

# Succincter text indexing with wildcards

Chris Thachuk  
University of British Columbia

CPM 2011  
June 27, 2011

**T:** \*aa\*aca\*a\*aa\*cacc\*ac

**T:** \*aa\*aca\*a\*aa\*cacc\*ac

↑  
wildcard

T: \*aa\*aca\*a\*aa\*cacc\*ac



text segment

**T:** \*aa\*aca\*a\*aa\*cacc\*ac

**P:** aca

T: \*aa\*aca\*a\*aa\*cacc\*ac

P: aca

**T:** \*aa\***aca**\*a\*aa\*cacc\*ac

**P:** aca

T: \*aa\*aca\*a\*aa\*cacc\*ac

P: aca



T: \*aa\*aca\*a\*aa\*cacc\*ac

P: aca

T: \*aa\*aca\*a\*aa\*cacc\*ac

4 matches

P: aca

## Problem

Create a **succinct** index for:

- a text **T** of length **n**
- over an alphabet of size  $\sigma$
- containing **d** **wildcard** positions

to support efficient pattern matching queries.

## Problem

Create a **succinct** index for:

- a text  $\mathbf{T}$  of length  $\mathbf{n}$
- over an alphabet of size  $\sigma$
- containing  $\mathbf{d}$  **wildcard** positions

to support efficient pattern matching queries.

## Assumption (in presentation)

No adjacent wildcards in  $\mathbf{T}$ .

## SNPs in the human genome

**T:** ...atgacagaggggtcaac...

                  ↑                  ↑                  ↑

                  a/t                  a/t                  g/c

$|\mathbf{T}| \approx 3$  billion bases  
 $\sigma = 4$   
 $\mathbf{d} \approx 3$  million positions

↪ also protein collections ( $\sigma = 20$ )

# Our results in context

## Non succinct indexes

Cole, Gottlieb, Lewenstein (2004)	$O(n \log^k n)$ words
Lam <i>et al.</i> , (2007)	$O(n)$ words

# Our results in context

## Non succinct indexes

Cole, Gottlieb, Lewenstein (2004)  $O(n \log^k n)$  words  
Lam *et al.*, (2007)  $O(n)$  words

## Succinct indexes

	Index space (bits)	Query working space (bits)
Tam <i>et al.</i> , 2009	$(3 + o(1))n \log \sigma$	$O(n \log d)$

# Our results in context

## Non succinct indexes

Cole, Gottlieb, Lewenstein (2004)  $O(n \log^k n)$  words  
Lam *et al.*, (2007)  $O(n)$  words

## Succinct indexes

	Index space (bits)	Query working space (bits)
Tam <i>et al.</i> , 2009	$(3 + o(1))n \log \sigma$	$O(n \log d)$
	(reduced working space by increasing query time)	
Tam <i>et al.</i> , 2009	$(3 + o(1))n \log \sigma$	$O(n \log \sigma)$



# Our results in context

## Non succinct indexes

Cole, Gottlieb, Lewenstein (2004)  $O(n \log^k n)$  words  
Lam *et al.*, (2007)  $O(n)$  words

## Succinct indexes

	Index space (bits)	Query working space (bits)
Tam <i>et al.</i> , 2009	$(3 + o(1))n \log \sigma$	$O(n \log d)$
	(reduced working space by increasing query time)	
Tam <i>et al.</i> , 2009	$(3 + o(1))n \log \sigma$	$O(n \log \sigma)$
<b>This work</b>	$(2 + o(1))n \log \sigma + O(n)$	$O(m \log n + md)$

(Ignoring smaller order terms)

### 3 Types of matches

Type 1: pattern *contained within* a text segment

T: \*aa\*aca\*a\*aa\*cacc\*ac

P: cac

[Lam *et al.*, 2007, Tam *et al.*, 2009]

### 3 Types of matches

Type 1: pattern *contained within* a text segment

T: \*aa\*aca\*a\*aa\*cacc\*ac

P: cac

Type 2: pattern *spanning one* wildcard

T: \*aa\*aca\*a\*aa\*cacc\*ac

P: caca

[Lam *et al.*, 2007, Tam *et al.*, 2009]

### 3 Types of matches

Type 1: pattern *contained within* a text segment

T: \*aa\*aca\*a\*aa\*cacc\*ac

P: cac

Type 2: pattern *spanning one* wildcard

T: \*aa\*aca\*a\*aa\*cacc\*ac

P: caca

Type 3: pattern *contains at least one* text segment

T: \*aa\*aca\*a\*aa\*cacc\*ac

P: acacaca

[Lam *et al.*, 2007, Tam *et al.*, 2009]

mississippi

Naïve BWT algorithm

For a string  $S$ :

# Review: SAs and the Burrows Wheeler transform

mississippi\$

## Naïve BWT algorithm

For a string  $S$ :

- 1 append \$ to  $S$  ( $\$ < c, \forall c \in \Sigma$ )

# Review: SAs and the Burrows Wheeler transform

```
mississippi$  
ississippi$m  
ssissippi$mi  
sissippi$mis  
issippi$miss  
ssippi$missi  
sippi$missis  
ippi$mississ  
ppi$mississi  
pi$mississip  
i$mississipp  
$mississippi
```

