

# Indexing the Bijective BWT

**Hideo Bannai (Kyushu University),**

**Juha Kärkkäinen (Helsinki Institute of Information Technology),**

***Dominik Köppl (TU Dortmund),***

**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
- #occ: number of occurrences
- time independent on #occ

# FM Index

on bijective BWT

## ingredients

- bijective BWT
- wavelet tree

## operation: backward search

- locate pattern
- #occ: number of occurrences
- time independent on #occ

# FM Index on bijective BWT

## ingredients

- bijective BWT
- wavelet tree

## operation: backward search

- locate pattern
- #occ: number of occurrences
- time independent on #occ

reason:  
bijective BWT is cool

# 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

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

(Chen-Fox-Lyndon theorem)



# properties [Duval' 88]

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

$\prec_{\omega}$

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

abababab...  
abaabaaba...

# bijjective BWT of senescence

s | enes | cen | ce

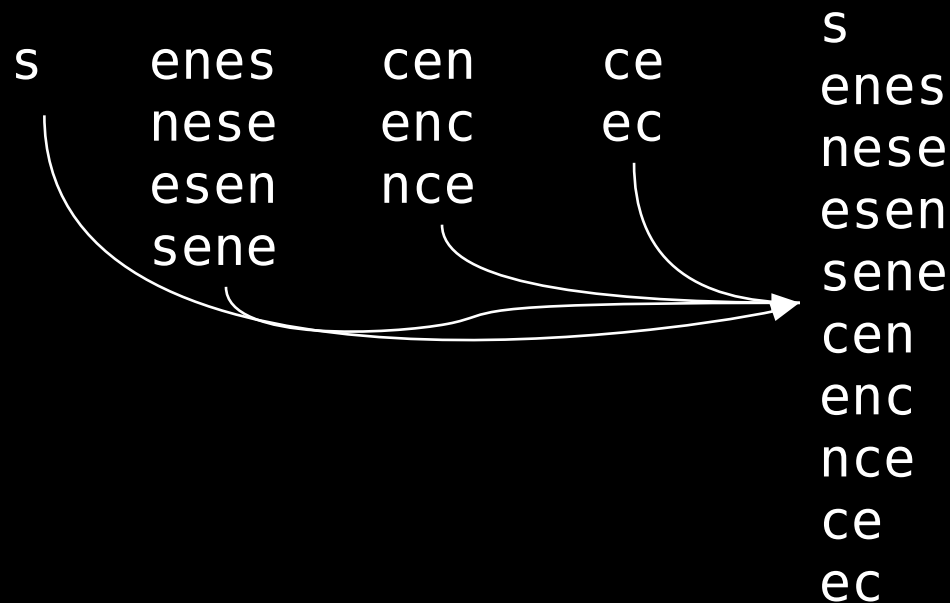
# bijjective BWT of senescence

s | enes | cen | ce

s	enes	cen	ce
	nese	enc	ec
	esen	nce	
	sene		

