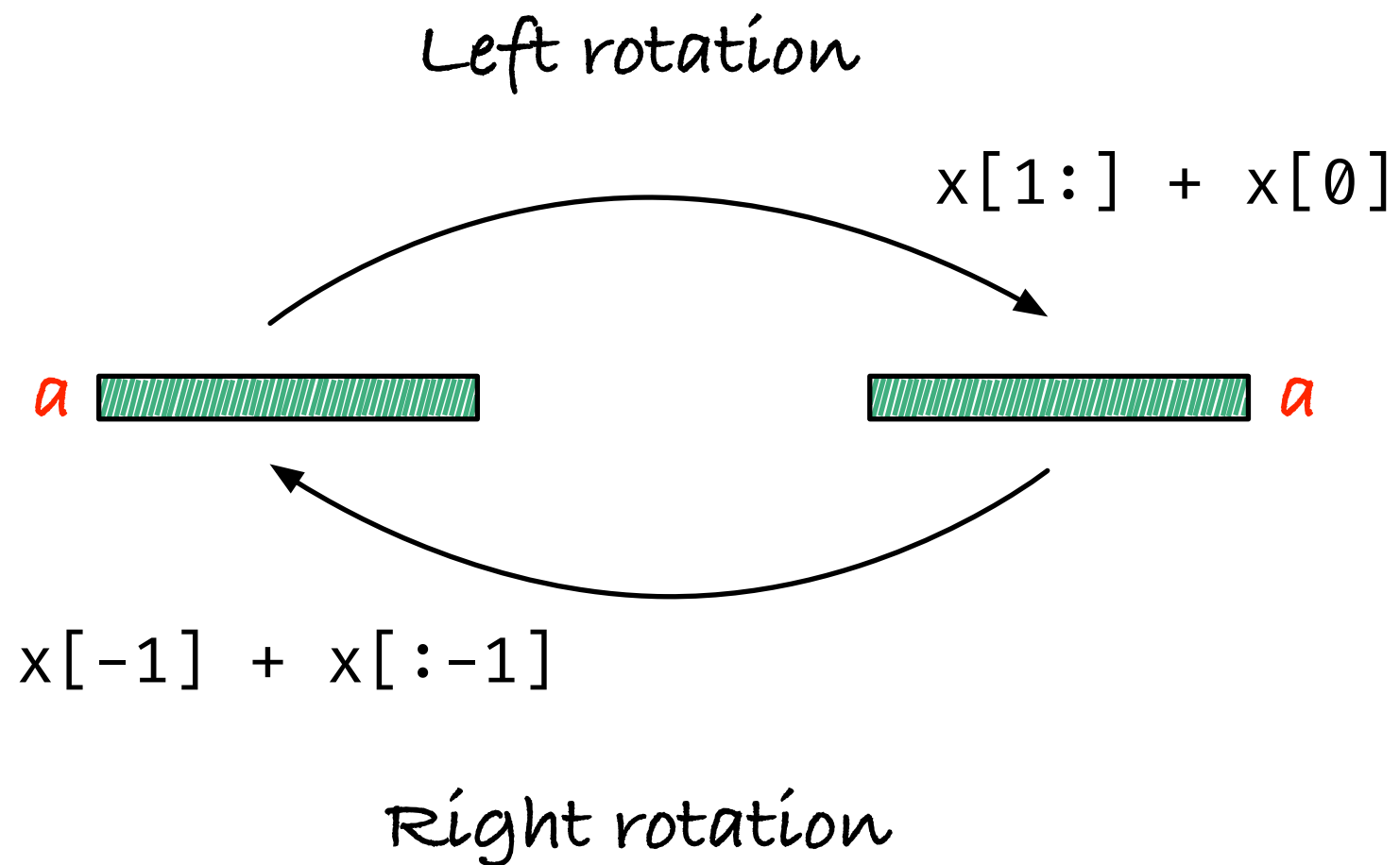


Burrows-Wheeler transform and FM-index search

(FM-index, after Paolo Ferragina and Giovanni
Manzini)

Burrows-Wheeler transform

String rotations



String rotations

-4: ppí\$missíssi

-3: p1\$mississip

-2: i\$mississippi

-1: \$mississippi

...síppí\$mississippi\$mississippi\$**mississippi**\$mississippi\$mississippi\$miss...

O: mississippi\$

+1: íssíssíppí\$m

+2: ssissippi\$mi

+3: síssíppí\$mis

+4: íssíppi\$miss

(infinite or just $x[i] == x[i \% n]$)

Burrows-Wheeler transform

$x = \text{mississippi}\$$

Take all rotations

0: mississippi\$	11: \$mississippi
1: ississippi\$m	10: i\$mississipp
2: ssissippi\$mi	7: ippi\$mississ
3: sissippi\$mis	4: issippi\$miss
4: issippi\$miss	1: ississippi\$m
5: ssippi\$missi	0: mississippi\$
6: sippi\$missis	9: pi\$mississip
7: ippi\$mississ	8: ppi\$mississi
8: ppi\$mississi	6: sippi\$missis
9: pi\$mississip	3: sissippi\$mis
10: i\$mississipp	5: ssippi\$missi
11: \$mississippi	2: ssissippi\$mi

Sort the rotations

Take the last column

$\text{bwt}(x) = \text{ipssm\$piissii}$

Burrows-Wheeler transform

$x = \text{mississippi\$}$

Take all rotations

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

Sort the rotations

Burrows-Wheeler matrix

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

Take the

bwt(x) is the last column in the Burrows-Wheeler matrix for x

$\text{bwt}(x) = \text{ipssm\$piissii}$

Burrows-Wheeler transform

- Burrows-Wheeler matrix is $O(n^2)$ space
 - Takes $O(n^2)$ time to fill in all rotations
 - Takes $O(n^2)$ time to radix sort it

Burrows-Wheeler transform

- Burrows-Wheeler matrix is $O(n^2)$ space
 - Takes $O(n^2)$ time to fill in all rotations
 - Takes $O(n^2)$ time to radix sort it
 - (We don't want any of this...)

BWT and suffix arrays

Sorted
rotations

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

sa(x)

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

BWT and suffix arrays

BWT and suffix arrays

$sa(x)$

$bwt(x)$

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

$sa(x)-1$ $bwt(x)$

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

Special case(?)

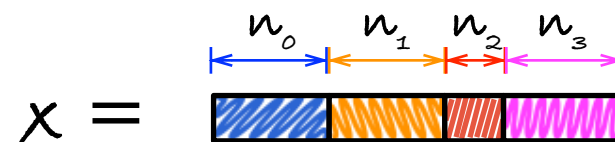
$bwt(x)[i] = \$$ if $sa[i] == 0$
 $bwt(x)[i] = x[sa[i]-1]$ otherwise

Why BWT?

- BWT is used in compression
 - e.g. in bzip2 (combined with many other things). It's the "b" in the name
- Compression likes low entropy (low randomness)
 - You get low randomness if some characters are much more frequent than others in a (sub-)string
 - We don't have time to talk compression and information theory, but I will give you a bit of an intuition (and nothing more)...

Run length encoding

- Given a string x , split it into blocks of consecutive identical characters, then encode each block as its length and one character.
- E.g. mississippi = (1,m),(1,i),(2,s),(1,i),(2,s),(1,i),(2,p),(1,i)



$$\text{RLE}(x) = [(n_0, \text{blue}), (n_1, \text{orange}), (n_2, \text{red}), (n_3, \text{purple})]$$

- This is exactly what we do with CIGAR strings

Run length encoding

- Run length encoding isn't always a good idea. Strings can get longer when you “compress them”.
- There is usually more magic in compression, but RLE is enough to illustrate why BWT is interesting
 - (and take my word that BWT is useful beyond RLE)

BWT tends to group characters

\$the_day_the_damned_dog_died
_damned_dog_died\$the_day_the
_day_the_damned_dog_died\$the
_died\$the_day_the_damned_dog
_dog_died\$the_day_the_damned
_the_damned_dog_died\$the_day
anned_dog_died\$the_day_the_d
ay_the_damned_dog_died\$the_d
d\$the_day_the_damned_dog_die
d_dog_died\$the_day_the_damne
damned_dog_died\$the_day_the_
day_the_damned_dog_died\$the_
died\$the_day_the_damned_dog_
dog_died\$the_day_the_damned_
e_damned_dog_died\$the_day_th
e_day_the_damned_dog_died\$th
ed\$the_day_the_damned_dog_di
ed_dog_died\$the_day_the_damn
g_died\$the_day_the_damned_do
he_damned_dog_died\$the_day_t
he_day_the_damned_dog_died\$t
ied\$the_day_the_damned_dog_d
mned_dog_died\$the_day_the_da
ned_dog_died\$the_day_the_dam
og_died\$the_day_the_damned_d
the_damned_dog_died\$the_day_
the_day_the_damned_dog_died\$
y_the_damned_dog_died\$the_da

$x = \text{"the_day_the_damned_dog_died\$"}$
 $\text{bwt}(x) = \text{"deegdyddee__hhinottdamd_\$a"}$

$\text{RLE}(x) = 1t1h1e1_1d1a1y1_1t1h1e1_1d1a1m1n1e1d1_1d1o1g1_1d1i1e1d1\$$
 $\text{RLE}(\text{bwt}(x)) = 1d2e1g1d1y2d2e4_2h1i1n1o2t1d1a1m1d1_1\$1a$

because BWT sorts by right-context

SA: $x[i]$ sorted wrt $x[i:]$

a 

BWT: $x[i]$ sorted wrt $x[i+1:] + x[i]$

 *a*

$x = \text{"the_day_the_damned_dog_died\$"}$

day_the_damned_dog_died\$

damned_dog_died\$

dog_died\$

died\$

...

d_dog_died\$the_day_the_damne
damned_dog_died\$the_day_the_
day_the_damned_dog_died\$the_
died\$the_day_the_damned_dog_
dog_died\$the_day_the_damned_
e_damned_dog_died\$the_day_th

...

...

