

EPR-dictionaries:

**A practical and fast data structure
for constant time searches
in unidirectional and bidirectional FM-indices**

Christopher Pockrandt, Marcel Ehrhardt, Knut Reinert

Free University of Berlin, Max-Planck-Institute for Molecular Genetics

FM-Index (by Ferragina and Manzini)

T = mississippi\$

i	S	F		L
1	12	\$	mississippi	i
2	11	i	\$mississip	p
3	8	i	ppi\$missis	s
4	5	i	ssippi\$mis	s
5	2	i	ssissippi\$	m
6	1	m	issippi	\$
7	10	p	i\$mississi	p
8	9	p	pi\$mississ	i
9	7	s	ippi\$missi	s
10	4	s	issippi\$mi	s
11	6	s	sippi\$miss	i
12	3	s	sissippi\$m	i

← sort

S_i = suffix array
F = First column
L = Last column = BWT

1	mississippi\$
2	issippi\$m
3	ssissippi\$mi
4	sissippi\$mis
5	issippi\$miss
6	ssippi\$missi
7	sippi\$missis
8	ippi\$mississ
9	ppi\$mississi
10	pi\$mississip
11	i\$mississipp
12	\$mississippi

FM-Index (by Ferragina and Manzini)



T = mississ**ip**pi\$

S = suffix array

F = First column

L = Last column = BWT

i	S	F		L
1	12	\$	mississippi	i
2	11	i	\$mississip	p
3	8	i	ppi\$missis	s
4	5	i	ssippi\$mis	s
5	2	i	ssissippi\$	m
6	1	m	issippi	\$
7	10	p	i\$mississi	p
8	9	p	pi\$mississ	i
9	7	s	ippi\$missi	s
10	4	s	issippi\$mi	s
11	6	s	sippi\$miss	i
12	3	s	sissippi\$m	i

1st p in F corresponds to 1st p in L
 2nd i in F corresponds to 2nd i in L
 ...

We also need a Count table C

c	Count (c)
\$	1
i	5
m	6
p	8
s	12

Backward search

i	S	F	L
1	12	\$ mississipp	i
2	11	i \$mississip	p
3	8	i ppi\$missis	s
4	5	i ssippi\$mis	s
5	2	i ssissippi\$	m
6	1	m ississippi	\$
7	10	p i\$mississi	p
8	9	p pi\$mississ	i
9	7	s ippi\$missi	s
10	4	s issippi\$mi	s
11	6	s sippi\$miss	i
12	3	s sissippi\$m	i

T = mississippi\$

Pattern P = "is"

[1, 12] corresponds to ϵ

[9, 12] corresponds to "s"

Backward search

i	S	F	L
1	12	\$	mississippi
2	11	i	\$mississip
3	8	i	ppi\$missis
4	5	i	ssippi\$mis
5	2	i	ssissippi\$
6	1	m	ississippi
7	10	p	i\$mississi
8	9	p	pi\$mississ
9	7	s	ippi\$missi
10	4	s	issippi\$mi
11	6	s	sippi\$miss
12	3	s	sissippi\$m

T = mississippi\$

Pattern P = "is"

[1, 12] corresponds to ϵ

[9, 12] corresponds to "s"

[4, 5] corresponds to "is"



$\text{rank}_i(8) = 2$

Bit vectors

BWT	i	p	s	s	m	\$	p	i	s	s	i	i
\$	0	0	0	0	0	1	0	0	0	0	0	0
i	1	0	0	0	0	0	0	1	0	0	1	1
m	0	0	0	0	1	0	0	0	0	0	0	0
p	0	1	0	0	0	0	1	0	0	0	0	0
s	0	0	1	1	0	0	0	0	1	1	0	0

Constant time rank support for each bit vector (Jacobson)

	block 1				block 2				block 3			
i	1	0	0	0	0	0	0	1	0	0	1	1

block #	rank
1	1
2	2
3	4

BWT	i	p	s	s	m	\$	p	i	s	s	i	i
\$	0	0	0	0	0	1	0	0	0	0	0	0
i	1	0	0	0	0	0	0	1	0	0	1	1
m	0	0	0	0	1	0	0	0	0	0	0	0
p	0	1	0	0	0	0	1	0	0	0	0	0
s	0	0	1	1	0	0	0	0	1	1	0	0

