

Space efficient locating with the r-index

Gagie, Navarro, Prezza
SODA 2017 & JACM 2018

Locate on a BWT-index

We consider the problem of finding the position in the text of all pattern occurrences.

Recall we only have the first and last column (F and L)

$T = \text{swiss miss miss missing}$

F	L
5 miss miss missingswiss	s
10 miss miss missingswiss miss	s
15 missingswiss miss miss	s
22 gswiss miss miss miss	n
20 ingswiss miss miss miss	s
2 iss miss miss missings	w
7 iss miss missingswiss	m
12 iss missingswiss miss	m
17 issingswiss miss miss	m
6 miss miss missingswiss	
11 miss missingswiss miss	
16 missingswiss miss miss	
21 ngswiss miss miss miss	i
4 s miss miss missingswi	s
9 s miss missingswiss mi	s
14 s missingswiss miss mi	s
19 singswiss miss miss mi	s
3 ss miss miss missingsw	i
8 ss miss missingswiss m	i
13 ss missingswiss miss m	i
18 ssingswiss miss miss m	i
0 swiss miss miss missin	g
1 wiss miss miss missing	s

Simple solution: uniform sampling
store one out of t SA values

To find the position of an occ
use te LF map to move backward
untill we reach a stored value

the parametr t induces a trade-off:

extra space: $(n/t)\log n$ bits

locate time: $O(t)$ per occurence

F	L
miss miss missingswiss	s
miss miss missingswiss miss	s
missingswiss miss miss	s
gswiss miss miss miss	n
ingswiss miss miss miss	s
iss miss miss missings	w
iss miss missingswiss	m
12 iss missingswiss miss	m
17 issingswiss miss miss	m
o miss miss missingswiss	i
11 miss missingswiss miss	s
missingswiss miss miss	s
ngswiss miss miss miss	i
s miss miss missingswi	s
s miss missingswiss mi	s
s missingswiss miss mi	s
singswiss miss miss mi	s
ss miss miss missingsw	i
ss miss missingswiss m	i
ss missingswiss miss m	i
-- ssingswiss miss miss m	i
0 swiss miss miss missin	g
1 wiss miss miss missing	s

When the input is highly compressible (for example consists of many variants of the same sequence) it is more convenient to use an index of size $O(r)$ words where r is the number of runs in the BWT.

In this setting storing (n/t) SA entries space dominates the index size: using BWT properties we can save space by storing only $2r$ SA entries

The resulting index is called the r-index
[Gagie, Prezza, Navarro 2018]

r-Index: locate 1st occurrence

Toehold Lemma:

to locate the lexicographically first occurrence of a pattern we only need the SA entries for rows containing the first occ of a run in L

Proof:

by induction on the backward search steps

Note: we also need the select operation on column L

F	L
5	miss miss missingswiss
-	miss missingswiss miss
-	missingswiss miss miss
22	gswiss miss miss miss
20	ingswiss miss miss miss
2	iss miss miss missings
7	iss miss missingswiss
12	iss missingswiss miss
-	iss missingswiss miss
6	miss miss missingswiss
11	miss missingswiss miss
-	missingswiss miss miss
21	ngswiss miss miss miss
4	s miss miss missingswi
9	s miss missingswiss mi
-	s missingswiss miss mi
-	singswiss miss miss mi
5	ss miss miss missingsw
8	ss miss missingswiss m
-	ss missingswiss miss m
-	ssingswiss miss miss m
6	swiss miss miss missin
1	wiss miss miss missing

