

Indexing the Bijective BWT

Hideo Bannai (Kyushu University),

Juha Kärkkäinen (Helsinki Institute of Information Technology),

Dominik Köppl (Kyushu University),

Marcin Piątkowski (Nicolaus Copernicus University)

This presentation received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 690941.

FM Index

ingredients

- BWT
- wavelet tree

FM Index

ingredients

- BWT
- wavelet tree

operation: backward search

- locate pattern
- time independent on number of occurrences
- $O(|P|)$ rank/select for pattern P

FM Index on bijective BWT

ingredients

- bijective BWT
- wavelet tree

operation: backward search

- locate pattern
- time independent on number of occurrences

$O(|P| \lg |P|)$ rank/select for pattern P

bijjective BWT is
the BWT of
the Lyndon factorization
of an input text
with respect to \prec_{ω}

bijjective BWT is
the BWT of

the Lyndon factorization

1.

of an input text

with respect to

\prec_{ω}

2.

Lyndon words

- a
- aabab

Lyndon word is smaller than

- any proper suffix
- any rotation

Lyndon words

- a
- aabab

Lyndon word is smaller than

- any proper suffix
- any rotation

not Lyndon words:

- abaab (rotation aabab smaller)
- abab (abab not smaller than suffix ab)

Lyndon factorization [Chen+ '58]

- input: text T
- output: factorization $T_1 \dots T_t$ with
 - T_i is Lyndon word
 - $T_x \geq_{\text{lex}} T_{x+1}$
 - factorization uniquely defined
 - linear time [Duval'88]

properties [Duval' 88]

- T_t :
 - smallest Lyndon word
 - smallest suffix of T
- T_x primitive
- T_1 longest Lyndon prefix of $T[1..]$
- T_{x+1} longest Lyndon prefix of $T[|T_1 \cdots T_x|+1..]$

\prec_{ω}

- $u \prec_{\omega} w \iff uuuu\dots \prec_{\text{lex}} wwww\dots$
- $ab \prec_{\text{lex}} aba$
- $aba \prec_{\omega} ab$

\prec_{ω}

• $u \prec_{\omega} w \iff uuuu\dots \prec_{\text{lex}} wwww\dots$

• $ab \prec_{\text{lex}} aba$

ab**a**babab...

• $aba \prec_{\omega} ab$

aba**a**baaba...

