

# Exact Matching, Part 2

photophorescent



Kyle Gorman (filling in for Steven Bedrick)  
CS/EE 5/655, 11/17/14

# Plan for today:

Z-algorithm review

Knuth-Morris-Pratt

Boyer-Moore

Definitions: for a string  $S$ :

$S[i,j]$  = contiguous substring starting at  $i$  and ending at  $j$ .  $S(i) = S[i,i]$

Definitions: for a string  $S$ :

$S[i,j]$  = contiguous substring starting at  $i$  and ending at  $j$ .  $S(i) = S[i,i]$

$S = \text{aardvark}$

Definitions: for a string  $S$ :

$S[i,j]$  = contiguous substring starting at  $i$  and ending at  $j$ .  $S(i) = S[i,i]$

$S = \text{aardvark}$

$S[2,4] = \text{ard}$

Definitions: for a string  $S$ :

$S[i,j]$  = contiguous substring starting at  $i$  and ending at  $j$ .  $S(i) = S[i,i]$

$S = \text{aardvark}$

$S[2,4] = \text{ard}$        $S(4) = \text{d}$

Definitions: for a string  $S$ :

$S[i,j]$  = contiguous substring starting at  $i$  and ending at  $j$ .  $S(i) = S[i,i]$

$S = \text{aardvark}$

$S[2,4] = \text{ard}$        $S(4) = \text{d}$

For  $i > 1$ ,  $Z_i(S)$  is the length of the longest prefix of  $S[i,|S|]$  that is also a prefix of  $S$ .

For  $i > 1$ ,  $Z_i(S)$  is the length of the longest prefix of  $S[i, |S|]$  that is also a prefix of  $S$ .



For  $i > 1$ ,  $Z_i(S)$  is the length of the longest prefix of  $S[i, |S|]$  that is also a prefix of  $S$ .

$S = \text{xtpxtd}$

For  $i > 1$ ,  $Z_i(S)$  is the length of the longest prefix of  $S[i, |S|]$  that is also a prefix of  $S$ .

$S = \text{xtpxtd}$

$S[4, |S|] = \text{xtd}$

$\text{xtp}\underline{\text{xtd}}$

For  $i > 1$ ,  $Z_i(S)$  is the length of the longest prefix of  $S[i, |S|]$  that is also a prefix of  $S$ .

$S = \text{xtpxtd}$

$S[4, |S|] = \text{xtd}$        $\text{xtp}\underline{\text{xtd}}$

$Z_4(S) = 2$        $\underline{\text{xtp}}\underline{\text{xtd}}$

S = aardvark

S = aardvark

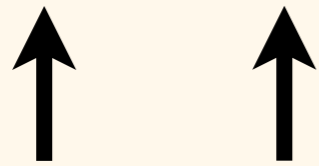


$S = \underline{a}ardv\underline{a}rk$



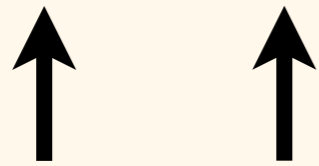
$$Z_6(S) = 1$$

$S = \underline{a}ardv\underline{a}rk$



$$Z_6(S) = 1$$

$S = \underline{a}ardv\underline{a}rk$

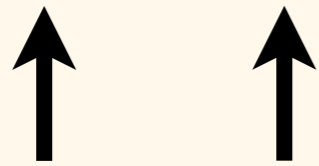


$$Z_6(S) = 1$$

$$Z_2(S) = 1$$



$S = \underline{a}ardv\underline{a}rk$



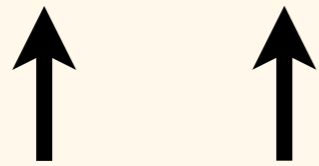
$S = \underline{alf}\underline{alfa}$



$$Z_6(S) = 1$$

$$Z_2(S) = 1$$

$S = \underline{a}ardv\underline{a}rk$



$$Z_6(S) = 1$$

$$Z_2(S) = 1$$

$S = \underline{alf}alfa$



$$Z_4(S) = 4$$

$S = \underline{a}ardv\underline{a}rk$



$$Z_6(S) = 1$$

$$Z_2(S) = 1$$

$S = \underline{a}l\underline{f}a\underline{l}f\underline{a}$



$$Z_4(S) = 4$$

$S = \underline{p}h\underline{o}t\underline{o}p\underline{h}o\underline{s}p\underline{h}o\underline{r}e\underline{s}c\underline{e}n\underline{t}$

$S = \underline{a}ardv\underline{a}rk$



$$Z_6(S) = 1$$

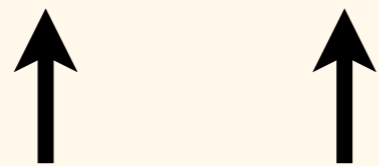
$$Z_2(S) = 1$$

$S = \underline{a}l\underline{f}a\underline{l}f\underline{a}$



$$Z_4(S) = 4$$

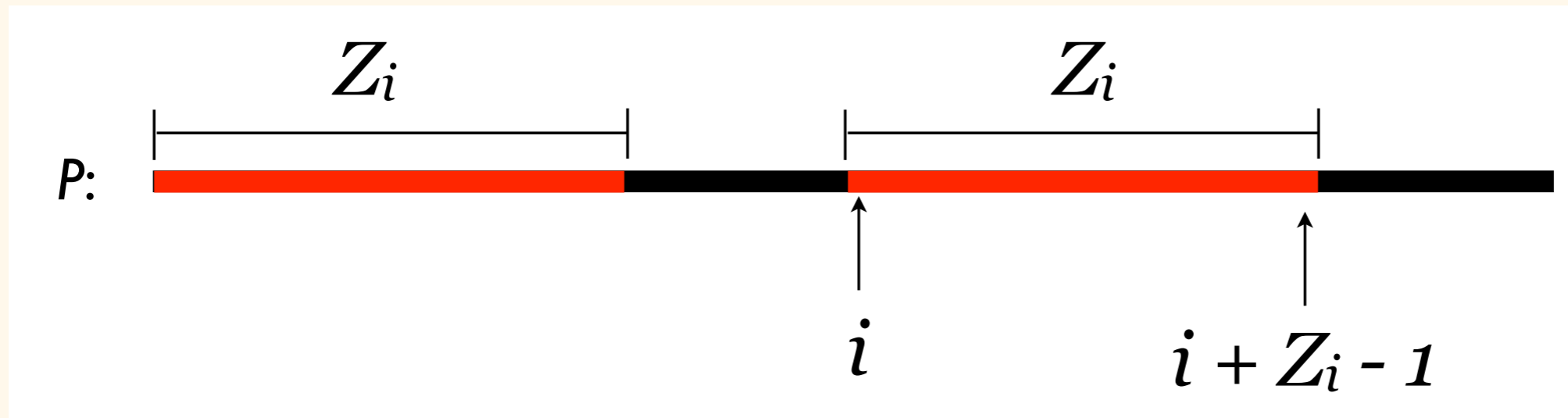
$S = \underline{p}h\underline{o}t\underline{o}p\underline{h}o\underline{s}p\underline{h}o\underline{r}e\underline{s}c\underline{e}n\underline{t}$



$$Z_6(S) = Z_{10}(S) = 3$$

These regions of prefix-overlap are called *z-boxes*.

P = photophosphorescent



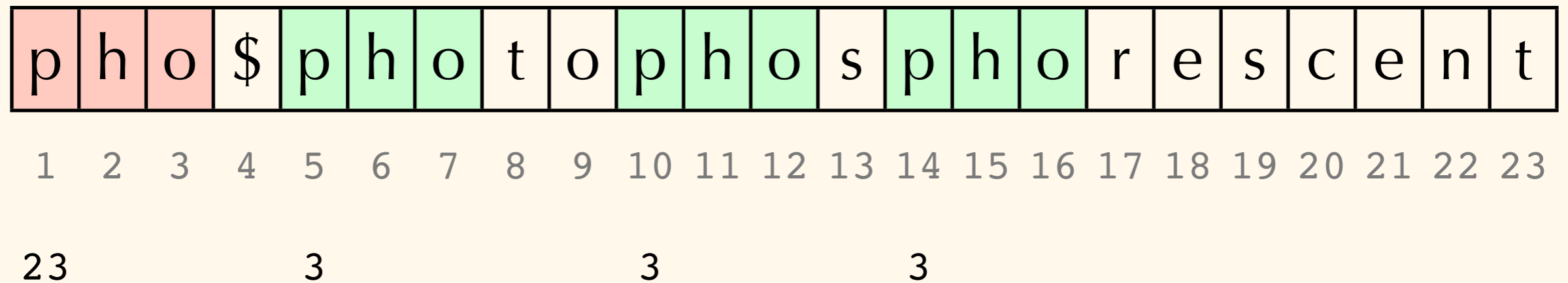
By being clever about constructing our string,  
we can easily find exact pattern matches:

p	h	o	\$	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e	n	t
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23
23				3					3				3									

Question: How to calculate  $Z_i$ ?

By being clever about constructing our string, we can easily find exact pattern matches:

Make the pattern ( $P$ ) the prefix...



Question: How to calculate  $Z_i$ ?

By being clever about constructing our string,  
we can easily find exact pattern matches:

Make the pattern ( $P$ ) the prefix...



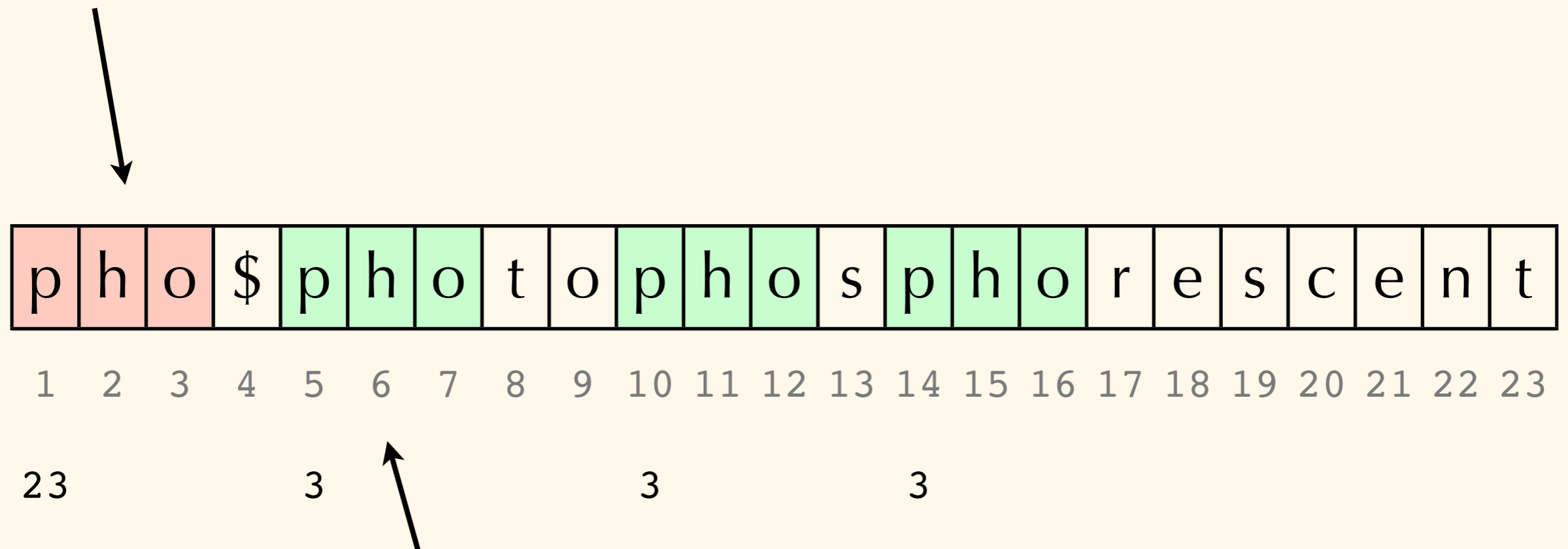
p	h	o	\$	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e	n	t
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23
23				3					3				3									

Question: How to calculate  $Z_i$ ?



By being clever about constructing our string, we can easily find exact pattern matches:

Make the pattern ( $P$ ) the prefix...



... and now, any occurrence of  $P$  is a repeat of the string's prefix, and so has  $Z_i = |P|$

Question: How to calculate  $Z_i$ ?

Question: How to calculate  $Z_i$ ?

Question: How to calculate  $Z_i$ ?

The naïve way:

Question: How to calculate  $Z_i$ ?

The naïve way:

For every position  $i$ , compute the longest common prefix between  $S$  and  $S[i, |S|]$

Question: How to calculate  $Z_i$ ?

The naïve way:

For every position  $i$ , compute the longest common prefix between  $S$  and  $S[i, |S|]$

Problem: This is  $O(n^2)$ !

Question: How to calculate  $Z_i$ ?

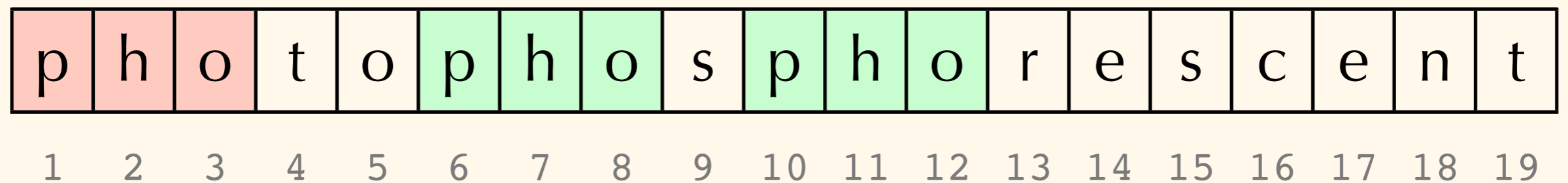
The naïve way:

For every position  $i$ , compute the longest common prefix between  $S$  and  $S[i, |S|]$

Problem: This is  $O(n^2)$ !

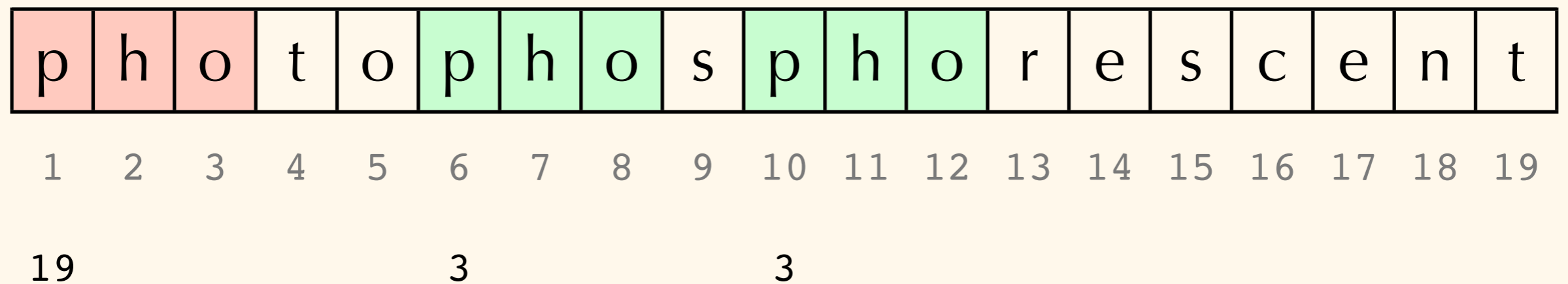
The solution involves thinking about the properties of Z-boxes.

At any given  $i$  with  $Z_i > 0$ , we can use  $Z_i$  to determine the length of its enclosing Z-box:



Also, we know that any position  $k$  inside a Z-box must correspond to a position  $k'$  somewhere in the prefix of  $S$ .

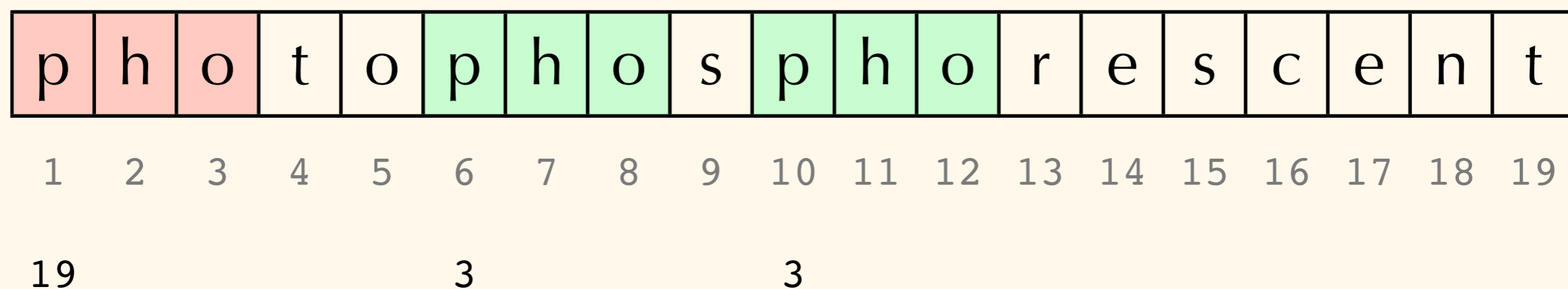
At any given  $i$  with  $Z_i > 0$ , we can use  $Z_i$  to determine the length of its enclosing Z-box:



Also, we know that any position  $k$  inside a Z-box must correspond to a position  $k'$  somewhere in the prefix of  $S$ .



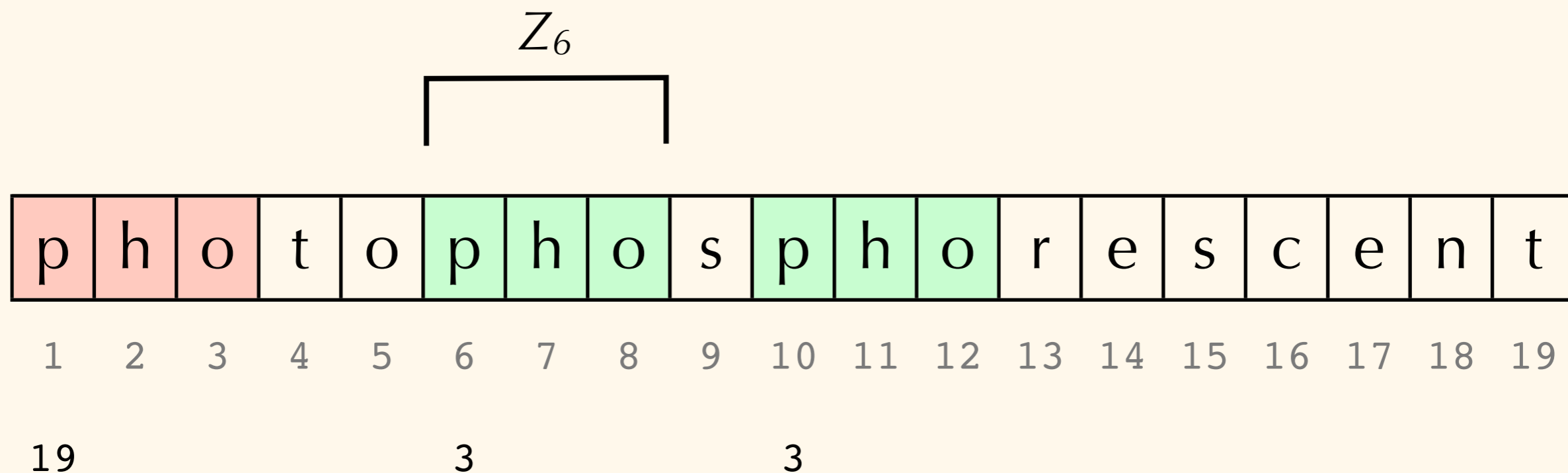
At any given  $i$  with  $Z_i > 0$ , we can use  $Z_i$  to determine the length of its enclosing Z-box:



For  $i = 6$ :  $Z_6(S) = 3$

Also, we know that any position  $k$  inside a Z-box must correspond to a position  $k'$  somewhere in the prefix of  $S$ .

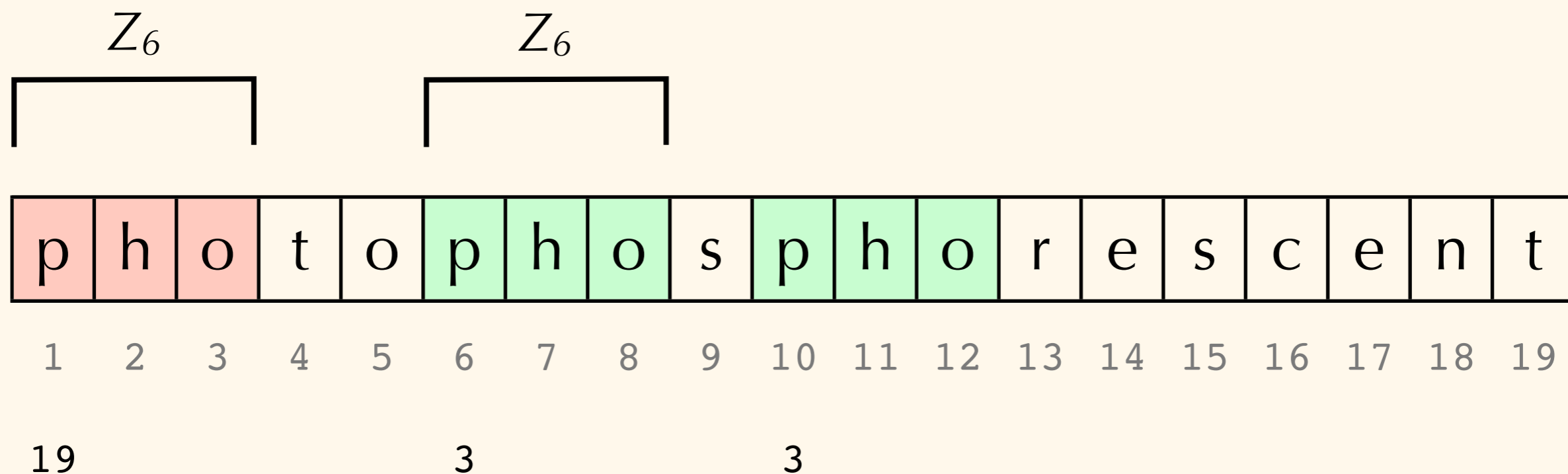
At any given  $i$  with  $Z_i > 0$ , we can use  $Z_i$  to determine the length of its enclosing Z-box:



For  $i = 6$ :  $Z_6(S) = 3$

Also, we know that any position  $k$  inside a Z-box must correspond to a position  $k'$  somewhere in the prefix of  $S$ .

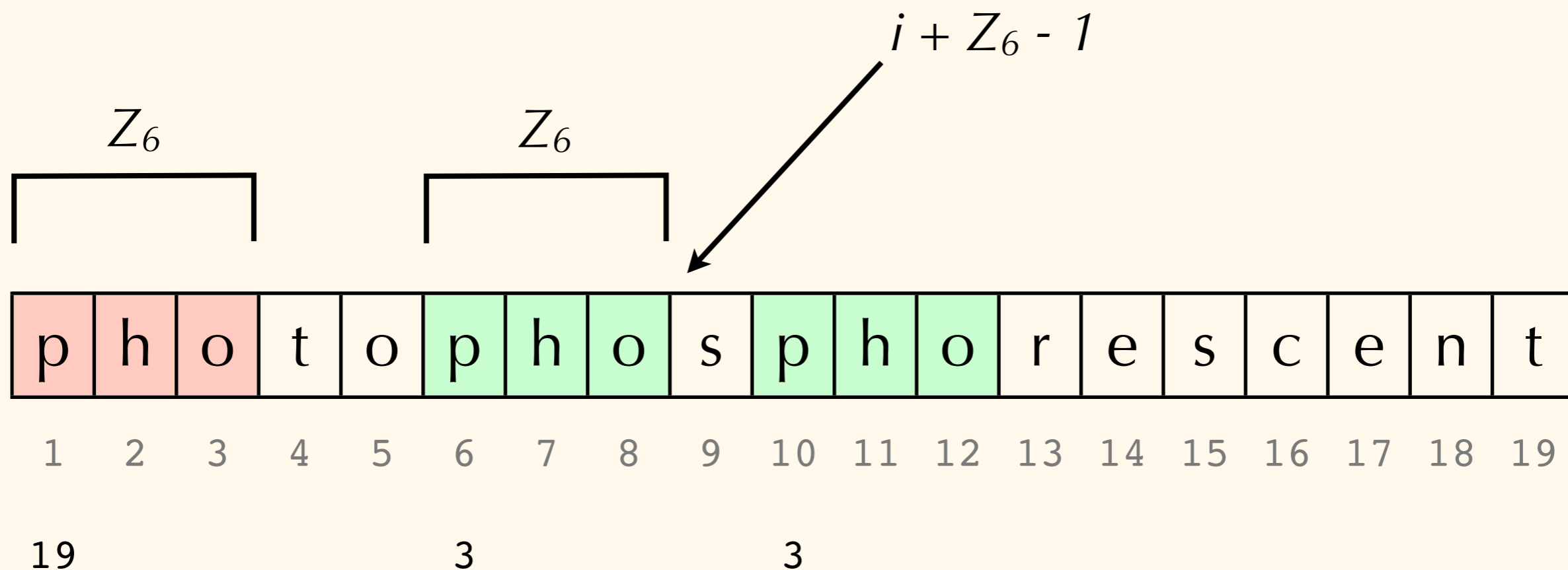
At any given  $i$  with  $Z_i > 0$ , we can use  $Z_i$  to determine the length of its enclosing Z-box:



For  $i = 6$ :  $Z_6(S) = 3$

Also, we know that any position  $k$  inside a Z-box must correspond to a position  $k'$  somewhere in the prefix of  $S$ .

At any given  $i$  with  $Z_i > 0$ , we can use  $Z_i$  to determine the length of its enclosing Z-box:



For  $i = 6$ :  $Z_6(S) = 3$

Also, we know that any position  $k$  inside a Z-box must correspond to a position  $k'$  somewhere in the prefix of  $S$ .