_the_damned_dog_died\$the_day
anned_dog_died\$the_day_the_d
ay_the_damned_dog_died\$the_d
d\$the_day_the_damned_dog_die
d_dog_died\$the_day_the_damne

...

he_day_the_damned_dog_died\$t
ied\$the_day_the_damned_dog_d
mned_dog_died\$the_day_the_da
ned_dog_died\$the_day_the_dam
og_died\$the_day_the_damned_d
the_damned_dog_died\$the_day

Why not use the first column?

\$the_day_the_damned_dog_died
 _damned_dog_died\$the_day_the
 _day_the_damned_dog_died\$the
 _died\$the_day_the_damned_dog
 _dog_died\$the_day_the_damned
 _the_damned_dog_died\$the_day
 anned_dog_died\$the_day_the_d
 ay_the_damned_dog_died\$the_d
 d\$the_day_the_damned_dog_die
 d_dog_died\$the_day_the_damne
 damned_dog_died\$the_day_the_
 day_the_damned_dog_died\$the_
 died\$the_day_the_damned_dog_
 dog_died\$the_day_the_damned_
 e_damned_dog_died\$the_day_th
 e_day_the_damned_dog_died\$th
 ed\$the_day_the_damned_dog_di
 ed_dog_died\$the_day_the_damn
 g_died\$the_day_the_damned_do
 he_damned_dog_died\$the_day_t
 he_day_the_damned_dog_died\$
 died\$the_day_the_damned_dog_d
 mned_dog_died\$the_day_the_da
 ned_dog_died\$the_day_the_dam
 og_died\$the_day_the_damned_d
 the_damned_dog_died\$the_day_
 the_day_the_damned_dog_died\$
 y_the_damned_dog_died\$the_da

$x = \text{"the_day_the_damned_dog_died\$"}$

$\text{bwt}(x) = \text{"deegdyddee__hhinottdama_\$a"}$

$\text{RLE}(x) = 1\text{t}1\text{h}1\text{e}1_1\text{d}1\text{a}1\text{y}1_1\text{t}1\text{h}1\text{e}1_1\text{d}1\text{a}1\text{m}1\text{n}1\text{e}1\text{d}1_1\text{d}1\text{o}1\text{g}1_1\text{d}1\text{i}1\text{e}1\text{d}1\$$

$\text{RLE}(\text{bwt}(x)) = 1\text{d}2\text{e}1\text{g}1\text{d}1\text{y}2\text{d}2\text{e}4_2\text{h}1\text{i}1\text{n}1\text{o}2\text{t}1\text{d}1\text{a}1\text{m}1\text{d}1_1\1a

$\text{RLE}(\text{1stCol}(x)) = 1\$5_2\text{a}6\text{d}4\text{e}1\text{g}2\text{h}1\text{i}1\text{m}1\text{n}1\text{o}2\text{t}1\text{y}$

Two things we want from compression:

- We want $\text{compr}(x)$ to be small
- We *also* want to get x back: $x = \text{inv}(\text{compr}(x))$
- The first column in the BWT matrix is optimal for RLE
- We cannot reconstruct x from it
- We can from the last column

Characters in first and last come in the same order...

(when identifying the characters by their position in the original string x\$)

$\$_0$ mississipp*i*₀
*i*₀ $\$$ mississip*p*₀
*i*₁ppi\$missis*s*₀
*i*₂ssippi\$mis*s*₁
*i*₃ssissippi\$m*m*₀
*m*₀ississippi\$*s*₀
*p*₀i\$mississip*p*₁
*p*₁pi\$mississ*i*₁
*s*₀ippi\$missi*s*₂
*s*₁issippi\$mi*s*₃
*s*₂sippi\$miss*i*₂
*s*₃sissippi\$m*i*₃

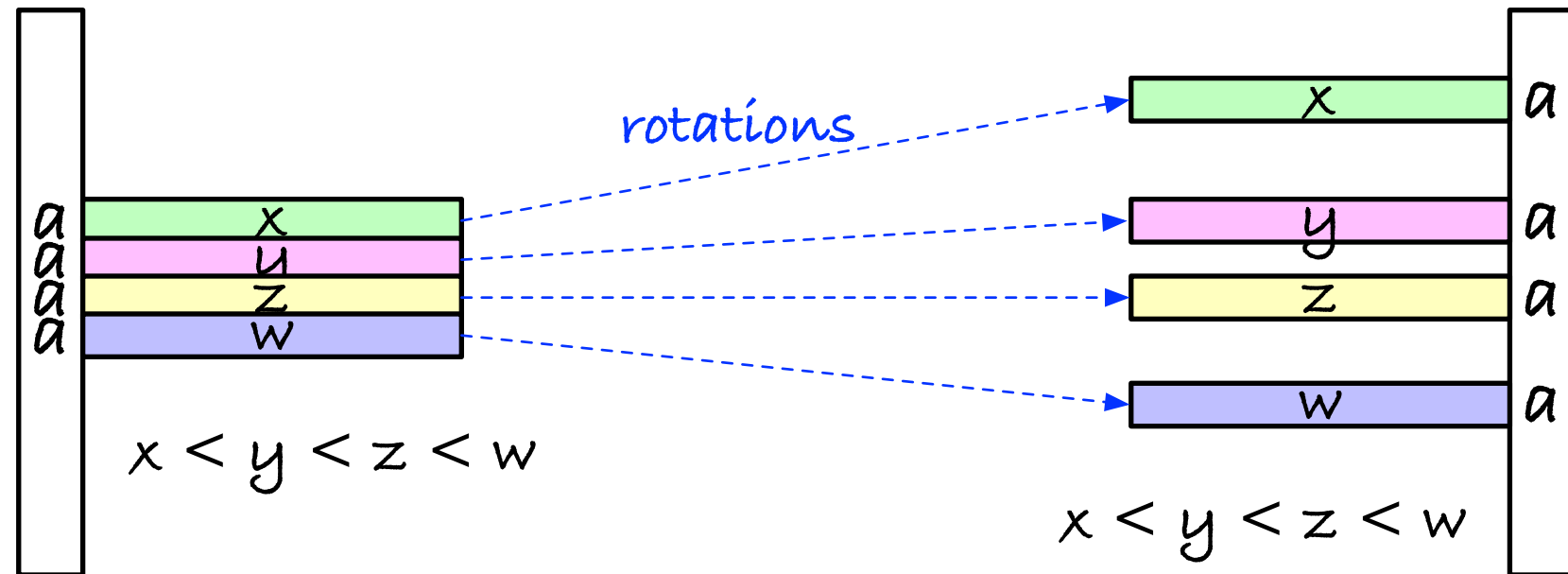
(when identifying the characters by their position in the original string x\$)



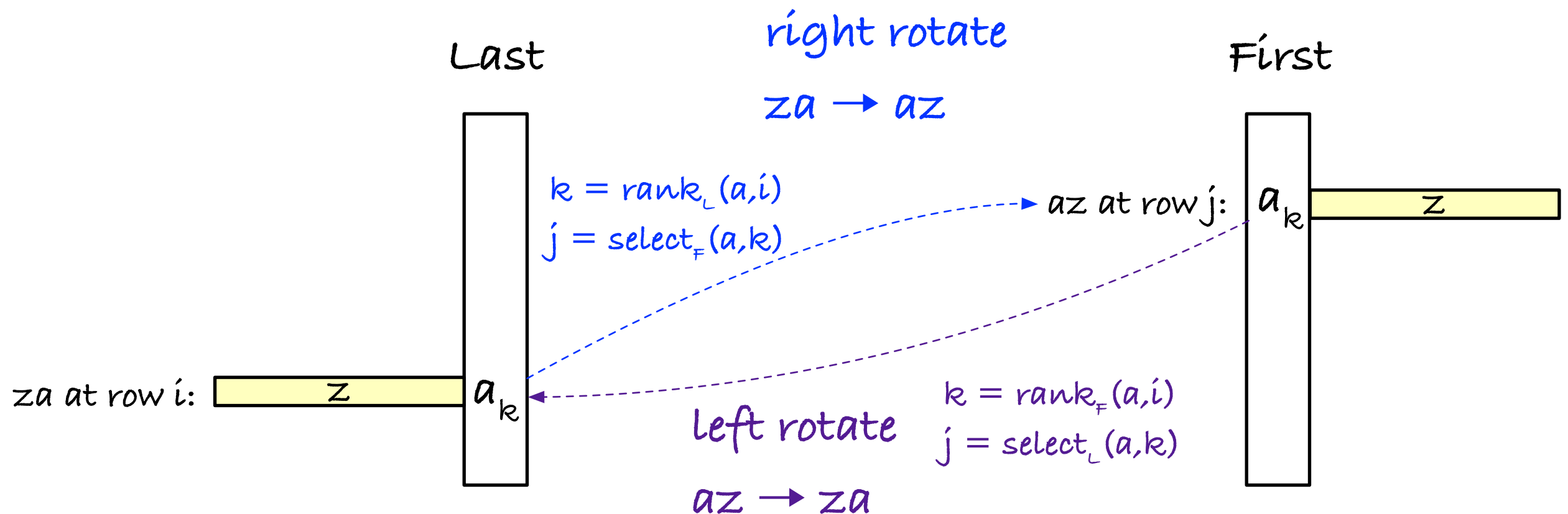
Characters come in the same order

a-bin in first column of BWT matrix

Placement in last column given by order of suffixes of the first column...



Rotating in the BWT matrix



These are common algorithmic operations. Here on strings (the columns are strings). We need to implement them somehow.