bijjective BWT of senescence

s | enes | cen | ce

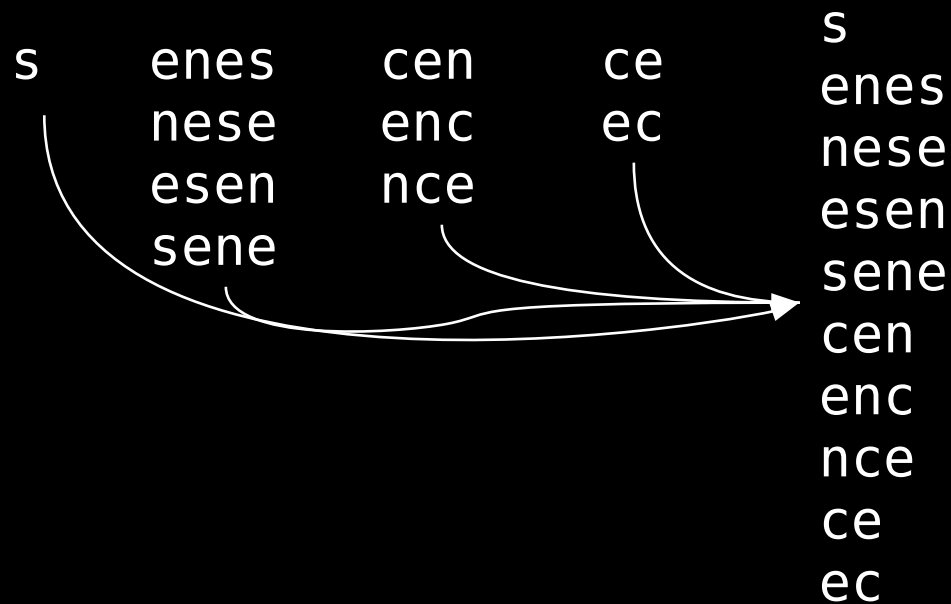
bijjective BWT of senescence

s | enes | cen | ce

s	enes	cen	ce
	nese	enc	ec
	esen	nce	
	sene		

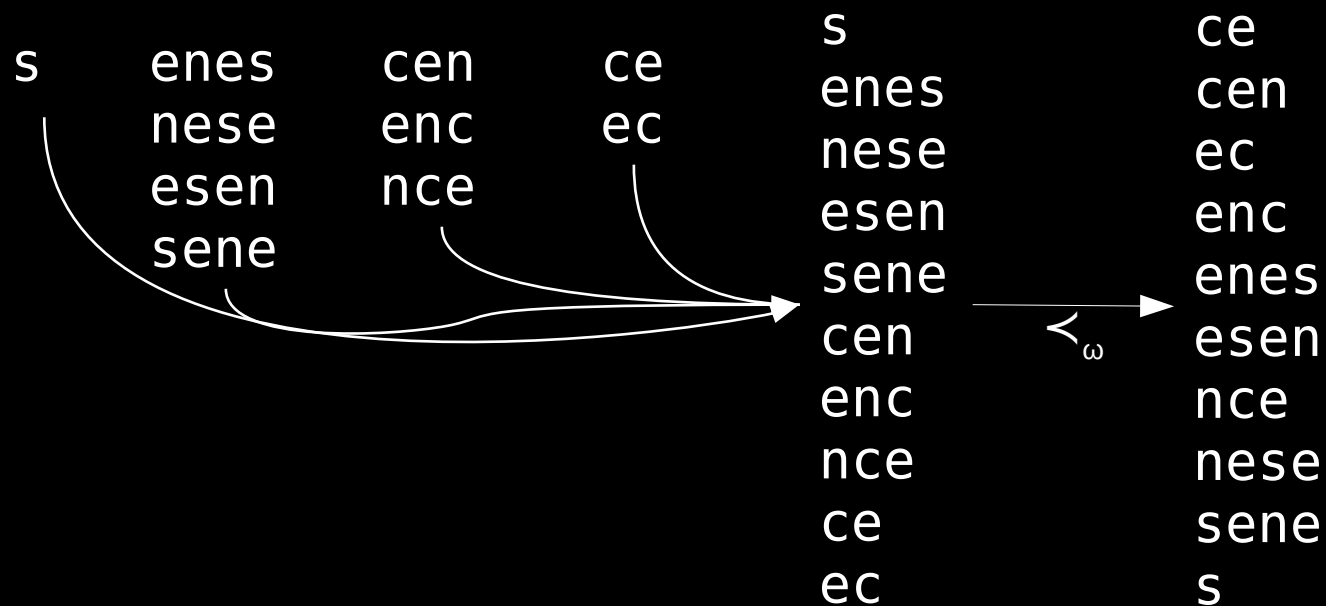
bijective BWT of senescence

s | enes | cen | ce



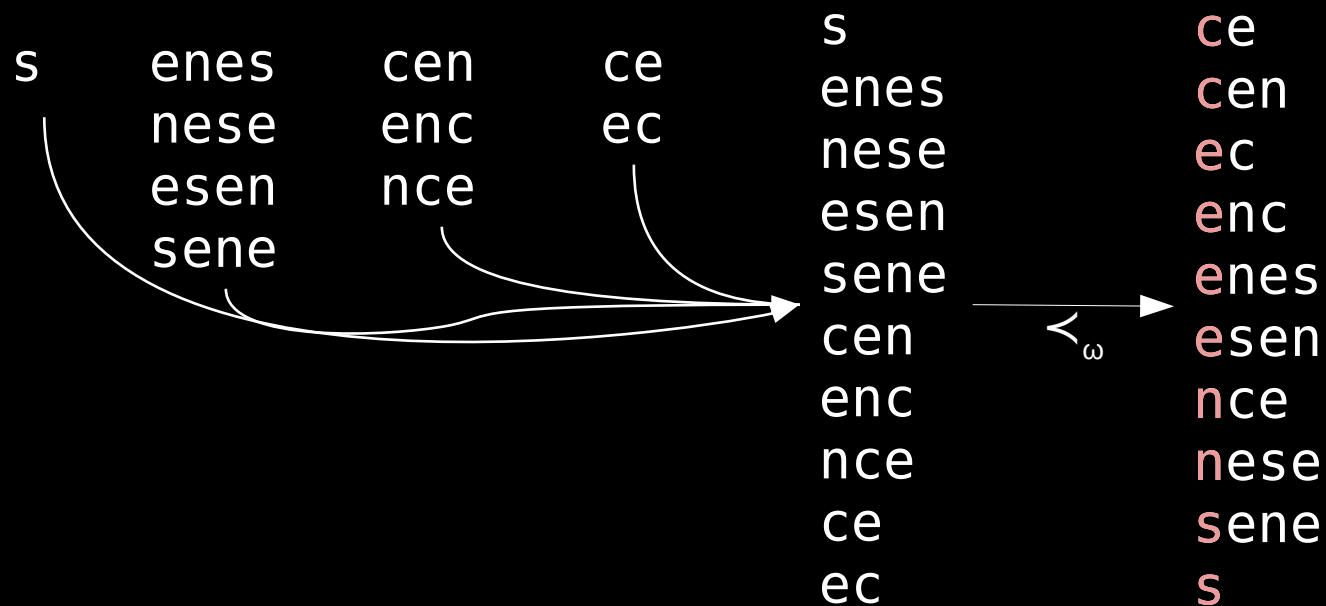
bijjective BWT of senescence

s | enes | cen | ce



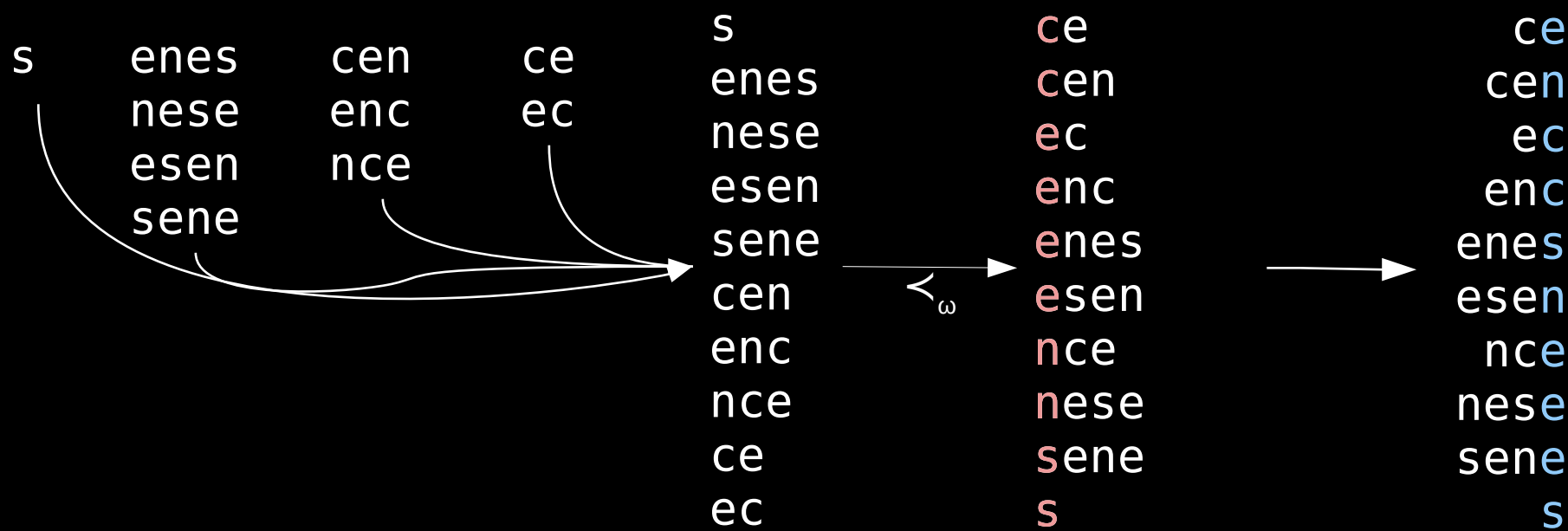
bijjective BWT of senescence

s | enes | cen | ce



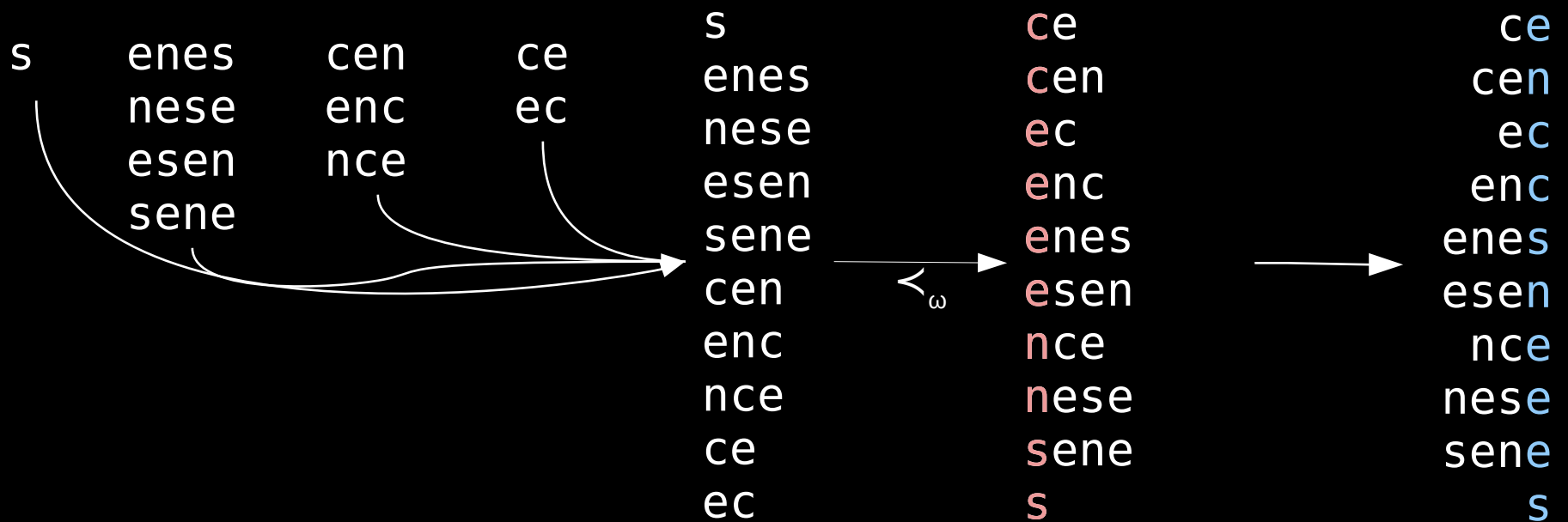
bijjective BWT of senescence

