

The BWT as a compressed index

Off-line text searching

We have a collection of documents (think of project Gutenberg or Wikipedia) and we want to search for information (substrings) in it.

Since the collection is known in advance, we speed-up the search building an index for the collection

The Suffix Array

The Suffix Array is the simplest indexing data structure supporting fast pattern searching

Consider for example

T = `swiss·miss·missing`

The Suffix Array

swiss·miss·missing

We consider all the
suffixes of the input text

1	swiss·miss·missing
2	wiss·miss·missing
3	iss·miss·missing
4	ss·miss·missing
5	s·miss·missing
6	·miss·missing
7	miss·missing
8	iss·missing
9	ss·missing
10	s·missing
11	·missing
12	missing
13	issing
14	ssing
15	sing
16	ing
17	ng
18	g

The Suffix Array

swiss·miss·missing

We consider all the
suffixes of the input text

... and we sort them
in lexicographic order

6	·miss·missing
11	·missing
18	g
16	ing
3	iss·miss·missing
8	iss·missing
13	issing
7	miss·missing
12	missing
17	ng
5	s·miss·missing
10	s·missing
15	sing
4	ss·miss·missing
9	ss·missing
14	ssing
1	swiss·miss·missing
2	wiss·miss·missing

The Suffix Array

swiss·miss·missing

Using binary search
we find that miss
appears starting at
positions 7 and 12

6	·miss·missing
11	·missing
18	g
16	ing
3	iss·miss·missing
8	iss·missing
13	issing
7	miss ·missing
12	miss ing
17	ng
5	s·miss·missing
10	s·missing
15	sing
4	ss·miss·missing
9	ss·missing
14	ssing
1	swiss·miss·missing
2	wiss·miss·missing

Using the suffix array we can find all the occurrences of any pattern P in T using binary search in $O(|P| \log |T|)$ time.

Enriching the suffix array with additional information we can reduce search time to $O(|P| + \log |T|)$.

The Suffix Array

swiss·miss·missing

To represent the Suffix
Array we use $|T|$ integers
in the range $1 \dots |T|$

$\Rightarrow |T| \log |T|$ bits

6	·miss·missing
11	·missing
18	g
16	ing
3	iss·miss·missing
8	iss·missing
13	issing
7	miss·missing
12	missing
17	ng
5	s·miss·missing
10	s·missing
15	sing
4	ss·miss·missing
9	ss·missing
14	ssing
1	swiss·miss·missing
2	wiss·miss·missing

Burrows-Wheeler Transform

swiss · miss · missing

To compute the **BWT** we go on transforming each suffix into a cyclic shift of **T**.

The last column of the resulting matrix is the **BWT**.

6	·miss·missingswiss	s
11	·missingswiss·miss	s
18	gswiss·miss·missi	n
16	ingswiss·miss·mis	s
3	iss·miss·missing	w
8	iss·missingswiss·	m
13	issingswiss·miss·	m
7	miss·missingswiss	·
12	missingswiss·miss	·
17	ngswiss·miss·miss	i
5	s·miss·missingswi	s
10	s·missingswiss·mi	s
15	singswiss·miss·mi	s
4	ss·miss·missingsw	i
9	ss·missingswiss·m	i
14	ssingswiss·miss·m	i
1	swiss·miss·missin	g
2	wiss·miss·missing	s

Burrows-Wheeler Transform

swiss · miss · missing

We discard the Suffix Array
but we keep track of the row
containing the original string



6	·miss·missing	swiss
11	·missing	swiss·miss
18	g	swiss·miss·missi
16	ing	swiss·miss·miss
3	iss·miss·missing	swiss
8	iss·missing	swiss·
13	issing	swiss·miss·
7	miss·missing	swiss·
12	missing	swiss·miss·
17	ng	swiss·miss·miss
5	s·miss·missing	swiss
10	s·missing	swiss·miss
15	sing	swiss·miss·miss
4	ss·miss·missing	swiss
9	ss·missing	swiss·miss
14	ssing	swiss·miss·miss
1	swiss·miss·missin	g
2	wiss·miss·missing	s

Search using the BWT

Suppose we want to count the occurrences of mis

Working with only F and L we successively find the range of rows prefixed by:
s, is, mis

F		L
•	miss·missingswiss	s
•	missingswiss·mis	s
g	swiss·miss·missi	n
i	ngswiss·miss·mis	s
i	ss·miss·missings	w
i	ss·missingswiss·	m
i	ssingswiss·miss·	m
m	iss·missingswiss	•
m	issingswiss·miss	•
n	gswiss·miss·miss	i
s	•miss·missingswi	s
s	•missingswiss·mi	s
s	ingswiss·miss·mi	s
s	s·miss·missingsw	i
s	s·missingswiss·m	i
s	singswiss·miss·m	i
s	wiss·miss·missin	g
w	iss·miss·missing	s

