

In-Place (Bijective) BWT Transforms

Kyushu University

Tohoku University

Dominik Köppl

Daiki Hashimoto

Diptarama

Ayumi Shinohara

data structures

Burrows-Wheeler Transform (BWT)

[Burrows,Wheeler '94]

Bijective BWT (BBWT)

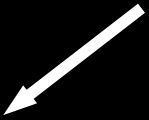
[Gil,Scott '12]

BWT of bacabbabb

$T = \text{bacabbabb\$}$

BWT of bacabbabb

$T = \text{bacabbabb\$}$

 all suffixes

- bacabbabb\$
- acabbabb\$
- cabbabb\$
- abbabb\$
- bbabb\$
- babb\$
- abb\$
- bb\$
- b\$
- \$

BWT of bacabbabb

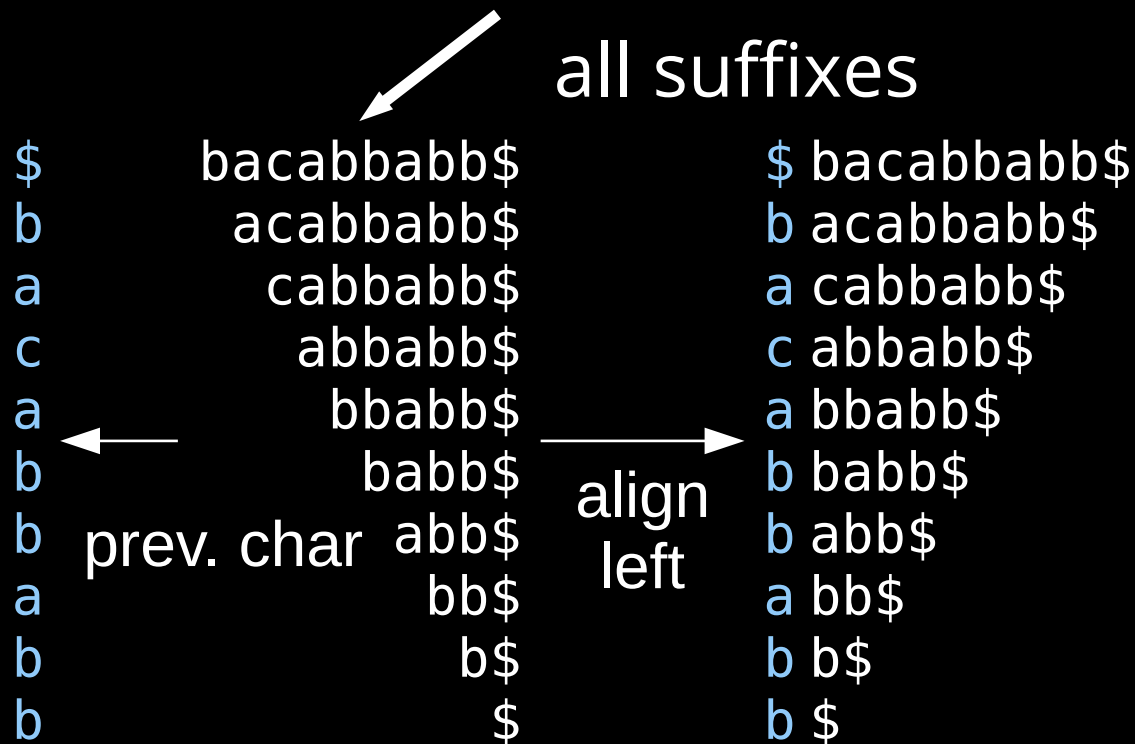
$T = \text{bacabbabb}\$$

all suffixes

\$		bacabbabb\$
b		acabbabb\$
a		cabbabb\$
c		abbabb\$
a		bbabb\$
b	←	babb\$
b	prev. char	abb\$
a		bb\$
b		b\$
b		\$

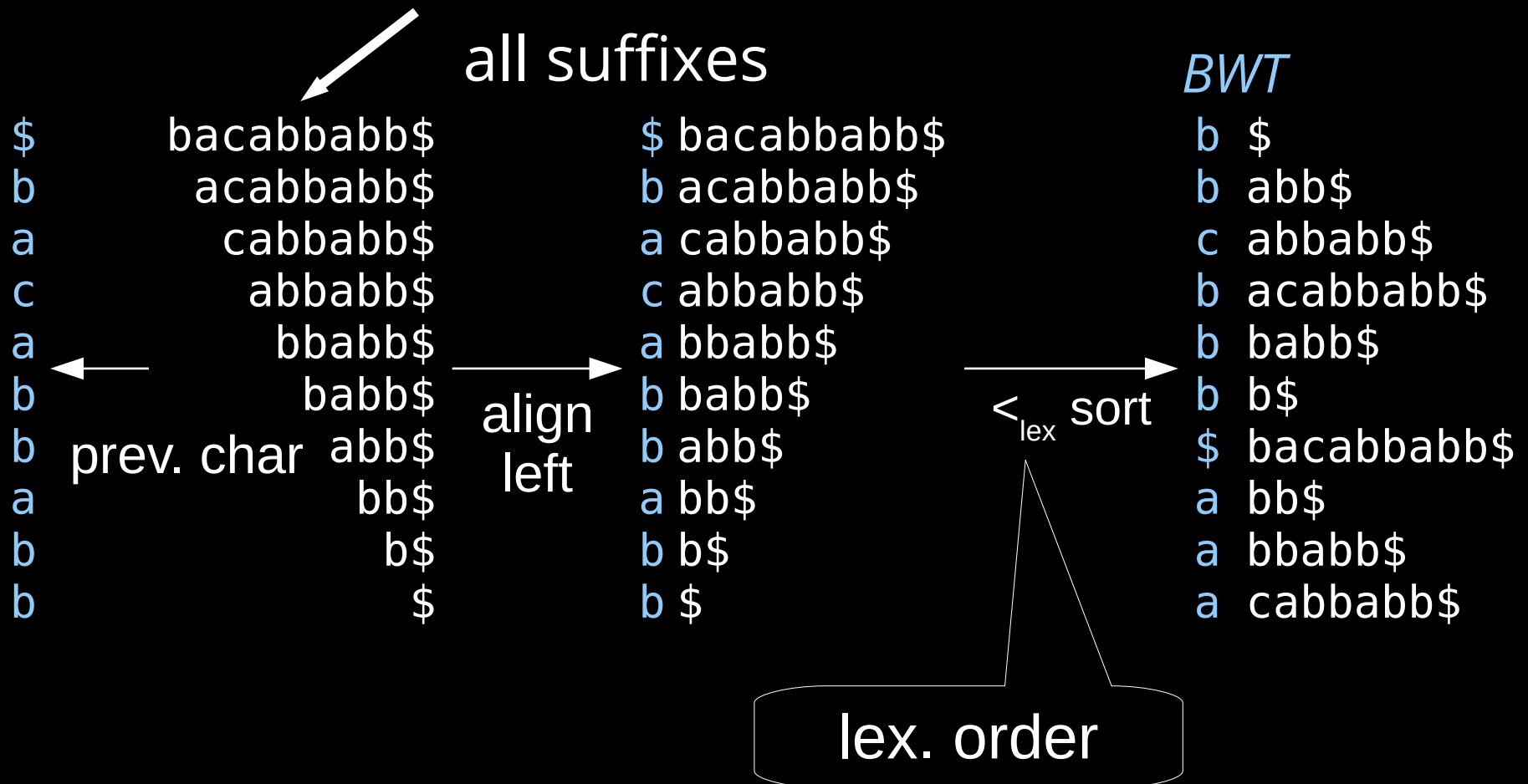
BWT of bacabbabb

$T = \text{bacabbabb}\$$



BWT of bacabbabb

$T = \text{bacabbabb}\$$



the BBWT is
the BWT of
the Lyndon factorization
of an input text
with respect to $<_{\omega}$

the BBWT is
the BWT of
the Lyndon factorization 1.
of an input text
with respect to $<_{\omega}$ 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 $T =$

T_1	T_2	\dots	T_t
-------	-------	---------	-------
- output: factorization $T_1 \dots T_t$ with
 - T_x is Lyndon word
 - $T_x \geq_{\text{lex}} T_{x+1}$
 - factorization uniquely defined
 - linear time [Duval'88]

(Chen-Fox-Lyndon Theorem)

example

$T = \text{bacabbabb}$

Lyndon factorization: $b|ac|abb|abb$

- b, ac, abb , and abb are Lyndon
- $b >_{\text{lex}} ac >_{\text{lex}} abb \geq_{\text{lex}} abb$

\prec_ω order

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

\prec_ω order

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

- $ab \prec_{\text{lex}} aba$

abababab...

- $aba \prec_\omega ab$

abaabaaba...