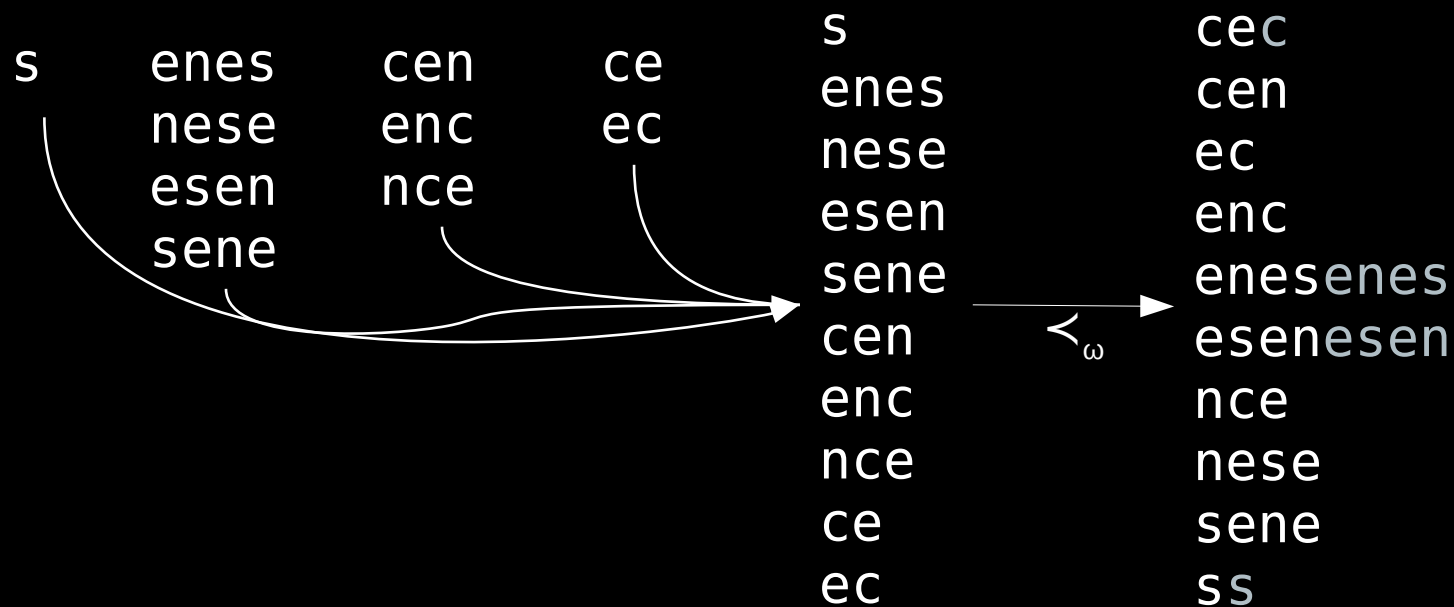
# bijection BWT of senescence

s | enes | cen | ce



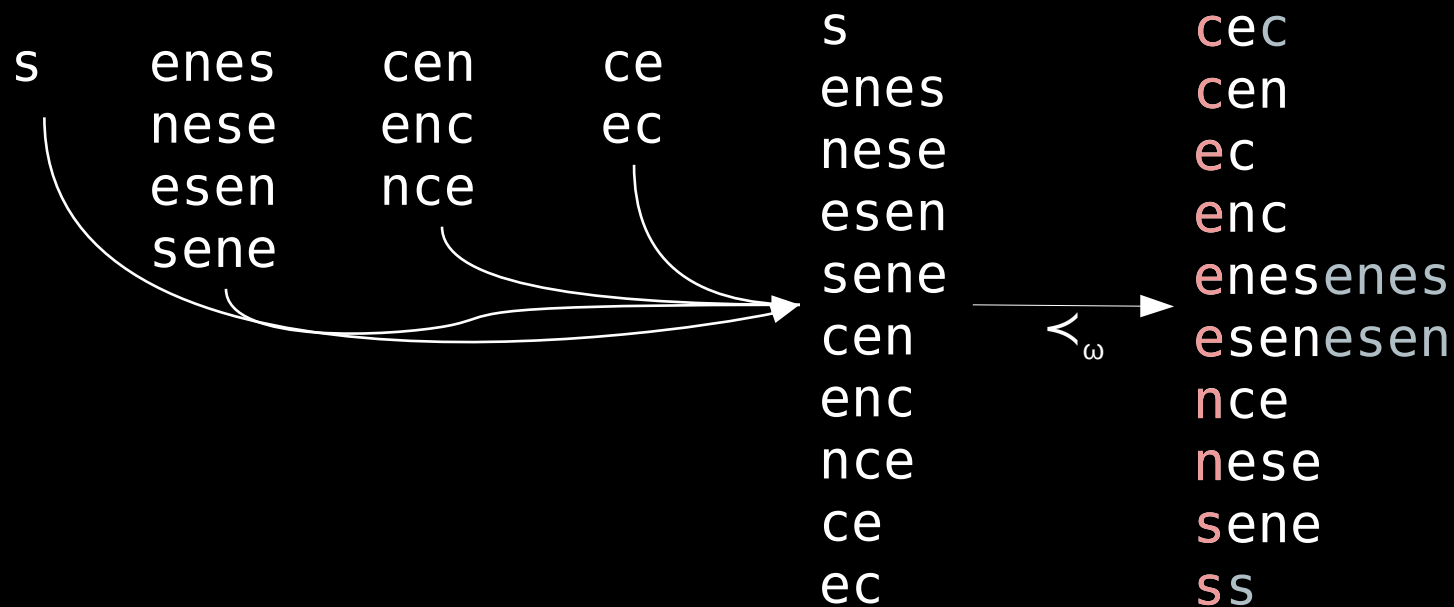
# bijjective BWT of senescence

s | enes | cen | ce



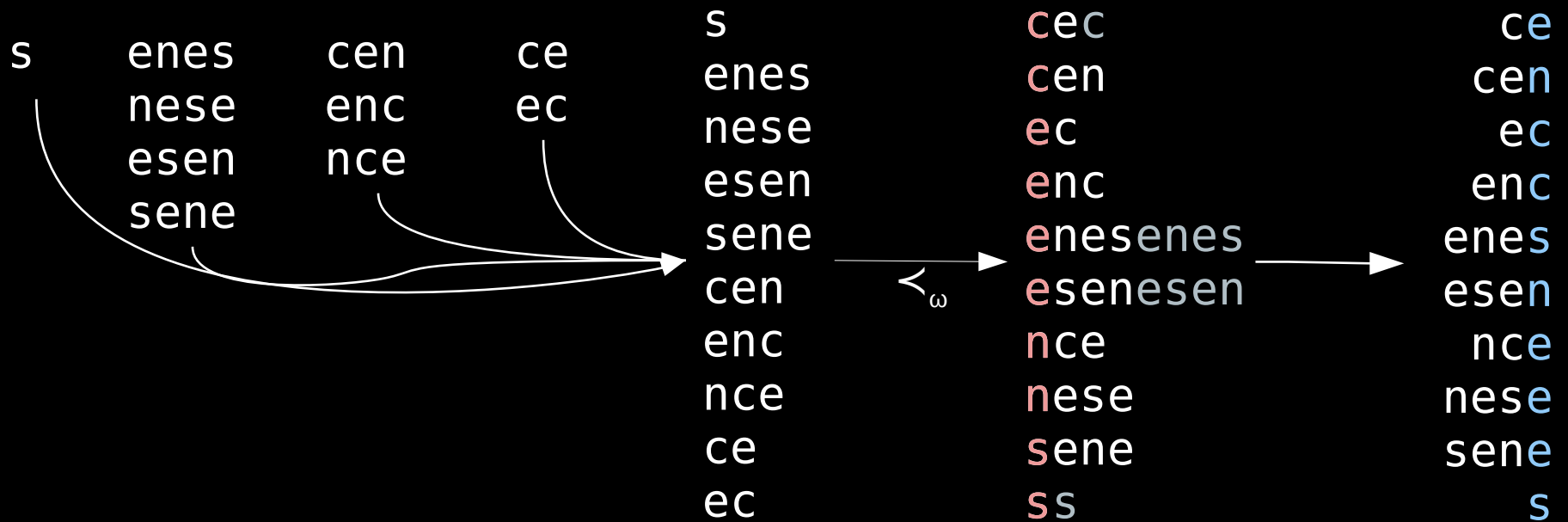
# bijjective BWT of senescence

s | enes | cen | ce



# bijjective BWT of senescence

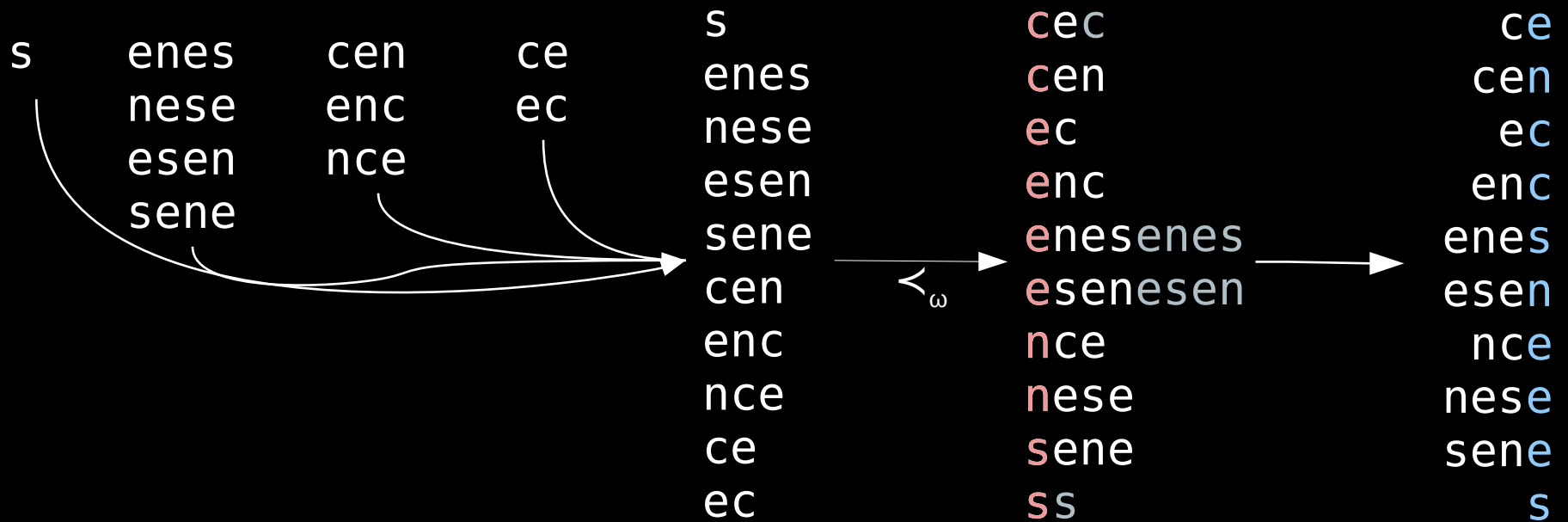
s | enes | cen | ce





# bijjective BWT of senescence

s | enes | cen | ce



result: encsneees

# *cycles*

*L*

*e*

*n*

*c*

*c*

*s*

*n*

*e*

*e*

*e*

*s*

*F*

*c*

*c*

*e*

*e*

*e*

*e*

*n*

*n*

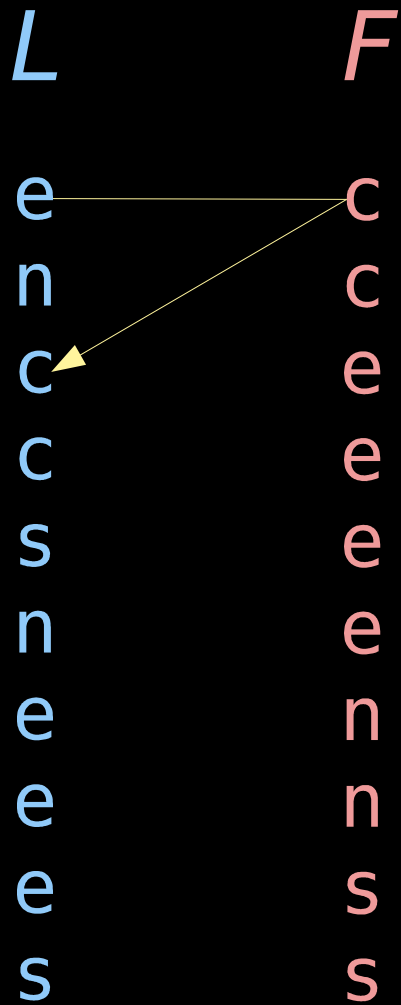