Backward search

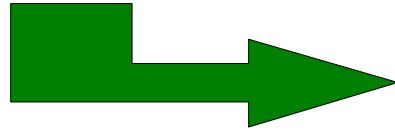
The rows prefixed
by **s** are easy to find

F
.
.
g
i
i
i
i
m
m
n
s
s
s
s
s
s
s
s
W

L
s
s
n
s
w
m
m
.
.
i
s
s
s
i
i
i
i
g
s

Backward search

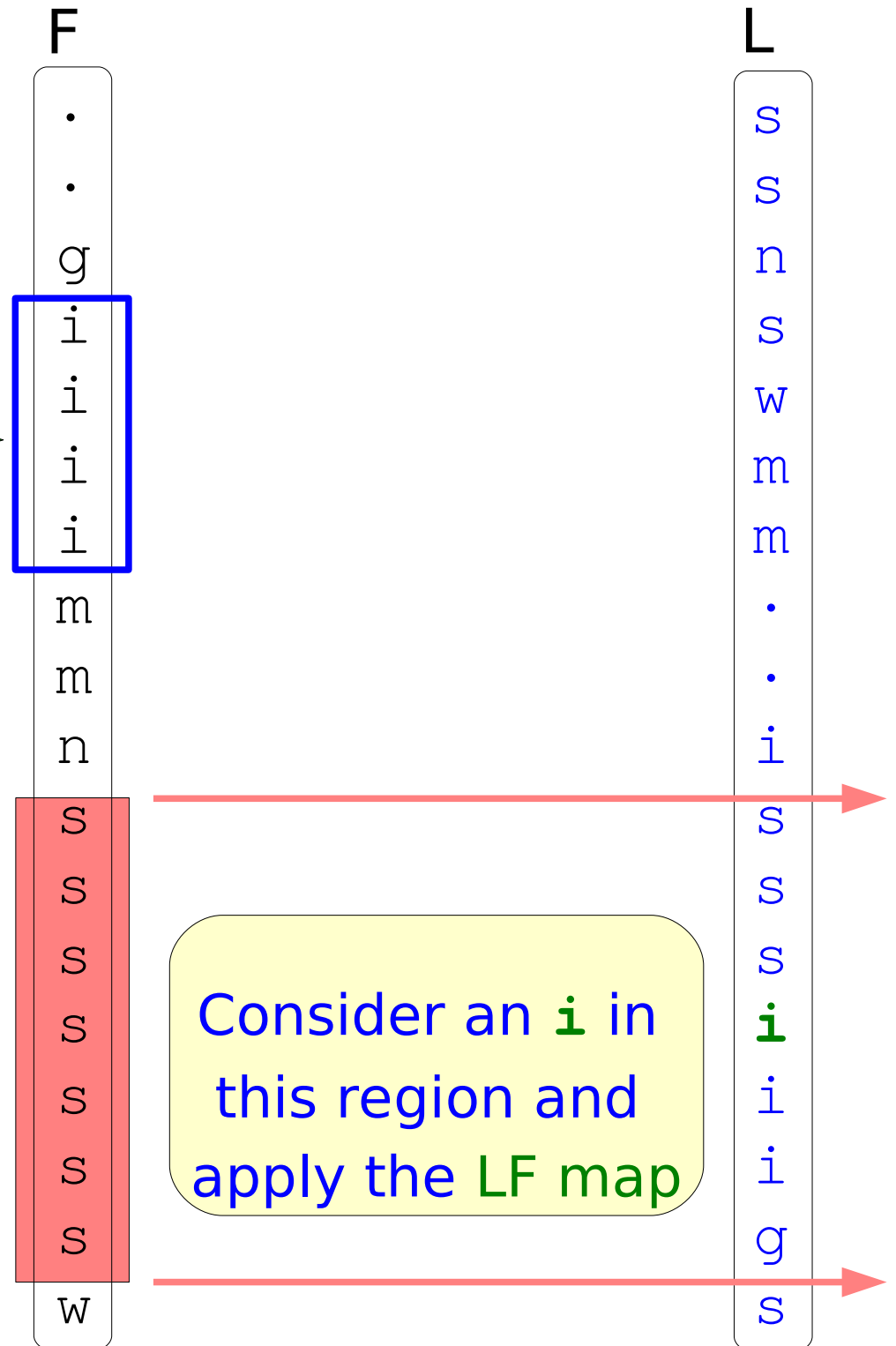
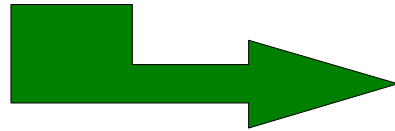
The rows prefixed by **is**
are a subset of **these** and
are consecutive.



