

BWT for repetitive texts, part 3: Toehold lemma setup

Ben Langmead



Please sign guestbook (www.langmead-lab.org/teaching-materials) to tell me briefly how you are using the slides. For original Keynote files, email me (ben.langmead@gmail.com).

Locate queries

	Count		Locate	
	Space	Time	Space	Time
FM Index (2000)	$O(n)$	$O(m)$	$O(n)$	$O(m + \text{occ})$
RLFM Index (2005)	$O(r)$	$O(m)$	$O(n)$	$O(m + \text{occ})$
r-index (2018)	$O(r)$	$O(m)$	$O(r)$	$O(m + \text{occ})$

n = reference length, m = query length,

r = # BWT runs

(bounds simplified)

Locate queries

How?

	Count		Locate	
	Space	Time	Space	Time
FM Index (2000)	$O(n)$	$O(m)$	$O(n)$	$O(m + \text{OCC})$
RLFM Index (2005)	$O(r)$	$O(m)$	$O(n)$	$O(m + \text{OCC})$
r-index (2018)	$O(r)$	$O(m)$	$O(r)$	$O(m + \text{OCC})$

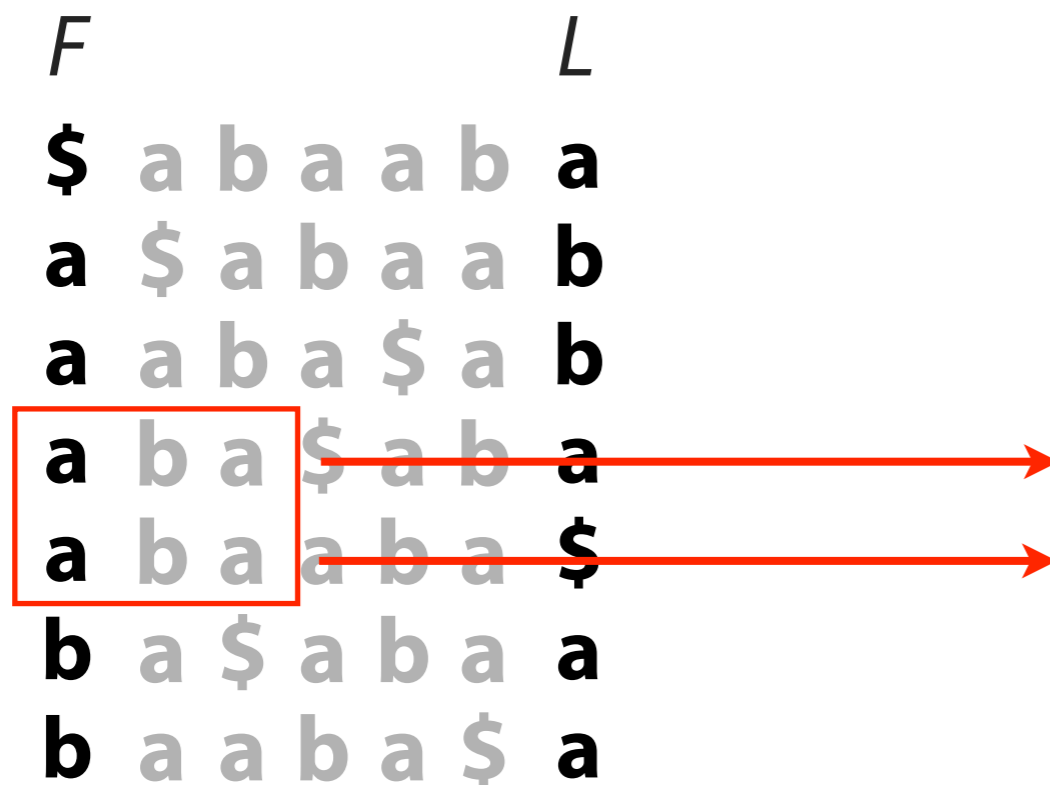
n = reference length, m = query length,

r = # BWT runs

(bounds simplified)

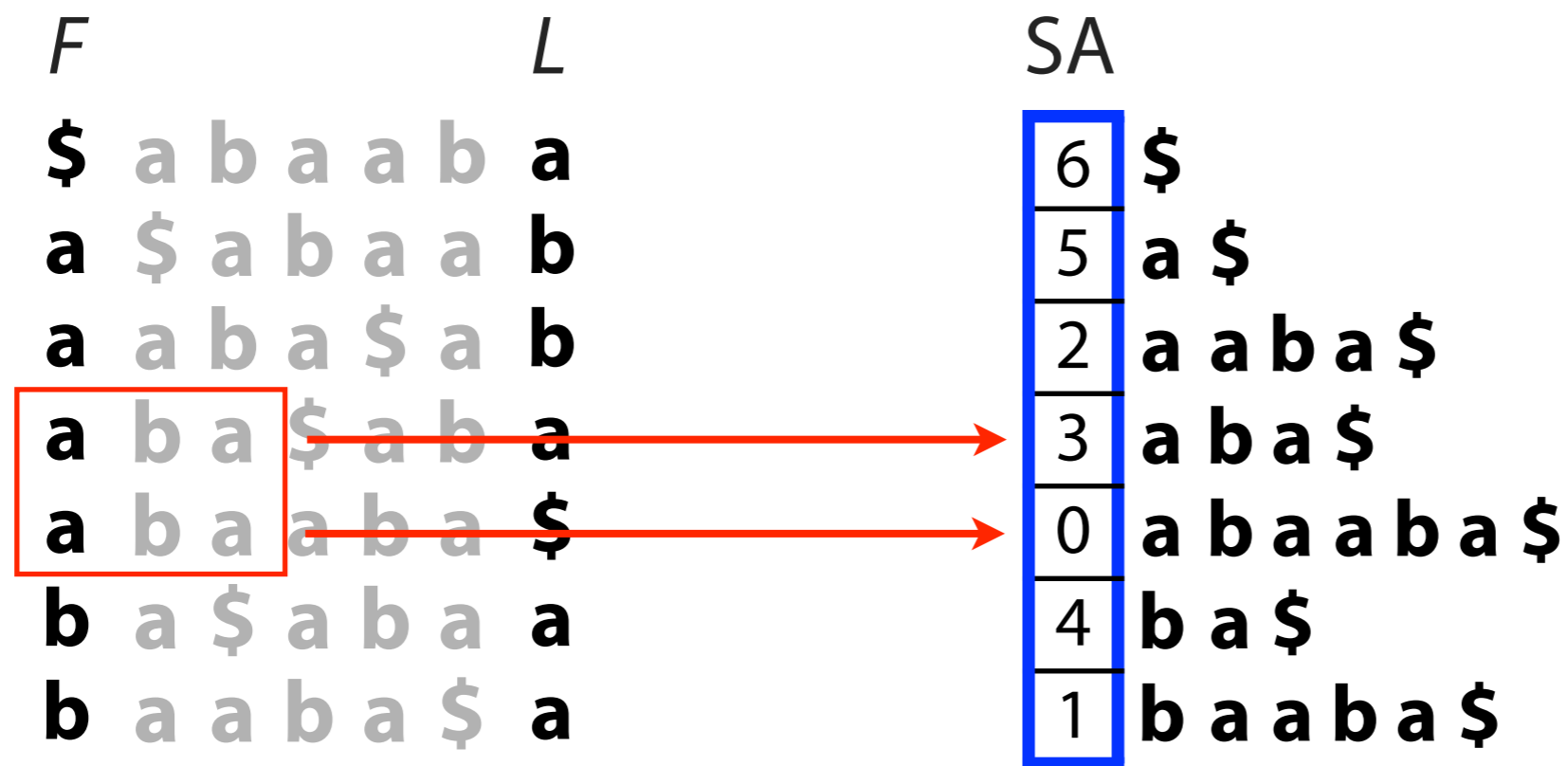
Locate query & SA samples

A precomputed "suffix array" stores the answers...



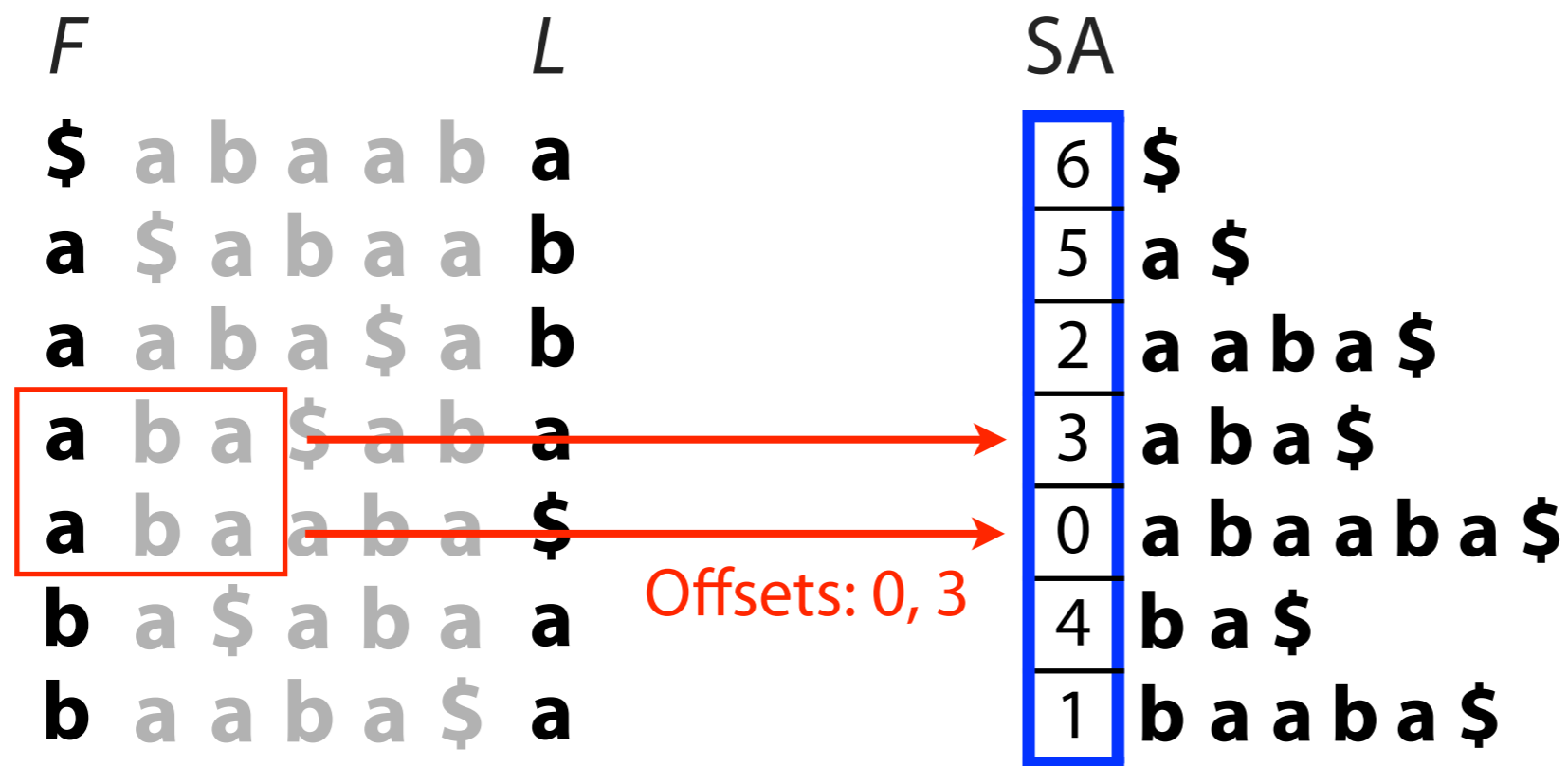
Locate query & SA samples

A precomputed "suffix array" stores the answers...



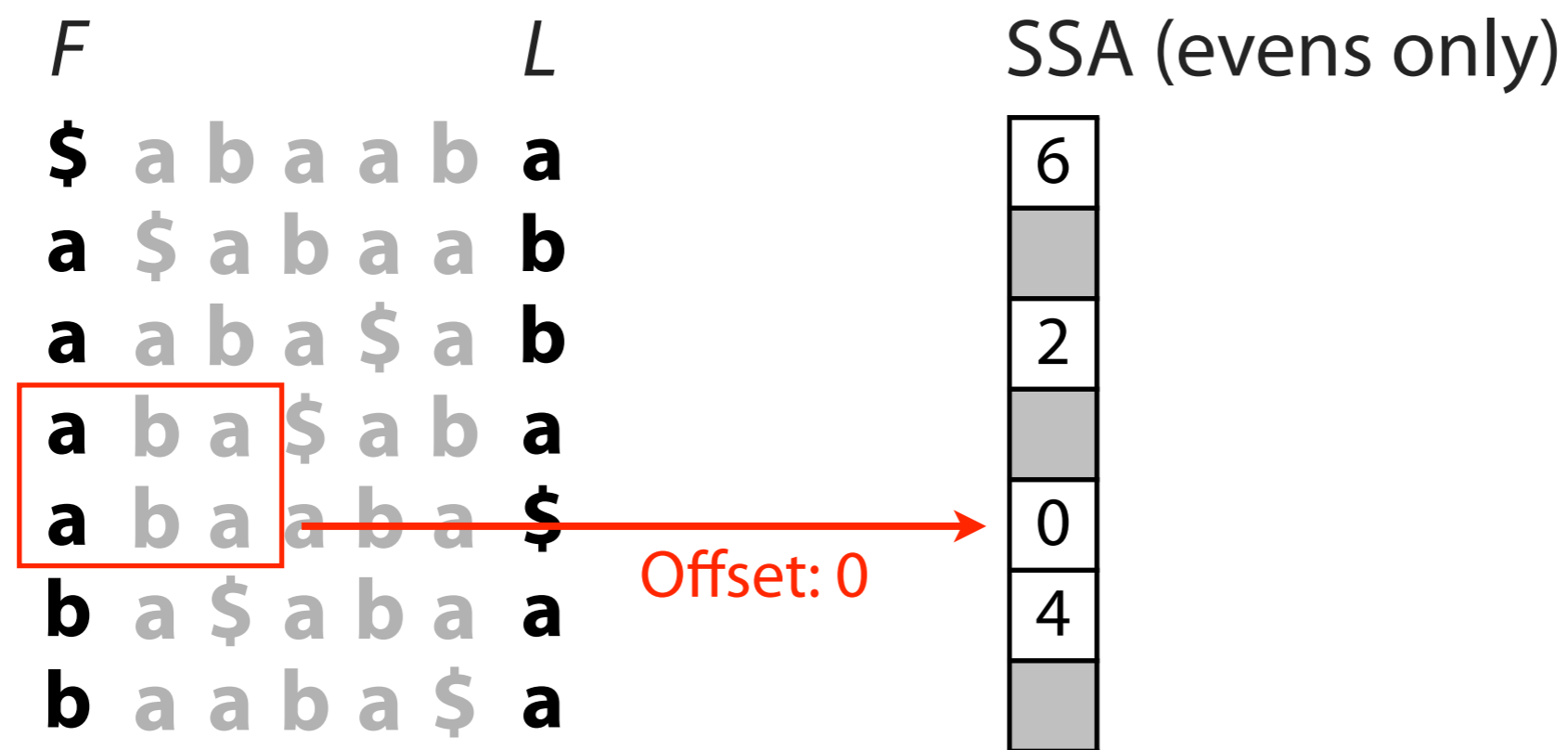
Locate query & SA samples

A precomputed "suffix array" stores the answers...