*s*

*s*

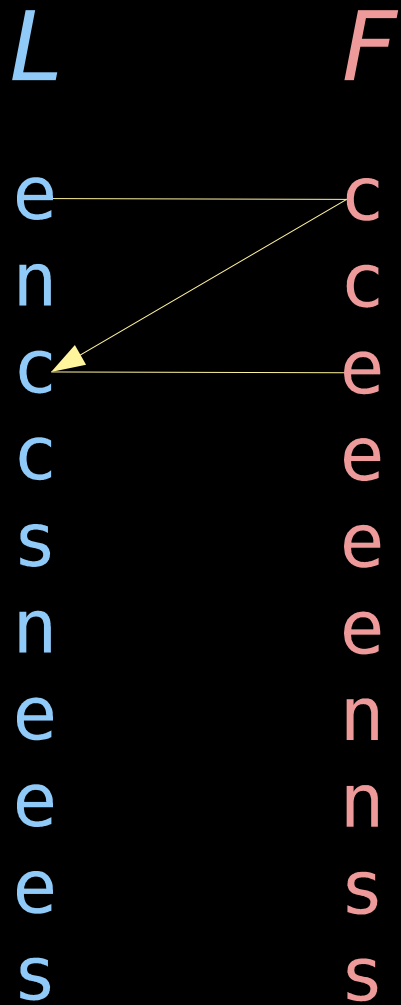
# cycles

<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

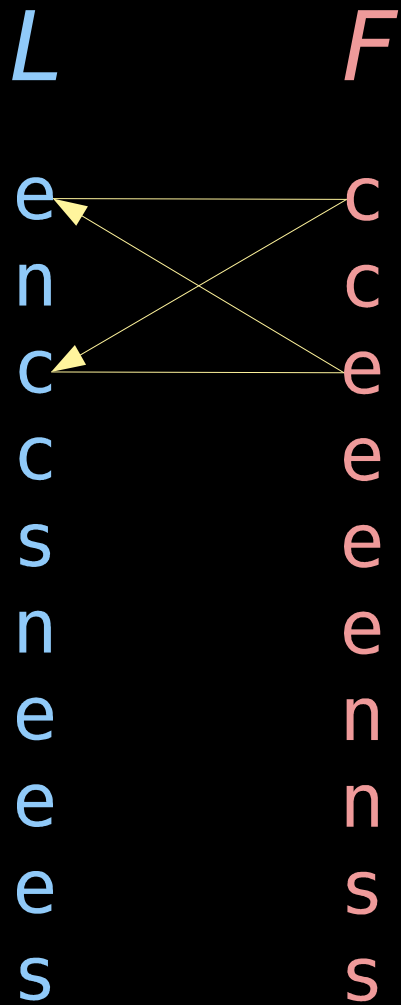
# *cycles*



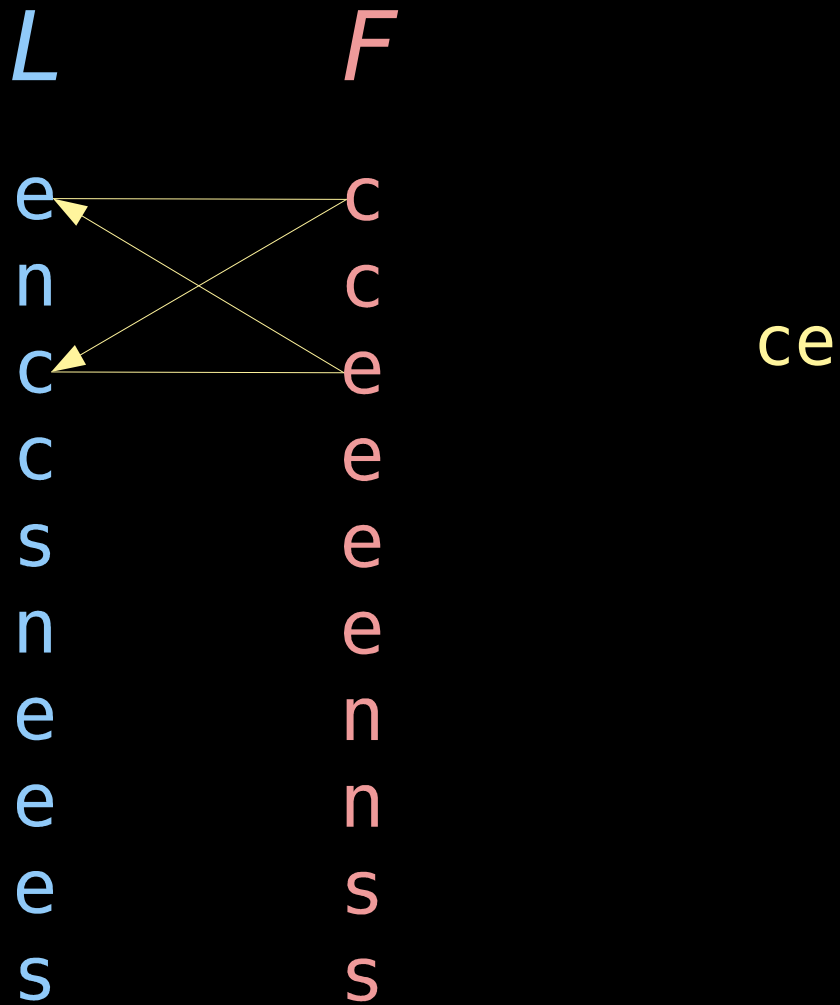
# cycles



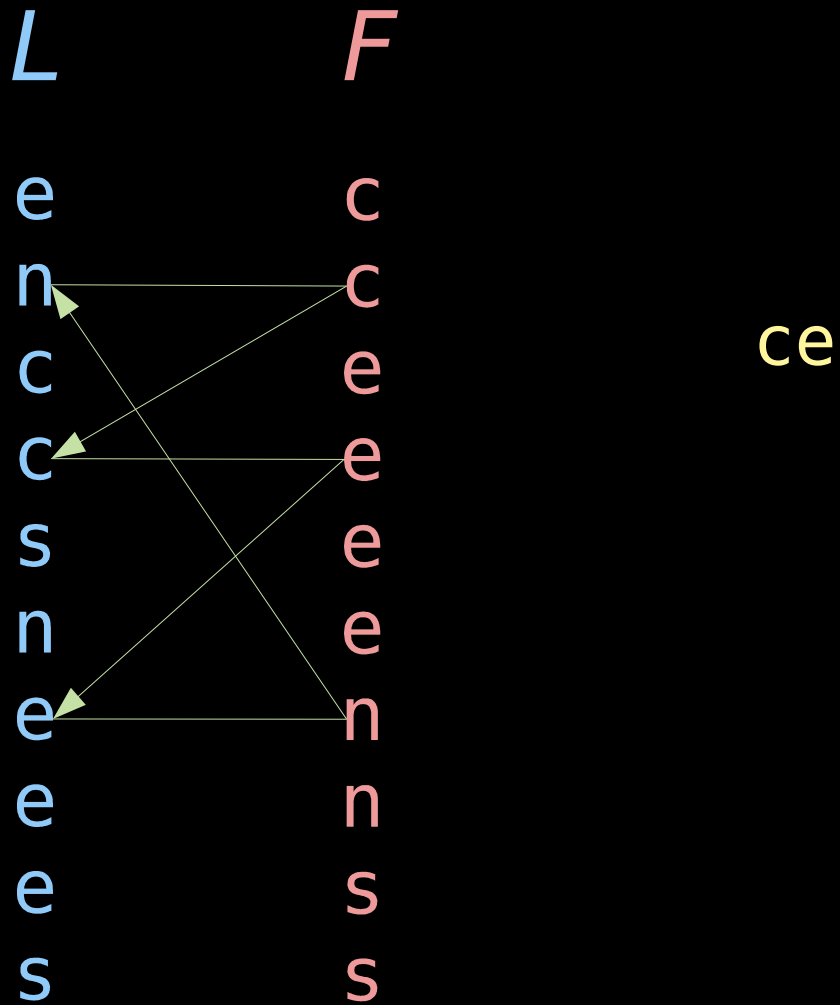
# cycles



# cycles

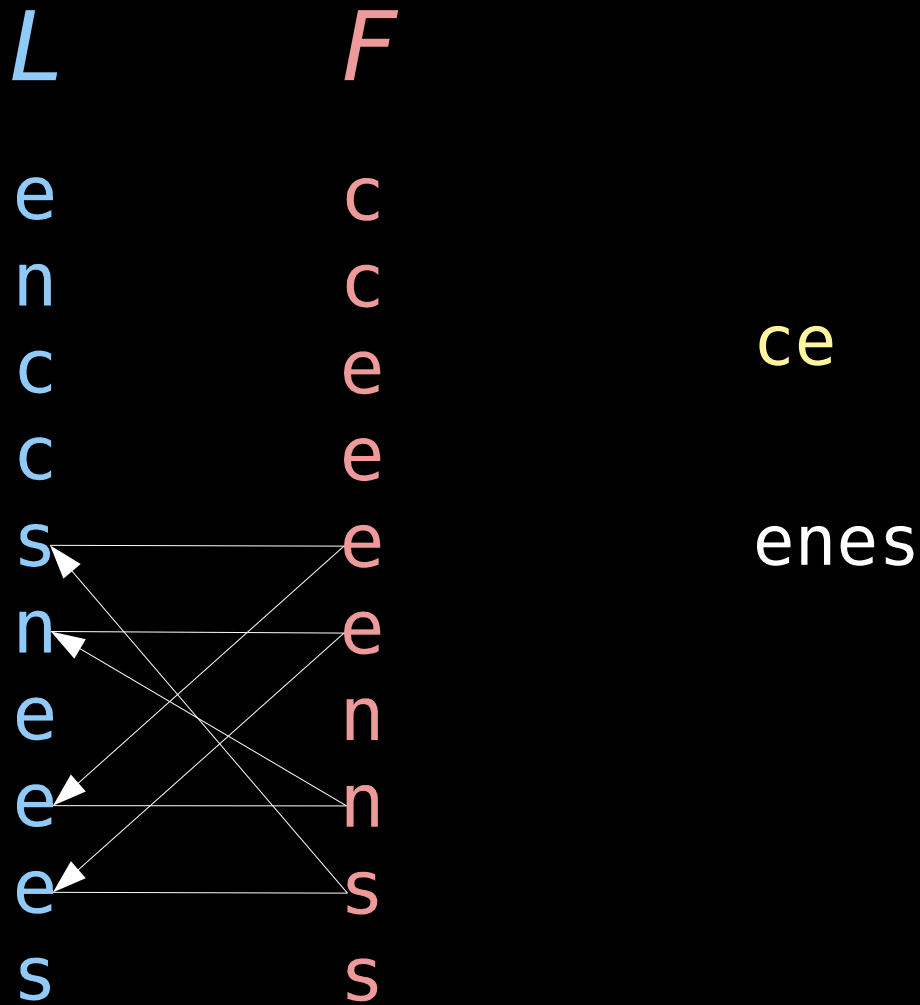


# cycles

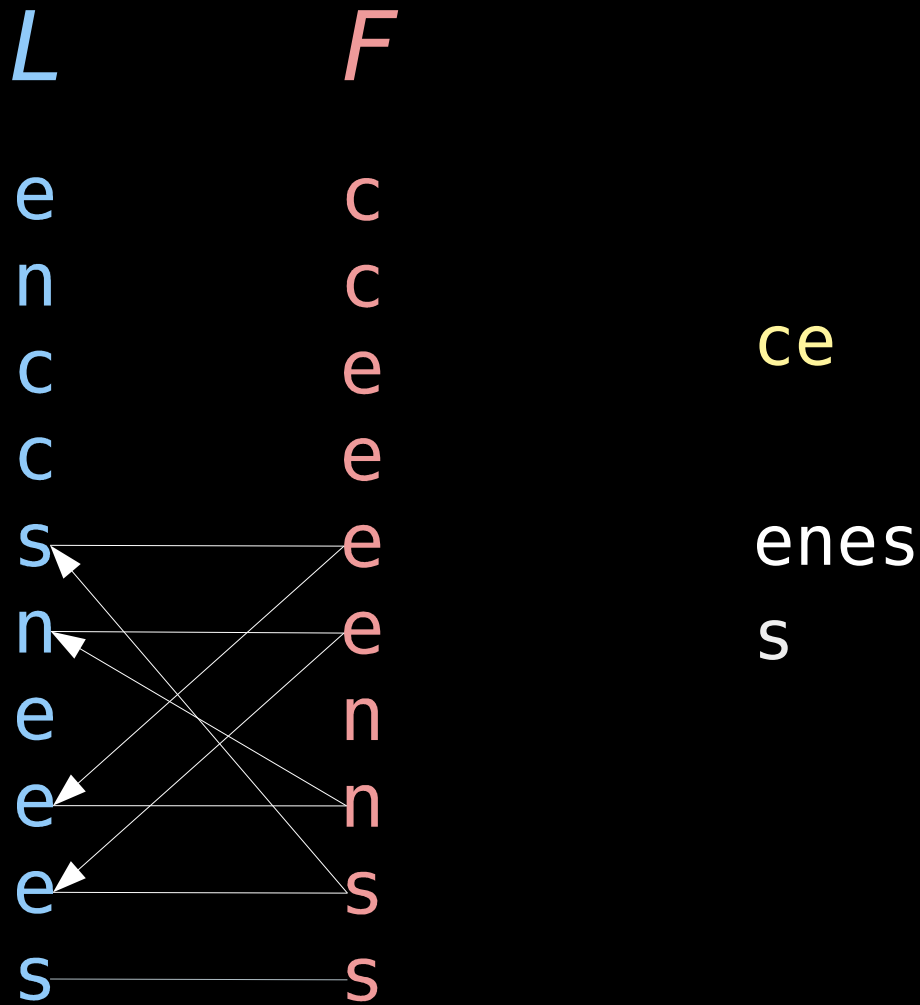




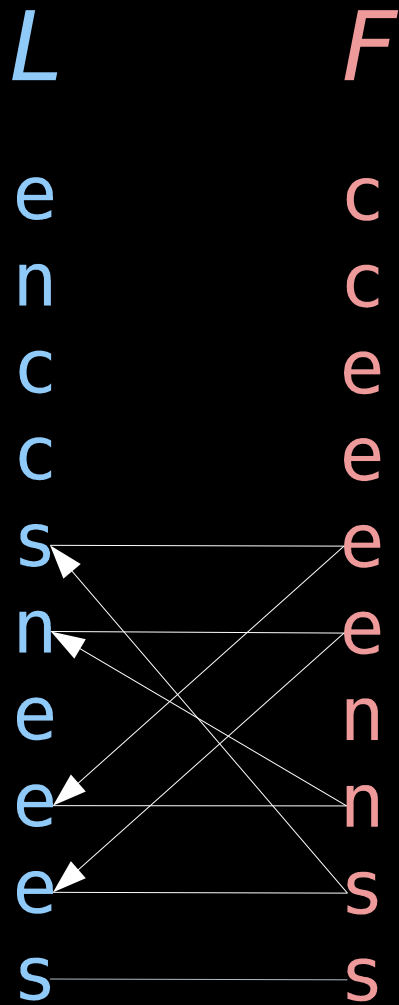
# cycles



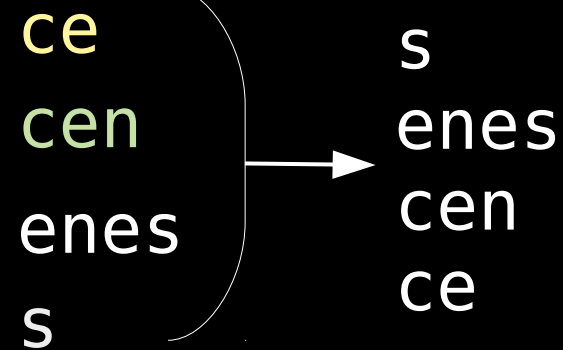
# cycles



# 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 )

enc

( e )

enes

e

esen

( n )

nce

( n )