BBWT of bacabbabb

b | ac | abb | abb

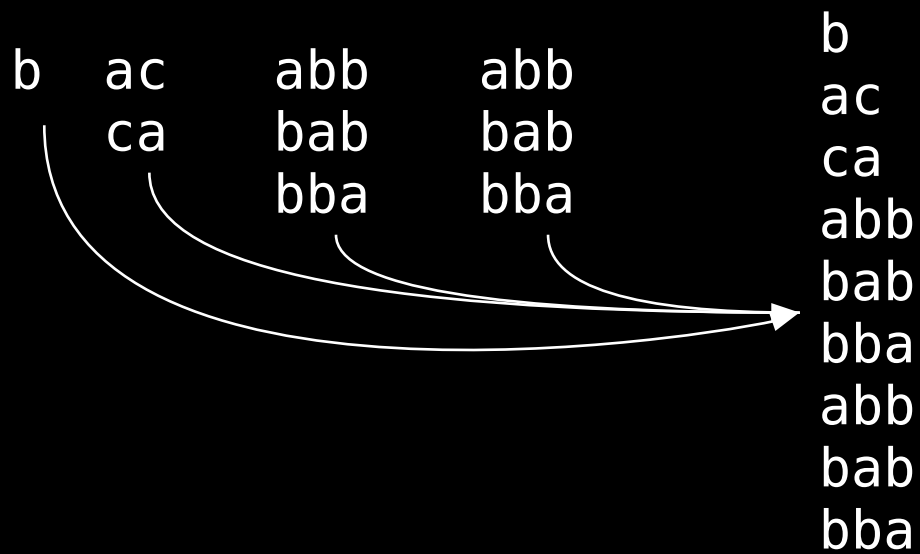
BBWT of bacabbabb

b | ac | abb | abb

b	ac	abb	abb
	ca	bab	bab
		bba	bba

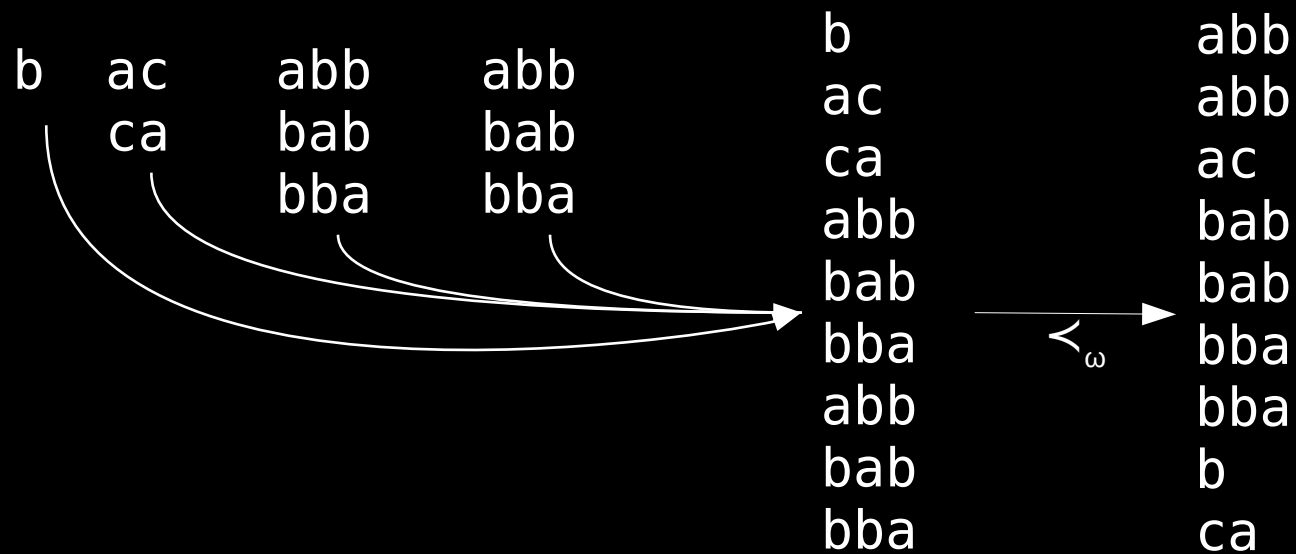
BBWT of bacabbabb

b | ac | abb | abb



BBWT of bacabbabb

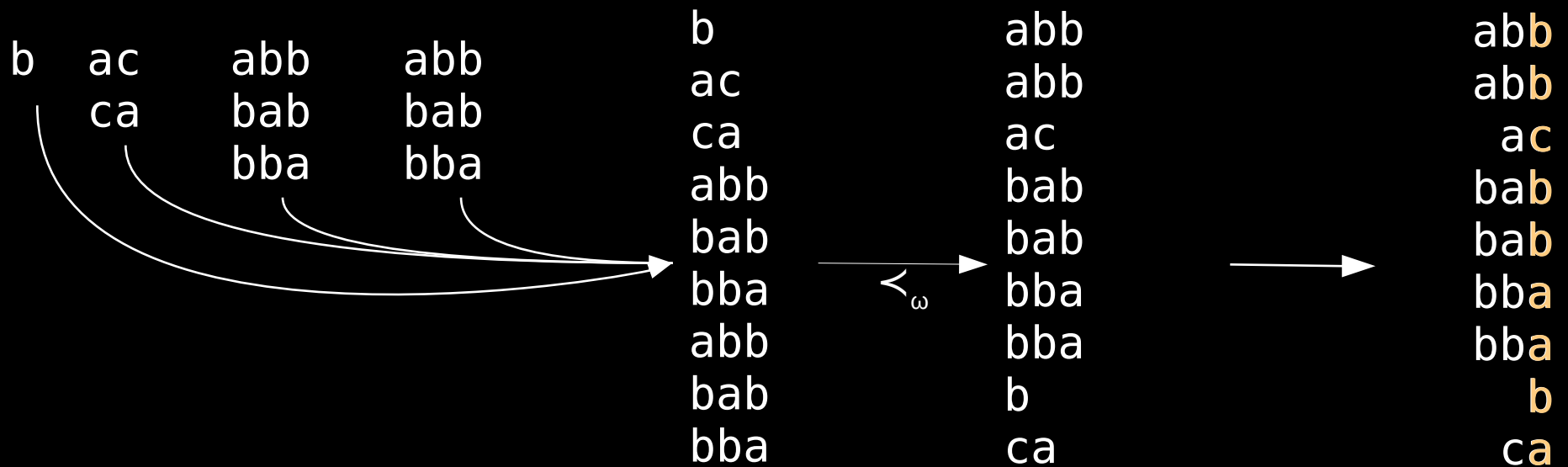
b | ac | abb | abb



BBWT of bacabbabb

b | ac | abb | abb

BBWT

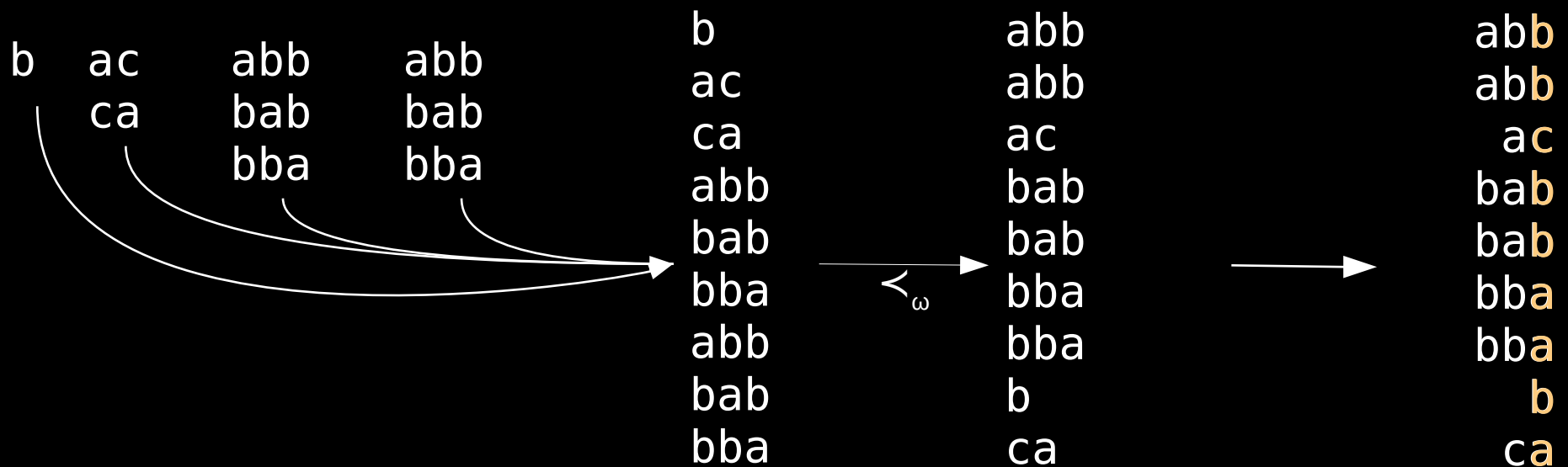


BBWT(*T*) = `bbcbbaaba`

BBWT of bacabbabb

b | ac | abb | abb

BBWT



$BBWT(T) = \text{bbcbbaaba}$

$BWT(T\$) = \text{bbcbbb\$aaa}$

motivation

properties of BBWT :

- no \$ necessary
- BBWT is more compressible than BWT for various inputs

[Scott and Gill '12]

- BBWT is indexable (full text index)
- is computable in $O(n)$ time with $O(n)$ words

[Bannai+ '19]

however, $O(n)$ words can be too much for large n

in-place computation

- Σ : alphabet, $\sigma := |\Sigma|$ alphabet size
- T : text, $n := |T|$
- $L := n \lg \sigma$ bits workspace
- aim : in-place computation

transform $T \leftrightarrow \text{BWT} \leftrightarrow \text{BBWT}$ with

$|L| + O(\lg n)$ bits of workspace L

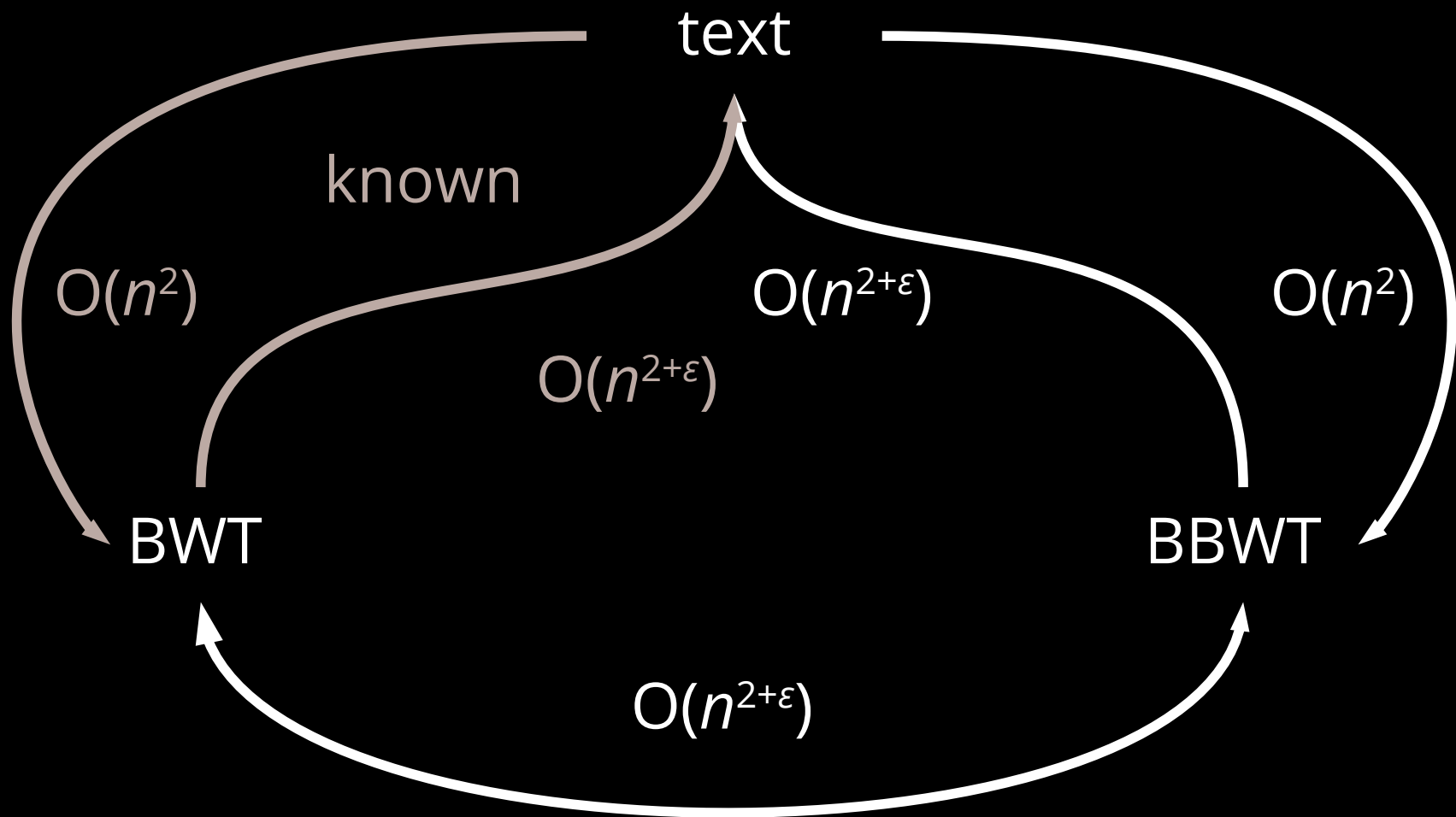


known solutions

input	output	work-space	time	reference
text	BWT	in-place	$O(n^2)$	Crochemore+ '15
BWT	text	in-place	$O(n^{2+\varepsilon})$	
text	BBWT	$O(n \lg \sigma)$ bits	$O(n \lg n / \lg \lg n)$	Bonomo+ '14

σ : alphabet size, n : text length,
 ε is a constant with $0 < \varepsilon < 1$

in-place conversions



working space: $n \lg \sigma + O(\lg n)$ bits (including text)

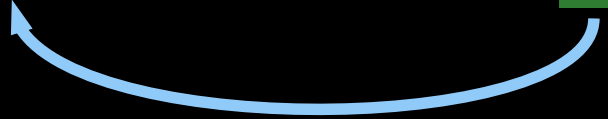
forward search

$T = \text{bacabbabb}\$$

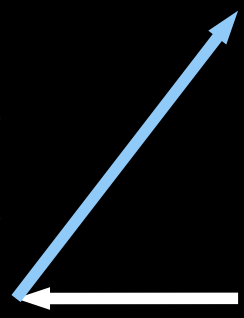
F	L
\$	b
a	b
a	c
a	b
b	b
b	b
b	$\$$
b	a
b	a
c	a

forward search

$T = \text{bacabbabb}\$$

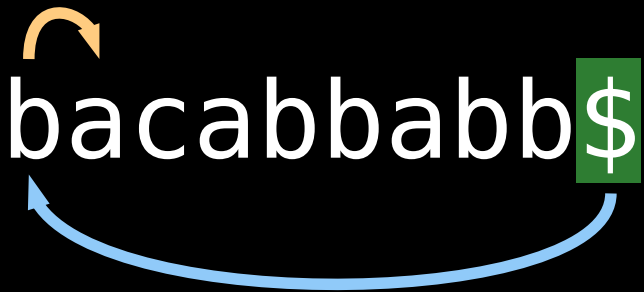
A blue curved arrow starts at the end of the string 'bacabbabb\$' and points back to the beginning, indicating a reverse traversal or search.