## Naïve BWT algorithm

For a string  $S$ :

- 1 append  $\$$  to  $S$  ( $\$ < c, \forall c \in \Sigma$ )
- 2 list all cyclic rotations of  $S$

# Review: SAs and the Burrows Wheeler transform

```
$mississippi  
i$mississipp  
ippi$mississ  
issippi$miss  
ississippi$m  
mississippi$  
pi$mississip  
ppi$mississi  
ssippi$missis  
sissippi$mis  
ssippi$missi  
ssissippi$mi
```

## Naïve BWT algorithm

For a string  $S$ :

- 1 append  $\$$  to  $S$  ( $\$ < c, \forall c \in \Sigma$ )
- 2 list all cyclic rotations of  $S$
- 3 sort rows of matrix (lexicographically)



# Review: SAs and the Burrows Wheeler transform

\$	i
i\$	p
ippi\$	s
issippi\$	s
ississippi\$m	
mississippi\$	
pi\$	p
ppi\$	i
sippi\$	s
sissippi\$	s
ssippi\$	i
ssissippi\$	i

## Naïve BWT algorithm

For a string  $S$ :

- 1 append  $\$$  to  $S$  ( $\$ < c, \forall c \in \Sigma$ )
- 2 list all cyclic rotations of  $S$
- 3 sort rows of matrix (lexicographically)
- 4 output last column of matrix

# Review: backwards search of pattern $P$ in string $S$

\$	i
i\$	p
ippi\$	s
issippi\$	s
ississippi\$m	
mississippi\$	
pi\$	p
ppi\$	i
sippi\$	s ←
sissippi\$	s
ssippi\$	i
ssissippi\$	i ←

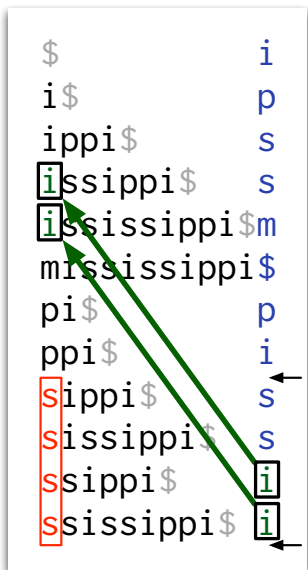
## Example

S S i S

↪ SA range: [8, 11]

[Ferragina & Manzini 2002]

# Review: backwards search of pattern $P$ in string $S$



## Example

S

S

i

S

↔ SA range: [8, 11]

[Ferragina & Manzini 2002]

# Review: backwards search of pattern $P$ in string $S$

\$	i
i\$	p
ippi\$	s
issippi\$	s ←
ississippi\$m	s ←
mississippi\$	
pi\$	p
ppi\$	i
sippi\$	s
sissippi\$	s
ssippi\$	i
ssissippi\$	i

Example

S S i S

↪ SA range: [3, 4]

[Ferragina & Manzini 2002]

# Review: backwards search of pattern $P$ in string $S$

\$	i
i\$	p
ippi\$	s
issippi\$	s ←
ississippi\$m	←
mississippi\$	
pi\$	p
ppi\$	i
sippi\$	s
issippi\$	s
ssippi\$	i
ssissippi\$	i

## Example

S      S      i      S

↪ SA range: [3, 4]

[Ferragina & Manzini 2002]

# Review: backwards search of pattern $P$ in string $S$

\$	i
i\$	p
ippi\$	s
issippi\$	s
ississippi\$m	
mississippi\$	
pi\$	p
ppi\$	i
sippi\$	s
<b>sis</b> sippi\$	s ←
ssippi\$	i ←
ssissippi\$	i

## Example

S            **s**            **i**            **s**

↪ SA range: [9, 9]

[Ferragina & Manzini 2002]

# Type 1 matching

Build **F**, a compressed suffix array of **T**.

	BWT
*a*aa*cacc*ac	a
*aa*aca*a*aa*cacc*ac	c
*aa*cacc*ac	a
*ac	c
*aca*a*aa*cacc*ac	a
*cacc*ac	a
a*a*aa*cacc*ac	c
a*aa*cacc*ac	*
a*aca*a*aa*cacc*ac	a
a*cacc*ac	a
aa*aca*a*aa*cacc*ac	*
aa*cacc*ac	*
ac	*
aca*a*aa*cacc*ac	*
acc*ac	c
c	a
c*ac	c
ca*a*aa*cacc*ac	a
cacc*ac	*
cc*ac	a

## Lemma

All  $occ_1$  Type 1 matches of  $P$  can be reported using:

Query time	$O( P  \log \sigma + occ_1 \log^{1+\epsilon} n)$
Index space	$(1 + o(1))n \log \sigma$ bits
Work space	$O(\log n)$ bits