s | enes | cen | ce



bijjective BWT of senescence

s | enes | cen | ce



result: encsneees

cycles

L

e

n

c

c

s

n

e

e

s

s

F

c

c

e

e

e

e

n

n

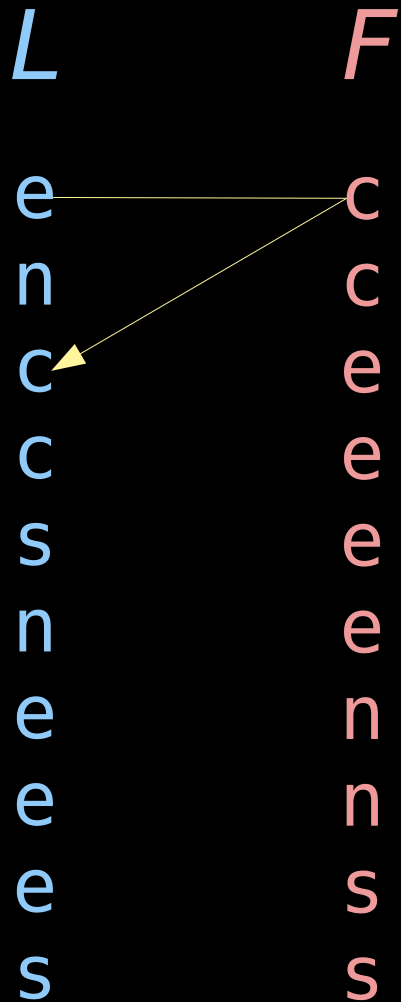
s

s

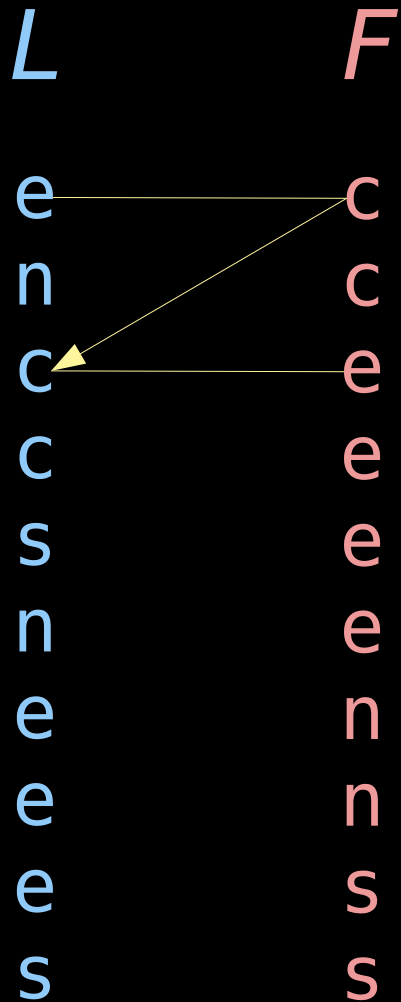
cycles

<i>L</i>	<i>F</i>
e	c
n	c
c	e
c	e
s	e
n	e
e	n
e	n
s	s
s	s

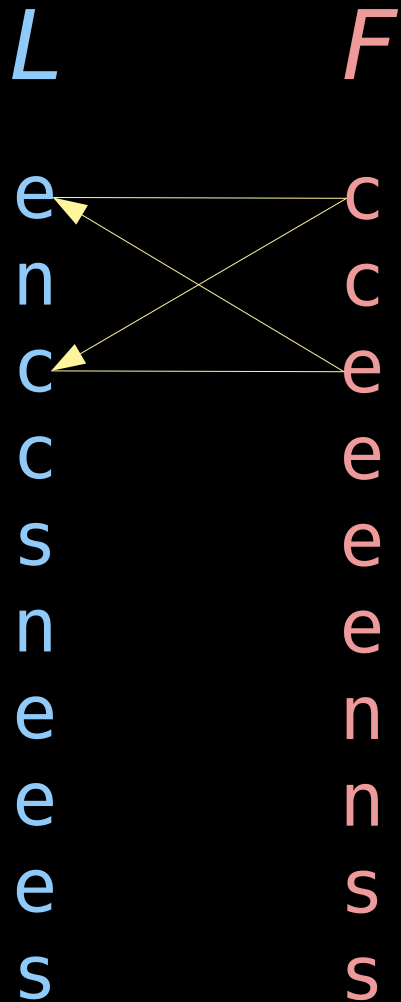
cycles



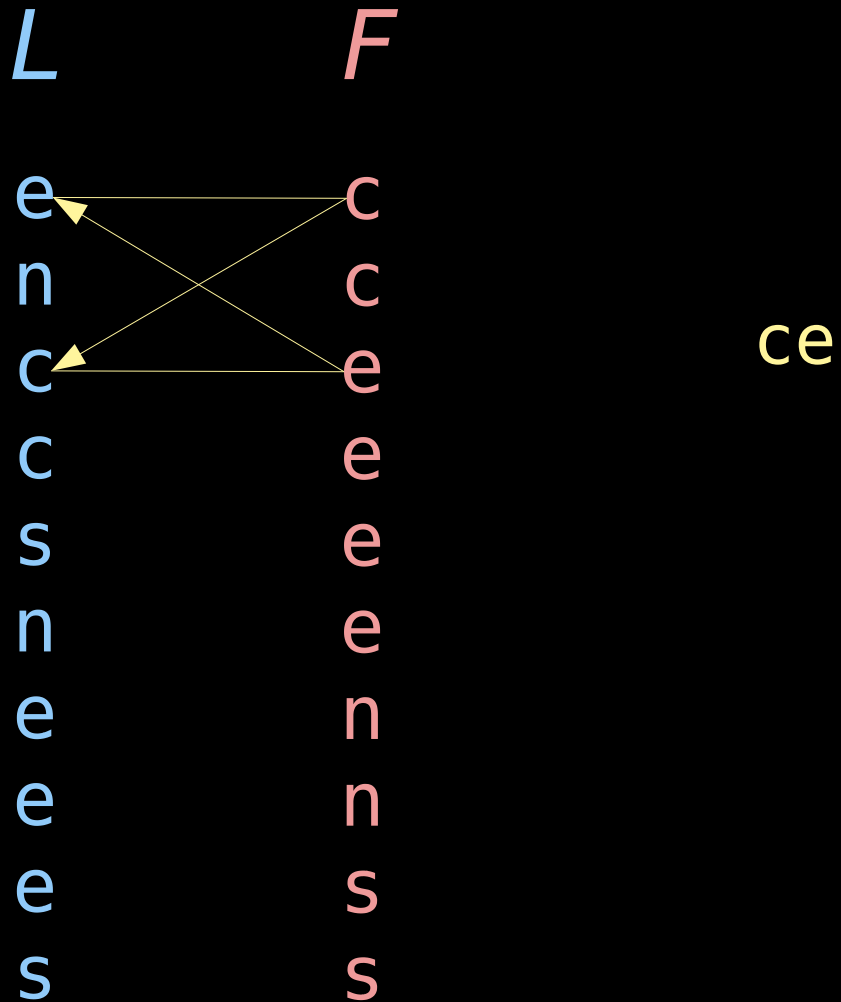
cycles



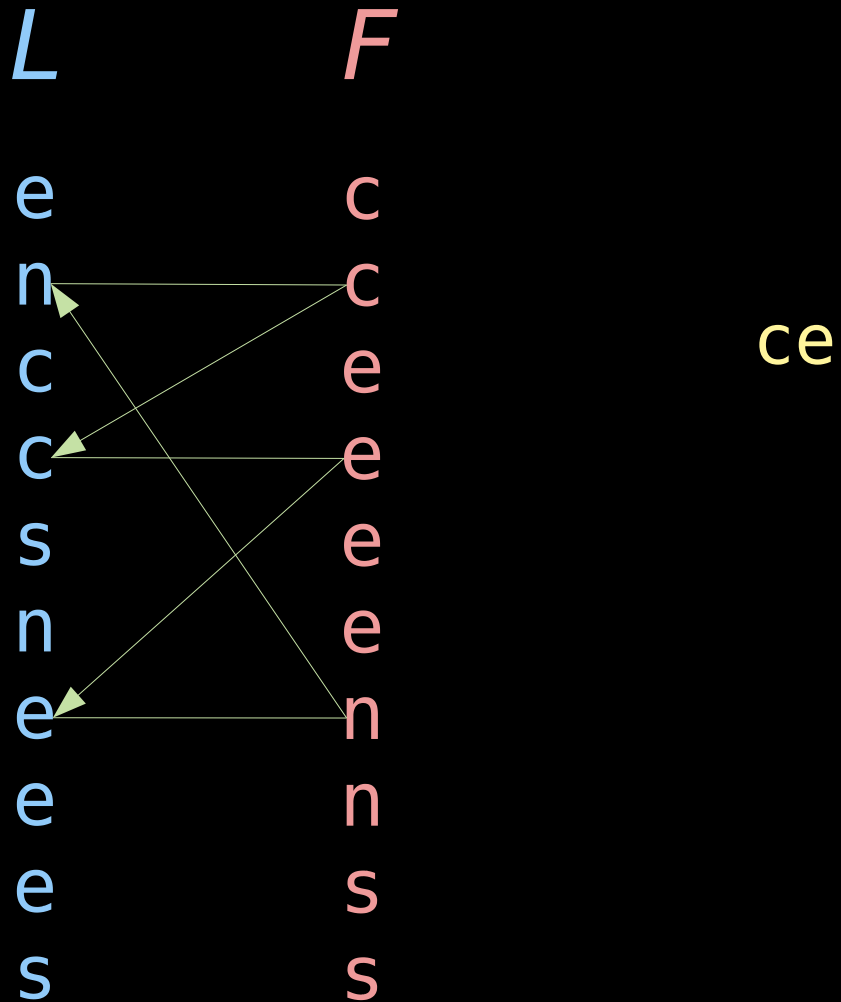
cycles



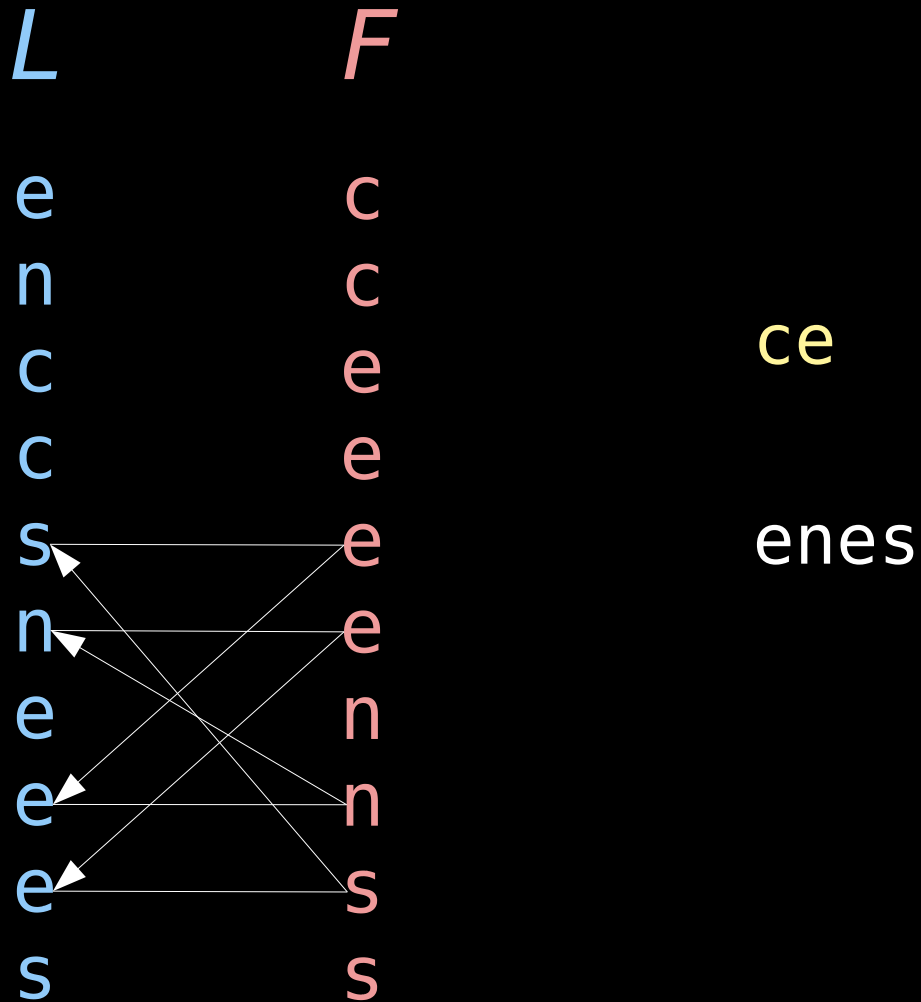
cycles