F		L
.		s
.		s
g		n
i		s
i		w
i		m
i		m
m		.
m		.
n		i
s		s
s		s
s		s
s		i
s		i
s		i
s		g
s		s
w		

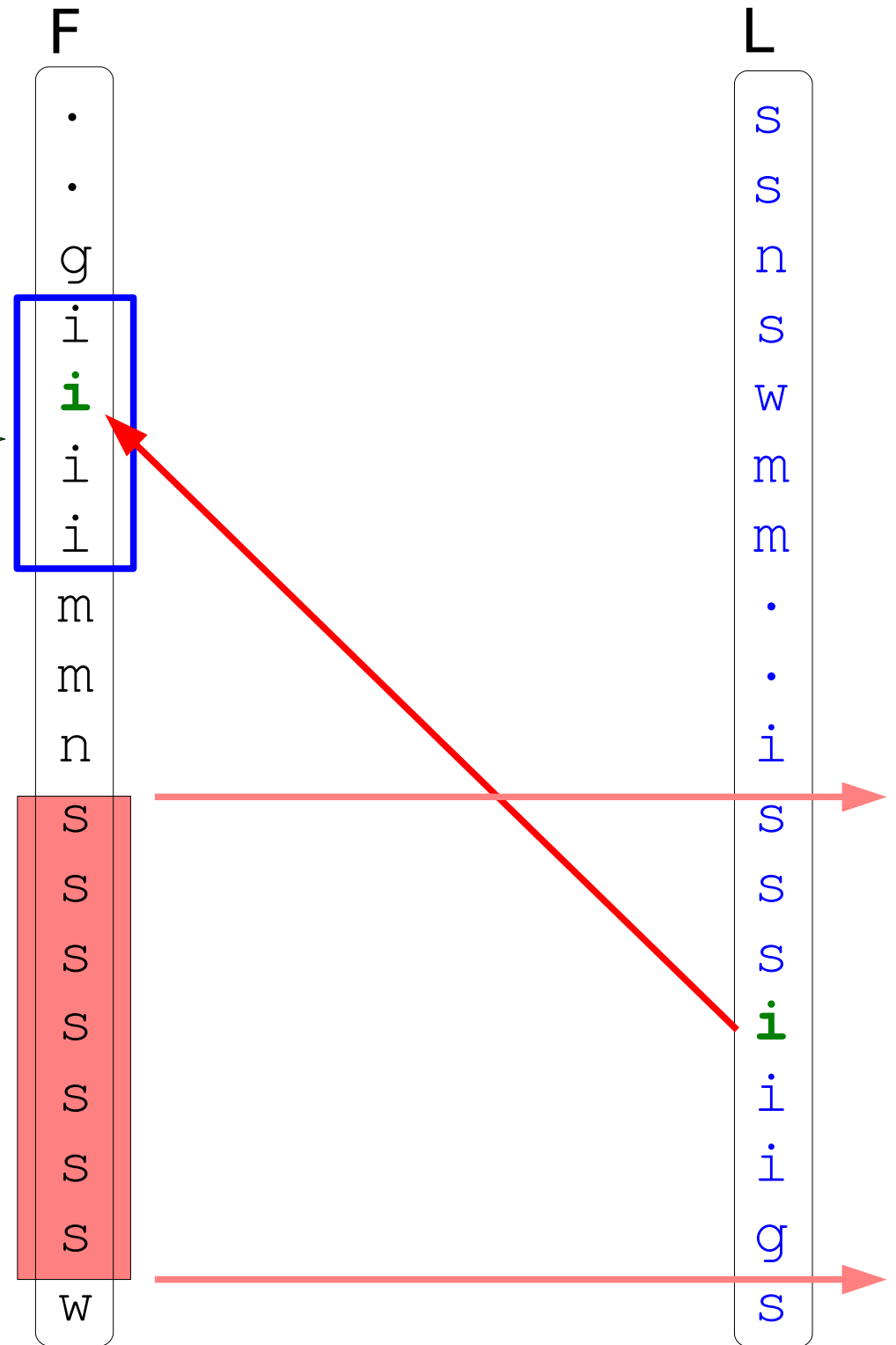
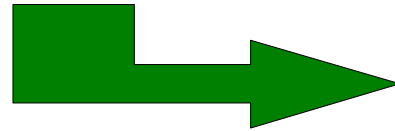
Backward search

The rows prefixed by **is**
are a subset of **these** and
are consecutive.