Constant time rank support for each bit vector (Jacobson)

	block 1				block 2				block 3			
i	1	0	0	0	0	0	0	1	0	0	1	1



block #	rank
1	1
2	2
3	4

$$\text{rank}_i(11) =$$

$$2 + 1$$



table lookup + in-block query

BWT	i	p	s	s	m	\$	p	i	s	s	i	i
\$	0	0	0	0	0	1	0	0	0	0	0	0
i	1	0	0	0	0	0	0	1	0	0	1	1
m	0	0	0	0	1	0	0	0	0	0	0	0
p	0	1	0	0	0	0	1	0	0	0	0	0
s	0	0	1	1	0	0	0	0	1	1	0	0

Constant time rank support for each bit vector (Jacobson)

	block 1				block 2				block 3			
i	1	0	0	0	0	0	0	1	0	0	1	1

block #	rank
1	1
2	2
3	4

Only $o(n)$ additional space per char

In total $n \cdot \sigma + o(n \cdot \sigma)$ space

Bidirectional search (Lam et al)

T = mississippi\$ P = is T^{rev} = ippississim\$

i	S	F	L
1	12	\$	mississipp i
2	11	i	\$mississip p
3	8	i	ppi\$missis s
4	5	i	ssippi\$mis s
5	2	i	ssissippi\$ m
6	1	m	ississippi \$
7	10	p	i\$mississi p
8	9	p	pi\$mississ i
9	7	s	ippi\$missi s
10	4	s	issippi\$mi s
11	6	s	sippi\$miss i
12	3	s	sissippi\$m i

i	S'	F'	L'
1	12	\$	ippississi m
2	10	i	m\$ippissis s
3	1	i	ppississim \$
4	7	i	ssim\$ippis s
5	4	i	ssissim\$ip p
6	11	m	\$ippississ i
7	3	p	ississim\$i p
8	2	p	pississim\$ i
9	9	s	im\$ippissi s
10	6	s	issim\$ippi s
11	8	s	sim\$ippiss i
12	5	s	sissim\$ipp i

Bidirectional search (Lam et al)

T = mississippi\$

P = is

T^{rev} = ippississim\$

i	S	F	L
1	12	\$	mississippi
2	11	i	\$mississip
3	8	i	ppi\$missis
4	5	i	ssippi\$mis
5	2	i	ssissippi\$
6	1	m	ississippi
7	10	p	i\$mississi
8	9	p	pi\$mississ
9	7	s	ippi\$missi
10	4	s	issippi\$mi
11	6	s	sippi\$miss
12	3	s	sissippi\$m

i	S'	F'	L'
1	12	\$	ippississi
2	10	i	m\$ippissis
3	1	i	ppississim
4	7	i	ssim\$ippis
5	4	i	ssissim\$ip
6	11	m	\$ippississ
7	3	p	ississim\$i
8	2	p	pississim\$
9	9	s	im\$ippissi
10	6	s	issim\$ippi
11	8	s	sim\$ippiss
12	5	s	sissim\$ipp

Bidirectional search (Lam et al)

T = mississippi\$

P = is

T^{rev} = ippississim\$

i	S	F		L
1	12	\$	mississipp	i
2	11	i	\$mississip	p
3	8	i	ppi\$missis	s
4	5	i	ssippi\$mis	s
5	2	i	ssissippi\$	m
6	1	m	ississippi	\$
7	10	p	i\$mississi	p
8	9	p	pi\$mississ	i
9	7	s	ippi\$missi	s
10	4	s	issippi\$mi	s
11	6	s	sippi\$miss	i
12	3	s	sissippi\$m	i

i	S'	F'		L'
1	12	\$	ippississi	m
2	10	i	m\$ippissis	s
3	1	i	ppississim	\$
4	7	i	ssim\$ippis	s
5	4	i	ssissim\$ip	p
6	11	m	\$ippississ	i
7	3	p	ississim\$i	p
8	2	p	pississim\$	i
9	9	s	im\$ippissi	s
10	6	s	issim\$ippi	s
11	8	s	sim\$ippiss	i
12	5	s	sissim\$ipp	i

Bidirectional search (Lam et al)

T = mississippi\$

P = is

T^{rev} = ippississim\$

i	S	F	L
1	12	\$	mississippi
2	11	i	\$ mississippi
3	8	i	p issippi
4	5	i	ssissippi\$
5	2	i	ssissippi\$
6	1	m	issippi\$
7	10	p	i\$mississippi
8	9	p	pi\$mississippi
9	7	s	ippi\$mississippi
10	4	s	issippi\$mississippi
11	6	s	sissippi\$mississippi
12	3	s	sissippi\$mississippi

i	S'	F'	L'
1	12	\$	ippississim
2	10	i	m\$ippississim
3	1	i	ppississim\$
4	7	i	ssim\$ippissim
5	4	i	ssissim\$ipp
6	11	m	\$ippississim
7	3	p	ississim\$ipp
8	2	p	pississim\$ipp
9	9	s	i m\$ippissim
10	6	s	i ssim\$ippissim
11	8	s	sim\$ippissim
12	5	s	sissim\$ippissim