<i>F</i>	<i>L</i>
\$	b
a	b
a	c
a	b
b	b
b	b
b	\$
b	a
b	a
c	a

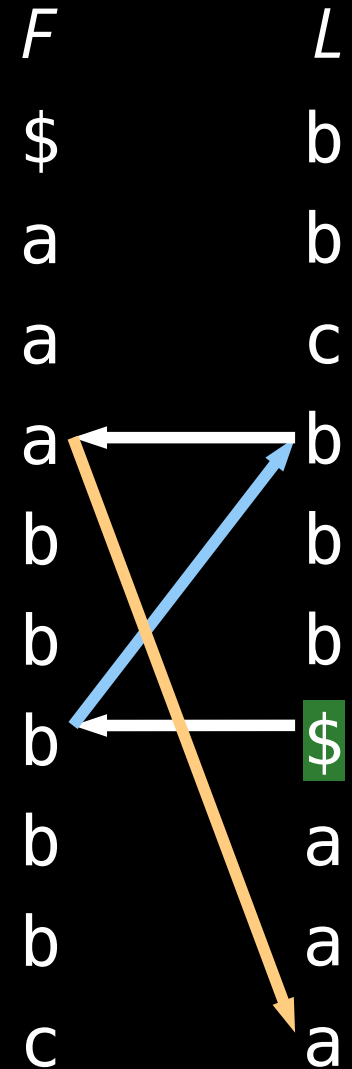
A blue arrow points from the '\$' in the 'L' column to the 'b' in the 'F' column. A white arrow points from the '\$' in the 'L' column to the 'b' in the 'F' column.

forward search

$T = \text{bacabbabb}\$$

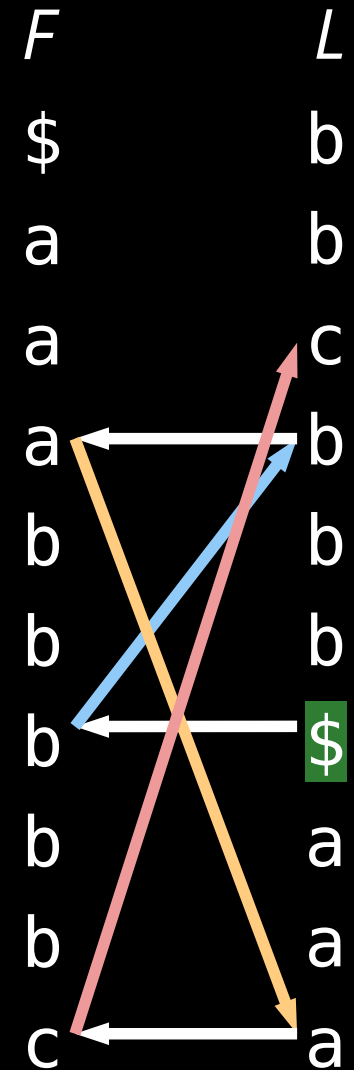
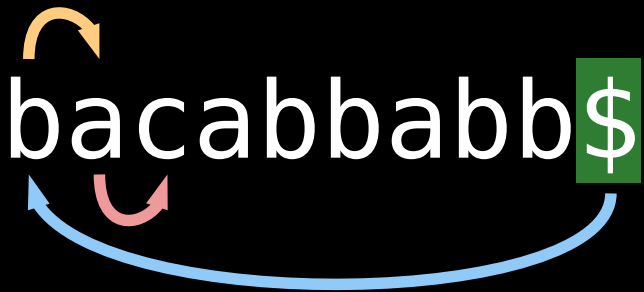


<i>F</i>	<i>L</i>
\$	b
a	b
a	c
a	b
b	b
b	b
b	\$
b	a
b	a
c	a



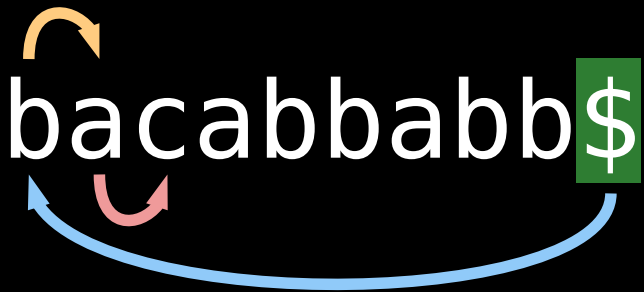
forward search

$T = \text{bacabbabb}\$$

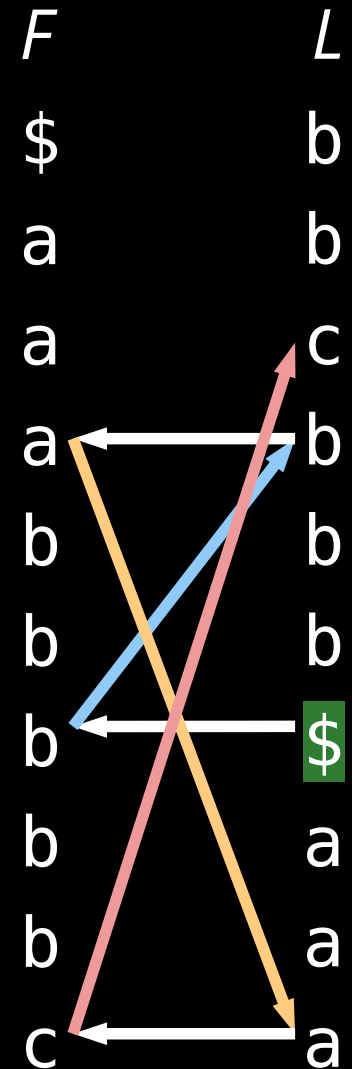


forward search

$T = \text{bacabbabb}\$$



can calculate with
rank and select on F and L



forward search

$T = \text{bacabbabb\$}$

FL mapping:

$$FL(i) = L.\text{select}_{F[i]}(F.\text{rank}_{F[i]}(F[i]))$$

$F.\text{rank}_{F[i]}(F[i])$

$L.\text{rank}_{L[i]}(L[i])$

	F		L
1	\$	b	1
1	a	b	2
2	a	c	1
3	a	b	3
1	b	b	4
2	b	b	5
3	b	\$	1
4	b	a	1
5	b	a	2
1	c	a	3

backward search

$T = \text{ba}\text{c}\text{abbabb}\$$

	F		L	
	1	\$	b	1
	1	a	b	2
	2	a	c	1
	3	a	b	3
	1	b	b	4
	2	b	b	5
	3	b	\$	1
	4	b	a	1
	5	b	a	2
	1	c	a	3

$F.\text{rank}_{F[i]}(F[i])$

$L.\text{rank}_{L[i]}(L[i])$

backward search

$T = \text{bacabbabb\$}$



$L.\text{rank}_{L[i]}(L[i])$

$F.\text{rank}_{F[i]}(F[i])$

	F		L
1	\$	b	1
1	a	b	2
2	a	c	1
3	a	b	3
1	b	b	4
2	b	b	5
3	b	\$	1
4	b	a	1
5	b	a	2
1	c	a	3

backward search

$T = \text{bacabbabb\$}$

$F.\text{rank}_{F[i]}(F[i])$

$L.\text{rank}_{L[i]}(L[i])$

	F		L
1	\$	b	1
1	a	b	2
2	a	c	1
3	a	b	3
1	b	b	4
2	b	b	5
3	b	\$	1
4	b	a	1
5	b	a	2
1	c	a	3

backward search

$T = \text{bacabbabb\$}$

$F.\text{rank}_{F[i]}(F[i])$

$L.\text{rank}_{L[i]}(L[i])$

	F		L
1	\$	b	1
1	a	b	2
2	a	c	1
3	a	b	3
1	b	b	4
2	b	b	5
3	b	\$	1
4	b	a	1
5	b	a	2
1	c	a	3

backward search

$T = \text{bacabbabb\$}$

LF mapping:

$\text{LF}(i) := F.\text{select}_{L[i]}(L.\text{rank}_{L[i]}(i))$

$L.\text{rank}_{L[i]}(L[i])$

$F.\text{rank}_{F[i]}(F[i])$

	F		L
1	\$	b	1
1	a	b	2
2	a	c	1
3	a	b	3
1	b	b	4
2	b	b	5
3	b	\$	1
4	b	a	1
5	b	a	2
1	c	a	3

backward search

$T = \text{bacabbabb\$}$

LF mapping:

$$\text{LF}(i) := F.\text{select}_{L[i]}(L.\text{rank}_{L[i]}(i))$$

$$= F.\text{select}_{L[i]}(1) + L.\text{rank}_{L[i]}(i) - 1$$

$F.\text{rank}_{F[i]}(F[i])$

$L.\text{rank}_{L[i]}(L[i])$

	F		L
1	\$	b	1
1	a	b	2
2	a	c	1
3	a	b	3
1	b	b	4
2	b	b	5
3	b	\$	1
4	b	a	1
5	b	a	2
1	c	a	3

backward search

$T = \text{bacabbabb\$}$

LF mapping:

$$\begin{aligned} \text{LF}(i) &:= F.\text{select}_{L[i]}(L.\text{rank}_{L[i]}(i)) \\ &= F.\text{select}_{L[i]}(1) + L.\text{rank}_{L[i]}(i) - 1 \\ &= |\{j : L[j] < L[i]\}| + L.\text{rank}_{L[i]}(i) \end{aligned}$$

$F.\text{rank}_{F[i]}(F[i])$

$L.\text{rank}_{L[i]}(L[i])$

	F		L
1	\$	b	1
1	a	b	2
2	a	c	1
3	a	b	3
1	b	b	4
2	b	b	5
3	b	\$	1
4	b	a	1
5	b	a	2
1	c	a	3

LF: time complexity

If we store $\text{BWT}(T)$ in L :

- $L[i] = \text{BWT}[i]$: $O(1)$ time

\Rightarrow for any c : $L.\text{rank}_c(i)$ in $O(n)$ time

- $\text{LF}(i) = \underbrace{|\{j : L[j] < L[i]\}|}_{O(n) \text{ time}} + \underbrace{L.\text{rank}_{L[i]}(i)}_{O(n) \text{ time}}$