# Additional definitions:

$$r_i = \max_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

$$l_i = \operatorname{argmax}_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$



# Additional definitions:

$$r_i = \max_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

$$l_i = \operatorname{argmax}_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$



$i$

$Z_i$

$r_i$

$l_i$

# Additional definitions:

$$r_i = \max_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

$$l_i = \operatorname{argmax}_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$



$i$  1

$Z_i$  19

$r_i$

$l_i$



# Additional definitions:

$$r_i = \max_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

$$l_i = \operatorname{argmax}_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$



$i$  1 2

$Z_i$  19

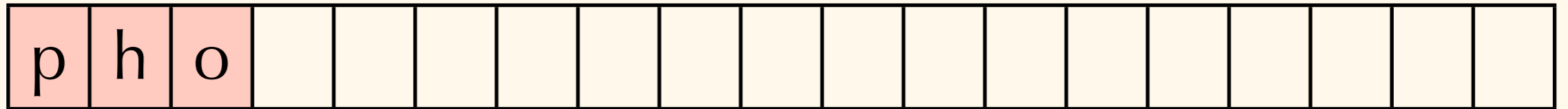
$r_i$  0

$l_i$  0

# Additional definitions:

$$r_i = \max_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

$$l_i = \operatorname{argmax}_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$



$i$  1 2 3

$Z_i$  19

$r_i$  0 0

$l_i$  0 0

# Additional definitions:

$$r_i = \max_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

$$l_i = \operatorname{argmax}_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

p	h	o	t															
---	---	---	---	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

$i$  1 2 3 4

$Z_i$  19

$r_i$  0 0 0

$l_i$  0 0 0

# Additional definitions:

$$r_i = \max_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

$$l_i = \operatorname{argmax}_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

p	h	o	t	o															
---	---	---	---	---	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

$i$     1    2    3    4    5

$Z_i$  19

$r_i$     0    0    0    0

$l_i$     0    0    0    0

# Additional definitions:

$$r_i = \max_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

$$l_i = \operatorname{argmax}_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

p	h	o	t	o	p													
---	---	---	---	---	---	--	--	--	--	--	--	--	--	--	--	--	--	--

$i$  1 2 3 4 5 6

$Z_i$  19 3

$r_i$  0 0 0 0 8

$l_i$  0 0 0 0 6





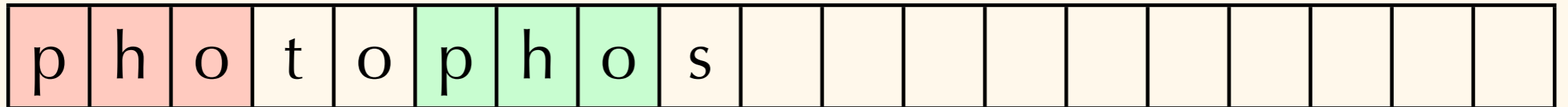




# Additional definitions:

$$r_i = \max_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

$$l_i = \operatorname{argmax}_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$



$i$  1 2 3 4 5 6 7 8 9

$Z_i$  19 3

$r_i$  0 0 0 0 8 8 8 8

$l_i$  0 0 0 0 6 6 6 6













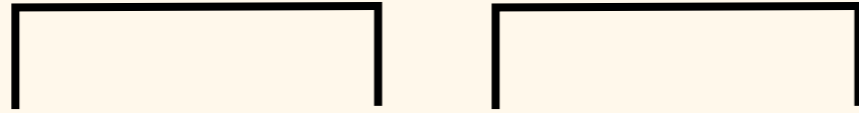




# Additional definitions:

$$r_i = \max_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

$$l_i = \operatorname{argmax}_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$



p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c			
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	--	--	--

$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
-----	---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----

$Z_i$	19					3				3						
-------	----	--	--	--	--	---	--	--	--	---	--	--	--	--	--	--

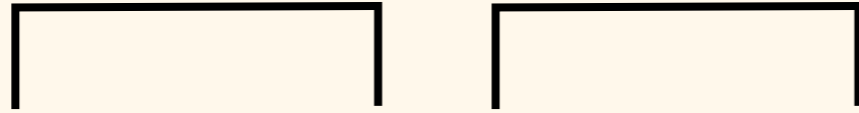
$r_i$		0	0	0	0	8	8	8	8	12	12	12	12	12	12	12
-------	--	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----

$l_i$		0	0	0	0	6	6	6	6	10	10	10	10	10	10	10
-------	--	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----

# Additional definitions:

$$r_i = \max_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$

$$l_i = \operatorname{argmax}_{\substack{1 < j \leq i \\ \text{s.t. } Z_j > 0}} (j + Z_j - 1)$$



p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e		
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	--	--

$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
-----	---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----

$Z_i$	19					3				3							
-------	----	--	--	--	--	---	--	--	--	---	--	--	--	--	--	--	--

$r_i$		0	0	0	0	8	8	8	8	12	12	12	12	12	12	12	12
-------	--	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----

$l_i$		0	0	0	0	6	6	6	6	10	10	10	10	10	10	10	10
-------	--	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----





What does this get us?

p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e	n	t
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

$i$  1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19

$Z_i$  19 3 3

$r_i$  0 0 0 0 8 8 8 8 12 12 12 12 12 12 12 12 12 12 12

$l_i$  0 0 0 0 6 6 6 6 10 10 10 10 10 10 10 10 10 10 10

For any position  $k$  where  $k > l_k$  (i.e., is inside a z-box), we know that there exists a position  $k'$  s.t.  $S(k) = S(k')$ .  $k' = k - l_k + 1$

E.g.:  $k = 7$

What does this get us?

	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e	n	t
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
$Z_i$	19					3				3									
$r_i$		0	0	0	0	8	8	8	8	12	12	12	12	12	12	12	12	12	12
$l_i$		0	0	0	0	6	6	6	6	10	10	10	10	10	10	10	10	10	10

For any position  $k$  where  $k > l_k$  (i.e., is inside a z-box), we know that there exists a position  $k'$  s.t.  $S(k) = S(k')$ .  $k' = k - l_k + 1$

E.g.:  $k = 7$      $S(k) = \text{"h"}$

What does this get us?

	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e	n	t
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
$Z_i$	19					3				3									
$r_i$		0	0	0	0	8	8	8	8	12	12	12	12	12	12	12	12	12	12
$l_i$		0	0	0	0	6	6	6	6	10	10	10	10	10	10	10	10	10	10

For any position  $k$  where  $k > l_k$  (i.e., is inside a z-box), we know that there exists a position  $k'$  s.t.  $S(k) = S(k')$ .  $k' = k - l_k + 1$

E.g.:  $k = 7$      $S(k) = \text{"h"}$      $l_k = 6$

What does this get us?

	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e	n	t
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
$Z_i$	19					3				3									
$r_i$		0	0	0	0	8	8	8	8	12	12	12	12	12	12	12	12	12	12
$l_i$		0	0	0	0	6	6	6	6	10	10	10	10	10	10	10	10	10	10

For any position  $k$  where  $k > l_k$  (i.e., is inside a z-box), we know that there exists a position  $k'$  s.t.  $S(k) = S(k')$ .  $k' = k - l_k + 1$

E.g.:  $k = 7$      $S(k) = \text{"h"}$      $l_k = 6$      $k' = 7 - 6 + 1 = 2$



What does this get us?

	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e	n	t
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
$Z_i$	19					3				3									
$r_i$		0	0	0	0	8	8	8	8	12	12	12	12	12	12	12	12	12	12
$l_i$		0	0	0	0	6	6	6	6	10	10	10	10	10	10	10	10	10	10

For any position  $k$  where  $k > l_k$  (i.e., is inside a z-box), we know that there exists a position  $k'$  s.t.  $S(k) = S(k')$ .  $k' = k - l_k + 1$

E.g.:  $k = 7$      $S(k) = \text{"h"}$      $l_k = 6$      $k' = 7 - 6 + 1 = 2$

$k'$  is the  $k$ 's equivalent position in the matching prefix of the string!

What does this get us?

	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e	n	t
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
$Z_i$	19					3				3									
$r_i$		0	0	0	0	8	8	8	8	12	12	12	12	12	12	12	12	12	12
$l_i$		0	0	0	0	6	6	6	6	10	10	10	10	10	10	10	10	10	10

$k'$  is the  $k$ 's equivalent position in the matching prefix of the string!

What does this get us?

	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e	n	t
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
$Z_i$	19					3				3									
$r_i$		0	0	0	0	8	8	8	8	12	12	12	12	12	12	12	12	12	12
$l_i$		0	0	0	0	6	6	6	6	10	10	10	10	10	10	10	10	10	10

$k'$  is the  $k$ 's equivalent position in the matching prefix of the string!

Therefore,  $Z_{k'}$  can tell us something about the structure of the prefix of the string.

What does this get us?

	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e	n	t
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
$Z_i$	19					3				3									
$r_i$		0	0	0	0	8	8	8	8	12	12	12	12	12	12	12	12	12	12
$l_i$		0	0	0	0	6	6	6	6	10	10	10	10	10	10	10	10	10	10

$k'$  is the  $k$ 's equivalent position in the matching prefix of the string!

Therefore,  $Z_{k'}$  can tell us something about the structure of the prefix of the string.

If  $Z_{k'} > 0$ , there must be repeating elements!

Putting it all together into an algorithm:

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :


Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:

Putting it all together into an algorithm:


Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*



Putting it all together into an algorithm:


Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :


If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way

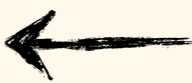
If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :


If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way

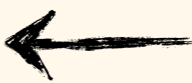
If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*


Calculate  $Z_k$  the normal way

If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

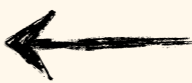
If  $k \leq r$ , we are inside of an already-found z-box, so:

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way


If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

If  $k \leq r$ , we are inside of an already-found z-box, so:

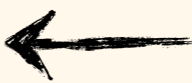
$$k' = k - l + 1 ; \beta = |S[k, r]|$$

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way

If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k, r]$$

We know that  $\beta$  must match  $S[k', Z_l]$ ...

If  $k \leq r$ , we *are* inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

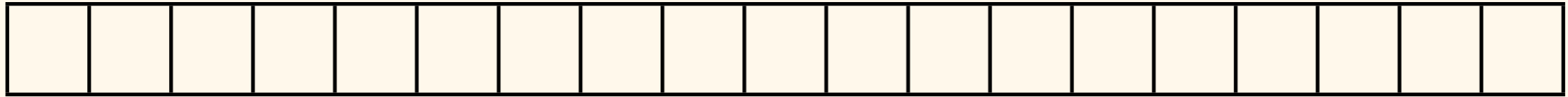
We know that  $\beta$  must match  $S[k',Z_l]...$

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we *are* inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



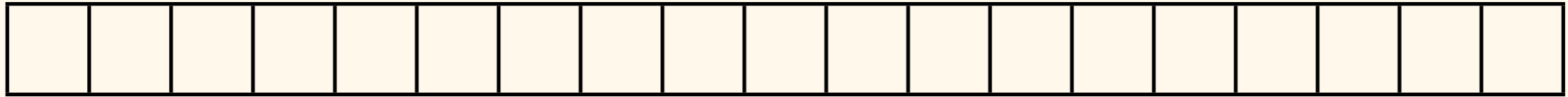
$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$



If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



$i$

$Z_i$

$r_i$

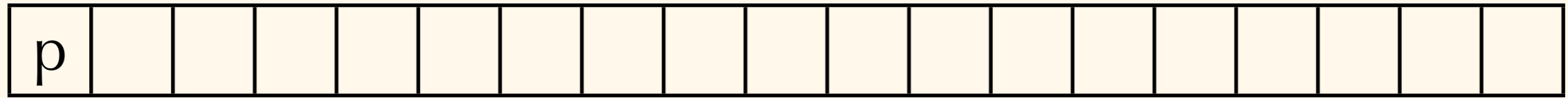
$l_i$

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



$i$  1

$Z_i$  19

$r_i$

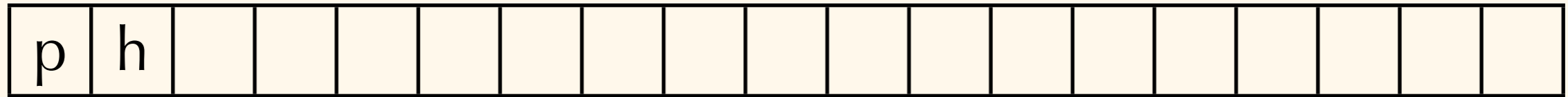
$l_i$

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



$i$     1    2

$Z_i$  19

$r_i$     0

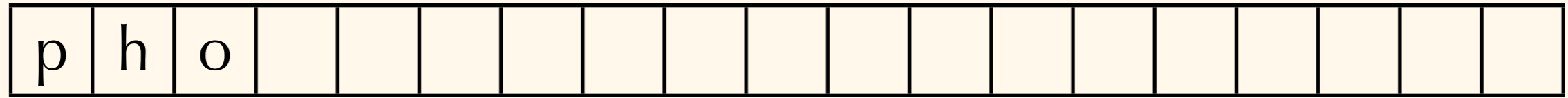
$l_i$     0

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



$i$     1    2    3

$Z_i$  19

$r_i$     0    0

$l_i$     0    0

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

p	h	o	t															
---	---	---	---	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

$i$     1    2    3    4

$Z_i$  19

$r_i$     0    0    0

$l_i$     0    0    0

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o	t	o														
$i$	1	2	3	4	5														
$Z_i$	19																		
$r_i$		0	0	0	0														
$l_i$		0	0	0	0														

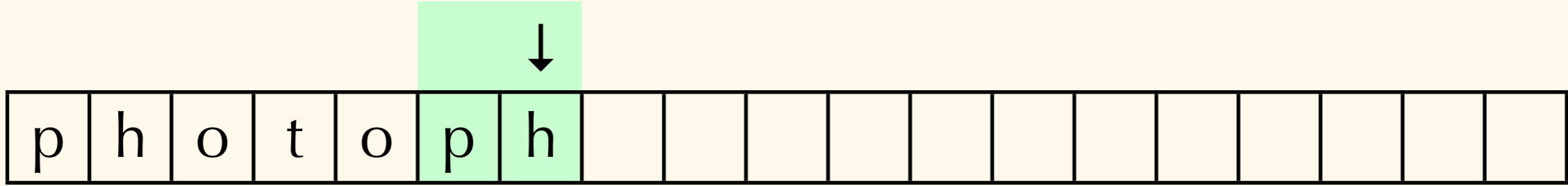
$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$



If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



$i$	1	2	3	4	5	6	7											
$Z_i$	19					3												
$r_i$		0	0	0	0	8												
$l_i$		0	0	0	0	6												

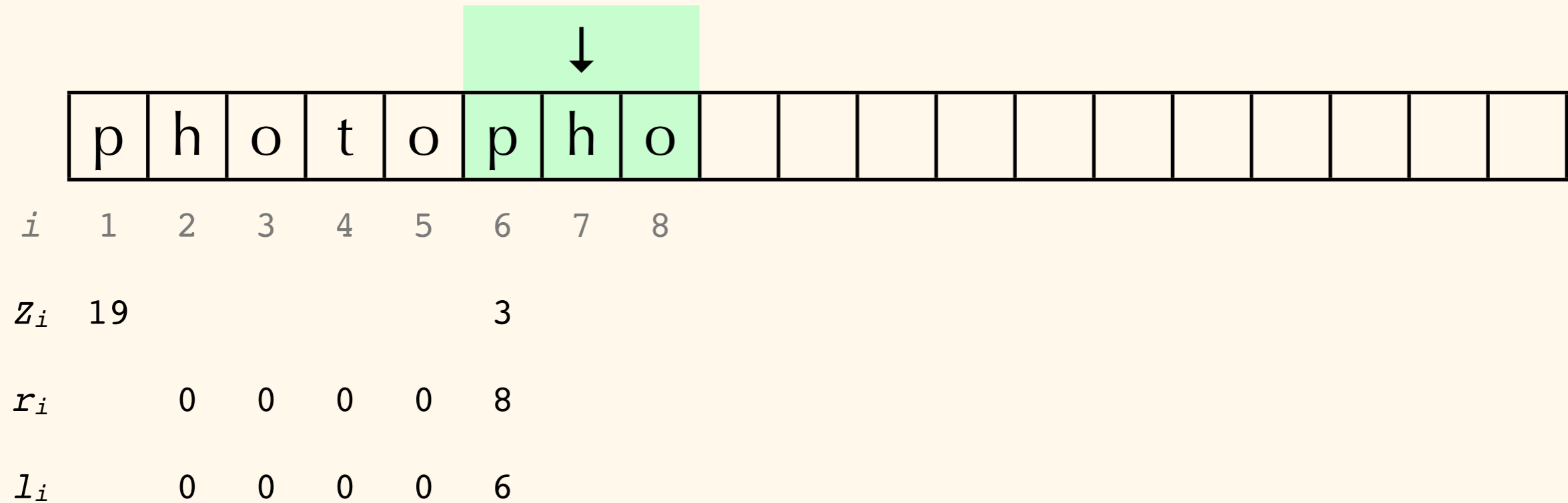
$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$



If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



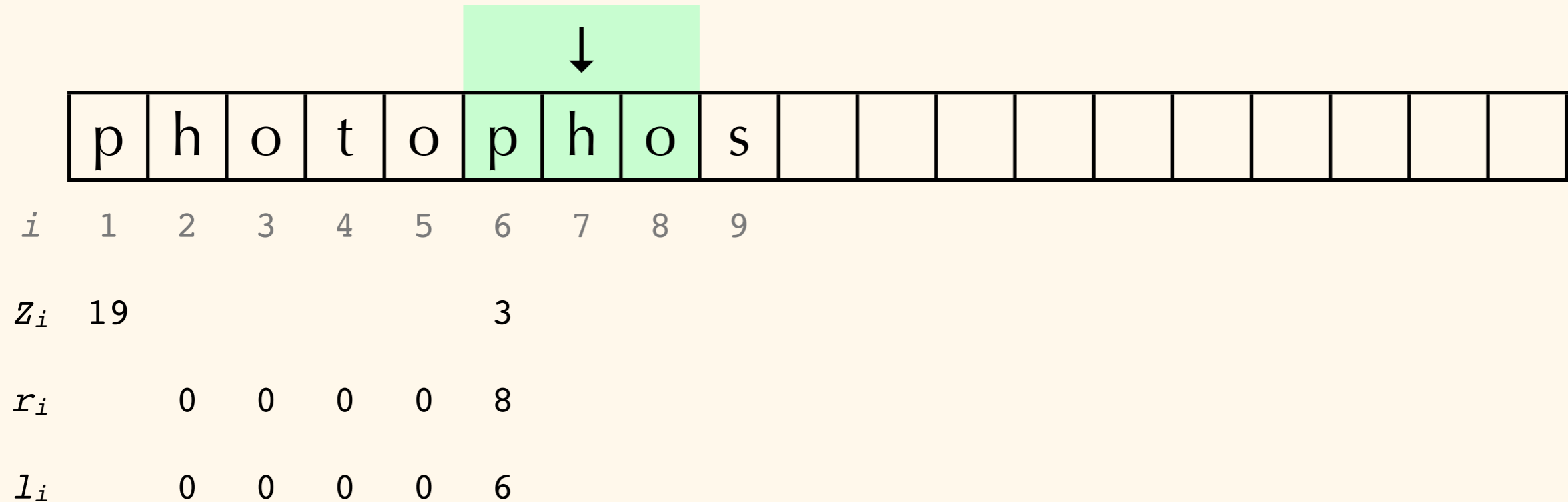
	p	h	o	t	o	p	h	o											
$i$	1	2	3	4	5	6	7	8											
$Z_i$	19					3													
$r_i$		0	0	0	0	8													
$l_i$		0	0	0	0	6													

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



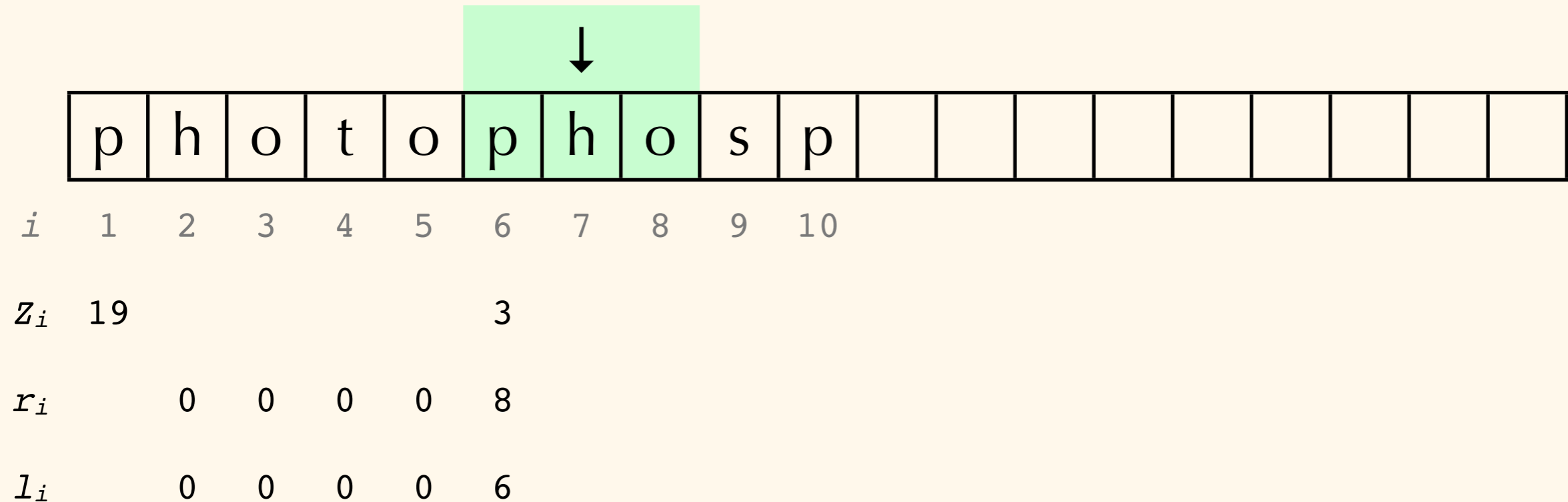
<i>i</i>	1	2	3	4	5	6	7	8	9										
<i>Z<sub>i</sub></i>	19					3													
<i>r<sub>i</sub></i>		0	0	0	0	8													
<i>l<sub>i</sub></i>		0	0	0	0	6													

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



	p	h	o	t	o	p	h	o	s	p									
$i$	1	2	3	4	5	6	7	8	9	10									
$Z_i$	19					3													
$r_i$		0	0	0	0	8													
$l_i$		0	0	0	0	6													

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o	t	o	p	h	o	s	p	h									
$i$	1	2	3	4	5	6	7	8	9	10	11									
$Z_i$	19					3														
$r_i$		0	0	0	0	8														
$l_i$		0	0	0	0	6														

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o	t	o	p	h	o	s	p	h	o							
$i$	1	2	3	4	5	6	7	8	9	10	11	12							
$Z_i$	19					3													
$r_i$		0	0	0	0	8													
$l_i$		0	0	0	0	6													

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$









If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c			
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16			
$Z_i$	19					3													
$r_i$		0	0	0	0	8													
$l_i$		0	0	0	0	6													

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e		
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17		
$Z_i$	19					3													
$r_i$		0	0	0	0	8													
$l_i$		0	0	0	0	6													

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e	n	
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	
$Z_i$	19					3													
$r_i$		0	0	0	0	8													
$l_i$		0	0	0	0	6													

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$



If  $k \leq r$ , we *are* inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

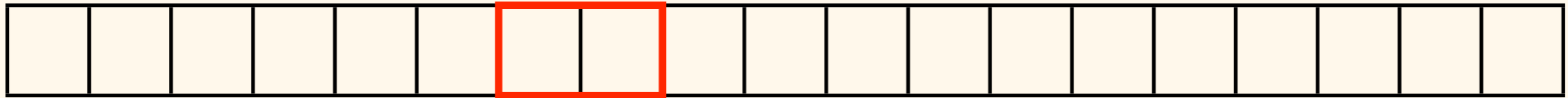
We know that  $\beta$  must match  $S[k',Z_l]...$

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

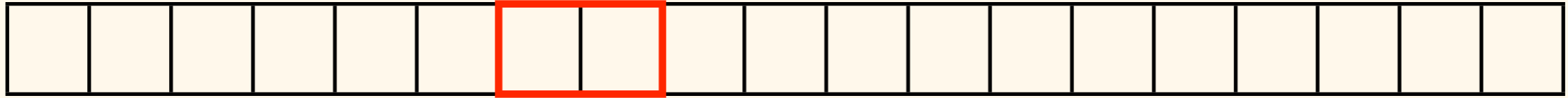


$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



$i$

$Z_i$

$r_i$

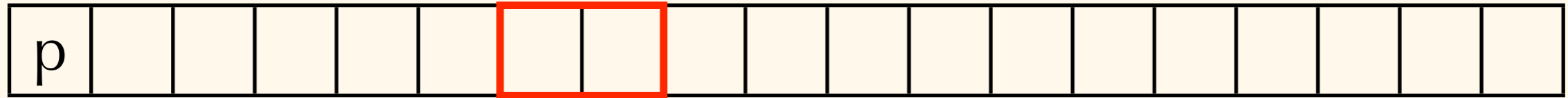
$l_i$

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



$i$  1

$Z_i$  19

$r_i$

$l_i$

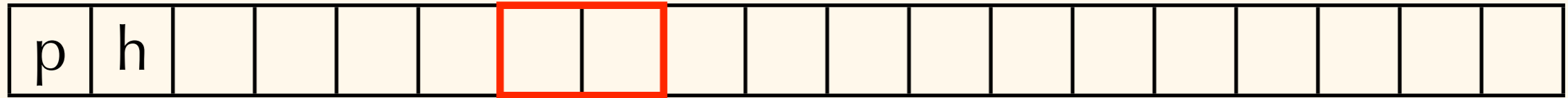
$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$



If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



$i$     1    2

$Z_i$  19

$r_i$     0

$l_i$     0

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o															
$i$	1	2	3															
$Z_i$	19																	
$r_i$		0	0															
$l_i$		0	0															

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o	t															
$i$	1	2	3	4															
$Z_i$	19																		
$r_i$		0	0	0															
$l_i$		0	0	0															

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o	t	o														
$i$	1	2	3	4	5														
$Z_i$	19																		
$r_i$		0	0	0	0														
$l_i$		0	0	0	0														