nese

s

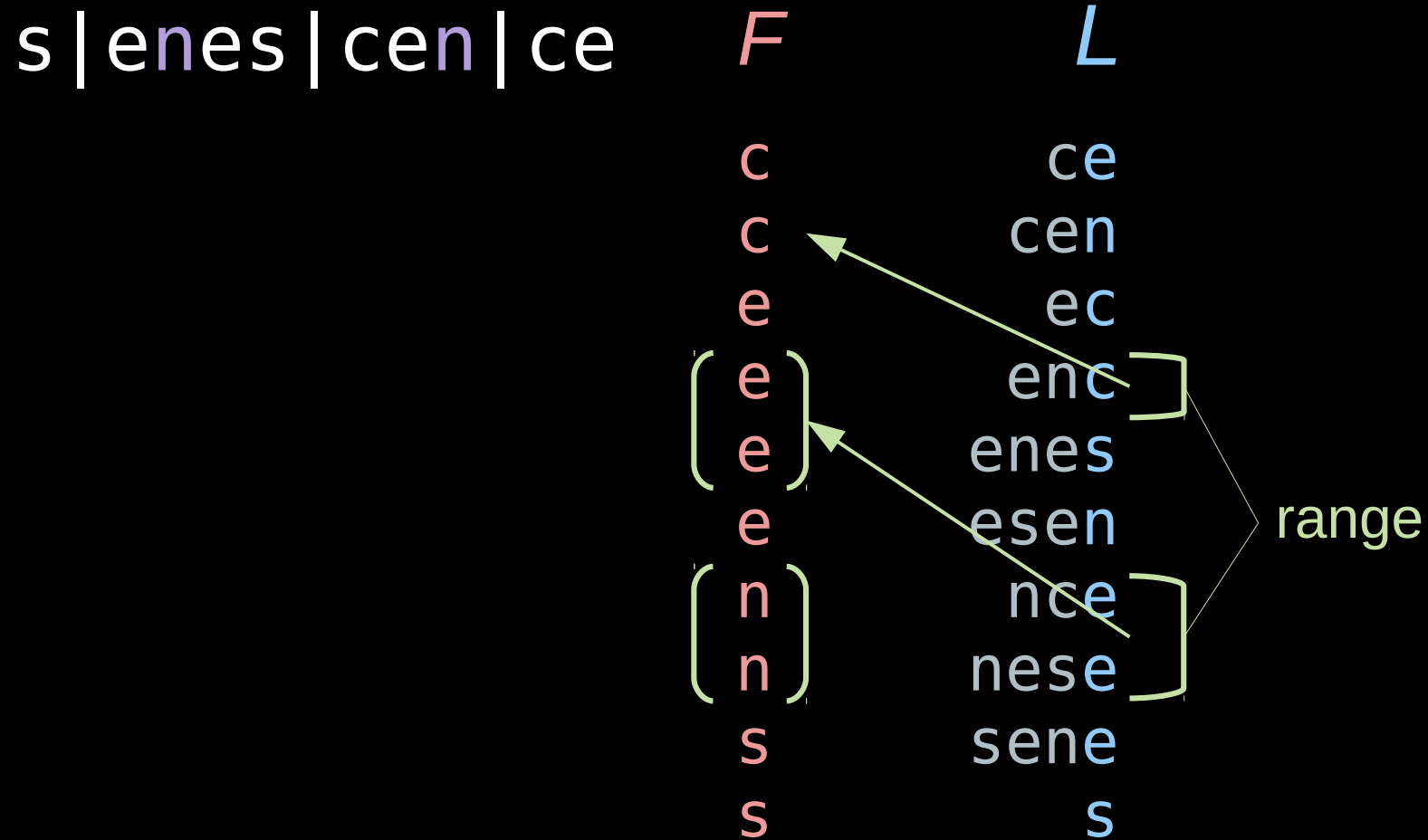
sene

s

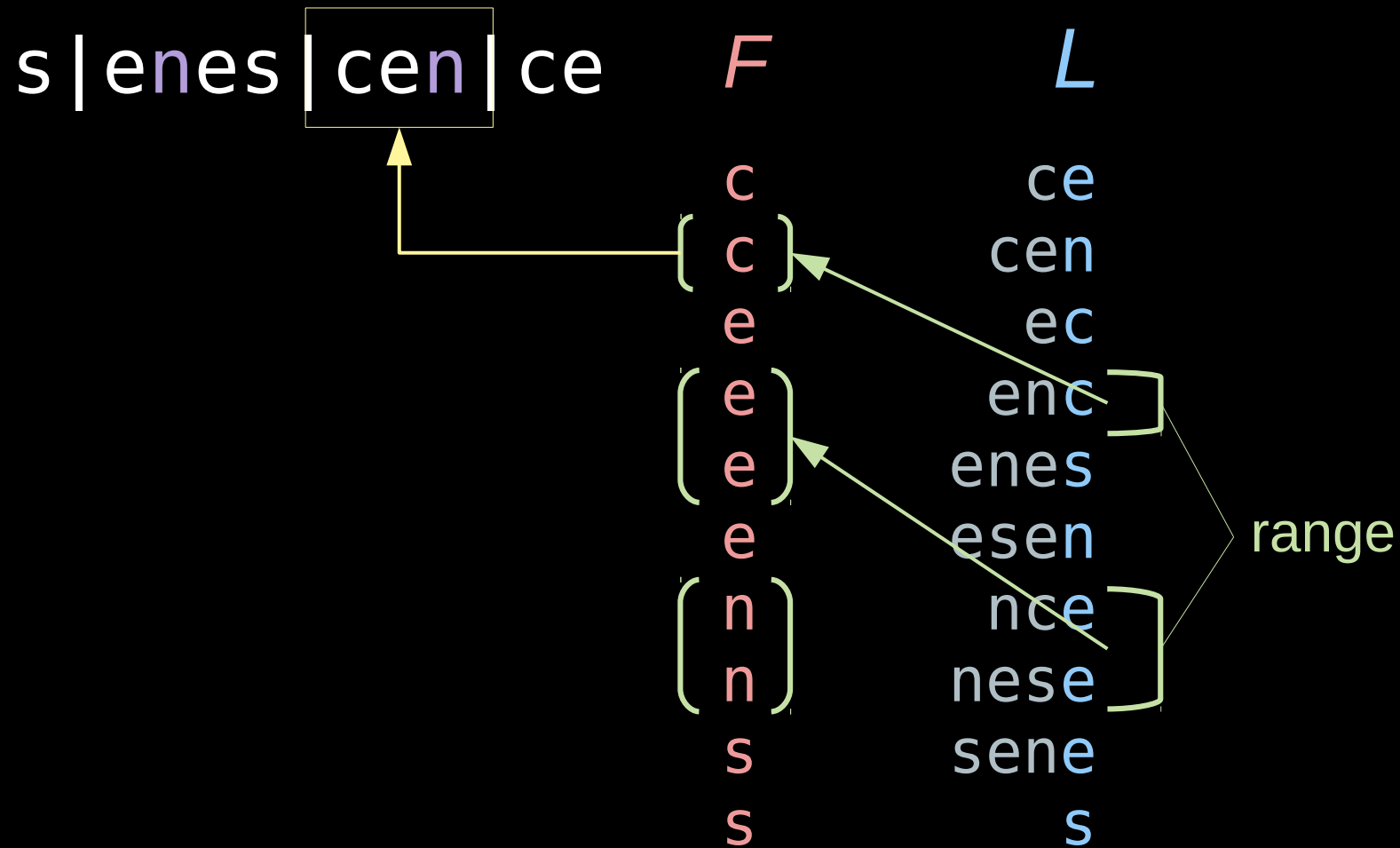
s

range

# backward search 'cen'



# backward search 'cen'





# 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

s

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	<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

⇒ occurrences inside factors

⇒ found within cycles

backward search  $\cong$  FM-index

# pattern $P$ is not a Lyndon word

- Lyndon factorization:  $P = P_1 \cdots P_m$
- $P_j$  substring of  $T_i$  or equal to  $T_i$
- search  $P_m$
- take care when starting with  $P_{m-1}$ !