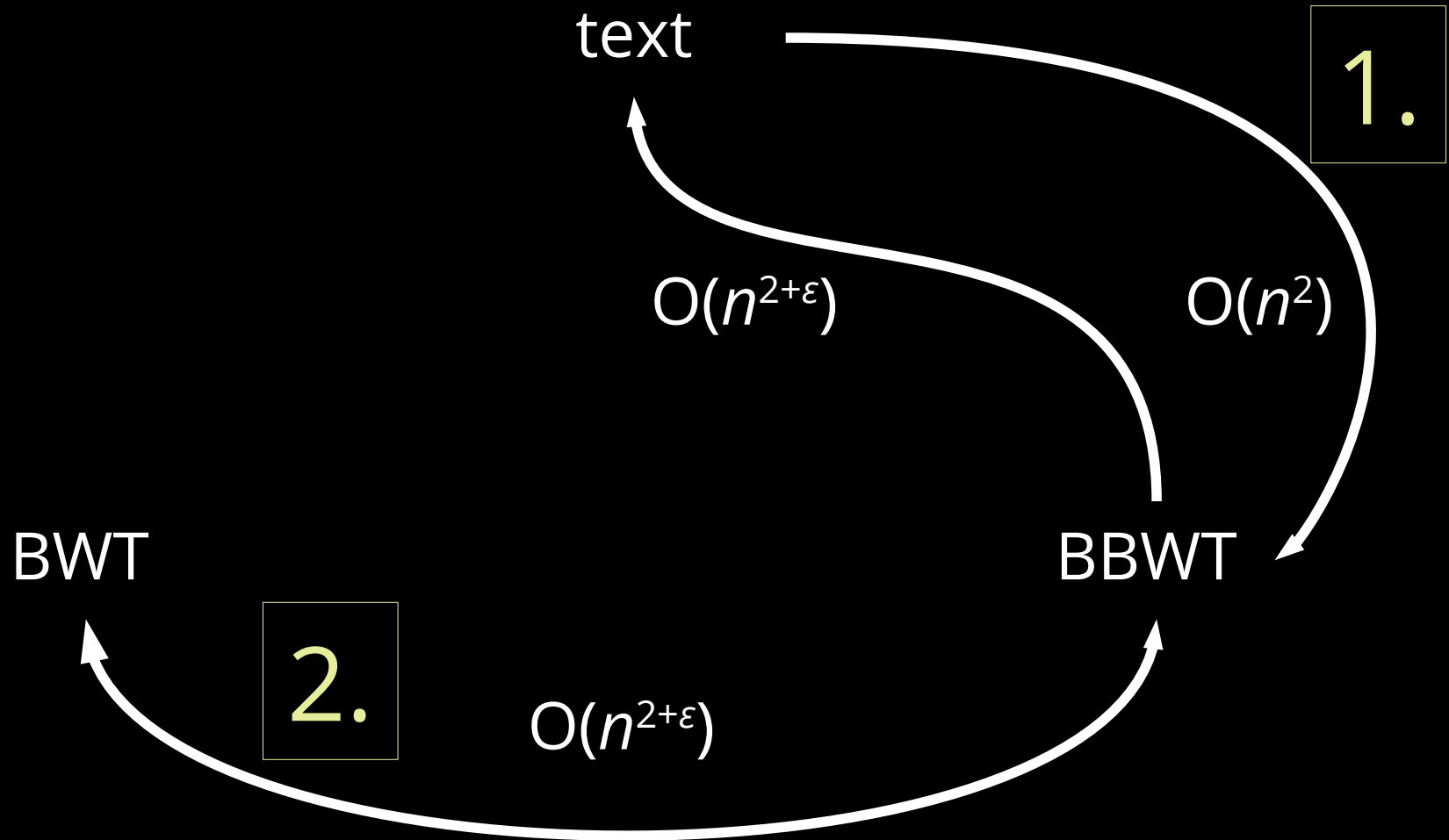
FL: time complexity

- $FL(i) = L.\text{select}_{F[i]}(F.\text{rank}_{F[i]}(F[i]))$
 $= L.\text{select}_{F[i]}(i - |\{j : L[j] < i\}|)$
- If we know $F[i]$: $FL(i)$ in $O(n)$ time
- however, the fastest in-place computation of $F[i]$ takes $O(n^{1+\varepsilon})$ time

[Munro,Raman '96]

for any constant ε with $0 < \varepsilon < 1$

road map



working space: $n \lg \sigma + O(\lg n)$ bits (including text)

text → BBWT

text \rightarrow BBWT

for each Lyndon factor T_x with $x = 1$ up to t :

prepend $T_x[|T_x|]$ to BBWT

$p \leftarrow 1$ (insert position in BBWT)

for each $i = |T_x|-1$ down to 1:

$p \leftarrow \text{LF}(p) + 1$

insert $T_x[i]$ at BBWT[p]

[Bonomo+ '14]

text \rightarrow BBWT

$T = \text{bacabbabb}$

- Lyndon factorization:

$b \mid ac \mid abb \mid abb$

- first: insert b

text \rightarrow BBWT

$T = \text{bacabbabb}$

- Lyndon factorization:

$b \mid ac \mid abb \mid abb$

- first: insert b

	F	L	
1	b	b	1

text \rightarrow BBWT

$T = \text{bacabbabb}$

- Lyndon factorization:

$b \mid ac \mid abb \mid abb$

- first: insert b

	F	L	
1	b	b	1

how to calculate?

	F	L	
1	a	b	1
2	a	b	2
3	a	c	1
1	b	b	3
2	b	b	4
3	b	a	1
4	b	a	2
5	b	b	5
1	c	a	3

BBWT($T_1 T_2$)

$T = b | ac | abb | abb = T_1 T_2 T_3 T_4$

- next Lyndon factor: ac

	F	L	
1	b	b	1

BBWT($T_1 T_2$)

$$T = b | ac | abb | abb = T_1 T_2 T_3 T_4$$

- next Lyndon factor: ac

	F	L	
1	b	b	1

	F	L	
1	b	c	1
1	c	b	1

BBWT($T_1 T_2$)

$$T = b|ac|abb|abb = T_1 T_2 T_3 T_4$$

- next Lyndon factor: ac

	F	L	
1	b	b	1

	F	L	
1	b	c	1
1	c	b	1

	F	L	
1	a	c	1
1	b	b	1
1	c	a	1

BBWT($T_1 T_2 T_3$)

$T = b \mid ac \mid abb \mid abb$

- next Lyndon factor: abb

	F	L	
1	a	c	1
1	b	b	1
1	c	a	1

$$\text{BBWT}(T_1 \ T_2 \ T_3)$$

$T = b \mid ac \mid abb \mid abb$

- next Lyndon factor: abb

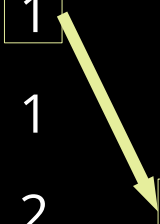
	F	L			F	L	
1	a	c	1	1	a	b	1
1	b	b	1	1	b	c	1
1	c	a	1	2	b	b	2
				1	c	a	1

BBWT($T_1 T_2 T_3$)

$T = b | ac | abb | abb$

- next Lyndon factor: abb

	F	L			F	L			F	L	
1	a	c	1	1	a	b	1	1	a	b	1
1	b	b	1	1	b	c	1	1	b	c	1
1	c	a	1	2	b	b	2	2	b	b	2
				1	c	a	1	3	b	b	3
								1	c	a	1



BBWT($T_1 T_2 T_3$)

$T = b | ac | abb | abb$

- next Lyndon factor: abb

F	L		F	L		F	L		F	L		
1	a	c	1	1	a	b	1	1	1	a	b	1
1	b	b	1	1	b	c	1	1	2	a	c	1
1	c	a	1	2	b	b	2	2	1	b	b	2
			1	1	c	a	1	3	2	b	b	3
								1	1	b	a	1
									1	c	a	2

text \rightarrow BBWT *in-place*

- |bacabbabb

text \rightarrow BBWT *in-place*

- |bacabbabb
- **b**|acabbabb

text \rightarrow BBWT *in-place*

- |bacabbabb
- **b**|acabbabb
- **ba****c**|abbabb

text \rightarrow BBWT *in-place*

- |bacabbabb
- **b**|acabbabb
- **b****ac**|abbabb
- **cba**|abbabb

text \rightarrow BBWT *in-place*

- |bacabbabb
- **b**|acabbabb
- **b****ac**|abbabb
- **cba**|abbabb
- **cba****abb**|abb

text \rightarrow BBWT *in-place*

- |bacabbabb
- **b**|acabbabb
- **b****ac**|abbabb
- **cba**|abbabb
- **cba****abb**|abb
- \vdots

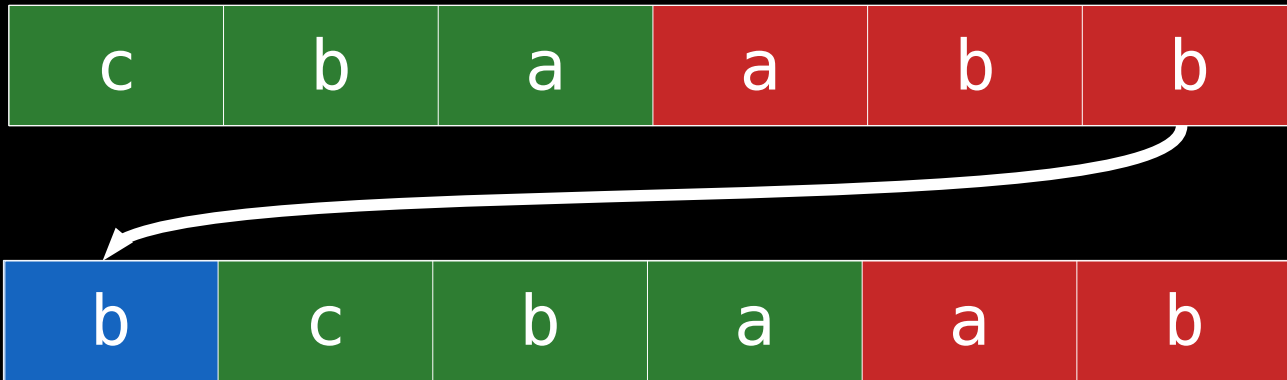
text \rightarrow BBWT *in-place*

- |bacabbabb
- **b**|acabbabb
- **b****ac**|abbabb
- **cba**|abbabb
- **cba****abb**|abb
- \vdots
- **bbcbbbaaba**|