We can rotate using only first and last (no middle columns needed)

We can always rebuild first from last (we only need to count characters)

$$\text{rank}_A(a, i) = |\{j < i : A[j] == a\}|$$

$$\text{select}_A(a, k) = \text{pos of the } k\text{'th } a \text{ in } A$$

Reversing transformation

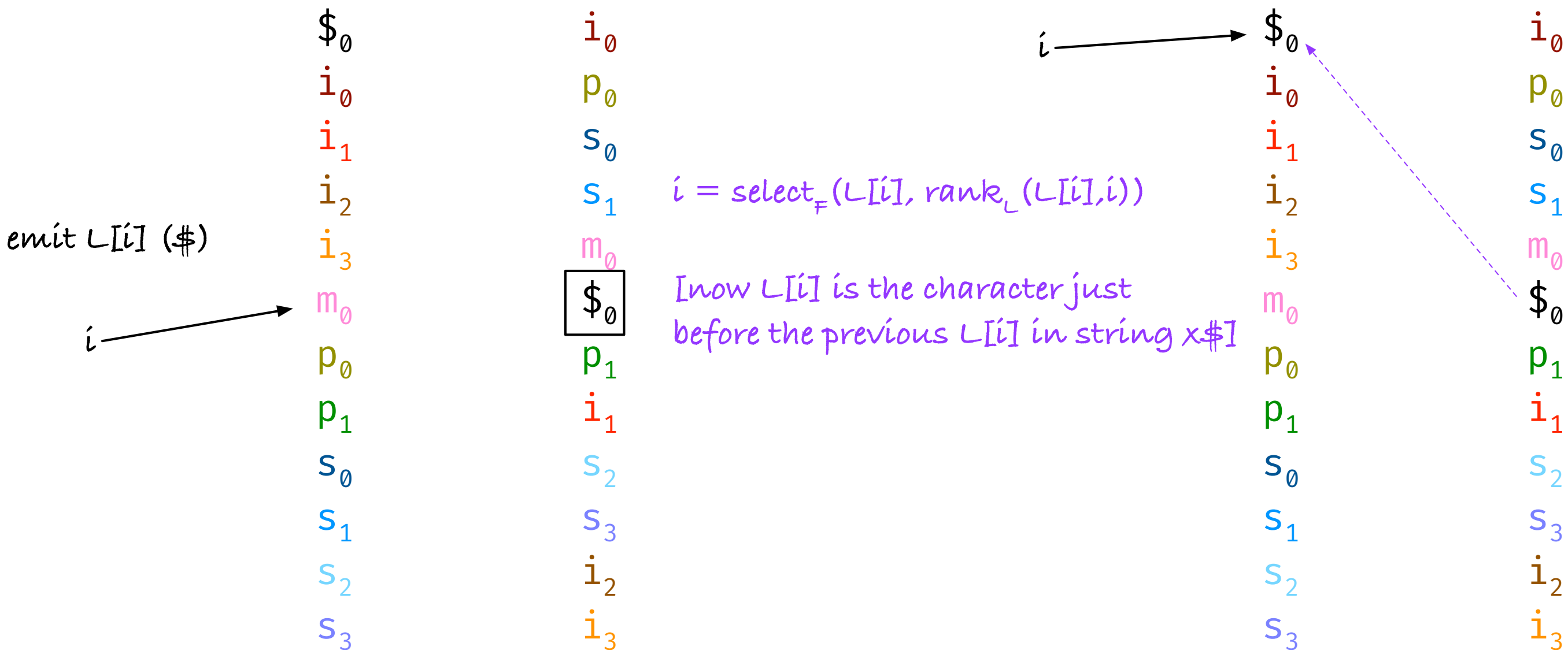
Input $L = \text{bwt}(\text{mississippi}\$) = \text{ipssm\$pissii}$

Preprocess as needed for rank and select...

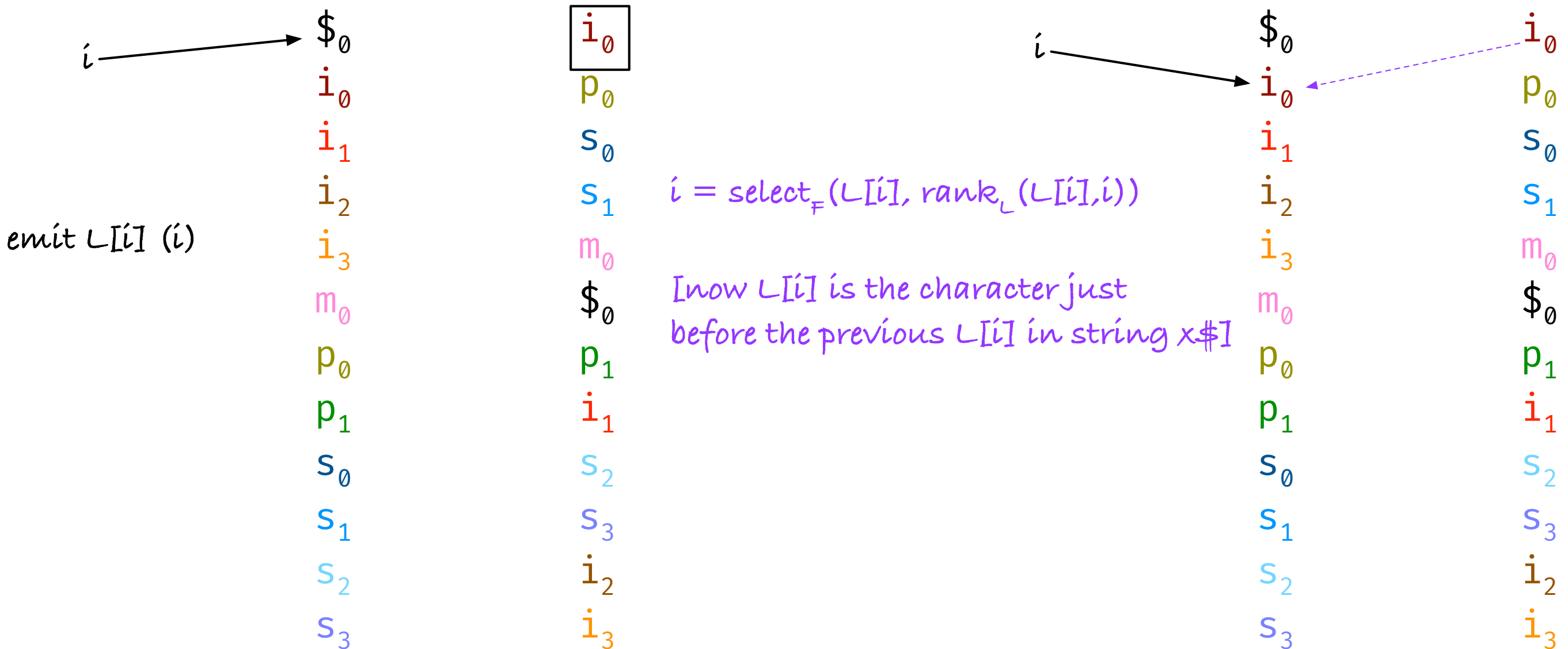
$i = \text{select}_L(\$, 0)$

$\$$ ₀mississipp*i*₀
*i*₀\$mississipp*p*₀
*i*₁ppi\$mississ*s*₀
*i*₂ssippi\$miss*s*₁
*i*₃ssissippi\$m*0*
i → *m*₀ississippi\$*0*
*p*₀i\$mississip*p*₁
*p*₁pi\$mississ*i*₁
*s*₀ippi\$missi*s*₂
*s*₁issippi\$mi*s*₃
*s*₂sippi\$miss*i*₂
*s*₃sissippi\$m*i*₃