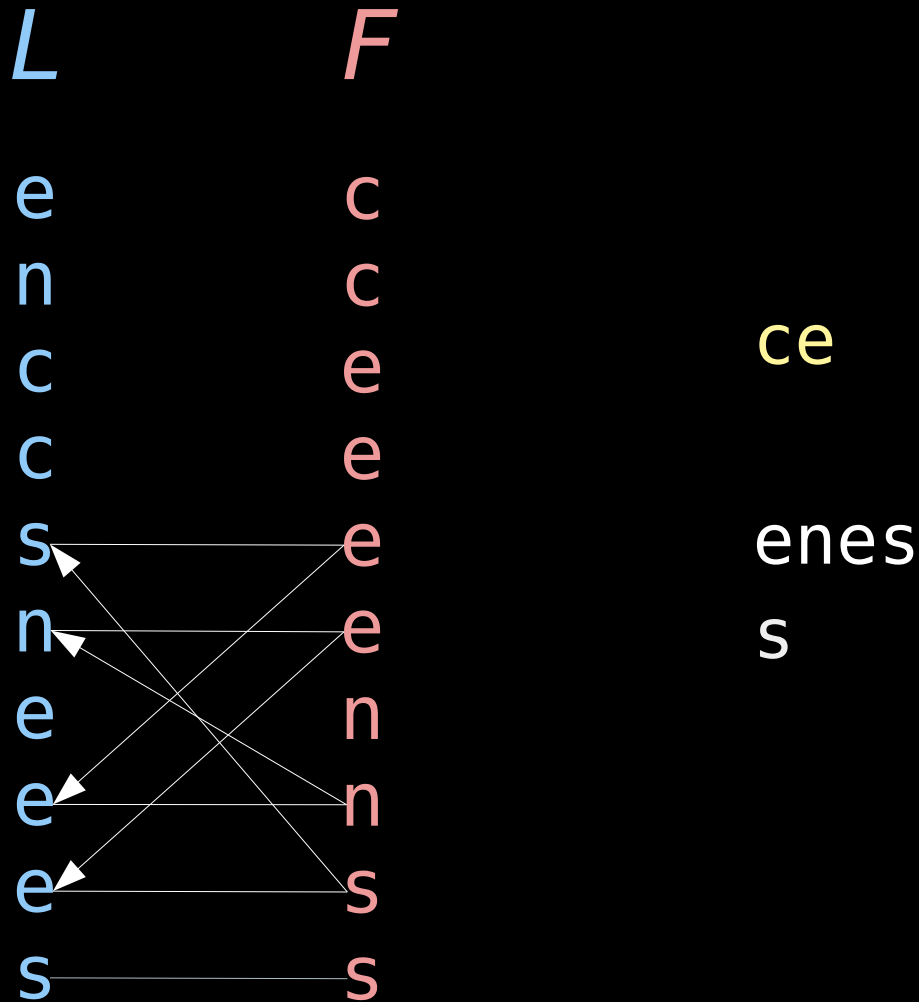
cycles



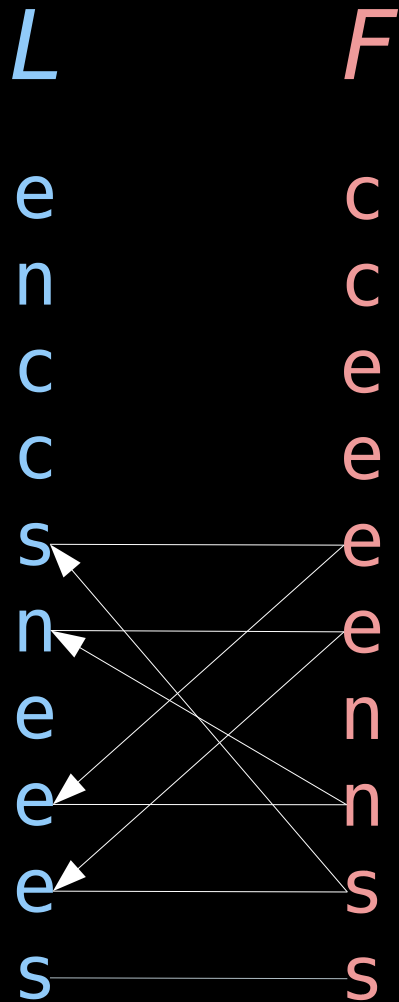
cycles



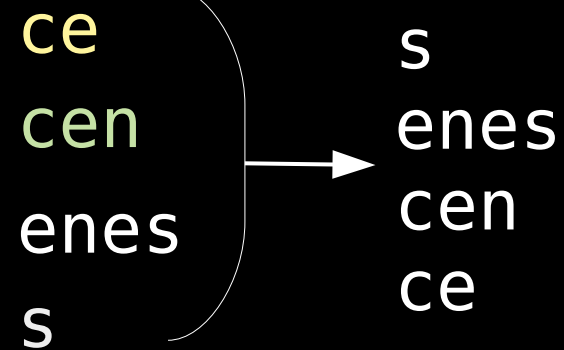
cycles



cycles



sort lexico. reversely



backward search 'cen'

s enes cen ce	<i>F</i>	<i>L</i>
	c	ce
	c	cen
	e	ec
	e	enc
	e	enes
	e	esen
	n	nce
	n	nese
	s	sene
	s	s

backward search 'cen'

s enes cen ce	<i>F</i>	<i>L</i>
	c	ce
	c	cen
	e	ec
	e	enc
	e	enes
	e	esen
	(n)	nce
	(n)	nese
	s	sene
	s	s

backward search 'cen'

s | enes | cen | ce

F

L

c

ce

c

cen

e

ec

(e)
(e)

enc

enes

e

esen

(n)
(n)