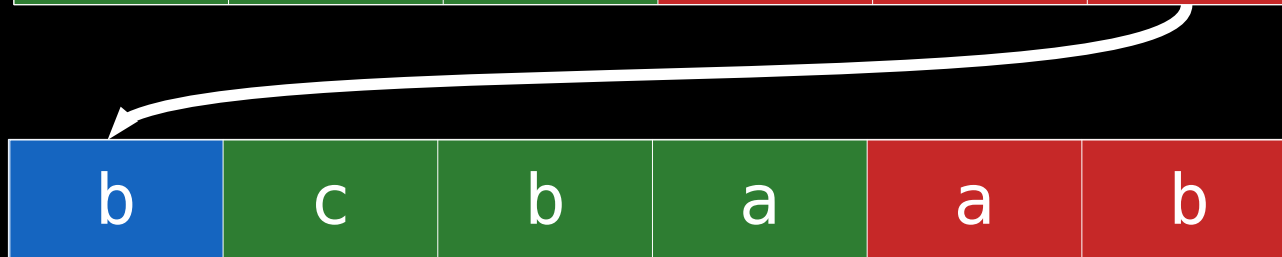
detailed transformation



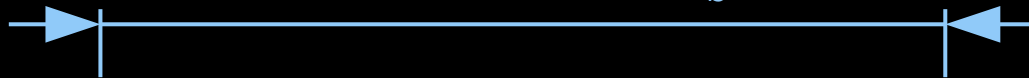
detailed transformation



detailed transformation

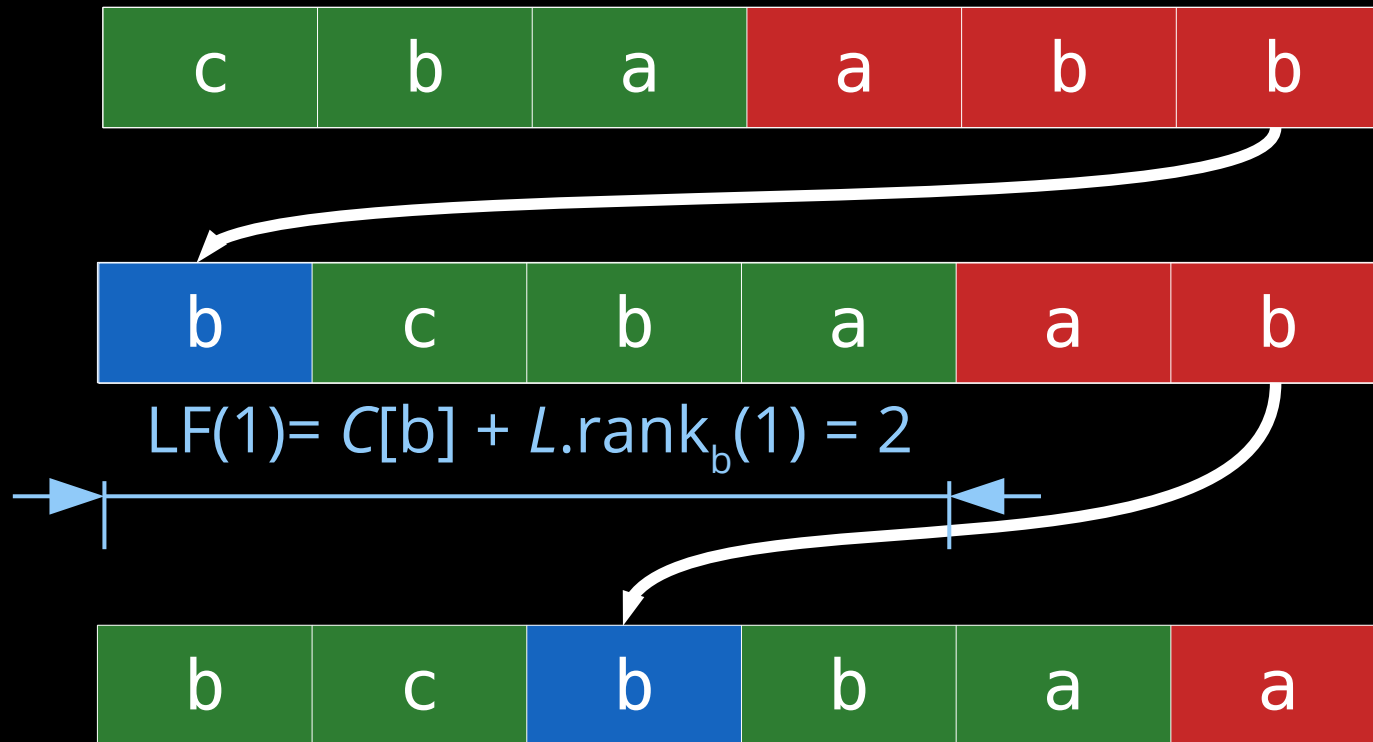


$$LF(1) = C[b] + L.rank_b(1) = 2$$



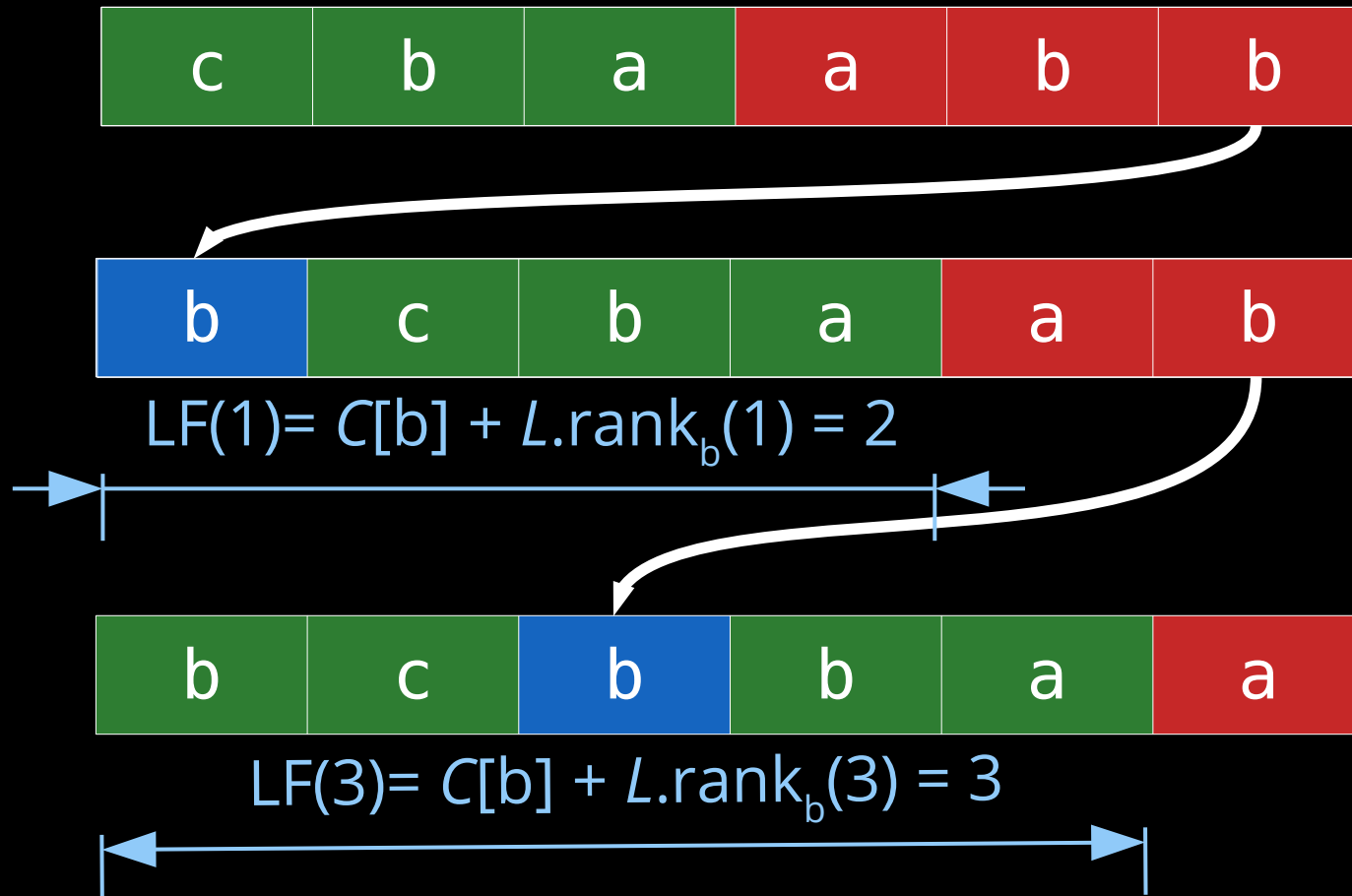
where $C[b] := |\{j : L[j] < b\}|$

detailed transformation



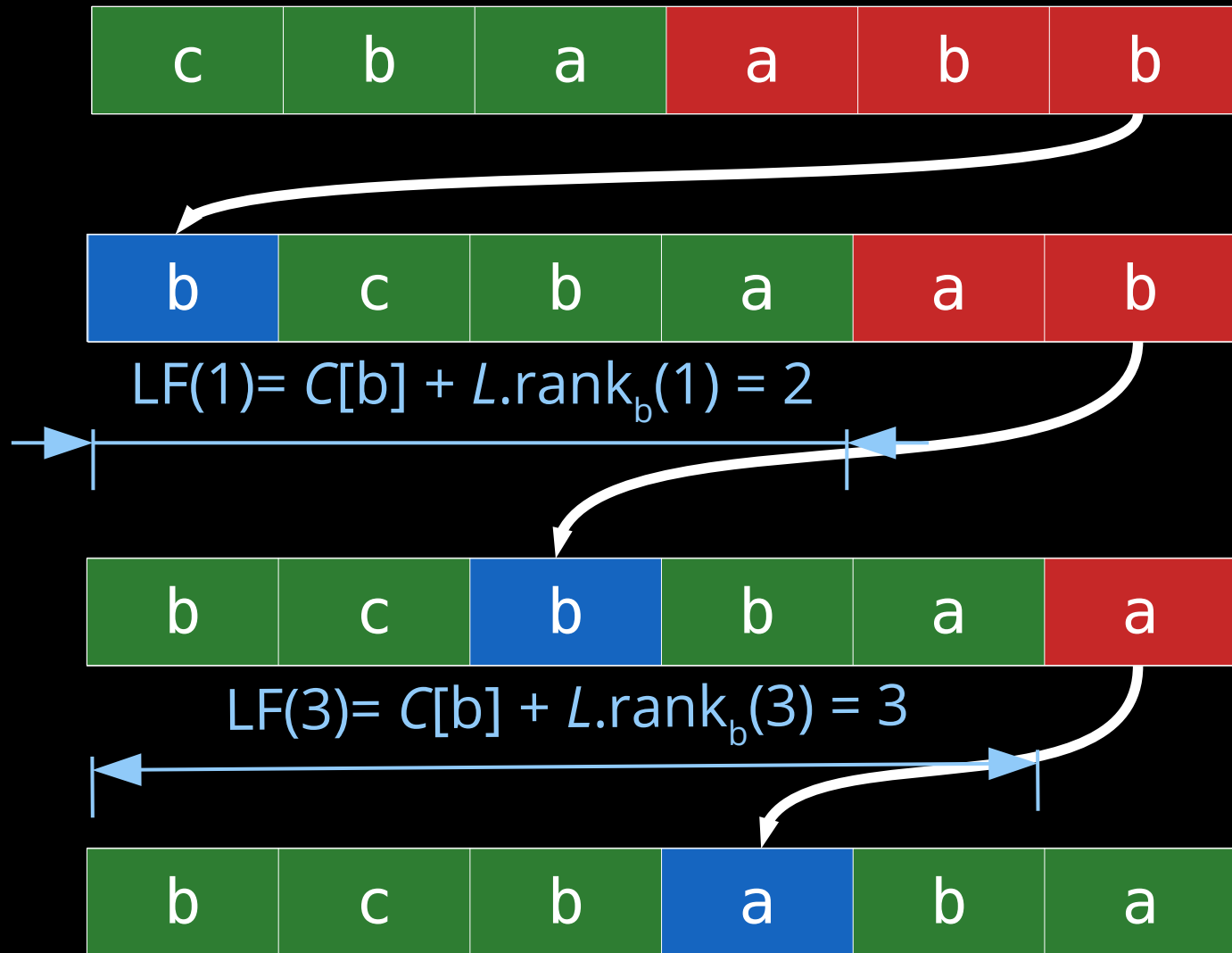
where $C[b] := |\{j : L[j] < b\}|$

detailed transformation



where $C[b] := |\{j : L[j] < b\}|$

detailed transformation



where $C[b] := |\{j : L[j] < b\}|$

BWT → BBWT

BWT \rightarrow BBWT *in-place*

- Duval's algorithm
 - computes Lyndon factorization
 - it runs in $O(n t_L)$ time,
where t_L is the time for accessing an entry of T
 - algorithm uses linear scans from any $T[i]$ to $T[i+1]$
- \Rightarrow emulate this with FL mapping
- \Rightarrow $O(n^{2+\varepsilon})$ time only with L storing BWT




BWT \rightarrow BBWT *in situ*

$T = b | ac | abb | abb$

F		L	
1	\$	b	1
1	a	b	2
2	a	c	1
3	a	b	3
1	b	b	4
2	b	b	5
3	b	\$	1
4	b	a	1
5	b	a	2
1	c	a	3

BWT \rightarrow BBWT *in situ*


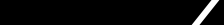

$T = b | ac | abb | abb$

	F		L
1	\$		b 1
1	a		b 2
2	a		c 1
3	a		b 3
1	b		b 4
2	b		b 5
3	b		\$ 1
4	b		a 1
5	b		a 2
1	c		a 3

BWT \rightarrow BBWT *in situ*

$T = b | ac | abb | abb$

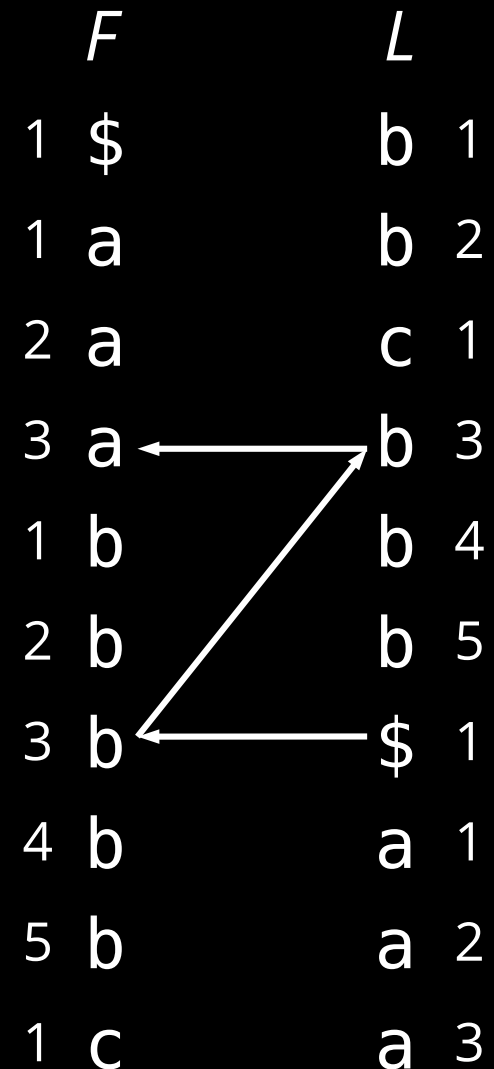
- with FL mapping + Duval
we detect the first Lyndon
factor $b | a \dots$