Bidirectional search (Lam et al)

T = mississippi\$

P = is

T^{rev} = ippississim\$

i	S	F		L
1	12	\$	mississipp	i
2	11	i	\$mississip	p
3	8	i	ppi\$missis	s
4	5	i	s sippi\$mis	s
5	2	i	s sissippi\$	m
6	1	m	ississippi	\$
7	10	p	i\$mississi	p
8	9	p	pi\$mississ	i
9	7	s	ippi\$missi	s
10	4	s	issippi\$mi	s
11	6	s	sippi\$miss	i
12	3	s	sissippi\$m	i

i	S'	F'		L'
1	12	\$	ippississi	m
2	10	i	m\$ippissis	s
3	1	i	ppississim	\$
4	7	i	ssim\$ippis	s
5	4	i	ssissim\$ip	p
6	11	m	\$ippississ	i
7	3	p	ississim\$i	p
8	2	p	pississim\$	i
9	9	s	i m\$ippissi	s
10	6	s	i ssim\$ippi	s
11	8	s	sim\$ippiss	i
12	5	s	sissim\$ipp	i

Bidirectional search (Lam et al)

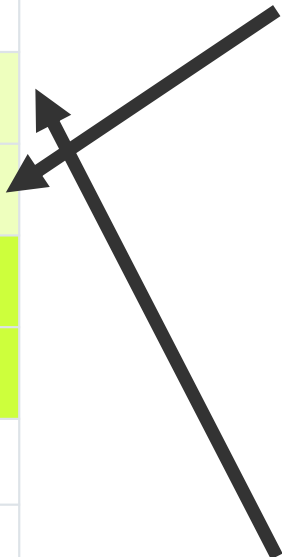
T = mississippi\$

P = is

T^{rev} = ippississim\$

i	S	F	L
1	12	\$	mississipp i
2	11	i	\$mississip p
3	8	i	ppi\$missis s
4	5	i	s sippi\$mis s
5	2	i	s sissippi\$ m
6	1	m	ississippi \$
7	10	p	i\$mississi p

i	S'	F'	L'
1	12	\$	i ppississi m
2	10	i	m\$ippissis s
3	1	i	ppississim \$
4	7	i	ssim\$ippis s
5	4	i	ssissim\$ip p
6	11	m	\$ippississ i
7	3	p	i ssissim\$ i p



**A forward search takes
O(Σ) backward searches / rank queries!**

8			
9			
10			
11	6	s	sippi\$miss i
12	3	s	sissippi\$m i

11	8	s	sim\$ippiss i
12	5	s	sissim\$ipp i

Bidirectional search (Lam et al)

T = mississippi\$

P = is

T^{rev} = ippississim\$

i	S	F	L
1	12	\$	mississippi
2	11	i	\$mississip
3	8	i	ppi\$missis
4	5	i	s issippi\$mis
5	2	i	s ississippi\$
6	1	m	ississippi
7	10	p	i\$mississi
8	9	p	pi\$mississ
9	7	s	ippi\$missi
10	4	s	issippi\$mi
11	6	s	sippi\$miss
12	3	s	sissippi\$m

i	S'	F'	L'
1	12	\$	ippississi
2	10	i	m\$ippissis
3	1	i	ppississim
4	7	i	ssim\$ippis
5	4	i	ssissim\$ip
6	11	m	\$ippississ
7	3	p	ississim\$i
8	2	p	pississim\$
9	9	s	i m\$ippissi
10	6	s	i ssim\$ippi
11	8	s	sim\$ippiss
12	5	s	sissim\$ipp

Observation: A forward search can be done with $O(1)$ backward-searches using less space (compared to Lam et al., 2009).

→ **Use prefix sums !**

BWT	m	s	\$	s	p	i	p	i	s	s	i	i
\$	0	0	1	0	0	0	0	0	0	0	0	0
i	0	0	1	0	0	1	0	1	0	0	1	1
m	1	0	1	0	0	1	0	1	0	0	1	1
p	1	0	1	0	1	1	1	1	0	0	1	1
s	1	1	1	1	1	1	1	1	1	1	1	1

Disadvantage : $n \cdot (\sigma - 1) + o(n \cdot \sigma)$ space

Use bit representation of BWT

	block 1				block 2				block 3			
BWT	A	C	A	T	G	T	A	C	A	G	C	A
	00	01	00	11	10	11	00	01	00	10	01	00