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

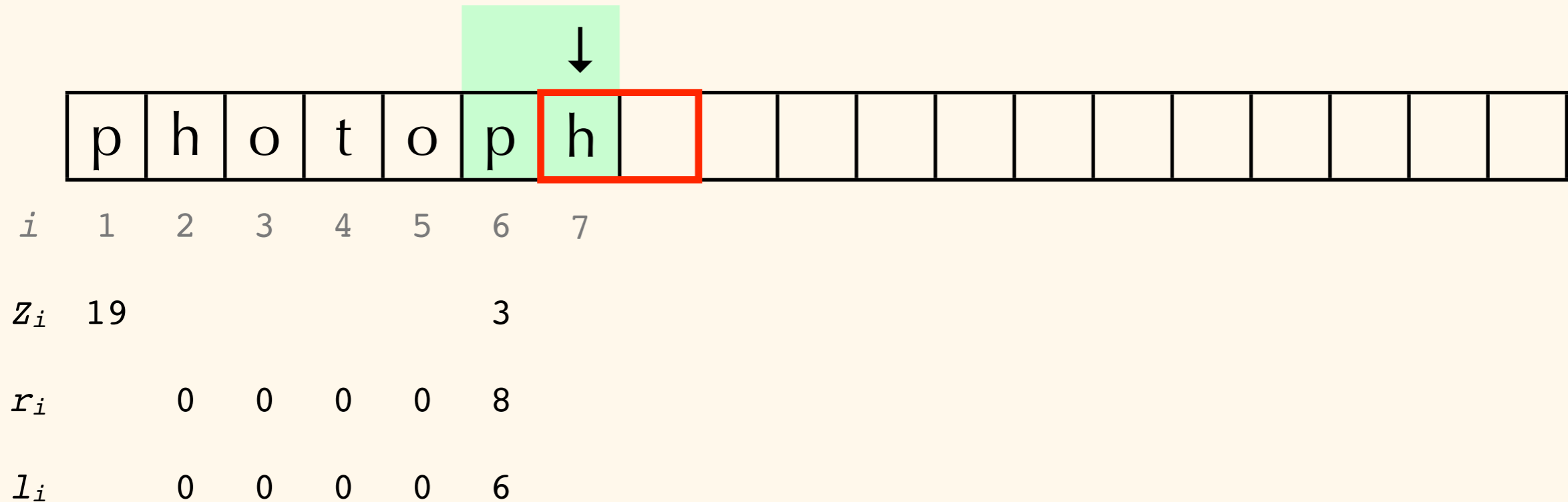
	p	h	o	t	o	p													
$i$	1	2	3	4	5	6													
$Z_i$	19					3													
$r_i$		0	0	0	0	8													
$l_i$		0	0	0	0	6													

$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

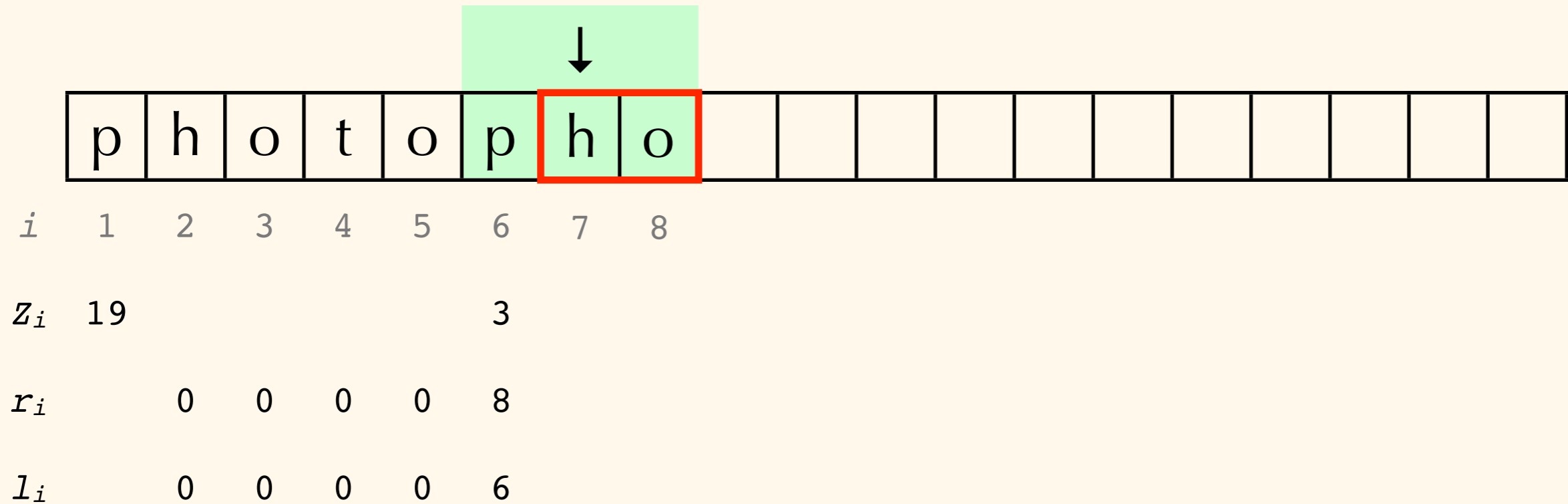


$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

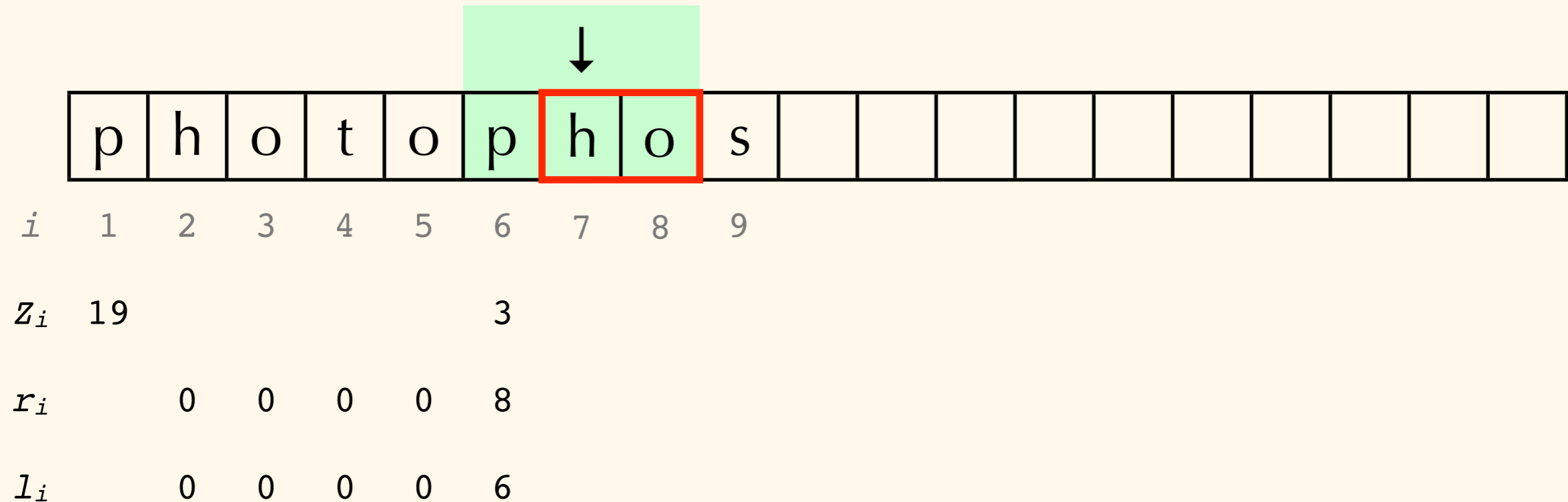


$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



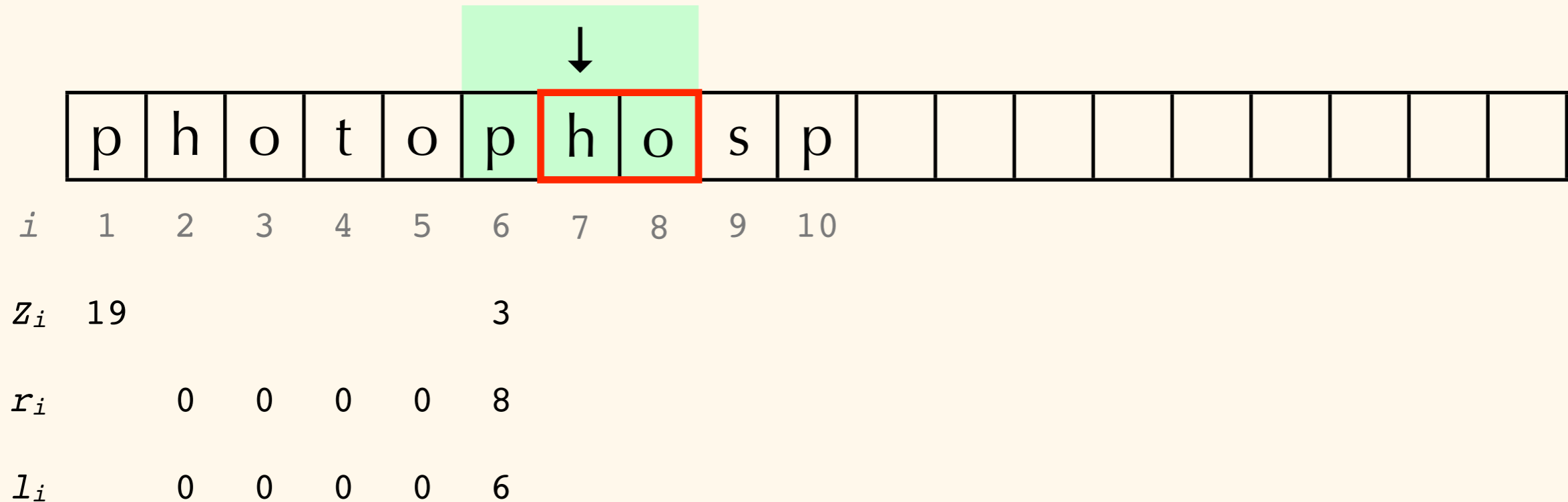
$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$



If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

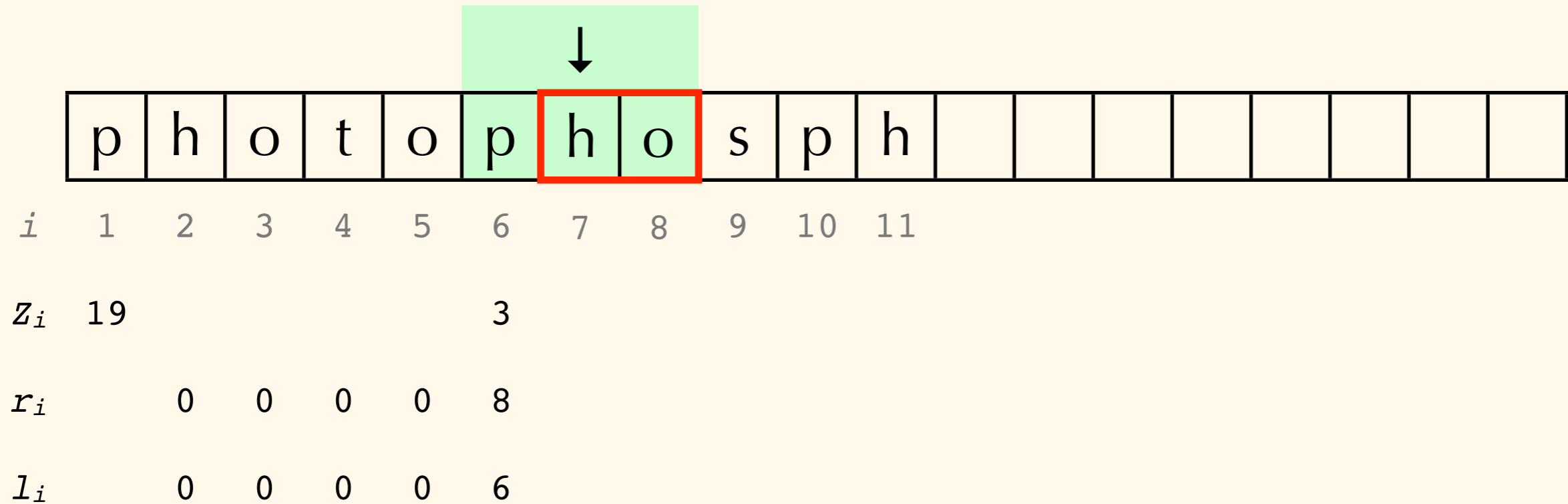


$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

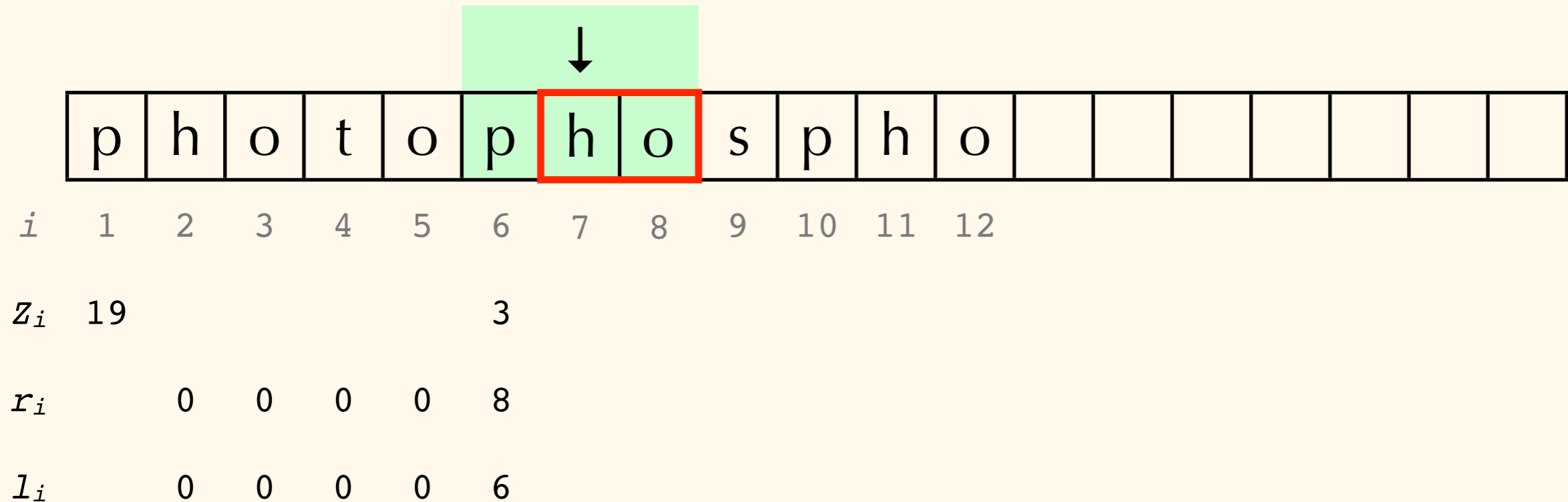


$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

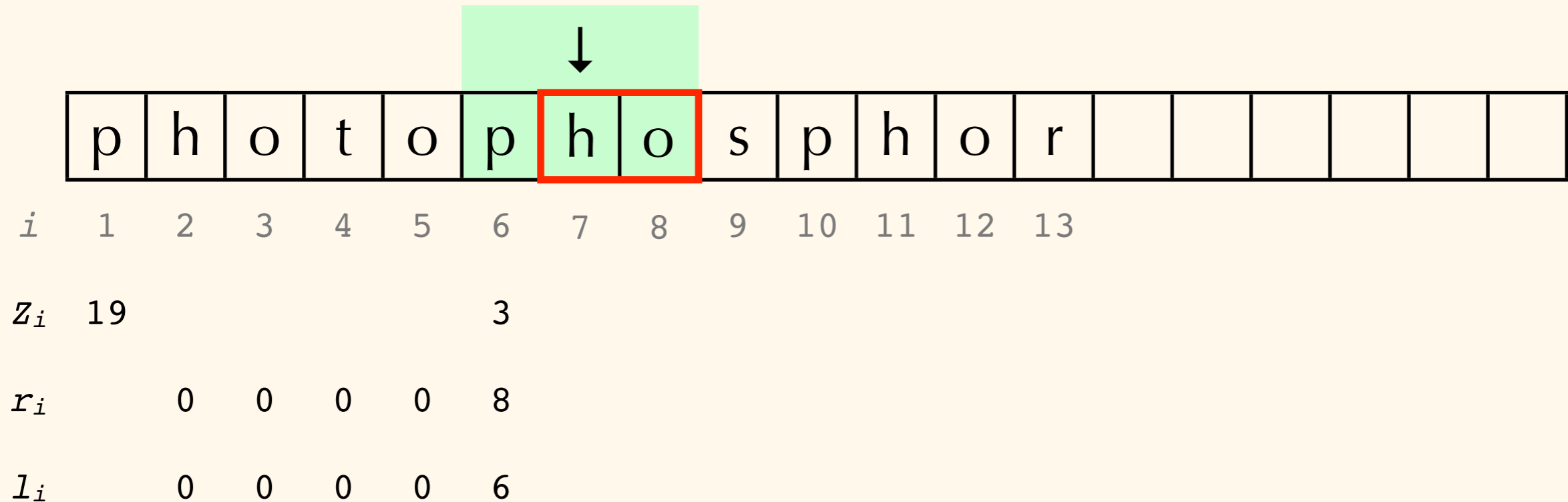


$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

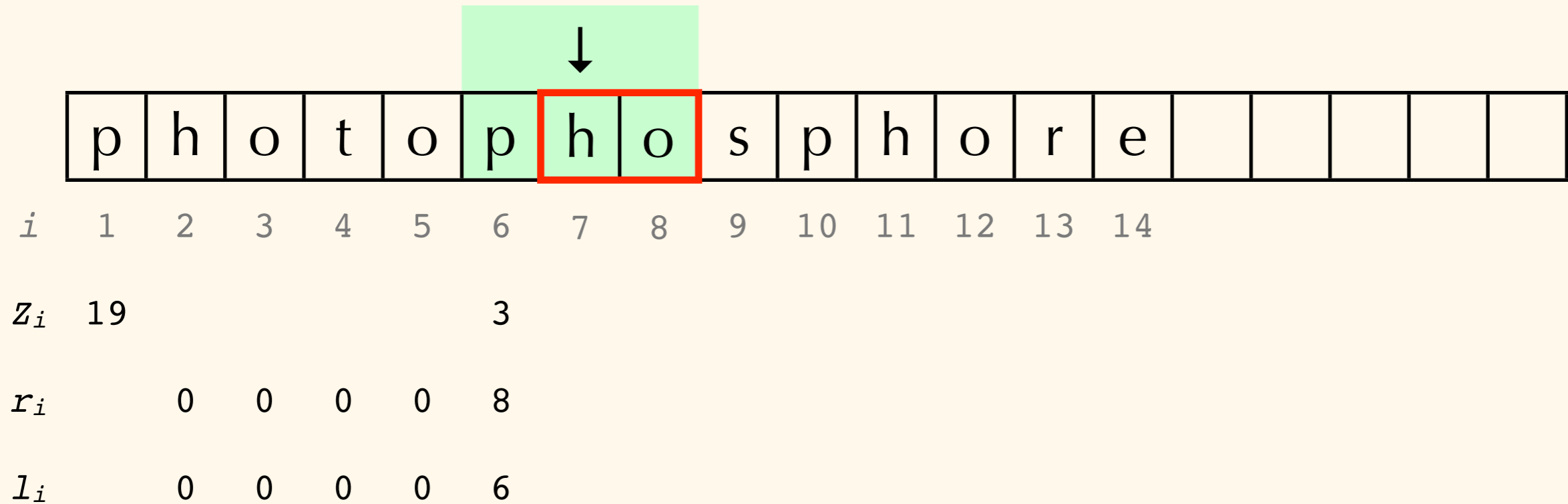


$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

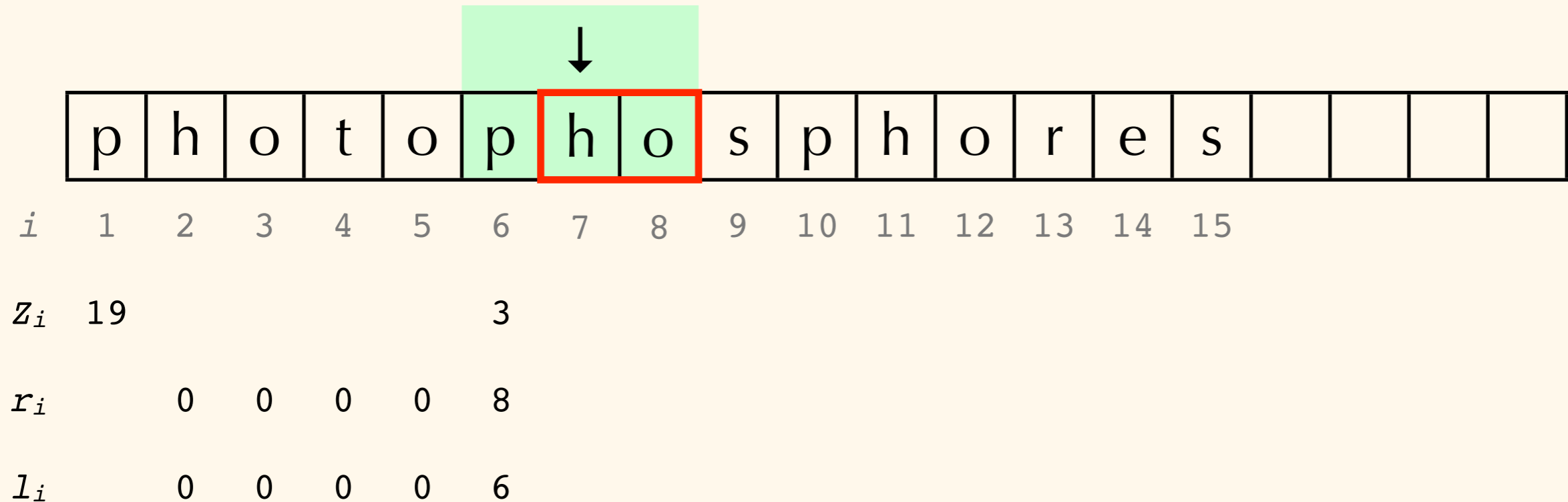


$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

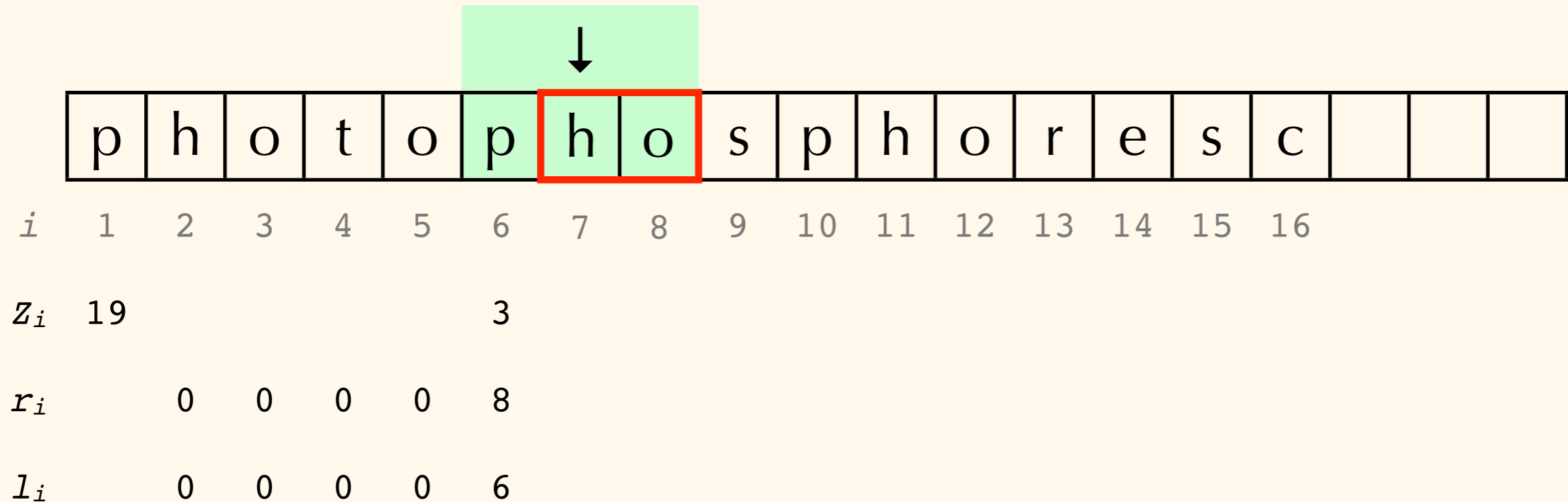


$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

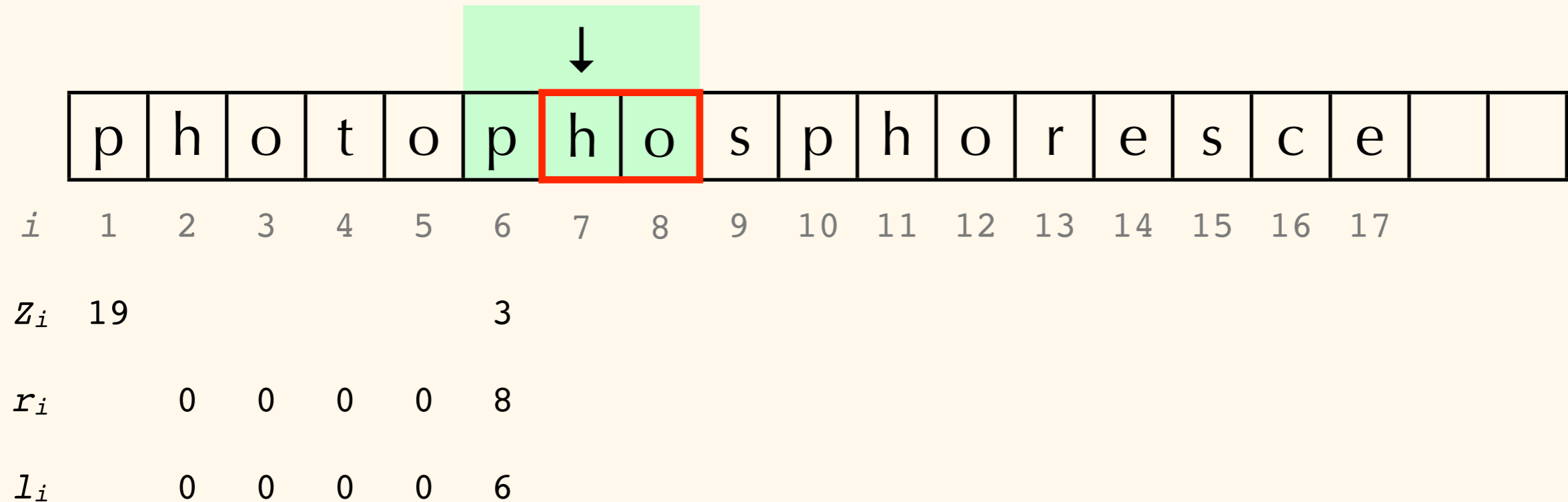


$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



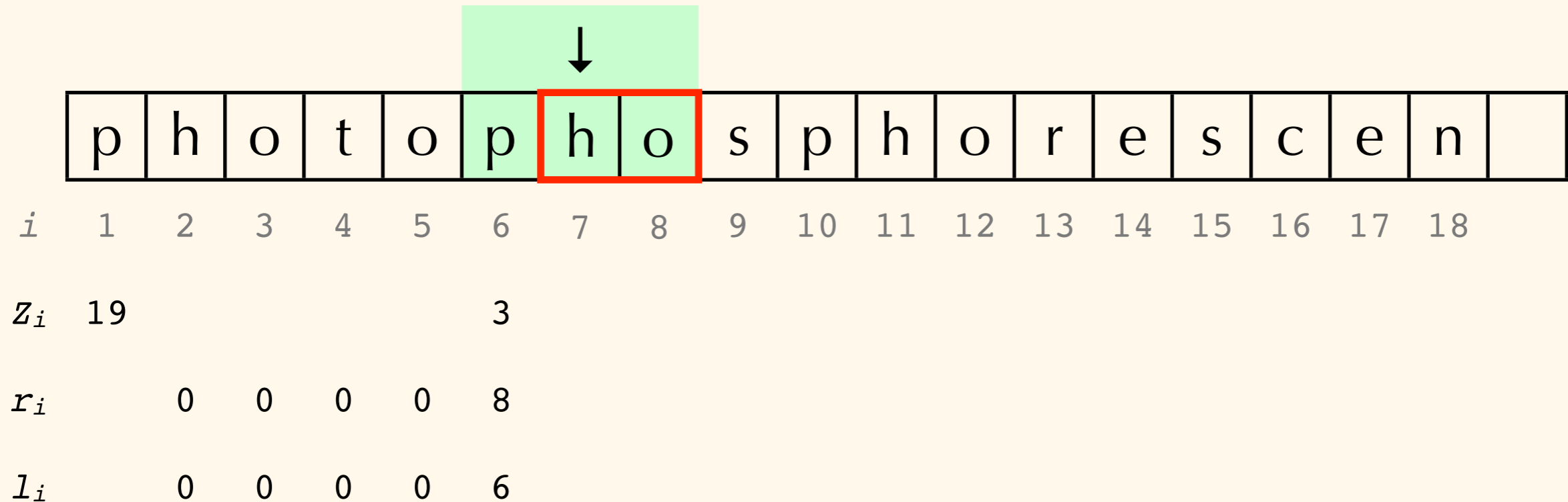
$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$