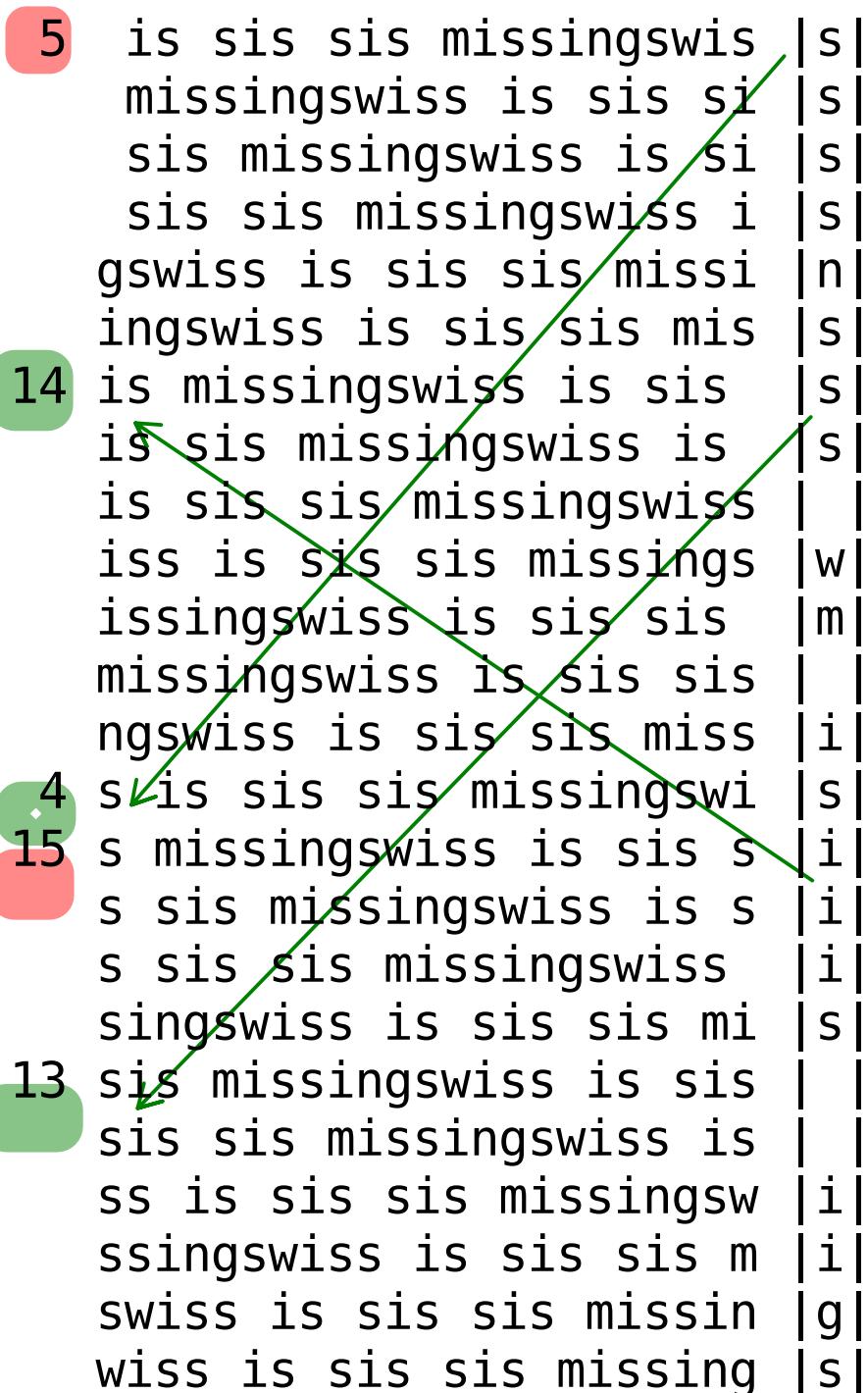
Example: searching "sis" in

T = swiss is sis sis missing

We only use the SA entries marked in red. The one in green are derived

Green arrows are applications of the LF map. Each LF application reduces the current position by 1

The first occurrence of "sis" is in text position 13



r-Index: locate next occurrence

 1 1 2
0 5 0 5 0
T = Swiss miss miss missing

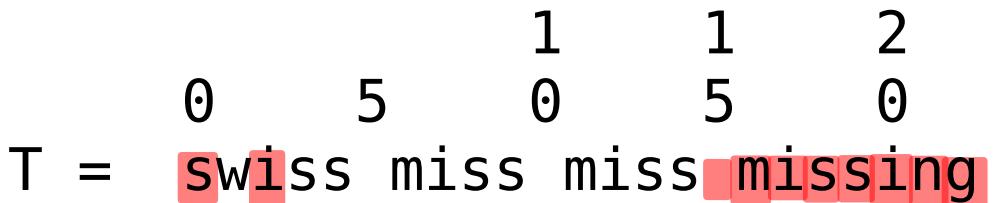
SA values at end/begin of runs
(15,22), (22,20), (20,2) (2,7)
(17,6) (16,21), (21,4) (19,3)
(18,0), (0,1)

Sorted pairs:

(0,1) (2,7), (15,22), (16,21)
(17,6) (18,0) (19,3) (20,2)
(21,4) (22,20)