# Matching prefixes of text segments

	BWT
*a*aa*cacc*ac	a
*aa*aca*a*aa*cacc*ac	c
*aa*cacc*ac	a
*ac	c
*aca*a*aa*cacc*ac	a
*cacc*ac	a
a*a*aa*cacc*ac	c
a*aa*cacc*ac	*
a*aca*a*aa*cacc*ac	a
a*cacc*ac	a
aa*aca*a*aa*cacc*ac	*
aa*cacc*ac	*
ac	* ←
aca*a*aa*cacc*ac	*
acc*ac	c ←
c	a
c*ac	c
ca*a*aa*cacc*ac	a
cacc*ac	*
cc*ac	a



# Matching prefixes of text segments

	BWT
*a*aa*cacc*ac	a
*aa*aca*a*aa*cacc*ac	c
*aa*cacc*ac	a
*ac	c
*aca*a*aa*cacc*ac	a
*cacc*ac	a
a*a*aa*cacc*ac	c
a*aa*cacc*ac	*
a*aca*a*aa*cacc*ac	a
a*cacc*ac	a
aa*aca*a*aa*cacc*ac	*
aa*cacc*ac	*
ac	*
aca*a*aa*cacc*ac	*
acc*ac	c
c	a
c*ac	c
ca*a*aa*cacc*ac	a
cacc*ac	*
cc*ac	a

$\text{rank}_*(BWT, i)$  queries determine a lexicographic range of text segments

# Type 2 matching

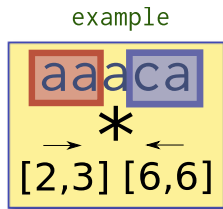
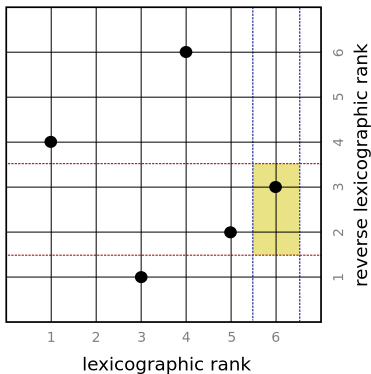
Use **F** and **R** to find candidate text segments

example		BWT		BWT <sup>R</sup>
	*a*aa*cacc*ac	a	*a*aca*aa*	a
	*aa*aca*a*aa*cacc*ac	c	*aa*a*aca*aa*	c
	*aa*cacc*ac	a	*aa*	a
	*ac	c	*aca*aa*	a
	*aca*a*aa*cacc*ac	a	*ca*ccac*aa*	a
	*cacc*ac	a	*ccac*aa*a*aca*aa*	a
	a*a*aa*cacc*ac	c	a*a*aca*aa*	a
	a*aa*cacc*ac	*	a*aa*	c
	a*aca*a*aa*cacc*ac	a	a*aca*aa*	*
	a*cacc*ac	a	a*	a
	aa*aca*a*aa*cacc*ac	*	a*ccac*aa*a*aca*aa*	c
	aa*cacc*ac	*	<b>aa</b> *a*aca*aa*	* ←
	ac	*	<b>aa</b> *	* ←
	aca*a*aa*cacc*ac	*	ac*aa*a*aca*aa*	c
	acc*ac	c	aca*aa*	*
	c	a	c*aa*a*aca*aa*	a
	c*ac	c	ca*aa*	a
	<b>ca</b> *a*aa*cacc*ac	a ←	ca*ccac*aa*a*aca*aa*	*
	<b>cacc</b> *ac	* ←	cac*aa*a*aca*aa*	c
	cc*ac	a	ccac*aa*a*aca*aa*	*

# Type 2 matching

Use  $Q$  to find compatible candidate [Bose *et al.*, 2009]

\*aa\*aca\*a\*aa\*cacc\*ac  
                  └───┬───┘  
                  match



## Lemma

All  $occ_2$  Type 2 matches of  $P$  can be reported using:

*Query time*  $O(|P| \log \sigma + (|P| + occ_2) \frac{\log d}{\log \log d})$

*Index space*  $(2 + o(1))n \log \sigma$  bits

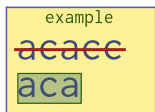
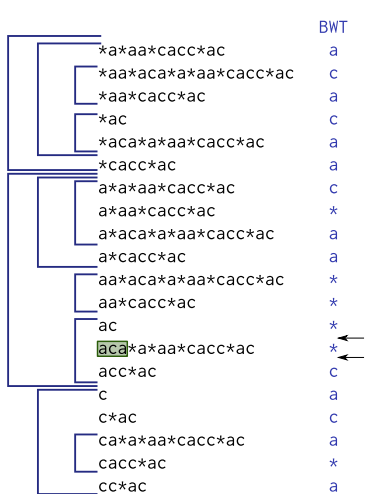
*Work space*  $O(m \log n)$  bits

T: \*aa\*aca\*a\*aa\*cacc\*ac  
P: acacaca

We will enhance **F** to support dictionary queries

- calculate matching statistics
- for each suffix, return text segments contained as a prefix

# Review: computing matching statistics



suffix	length	SA range
<b>c</b>	1	[15,19]
<b>cc</b>	2	[19,19]
<b>acc</b>	3	[14,14]
<b>cacc</b>	4	[18,18]
<b>acacc</b>	3	[13,13]

[Ohlebusch, Gog & Kügel 2010]