If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

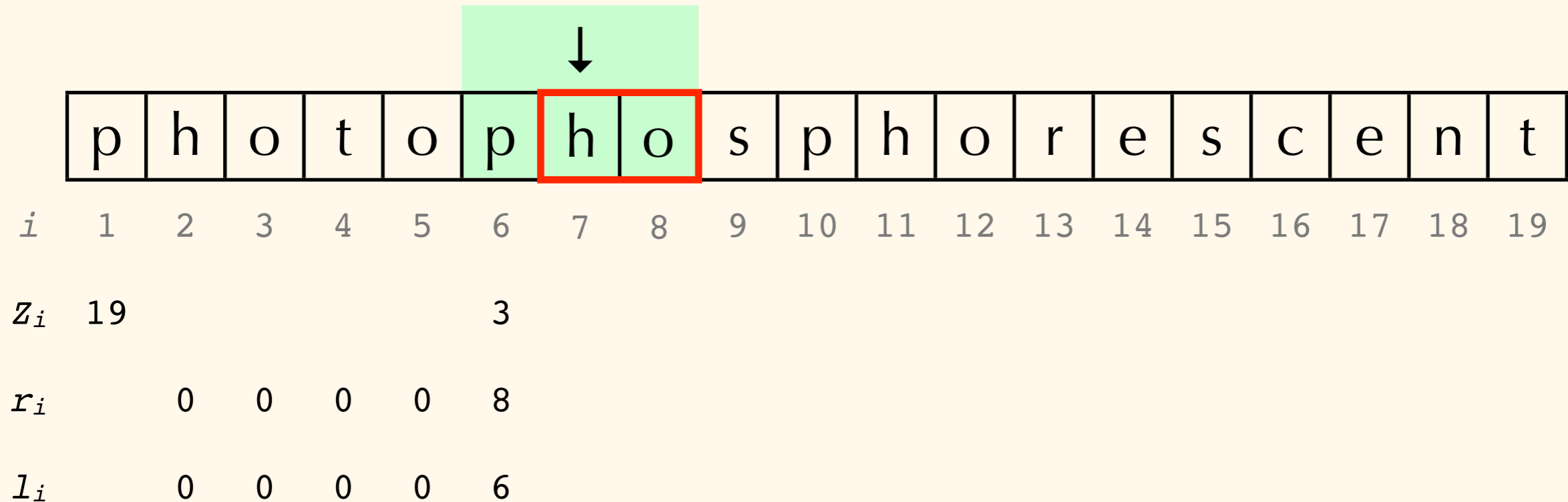


$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



$$\beta = S[k,r] = S[7,8] = \text{"ho"}$$

If  $k \leq r$ , we *are* inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

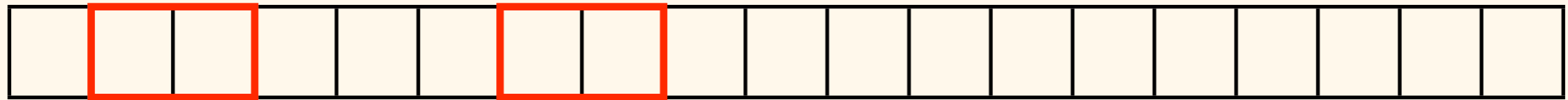
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



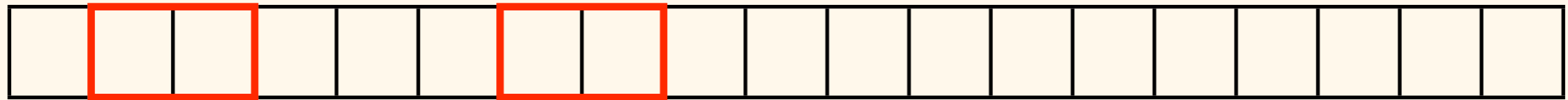
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



$i$

$Z_i$

$r_i$

$l_i$

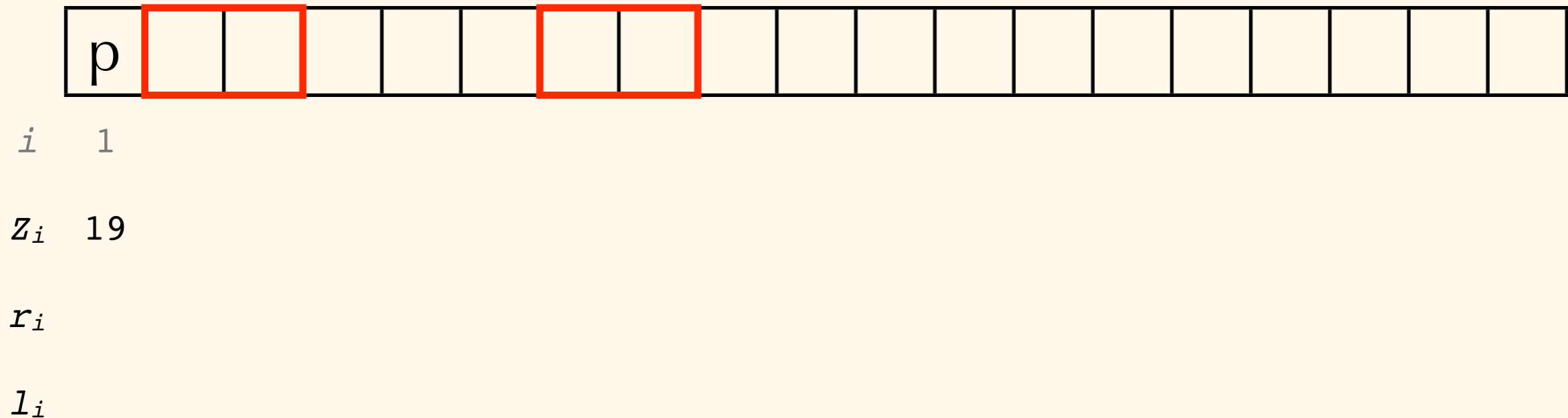
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



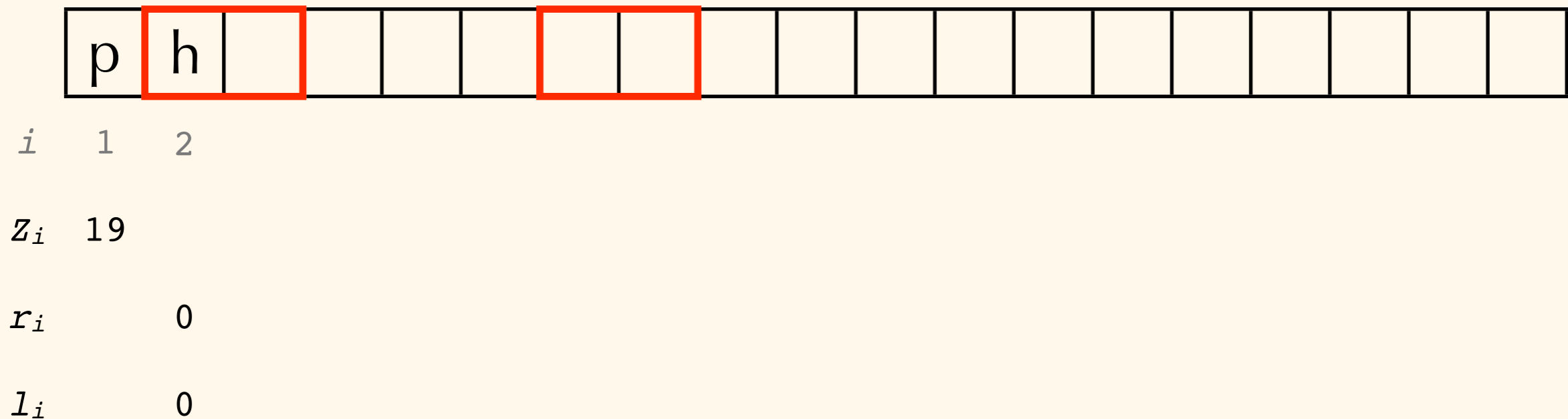
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o															
$i$	1	2	3															
$Z_i$	19																	
$r_i$		0	0															
$l_i$		0	0															

$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$



If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o	t														
$i$	1	2	3	4														
$Z_i$	19																	
$r_i$		0	0	0														
$l_i$		0	0	0														

$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o	t	o														
$i$	1	2	3	4	5														
$Z_i$	19																		
$r_i$		0	0	0	0														
$l_i$		0	0	0	0														

$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o	t	o	p												
$i$	1	2	3	4	5	6												
$Z_i$	19					3												
$r_i$		0	0	0	0	8												
$l_i$		0	0	0	0	6												

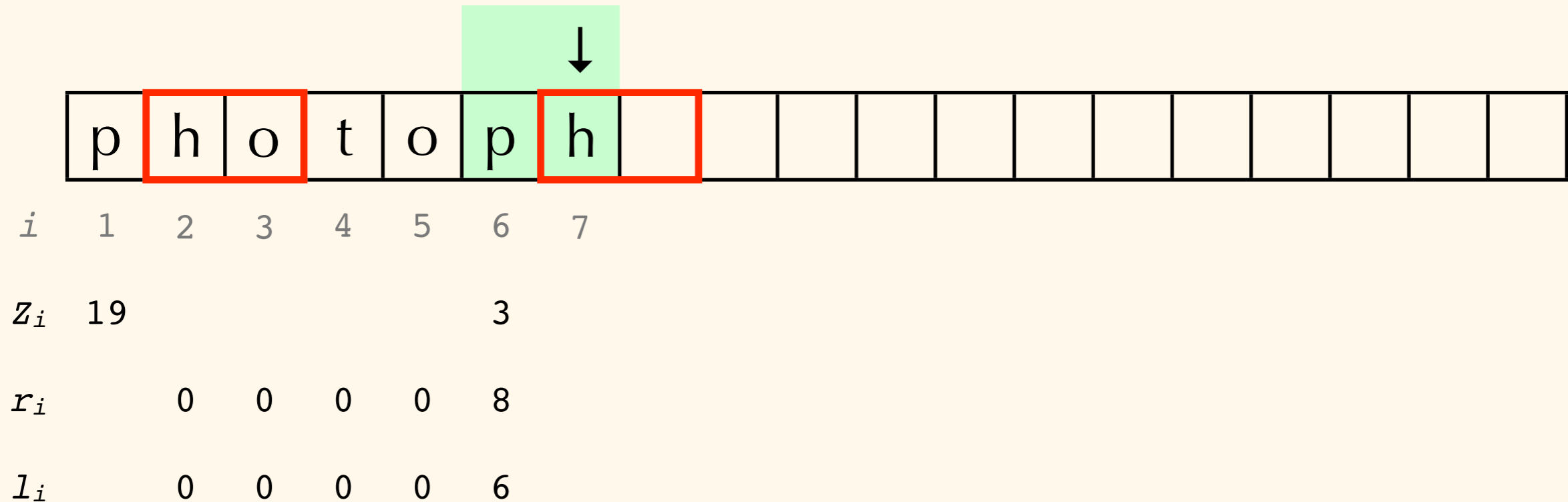
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



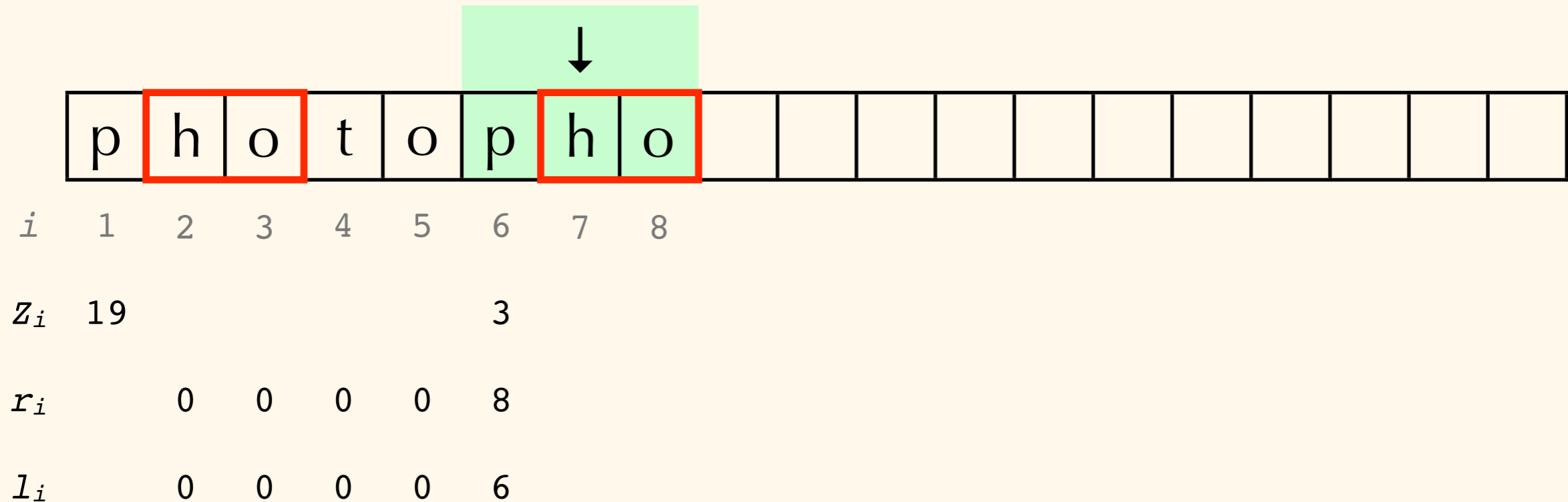
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



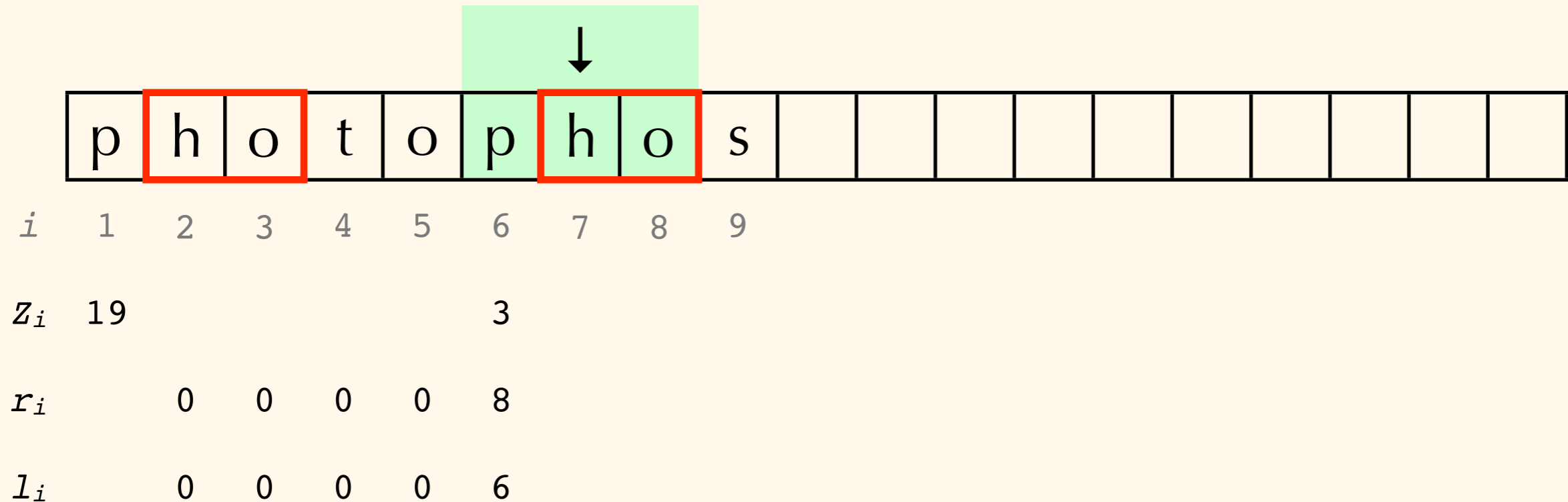
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



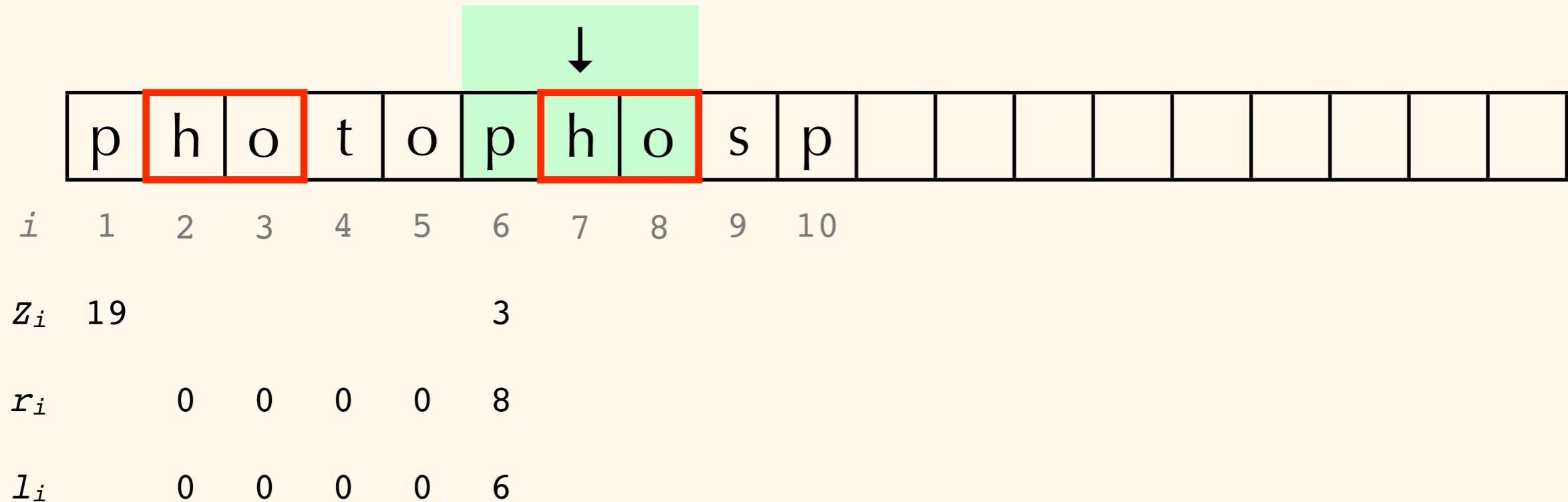
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



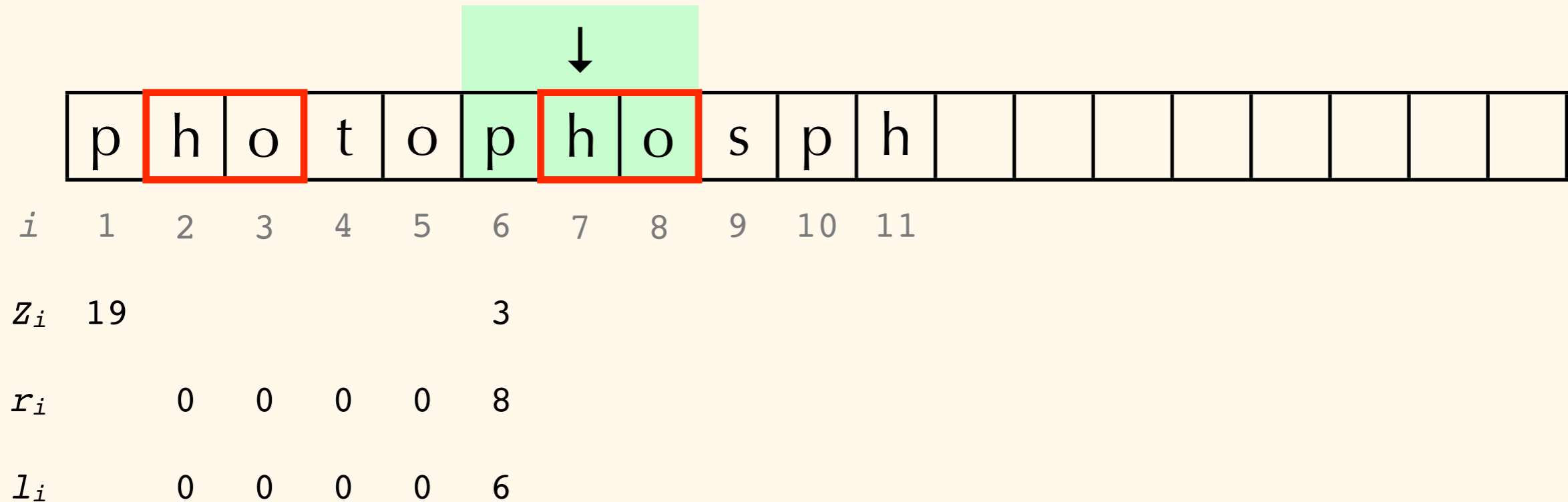
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



$$\beta = S[k,r] = \text{"ho"}$$

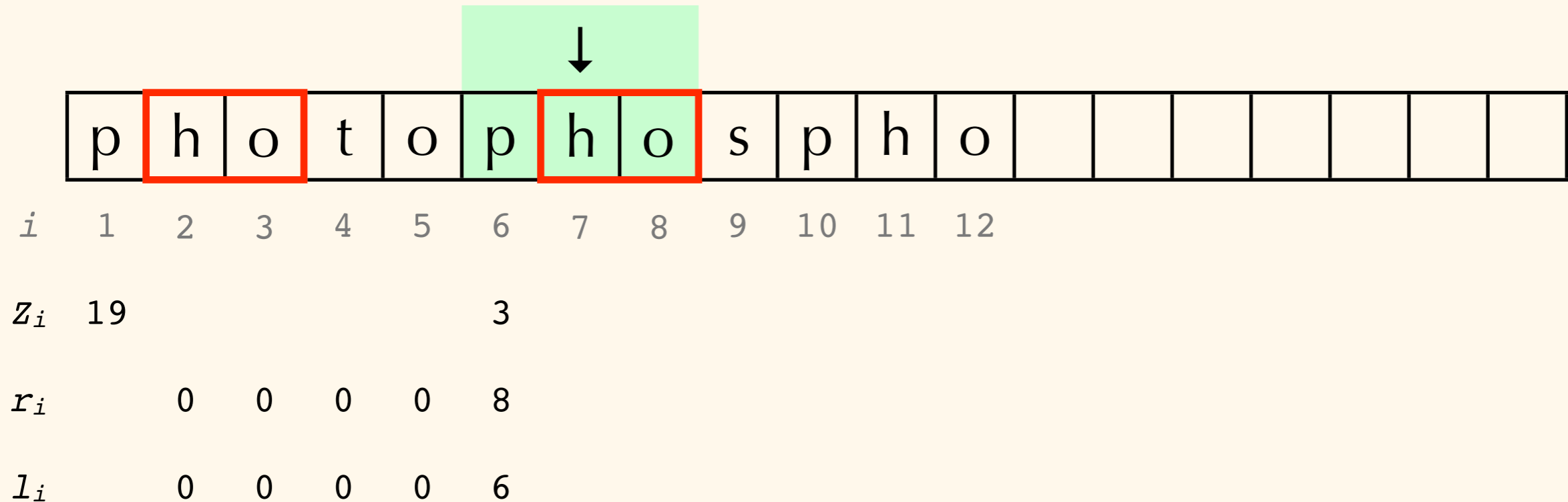
$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$



If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



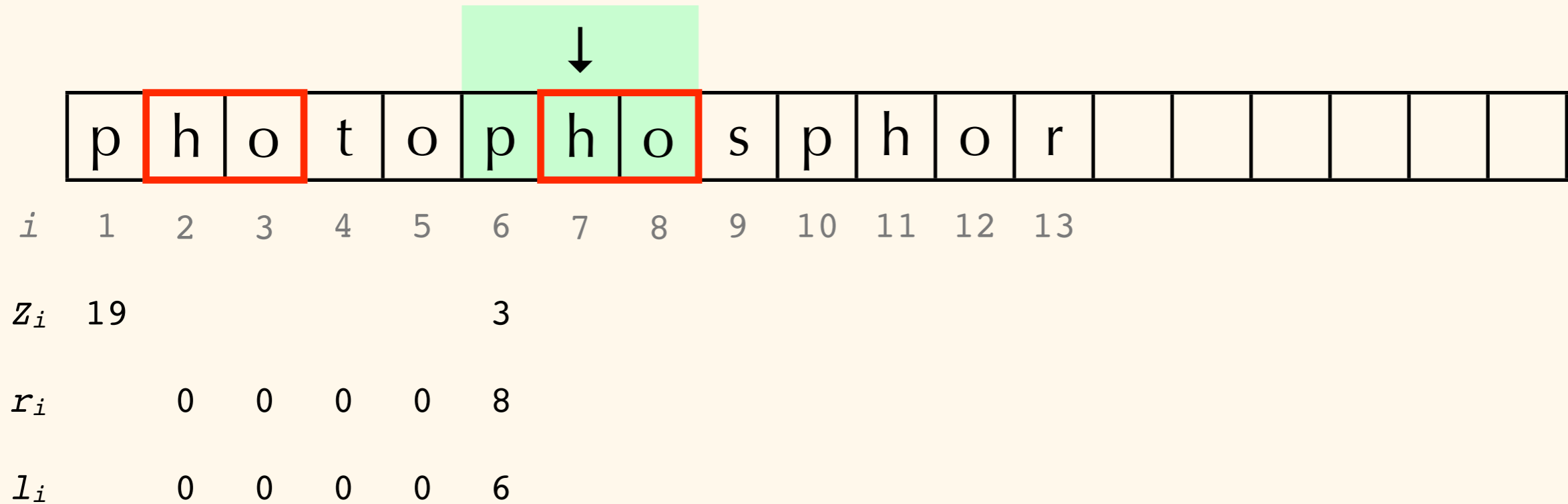
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