F	L
5 miss miss missingswiss	S
10 miss missingswiss miss	S
15 missingswiss miss miss	S
22 gswiss miss miss miss	n
20 ingswiss miss miss miss	s
2 iss miss miss missingswiss	w
7 iss miss missingswiss	m
12 iss missingswiss miss	m
17 issingswiss miss miss	m
6 miss miss missingswiss	
11 miss missingswiss miss	
16 missingswiss miss miss	
21 ngswiss miss miss miss	i
4 s miss miss missingswiss	s
9 s miss missingswiss miss	s
14 s missingswiss miss miss	s
19 singswiss miss miss miss	s
3 ss miss miss missingswiss	i
8 ss miss missingswiss miss	i
13 ss missingswiss miss miss	i
18 ssingswiss miss miss miss	i
0 swiss miss miss missings	g
1 wiss miss miss missing	s

r-Index: locate next occurrence

T = 
0 5 0 5 0
1 1 2

Sorted pairs:

(0,1) (2,7), (15,22), (16,21)
(17,6) (18,0) (19,3) (20,2)
(21,4) (22,20)

Lemma: $L(p)=L(p+1)$ $q=LF(p)$
 $q+1=L(p+1)$

rows ending with the same symbol stay together!

F L
miss miss missingswiss
miss missingswiss miss
missingswiss miss miss
gswiss miss miss miss
ingswiss miss miss miss
iss miss miss miss
iss miss miss missings
iss miss missingswiss
iss missingswiss miss
iss missingswiss miss
miss miss missingswiss
miss missingswiss miss
missingswiss miss miss
ngswiss miss miss miss
s miss miss missingswi
s miss missingswiss mi
s missingswiss miss mi
singswiss miss miss mi
ss miss miss missingsw
ss miss missingswiss m
ss missingswiss miss mi
ssingswiss miss miss mi
swiss miss miss missin
wiss miss miss missing

r-Index: locate next occurrence

		1	1	2
0	5	0	5	0
T =	swiss miss miss	missing		

Sorted pairs:

- (0,1) (2,7), (15,22), (16,21)
- (17,6) (18,0) (19,3) (20,2)
- (21,4) (22,20)

Given the text position of a row using a predecessor query of the sorted pairs we can retrieve the text position of the next row

F	L
miss	s
miss	s
missingswiss	s
missingswiss	s
miss	s
miss	s
miss	n
miss	s
miss	s
miss	s
miss	w
miss	s
miss	m
miss	w
miss	s
miss	m
miss	m
miss	i
miss	s
miss	i
miss	g
miss	s

r-Index: locate next occurrence

		1	1	2
	0	5	0	5
T =	swiss	miss	miss	missing

Sorted pairs:

$(0, 1)$ $(2, 7)$, $(15, 22)$, $(16, 21)$
 $(17, 6)$ $(18, 0)$ $(19, 3)$ $(20, 2)$
 $(21, 4)$ $(22, 20)$

Examples :

16->21 (from the 4th pair)

3->2(LF) ->7(2nd pair)->8

$$10 \rightarrow 2(\text{pred}) \rightarrow 7 + (10 - 2) \rightarrow 15$$

	F	L
10	miss miss missingswiss miss miss missingswiss mis missingswiss miss mis gswiss miss miss missi ingswiss miss miss mis iss miss miss missings iss miss missingswiss iss missingswiss miss issingswiss miss miss miss miss missingswiss miss missingswiss miss missingswiss miss miss ngswiss miss miss miss s miss miss missingswi s miss missingswiss mi s missingswiss miss mi singswiss miss miss mi ss miss miss missingsw ss miss missingswiss m ss missingswiss miss m ssingswiss miss miss m swiss miss miss missin wiss miss miss missing	s s s n s w m m m i s s s s s i s s s s i i i i i g s
?		
2		
7		
16		
?		
3		
?		

The space/time bounds for the r-index are:

$O(r)$ words $O((|p| + occ) \log \log n)$ time

$O(r \log \log n)$ words $O(|p| + occ)$ time (optimal)

Timeline of BWT based indexing

1994 BWT (Burrows, Wheeler)

1997 bzip2 (Seward)

2000 Backward search (Ferragina, GM)

2003 Wavelet Trees (Gupta, Grossi, Vitter)

2017 r-index (Gagie, Navarro, Prezza)

2017 Wheeler-Graphs (Gagie, GM, Siren)

It took 17 years
to devise a space
efficient locate

Next topic! →