# Finding dictionary matches

## Case 1

example 'acc'  
sp = 15

	BWT
\$	c
*a*aa*cacc*ac\$	a
*aa*aca*a*aa*cacc*ac\$	\$
*aa*cacc*ac\$	a
*ac\$	c
*aca*a*aa*cacc*ac\$	a
*cacc*ac\$	a
a*a*aa*cacc*ac\$	c
a*aa*cacc*ac\$	*
a*aca*a*aa*cacc*ac\$	a
a*cacc*ac\$	a
aa*aca*a*aa*cacc*ac\$	*
aa*cacc*ac\$	*
ac\$	*
aca*a*aa*cacc*ac\$	*
acc*ac\$	c ←
	a ←
c\$	a
c*ac\$	c
ca*a*aa*cacc*ac\$	a
cacc*ac\$	*
cc*ac\$	a

# Finding dictionary matches

## Case 1

example 'acc'  
sp = 15

		BWT
	\$	c
	*a*aa*cacc*ac\$	a
	*aa*aca*a*aa*cacc*ac\$	\$
	*aa*cacc*ac\$	a
	*ac\$	c
	*aca*a*aa*cacc*ac\$	a
	*cacc*ac\$	a
	a*a*aa*cacc*ac\$	c
	a*aa*cacc*ac\$	*
	a*aca*a*aa*cacc*ac\$	a
	a*cacc*ac\$	a
(1, a)	aa*aca*a*aa*cacc*ac\$	*
(2, aa)	aa*cacc*ac\$	*
(3, aa)	ac\$	*
(4, ac)	aca*a*aa*cacc*ac\$	*
(5, aca)	acc*ac\$	c ←
	c\$	c ←
	c*ac\$	a
	ca*a*aa*cacc*ac\$	c
(6, caca)	cacc*ac\$	a
	cc*ac\$	*
		a



# Finding dictionary matches

## Case 1

example 'acc'  
sp = 15

	\$		BWT	c
	*a*aa*cacc*ac\$			a
	*aa*aca*a*aa*cacc*ac\$			\$
	*aa*cacc*ac\$			a
	*ac\$			c
	*aca*a*aa*cacc*ac\$			a
	*cacc*ac\$			a
	a*a*aa*cacc*ac\$			c
	a*aa*cacc*ac\$			*
	a*aca*a*aa*cacc*ac\$			a
	a*cacc*ac\$			a
(1, a)	aa*aca*a*aa*cacc*ac\$			*
(2, aa)	aa*cacc*ac\$			*
(3, aa)	ac\$			*
(4, ac)	aca*a*aa*cacc*ac\$			*
(5, aca)	acc*ac\$			c ←
	c\$			a ←
	c*ac\$			a
	ca*a*aa*cacc*ac\$			a
(6, caca)	cacc*ac\$			*
	cc*ac\$			a

BP  
1 2 3  
( ( ( ) ) )  
4 5  
( ( ) ) )  
6  
( )

[Munro & Ramen 2002]

# Finding dictionary matches

## Case 1

example 'acc'  
sp = 15

	BWT	B	BP
\$	c	0	
*a*aa*cacc*ac\$	a	0	
*aa*aca*a*aa*cacc*ac\$	\$	0	
*aa*cacc*ac\$	a	0	
*ac\$	c	0	
*aca*a*aa*cacc*ac\$	a	0	
*cacc*ac\$	a	0	
a*a*aa*cacc*ac\$	c	1	
a*aa*cacc*ac\$	*	0	
a*aca*a*aa*cacc*ac\$	a	0	
a*cacc*ac\$	a	0	
aa*aca*a*aa*cacc*ac\$	*	1	
aa*cacc*ac\$	*	0	
ac\$	*	1	
aca*a*aa*cacc*ac\$	*	1	
<b>acc*ac\$</b>	c	1	(
c\$	a	1	(
c*ac\$	c	0	(
ca*a*aa*cacc*ac\$	a	0	(
(6, caca) cacc*ac\$	*	1	(
cc*ac\$	a	1	(

[Ramen et al., 2002]

# Finding dictionary matches

## Case 1

example 'acc'  
sp = 15

	BWT	B	B'	
\$	c	0	0	
*a*aa*cacc*ac\$	a	0	1	
*aa*aca*a*aa*cacc*ac\$	\$	0	0	
*aa*cacc*ac\$	a	0	0	
*ac\$	c	0	1	
*aca*a*aa*cacc*ac\$	a	0	0	
*cacc*ac\$	a	0	0	
a*a*aa*cacc*ac\$	c	1	0	
a*aa*cacc*ac\$	*	0	1	
a*aca*a*aa*cacc*ac\$	a	0	0	
a*cacc*ac\$	a	0	1	
aa*aca*a*aa*cacc*ac\$	*	1	0	
aa*cacc*ac\$	*	0	1	
ac\$	*	1	0	
aca*a*aa*cacc*ac\$	*	1	0	
<b>acc*ac\$</b>	<b>c</b> ←	<b>1</b>	<b>1</b>	
c\$	a ←	1	0	
c*ac\$	c	0	1	
ca*a*aa*cacc*ac\$	a	0	0	
(6, caca) cacc*ac\$	*	1	1	
cc*ac\$	a	1		