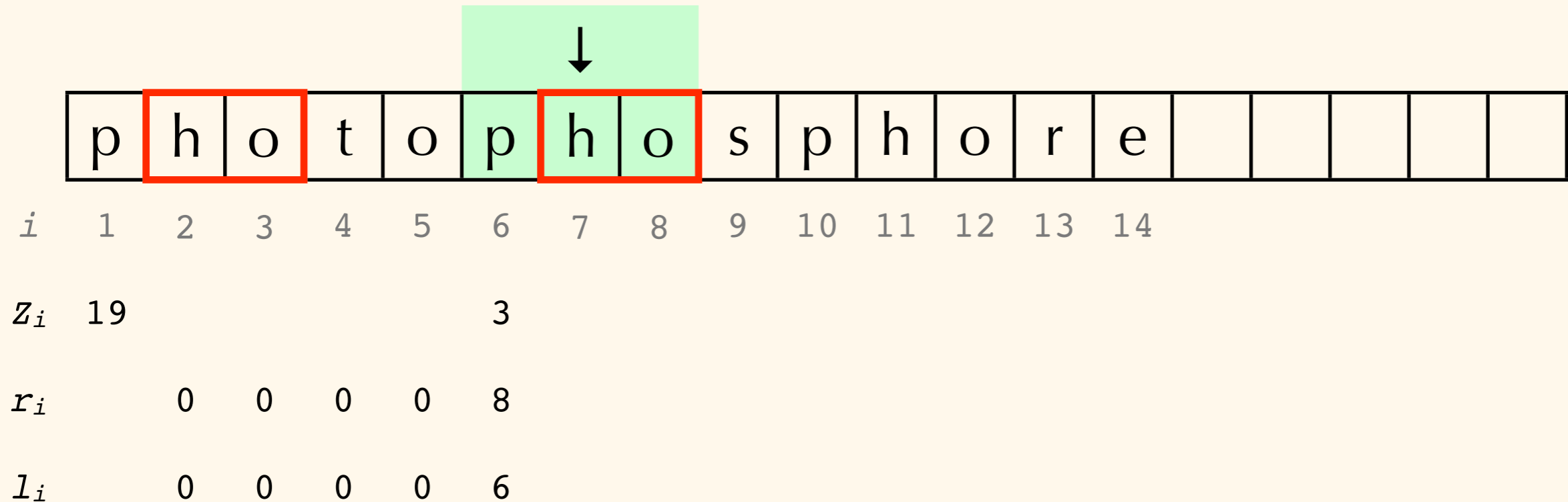
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



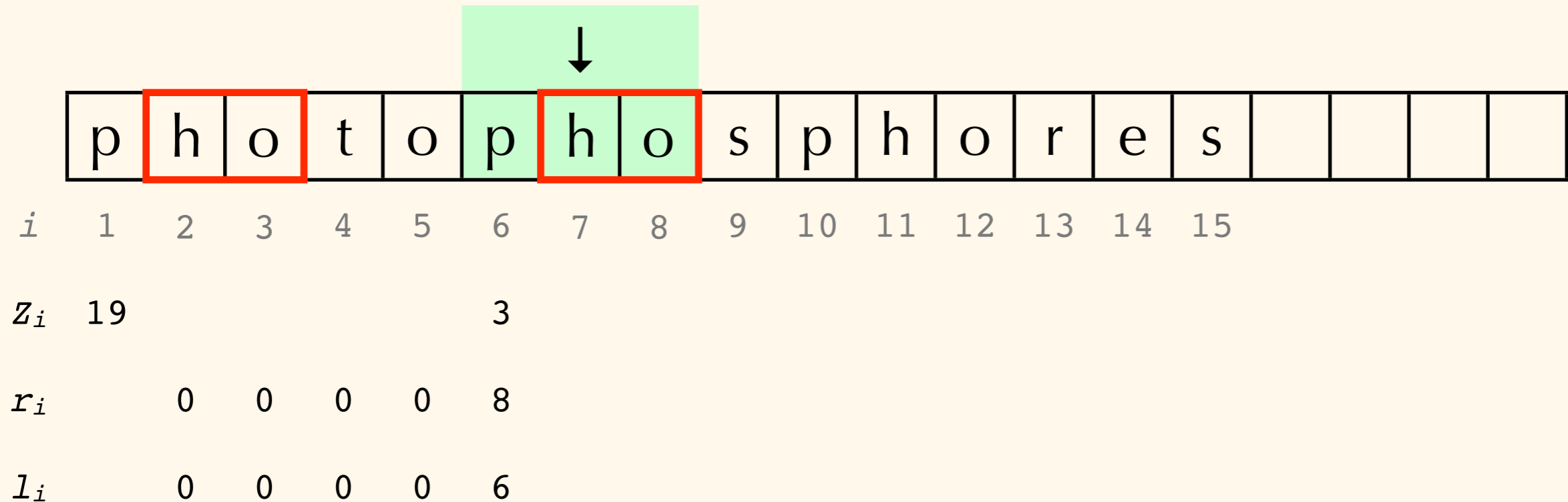
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



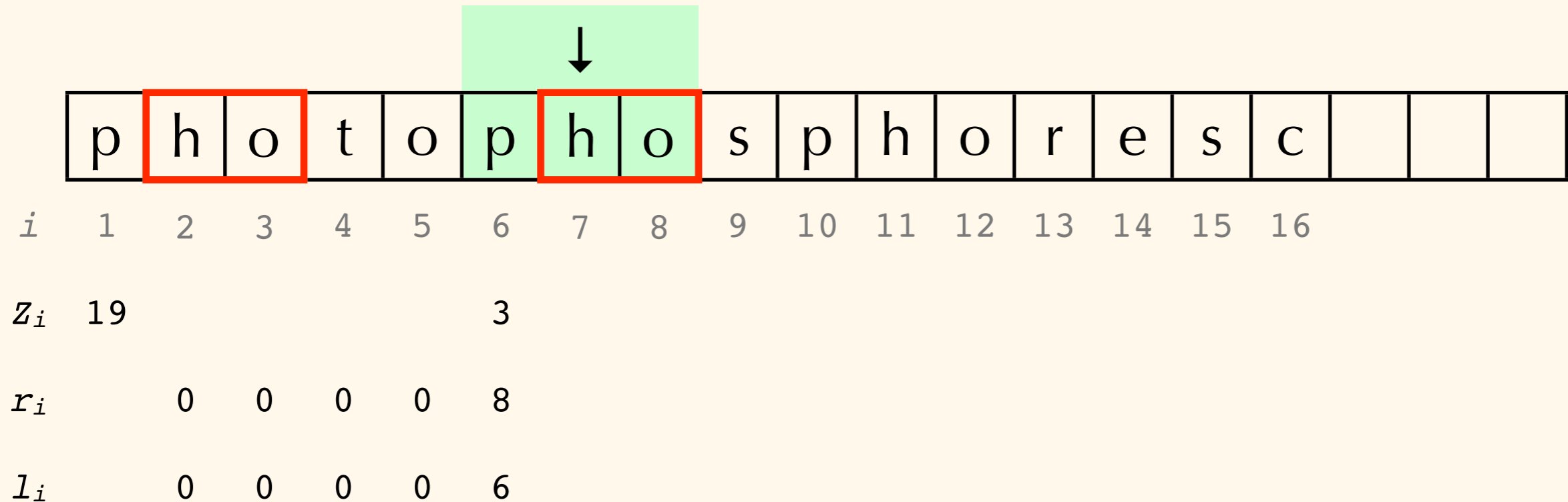
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



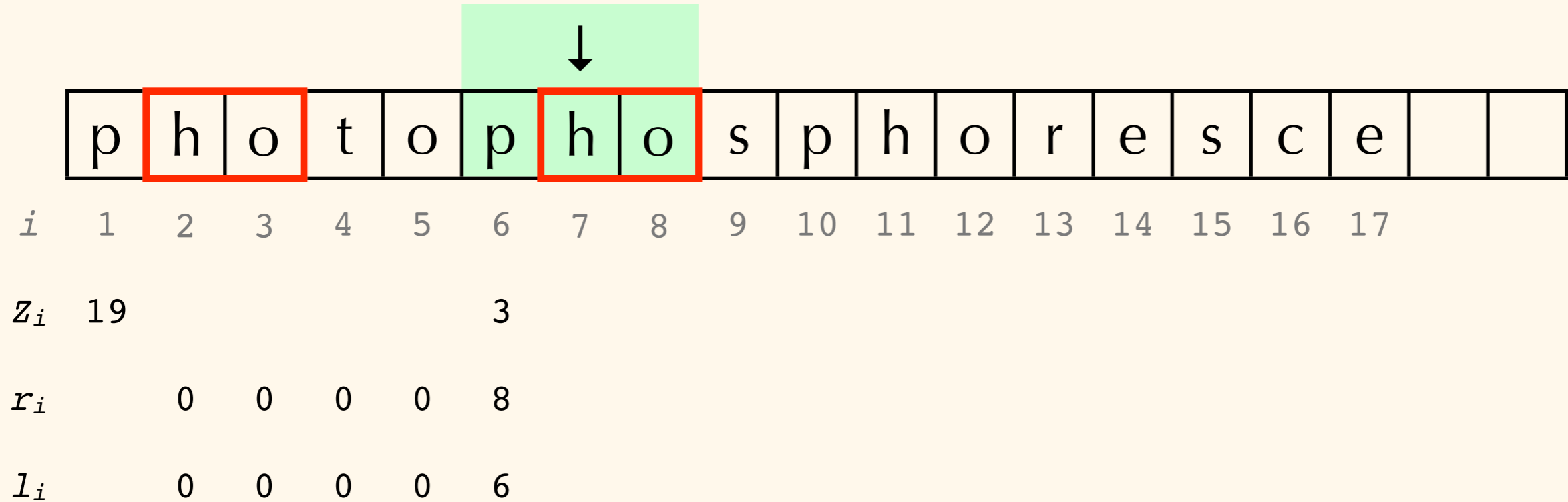
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$



$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$

	p	h	o	t	o	p	h	o	s	p	h	o	r	e	s	c	e	n	
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	
$Z_i$	19					3													
$r_i$		0	0	0	0	8													
$l_i$		0	0	0	0	6													

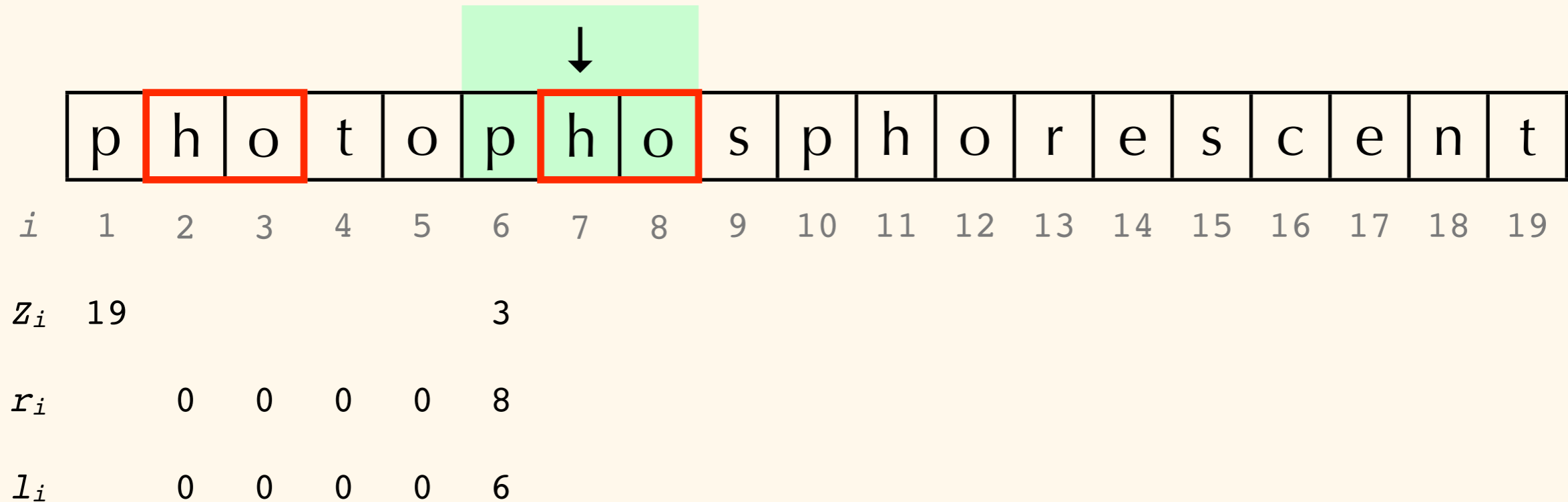
$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

We know that  $\beta$  must match  $S[k',Z_l]...$




$$\beta = S[k,r] = \text{"ho"}$$

$$S[k',Z_l] = S[7-6+1, 3] = S[2,3] = \text{"ho"}$$

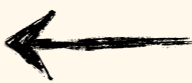


Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way


If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

If  $k \leq r$ , we are inside of an already-found z-box, so:

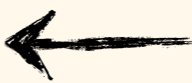
$$k' = k - l + 1 ; \beta = S[k,r]$$

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way

If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*


If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = |S[k, r]|$$

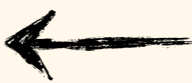
Two possibilities:

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way

If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

If  $k \leq r$ , we are inside of an already-found z-box, so:


$$k' = k - l + 1 ; \beta = |S[k, r]|$$

Two possibilities:

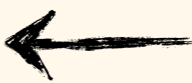
$S(k)$  could simply be part of a matching substring...

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way

If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = |S[k,r]|$$

Two possibilities:

$S(k)$  could simply be part of a matching substring...

Alternatively,  $S[k,r]$  could itself be a matching prefix of  $S$ !

Two possibilities:

$S(k)$  could simply be part of a matching substring...

Two possibilities:

$S(k)$  could simply be part of a matching substring...



Two possibilities:

$S(k)$  could simply be part of a matching substring...



$i$

$Z_i$

Two possibilities:

$S(k)$  could simply be part of a matching substring...



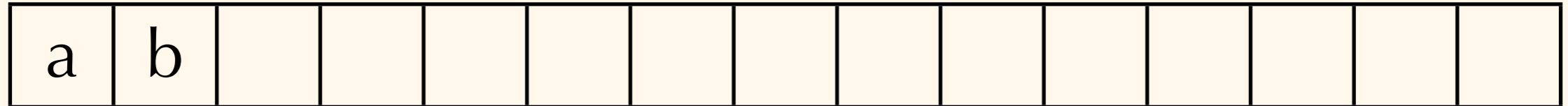
$i$  1

$Z_i$  15



Two possibilities:

$S(k)$  could simply be part of a matching substring...

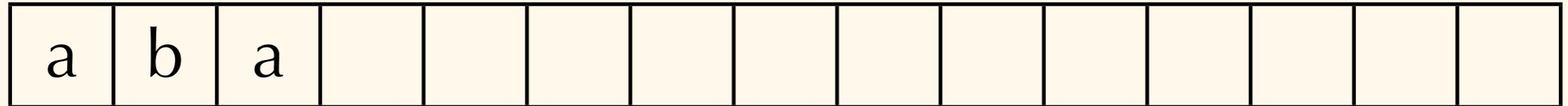


$i$     1    2

$Z_i$    15

Two possibilities:

$S(k)$  could simply be part of a matching substring...

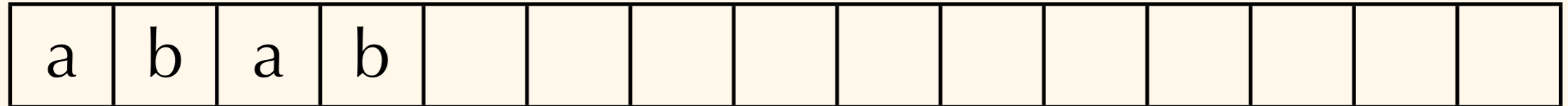


$i$     1    2    3

$Z_i$    15        2

Two possibilities:

$S(k)$  could simply be part of a matching substring...



$i$     1    2    3    4

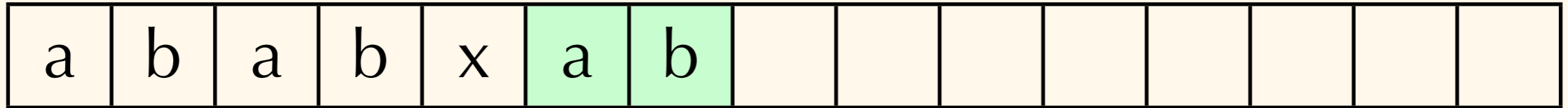
$Z_i$    15        2





Two possibilities:

$S(k)$  could simply be part of a matching substring...



$i$     1    2    3    4    5    6    7

$Z_i$    15        2                4







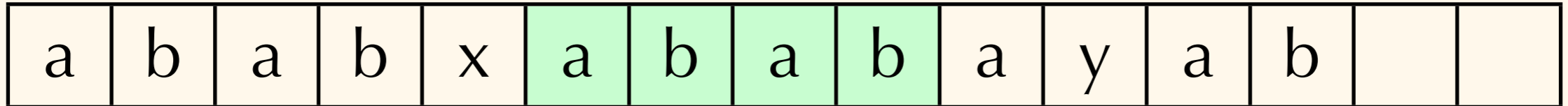






Two possibilities:

$S(k)$  could simply be part of a matching substring...



$i$     1    2    3    4    5    6    7    8    9    10    11    12    13

$Z_i$    15        2                4







Two possibilities:

$S(k)$  could simply be part of a matching substring...



	a	b	a	b	x	a	b	a	b	a	y	a	b	a	b
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$Z_i$	15		2			4									

In this case, we know that the characters between  $S(k)$  and  $S(r)$  are all part of a single match...



Two possibilities:

$S(k)$  could simply be part of a matching substring...



	a	b	a	b	x	a	b	a	b	a	y	a	b	a	b
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$Z_i$	15		2			4									

In this case, we know that the characters between  $S(k)$  and  $S(r)$  are all part of a single match...

... and so:

- $Z_k$  must be equal to  $Z_{k'}$ , and therefore...
- we don't need to update  $l$  and  $r$



Two possibilities:

$S(k)$  could simply be part of a matching substring...

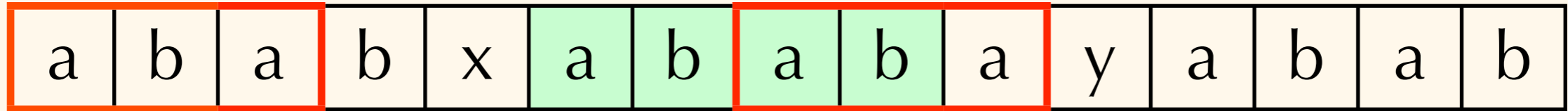


	a	b	a	b	x	a	b	a	b	a	y	a	b	a	b
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$Z_i$	15		2			4									

Alternatively,  $S[k,r]$  could itself be a matching prefix of  $S$ !



Alternatively,  $S[k,r]$  could itself be a matching prefix of  $S$ !



$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$Z_i$	15		2			4									

In this case, we are in an *overlapping* z-box...

Alternatively,  $S[k,r]$  could itself be a matching prefix of  $S$ !



	a	b	a	b	x	a	b	a	b	a	y	a	b	a	b
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$Z_i$	15		2			4									

In this case, we are in an *overlapping* z-box...

... so  $Z_k$  might be potentially end up being different from  $Z_{k'}$ ...

Alternatively,  $S[k,r]$  could itself be a matching prefix of  $S$ !



	a	b	a	b	x	a	b	a	b	a	y	a	b	a	b
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$Z_i$	15		2			4									


In this case, we are in an *overlapping* z-box...

... so  $Z_k$  might be potentially end up being different from  $Z_{k'}$ ...

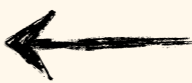
... and we might need to update  $l$  and  $r$  (to reflect the boundaries of the *new* z-box).

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way

If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

If  $k \leq r$ , we are inside of an already-found z-box, so:

$$k' = k - l + 1 ; \beta = S[k,r]$$

Two possibilities:

$S(k)$  could simply be part of a matching substring...


Alternatively,  $S[k,r]$  could itself be a matching prefix of  $S$ !

How to tell which condition? Compare  $Z_{k'}$  to  $|\beta|$ :

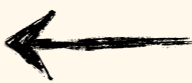


Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way


If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

If  $k \leq r$ , we are inside of an already-found z-box, so:

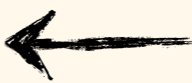
$$k' = k - l + 1 ; \beta = S[k,r]$$

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way

If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*


If  $k \leq r$ , we are inside of an already-found z-box, so:

$k' = k - l + 1$  ;  $\beta = S[k,r]$

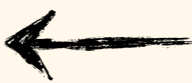
If  $Z_{k'} < |\beta|$ ,  $Z_k = Z_{k'}$  and  $l, r$  are unchanged;

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way

If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

If  $k \leq r$ , we are inside of an already-found z-box, so:


$k' = k - l + 1$  ;  $\beta = S[k,r]$

If  $Z_{k'} < |\beta|$ ,  $Z_k = Z_{k'}$  and  $l, r$  are unchanged;

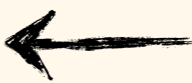
If  $Z_{k'} \geq |\beta|$ :

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way

If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

If  $k \leq r$ , we are inside of an already-found z-box, so:

$k' = k - l + 1$  ;  $\beta = S[k,r]$


If  $Z_{k'} < |\beta|$ ,  $Z_k = Z_{k'}$  and  $l, r$  are unchanged;

If  $Z_{k'} \geq |\beta|$ :

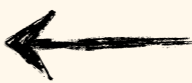
We know that  $Z_k$  must be at least  $Z_{k'}$  - but it could be longer...

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way

If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

If  $k \leq r$ , we are inside of an already-found z-box, so:

$k' = k - l + 1$  ;  $\beta = S[k,r]$

If  $Z_{k'} < |\beta|$ ,  $Z_k = Z_{k'}$  and  $l, r$  are unchanged;


If  $Z_{k'} \geq |\beta|$ :

We know that  $Z_k$  must be at least  $Z_{k'}$  - but it could be longer...

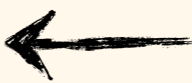
Start looking for a match between  $S[r + 1, ]$  and  $S[|\beta| + 1, ]$

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way

If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

If  $k \leq r$ , we are inside of an already-found z-box, so:

$k' = k - l + 1$  ;  $\beta = S[k,r]$

*k is inside a z-box, but is not the start of an overlapping box*

If  $Z_{k'} < |\beta|$ ,  $Z_k = Z_{k'}$  and  $l, r$  are unchanged; 


If  $Z_{k'} \geq |\beta|$ :

We know that  $Z_k$  must be at least  $Z_{k'}$  - but it could be longer...

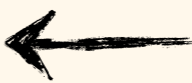
Start looking for a match between  $S[r + 1, ]$  and  $S[|\beta| + 1, ]$

Putting it all together into an algorithm:

Initialize  $l$  and  $r$  to 0; for each  $k$ ,  $1 < k \leq |S|$  :

If  $k > r$ , we are not in a z-box, so:  *i.e., not inside a previously-found matching region*

Calculate  $Z_k$  the normal way


If  $Z_k > 0$ , set  $l = k$  and  $r = k + Z_k - 1$   *k is the beginning of a match of length  $Z_k$*

If  $k \leq r$ , we are inside of an already-found z-box, so:

$k' = k - l + 1$  ;  $\beta = S[k, r]$

*k is inside a z-box, but is not the start of an overlapping box*

If  $Z_{k'} < |\beta|$ ,  $Z_k = Z_{k'}$  and  $l, r$  are unchanged; 

If  $Z_{k'} \geq |\beta|$ :  *k must itself be the beginning of a new, overlapping z-box*

We know that  $Z_k$  must be at least  $Z_{k'}$  - but it could be longer...

Start looking for a match between  $S[r + 1, ]$  and  $S[|\beta| + 1, ]$









Alternatively,  $S[k,r]$  could itself be a matching prefix of  $S$ !

		↓					↓								
	a	b	a	b	x	a	b	a	b	a	y	a	b	a	b
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$Z_i$	15		<u>2</u>			4									

In this case, we know that  $Z_k$  will be at least 2...

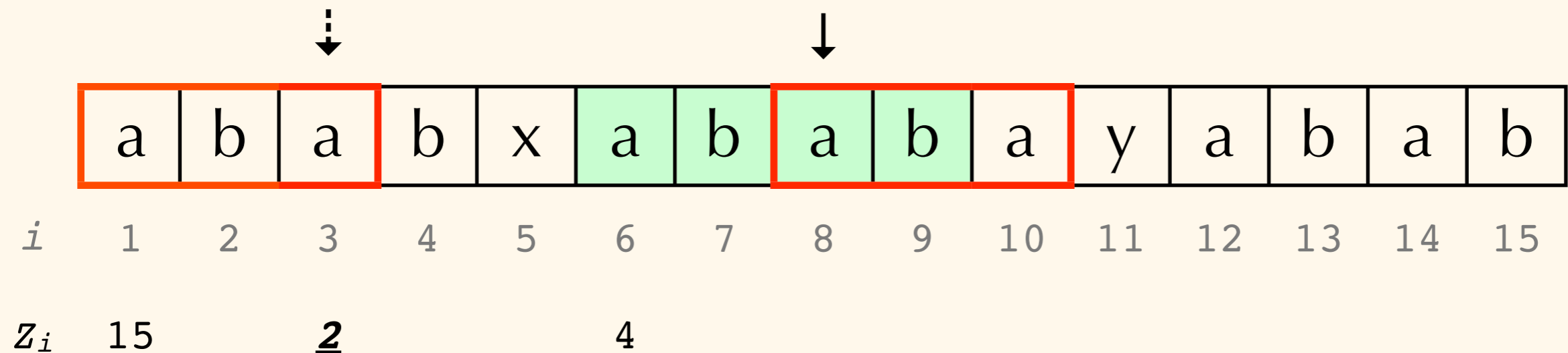
Alternatively,  $S[k,r]$  could itself be a matching prefix of  $S$ !