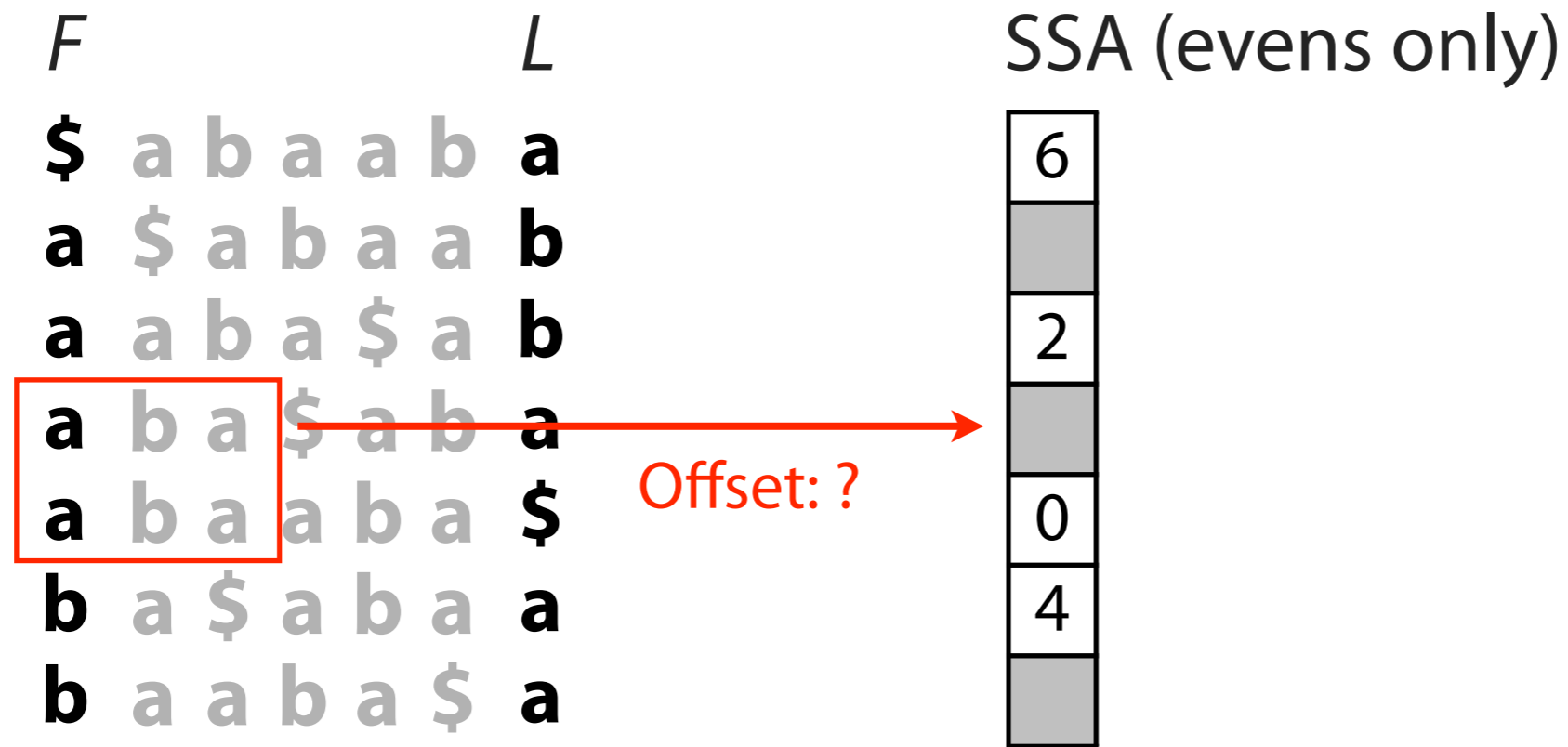


n integers is too big, so we **sample** instead

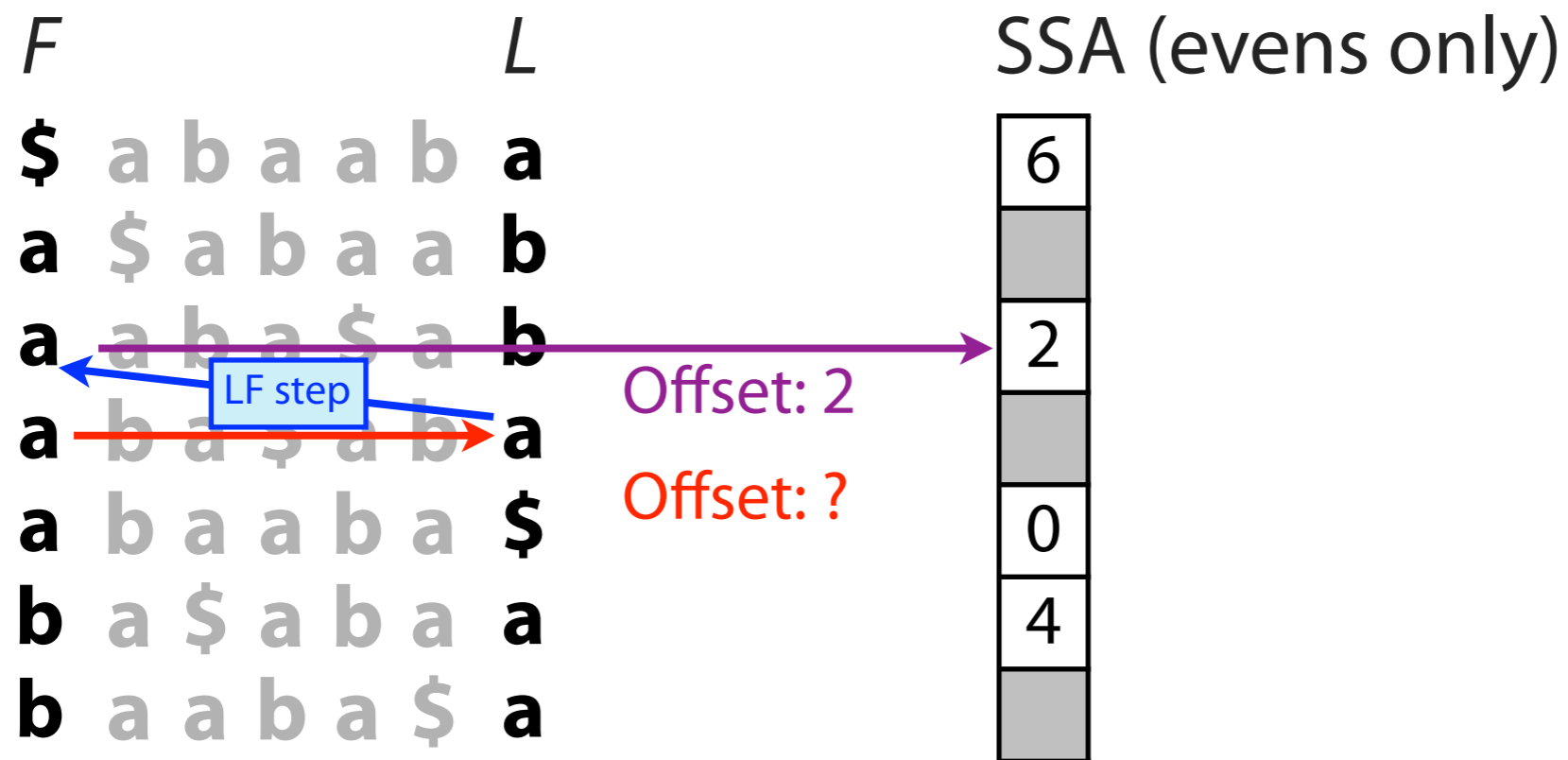
Locate query & SA samples



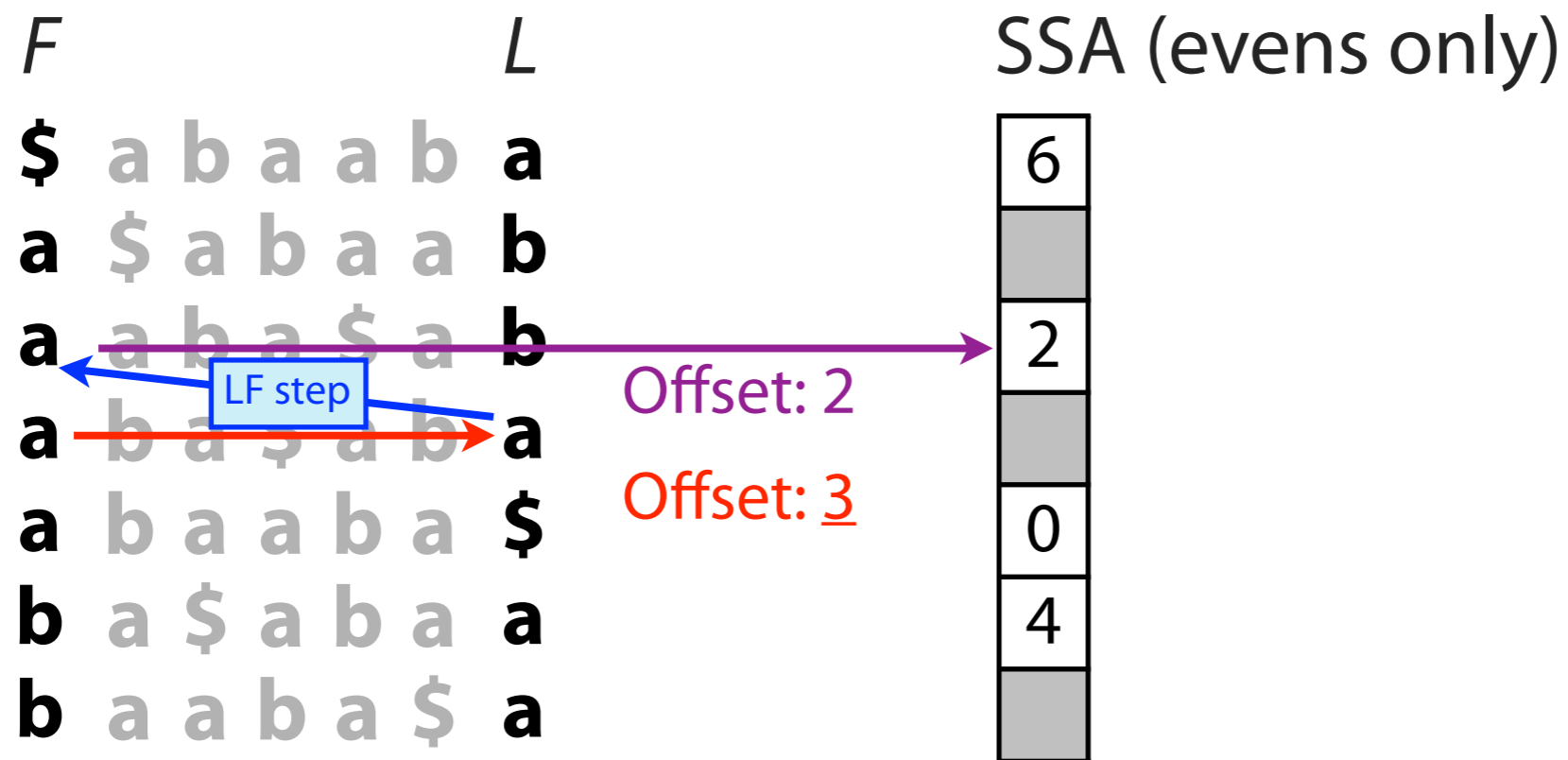
Locate query & SA samples



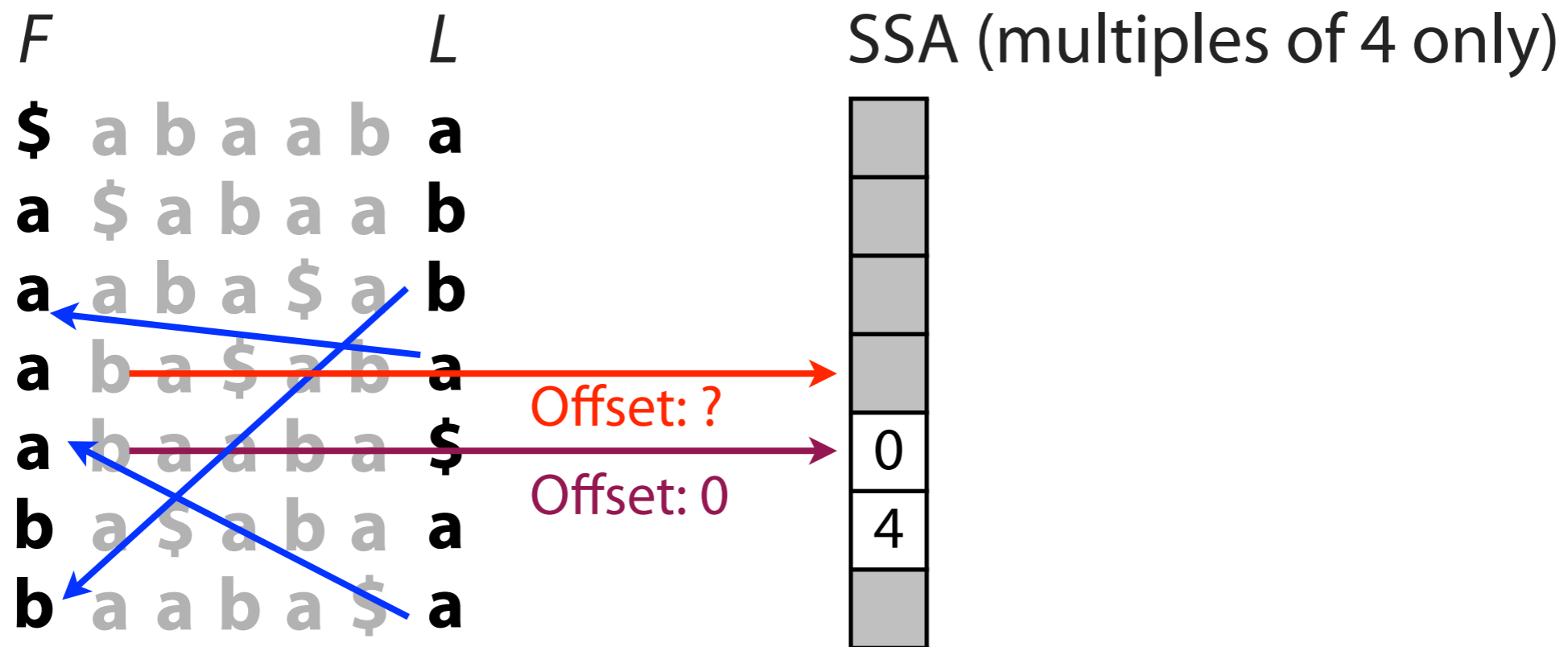
Locate query & SA samples



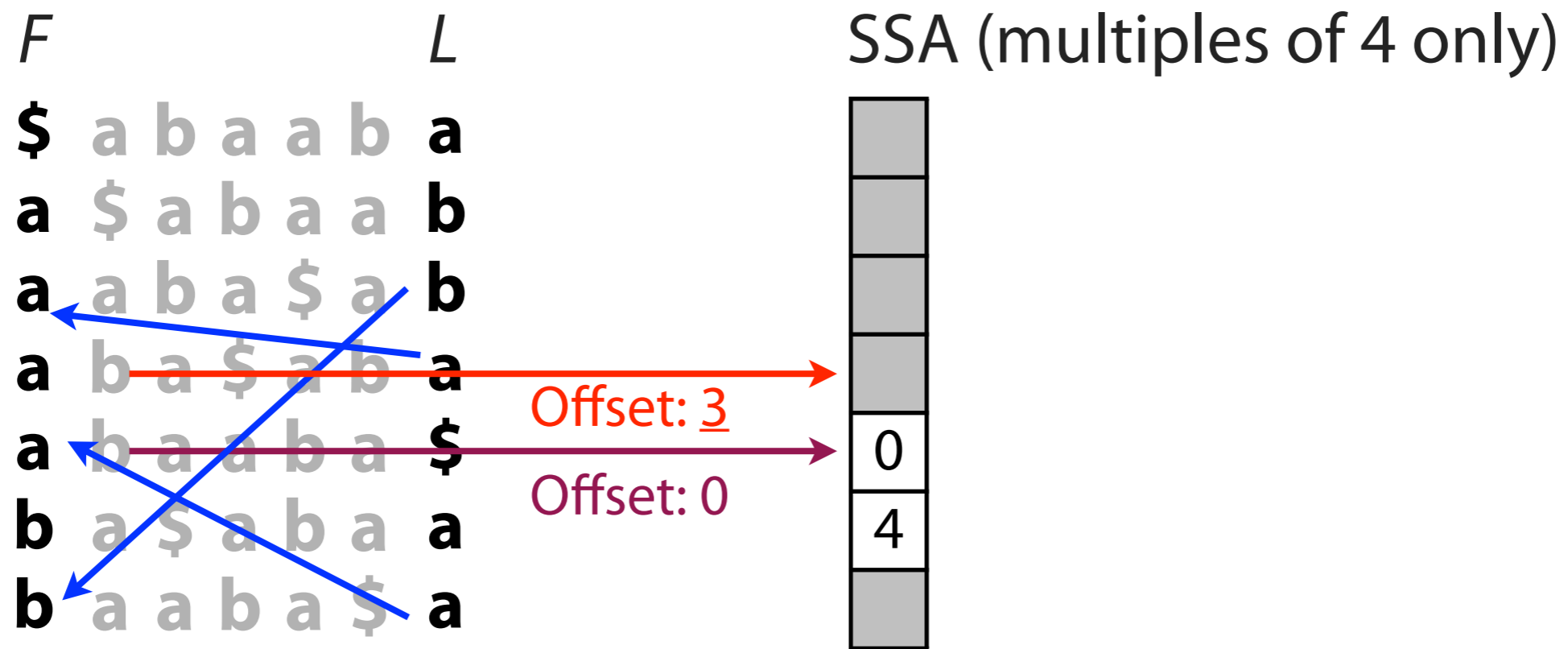
