 4o w 2h 3o w h o 1?
2	4o w 2h 3o w h o 1?
3	w 2h 3o w h o 1?
4	2h 3o w h o 1?
5	3o w h o 1?
6	w h o 1?
7	h o 1?
8	o 1?
9	1?

Step 3: Sort all suffixes

Order two suffixes by the first pair of

- ▶ **mismatching** characters or

$\boxed{6}$	w h o 1?	$\boxed{1}$	h 4o w 2h 3o w h o 1?
$\boxed{3}$	w 2h 3o w h o 1?	$\boxed{0}$	w h 4o w 2h 3o w h o 1?

- ▶ aligned **sample suffixes** (**anchor pair**)

$\boxed{3}$	w 2h 3o w h o 1?	$\boxed{4}$	2h 3o w h o 1?
$\boxed{0}$	w h 4o w 2h 3o w h o 1?	$\boxed{5}$	3o w h o 1?

⇒ choose the sample to ensure **early anchor pairs**

How to choose the sample?

Regular sample (every k -th suffix)

- ▶ regularity \implies fast sorting
- ▶ **no anchor pairs** for most pairs of suffixes

w³h o⁷w h⁵o w²h o⁶w h⁴o w¹h o
7w h⁵o w²h o⁶w h⁴o w¹h o

Random sample (of size n/k)

- ▶ expected distance to an **anchor pair** is k^2
- ▶ **sorting** time $\Omega(n^2/k)$ in the worst case
- ▶ not good for any k

Difference cover sample

Choose a v -periodic sample $\{i \in [0, n) \mid i \bmod v \in D\}$
for a difference cover D modulo v

A difference cover D modulo v is a
subset of $[0, v)$ such that for all $i \in [0, v)$
there exist $j, k \in D$ with $i \equiv k - j \pmod{v}$

Example:

$D = \{2, 4, 5\}$ is a difference cover modulo 7

i	$k - j$
0	2 - 2
1	5 - 4
2	4 - 2
3	5 - 2
4	2 - 5
5	2 - 4
6	4 - 5

Difference cover sample

Choose a v -periodic sample $\{i \in [0, n) \mid i \bmod v \in D\}$
for a difference cover D modulo v

A difference cover D modulo v is a subset of $[0, v)$ such that for all $i \in [0, v)$ there exist $j, k \in D$ with $i \equiv k - j \pmod{v}$

i	$k - j$
0	2 - 2
1	5 - 4
2	4 - 2
3	5 - 2
4	2 - 5
5	2 - 4
6	4 - 5

Example:

$D = \{2, 4, 5\}$ is a difference cover modulo 7

D -sample:

| 0 1 2 3 4 5 6 | 7 8 9 10 11 12 13 | 14 15 16 17 18 19 20 | ...

Difference cover sample

Choose a v -periodic sample $\{i \in [0, n) \mid i \bmod v \in D\}$
for a difference cover D modulo v

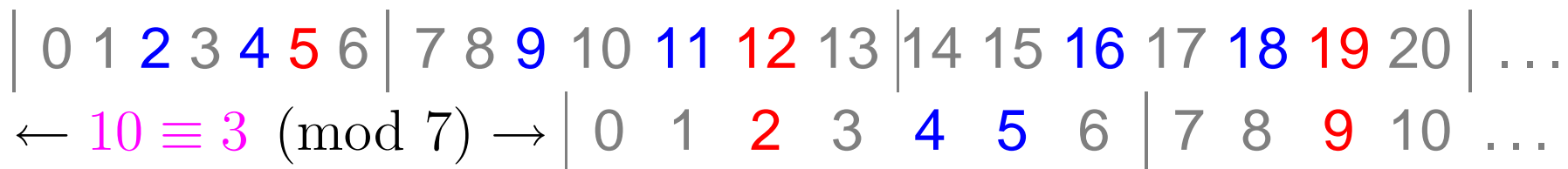
A difference cover D modulo v is a subset of $[0, v)$ such that for all $i \in [0, v)$ there exist $j, k \in D$ with $i \equiv k - j \pmod{v}$

i	$k - j$
0	2 - 2
1	5 - 4
2	4 - 2
3	5 - 2
4	2 - 5
5	2 - 4
6	4 - 5

Example:

$D = \{2, 4, 5\}$ is a difference cover modulo 7

D -sample:



Difference cover sample

Key properties

- ▶ ensures an **anchor pair within distance v**
- ▶ **closest anchor pair** can be found **in constant time** after $\mathcal{O}(v)$ **time and space preprocessing**
- ▶ **[Colbourn & Ling '00]**: difference cover of size $\mathcal{O}(\sqrt{v})$ is easy to compute for all $v \implies$ **sample size $\mathcal{O}(n/\sqrt{v})$**
- ▶ v -periodic \implies sample sorting in $\mathcal{O}(\sqrt{v}n + (n/\sqrt{v}) \log n)$ **time** and $\mathcal{O}(v + n/\sqrt{v})$ **space**

Difference cover sample algorithm

1. Compute and preprocess the difference cover
 $\mathcal{O}(v)$ time and space
2. Sort the sample
 $\mathcal{O}(\sqrt{v}n + (n/\sqrt{v}) \log n)$ time and $\mathcal{O}(v + n/\sqrt{v})$ space
3. Sort all suffixes in place
 - (a) Sort suffixes by first v characters
 $\mathcal{O}(vn + n \log n)$ time and $\mathcal{O}(\log n)$ extra space
 - (b) Finish sorting using anchor pairs
 $\mathcal{O}(n \log n)$ time and $\mathcal{O}(v + n/\sqrt{v})$ extra space