→ Needs only $n \cdot \log \sigma + o(n \cdot \log \sigma \cdot \sigma)$

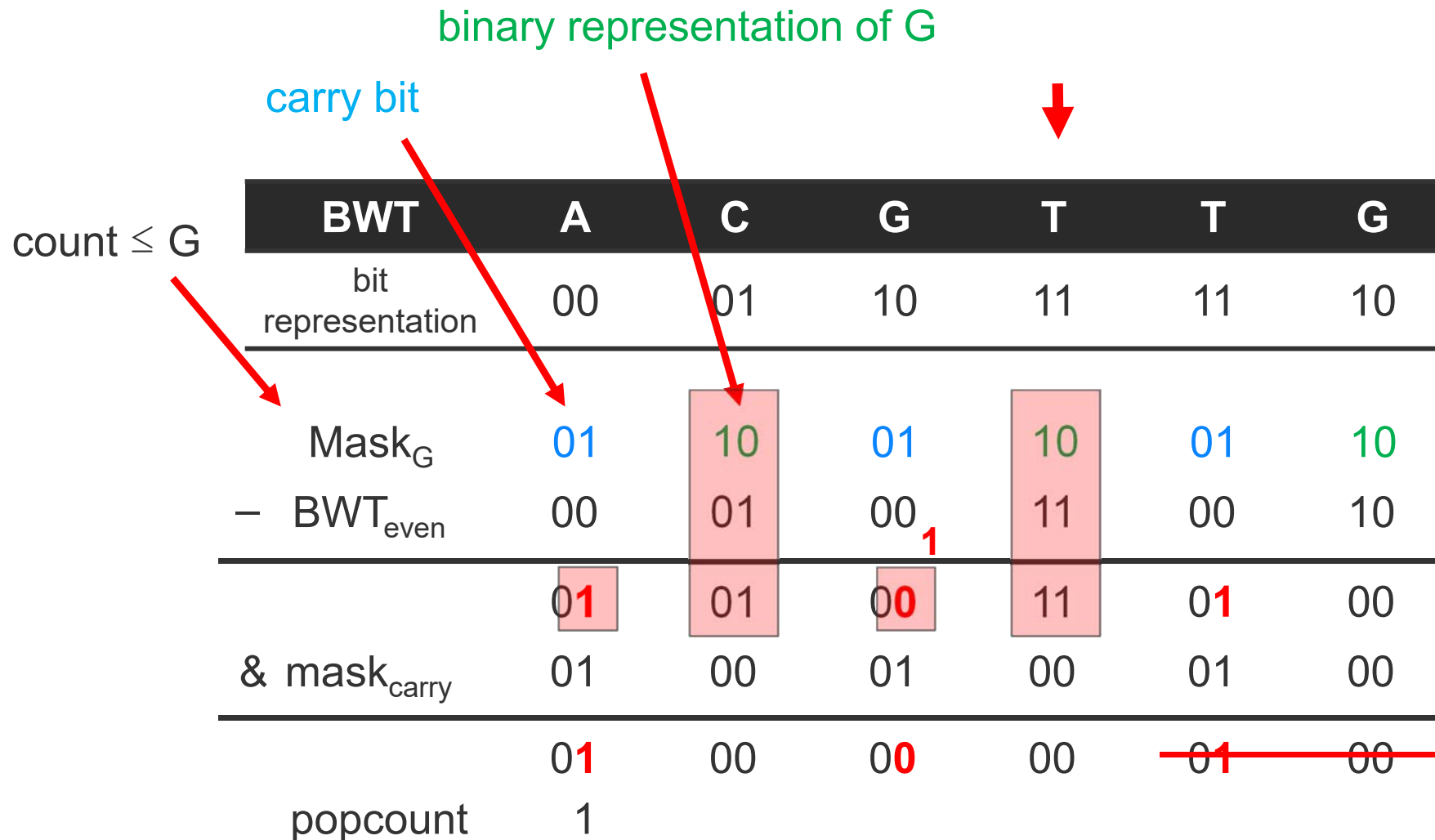
count $\leq G$

id	rank $\leq A$	rank $\leq C$	rank $\leq G$	rank $\leq T$
1	2	3	3	4
2	3	5	6	8
3	5	8	9	12

How to count in-block query in constant time?