Locate query & SA samples



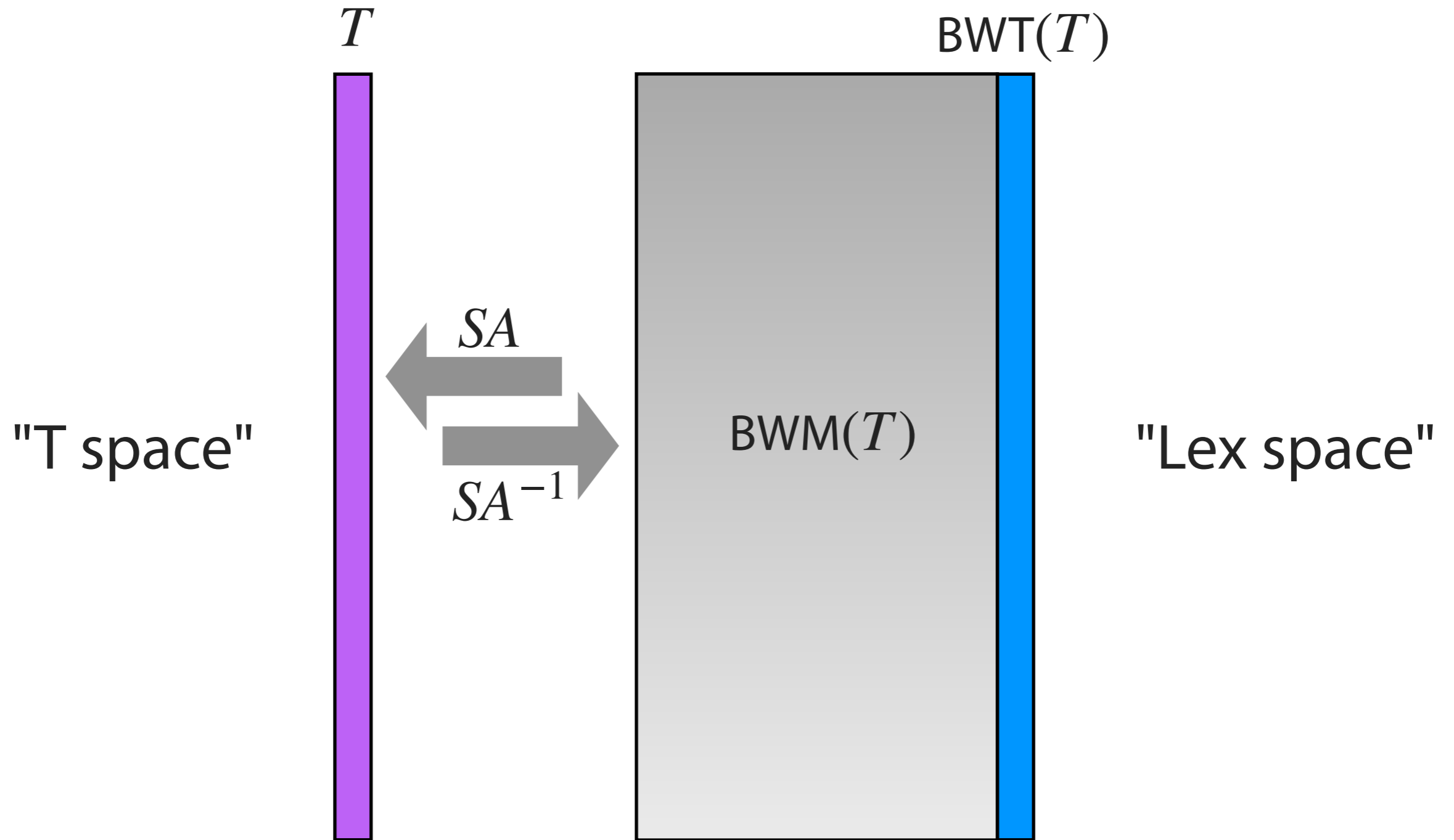
Locate query & SA samples



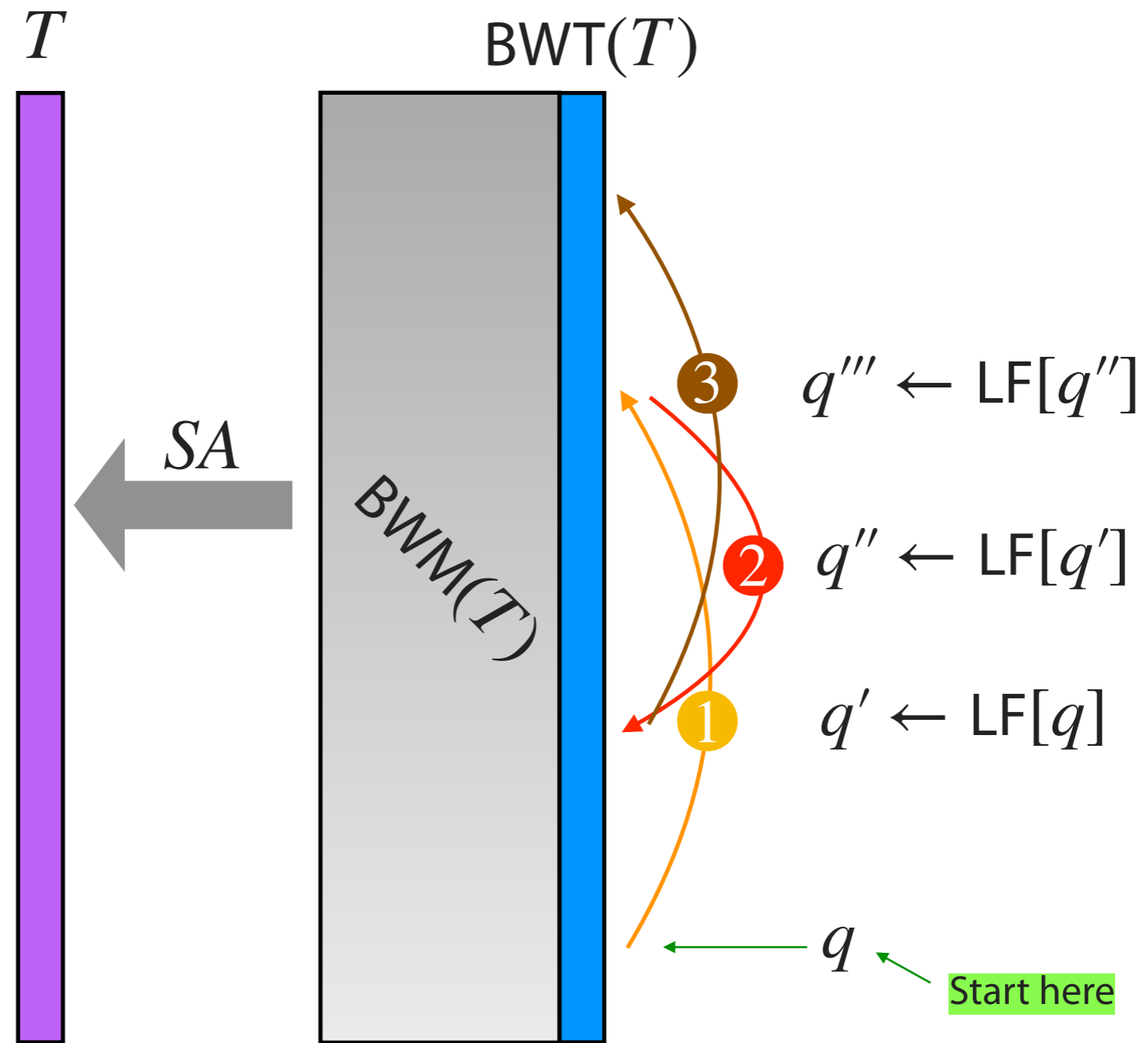
Locate query & SA samples



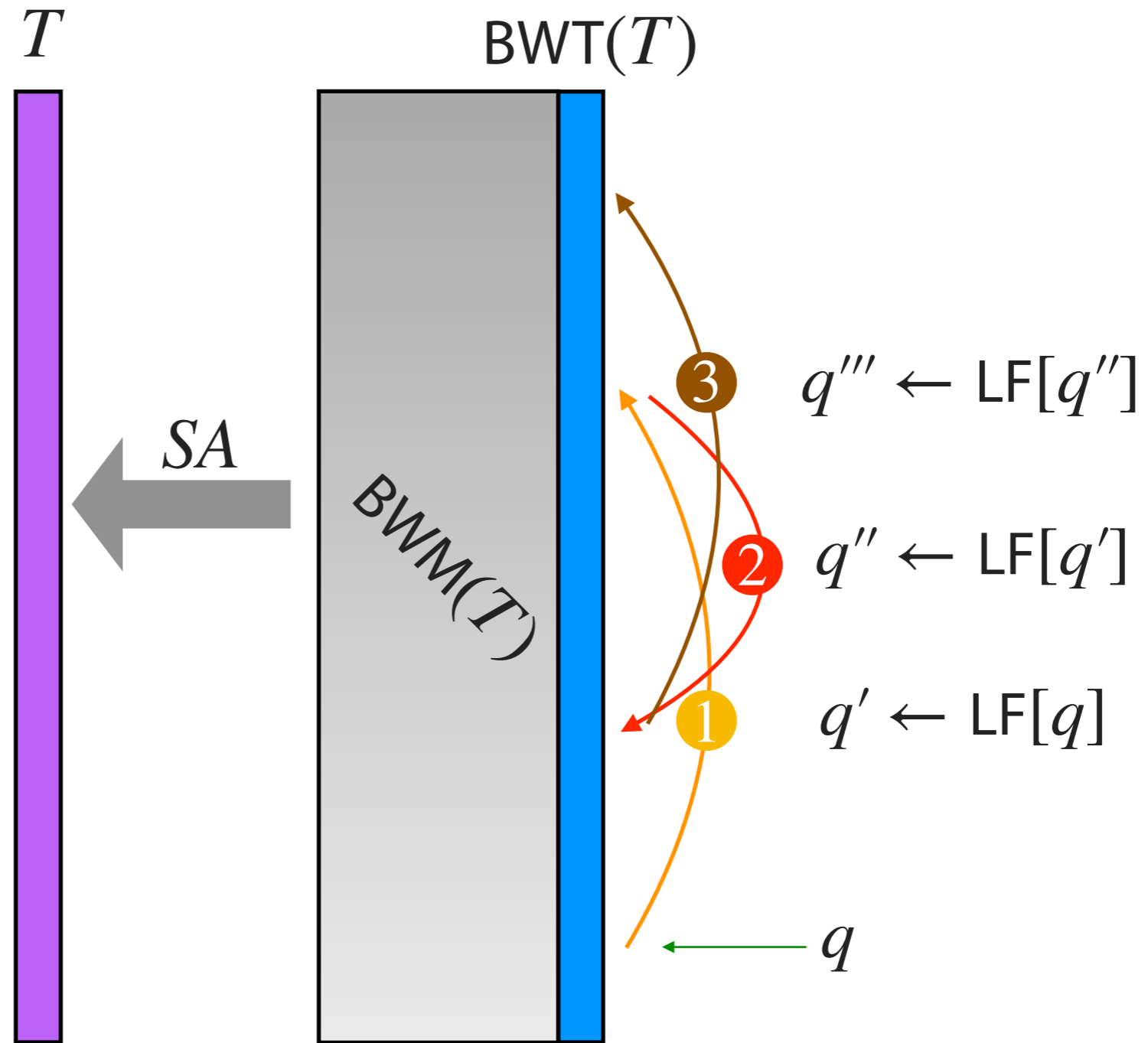
Fundamental movements



Fundamental movements: LF



Fundamental movements: LF

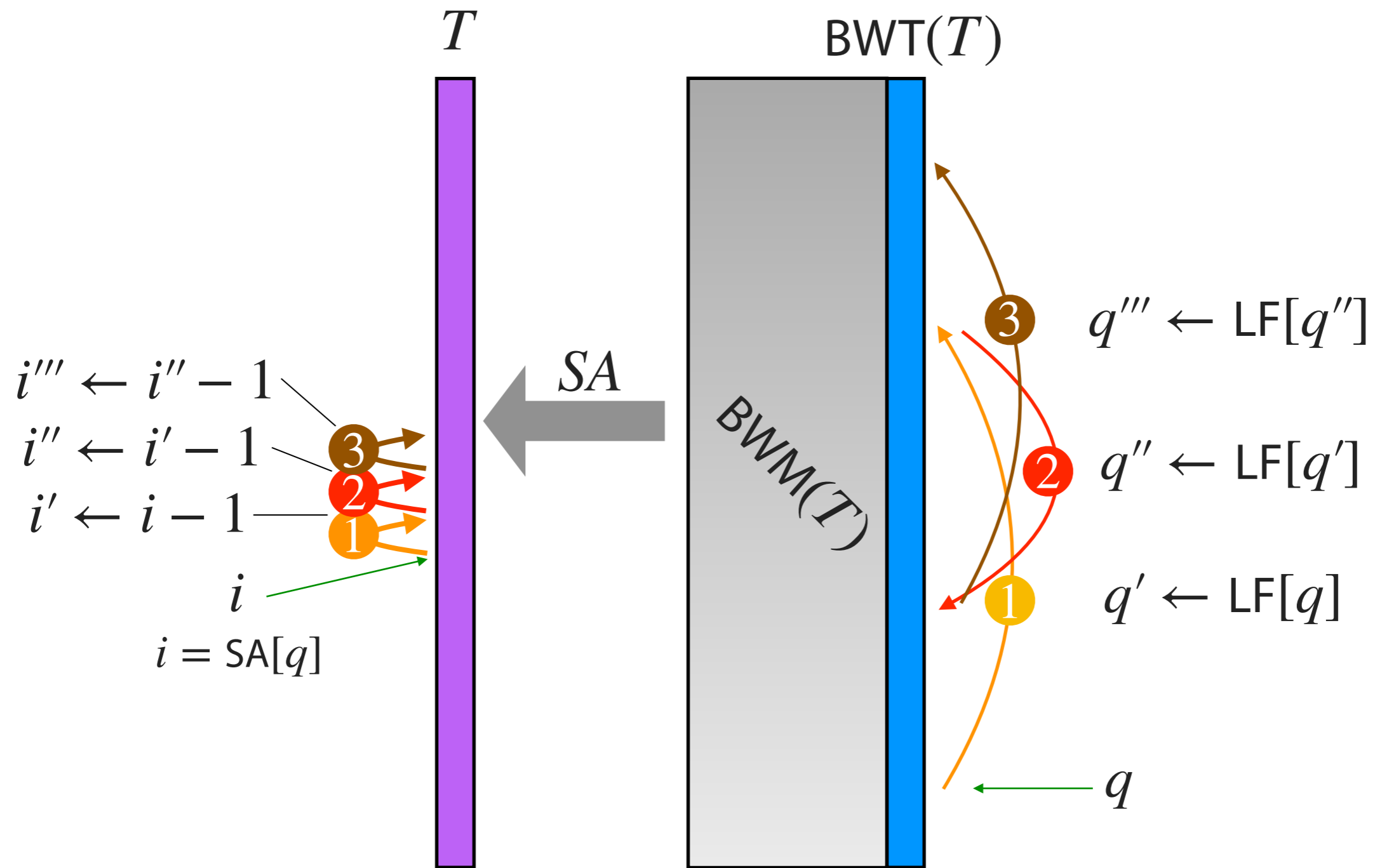


An LF step
(in lex space)...