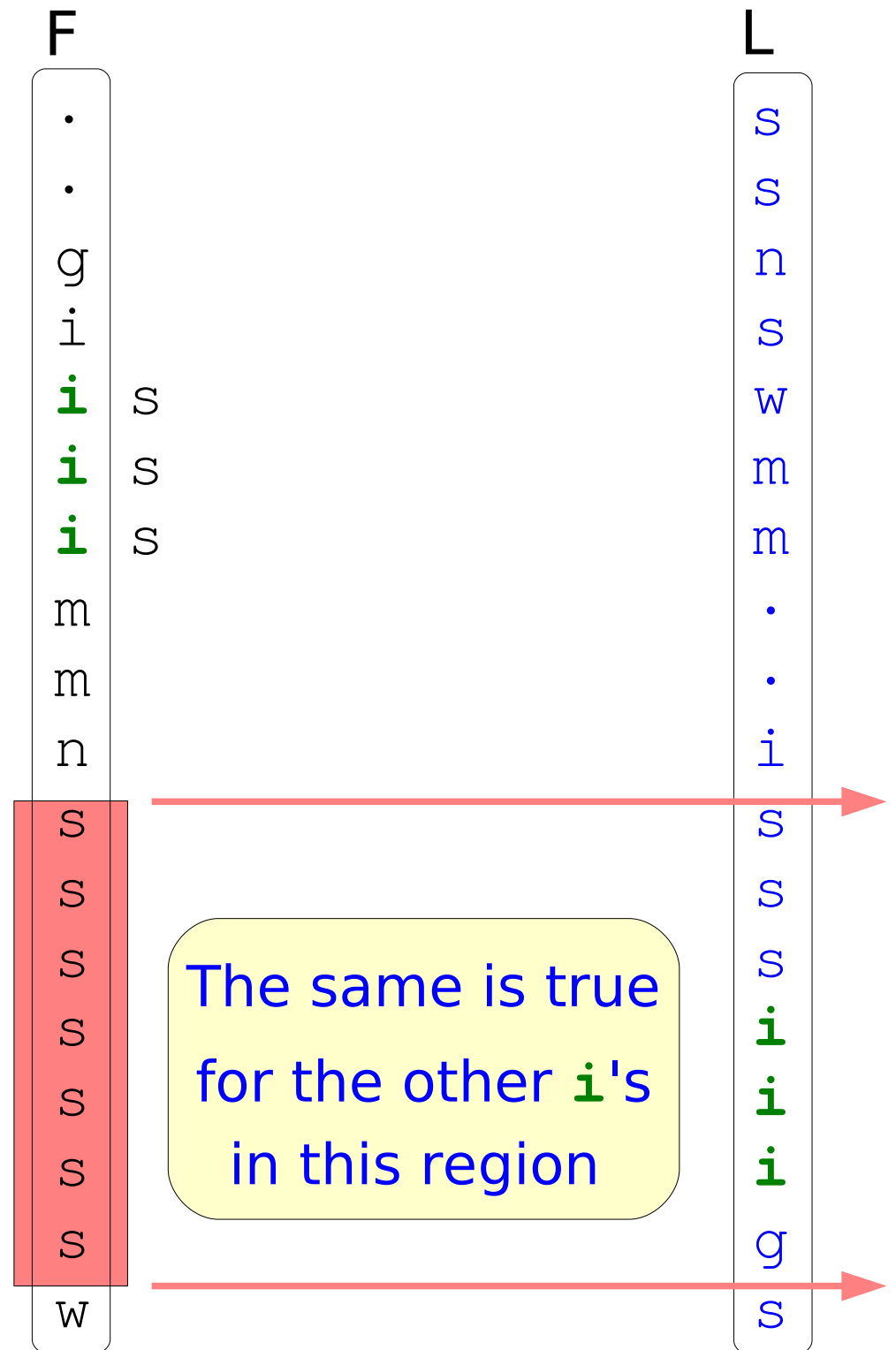


Backward search

The rows prefixed by **is**
are a subset of **these** and
are consecutive.



Backward search



Backward search

We have found the rows
prefixed by **is**!