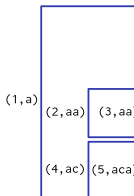
BP  
( ( ( ( ( ( ( ( ( ( ( ) ) ) ) ) ) ) ) )

[Grossi & Vitter 2000]

# Finding dictionary matches

## Case 1

```
example 'acc'  
sp = 15  
tmp1 = rank1(B, 15)
```



	BWT	B	B'
\$	c	0	0
*a*aa*cacc*ac\$	a	0	1
aa*aca*a*aa*cacc*ac\$	\$	0	0
aa*cacc*ac\$	a	0	0
*ac\$	c	0	1
*aca*a*aa*cacc*ac\$	a	0	0
*cacc*ac\$	a	0	0
a*a*aa*cacc*ac\$	c	1	0
a*aa*cacc*ac\$	*	0	1
a*aca*a*aa*cacc*ac\$	a	0	0
a*cacc*ac\$	a	0	1
aa*aca*a*aa*cacc*ac\$	*	1	0
aa*cacc*ac\$	*	0	1
ac\$	*	1	0
aca*a*aa*cacc*ac\$	*	1	0
<b>acc*ac\$</b>	<b>c</b> ←	<b>1</b>	<b>1</b>
c\$	a	1	0
c*ac\$	c	0	1
ca*a*aa*cacc*ac\$	a	0	0
(6, caca) cacc*ac\$	*	1	1
cc*ac\$	a	1	1

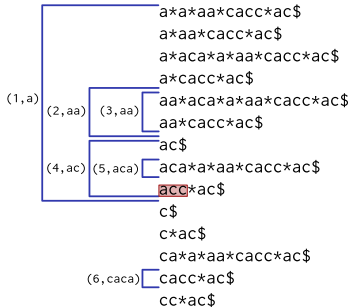
BP  
((((( ( ( ( ( ( ( ( (

1 2 3 4 5 6

# Finding dictionary matches

## Case 1

```
example 'acc'
sp = 15
tmp1 = rank1(B, 15)
tmp2 = select1(B', tmp1)
```



	BWT	B	B'
\$	c	0	0
*a*aa*cacc*ac\$	a	0	1
*aa*aca*a*aa*cacc*ac\$	\$	0	0
*aa*cacc*ac\$	a	0	0
*ac\$	c	0	1
*aca*a*aa*cacc*ac\$	a	0	0
*cacc*ac\$	a	0	0
a*a*aa*cacc*ac\$	c	1	0
a*aa*cacc*ac\$	*	0	1
a*aca*a*aa*cacc*ac\$	a	0	0
a*cacc*ac\$	a	0	1
aa*aca*a*aa*cacc*ac\$	*	1	0
aa*cacc*ac\$	*	0	1
ac\$	*	1	0
aca*a*aa*cacc*ac\$	*	1	0
acc*ac\$	c	1	1
c\$	a	1	0
c*ac\$	c	0	1
ca*a*aa*cacc*ac\$	a	0	0
(6, caca) cacc*ac\$	*	1	1
cc*ac\$	a	1	

BP  
 ( ( ( ( ( ( ( ( ( ( ( ( ) ) ) ) ) ) ) ) ) ) ) )

# Finding dictionary matches

## Case 1

```
example 'acc'  
sp = 15  
tmp1 = rank1(B, 15)  
tmp2 = select1(B', tmp1)  
pos = rank0(B', tmp2)+1
```

	BWT	B	B'
\$	c	0	0
*a*aa*cacc*ac\$	a	0	1
*aa*aca*a*aa*cacc*ac\$	\$	0	0
*aa*cacc*ac\$	a	0	0
*ac\$	c	0	1
*aca*a*aa*cacc*ac\$	a	0	0
*cacc*ac\$	a	0	0
a*a*aa*cacc*ac\$	c	1	0
a*aa*cacc*ac\$	*	0	1
a*aca*a*aa*cacc*ac\$	a	0	0
a*cacc*ac\$	a	0	1
aa*aca*a*aa*cacc*ac\$	*	1	0
aa*cacc*ac\$	*	0	0
ac\$	*	1	0
aca*a*aa*cacc*ac\$	*	1	0
acc*ac\$	c ←	1	1
c\$	a ←	1	0
c*ac\$	c	0	1
ca*a*aa*cacc*ac\$	a	0	0
(6, caca) cacc*ac\$	*	1	1
cc*ac\$	a	1	

BP  
1 2 3  
( ( ( ) ) )  
4 5  
( ( ) )  
6

if BP[pos] = ')' then find matching '('