$$SA[LF[q]] = SA[q] - 1$$

...moves to the left
by 1 in T space

Fundamental movements: LF

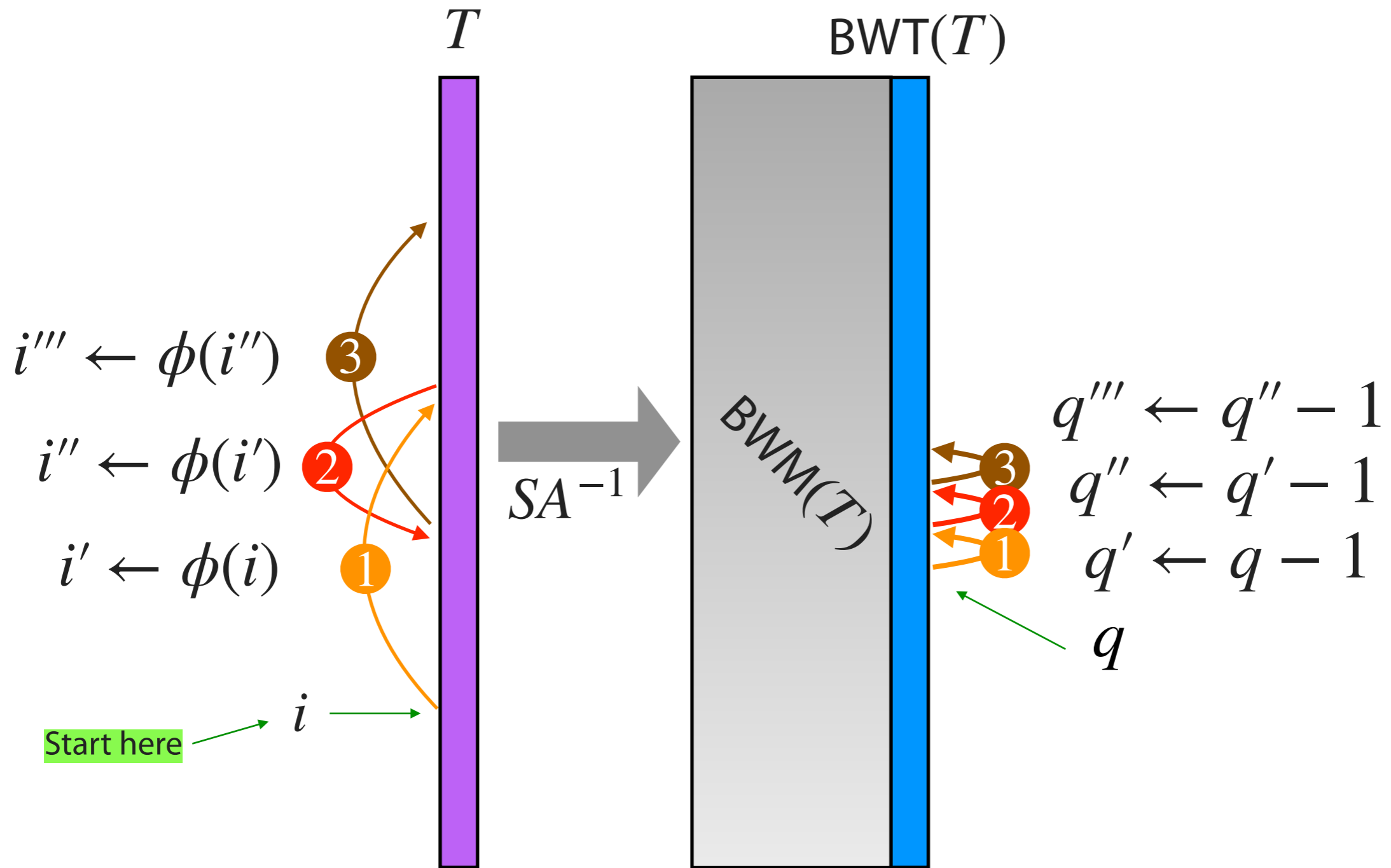


An LF step
(in lex space)...

$$SA[LF[q]] = SA[q] - 1$$

...moves to the left
by 1 in T space

Fundamental movements: ϕ



$$SA^{-1}[\phi(i)] = SA^{-1}[i] - 1$$

Fundamental movements: "Runny LF"

When T is repetitive, LF can be thought of as permuting **runs**

I.e. LF mapping is "runny"

T row_row_row_your_boatrow_row_row_your_boatrow_row_row_your_boat\$

L trrrwwwwwwwooo__bbbyyrrrrrrrrrruuutt\$_____aaaoooooooooooo__

F \$_____aaabbboooooooooooooooooorrrrrrrrrrrtttuuwwwwwwyyy



Fundamental movements: "Runny LF"

When T is repetitive, LF can be thought of as permuting **runs**

I.e. LF mapping is "runny"

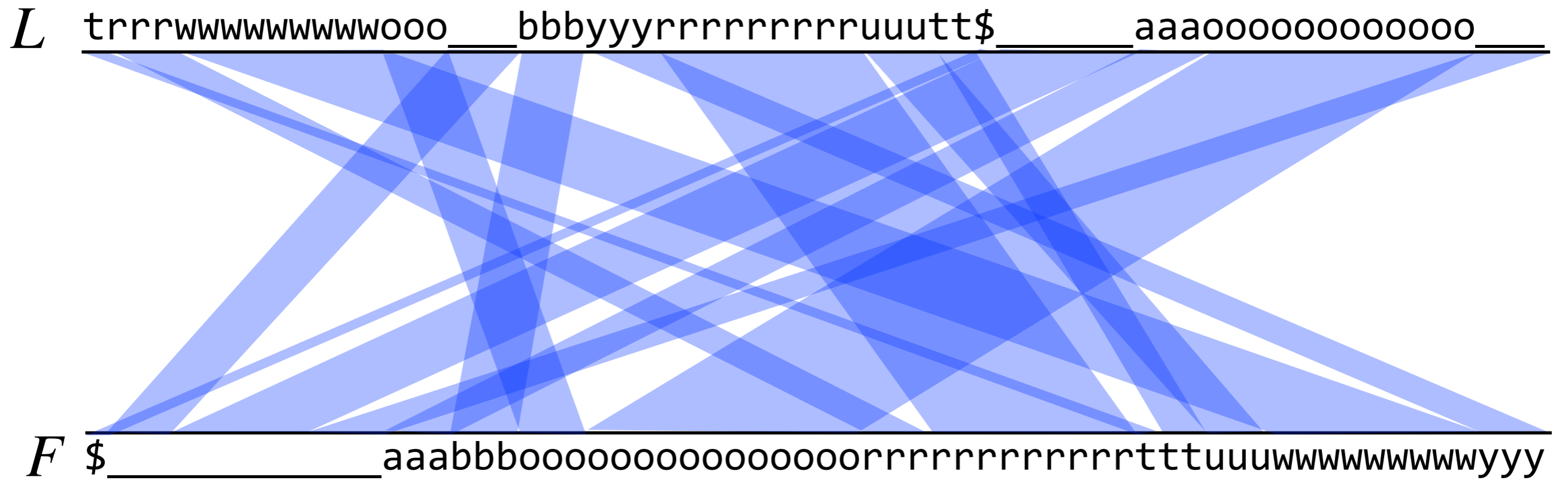
T row_row_row_your_boatrow_row_row_your_boatrow_row_row_your_boat\$

L trrrwwwwwwwooo__bbbyyrrrrrrrrrruuutt\$__aaoooooooooooo__

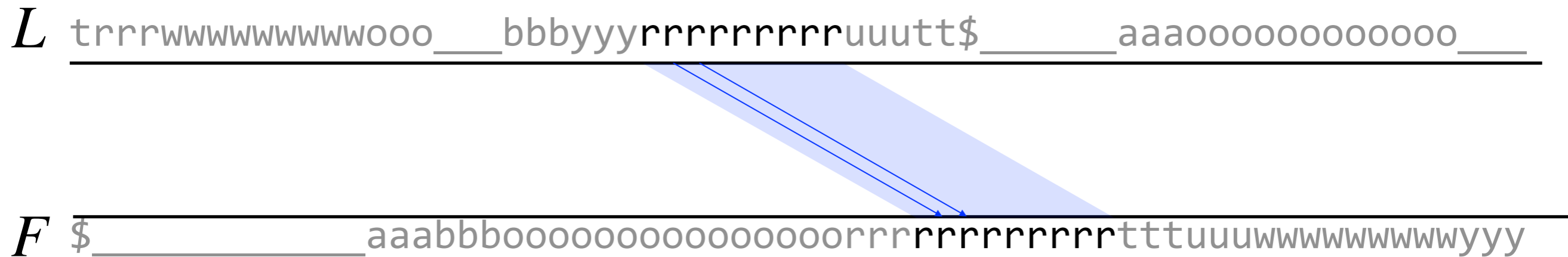
F \$____aaabbboooooooooooooooooorrrrrrrrrrrrrrrrrtttuuwwwwwwyyy



Runny LF



Runny LF principle

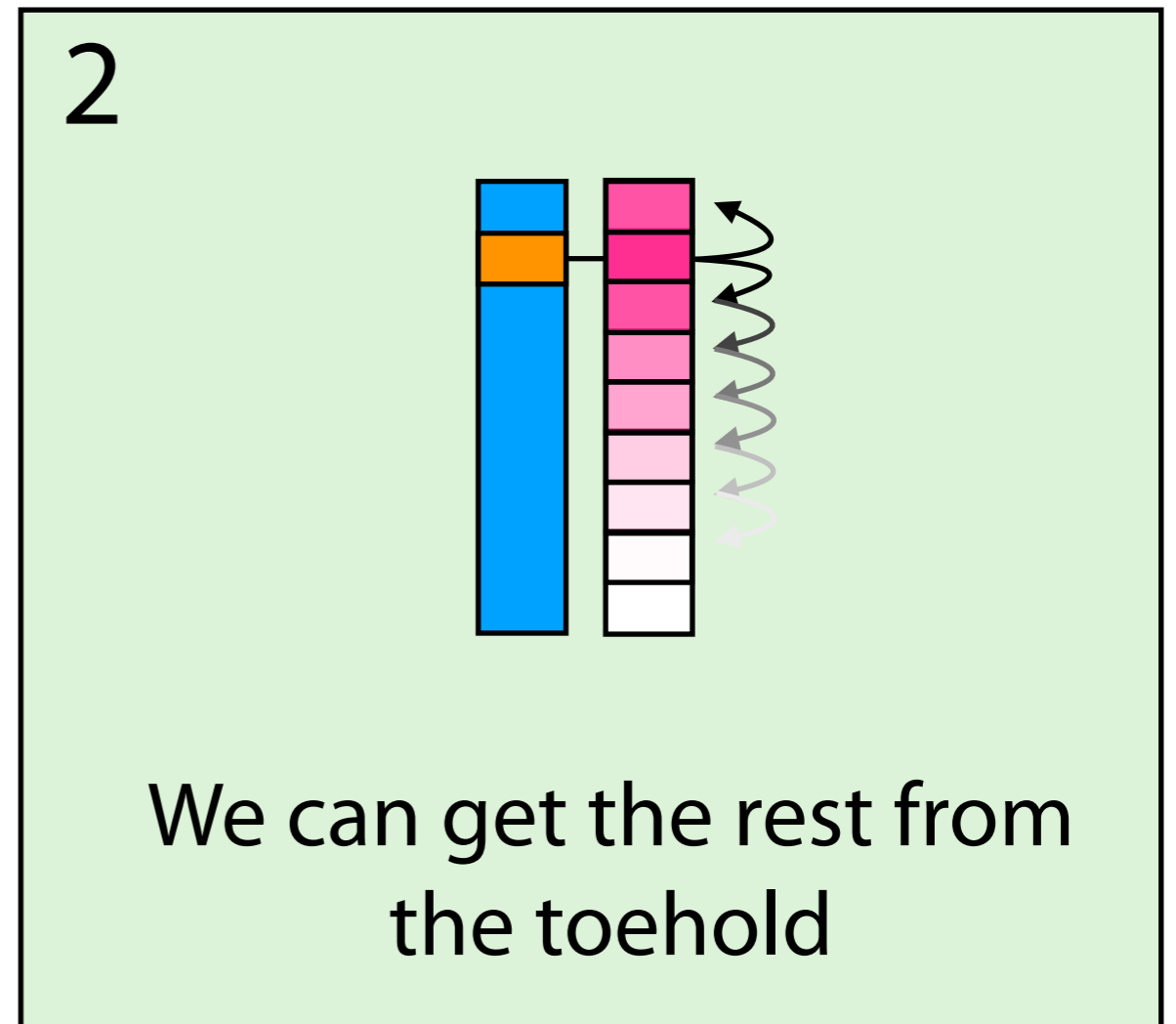
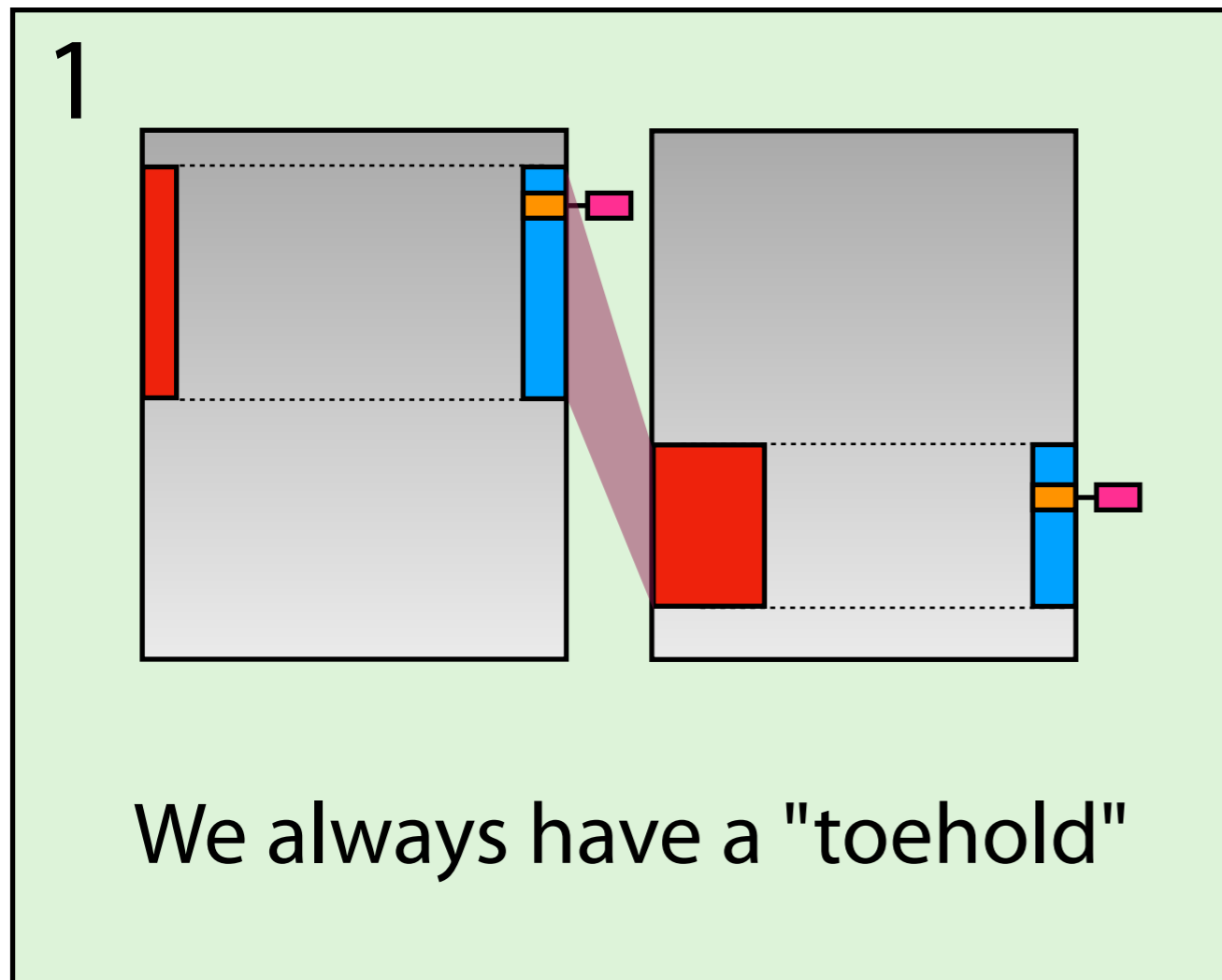