	<i>F</i>		<i>L</i>
1	\$		b 1
1	a		b 2
2	a		c 1
3	a		b 3
1	b		b 4
2	b		b 5
3	b		\$ 1
4	b		a 1
5	b		a 2
1	c		a 3

construction of a cycle

$T = b | ac | abb | abb$

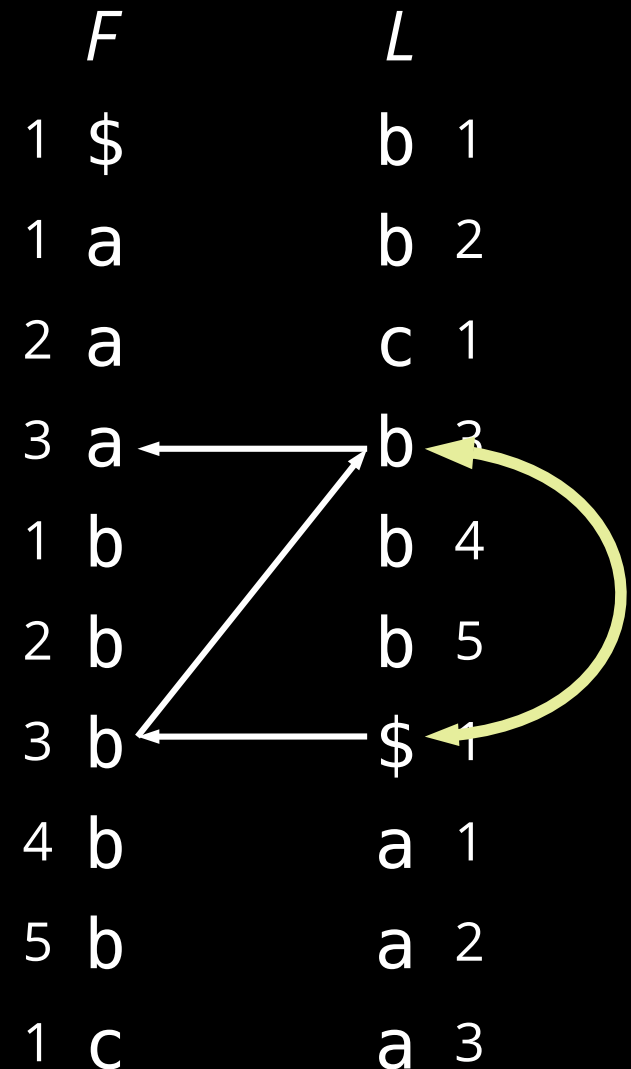
- aim: create cycle $b \rightarrow b$



construction of a cycle

$T = b | ac | abb | abb$

- aim: create cycle $b \rightarrow b$
- since FL maps \$ to $\pi[1]$ we want to exchange \$ and **b**



construction of a cycle

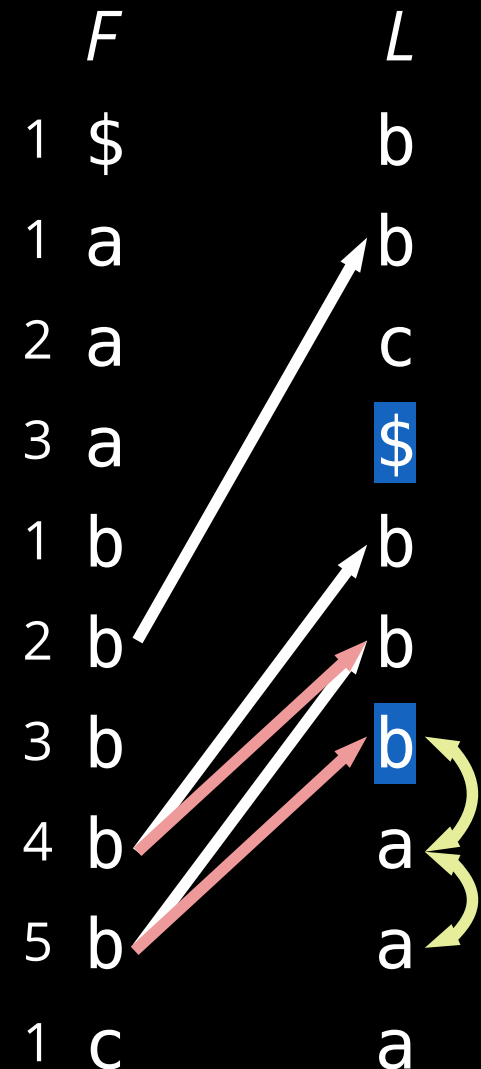
$T = b | ac | abb | abb$

- aim: create cycle $b \rightarrow b$
- since FL maps \$ to $\pi[1]$ we want to exchange \$ and b
- however: might not work
- need to fix red arrows

	<i>F</i>		<i>L</i>		
1	\$		b	1	1
1	a		b	2	2
2	a		c	1	1
3	a		\$	3	1
1	b		b	4	3
2	b		b	5	4
3	b		b	1	5
4	b		a	1	1
5	b		a	2	2
1	c		a	3	3

construction of a cycle

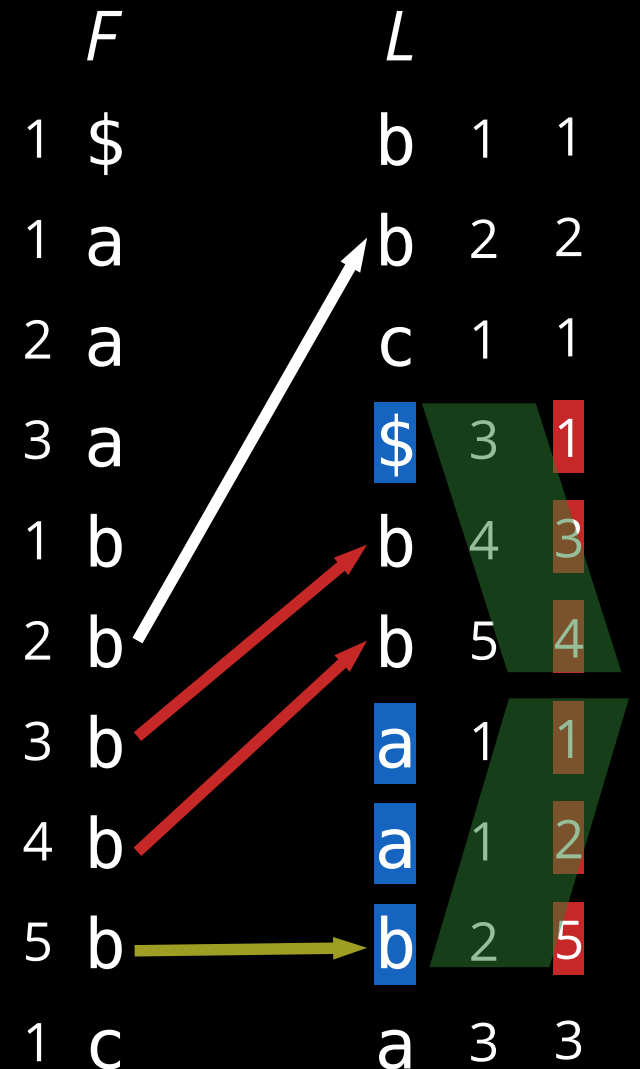
- since there are **two red arrows**:
- switch below the exchanged **b** the next **two entries**



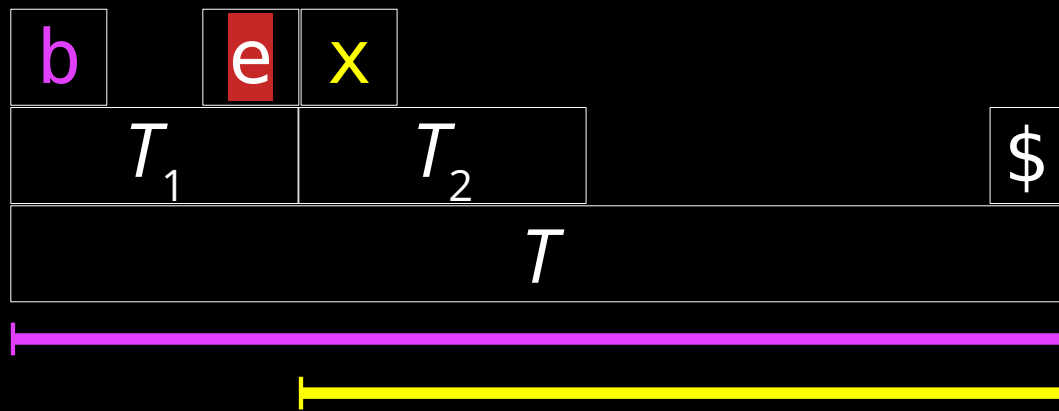
construction of a cycle

- the cycle moved below the exchange

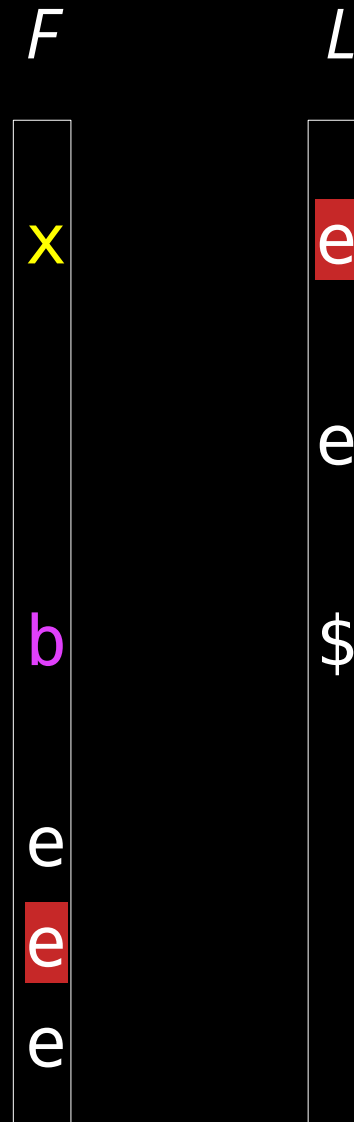
⇒ modified LF mapping just “moved”