F

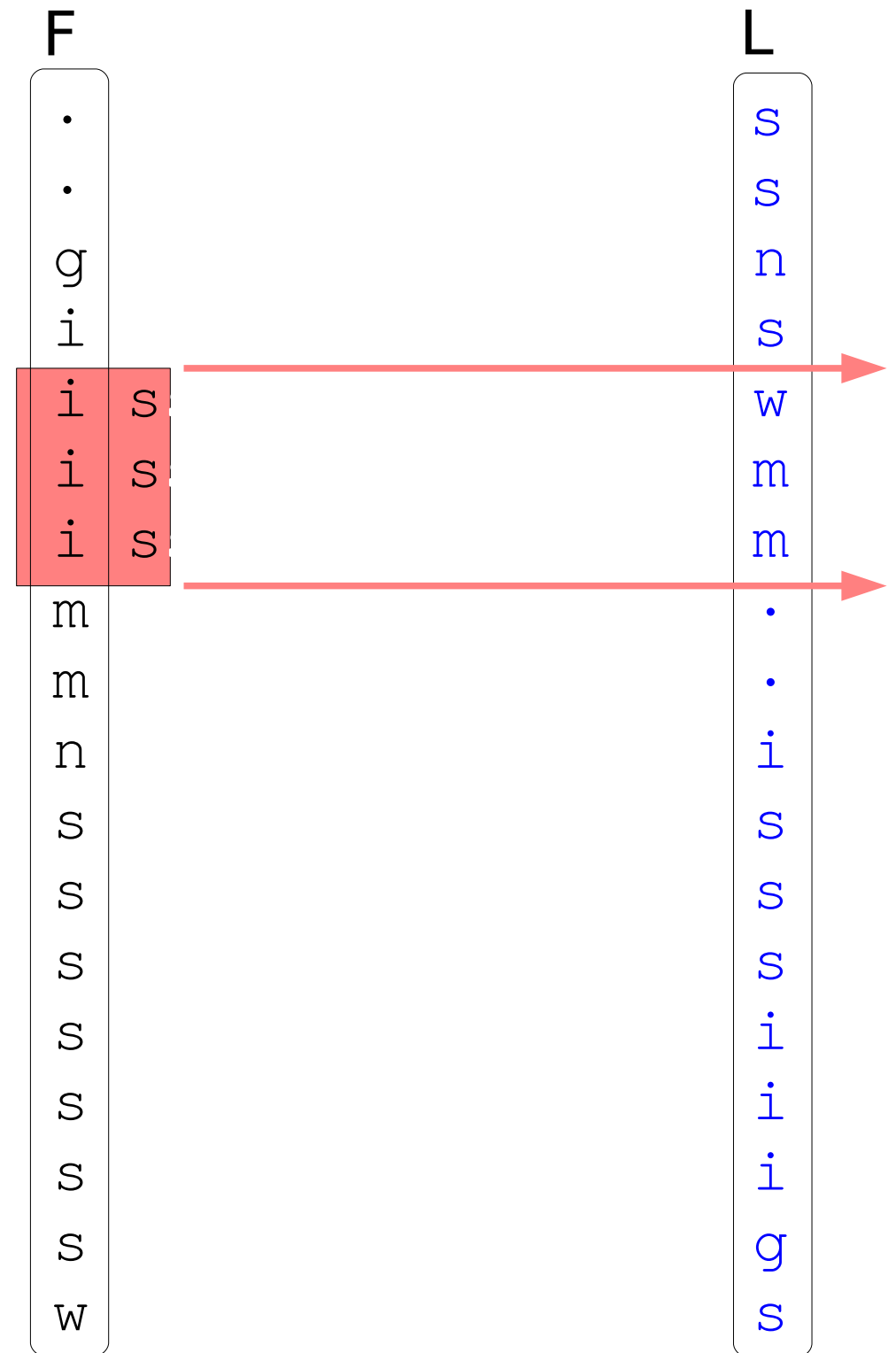
.
.
g
i
i
i
i
i
m
m
n
s
s
s
s
s
s
s
s
s
W