If $BWT[q] = BWT[q + 1]$, then $LF[q + 1] = LF[q] + 1$

Two positions in the same BWT run move through LF in a **parallel** fashion

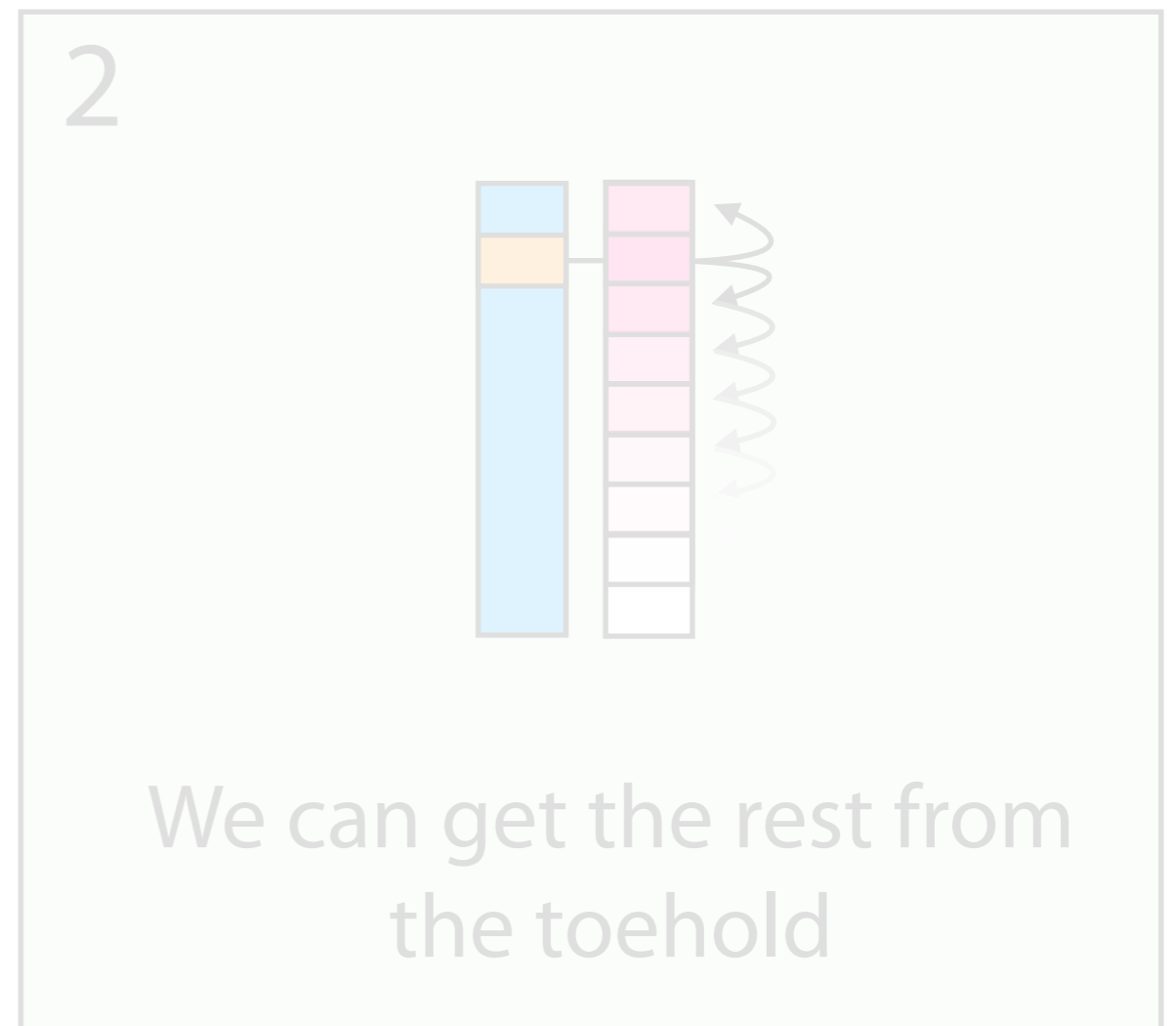
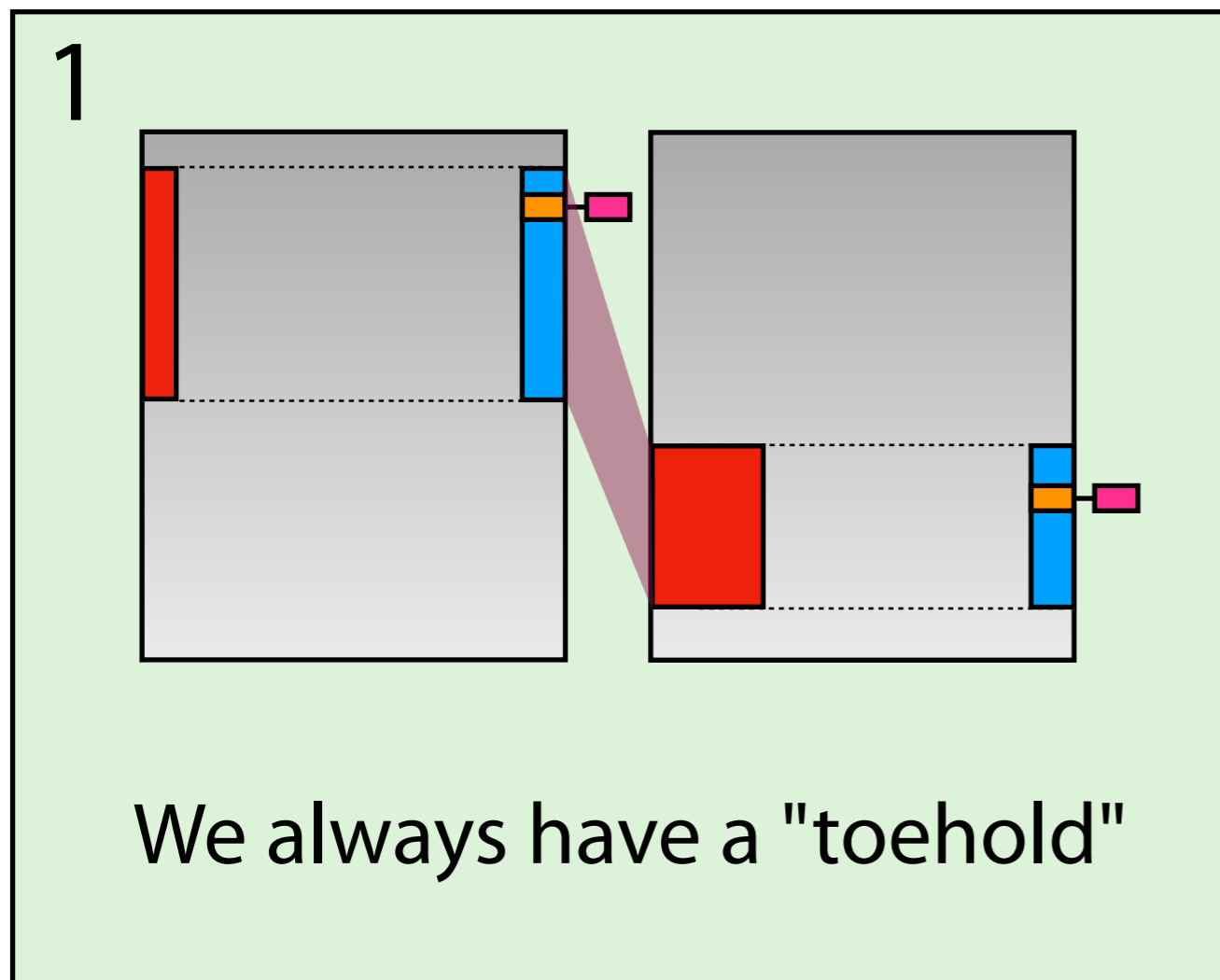
$O(r)$ locate queries

Algorithm has two components:



$O(r)$ locate queries

Algorithm has two components:



Locate query & SA samples

T row_row_row_your_boatrow_row_row_your_boatrow_row_row_your_boat\$
L trrrwwwwwwwooo__bbbyyrrrrrrrrruutt\$__aaoooooooooooo__

Samples at run heads

Locate query & SA samples

T row_row_row_your_boatrow_row_row_your_boatrow_row_row_your_boat\$

L trrrwwwwwwwooo__bbbyyrrrrrrrrrrruutt\$__aaoooooooooooo__

62

Samples at run heads

Locate query & SA samples

T row_row_row_your_boatrow_row_row_your_boatrow_row_row_your_boat\$

L trrrwwwwwwwooo__bbbyyrrrrrrrrrrruutt\$__aaooooooooooooo__

62 57



Samples at run heads

Locate query & SA samples

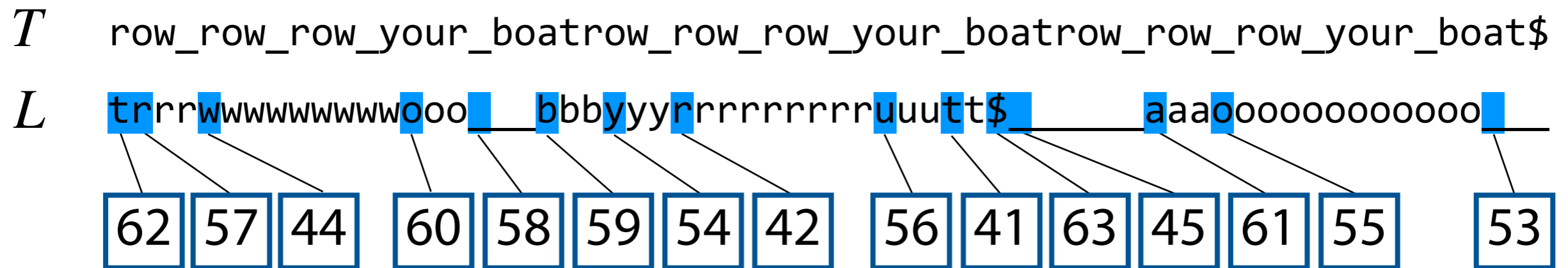
T row_row_row_your_boatrow_row_row_your_boatrow_row_row_your_boat\$