nce

nese

s

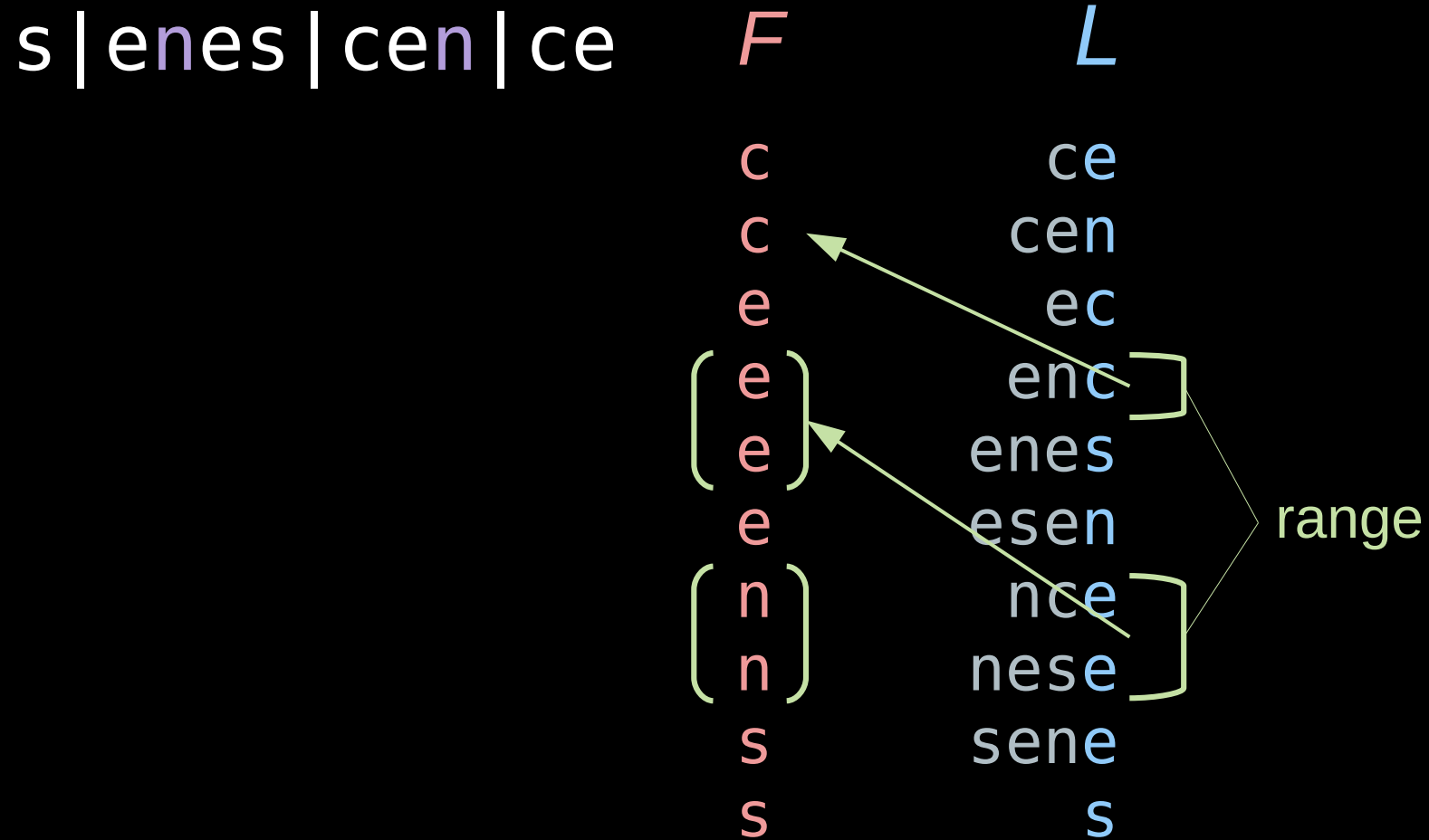
sene

s

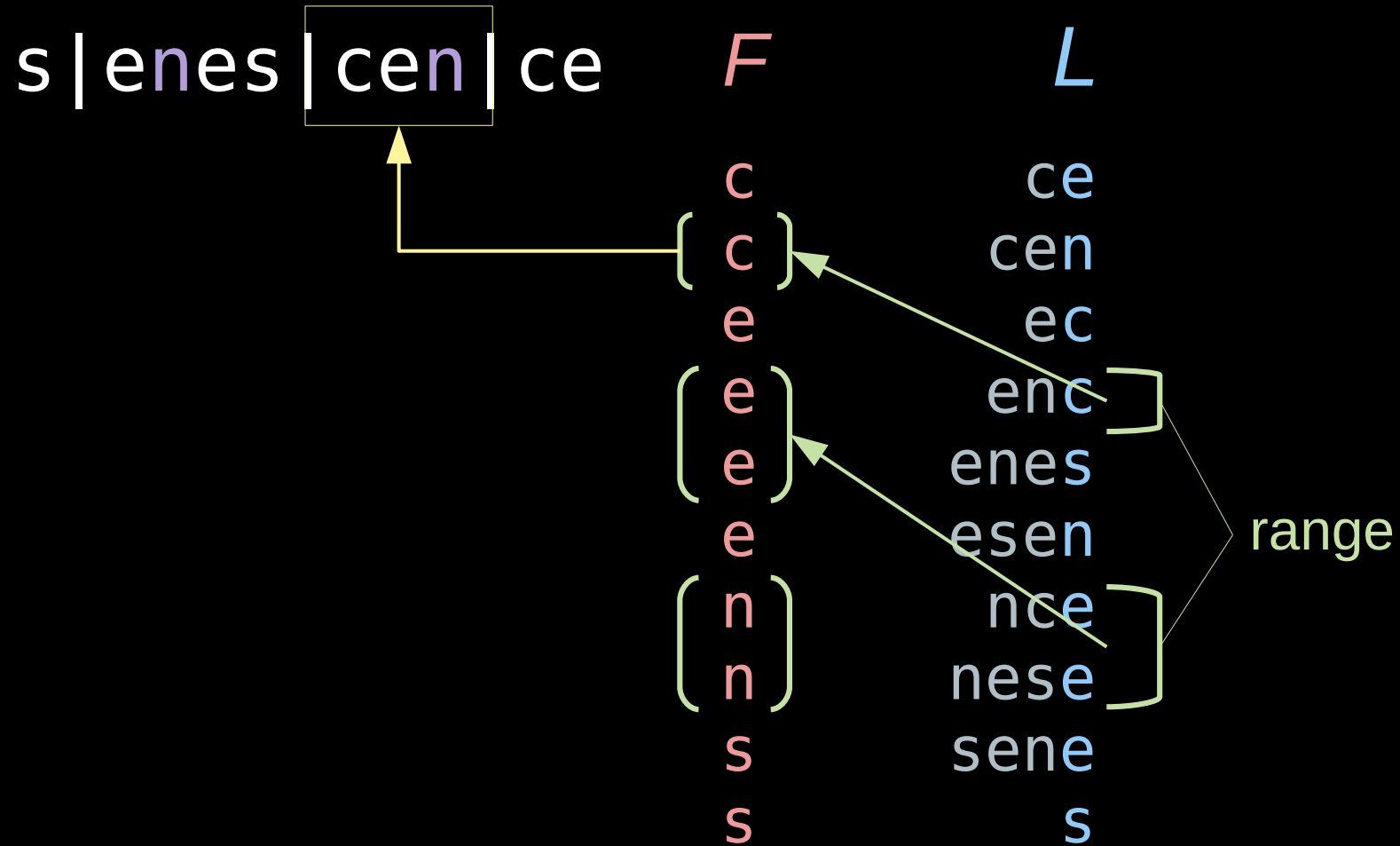
s

range

backward search 'cen'



backward search 'cen'



backward search 'ss'

<i>s</i> <i>enes</i> <i>cen</i> <i>ce</i>	<i>F</i>	<i>L</i>
	<i>c</i>	<i>ce</i>
	<i>c</i>	<i>cen</i>
	<i>e</i>	<i>ec</i>
	<i>e</i>	<i>enc</i>
	<i>e</i>	<i>enes</i>
	<i>e</i>	<i>esen</i>
	<i>n</i>	<i>nce</i>
	<i>n</i>	<i>nese</i>
	<i>s</i>	<i>sene</i>
	<i>s</i>	<i>s</i>

backward search 'ss'

<i>s</i> <i>enes</i> <i>cen</i> <i>ce</i>	<i>F</i>	<i>L</i>
	c	ce
	c	cen
	e	ec
	e	enc
	e	enes
	e	esen
	n	nce
	n	nese
	(<i>s</i>)	sene
	(<i>s</i>)	s

backward search 'ss'

s enes cen ce	<i>F</i>	<i>L</i>
	c	ce
	c	cen
	e	ec
	e	enc
	e	enes
	e	esen
	n	nce
	n	nese
	(s)	sene
	(s)	s]

backward search 'ss'

s | enes | cen | ce

F

L

c

ce

c

cen

e

ec

e

enc

e

enes

e

esen

n

nce

n

nese

s

sene



backward search 'ss'

s | enes | cen | ce

F

L

c

ce

c

cen

e

ec

e

enc

e

enes

e

esen

n

nce

n

nese

s

sene



- cen is Lyndon word
- ss is **not**

pattern is Lyndon word

$T =$

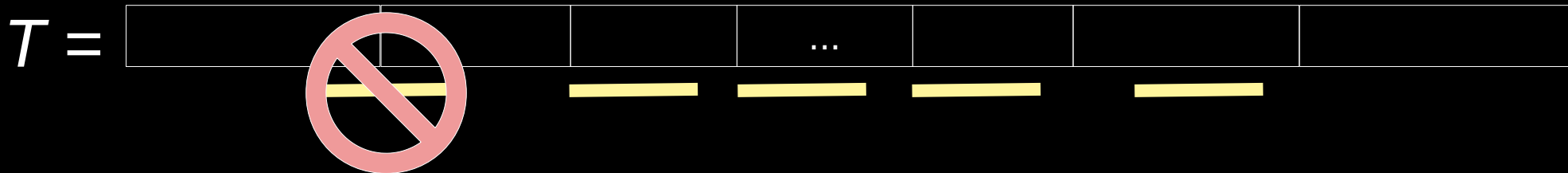
			...			
--	--	--	-----	--	--	--

pattern is Lyndon word



pattern is Lyndon word

cannot cross Lyndon factor border



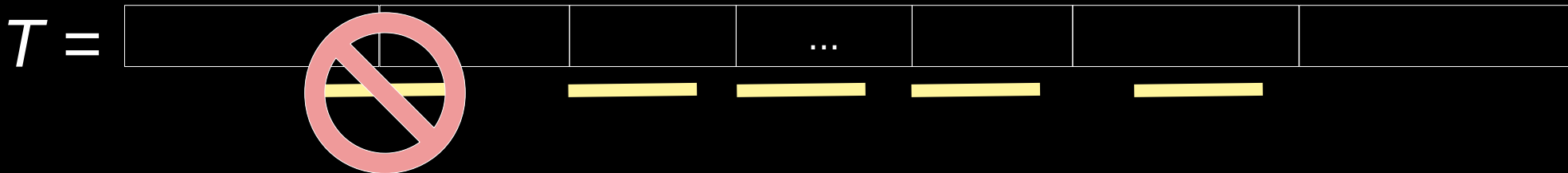
pattern is Lyndon word

cannot cross Lyndon factor border

⇒ occur inside factors

⇒ found within cycles

backward search \cong FM-index



pattern P is not a Lyndon word

- Lyndon factorization: $P = P_1 \cdots P_m$
- P_y substring of T_x or equal to T_x

algorithm:

- search P_m
- take care when starting with P_{m-1} !