		↓					↓								
	a	b	a	b	x	a	b	a	b	a	y	a	b	a	b
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$Z_i$	15		<u>2</u>			4									

In this case, we know that  $Z_k$  will be at least 2...

... and so we'll need to do prefix matching, but we *don't* have to start at the beginning of the string.

Alternatively,  $S[k,r]$  could itself be a matching prefix of  $S$ !

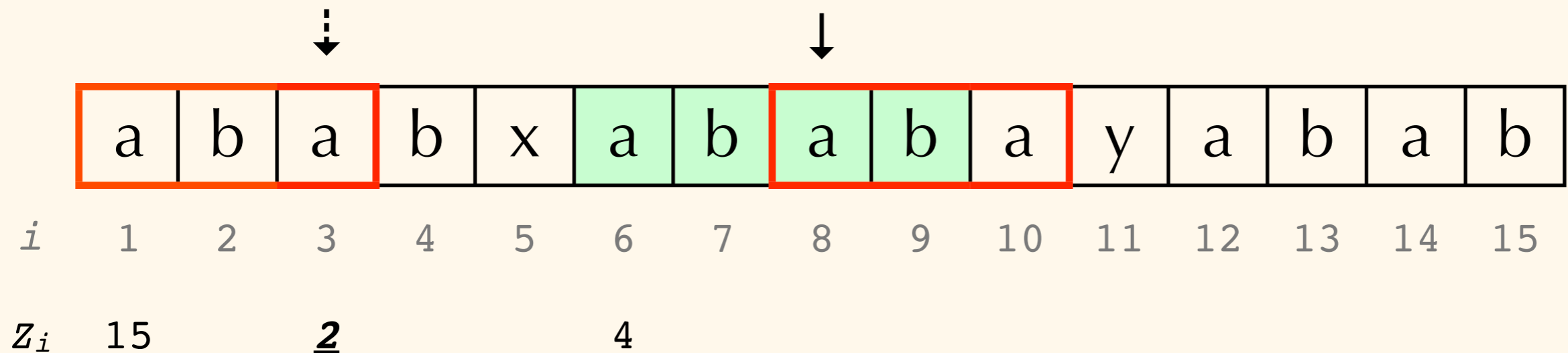


In this case, we know that  $Z_k$  will be at least 2...

... and so we'll need to do prefix matching, but we *don't* have to start at the beginning of the string.

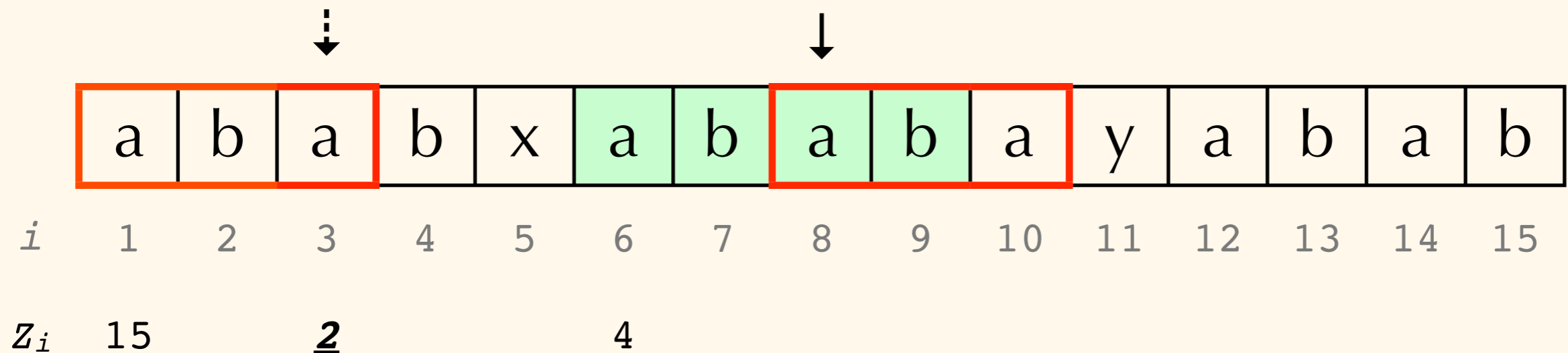
We know that  $S[k,r]$  *must* match  $S[1,Z_k]$ , so can use our current  $r$  offset to tell us where to start looking for a continued match.

Alternatively,  $S[k,r]$  could itself be a matching prefix of  $S$ !



We know that  $S[k,r]$  *must* match  $S[1,Z_k]$ , so can use our current  $r$  offset to tell us where to start looking for a continued match.

Alternatively,  $S[k,r]$  could itself be a matching prefix of  $S$ !



We know that  $S[k,r]$  *must* match  $S[1,Z_k]$ , so can use our current  $r$  offset to tell us where to start looking for a continued match.

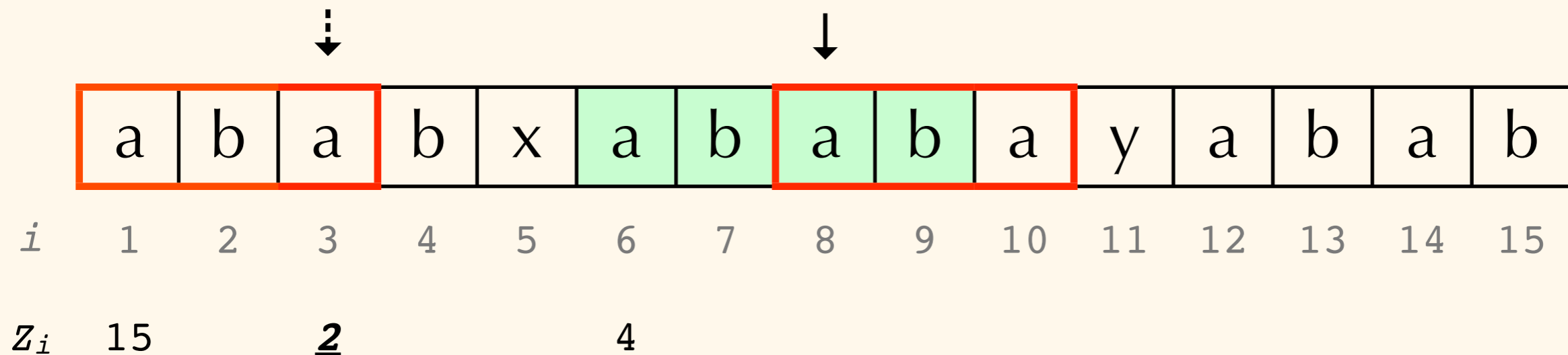
Start comparing  $S[r+1,]$  and  $S[|\beta|+1,]$ ; call the position in  $S$  where the first mismatch occurs  $q$ .





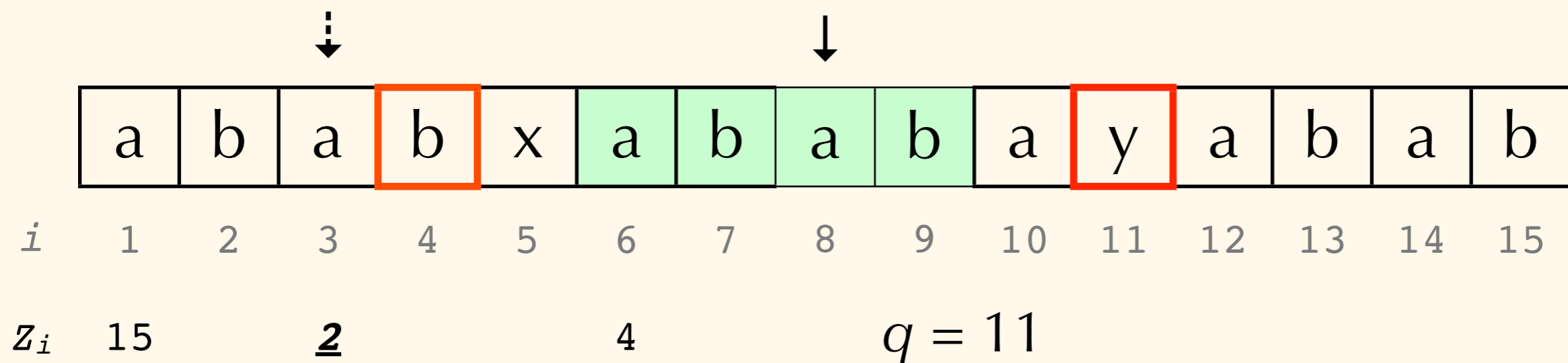


Alternatively,  $S[k,r]$  could itself be a matching prefix of  $S$ !

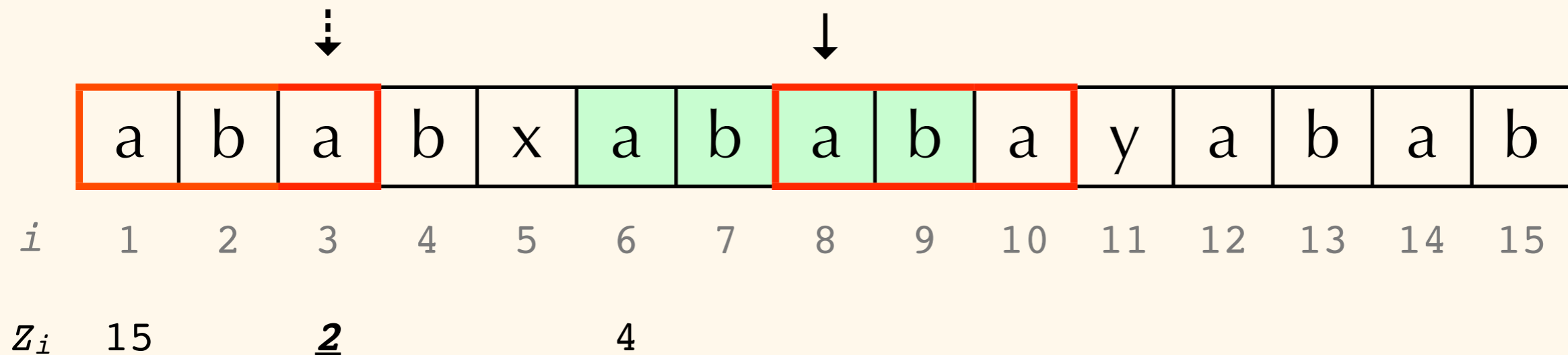


We know that  $S[k,r]$  must match  $S[1,Z_k]$ , so can use our current  $r$  offset to tell us where to start looking for a continued match.

Start comparing  $S[r+1,]$  and  $S[|\beta|+1,]$ ; call the position in  $S$  where the first mismatch occurs  $q$ .



Alternatively,  $S[k,r]$  could itself be a matching prefix of  $S$ !



We know that  $S[k,r]$  *must* match  $S[1,Z_k]$ , so can use our current  $r$  offset to tell us where to start looking for a continued match.

Start comparing  $S[r+1,]$  and  $S[|\beta|+1,]$ ; call the position in  $S$  where the first mismatch occurs  $q$ .

$$q = 11$$

We can now set  $Z_k = q - k$ ; the new  $r = q - 1$ ; and  $l = k$

			↓					↓							
	a	b	a	b	x	a	b	a	b	a	y	a	b	a	b
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$Z_i$	15		<u>2</u>			4									

Start comparing  $S[r+1, ]$  and  $S[|\beta|+1, ]$ ; call the position in  $S$  where the first mismatch occurs  $q$ .

$$q = 11$$

We can now set  $Z_k = q - k$ ; the new  $r = q - 1$ ; and  $l = k$

			↓					↓							
	a   b   a			b	x	a	b	a   b   a			y	a	b	a	b
$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$Z_i$	15		<u>2</u>			4		3							

Start comparing  $S[r+1, ]$  and  $S[|\beta|+1, ]$ ; call the position in  $S$  where the first mismatch occurs  $q$ .

$$q = 11$$

We can now set  $Z_k = q - k$ ; the new  $r = q - 1$ ; and  $l = k$

Now, continue on to the next  $k$ ...

What does this get us?

What does this get us?

A lot, in the *worst* case:

What does this get us?

A lot, in the *worst* case:

(a long, repetitive pattern, with a string full of matches or near-matches)



What does this get us?

A lot, in the *worst* case:

(a long, repetitive pattern, with a string full of matches or near-matches)

We get to “skip ahead” quite often in that case.

What does this get us?

A lot, in the *worst* case:

(a long, repetitive pattern, with a string full of matches or near-matches)

We get to “skip ahead” quite often in that case.

In a more “normal” case.... meh.

What does this get us?

A lot, in the *worst* case:

(a long, repetitive pattern, with a string full of matches or near-matches)

We get to “skip ahead” quite often in that case.

In a more “normal” case.... meh.

Less skipping == less (relative) benefit

# Plan for today:

Z-algorithm review

Knuth-Morris-Pratt

Boyer-Moore

# Historical notes:



Donald Knuth  
1938 –



Vaughan Pratt  
1944 –



James Morris  
1941 –

The KMP algorithm was discovered in 1974 by K & P at Stanford, and independently by M at CMU in that same year.

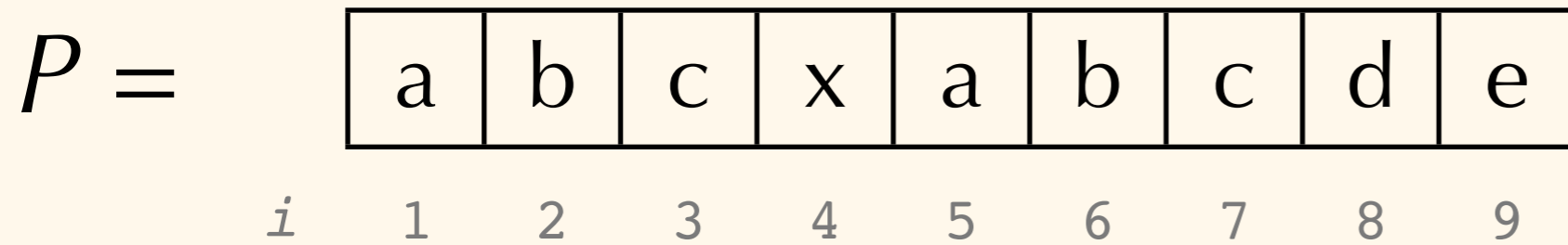
The three authors formally published together in 1977.

The Knuth-Morris-Pratt algorithm builds on top of the Z-algorithm...

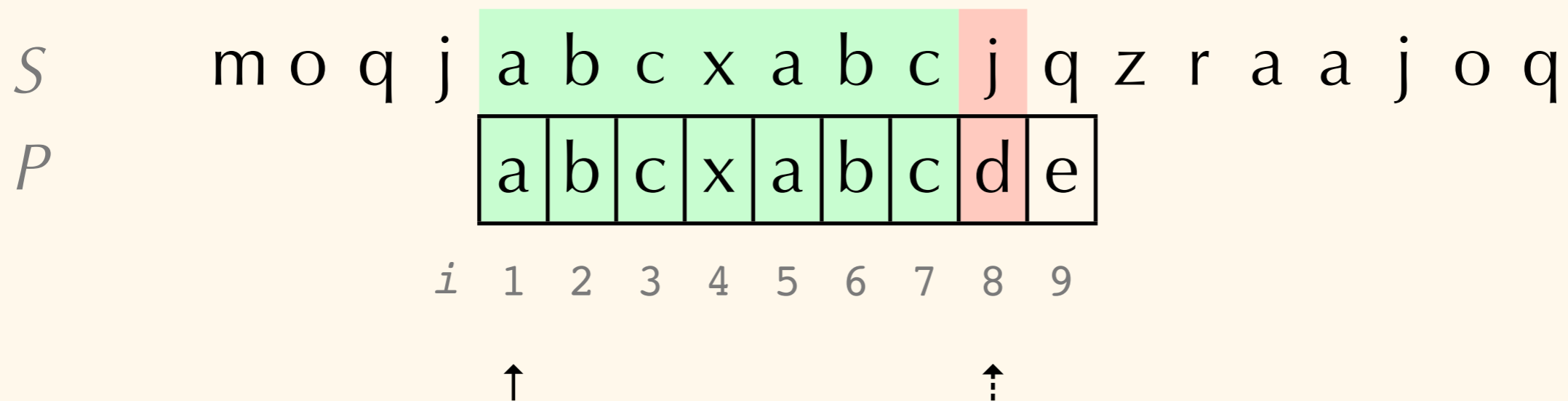
... the main innovation being that rather than moving through  $S$  one character at a time...

... we use information about repetitive segments of  $P$  to help us move more quickly.

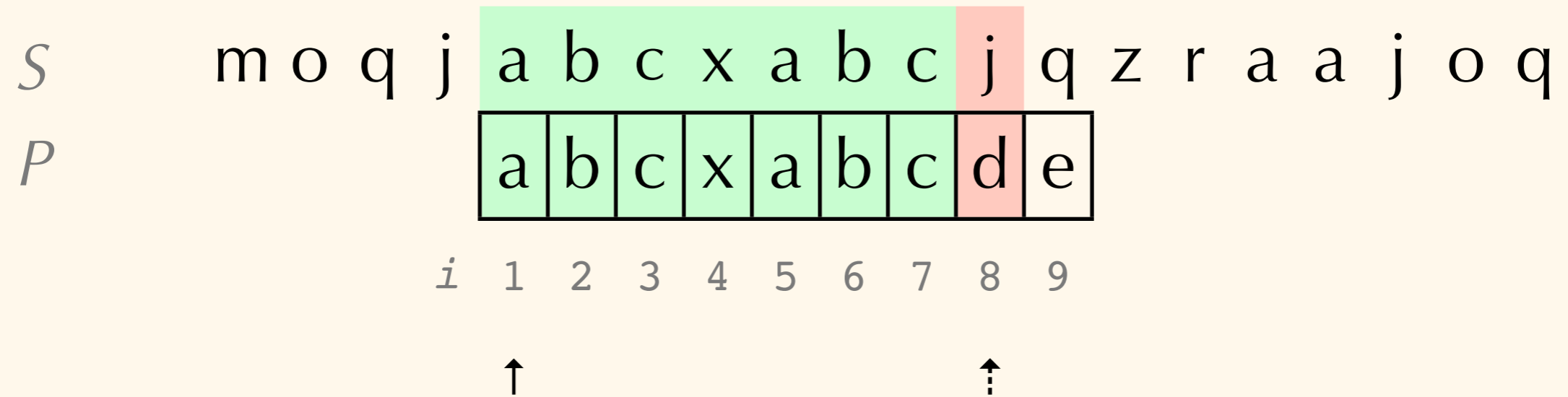
For example, (from Gusfield):



Consider a mismatch between  $P$  and  $S$  occurring at position 8 in  $P$ :

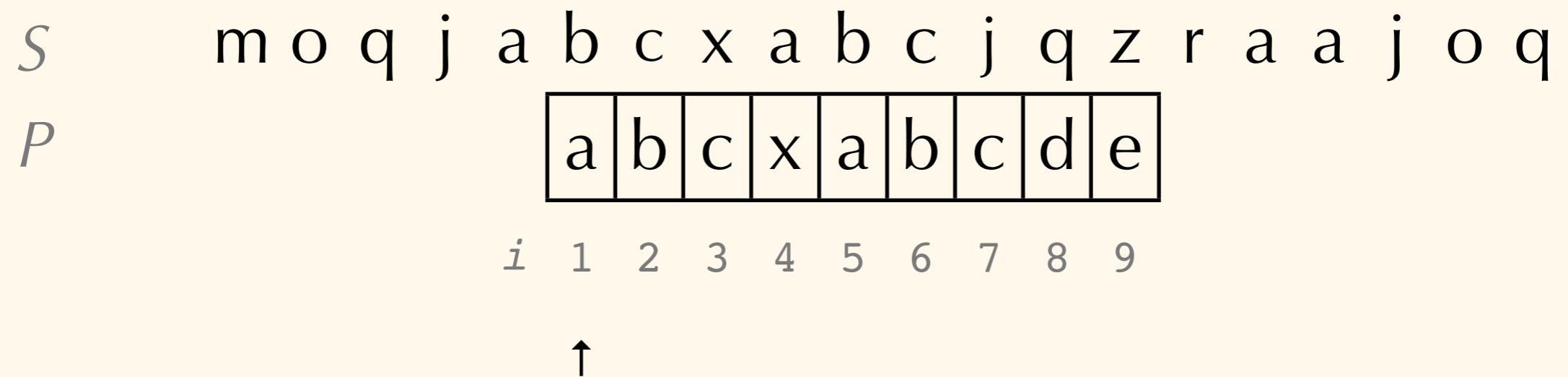


The previous algorithm would shift  $P$  down by one position, and then resume searching for a match:



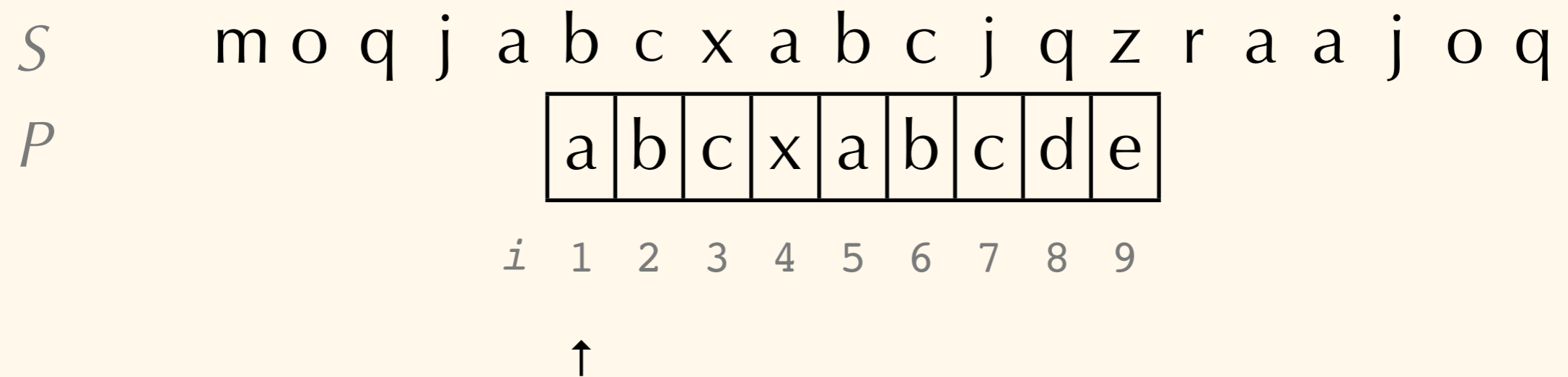


The previous algorithm would shift  $P$  down by one position, and then resume searching for a match:



$S$   
 $P$

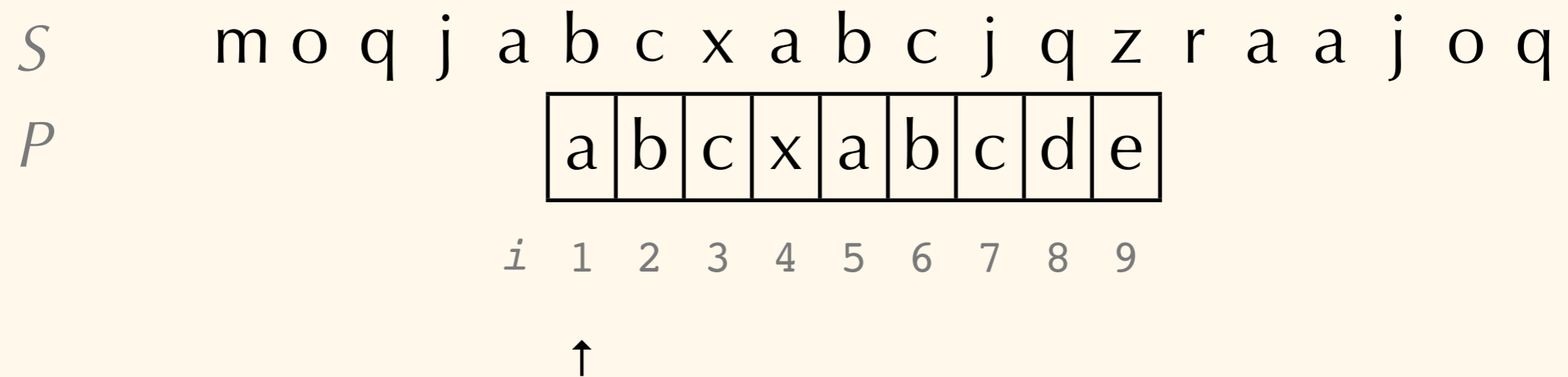
The previous algorithm would shift  $P$  down by one position, and then resume searching for a match:



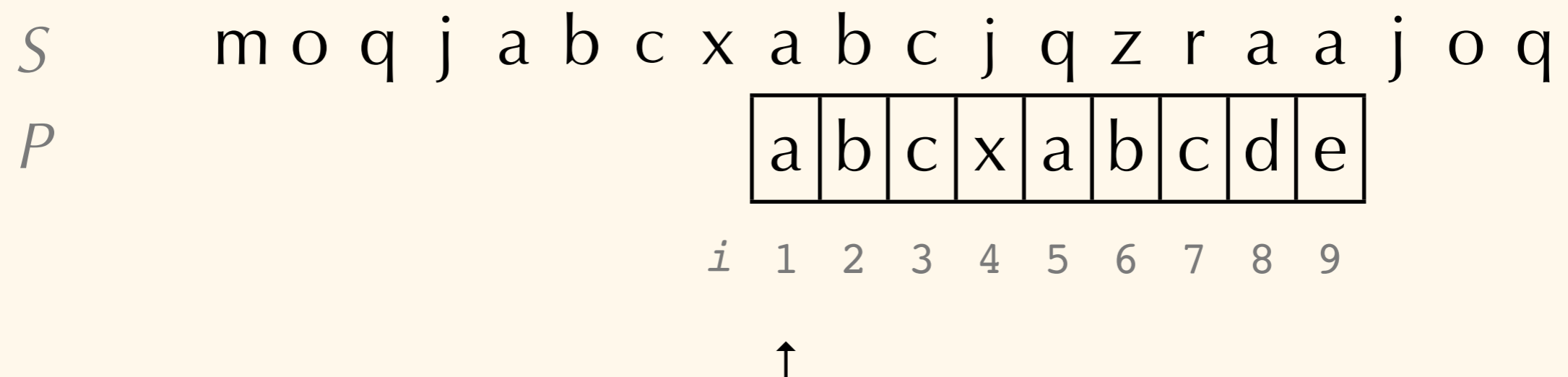
However, looking at the pattern, we can see that we *could* have shifted further without missing anything!