L trrrwwwwwwwooo__bbbyyrrrrrrrrrrruutt\$__aaooooooooooooo__



Samples at run heads

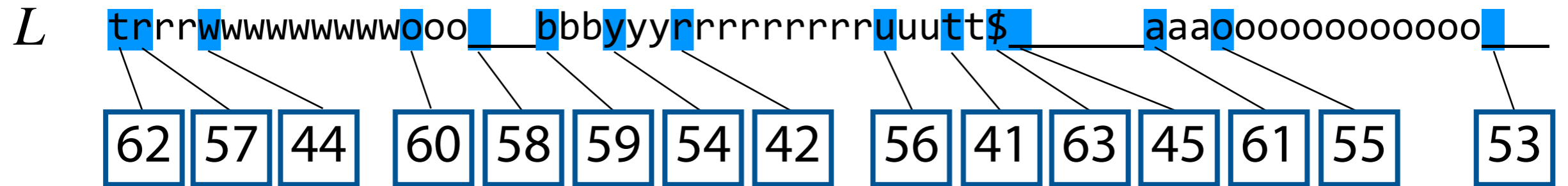
Locate query & SA samples



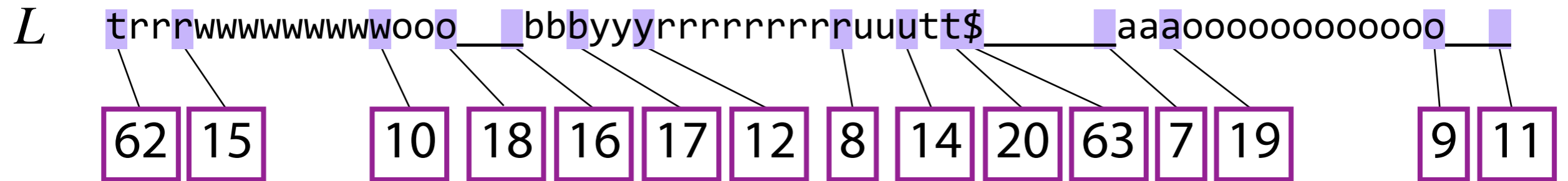
Samples at run heads

Locate query & SA samples

T row_row_row_your_boatrow_row_row_your_boatrow_row_row_your_boat\$

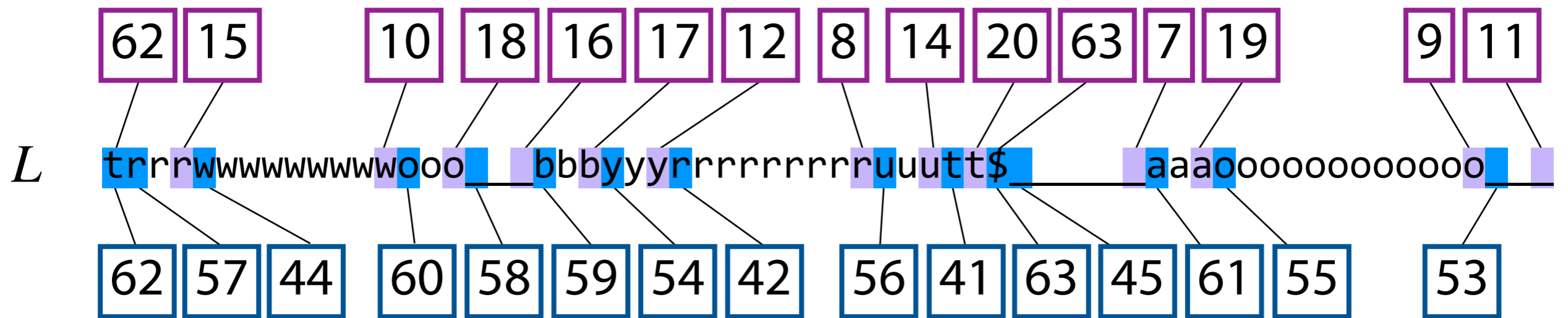


Samples at run heads



Samples at run tails

Locate query & SA samples

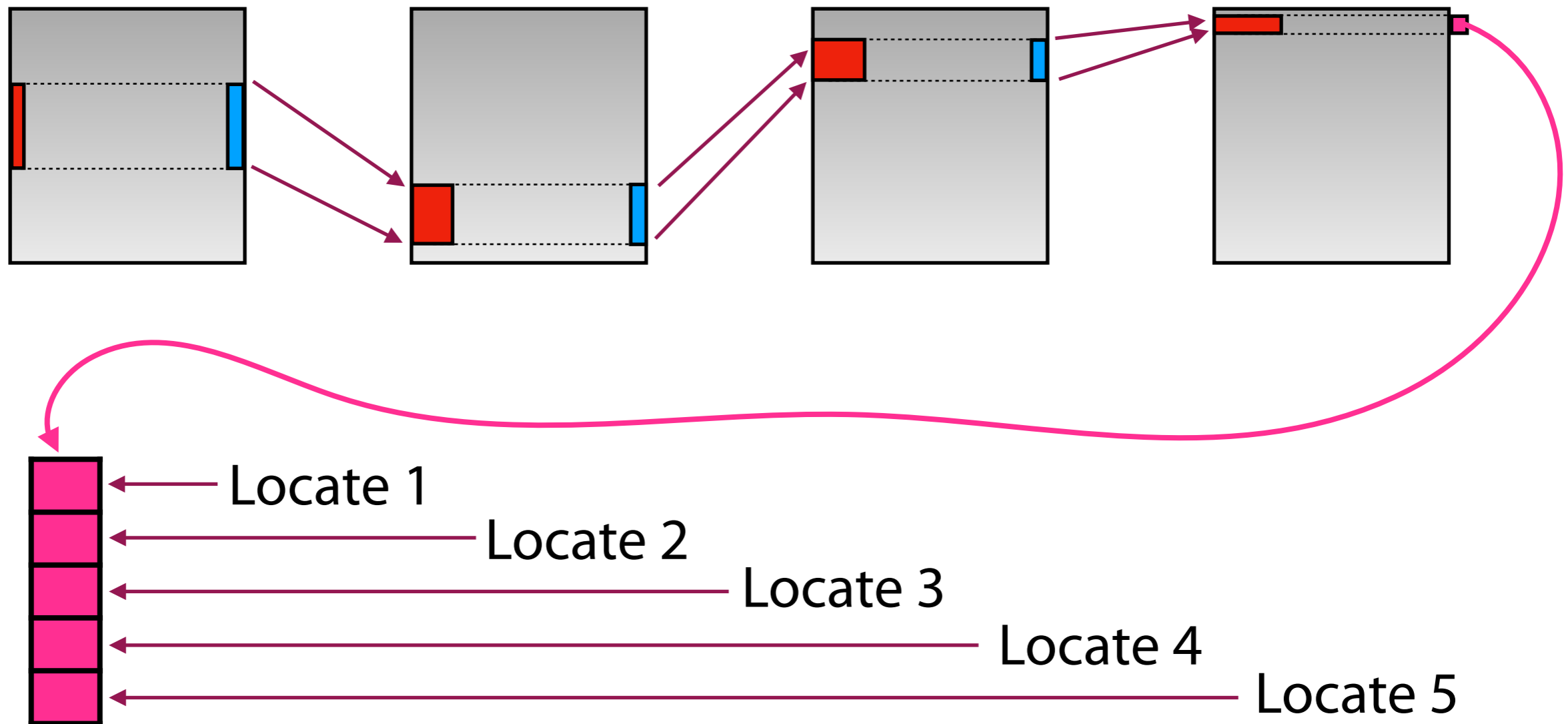


Samples at run heads & tails

Standard backward search

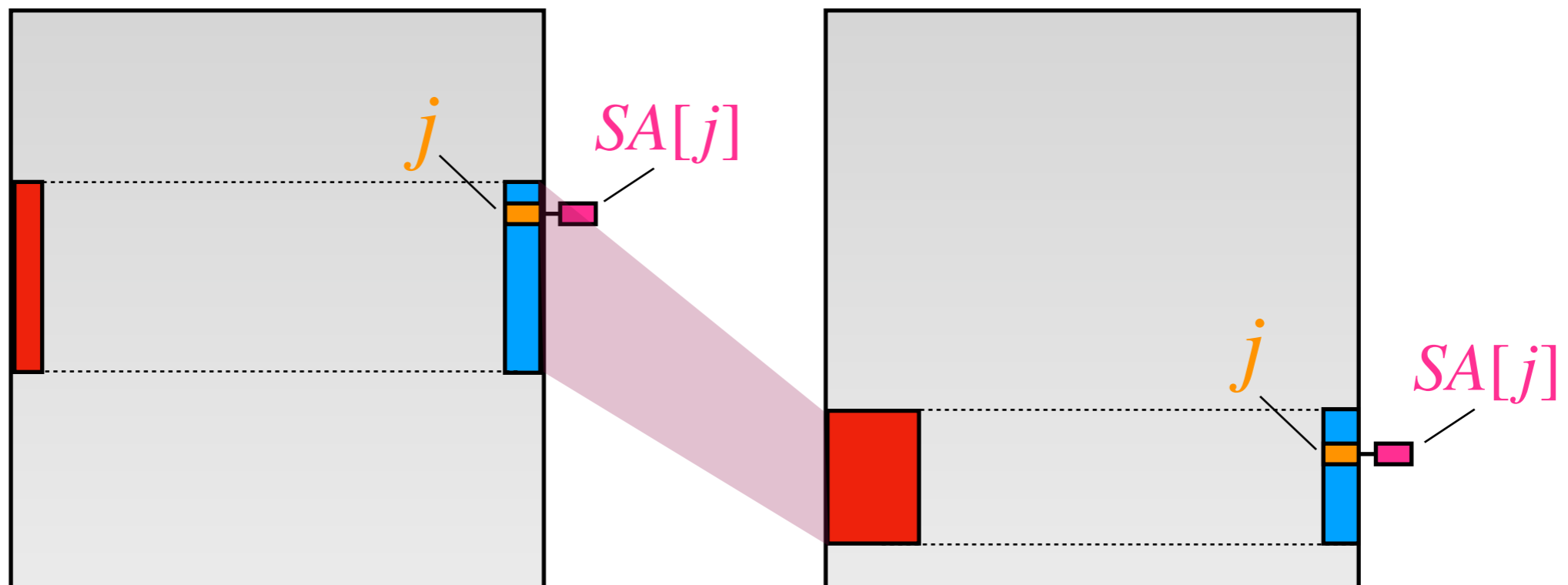
Usually we do backward search first, locate queries second:

Backward search



Backward search with toehold

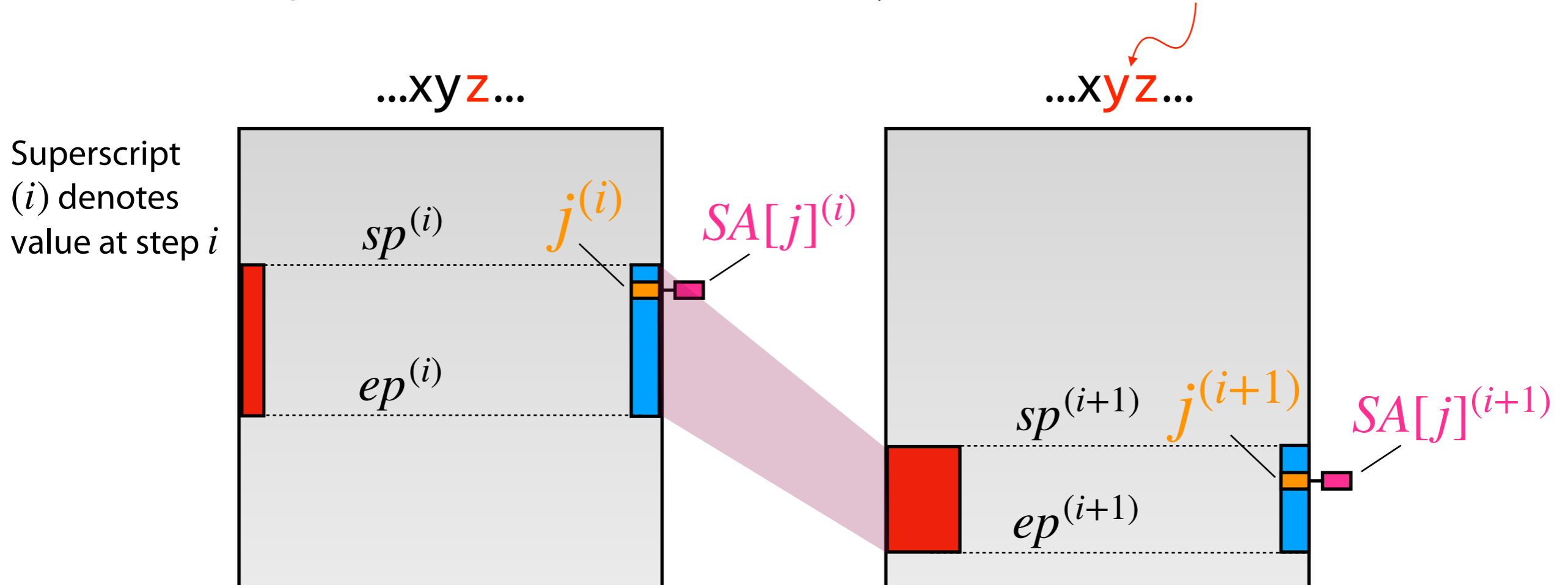
Augment backward search so that *at every step*, we have one offset j in our range $[sp, ep]$ for which we know $SA[j]$



How to update j and $SA[j]$ at each step?

Backward search with toehold

Say we have a j , $[sp, ep]$ and $SA[j]$ for step i and we're updating for step $i + 1$. Next query character is y .