L

s
s
n
s
w
m
m
.
.
i
s
s
s
i
i
i
g
s

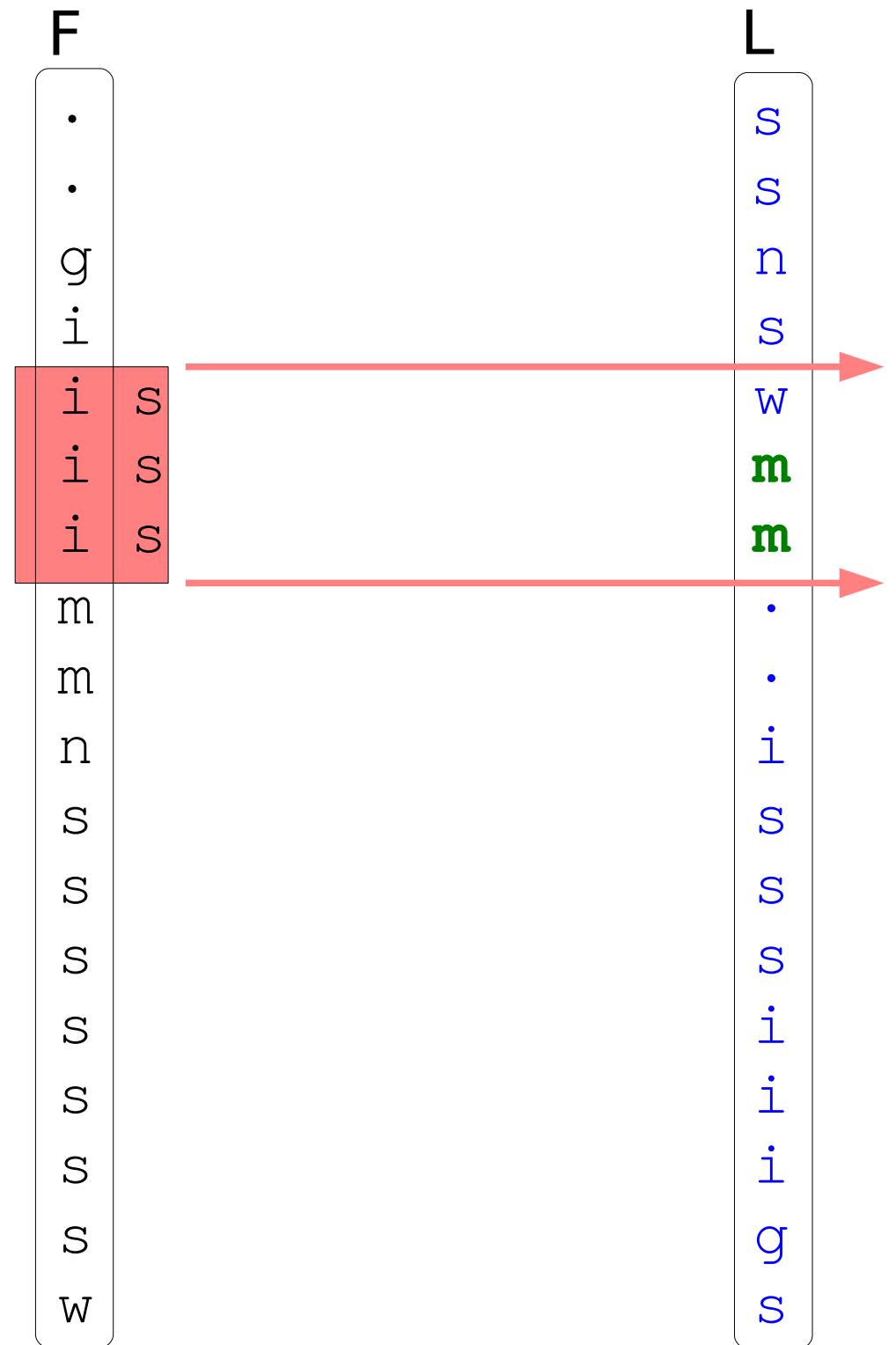
Backward search

To find the rows prefixed
by **mis** we proceed in
the same way



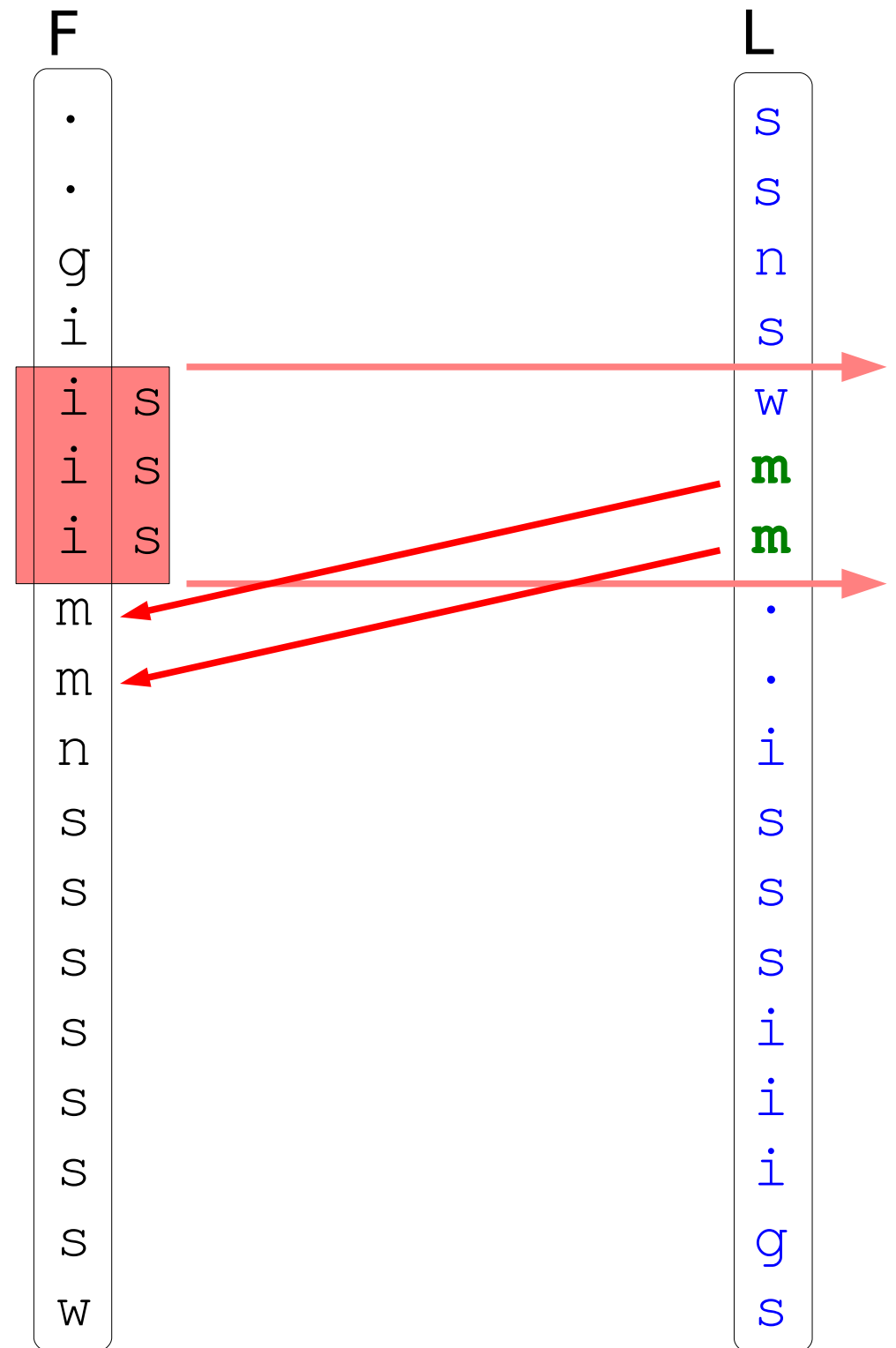
Backward search

To find the rows prefixed
by **mis** we proceed in
the same way



Backward search

To find the rows prefixed
by **mis** we proceed in
the same way



Backward search

We have found the rows
prefixed by **mis**!

F	L
.	s
.	s
g	n
i	s
i	w
i	m
i	m
m	.
m	.
n	i
s	s
s	s
s	s
s	i
s	i
s	i
s	g
w	s

Backward search

We have found the rows
prefixed by **mis**!

Even if we only have
F and L we can work
as if we had the SA

F		L
·	miss·missingswiss	s
·	missingswiss·mis	s