# Finding dictionary matches

## Case 1

```
example 'acc'
sp = 15
tmp1 = rank1(B, 15)
tmp2 = select1(B', tmp1)
pos = rank0(B', tmp2)+1
```

	BWT	B	B'	BP
\$	c	0	0	
*a*aa*cacc*ac\$	a	0	1	
*aa*aca*a*aa*cacc*ac\$	\$	0	0	
*aa*cacc*ac\$	a	0	0	
*ac\$	c	0	1	
*aca*a*aa*cacc*ac\$	a	0	0	
*cacc*ac\$	a	0	0	
a*a*aa*cacc*ac\$	c	1	0	
a*aa*cacc*ac\$	*	0	1	
a*aca*a*aa*cacc*ac\$	a	0	0	
a*cacc*ac\$	a	0	1	
aa*aca*a*aa*cacc*ac\$	*	1	0	
aa*cacc*ac\$	*	0	1	
ac\$	*	1	0	
aca*a*aa*cacc*ac\$	*	1	0	
<b>acc*ac\$</b>	<b>c</b> ←	1	1	
c\$	<b>a</b> ←	1	0	
c*ac\$	c	0	1	
ca*a*aa*cacc*ac\$	a	0	0	
(6, caca) cacc*ac\$	*	1	1	
cc*ac\$	a	1		

(1, a)				
(2, aa)				
(3, aa)				
(4, ac)				
(5, aca)				
(6, caca)				

(	1
(	2
(	3
)	4
)	5
)	6

if BP[pos] = ')' then find matching '('

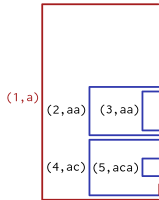
# Finding dictionary matches

## Case 1

```

example 'acc'
sp = 15
tmp1 = rank1(B, 15)
tmp2 = select1(B', tmp1)
pos = rank0(B', tmp2)+1
    
```

	BWT	B	B'	BP
\$	c	0	0	
*a*aa*cacc*ac\$	a	0	1	
*aa*aca*a*aa*cacc*ac\$	\$	0	0	
*aa*cacc*ac\$	a	0	0	
*ac\$	c	0	1	
*aca*a*aa*cacc*ac\$	a	0	0	
*cacc*ac\$	a	0	0	
a*a*aa*cacc*ac\$	c	1	0	
a*aa*cacc*ac\$	*	0	1	
a*aca*a*aa*cacc*ac\$	a	0	0	
a*cacc*ac\$	a	0	1	
aa*aca*a*aa*cacc*ac\$	*	1	0	
aa*cacc*ac\$	*	0	1	
ac\$	*	1	0	
aca*a*aa*cacc*ac\$	*	1	0	
acc*ac\$	c ←	1	1	
c\$	a ←	1	0	
c*ac\$	c	0	1	
ca*a*aa*cacc*ac\$	a	0	0	
(6, caca) cacc*ac\$	*	1	1	
cc*ac\$	a	1		



BP  
 1 2 3  
 ( ( ( ) ) )  
 4 5  
 ( ( ) )  
 6  
 ( )

if BP[pos] = ')' then find matching '('



# Finding dictionary matches

## Case 2

example 'acc'  
sp = 16

	BWT	B	B'	
	c	0	0	
	a	0	1	
	\$	0	0	
	a	0	0	
	c	0	1	
	a	0	0	
	a	0	0	
	c	0	1	
	a	0	0	
	a	0	0	
	c	1	0	
	*	0	1	
	a	0	0	
	a	0	1	
	a	0	0	
	*	1	0	
	*	0	1	
	*	1	0	
	*	1	0	
	c	1	1	
	a ←	1	0	
	c	0	1	
	a	0	0	
	a	0	0	
	*	1	1	
	a ←	1		

	B	B'	BP
	0	0	
	0	1	
	0	0	
	0	0	
	0	1	
	0	0	
	0	0	
	0	1	
	1	0	
	1	0	
	1	1	
	1	0	
	0	1	
	0	0	
	1	1	
	1		

	BP
	1 2 3
	( ( ( ) ) ( ) ) ( )
	4 5
	( ( ) ) ( )
	6
	( )

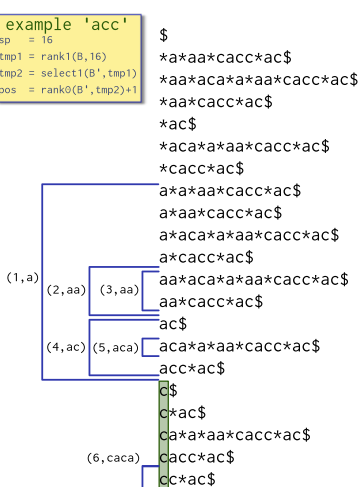
  

Diagram Labels	Row
(1, a)	a*a*aa*cacc*ac\$
	a*aa*cacc*ac\$
	a*aca*a*aa*cacc*ac\$
	a*cacc*ac\$
(2, aa)	aa*aca*a*aa*cacc*ac\$
(3, aa)	aa*cacc*ac\$
	ac\$
(4, ac)	aca*a*aa*cacc*ac\$
(5, aca)	acc*ac\$
	c\$
	c*ac\$
	ca*a*aa*cacc*ac\$
(6, caca)	cacc*ac\$
	cc*ac\$