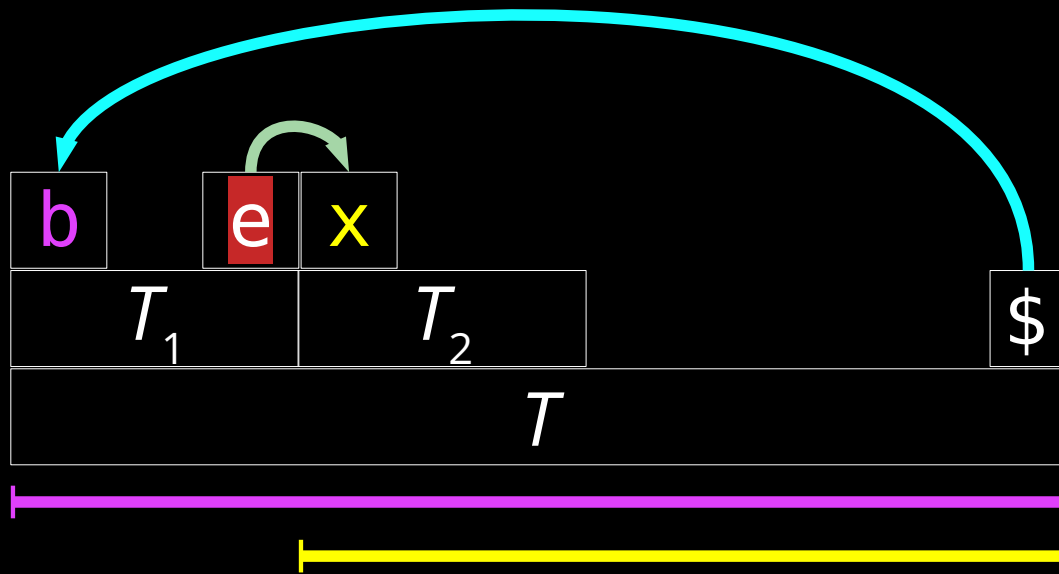
abstract idea



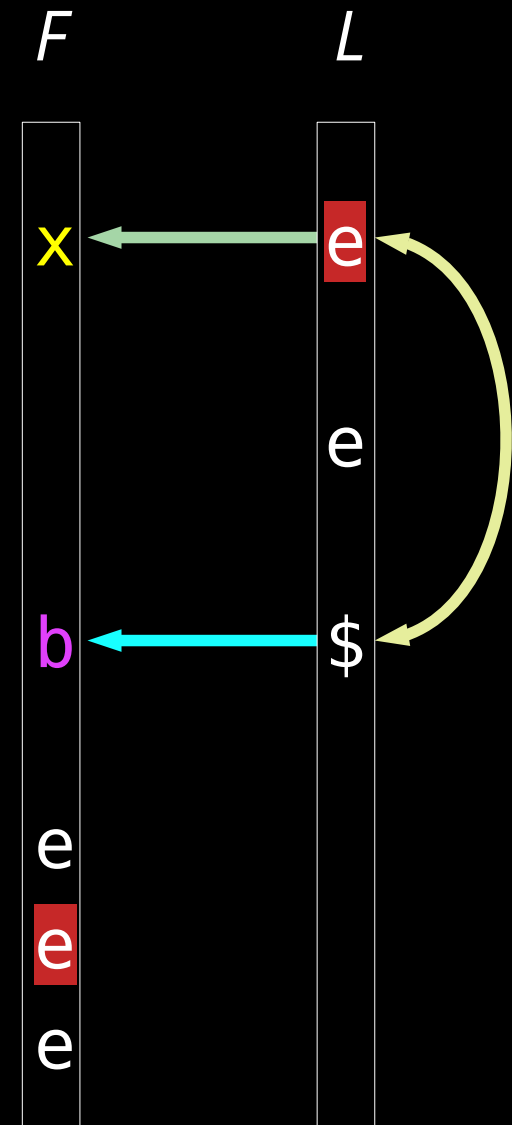
- $T_1 \geq_{\text{lex}} T_2 \geq_{\text{lex}} \dots \geq_{\text{lex}} T_t$
 $\Rightarrow \pi[1..] >_{\text{lex}} \pi[|T_1|..]$



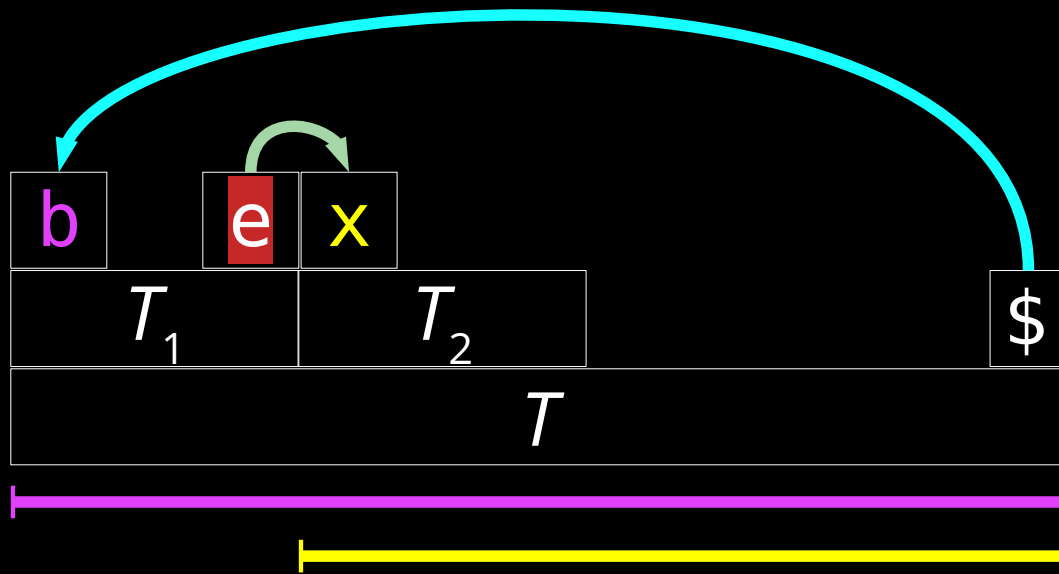
abstract idea



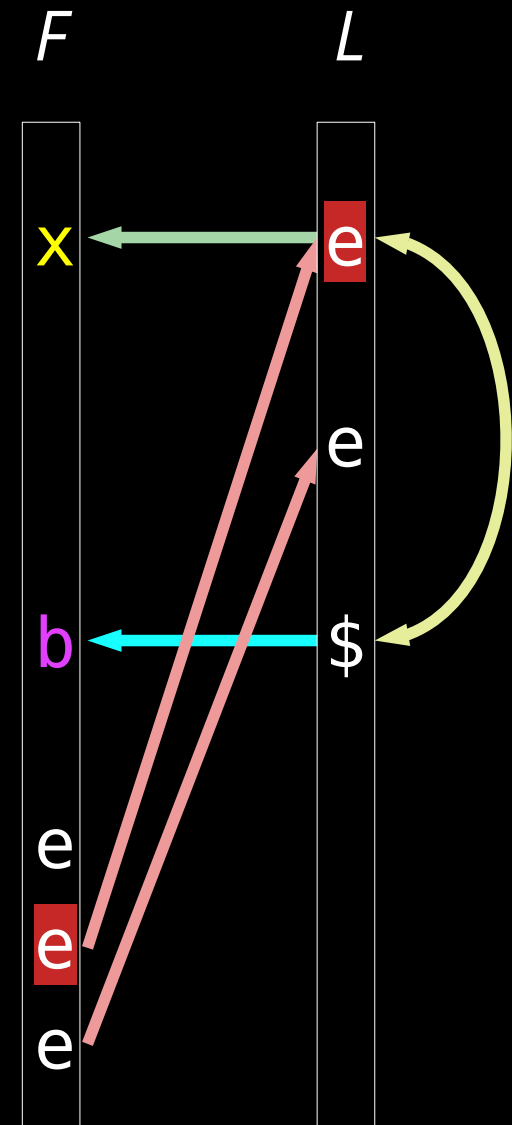
- $T_1 \geq_{\text{lex}} T_2 \geq_{\text{lex}} \dots \geq_{\text{lex}} T_t$
 $\Rightarrow \pi[1..] >_{\text{lex}} \pi[|T_1|..]$



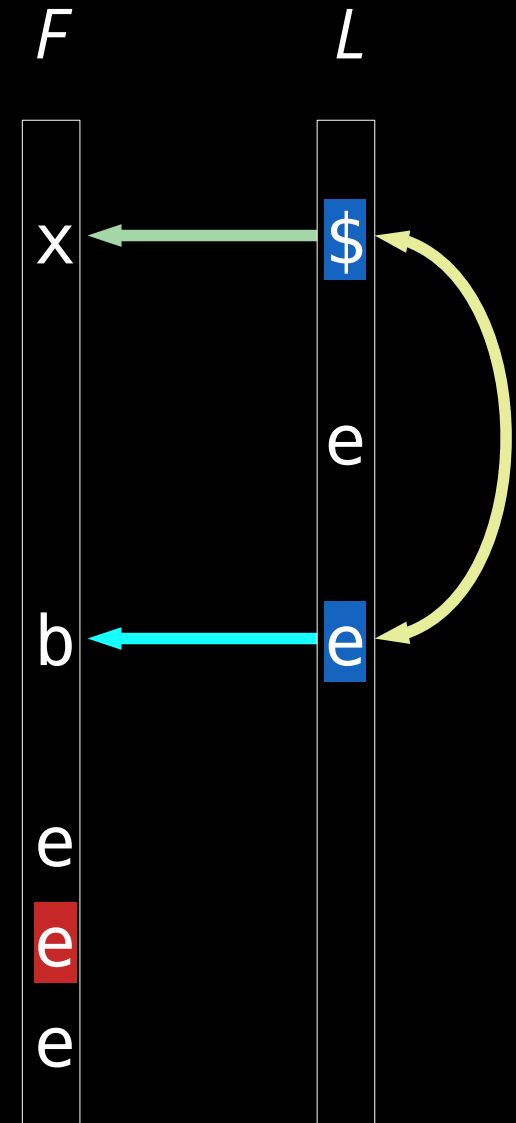
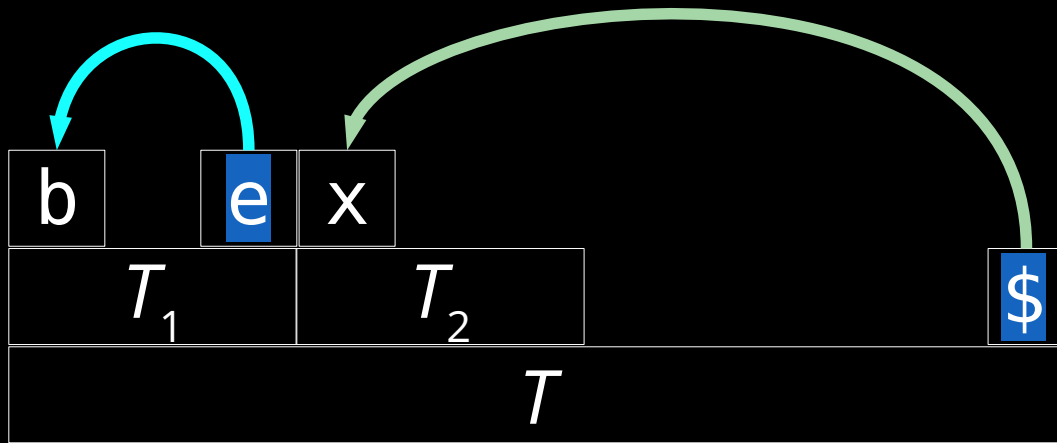
abstract idea



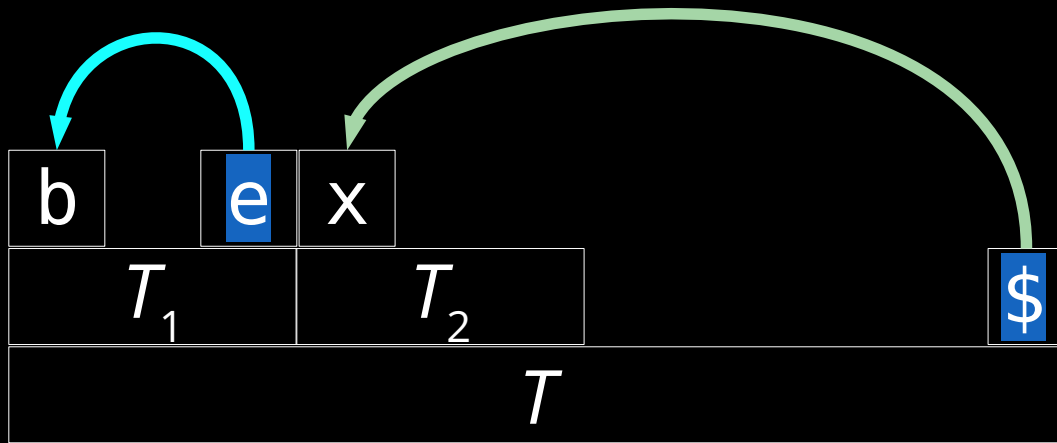
- $T_1 \geq_{\text{lex}} T_2 \geq_{\text{lex}} \dots \geq_{\text{lex}} T_t$
 $\Rightarrow \pi[1..] >_{\text{lex}} \pi[|T_1|..]$
- need to change red arrows