- backward search $P = se$

$s | enes | cen | ce$

<i>F</i>	<i>L</i>
c	ce
c	cen
e	ec
e	enc
e	enes
e	esen
n	nce
n	nese
s	sene
s	s

• backward search $P = se$

$s | enes | cen | ce$

• $P_2 = e$

F

L

c ce

c cen

e ec

e enc

e enes

e esen

n nce

n nese

s sene

s s

• backward search $P = se$

$s | enes | cen | ce$

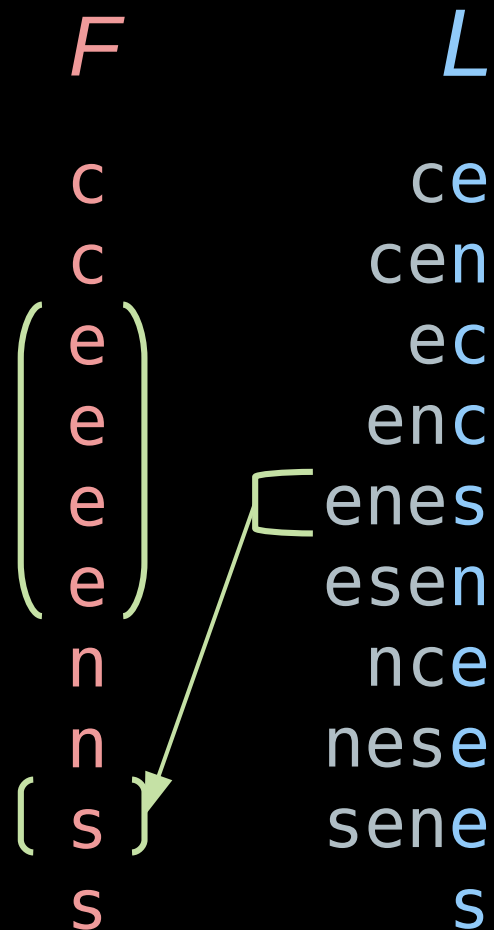
• $P_2 = e$

F	L
c	ce
c	cen
e	ec
e	enc
e	enes
e	esen
n	nce
n	nese
s	sene
s	s

• backward search $P = se$

• $P_2 = e$

• $P_1 = s$

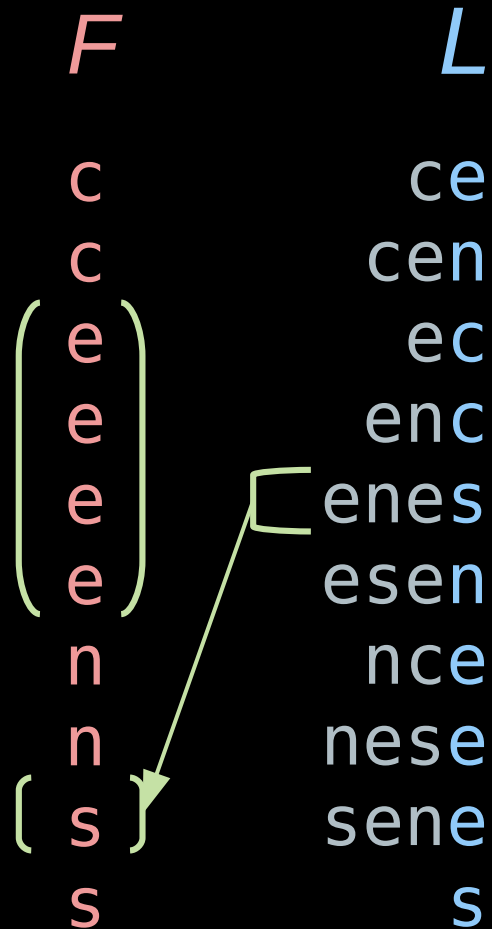


s | enes | cen | ce

• backward search $P = se$

• $P_2 = e$

• $P_1 = s$



s | enes | cen | ce

found

s | enes | cen | ce

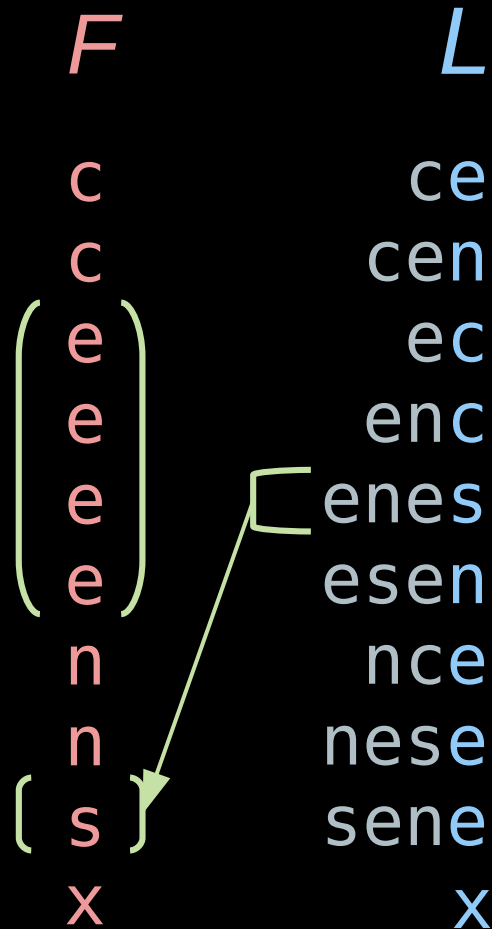
not counted

false occurrence

• backward search $P = se$

• $P_2 = e$

• $P_1 = s$



x | enes | cen | ce

found

x | enes | cen | ce

not counted

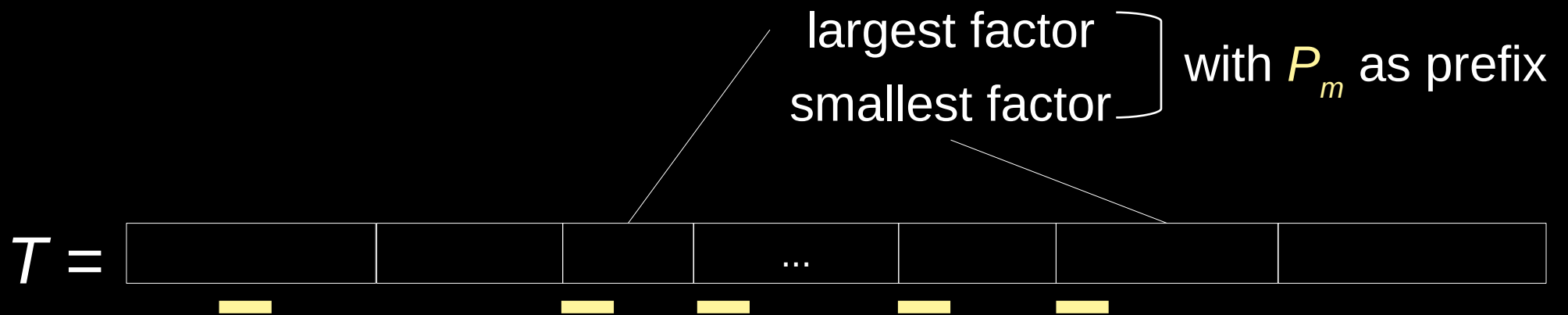