$S$   
 $P$

The previous algorithm would shift  $P$  down by one position, and then resume searching for a match:



However, looking at the pattern, we can see that we *could* have shifted further without missing anything!



# Definitions:

$sp_i(P)$ : The length of the longest *suffix* of  $P[1, i]$  that is also a *prefix* of  $P$ . (and let  $sp_0 = 0$ )

	x	t	p	x	t	d
$i$	1	2	3	4	5	6
$Z_i$	19	0	0	2	0	0
$sp_i$	0	0	0	0	2	0

“Failure Function”:  $F(i): sp_{i-1} + 1$

# Definitions:

$sp_i(P)$ : The length of the longest *suffix* of  $P[1, i]$  that is also a *prefix* of  $P$ . (and let  $sp_0 = 0$ )

	x	t	p	x	t	d	
$i$	1	2	3	4	5	6	7
$Z_i$	19	0	0	2	0	0	
$sp_i$	0	0	0	0	2	0	
$F(i)$	1	1	1	1	1	3	1

“Failure Function”:  $F(i): sp_{i-1} + 1$

# Computing $sp_i(P)$ can be done in linear time, using a modification of the Z-algorithm:

Initialize  $sp_i(P) = 0$  for all  $P(i)$

From left-to-right over  $P$ , hence at position  $k$  we have already calculated  $Z_{k-1}$  and have current values for  $l$  and  $r$

- If  $k > r$ , begin comparing with beginning of  $P$ . Length of match is  $Z_k$ . If  $Z_k > 0$ , then  $r = k + Z_k - 1$  and  $l = k$ .

\* If  $sp_r(P) = 0$  then set  $sp_r(P) = Z_k$

- If  $k \leq r$ , then  $P(k) = P(k')$  where  $k' = k - l - 1$

Further,  $P[k, r] = P[k', Z_l]$

Thus,  $Z_k \geq \min(Z_{k'}, |P[k, r]|)$

- If  $Z_{k'} < |P[k, r]|$ , then  $Z_k = Z_{k'}$  and  $r, l$  unchanged
- If  $Z_{k'} > |P[k, r]|$ , then  $Z_k = |P[k, r]|$  and  $r, l$  unchanged
- Otherwise, begin comparing position  $r + 1$  with  $|P[k, r]| + 1$

If mismatch at position  $q$ , then  $Z_k = q - k$ ,  $l = k$ ,  $r = q - 1$

\* If  $sp_r(P) = 0$  then set  $sp_r(P) = Z_k$

(See Gusfield for more gory details)

	x	t	p	x	t	d	
$i$	1	2	3	4	5	6	7
$Z_i$	19	0	0	2	0	0	
$sp_i$	0	0	0	0	2	0	
$F(i)$	1	1	1	1	1	3	1

We will use  $F(i)$  to tell us how far we can safely shift  $P$  along  $S$  when we encounter a mis-match.

	x	t	p	x	t	d	
$i$	1	2	3	4	5	6	7
$Z_i$	19	0	0	2	0	0	
$sp_i$	0	0	0	0	2	0	
$F(i)$	1	1	1	1	1	3	1

We will use  $F(i)$  to tell us how far we can safely shift  $P$  along  $S$  when we encounter a mis-match.

Basic idea: if we encounter a mismatch at character  $i$  of  $P$ , we can shift  $P$  down  $F(i)$  positions along  $S$ .



Demo (roark\_kmp.pdf)

# Notes on KMP:

Notes on KMP:

KMP is a classic and widely-known linear time exact-match algorithm...

# Notes on KMP:

KMP is a classic and widely-known linear time exact-match algorithm...

... and gives a nice speedup over the “simple” linear-time algorithm...

## Notes on KMP:

KMP is a classic and widely-known linear time exact-match algorithm...

... and gives a nice speedup over the “simple” linear-time algorithm...

... but for many situations, it is *not* the method of choice.

## Notes on KMP:

KMP is a classic and widely-known linear time exact-match algorithm...

... and gives a nice speedup over the “simple” linear-time algorithm...

... but for many situations, it is *not* the method of choice.

The Boyer-Moore algorithm gives better *typical* performance.

# Historical Notes:



Bob Boyer



J Strother Moore

The Boyer-Moore algorithm was developed while BB was at SRI and JSM was at Xerox PARC, and was published in 1977.

*Fun fact: Moore's first name is, in fact, the alphabetic letter "J" – it's not an abbreviation.*

Boyer-Moore improves over KMP to give sub-linear typical performance, and linear worst-case.

Key ideas:



Boyer-Moore improves over KMP to give sub-linear typical performance, and linear worst-case.

Key ideas:

1. We still move the pattern (a.k.a. the “needle”) from *left to right* through the string (the “haystack”)...

Boyer-Moore improves over KMP to give sub-linear typical performance, and linear worst-case.

Key ideas:

1. We still move the pattern (a.k.a. the “needle”) from *left to right* through the string (the “haystack”)...  
... but we start searching for matching characters from the *right* of the pattern!

Boyer-Moore improves over KMP to give sub-linear typical performance, and linear worst-case.

Key ideas:

1. We still move the pattern (a.k.a. the “needle”) from *left to right* through the string (the “haystack”)...  
... but we start searching for matching characters from the *right* of the pattern!
2. If a mis-matched character  $T$  never occurs in  $P$ , we shift  $P$  completely past that character.

Boyer-Moore improves over KMP to give sub-linear typical performance, and linear worst-case.

Key ideas:

1. We still move the pattern (a.k.a. the “needle”) from *left to right* through the string (the “haystack”)...  
... but we start searching for matching characters from the *right* of the pattern!
2. If a mis-matched character  $T$  never occurs in  $P$ , we shift  $P$  completely past that character.
3. We calculate an optimal shift amount (as in KMP), but we use suffixes rather than prefixes.

Caveat: Some of the Boyer-Moore pre-processing steps can be tricky to get one's head around.

The explanation in Chapter 2 of Gusfield is very decent, and you will ***need*** to spend some time working through it to fully grok the algorithm.

Definitions:

$R(x) =$

Definitions:

“Bad character rule”:

$R(x) =$

Definitions:

“Bad character rule”:

For each character  $x$  in  $P...$

$R(x) =$



Definitions:

“Bad character rule”:

For each character  $x$  in  $P$ ...

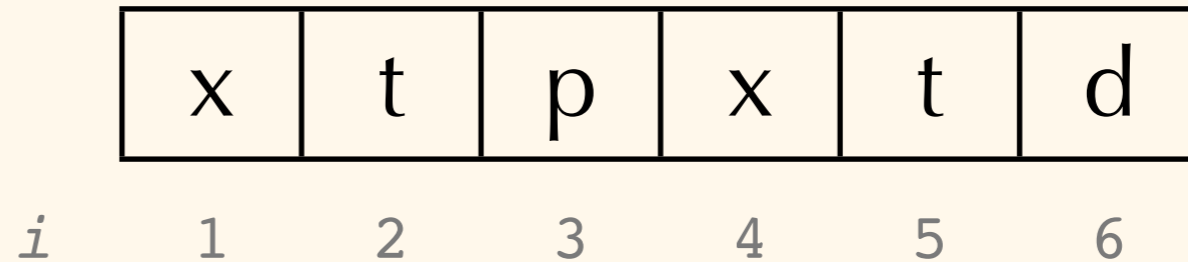
$R(x)$  = the position of the right-most  
occurrence of  $x$  in  $P$

Definitions:

“Bad character rule”:

For each character  $x$  in  $P$ ...

$R(x)$  = the position of the right-most occurrence of  $x$  in  $P$



Definitions:

“Bad character rule”:

For each character  $x$  in  $P$ ...

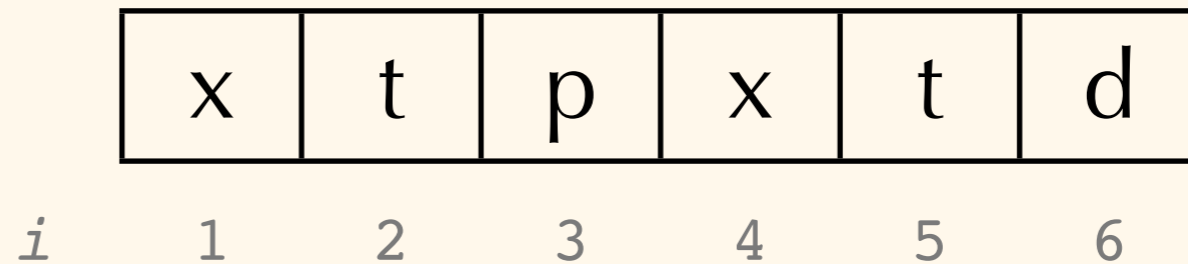
$R(x)$  = the position of the right-most occurrence of  $x$  in  $P$

$$R(x) = 4$$

$$R(t) = 5$$

$$R(p) = 3$$

$$R(d) = 6$$



“Bad character rule”:

$$R(x) = 4$$

$$R(t) = 5$$

$$R(p) = 3$$

$$R(d) = 6$$

	x	t	p	x	t	d
<i>i</i>	1	2	3	4	5	6

“Bad character rule”:

$$R(x) = 4$$

$$R(t) = 5$$

$$R(p) = 3$$

$$R(d) = 6$$

	x	t	p	x	t	d
<i>i</i>	1	2	3	4	5	6

If a mismatch occurs between character  $i$  of  $P$  and  $k$  of  $T$ , shift  $P$  by  $\max(1, i - R(T(k)))$

“Bad character rule”:

$$R(x) = 4$$

$$R(t) = 5$$

$$R(p) = 3$$

$$R(d) = 6$$

	x	t	p	x	t	d
<i>i</i>	1	2	3	4	5	6

If a mismatch occurs between character  $i$  of  $P$  and  $k$  of  $T$ , shift  $P$  by  $\max(1, i - R(T(k)))$

Furthermore, if  $T(k)$  does not appear in  $P$ , shift  $P$  by  $|P|$  (since there's no way that a match could involve  $k$ ).

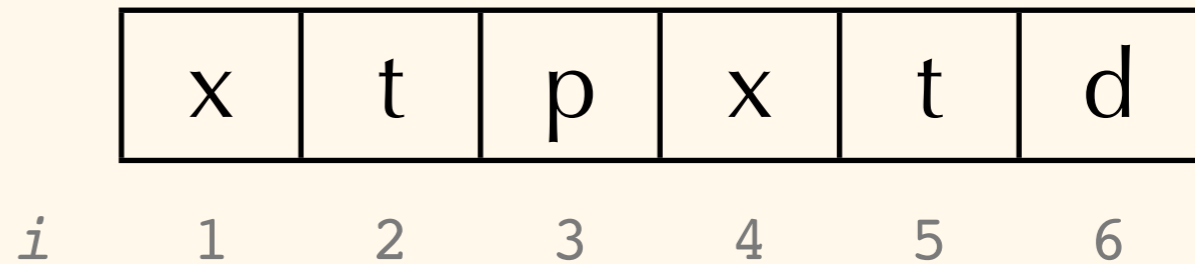
“Bad character rule”:

$$R(x) = 4$$

$$R(t) = 5$$

$$R(p) = 3$$

$$R(d) = 6$$



If a mismatch occurs between character  $i$  of  $P$  and  $k$  of  $T$ , shift  $P$  by  $\max(1, i - R(T(k)))$

Furthermore, if  $T(k)$  does not appear in  $P$ , shift  $P$  by  $|P|$  (since there's no way that a match could involve  $k$ ).

Intuition: Since we're matching from right-to-left,  $R(x)$  tells us the first place there could *possibly* be a match.

$$R(a) = 5$$

$$R(b) = 6$$

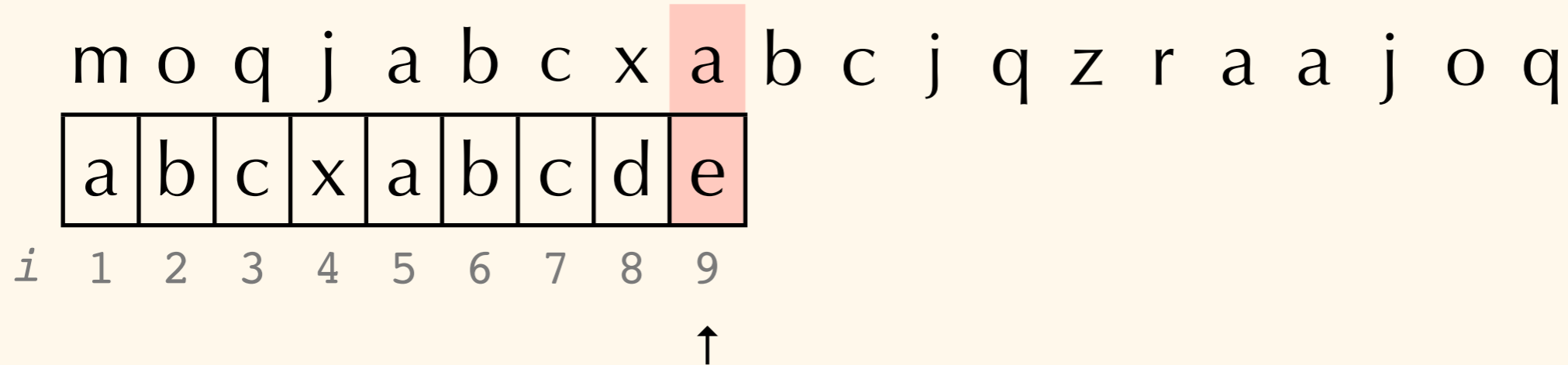
$$R(c) = 7$$

$$R(x) = 4$$

$$R(d) = 8$$

$$R(e) = 9$$

If a mismatch occurs between character  $i$  of  $P$  and  $k$  of  $T$ , shift  $P$  by  $\max(1, i - R(T(k)))$



“a” and “e” don’t match; the next *possible* location in  $S$  that a match *could* occur would involve position  $R(a)$  of the pattern being aligned with the current index.



$$R(a) = 5$$

$$R(b) = 6$$

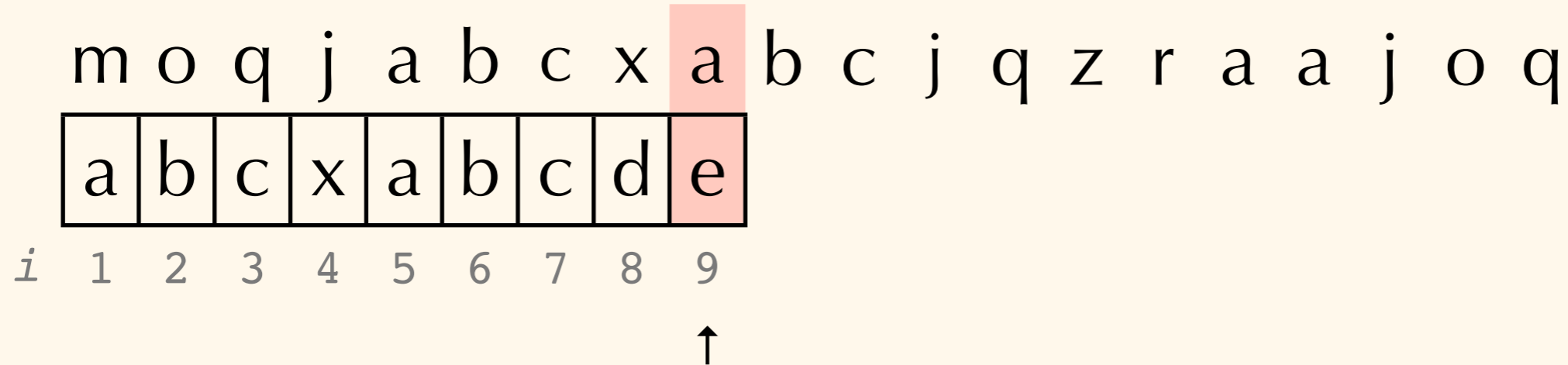
$$R(c) = 7$$

$$R(x) = 4$$

$$R(d) = 8$$

$$R(e) = 9$$

If a mismatch occurs between character  $i$  of  $P$  and  $k$  of  $T$ , shift  $P$  by  $\max(1, i - R(T(k)))$



$$R(T(k)) = R(T(9)) = R(a) = 5$$

We shift the pattern (and the search index!) by  $i - 5 = 9 - 5 = 4$  positions.

$$R(a) = 5$$

$$R(b) = 6$$

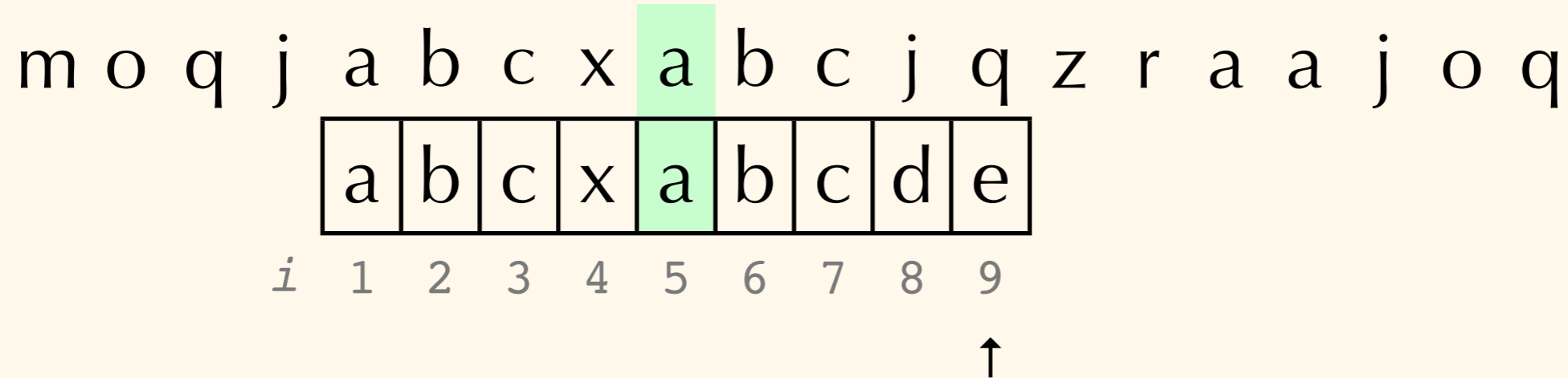
$$R(c) = 7$$

$$R(x) = 4$$

$$R(d) = 8$$

$$R(e) = 9$$

If a mismatch occurs between character  $i$  of  $P$  and  $k$  of  $T$ , shift  $P$  by  $\max(1, i - R(T(k)))$



$$R(T(k)) = R(T(9)) = R(a) = 5$$

We shift the pattern (and the search index!) by  $i - 5 = 9 - 5 = 4$  positions.

$$R(a) = 5$$

$$R(b) = 6$$

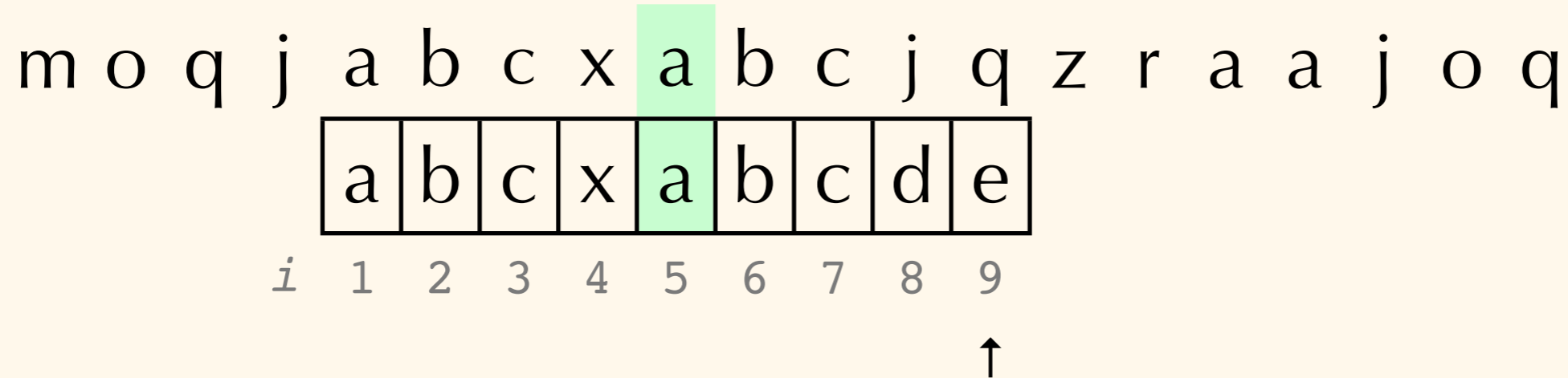
$$R(c) = 7$$

$$R(x) = 4$$

$$R(d) = 8$$

$$R(e) = 9$$

If a mismatch occurs between character  $i$  of  $P$  and  $k$  of  $T$ , shift  $P$  by  $\max(1, i - R(T(k)))$



Then, resume searching from the right-hand side of the pattern moving to the left.

$$R(a) = 5$$

$$R(b) = 6$$

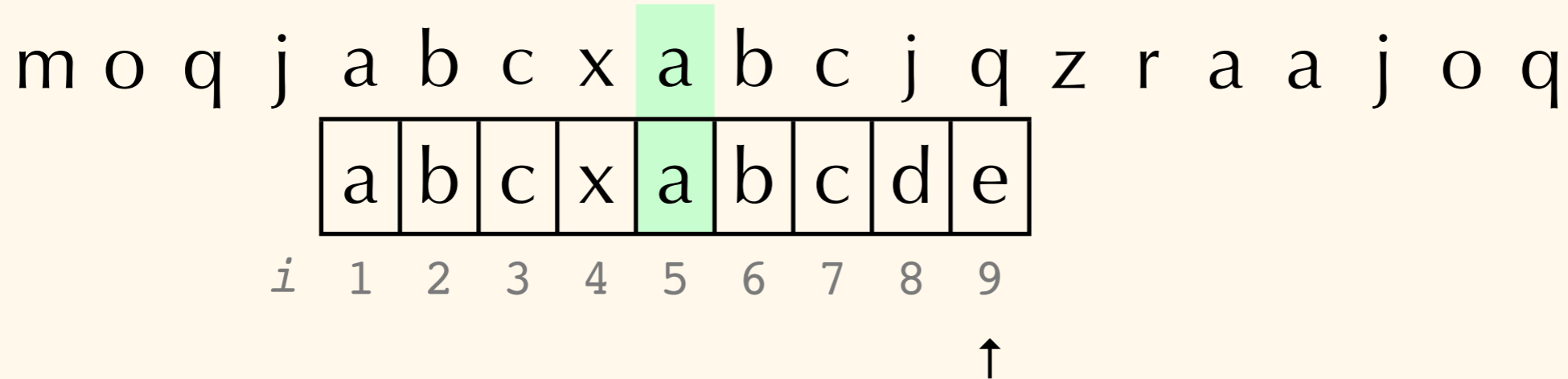
$$R(c) = 7$$

$$R(x) = 4$$

$$R(d) = 8$$

$$R(e) = 9$$

If a mismatch occurs between character  $i$  of  $P$  and  $k$  of  $T$ , shift  $P$  by  $\max(1, i - R(T(k)))$



Note that, in this case, our next shift would be quite large, since “q” does not appear in the pattern!

“Bad character rule”:

According to Gusfield, just using the bad suffix rule on its own yields good performance on “normal” English...

“Bad character rule”:

According to Gusfield, just using the bad suffix rule on its own yields good performance on “normal” English...

... but in some situations (e.g. small alphabets), it is less effective. Question for discussion: why?

“Bad character rule”:

According to Gusfield, just using the bad suffix rule on its own yields good performance on “normal” English...

... but in some situations (e.g. small alphabets), it is less effective. Question for discussion: why?

To assist, Boyer-Moore defines the “Good suffix” rule.

*“The original preprocessing method [278] for the strong good suffix rule is generally considered quite difficult and somewhat mysterious (although a weaker version of it is easy to understand). In fact, the preprocessing for the strong rule was given incorrectly in [278] and corrected, without much explanation, in [384]. Code based on [384] is given without real explanation in the text by Baase [32], but there are no published sources that try to fully explain the method.” ...*



*“The original preprocessing method [278] for the strong good suffix rule is generally considered quite difficult and somewhat mysterious (although a weaker version of it is easy to understand). In fact, the preprocessing for the strong rule was given incorrectly in [278] and corrected, without much explanation, in [384]. Code based on [384] is given without real explanation in the text by Baase [32], but there are no published sources that try to fully explain the method.” ...*

*“In contrast, the fundamental preprocessing of P discussed in Chapter 1 makes the needed preprocessing very simple.”*

Intuition of good suffix rule:

Intuition of good suffix rule:

When we hit a mismatch, we want to shift our pattern further along the string.

Intuition of good suffix rule:

When we hit a mismatch, we want to shift our pattern further along the string.

If we've found a partial match, and our pattern contains repeated elements...

Intuition of good suffix rule:

When we hit a mismatch, we want to shift our pattern further along the string.

If we've found a partial match, and our pattern contains repeated elements...

... we can shift our pattern down until the right-most repeated element\* to the left of our current position aligns with the current partial match.

Intuition of good suffix rule:

When we hit a mismatch, we want to shift our pattern further along the string.

If we've found a partial match, and our pattern contains repeated elements...

... we can shift our pattern down until the right-most repeated element\* to the left of our current position aligns with the current partial match.

In many cases, this will be a further shift than the bad-character rule would have given us!

Intuition of good suffix rule:

When we hit a mismatch, we want to shift our pattern further along the string.

If we've found a partial match, and our pattern contains repeated elements...

... we can shift our pattern down until the right-most repeated element\* to the left of our current position aligns with the current partial match.

In many cases, this will be a further shift than the bad-character rule would have given us!

\*: Some conditions apply, stay tuned...

# Illustration from Gusfield

$T$   $x$  |  $t$  |

---

$P$  before shift  $z$  |  $t'$  |  $y$  |  $t$  |

$P$  after shift  $z$  |  $t'$  |  $y$  |  $t$  |

Region of partial match





# Illustration from Gusfield

$T$   $x$  |  $t$  |

---

$P$  before shift  $z$  |  $t'$  |  $y$  |  $t$  |

$P$  after shift  $z$  |  $t'$  |  $y$  |  $t$  |

Point of mis-match



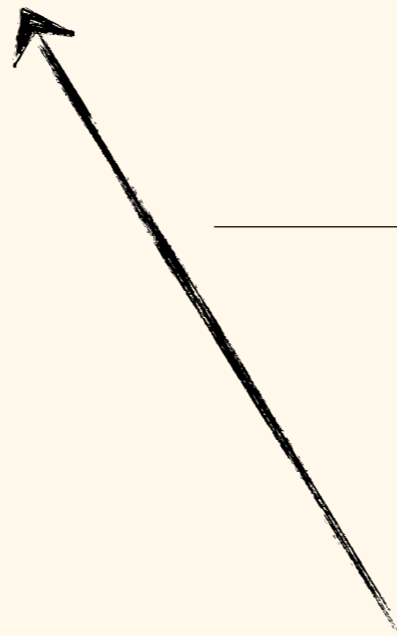
# Illustration from Gusfield

$T$   $x$  |  $t$  |

---

$P$  before shift  $z$  |  $t'$  |  $y$  |  $t$  |

$P$  after shift  $z$  |  $t'$  |  $y$  |  $t$  |



Repetition of partial-match region in pattern

# Illustration from Gusfield

$T$   $x$  |  $t$  |

---

$P$  before shift  $z$  |  $t'$  |  $y$  |  $t$  |

$P$  after shift  $z$  |  $t'$  |  $y$  |  $t$  |



Closest *possible* place for match to occur.