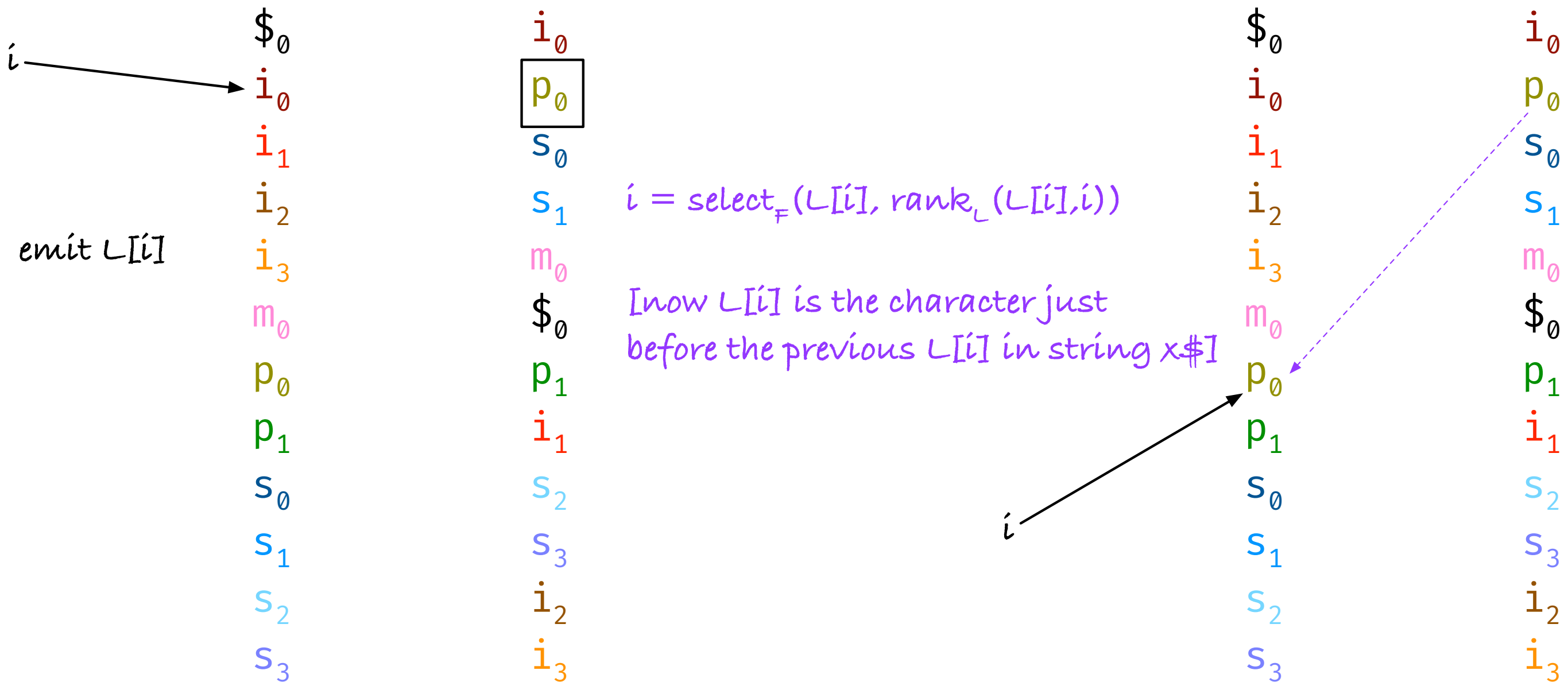
Reversing transformation



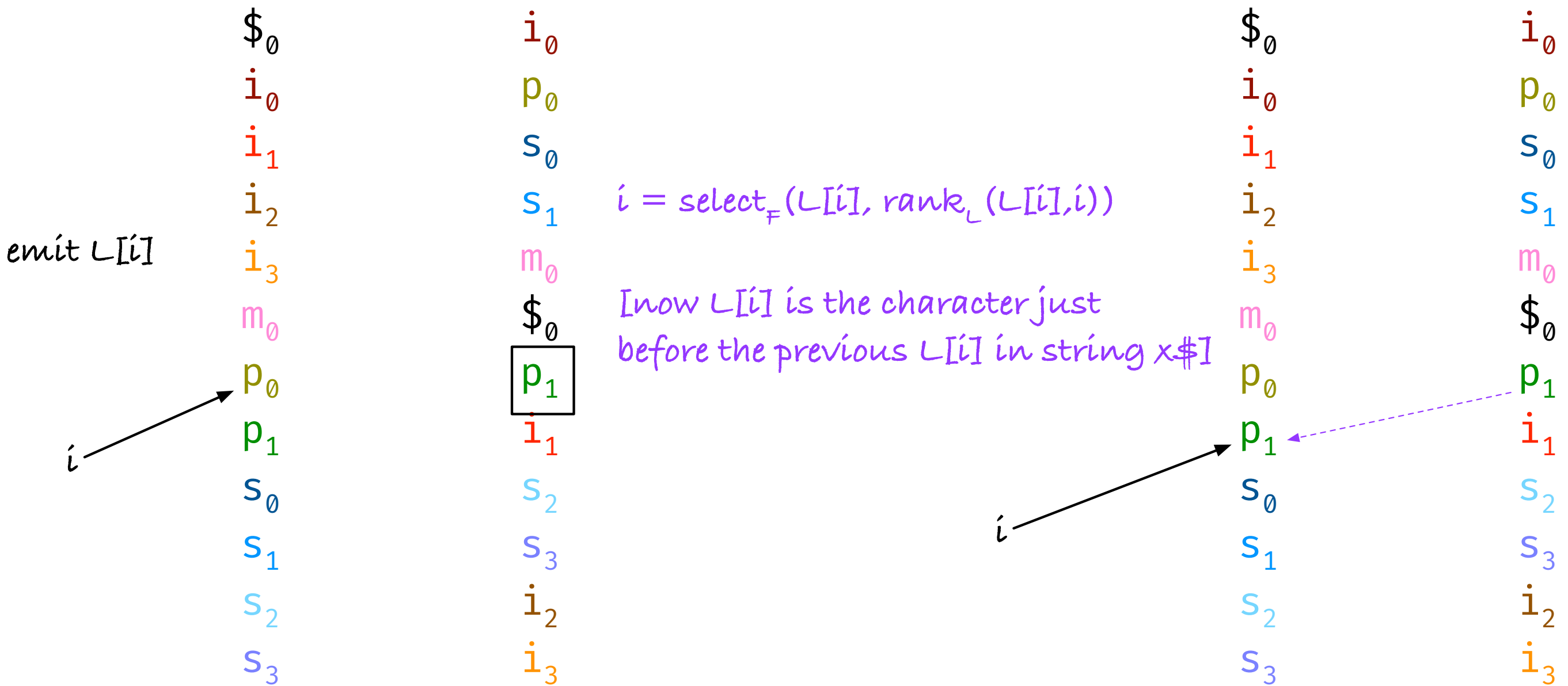
Reversing transformation



Reversing transformation



Reversing transformation



Emitted: $ppi\#$

Reversal as (implicit) rotations

We don't `x` or its rotations,
but we can see the last character
for all of them...

```
row select_($,0)
```

```
rotate right  
rotate right  
rotate right  
rotate right  
rotate right  
rotate right  
rotate right  
rotate right  
rotate right  
rotate right  
rotate right
```

↓
L[i]

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

You start with the rotation that ends in \$. Then you keep rotating to get the letter to the left of the current in `x`, until you have seen `n` characters. Then you have seen `x` reversed (so reverse the result or store it backwards).

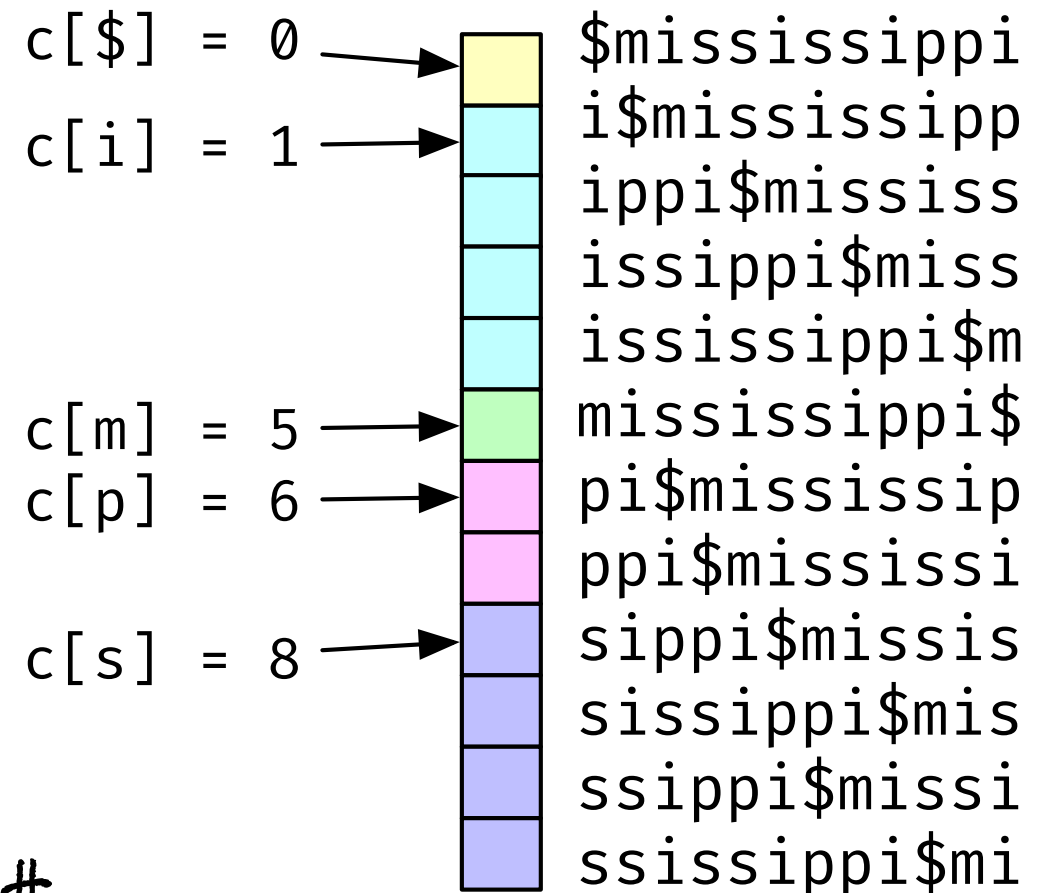
Running time

- You do $\text{rank}_L(L[i], i)$ and $\text{select}_F(L[i], k)$ n times.
- We shall see that we can do both operations in $O(1)$.
- The running time for reversal: $\text{bwt}(x) \rightarrow x$ is $O(n)$.

Select in first column

$c[a] = \text{index of bin } a$

Just bucketing as we have
seen a thousand times before.