false occurrence

$$T = \left[\begin{array}{ccccccc} & & & \dots & & & \end{array} \right]$$

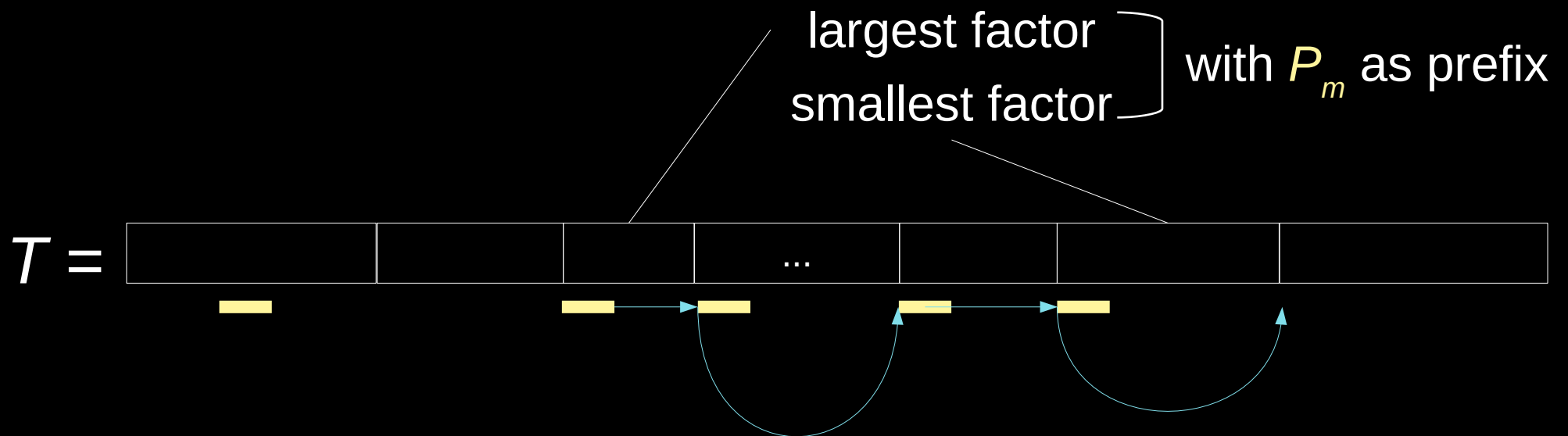
- backward search P_m



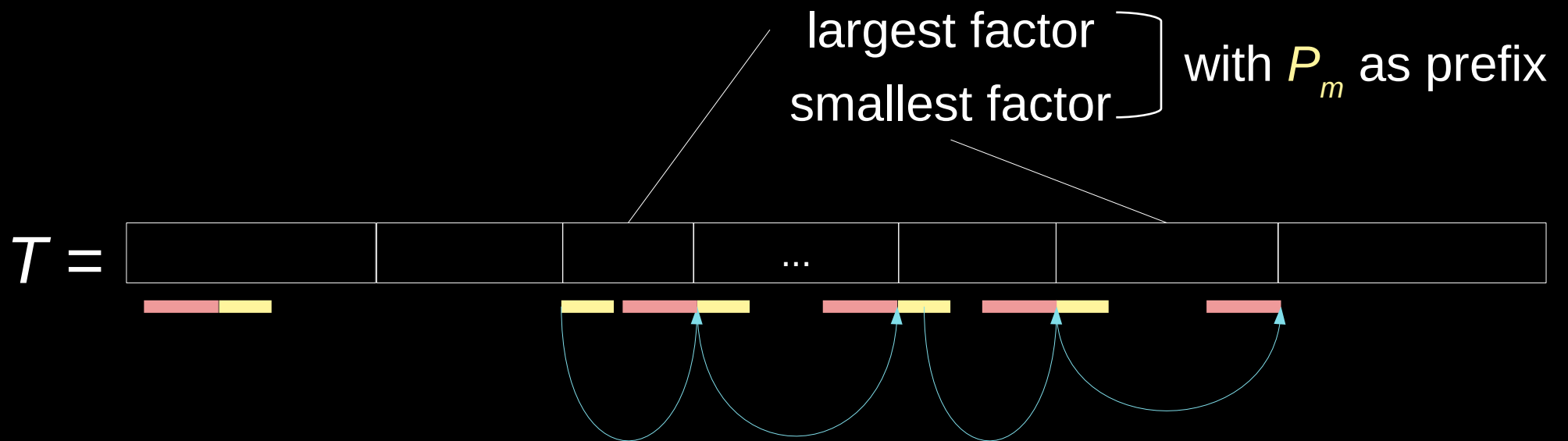
- backward search P_m



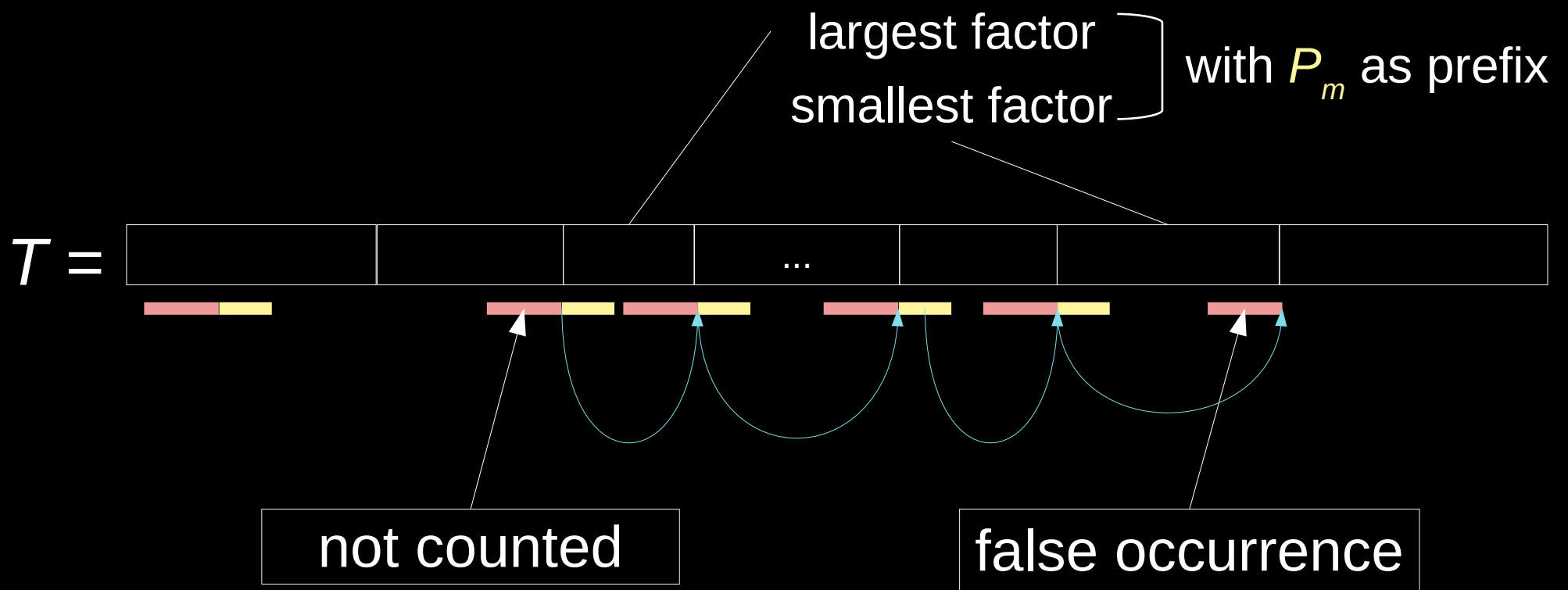
- backward search P_m
- continue search $P_{m-1}P_m$



- backward search P_m
- continue search $P_{m-1}P_m$



- backward search P_m
- continue search $P_{m-1}P_m$



- after finding range of P_m :
 - for border $P_{m-1}P_m$ maintain
 - pointer to not-counted occurrence
 - pointer to false occurrence
- in total backward search on
 - range
 - at most $2m$ individual values
- smallest/largest factor with P_m as prefix = ?

location of factors T_i

s | enes | cen | ce

L

ce

cen

ec

enc

enes

esen

nce

nese

sene

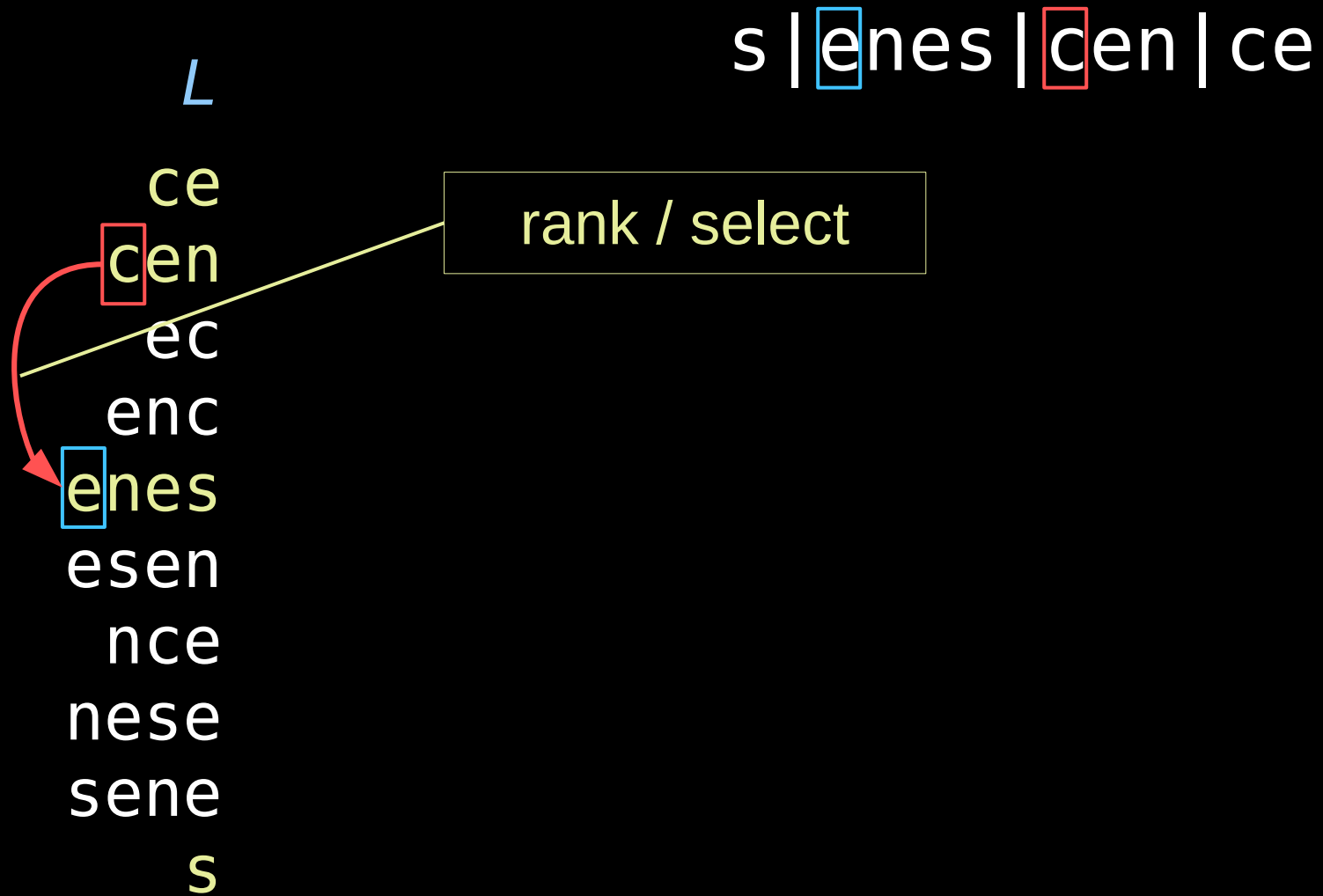
S

Lyndon factors T_t, \dots, T_1

u, v Lyndon words:

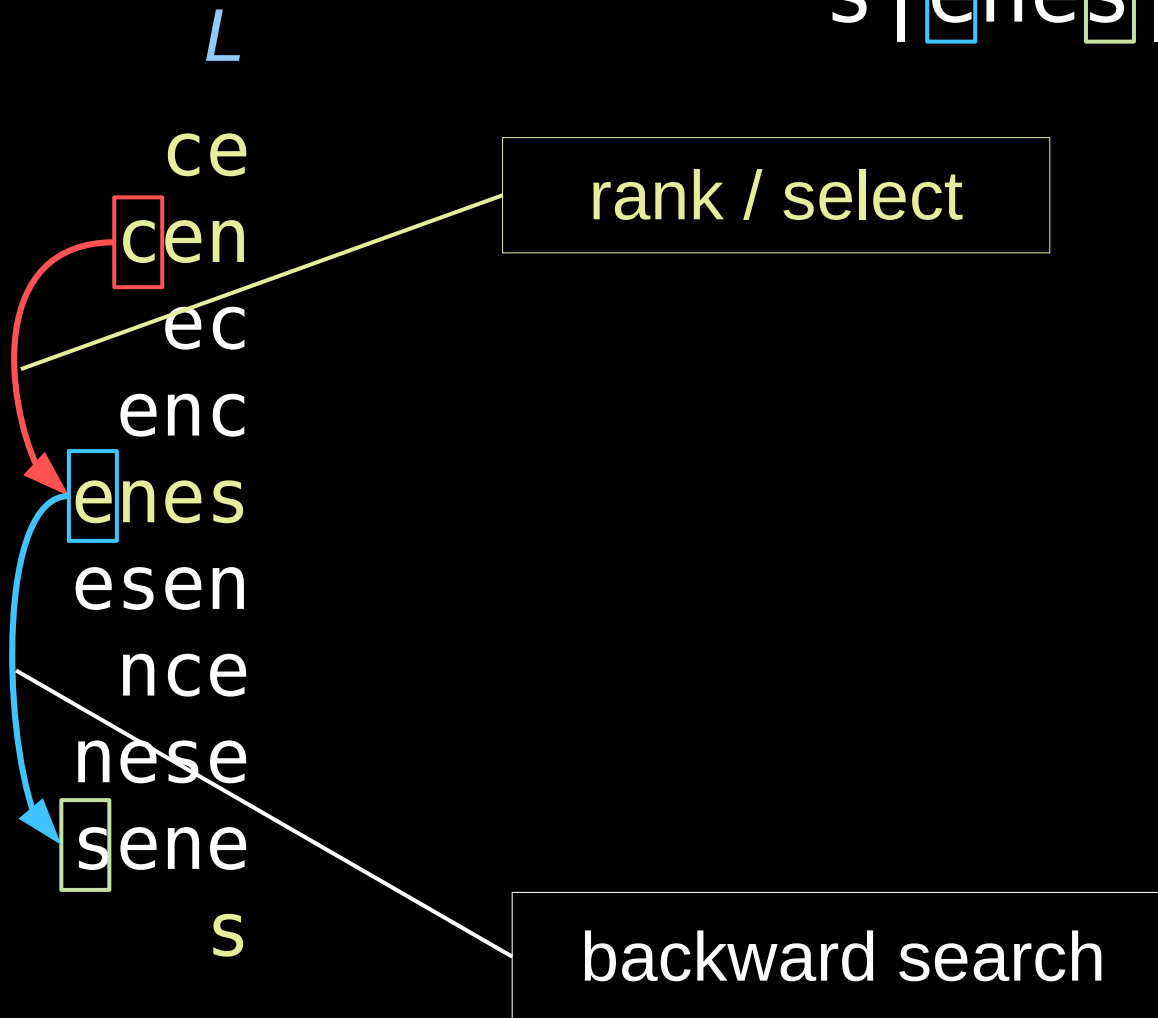
$$u <_{\text{lex}} v \iff u <_{\omega} v$$

location of factors T_i

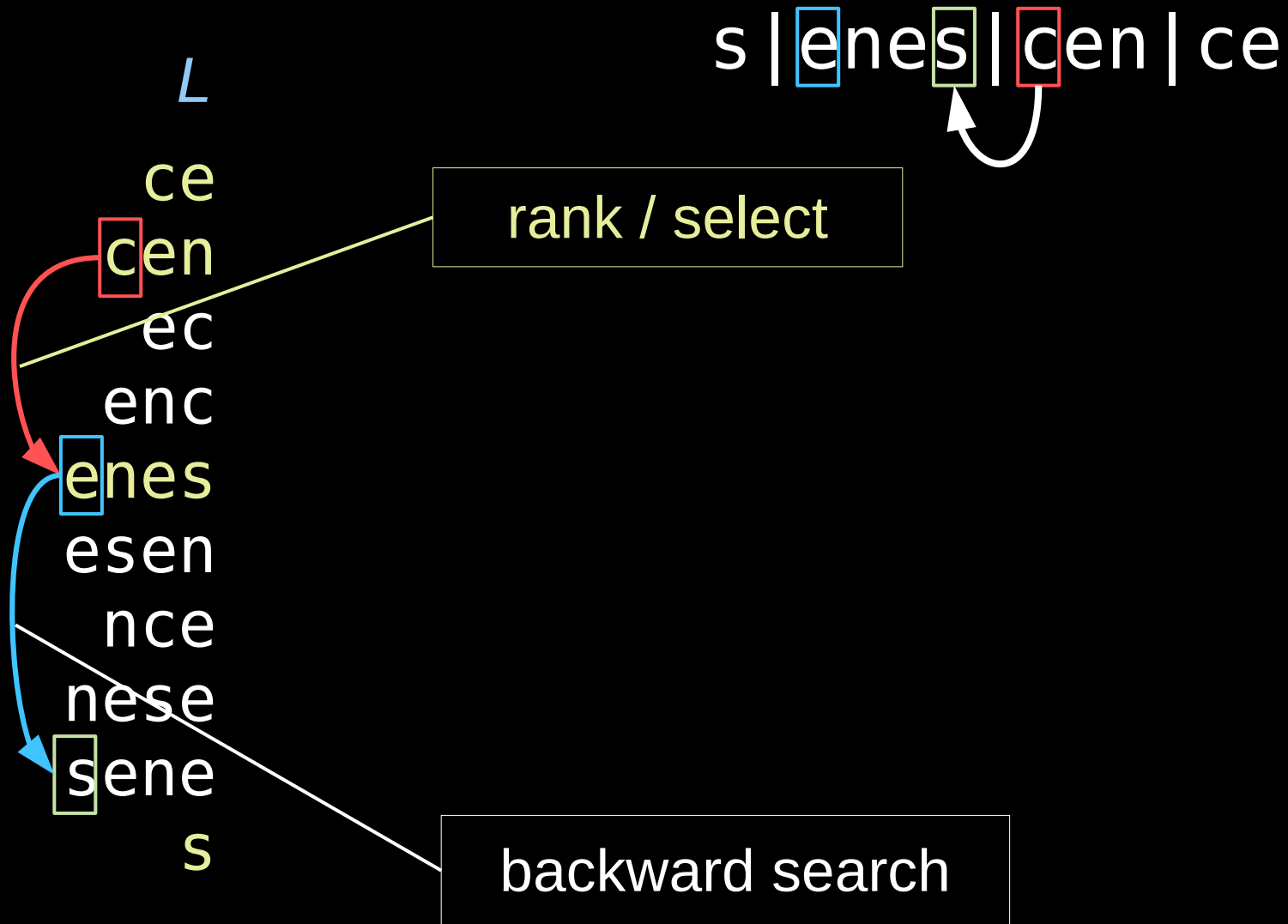


location of factors T_i

s | enes | cen | ce



location of factors T_i



conclusion

- FM index with bijective BWT
 - for each pattern character $O(\lg |P|)$ additional rank/selects
 - $\Rightarrow O(\lg |P|)$ times slower than FM index
- uses properties of Lyndon factorization on
 - text
 - pattern P

conclusion

- FM index with bijective BWT
 - for each pattern character $O(\lg |P|)$ additional rank/selects
 - $\Rightarrow O(\lg |P|)$ times slower than FM index
- uses properties of Lyndon factorization on
 - text
 - pattern P

Thank you for your attention. Any questions are welcome!