# Finding dictionary matches

## Case 2

```
example 'acc'
sp = 16
tmp1 = rank1(B, 16)
tmp2 = select1(B', tmp1)
pos = rank0(B', tmp2) + 1
```



BWT	B	B'	BP
c	0	0	
a	0	1	
\$	0	0	
a	0	0	
c	0	1	
a	0	0	
a	0	0	
c	1	0	
*	0	1	
a	0	0	
a	0	1	
*	1	0	
*	0	1	
*	1	0	
*	1	0	
c	1	1	
a	1	0	
c	0	1	
a	0	0	
*	1	1	
a	1		

if BP[pos] = ' ( ' then find parent interval

# Finding dictionary matches

## Case 2

```
example 'acc'
sp = 16
tmp1 = rank1(B, 16)
tmp2 = select1(B', tmp1)
pos = rank0(B', tmp2) + 1
```

	BWT	B	B'	
\$	c	0	0	
*a*aa*cacc*ac\$	a	0	1	
*aa*aca*a*aa*cacc*ac\$	\$	0	0	
*aa*cacc*ac\$	a	0	0	
*ac\$	c	0	1	
*aca*a*aa*cacc*ac\$	a	0	0	
*cacc*ac\$	a	0	0	
a*a*aa*cacc*ac\$	c	1	0	
a*aa*cacc*ac\$	*	0	1	
a*aca*a*aa*cacc*ac\$	a	0	0	
a*cacc*ac\$	a	0	1	
aa*aca*a*aa*cacc*ac\$	*	1	0	
aa*cacc*ac\$	*	0	1	
ac\$	*	1	0	
aca*a*aa*cacc*ac\$	*	1	0	
acc*ac\$	c	1	1	
c\$	a ←	1	0	
c*ac\$	c	0	1	
ca*a*aa*cacc*ac\$	a	0	0	
cacc*ac\$	*	1	1	
cc*ac\$	a ←	1		

	BP
	1 2 3
	(( ( ( ) ) ) )
	4 5
	(( ( ) ) )
	6
	(( ) )

NO MATCH

if BP[pos] = ' ( ' then find parent interval

## Theorem

A full-text dictionary using  $(1 + o(1))n \log \sigma + O(n) + O(d \log n)$  bits supports the following operations efficiently:

- report/count text segments *containing*  $P$
- report/count text segments *prefixed by*  $P$
- report/count text segments *contained in*  $P$

# Type 3 matching

## Preliminaries

**T:** \*aa\*aca\*acaa\*cac\*ac

**P:** acacacac

Every text segment contained within **P** is a [candidate](#)

# Type 3 matching

## Preliminaries

**T:** \*aa\*aca\*acaacac\*ac

          └──┬──┘   └──┬──┘  
          prefix   suffix  
          condition condition

**P:**       acacacac

Every text segment contained within **P** is a **candidate**

A candidate is **valid** if it satisfies:

- 1 suffix condition
- 2 prefix condition

# Algorithm overview

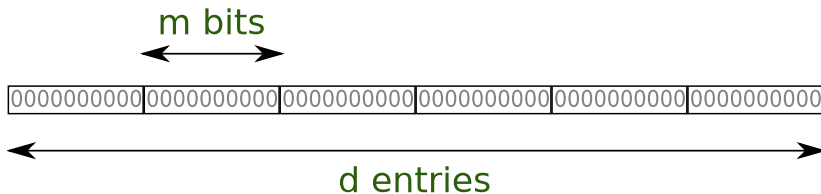
We design a dynamic programming algorithm:

- $|\mathbf{P}|$  phases of algorithm ( $i = |\mathbf{P}|, \dots, 2, 1$ )
- phase  $i$  considers candidates prefixing  $\mathbf{P}[i..|\mathbf{P}|]$

# Algorithm overview

We design a dynamic programming algorithm:

- $|\mathbf{P}|$  phases of algorithm ( $i = |\mathbf{P}|, \dots, 2, 1$ )
- phase  $i$  considers candidates prefixing  $\mathbf{P}[i..|\mathbf{P}|]$
- potential matches are tracked in a working array  $\mathbf{W}$





# Type 3 matching

## Case 1

- ↪ Text segment  $j$  matches  $\mathbf{P}[j..|\mathbf{P}|]$
- ↪ **(case 1)**  $\mathbf{P}$  must end within text segment  $j + 1$

T: ...\*aa\*aca\*acaa\*cac\*ac

P:       acacacac

0000000000 0000000000 0000000000 0000000000 0000000000 0000000000

# Type 3 matching

## Case 1

↪ Text segment  $j$  matches  $\mathbf{P}[j..|\mathbf{P}|]$

↪ **(case 1)**  $\mathbf{P}$  must end within text segment  $j + 1$

- 1 verify suffix condition

suffix  
condition



T: ...\*aa\*aca\*acaa\*cac\*ac

P: acacacac

0000000000 0000000000 0000000000 0000000000 0000000000 0000000000

# Type 3 matching

## Case 1

- ↪ Text segment  $j$  matches  $\mathbf{P}[j..|\mathbf{P}|]$
- ↪ **(case 1)**  $\mathbf{P}$  must end within text segment  $j + 1$