$x = \text{mississippi\$}$

$\text{counts} = \{\$: 1, i: 4, m: 1, p: 2, s: 4\}$

$\text{buckets} = \{\$: 0, i: 1, m: 5, p: 6, s: 8\}$

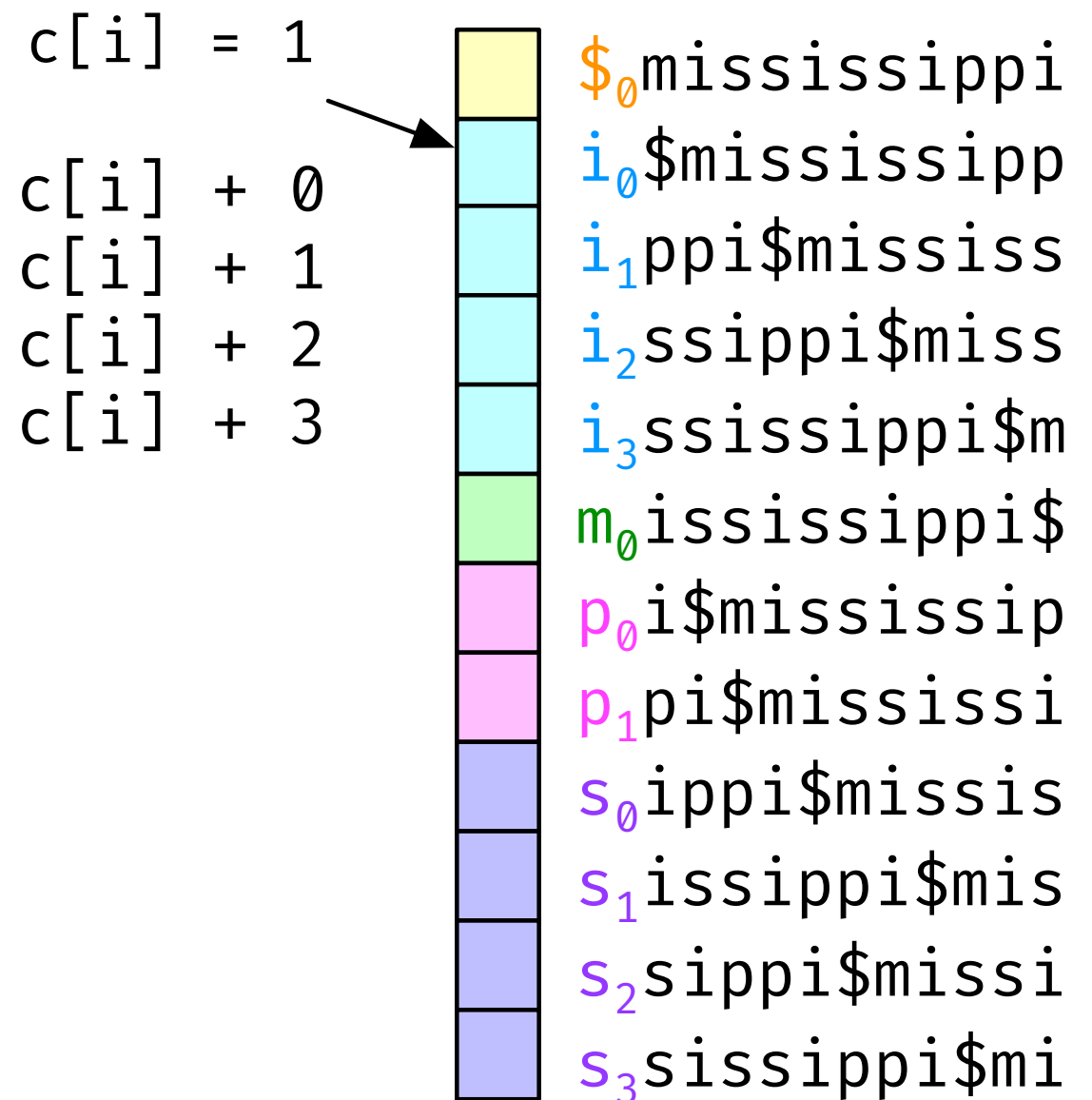
Notice: we don't need the first column for this. We can count in the last, that we already have.

Select in first column

$\text{select}_F(a, k) = \text{index of } k\text{'th } a$
 $= \text{a number } k \text{ in the bucket}$

$$\text{select}_F(a, k) = C[a] + k$$

We can build the C table in $O(n)$
 and then we can select in $O(1)$



Rank in last column

O table just brute force rank...

$$O[0,:] = 0$$

$$O[i,a] = O[i-1,a] + (L[i-1] == a)$$

$$\text{rank}_L(a,k) = O[k,a]$$

We can build the O table in $O(n)$
 (because σ is $O(1)$)
 and then we can select in $O(1)$

(add an extra row, $O[n+1,:]$, it pays off later) →

	$O[\$, i, m, p, s]$
$\$ _0 \text{mississippi} i_0$	0, 0, 0, 0, 0
$i_0 \$ \text{mississippi} p_0$	0, 1, 0, 0, 0
$i_1 p p i \$ \text{mississ} s_0$	0, 1, 0, 1, 0
$i_2 s s i p p i \$ \text{miss} s_1$	0, 1, 0, 1, 1
$i_3 s s i s s i p p i \$ m_0$	0, 1, 0, 1, 2
$m_0 i s s i s s i p p i \$$	0, 1, 1, 1, 2
$p_0 i \$ \text{missississ} i_1$	1, 1, 1, 1, 2
$p_1 p i \$ \text{missississ} i_1$	1, 1, 1, 2, 2
$s_0 i p p i \$ \text{mississ} s_2$	1, 2, 1, 2, 2
$s_1 i s s i p p i \$ \text{miss} s_3$	1, 2, 1, 2, 3
$s_2 s i p p i \$ \text{mississ} i_2$	1, 2, 1, 2, 4
$s_3 s i s s i p p i \$ m i_3$	1, 3, 1, 2, 4
	1, 4, 1, 2, 4

The O table is $O(\sigma n)$ integers (or $O(\sigma n \log n)$ bits). We call that $O(n)$ with constant alphabet and the RAM model.

We can fill the O table in $O(\sigma n)$ by running line by line or column by column.

Rank in last column

$O[i, a]$ is the number of a 's
in $\text{bwt}(x)$ before index i

(which is what $\text{rank}_{\text{bwt}(x)}(a, i)$ is)

$$O[7, \$] = 1$$

	$O[\$, i, m, p, s]$
0:	0, 0, 0, 0, 0
1:	0, 1, 0, 0, 0
2:	0, 1, 0, 1, 0
3:	0, 1, 0, 1, 1
4:	0, 1, 0, 1, 2
5:	0, 1, 1, 1, 2
6:	1, 1, 1, 1, 2
7:	1, 1, 1, 2, 2
8:	1, 2, 1, 2, 2
9:	1, 2, 1, 2, 3
10:	1, 2, 1, 2, 4
11:	1, 3, 1, 2, 4
12:	1, 4, 1, 2, 4

$\text{bwt}(x)$

i
 p
 s
 s
 m
 $\$$
 p
 i ← 7
 s
 s
 i
 i

Rank in last column

$O[i, a]$ is the number of a 's
in $\text{bwt}(x)$ before index i

(which is what $\text{rank}_{\text{bwt}(x)}(a, i)$ is)

$$O[7, i] = 1$$

	O[\$, i, m, p, s]
0:	0, 0, 0, 0, 0
1:	0, 1, 0, 0, 0
2:	0, 1, 0, 1, 0
3:	0, 1, 0, 1, 1
4:	0, 1, 0, 1, 2
5:	0, 1, 1, 1, 2
6:	1, 1, 1, 1, 2
7:	1, 1, 1, 2, 2
8:	1, 2, 1, 2, 2
9:	1, 2, 1, 2, 3
10:	1, 2, 1, 2, 4
11:	1, 3, 1, 2, 4
12:	1, 4, 1, 2, 4

$\text{bwt}(x)$

i

p

s

s

m

\$

p

i

s

s

i

i

7

(we don't count
the i at index 7)

Rank in last column

$O[i, a]$ is the number of a 's
in $\text{bwt}(x)$ before index i

(which is what $\text{rank}_{\text{bwt}(x)}(a, i)$ is)

$$O[7, m] = 1$$

	$O[\$, i, m, p, s]$
0:	0, 0, 0, 0, 0
1:	0, 1, 0, 0, 0
2:	0, 1, 0, 1, 0
3:	0, 1, 0, 1, 1
4:	0, 1, 0, 1, 2
5:	0, 1, 1, 1, 2
6:	1, 1, 1, 1, 2
7:	1, 1, 1, 2, 2
8:	1, 2, 1, 2, 2
9:	1, 2, 1, 2, 3
10:	1, 2, 1, 2, 4
11:	1, 3, 1, 2, 4
12:	1, 4, 1, 2, 4

$\text{bwt}(x)$

i
 p
 s
 s
 m
 $\$$
 p
 i
 s
 s
 i
 i

7

Rank in last column

$O[i, a]$ is the number of a 's
in $\text{bwt}(x)$ before index i

(which is what $\text{rank}_{\text{bwt}(x)}(a, i)$ is)

$$O[7, p] = 2$$

	$O[\$, i, m, p, s]$
0:	0, 0, 0, 0, 0
1:	0, 1, 0, 0, 0
2:	0, 1, 0, 1, 0
3:	0, 1, 0, 1, 1
4:	0, 1, 0, 1, 2
5:	0, 1, 1, 1, 2
6:	1, 1, 1, 1, 2
7:	1, 1, 1, 2, 2
8:	1, 2, 1, 2, 2
9:	1, 2, 1, 2, 3
10:	1, 2, 1, 2, 4
11:	1, 3, 1, 2, 4
12:	1, 4, 1, 2, 4

$\text{bwt}(x)$

i
p
 s
 s
 m
 $\$$
p
 i
 s
 s
 i
 i

7

Rank in last column

$O[i, a]$ is the number of a 's
in $\text{bwt}(x)$ before index i

(which is what $\text{rank}_{\text{bwt}(x)}(a, i)$ is)

$$O[7, s] = 2$$

	$O[\$, i, m, p, s]$
0:	0, 0, 0, 0, 0
1:	0, 1, 0, 0, 0
2:	0, 1, 0, 1, 0
3:	0, 1, 0, 1, 1
4:	0, 1, 0, 1, 2
5:	0, 1, 1, 1, 2
6:	1, 1, 1, 1, 2
7:	1, 1, 1, 2, 2
8:	1, 2, 1, 2, 2
9:	1, 2, 1, 2, 3
10:	1, 2, 1, 2, 4
11:	1, 3, 1, 2, 4
12:	1, 4, 1, 2, 4

$\text{bwt}(x)$

i
 p
 s
 s
 m
 $\$$
 p
 i
 s
 s
 i
 i

7

Rank in last column

bwt(x)

i
p
s
s
m
\$
p
i
s
s
i
i

O[\$, i, m, p, s]

0:	0	0	0	0	0
1:	0	1	0	0	0
2:	0	1	0	1	0
3:	0	1	0	1	1
4:	0	1	0	1	2
5:	0	1	1	1	2
6:	1	1	1	1	2
7:	1	1	1	2	2
8:	1	2	1	2	2
9:	1	2	1	2	3
10:	1	2	1	2	4
11:	1	3	1	2	4
12:	1	4	1	2	4

$O[i, a]$ is the number of a 's
in $\text{bwt}(x)$ before index i

(which is what $\text{rank}_{\text{bwt}(x)}(a, i)$ is)

$O[16, :] == \text{character counts in } x$

12

Rank in last column

You can easily represent row zero implicitly

$$O[i,a] = 0 \text{ if } i = 0 \text{ else } O'[i-1,a]$$

In the LM index search you never look up in column \$, so you don't need it.

(It is still a big table, though...)