# Illustration from Gusfield

$T$   $x$  |  $t$  |

---

$P$  before shift  $z$  |  $t'$  |  $y$  |  $t$  |

$P$  after shift  $z$  |  $t'$  |  $y$  |  $t$  |



Point from which we start our next iteration of the algorithm.

More formally:

$t$ : A substring of  $T$  and  $P$  matching at a particular alignment.

More formally:

$t$ : A substring of  $T$  and  $P$  matching at a particular alignment.

$t'$ : The right-most copy of  $t$  in  $P$  s.t.:

More formally:

$t$ : A substring of  $T$  and  $P$  matching at a particular alignment.

$t'$ : The right-most copy of  $t$  in  $P$  s.t.:

$t'$  is *not* a suffix of  $P$  **and**:

More formally:

$t$ : A substring of  $T$  and  $P$  matching at a particular alignment.

$t'$ : The right-most copy of  $t$  in  $P$  s.t.:

$t'$  is *not* a suffix of  $P$  **and**:

the character to the left of  $t'$  in  $P$  differs from the character to the left of  $t$  in  $P$ .



More formally:

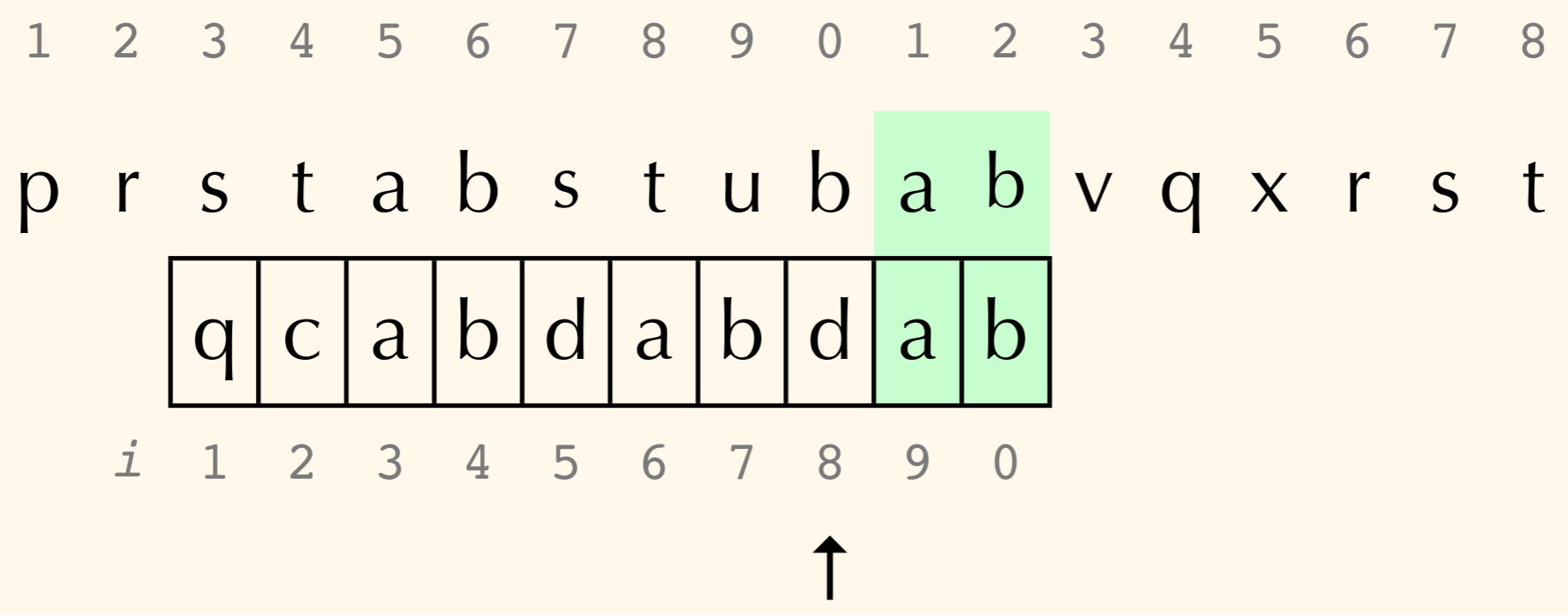
$t$ : A substring of  $T$  and  $P$  matching at a particular alignment.

$t'$ : The right-most copy of  $t$  in  $P$  s.t.:

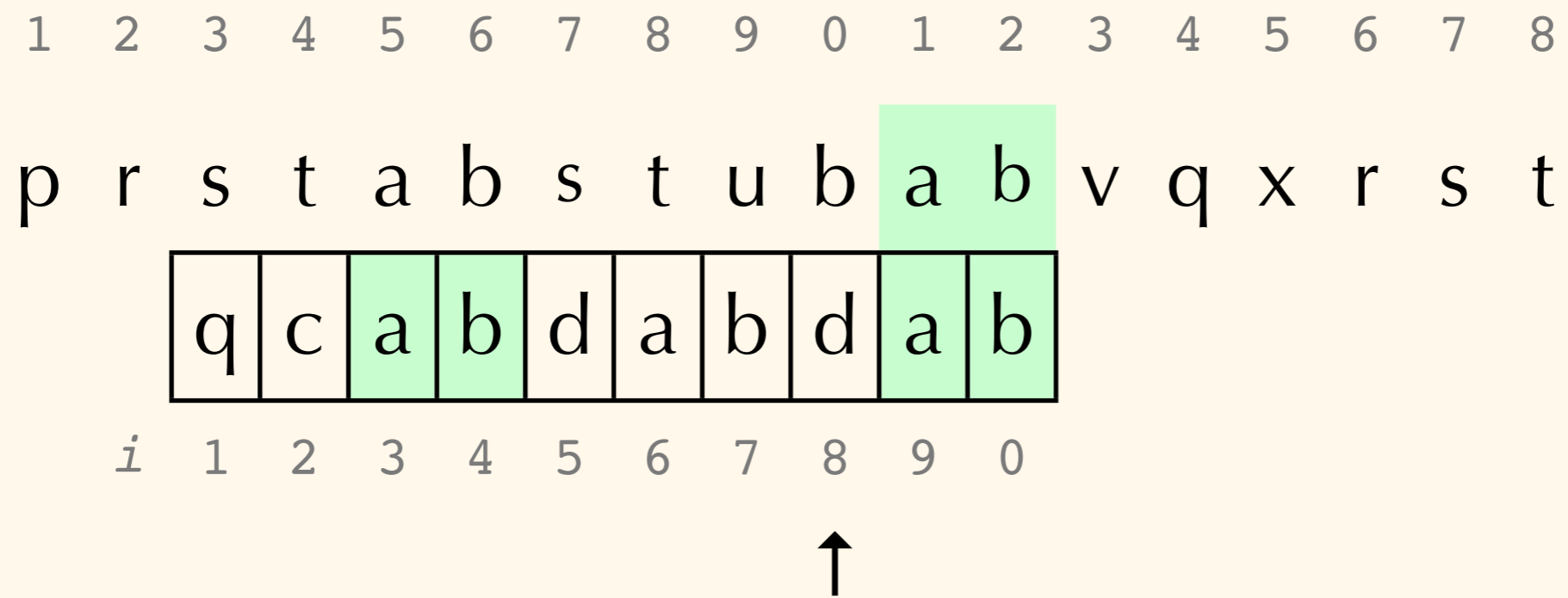
$t'$  is *not* a suffix of  $P$  **and**:

the character to the left of  $t'$  in  $P$  differs from the character to the left of  $t$  in  $P$ .

If  $t'$  exists, shift  $P$  so that  $t'$  in  $P$  is below substring  $t$  in  $T$ .

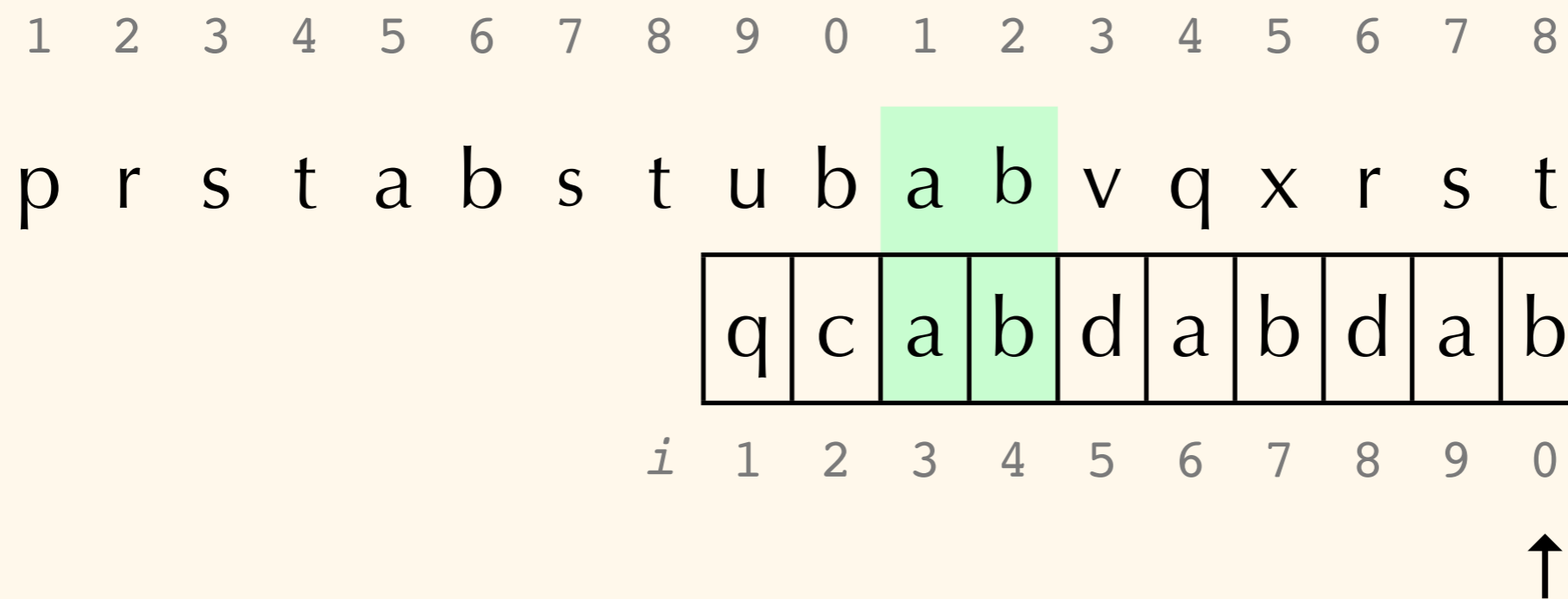


*t*: ab



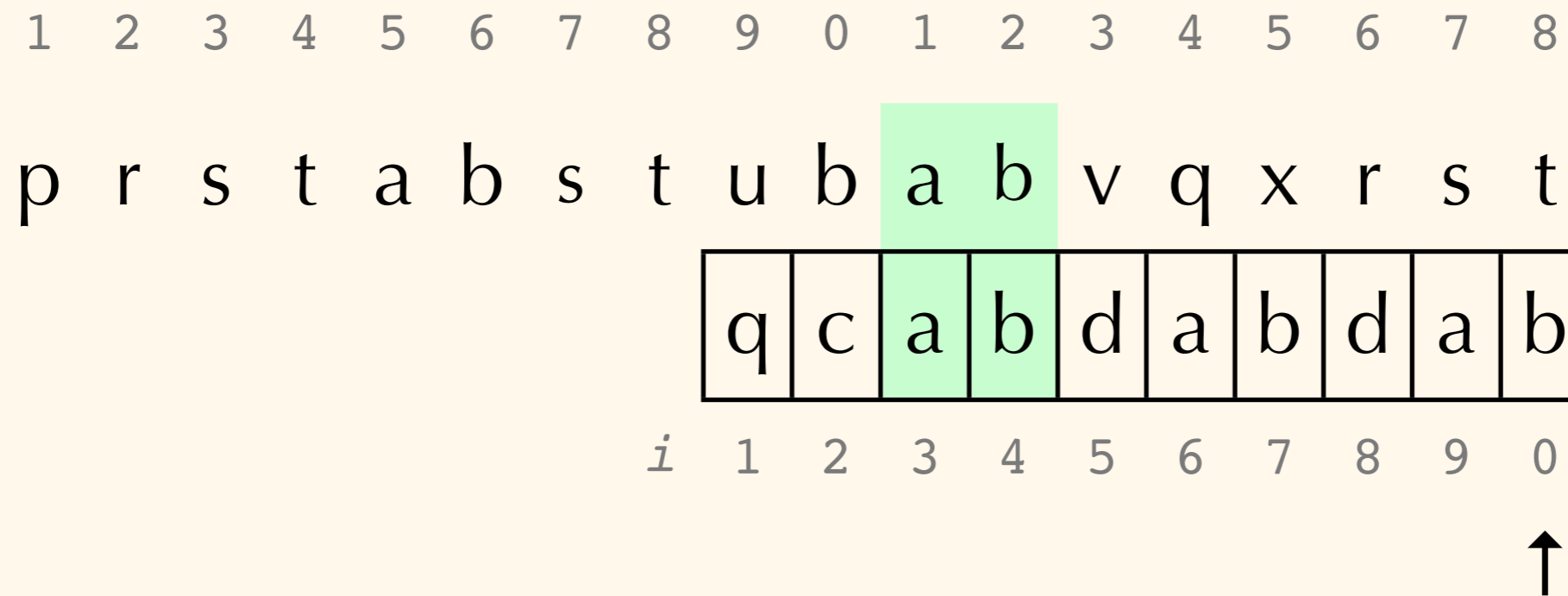
$t$ : ab

$t'$ : occurs at position 3 in  $P$



$t$ : ab

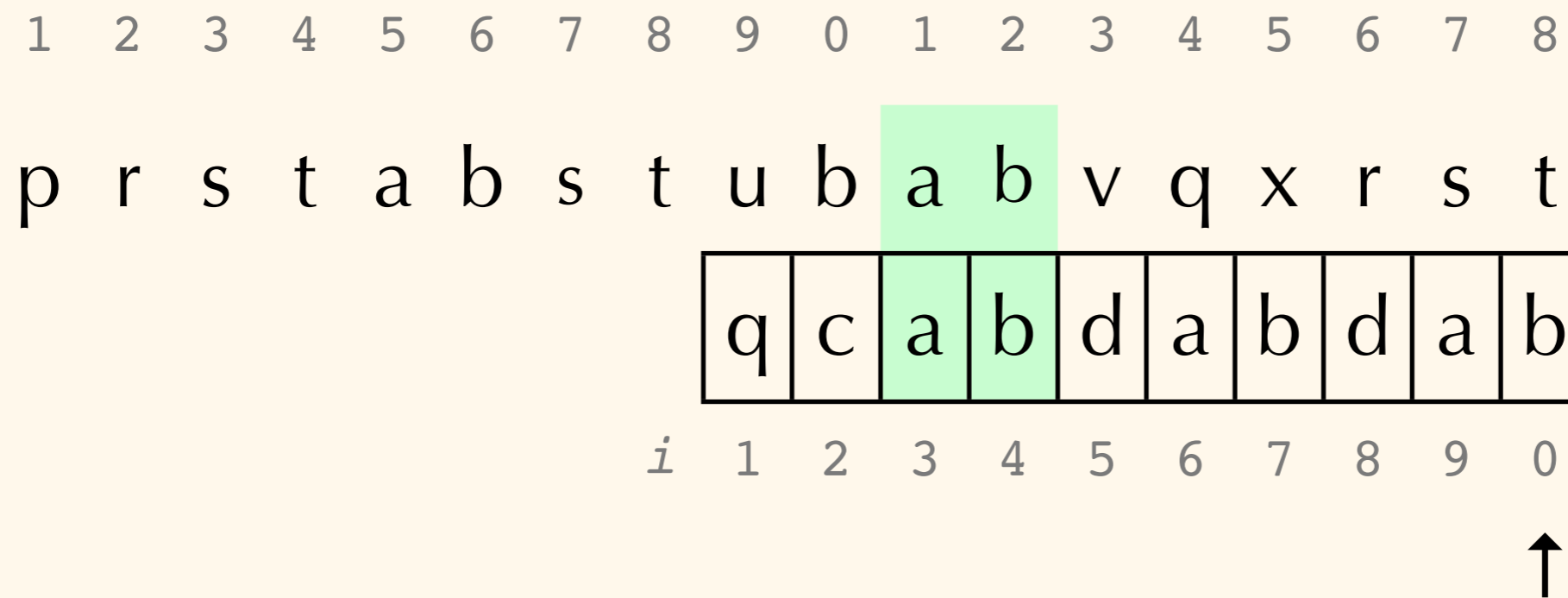
$t'$ : occurs at position 3 in  $P$



$t$ : ab

$t'$ : occurs at position 3 in  $P$

Note that had we relied on the “bad character” rule, we only would have been able to shift down 1 position!



$t$ : ab

$t'$ : occurs at position 3 in  $P$

Great! How do we find  $t'$ ?

# Definitions:

For each position  $i$ , let  $k = |P| - i + 1$

(i.e.,  $k = |P[i, |P|]|$ )

(i.e.,  $k =$  the length of the suffix of  $P$  starting at  $i$ .)

# Definitions:

For each position  $i$ , let  $k = |P| - i + 1$



# Definitions:

For each position  $i$ , let  $k = |P| - i + 1$

Let  $L(i)$  = the largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$ .

# Definitions:

For each position  $i$ , let  $k = |P| - i + 1$

Let  $L(i)$  = the largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$ .

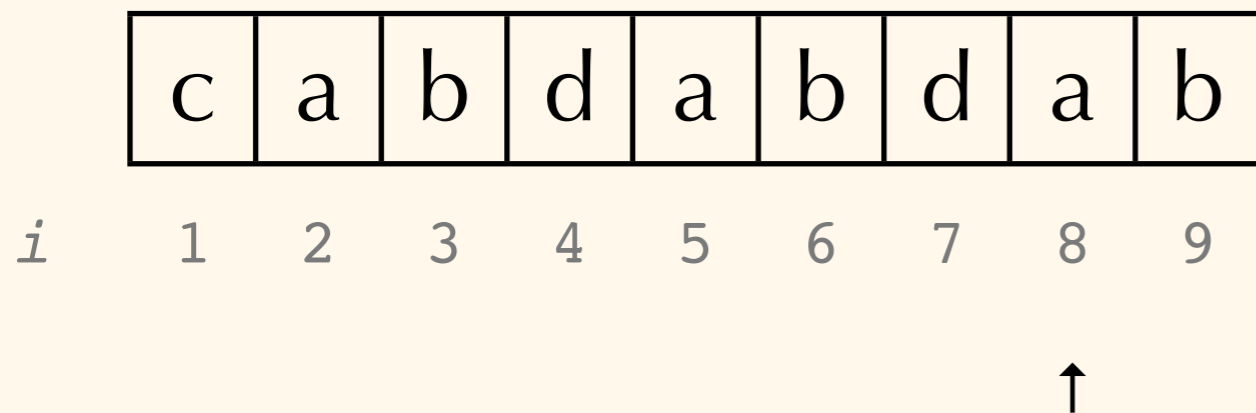
Let  $L'(i)$  = largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$  *and* s.t. the character preceding that suffix is not equal to  $P(i - 1)$ .

# Definitions:

For each position  $i$ , let  $k = |P| - i + 1$

Let  $L(i)$  = the largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$ .

Let  $L'(i)$  = largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$  *and* s.t. the character preceding that suffix is not equal to  $P(i - 1)$ .

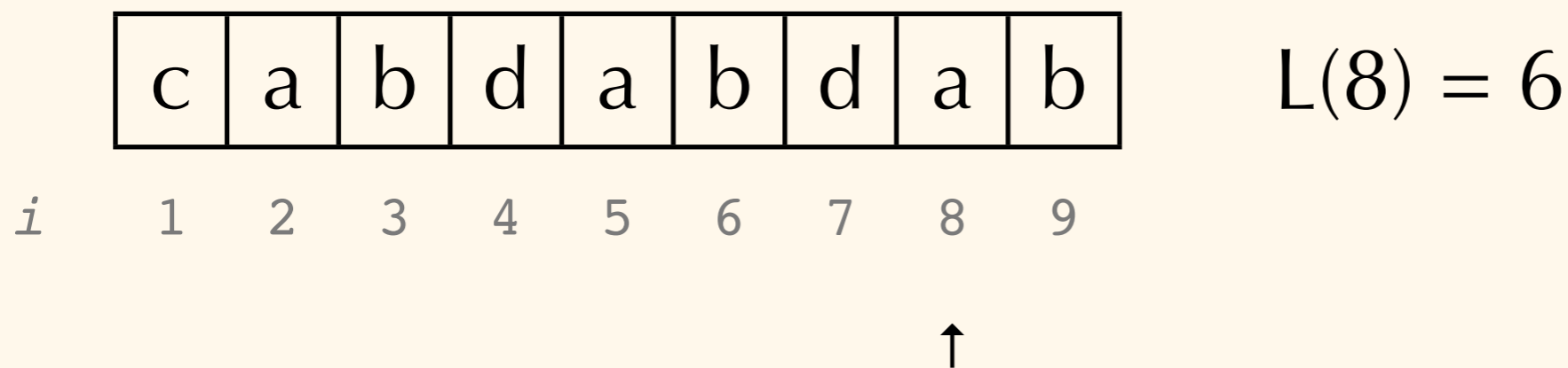


# Definitions:

For each position  $i$ , let  $k = |P| - i + 1$

Let  $L(i)$  = the largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$ .

Let  $L'(i)$  = largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$  *and* s.t. the character preceding that suffix is not equal to  $P(i - 1)$ .

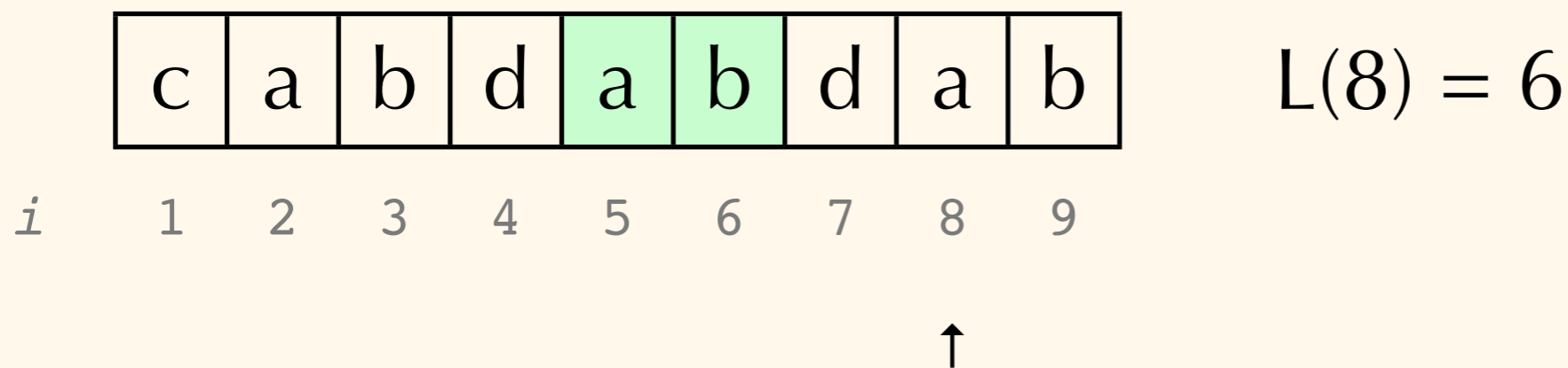


# Definitions:

For each position  $i$ , let  $k = |P| - i + 1$

Let  $L(i)$  = the largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$ .

Let  $L'(i)$  = largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$  *and* s.t. the character preceding that suffix is not equal to  $P(i - 1)$ .





## Definitions:

For each position  $i$ , let  $k = |P| - i + 1$

Let  $L(i)$  = the largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$ .

Let  $L'(i)$  = largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$  *and* s.t. the character preceding that suffix is not equal to  $P(i - 1)$ .

We'll use  $L(i)$  and  $L'(i)$  for detailed steps in the algorithm...

## Definitions:

For each position  $i$ , let  $k = |P| - i + 1$

Let  $L(i)$  = the largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$ .

Let  $L'(i)$  = largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$  *and* s.t. the character preceding that suffix is not equal to  $P(i - 1)$ .

We'll use  $L(i)$  and  $L'(i)$  for detailed steps in the algorithm...

... but note that any position  $i$  in  $P$  where  $L_i(P) > 0$  must have a corresponding repeating segment! ( $t'$  in previous slides)



## Definitions:

For each position  $i$ , let  $k = |P| - i + 1$

Let  $L(i)$  = the largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$ .

Let  $L'(i)$  = largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$  *and* s.t. the character preceding that suffix is not equal to  $P(i - 1)$ .

We'll use  $L(i)$  and  $L'(i)$  for detailed steps in the algorithm...

Gusfield chap. 2 gives a lovely algorithm for computing  $L(i)$  and  $L'(i)$  in  $O(|P|)$  time by using Z-boxes!

## Definitions:

For each position  $i$ , let  $k = |P| - i + 1$

Let  $L(i) =$  the largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$ .

Let  $L'(i) =$  largest position from which  $P[i,]$  matches a suffix of  $P[1, L(i)]$  *and* s.t. the character preceding that suffix is not equal to  $P(i - 1)$ .

One more definition:  $l(i)$ .

$l(i)$  is the length of the largest suffix of  $P(i,)$  that is also a prefix of  $P$ , and can also be calculated in linear time (see Gusfield).

Putting it all together:

1. Pre-calculate bad character table;
2. Pre-calculate good-suffix table;
3. Start at position  $|P|$  in  $T$ , move from right to left.
  - a. Look for matching characters;
  - b. If no match, skip  $\max(\text{bad character, good suffix})$  positions further down in the string.
4. Wash, Rinse, Repeat!

Putting it all together: [roark\\_boyer\\_moore.pdf](#)