Rank queries (in-block query)

Example: Count **even** occurrences of characters $\leq G$



Enhanced Prefixsum Rank dictionaries

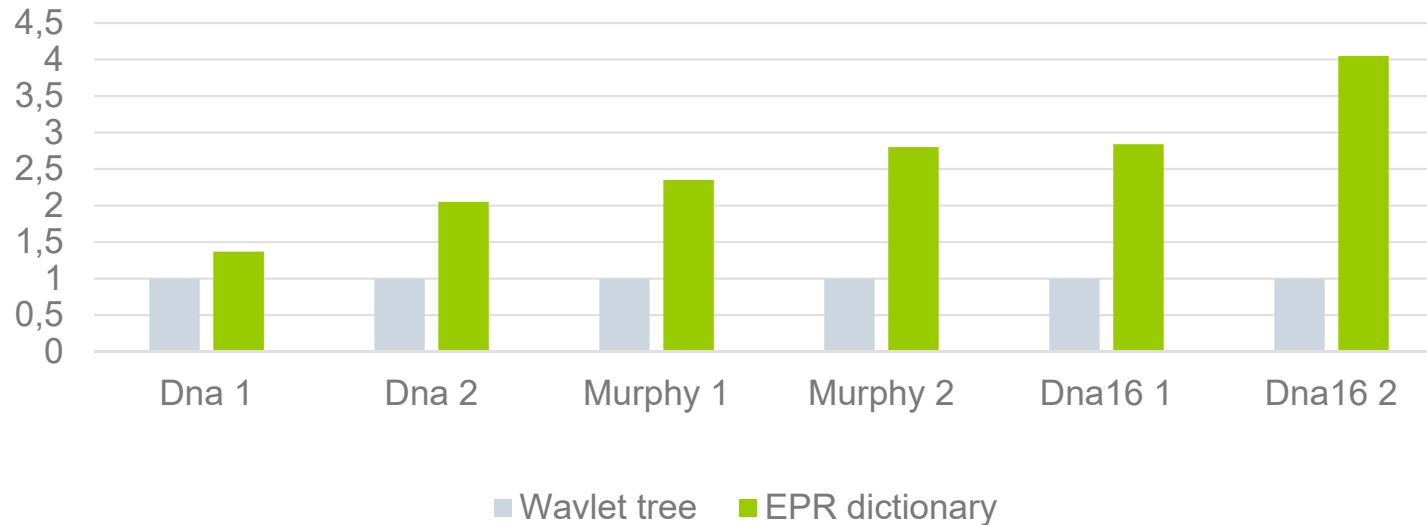
- backward and forward searches: $O(1)$
- $O(n \cdot \log \sigma) + o(n \cdot \sigma \cdot \log \sigma)$
space consumption

Wavelet trees

- backward and forward searches: $O(\log \sigma)$
- $O(n \cdot \log \sigma)$ space consumption

Runtime comparisons

Speedup for 1 and 2 directional searches



Index	Dna		Murphy10		Dna16	
	time	factor	time	factor	time	factor
WaveletTree	20.73s	1.00	52.35s	1.00	66.45s	1.00
EPR	15.09s	1.37	22.27s	2.35	23.39s	2.84
2WaveletTree	41.21s	1.00	66.59s	1.00	98.74s	1.00
2EPR	20.09s	2.05	23.80s	2.80	24.39s	4.05

1 million queries of length 200, sampled from a text of length 10^8

Space consumption

Index	Dna	Murphy10	Dna16
WaveletTree	30	51	60
EPR	42	156	227

rank data structure in Megabytes

We presented a practical rank dictionary that

- has the **same big O space consumption as wavelet trees**
- allows for **constant time** backward and forward searches
- has a speed-up factor between 2 and 4
- ~17% more space for DNA (CSA with sampling rate of 10)
- Available in the SeqAn library (www.seqan.de)
- Funding from the MPI for molecular genetics

Thank you for your attention



Runtime comparisons

Index	Dna		Murphy10		Dna16	
	time	factor	time	factor	time	factor
WaveletTree	20.73s	1.00	52.35s	1.00	66.45s	1.00
EPR	15.09s	1.37	22.27s	2.35	23.39s	2.84
2WaveletTree	41.21s	1.00	66.59s	1.00	98.74s	1.00
2SDSL	43.54s	0.95	74.69s	0.89	89.07s	1.11
2Schnattinger	59.59s	0.69	91.30s	0.73	107.04s	0.92
2EPR	20.09s	2.05	23.80s	2.80	24.39s	4.05

1 million queries of length 200, sampled from a text of length 10^8

Space consumption

Index	Dna	Murphy10	Dna16	Protein
EPR	42	156	227	478
WaveletTree	30	51	60	72
2EPR	84	311	454	955
2WaveletTree	60	102	120	144
2SDSL	68	105	122	145
Schnattinger	75	108	123	146

rank data structure in Megabytes