	O[\$, i, m, p, s]				
0:	0	0	0	0	0
1:	0	1	0	0	0
2:	0	1	0	1	0
3:	0	1	0	1	1
4:	0	1	0	1	2
5:	0	1	1	1	2
6:	1	1	1	1	2
7:	1	1	1	2	2
8:	1	2	1	2	2
9:	1	2	1	2	3
10:	1	2	1	2	4
11:	1	3	1	2	4
12:	1	4	1	2	4

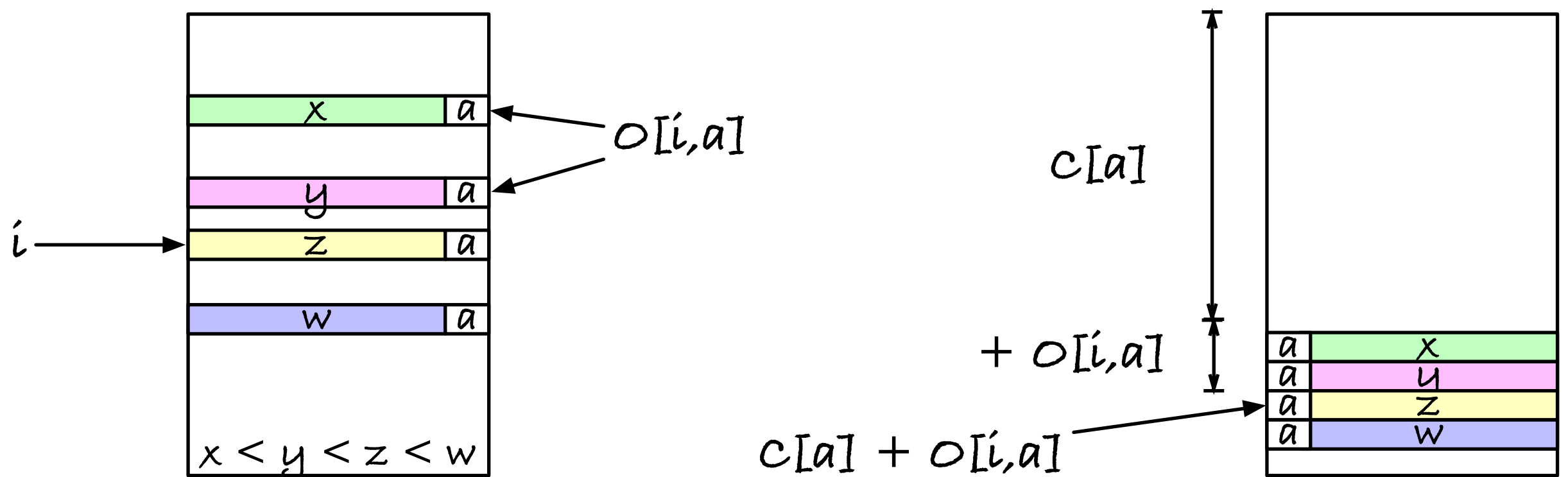
We can talk about ways to reducing the size in one of the mystery lectures at the end of the semester.

FM index

Ferragina and Manzini index

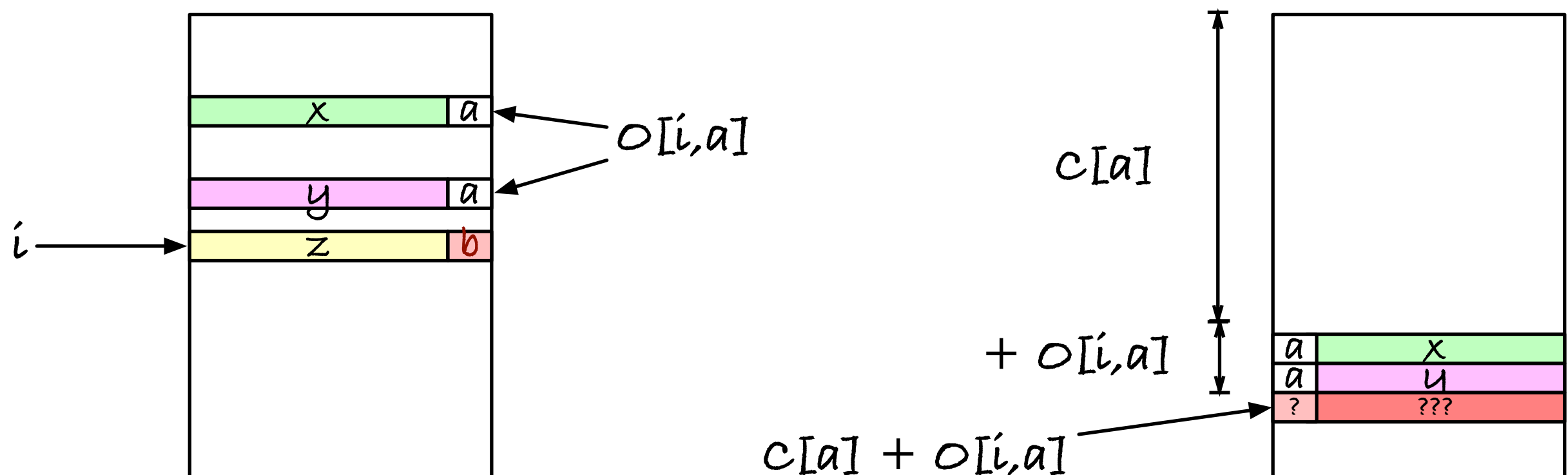
General jump

Rotating the last character to the front...



General jump

Rotating the last character to the front...



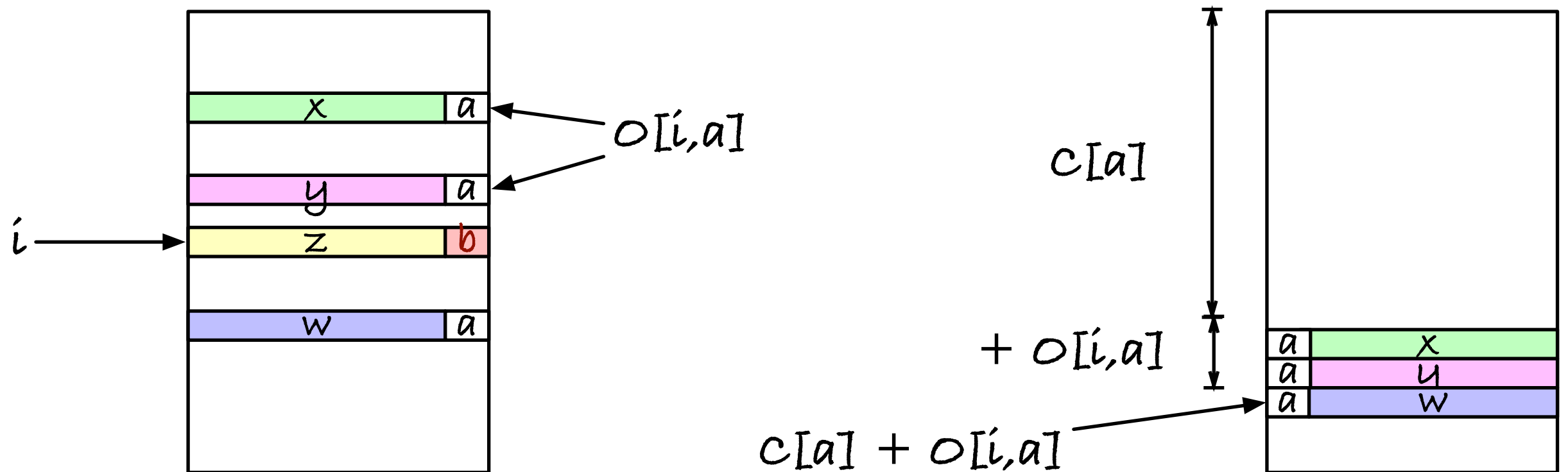
?	???
---	-----

 is the smallest rotation \geq

a	z
---	---

General jump

Rotating the last character to the front...

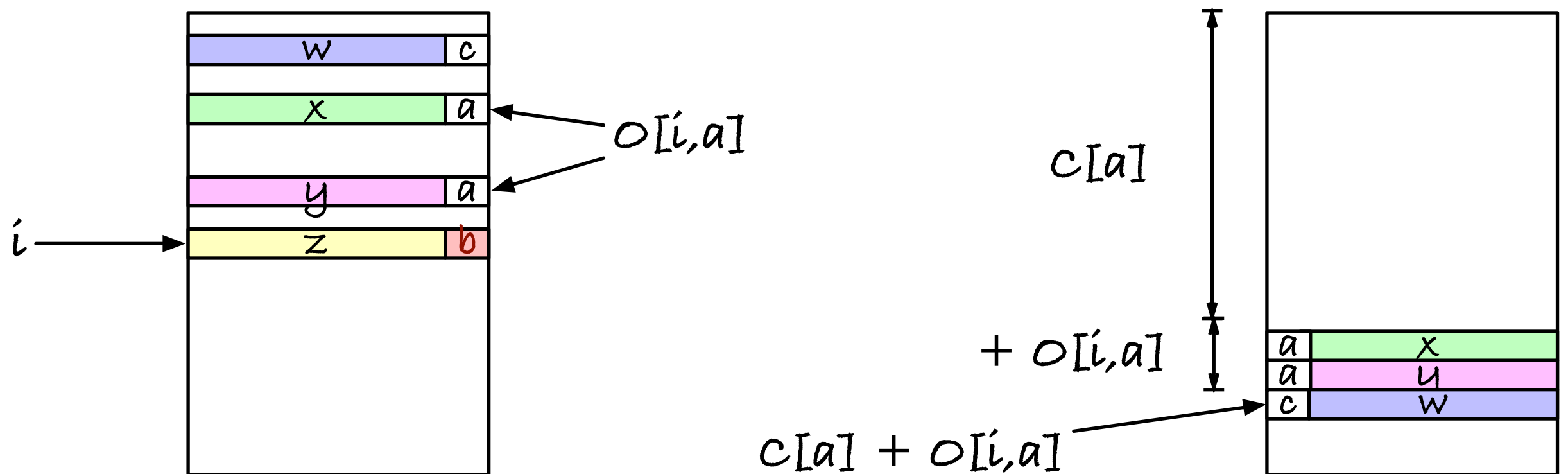


Assume $\boxed{?} \boxed{???} == \boxed{a} \boxed{w}$

Then $w \geq z$ (or it would go above $c[a] + o[i, a]$), and it is the smallest aw since the a -bin is sorted.