g	swiss·miss·missi	n
i	ngswiss·miss·mis	s
i	ss·miss·missing	w
i	ss·missingswiss·	m
i	ssingswiss·miss·	m
m	iss·missingswiss	·
m	issingswiss·miss	·
n	gswiss·miss·miss	i
s	·miss·missingswi	s
s	·missingswiss·mi	s
s	ingswiss·miss·mi	s
s	s·miss·missingsw	i
s	s·missingswiss·m	i
s	singswiss·miss·m	i
s	wiss·miss·missin	g
w	iss·miss·missing	s

A closer look

The basic step is going from the rows prefixed by **p** to the rows prefixed by **cp**

113
114
115
116
117

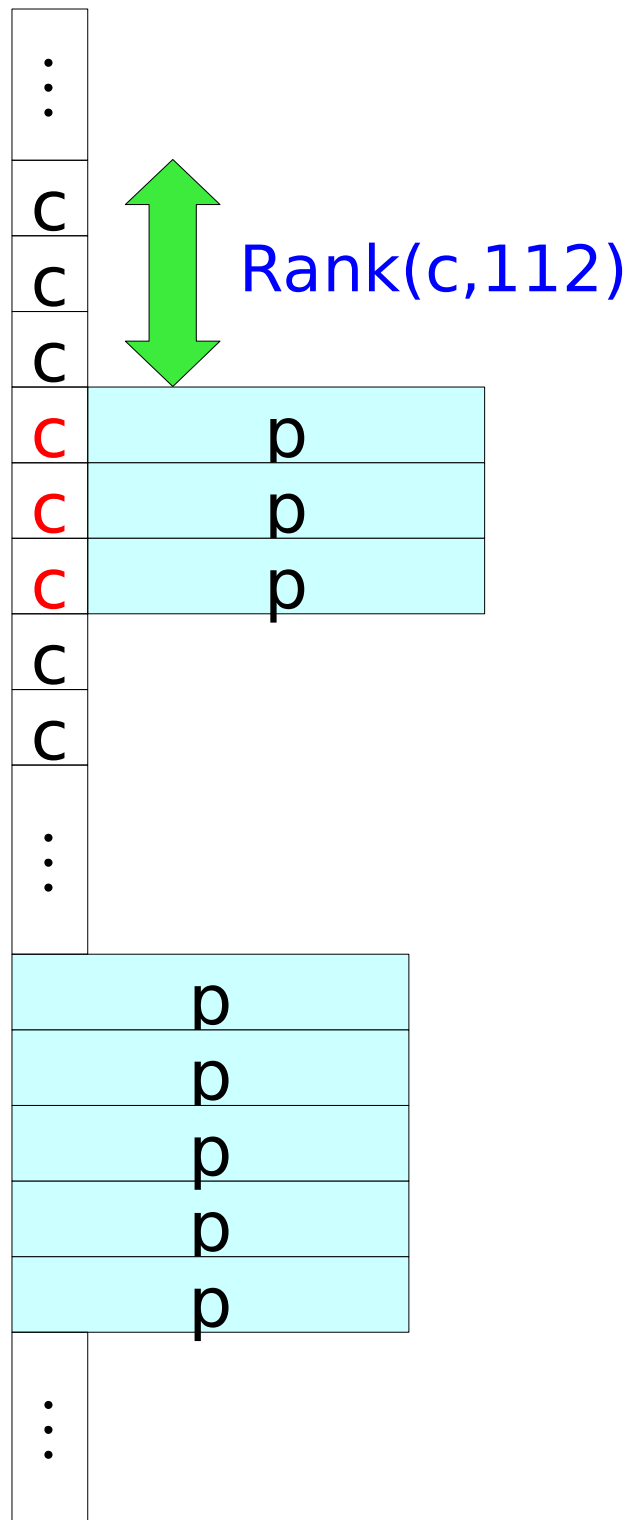
⋮	
c	
c	
c	
c	p
c	p
c	p
c	
c	
⋮	
	p
	p
	p
	p
	p
⋮	

c
c
c
c
c
y
x
c
c

A closer look

The basic step is going from the rows prefixed by **p** to the rows prefixed by **cp**