Total: $\mathcal{O}(vn + n \log n)$ time and $\mathcal{O}(v + n/\sqrt{v})$ extra space
If $v = \log n$: $\mathcal{O}(n \log n)$ time and $\mathcal{O}(n/\sqrt{\log n})$ extra space

Experimental results

LS = Larsson & Sadakane '99

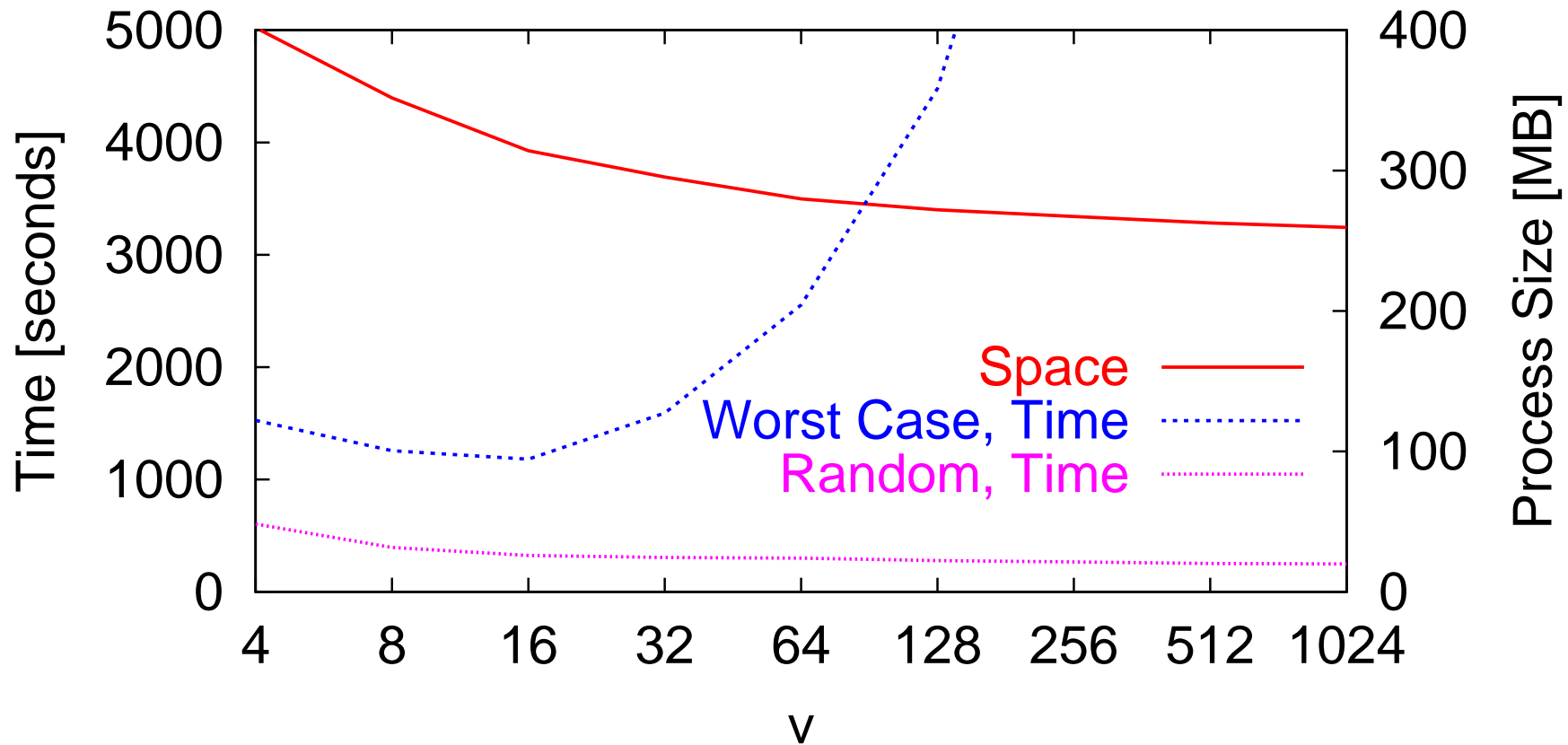
MF = Manzini & Ferragina '02

DC = difference cover algorithm with $v = 64$

text	time [sec]			space [MB]		
	LS	DC	MF	LS	DC	MF
bible (4.0M)	12	5.5	2.1	33	24	22
e.coli (4.6M)	14	5.6	2.1	38	27	25
gcc source (86M)	529	512	299	694	484	439
(A) ^{50M}	160	494	8.3	401	280	252
(1K random) ^{50K}	1222	2251	1M ^a	401	280	252
(500K random) ¹⁰⁰	1311	581	1M ^a	401	280	252

^aRough estimate by extrapolation

Experimental results



text length = 50M
alphabet size = 26

Worst Case = (1K random)^{50K}

Conclusion

Summary

- ▶ fast and lightweight suffix array construction in **theory and practice**
- ▶ based on **difference cover samples**

Extensions

- ▶ LCP computation
- ▶ sorting arbitrary subsets of suffixes

Future work

- ▶ still faster implementation