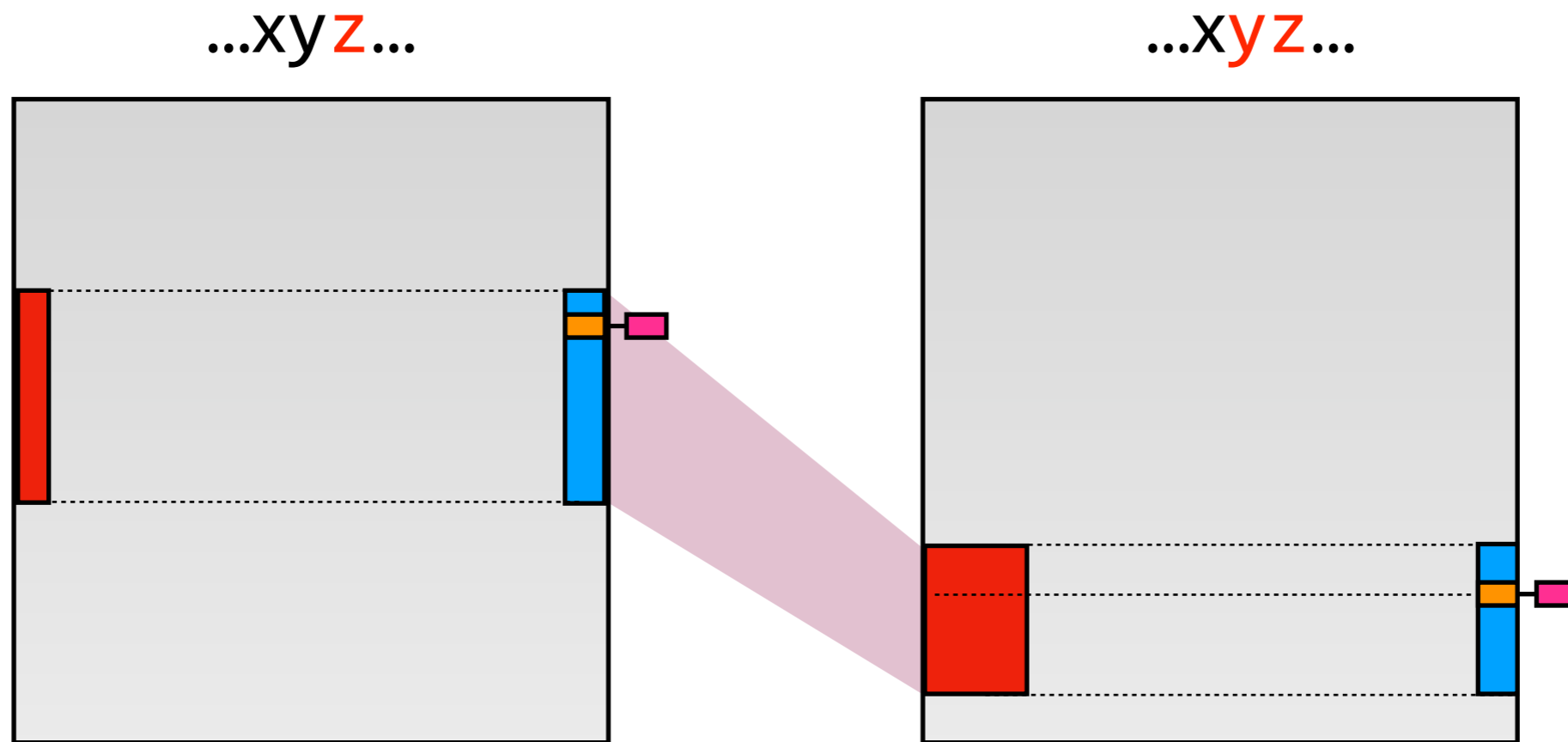


Two cases: Case 1: $BWT[j^{(i)}] = y$ Case 2: $BWT[j^{(i)}] \neq y$

Backward search with toehold

Case 1: $\text{BWT}[j^{(i)}] = y$

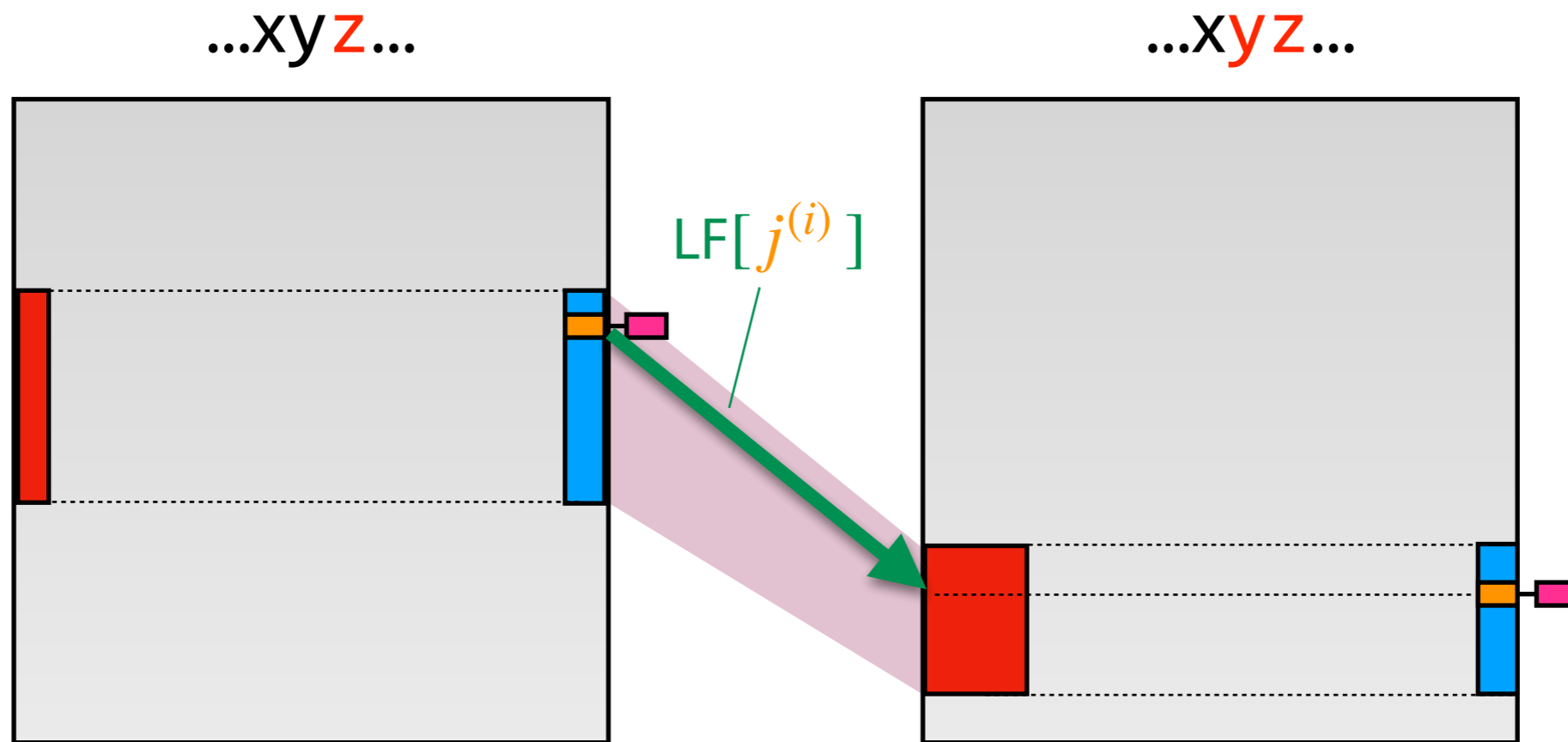
LF does all the work



Backward search with toehold

Case 1: $\text{BWT}[j^{(i)}] = y$

LF does all the work



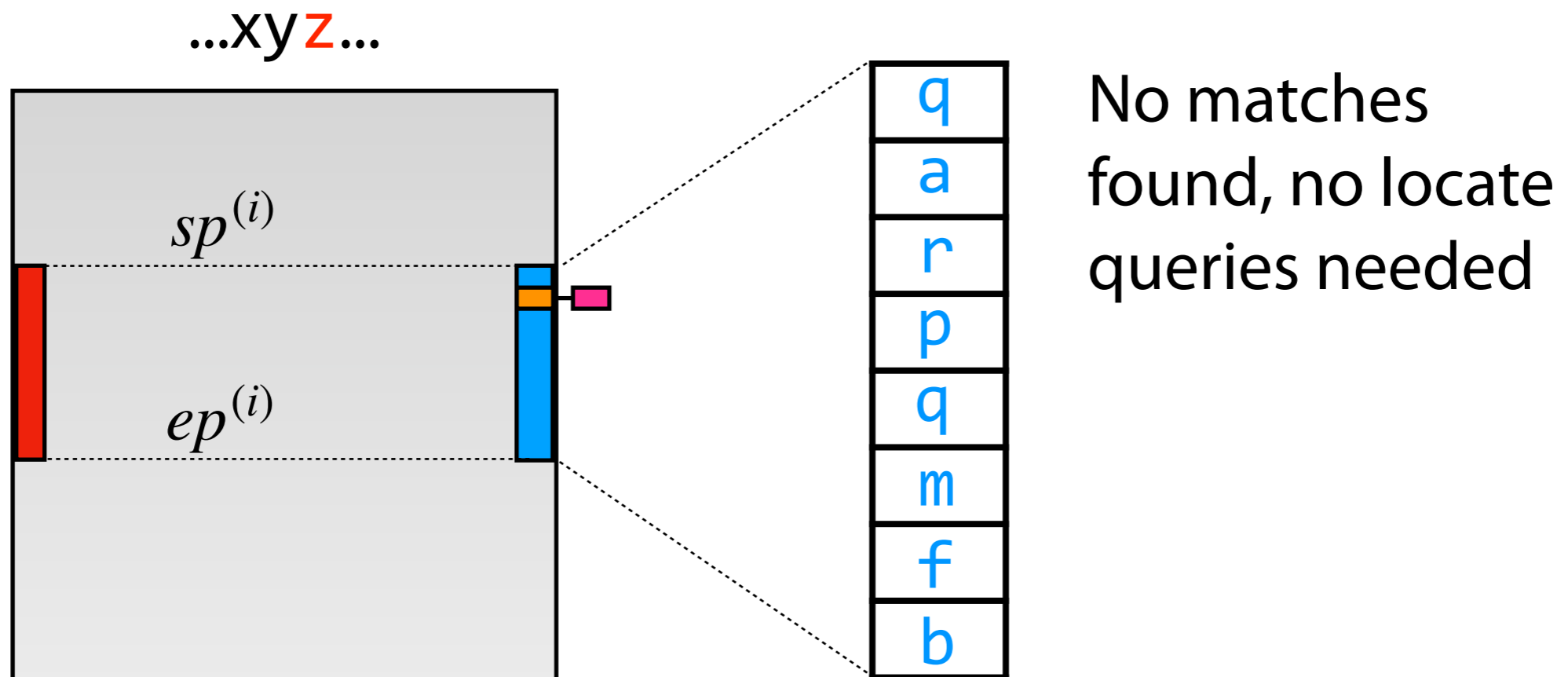
$$j^{(i+1)} = \text{LF}[j^{(i)}]$$

$$\text{SA}[j]^{(i+1)} = \text{SA}[j]^{(i)} - 1$$

Backward search with toehold

Case 2: $\text{BWT}[j^{(i)}] \neq y$

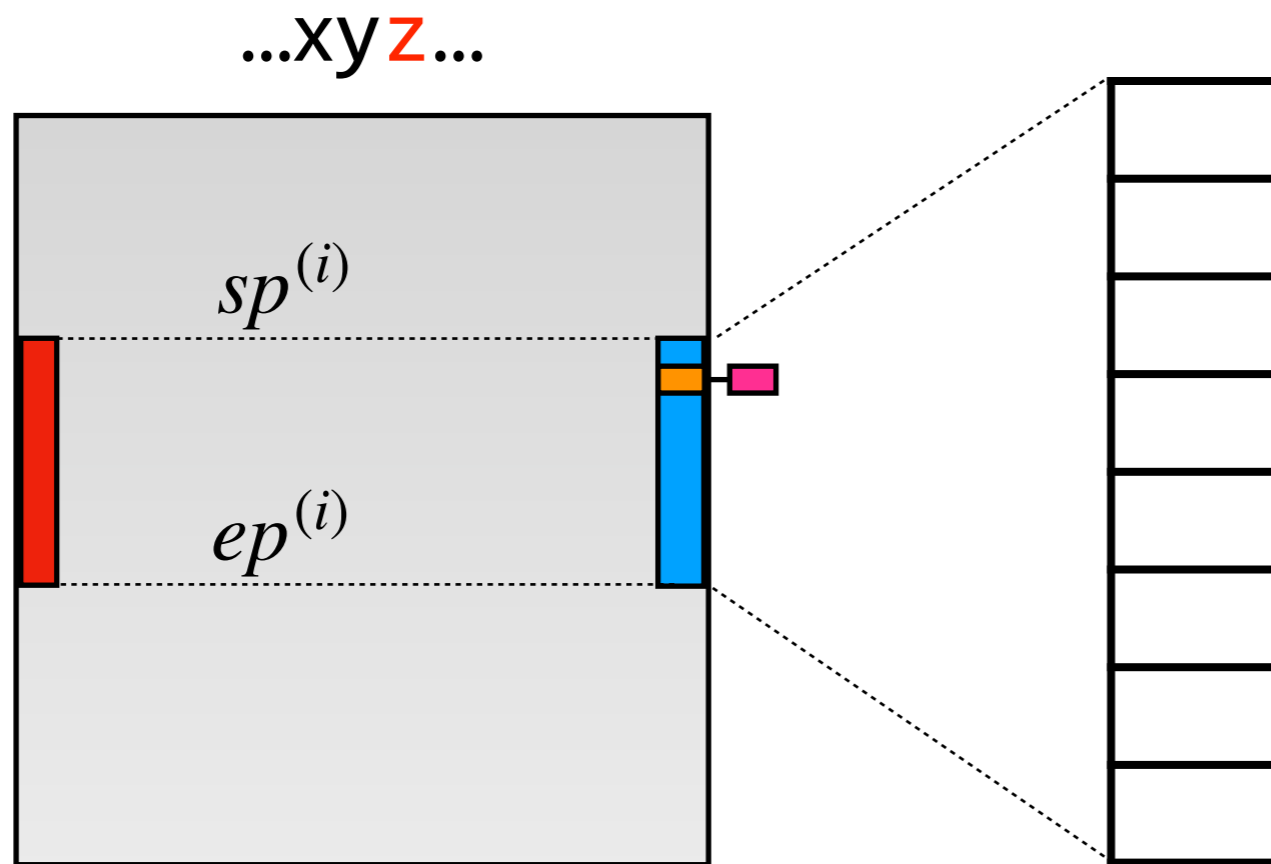
Case 2a: **No** element in $\text{BWT}[sp^{(i)} .. ep^{(i)}]$ equals y



Case 2b (next slide): Some elements in $\text{BWT}[sp^{(i)} .. ep^{(i)}]$ match y , but $\text{BWT}[j^{(i)}]$ doesn't

Backward search with toehold

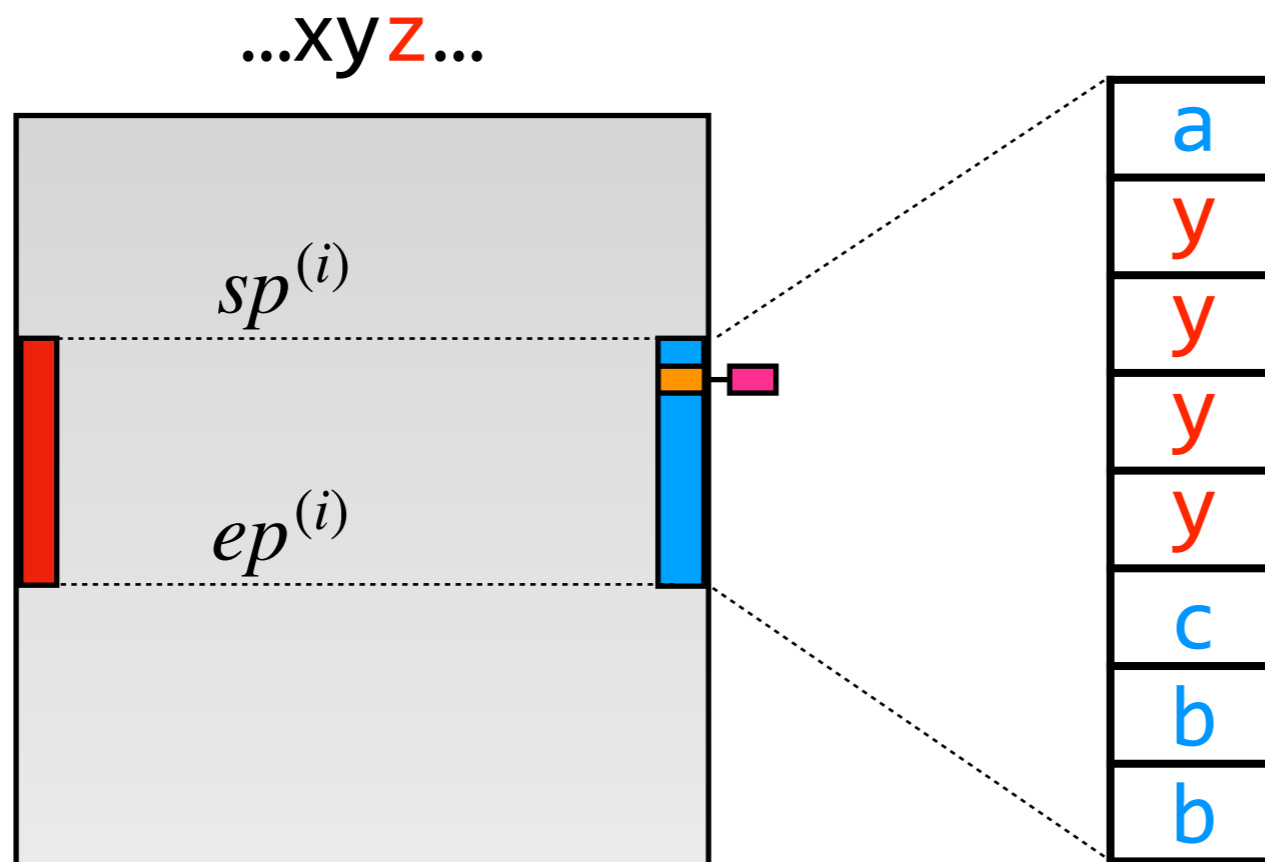
Case 2b: Some elements in $\text{BWT}[sp^{(i)} .. ep^{(i)}]$ match y ,
but $\text{BWT}[j^{(i)}]$ doesn't



$\text{BWT}[sp^{(i)} .. ep^{(i)}]$
must contain a
 y -run head or tail.

Backward search with toehold

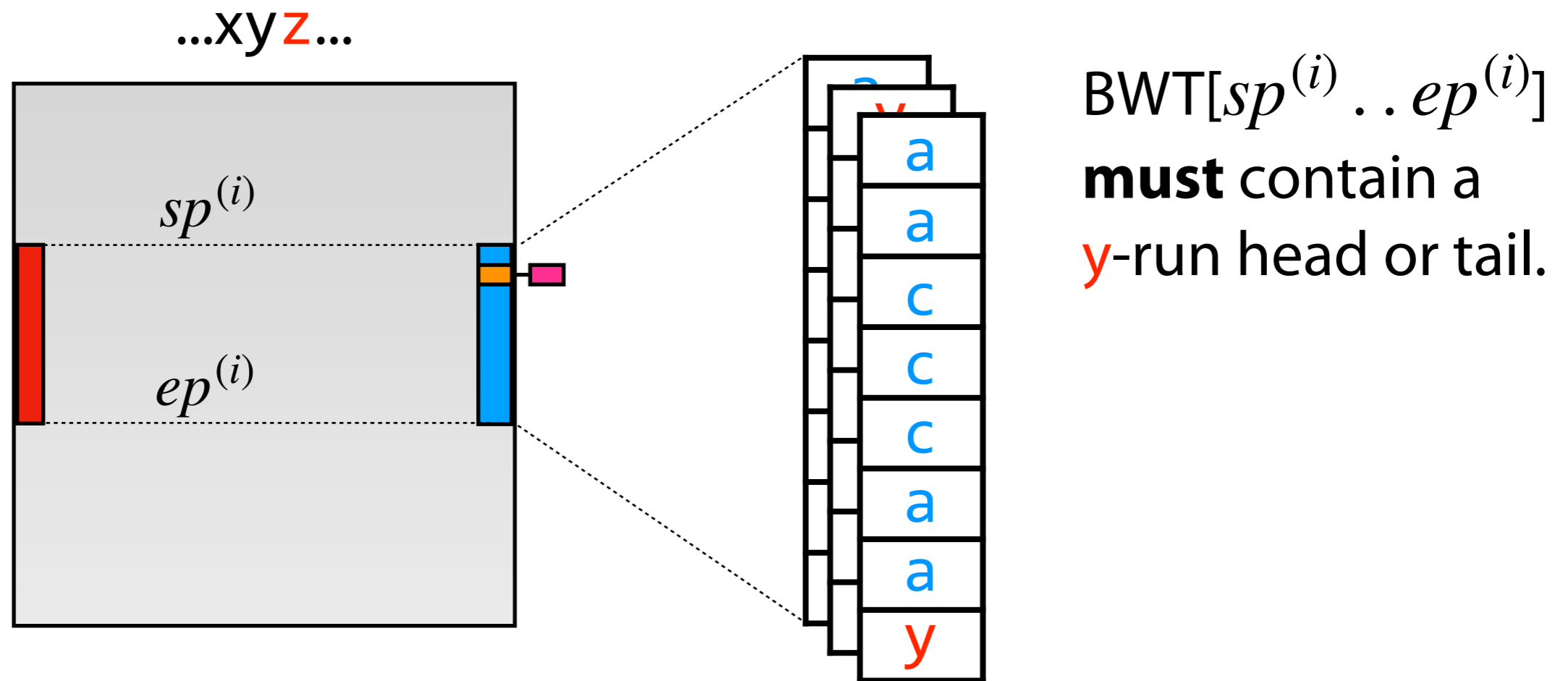
Case 2b: Some elements in $\text{BWT}[sp^{(i)} .. ep^{(i)}]$ match y ,
but $\text{BWT}[j^{(i)}]$ doesn't



$\text{BWT}[sp^{(i)} .. ep^{(i)}]$
must contain a
 y -run head or tail.