- 1 verify suffix condition
- 2 **if**  $\mathbf{P}$  must begin in text segment  $j - 1$  **then**  
verify prefix condition

prefix  
condition



T: ...\*aa\*aca\*acaa\*cac\*ac

P: acacacac



# Type 3 matching

## Case 1

↪ Text segment  $j$  matches  $\mathbf{P}[j..|\mathbf{P}|]$

↪ **(case 1)**  $\mathbf{P}$  must end within text segment  $j + 1$

- 1 verify suffix condition
- 2 **if**  $\mathbf{P}$  must begin in text segment  $j - 1$  **then**  
verify prefix condition

T: ...\*aa\*aca\*acaa\*cac\*ac

P: acacacac

0000000000 0000000000 0000000000 0000000000 0000000000 0000000000

# Type 3 matching

## Case 1

- ↪ Text segment  $j$  matches  $\mathbf{P}[j..|\mathbf{P}|]$
- ↪ **(case 1)**  $\mathbf{P}$  must end within text segment  $j + 1$

- 1 verify suffix condition
- 2 **if**  $\mathbf{P}$  must begin in text segment  $j - 1$  **then**  
verify prefix condition
- 3 **else**  
record partial match

set  $i$  bit of  $W[j]$  to 1

T: ...\*aa\*aca\*acaa\*cac\*ac

P: aaaaacacac

0000000000 0000000000 0001000000 0000000000 0000000000 0000000000

# Type 3 matching

## Case 2

- ↪ Text segment  $j$  matches  $\mathbf{P}[i..|\mathbf{P}|]$
- ↪ **(case 2)**  $\mathbf{P}$  must contain text segment  $j + 1$

**T:** aaa\*aa\*aca\*acaa\*cac\*ac

**P:** aaaaaacacac

0000000000|0000000000|0001000000|0000000000|0000000000|0000000000

# Type 3 matching

## Case 2

- ↪ Text segment  $j$  matches  $\mathbf{P}[i..|\mathbf{P}|]$
- ↪ **(case 2)**  $\mathbf{P}$  must contain text segment  $j + 1$

- 1 verify suffix condition by checking  $\mathbf{W}$

verify bit  $(i-L[j]-1)$  bit of  $\mathbf{W}[j-1]$

T: aaa\*aa\*aca\*acaa\*cac\*ac

P: aaaaaacacac

0000000000|0000000000|0001000000|0000000000|0000000000|0000000000

# Type 3 matching

## Case 2

↪ Text segment  $j$  matches  $\mathbf{P}[j..|\mathbf{P}|]$   
↪ **(case 2)**  $\mathbf{P}$  must contain text segment  $j + 1$

- 1 verify suffix condition by checking  $\mathbf{W}$
- 2 **if**  $\mathbf{P}$  must begin in text segment  $j - 1$  **then**  
verify prefix condition
- 3 **else**  
record partial match

**T:** aaa\*aa\*aca\*acaa\*cac\*ac

**P:** aaaaaacacac

0000000000 0000000000 0001000000 0000000000 0000000000 0000000000



# Type 3 matching

## Lemma

All  $occ_3$  Type 3 matches of  $P$  can be reported using:

<i>Query time</i>	$O( P  \log \sigma + \gamma)$
<i>Index space</i>	$(2 + o(1))n \log \sigma$ bits
<i>Work space</i>	$O(m \log n + md)$ bits

where  $\gamma$  is number of text segments contained in  $P$

## Theorem

*Given a text  $T$  of length  $n$  containing  $d$  wildcards, all matches of a pattern  $P$  can be reported using:*

*Query time*  $O(|P| \log \sigma + m \frac{\log d}{\log \log d} + occ_1 \log n + occ_2 \frac{\log d}{\log \log d} + \gamma)$

*Index space*  $(2 + o(1))n \log \sigma + O(n) + O(m \log n)$  bits

*Work space*  $O(m \log n + md)$  bits

## Theorem

*Given a text  $T$  of length  $n$  containing  $d$  wildcards, all matches of a pattern  $P$  can be reported using:*

*Query time*  $O(|P| \log \sigma + m \frac{\log d}{\log \log d} + occ_1 \log n + occ_2 \frac{\log d}{\log \log d} + \gamma)$

*Index space*  $(2 + o(1))n \log \sigma + O(n) + O(m \log n)$  bits

*Work space*  $O(m \log n + md)$  bits

## Theorem

Given a text  $T$  of length  $n$  containing  $d$  wildcards, all matches of a pattern  $P$  can be reported using:

*Query time*  $O(|P| \log \sigma + m \frac{\log d}{\log \log d} + occ_1 \log n + occ_2 \frac{\log d}{\log \log d} + \gamma)$

*Index space*  $(2 + o(1))n \log \sigma + O(n) + O(m \log n)$  bits

*Work space*  $O(m \log n + md)$  bits

# Significance of working space reduction

## Short read alignment to human genome

**|T|**  $\approx$  3 billion bases  
**d**  $\approx$  3 million bases  
**m**  $\approx$  32 – 64

Working space reduced from GBs [Tam *et al.*, 2009] to 10's of MBs

## Future Work

Can index space be reduced to  $(1 + o(1))n \log \sigma$  bits?

$\rightsquigarrow$  **YES, but** with query time:

$O(\max(m^2 \frac{\log \sigma}{\log \log \sigma}, \text{current query time}))$

Thanks to Anne Condon and anonymous reviewers.