- backward search  $P = se$

<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  $P = se$

- $P_2 = e$

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

• backward search  $P = se$

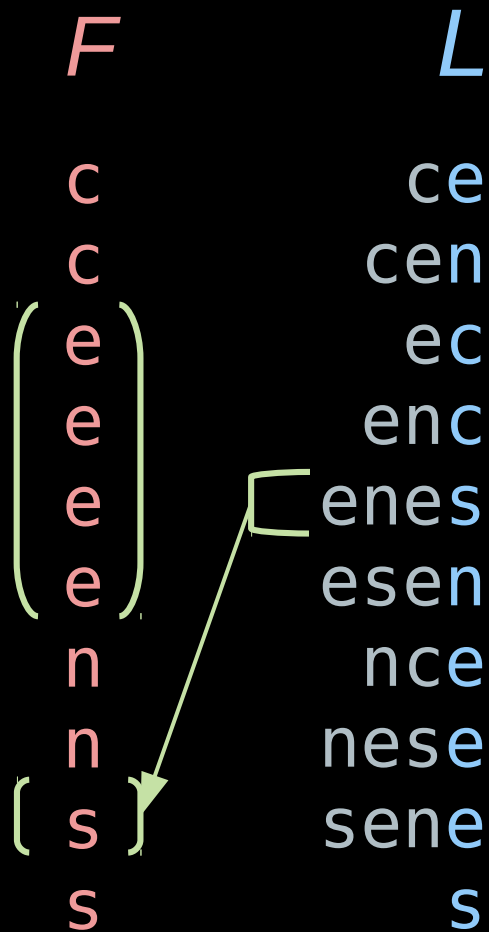
•  $P_2 = e$

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

- backward search  $P = se$

- $P_2 = e$

- $P_1 = s$



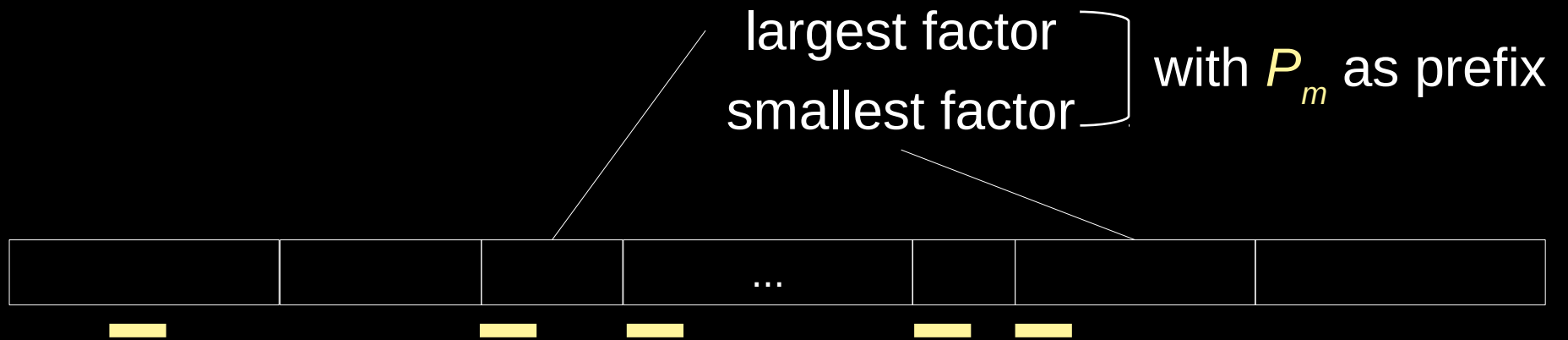




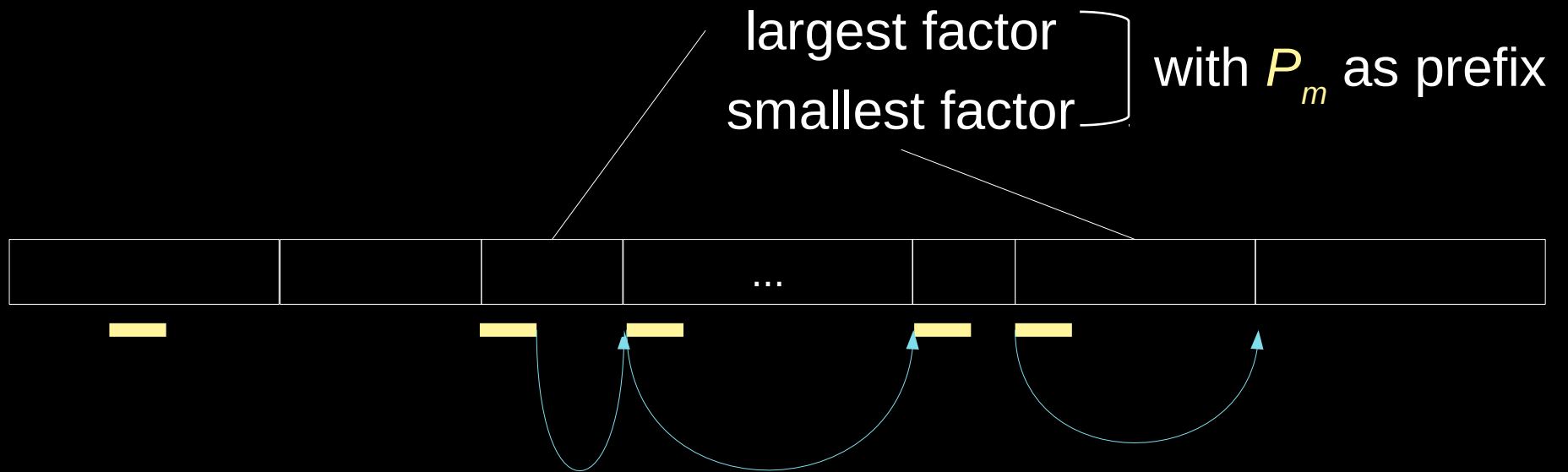
- backward search  $P_m$



- backward search  $P_m$



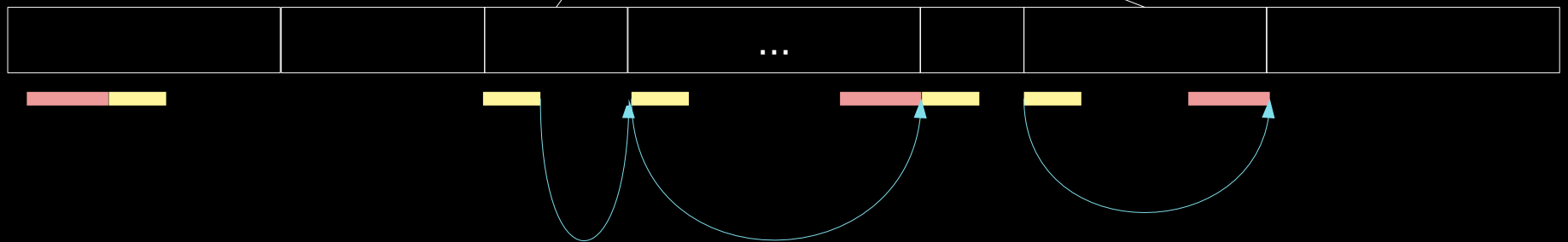
- backward search  $P_m$
- continue search  $P_{m-1}P_m$