Backward search with toehold

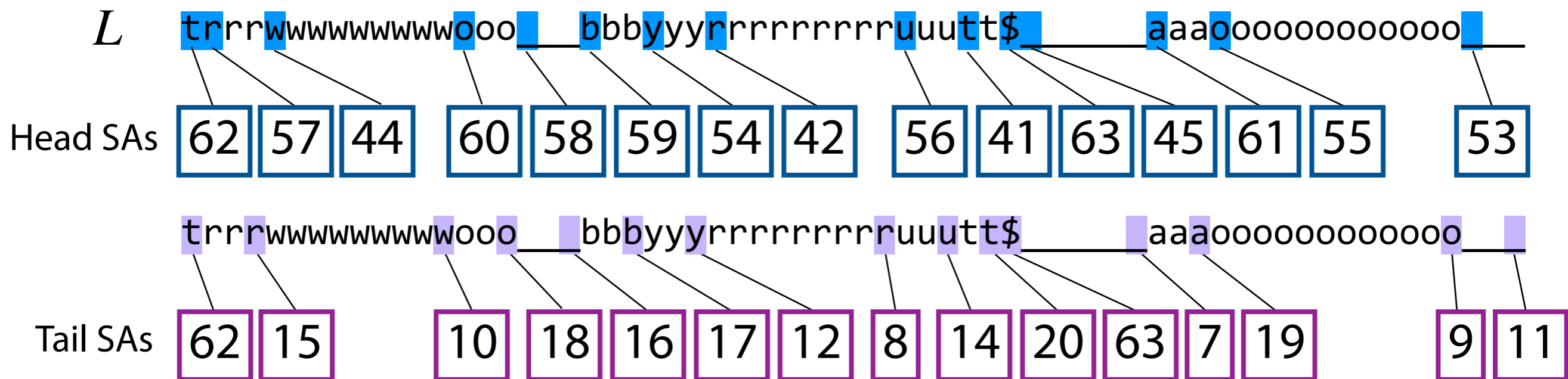
Case 2b: Some elements in $\text{BWT}[sp^{(i)} .. ep^{(i)}]$ match y ,
but $\text{BWT}[j^{(i)}]$ doesn't



Let $j^{(i+1)}$ be the offset of a y -head or y -tail. Look up $SA[j^{(i+1)}]$ in our structure holding SA items at heads/tails.

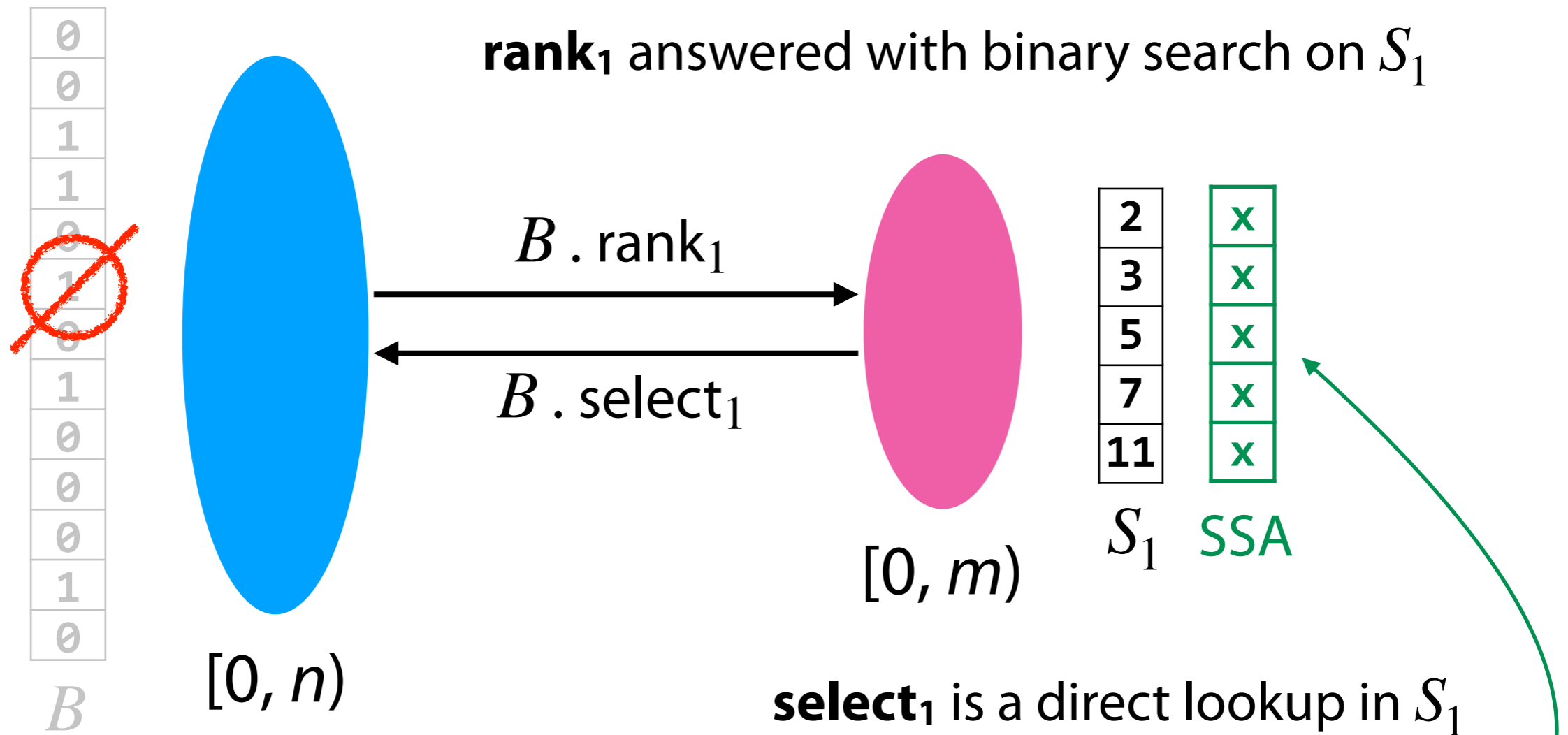
Toehold predecessor structure

How exactly to look up $SA[j^{(i+1)}]$ in our structure holding SA items at heads/tails?



Recall our discussion of the **predecessor structure** from the last video...

Run-length FM Index



Associating values with set bits is as easy as adding **more arrays** parallel to S_1

Toehold predecessor structure

Similarly, we can add associativity to the structure we referenced here:

Run-length FM Index

Predecessor queries on B

$O(\log \log n/r)$ time, $O(r)$ space
Belazzougui & Navarro, Thm A.1

Rank queries on S

$O(\log \log \sigma)$ time, $O(r)$ space
Belazzougui & Navarro, Thm 5.4

Select queries on C & B'

$O(1)$ time, $O(r)$ space
Array of cumulative counts

RLFM index:

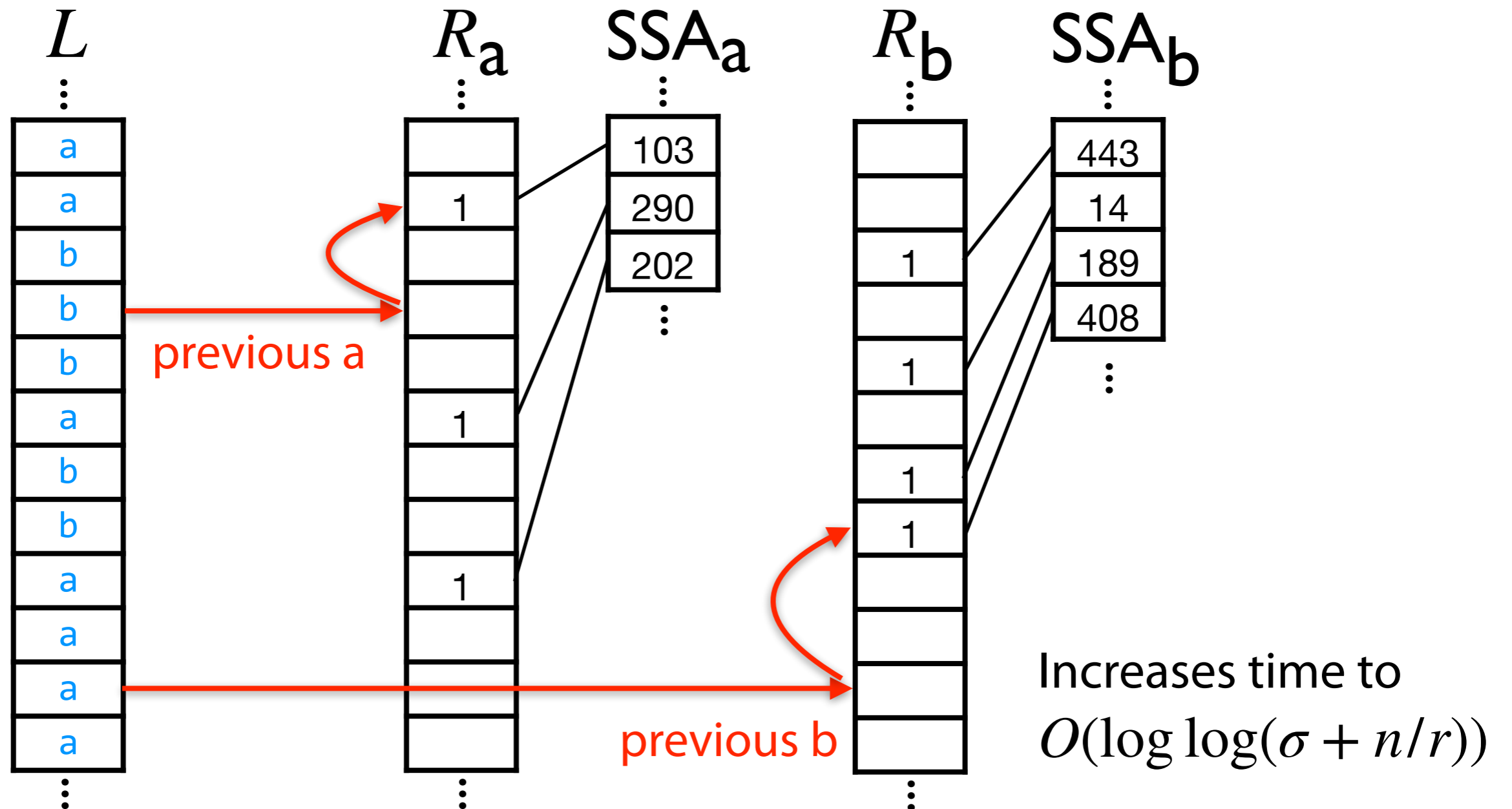
```
 $r \leftarrow B . \text{rank}_1(\text{off} + 1)$   
 $o \leftarrow B . \text{select}_1(r)$   
 $pv \leftarrow S . \text{rank}_c(r - 1)$   
 $s_1 \leftarrow C[c]$   
 $s_2 \leftarrow B'[c] . \text{select}_1(pv)$   
 $\text{next\_off} \leftarrow s_1 + s_2 + (\text{off} - o)$ 
```

Gagie T, Navarro G, and Prezza P. Optimal-time text indexing in BWT-runs bounded space. Proceedings of 29th SODA, ACM-SIAM. 2018. pp1459—1477.

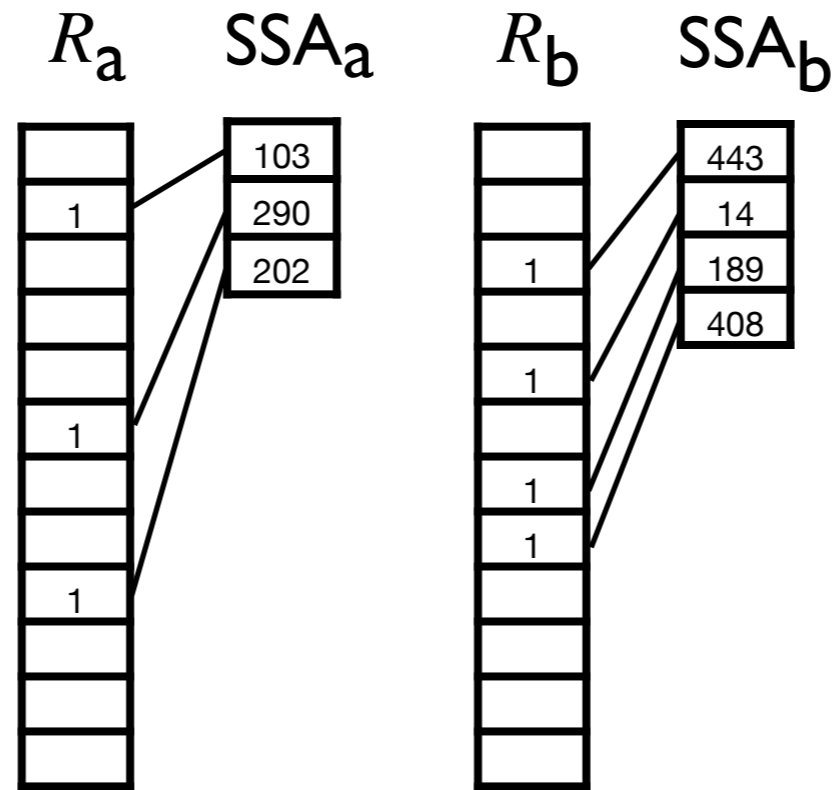
Belazzougui, Djamel, and Gonzalo Navarro. "Optimal lower and upper bounds for representing sequences." ACM Transactions on Algorithms (TALG) 11.4 (2015): 1-21.

Toehold predecessor structure

Final point: structure needs to be stratified by character



Toehold predecessor structure



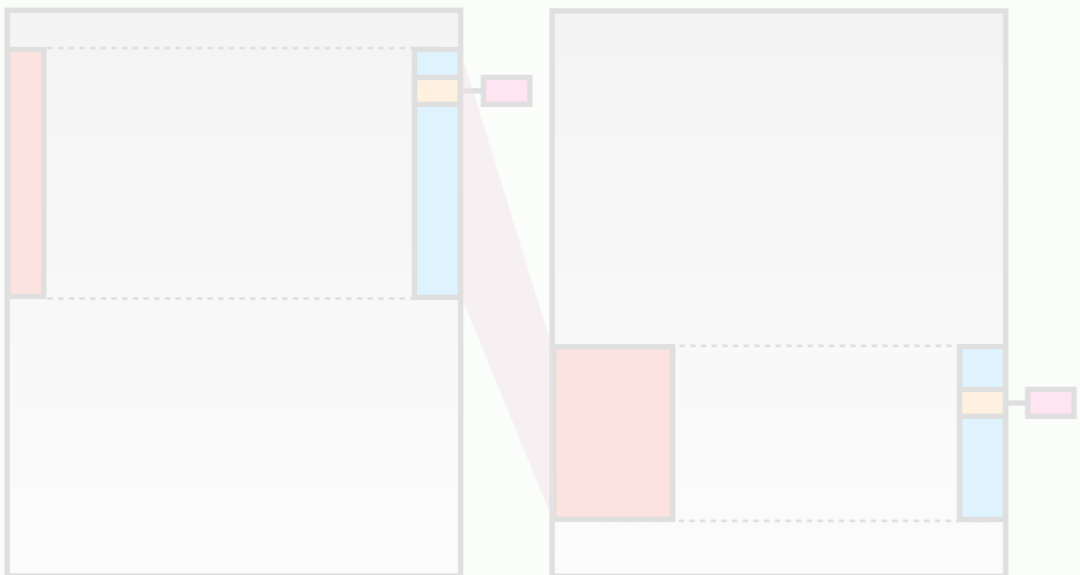
L trrrwwwwwwwooo__bbbyyrrrrrrrrruutt\$_____aaaoooooooooooo__

\$	39					
_	16	18	40	45	61	63
a	46	48				
b	19	21				
o	13	15	49	60		

r	1	3	25	33
t	0	37	38	
u	34	36		
w	4	12		
y	22	24		

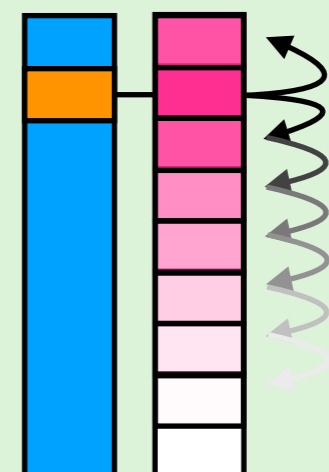
Next video

1



We always have a "toehold"

2



We can get the rest from the toehold