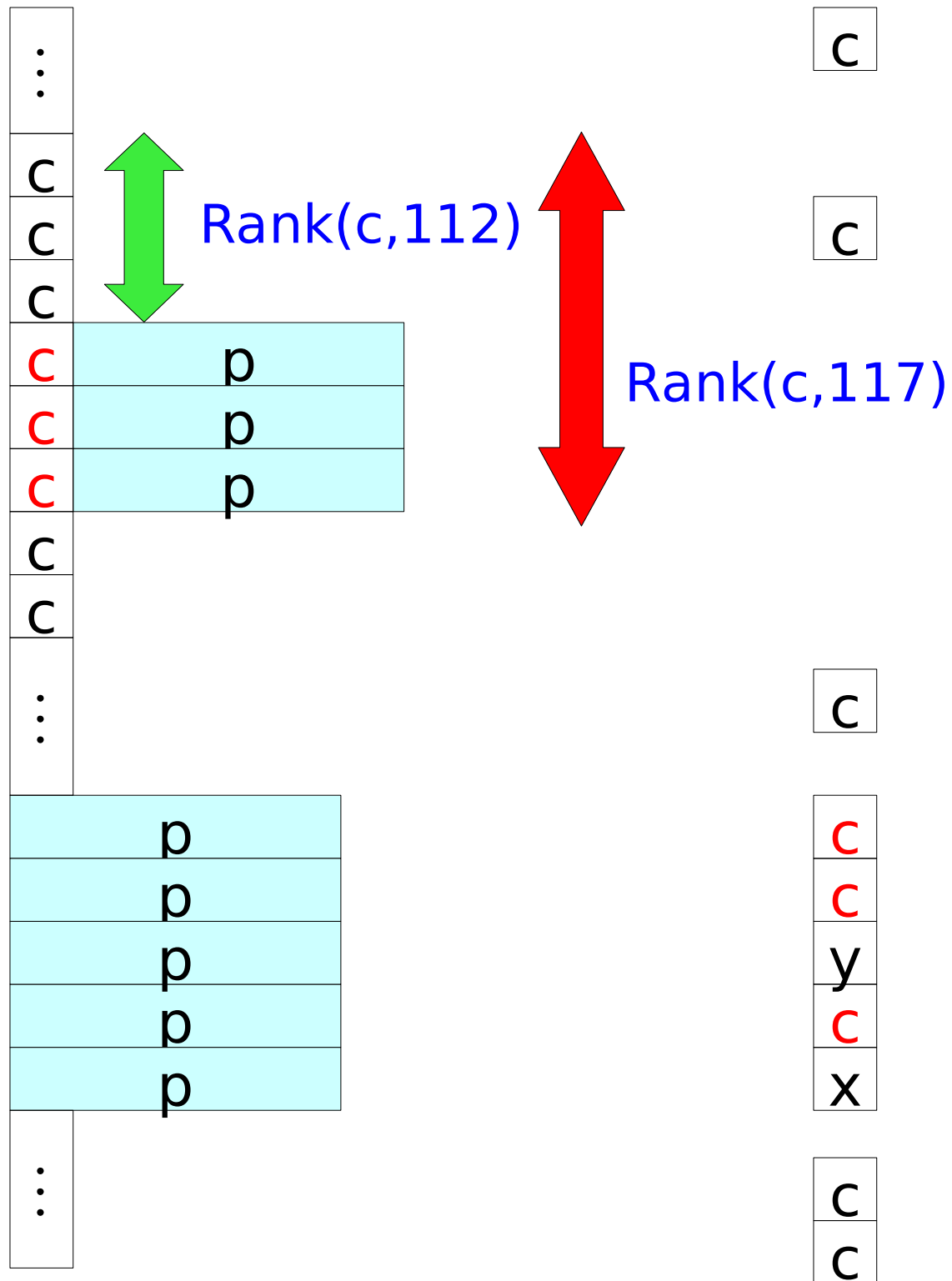
113
114
115
116
117



A closer look

The basic step is going from the rows prefixed by **p** to the rows prefixed by **cp**

113
114
115
116
117



Summing up

Each basic step requires two rank queries on L

We can do a rank query in $O(\log|A|)$ time on a compressed sequence.

Finding the range of rows prefixed by a pattern P takes $O(|P| \log|A|)$ time (no dependency on $|T|!$).

We have a compressed representation of T supporting fast queries.