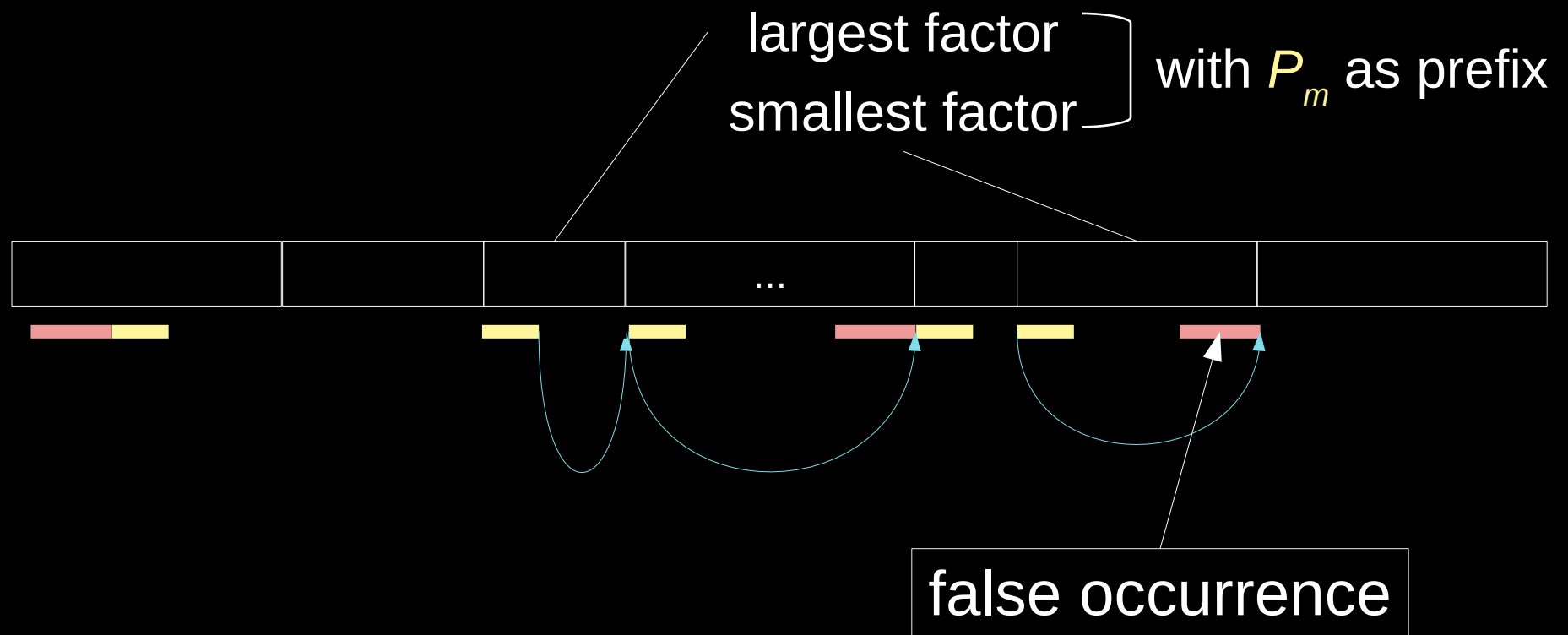


- backward search  $P_m$
- continue search  $P_{m-1}P_m$

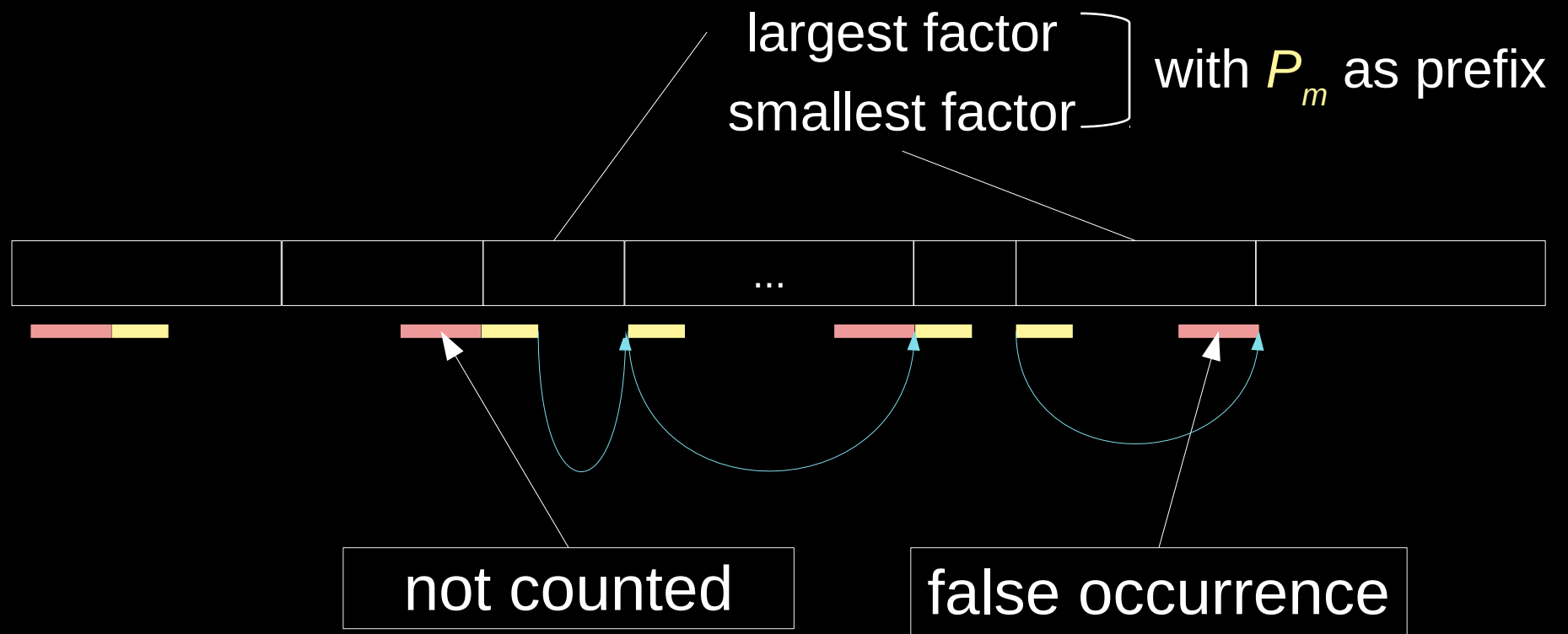
largest factor  
smallest factor } with  $P_m$  as prefix



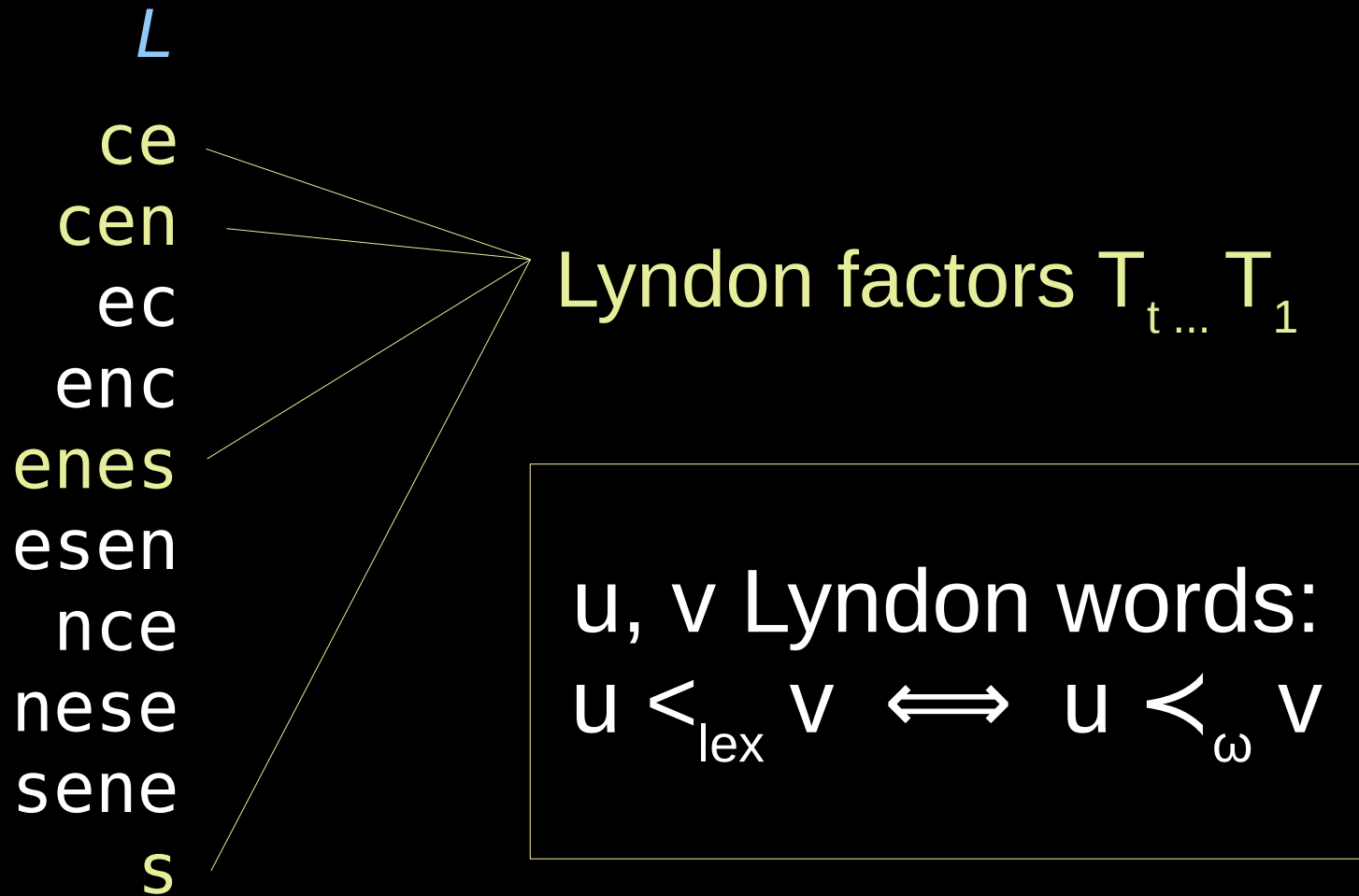
- backward search  $P_m$
- continue search  $P_{m-1}P_m$



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



# location of factors $T_i$



- 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

# conclusion

- FM index with bijective BWT
- uses properties of Lyndon factorization on
  - text
  - pattern  $P = P_1 \cdots P_m$
- currently  $O(m)$  times slower than FM index
- extended BWT does not seem to be a good candidate (no Lyndon word properties)

# conclusion

- FM index with bijective BWT
- uses properties of Lyndon factorization on
  - text
  - pattern  $P = P_1 \cdots P_m$
- currently  $O(m)$  times slower than FM index
- extended BWT does not seem to be a good candidate (no Lyndon word properties)

Thank you for your attention. Any questions are welcome!