General jump

Rotating the last character to the front...



Assume

?	???
---	-----

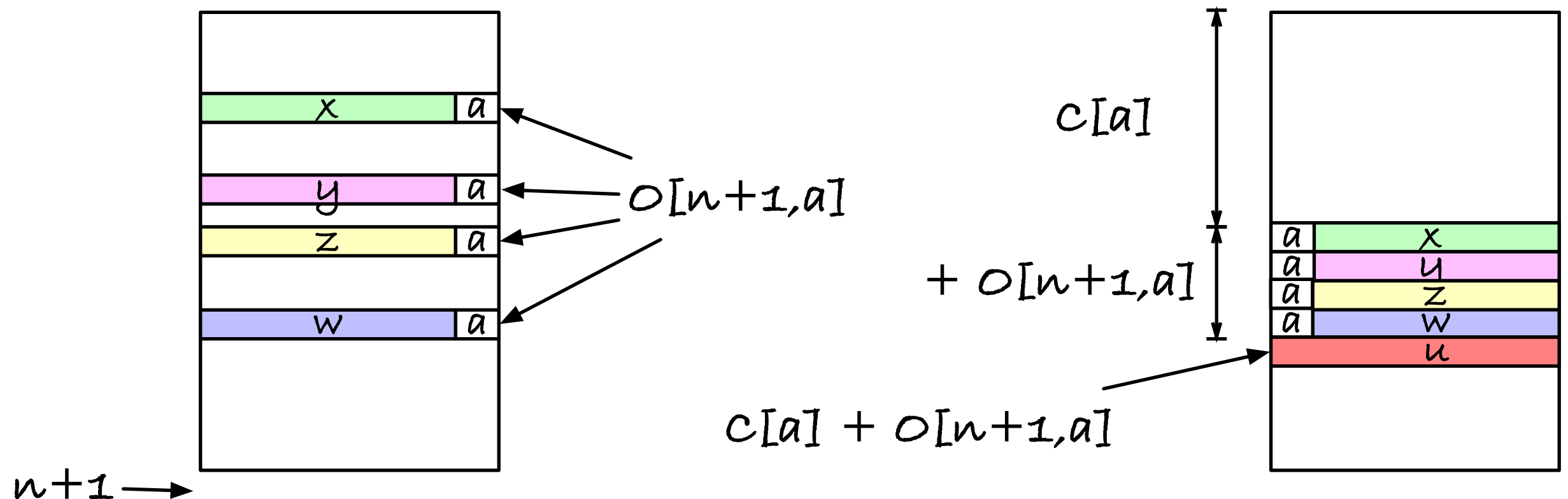
 ==

c	w
---	---

Then cw is the smallest rotation after those that start with a , and since all rotations that start with a , ax , are smaller than az (or we wouldn't point beyond the a -bin), it is the smallest greater than az .

General jump

Special case for row $n - 1$



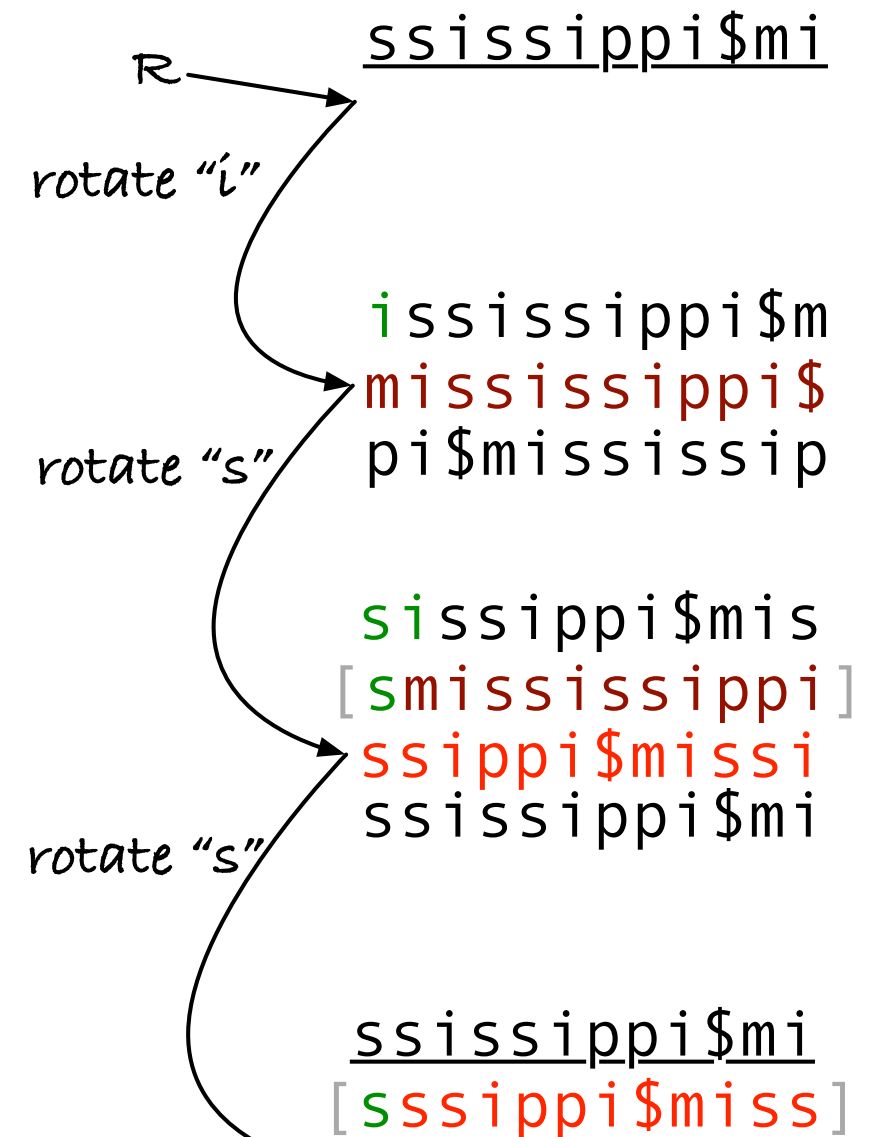
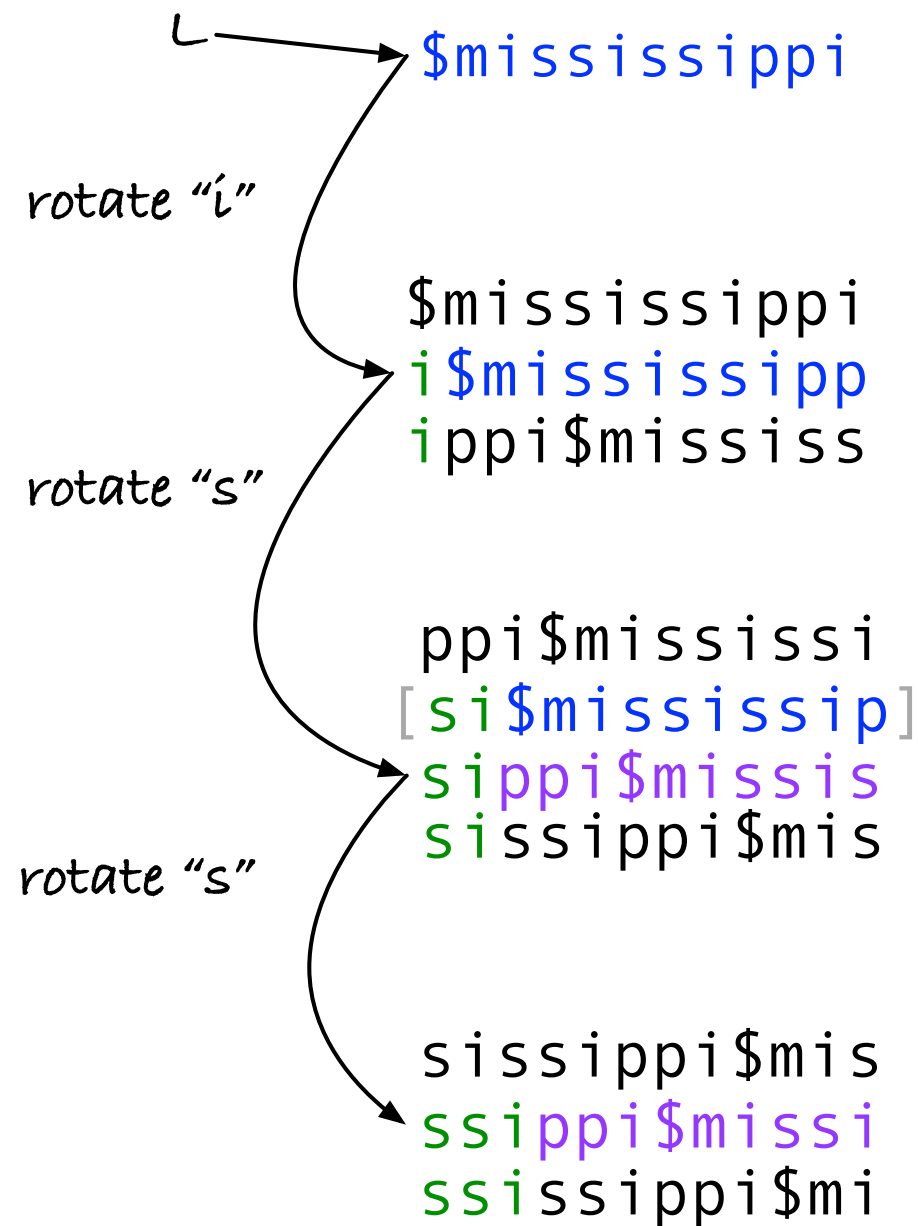
u is the smallest string after the a -bin

FM index search

Start with the smallest and one-past-the-largest rotations.
 "Rotate" according to reverse p.

```
L, R = 0, n+1
for a in reversed(p):
    if L == R: break
    L = C[a] + O[L,a]
    R = C[a] + O[R,a]
```

$x = \text{"mississippi"}$
 $p = \text{"ssi"}$

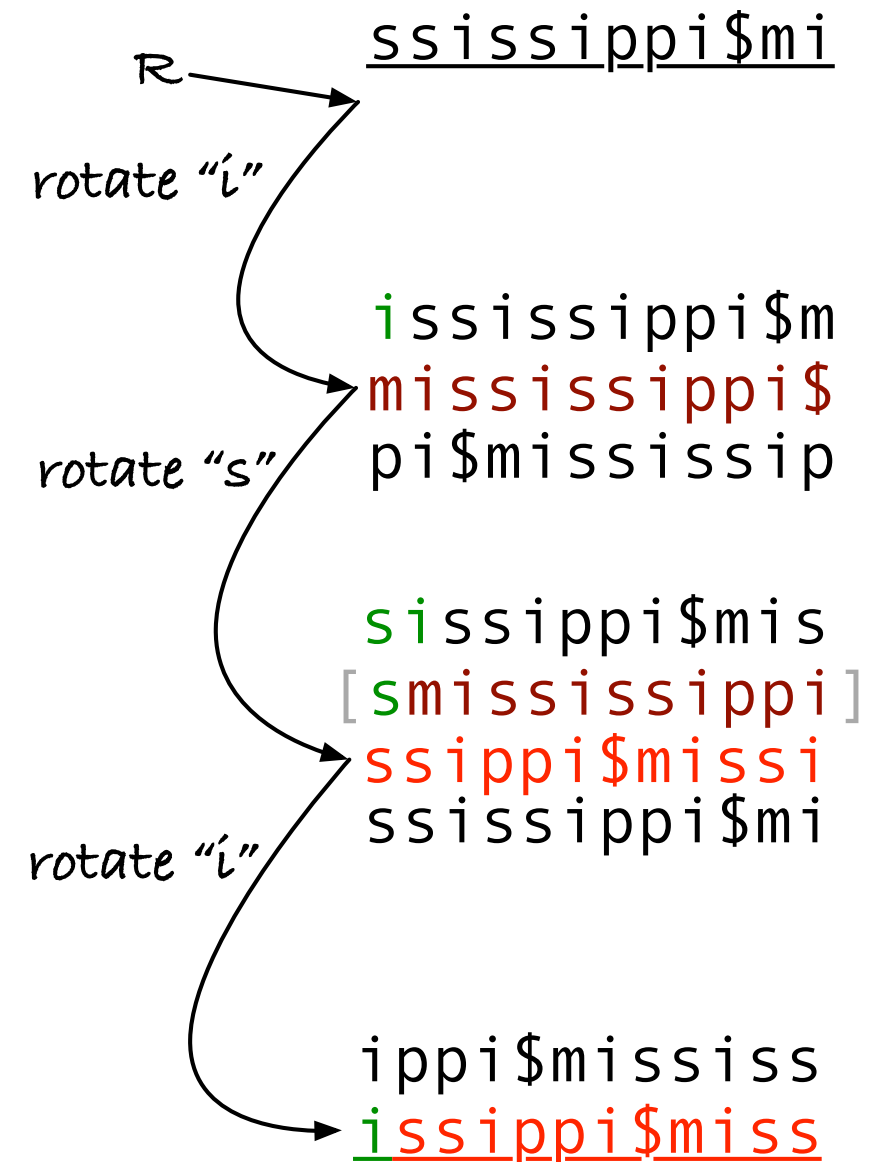
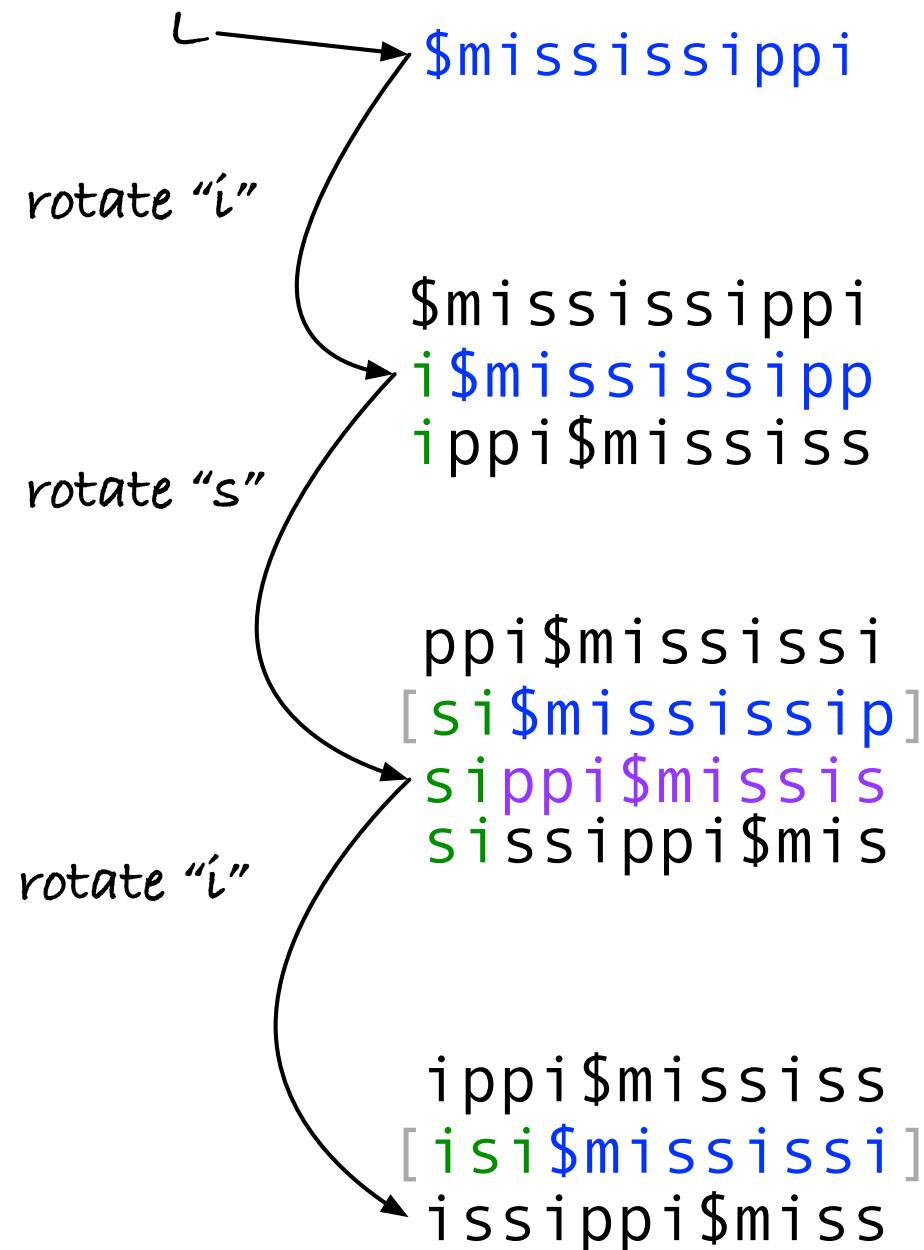


FM index search

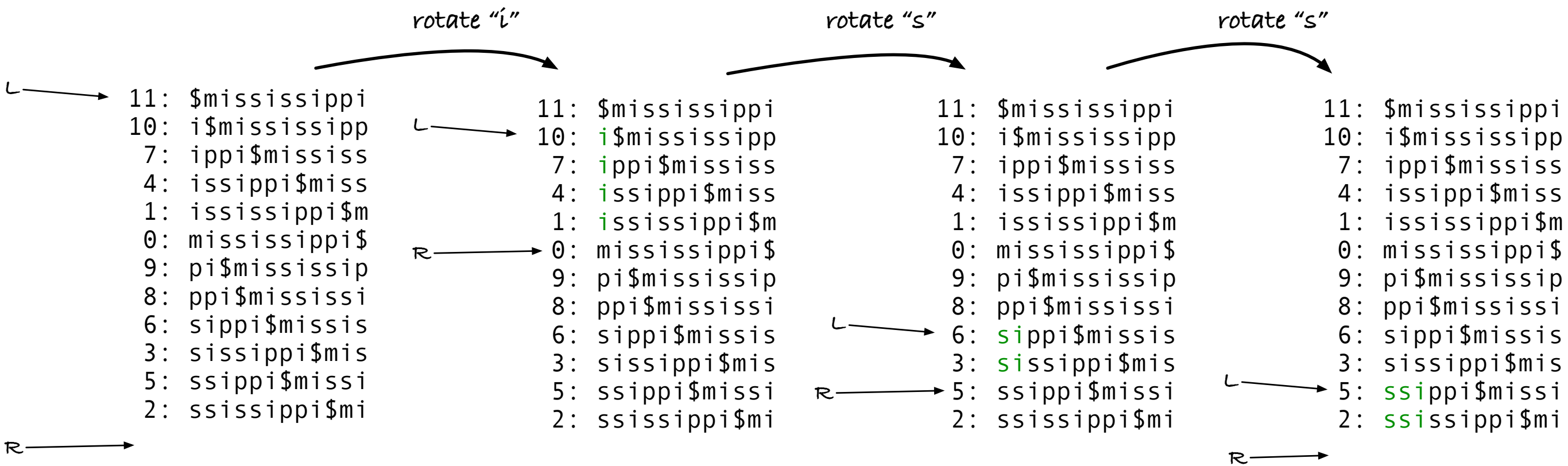
Start with the smallest and one-past-the-largest rotations.
 "Rotate" according to reverse p.

```
L, R = 0, n+1
for a in reversed(p):
    if L == R: break
    L = C[a] + O[L,a]
    R = C[a] + O[R,a]
```

x = "mississippi"
 p = "isi"



FM index search



Matches at $sa[i]$ for i in range(L,R)

Running time

- Preprocessing: $O(n)$
 - Suffix array [to get bwt(x)]
 - C and O tables
- Search: $O(m)$



That's all Folks!