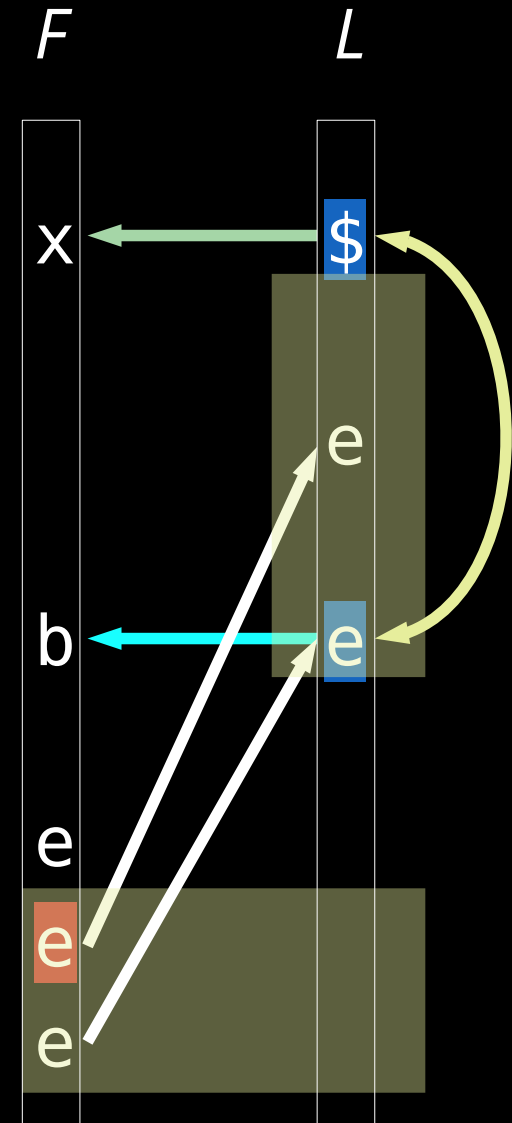
abstract idea



abstract idea



the number of e 's between the exchanged $\$$ and e =
 the number of entries to switch
 after the e in F that mapped to the
 exchanged e



carrying on with example

<i>F</i>	<i>L</i>
1 \$	b 1
1 a	b 2
2 a	c 1
3 a	\$ 1
1 b	b 3
2 b	b 4
3 b	a 1
4 b	a 2
5 b	b 5
1 c	a 3

carrying on with example

	F		L	
1	\$		b	1
1	a		b	2
2	a		c	1
3	a		\$	1
1	b		b	3
2	b		b	4
3	b		a	1
4	b		a	2
5	b		b	5
1	c		a	3

The diagram illustrates a matching between two sets, F and L . The elements of F are listed on the left, and the elements of L are listed on the right. Arrows indicate the following matches:

- From $F=3, a$ to $L=2, c$
- From $F=1, c$ to $L=3, a$
- From $F=3, a$ to $L=1, \$$

carrying on with example

	<i>F</i>		<i>L</i>
1	\$	b	1
1	a	b	2
2	a	c	1
3	a	\$	1
1	b	b	3
2	b	b	4
3	b	a	1
4	b	a	2
5	b	b	5
1	c	a	3

	<i>F</i>		<i>L</i>
1	\$	b	1
1	a	b	2
2	a	c	1
3	a	\$	1
1	b	b	3
2	b	b	4
3	b	a	1
4	b	a	2
5	b	b	5
1	c	a	3

carrying on with example

	<i>F</i>		<i>L</i>	
1	\$		b	1
1	a		b	2
2	a		c	1
3	a		\$	1
1	b		b	3
2	b		b	4
3	b		a	1
4	b		a	2
5	b		b	5
1	c		a	3

	<i>F</i>		<i>L</i>	
1	\$		b	1
1	a		b	2
2	a		c	1
3	a		\$	1
1	b		b	3
2	b		b	4
3	b		a	1
4	b		a	2
5	b		b	5
1	c		a	3

	<i>F</i>		<i>L</i>	
1	\$		b	1
1	a		b	2
2	a		\$	1
3	a		c	1
1	b		b	3
2	b		b	4
3	b		a	1
4	b		a	2
5	b		b	5
1	c		a	3

open problems

$O(n^{1+\varepsilon})$ time

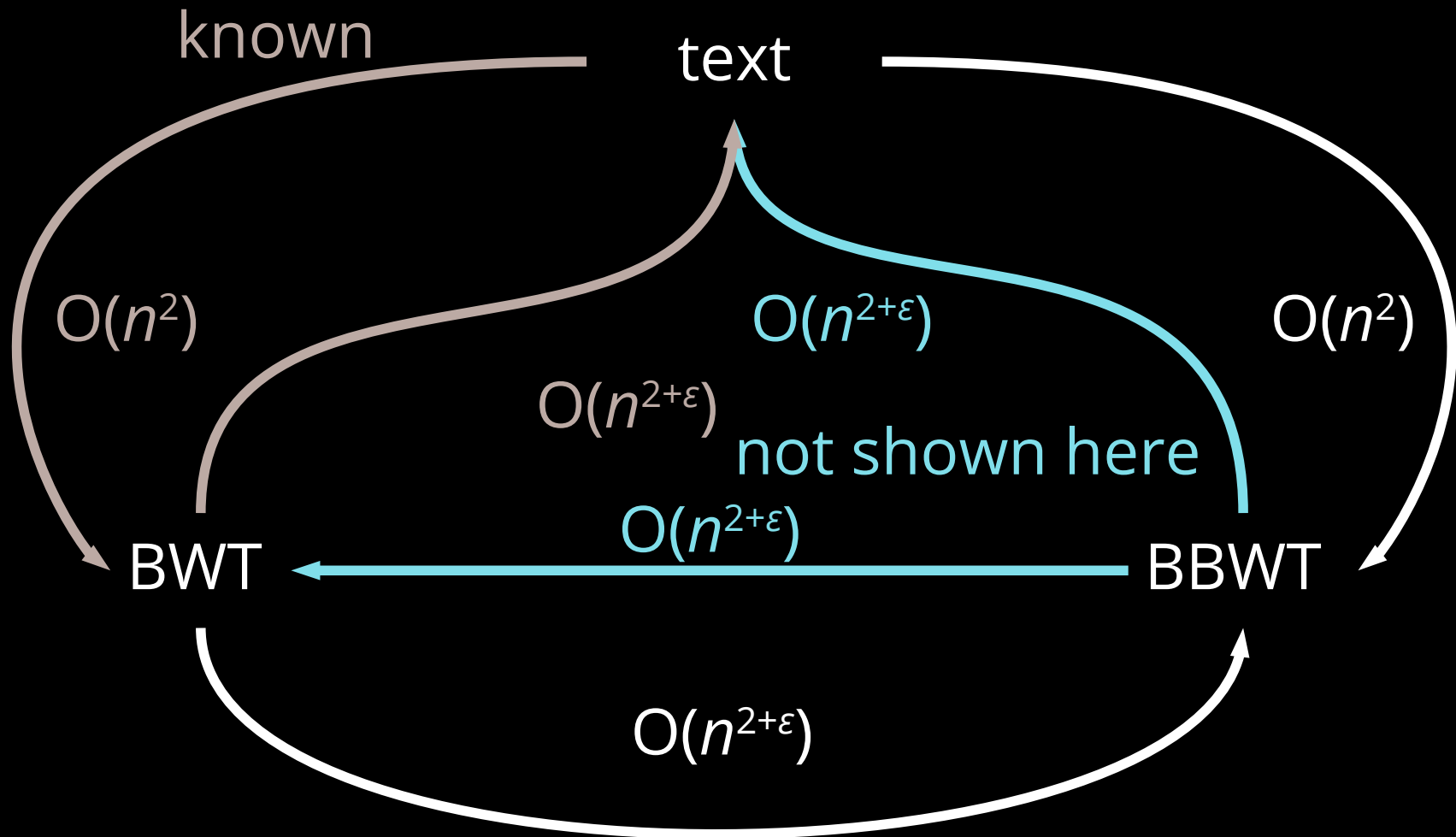
- can we get rid of the FL mapping?
(use only LF mapping)
- trade-off algorithm for time \leftrightarrow space
- Is the number of distinct Lyndon words of T bounded by the runs in the BBWT of T ?

if so:

$O(r)$ words run-length compressed BBWT-index

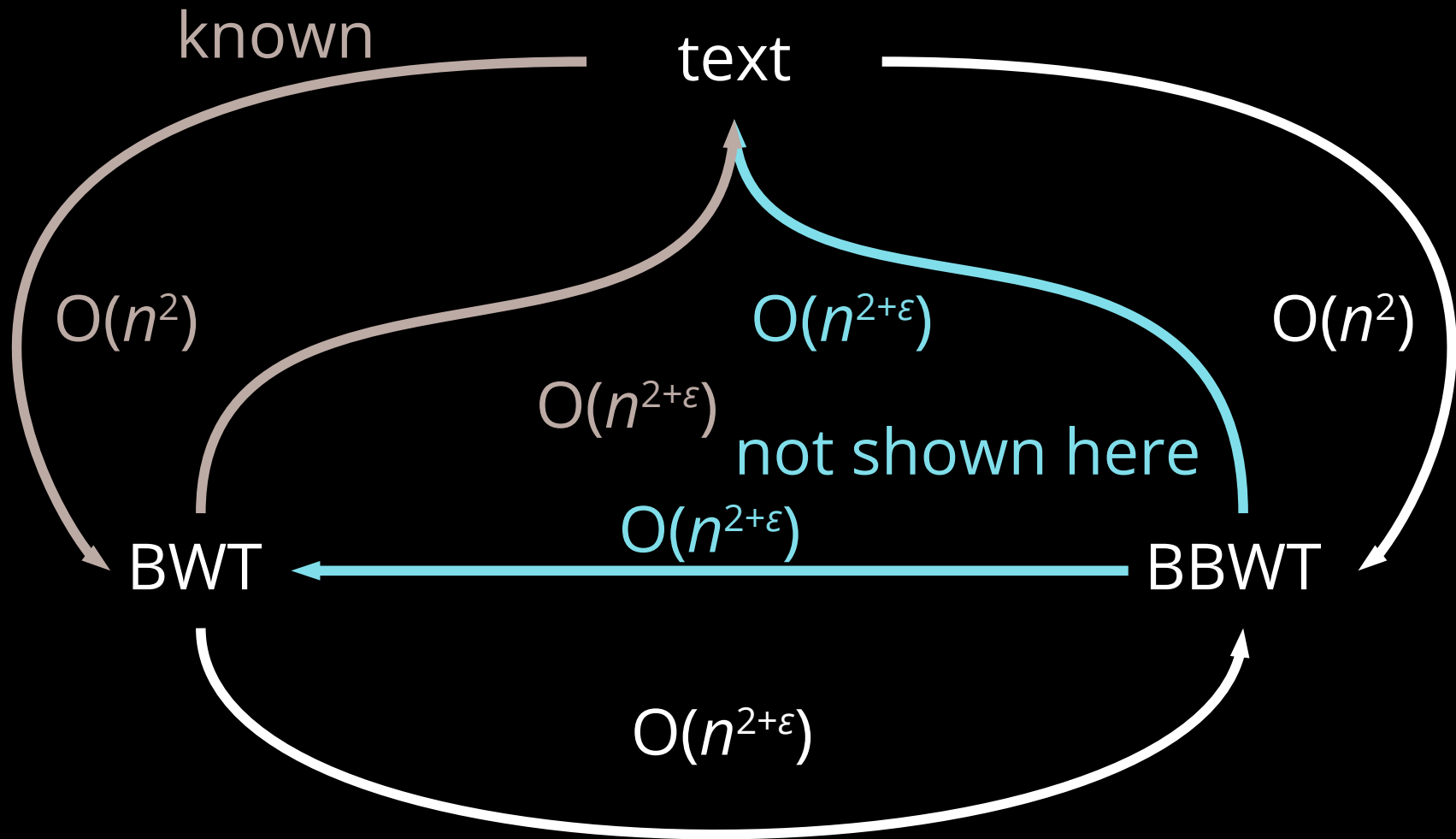
(r : runs in BBWT)

in-place conversions



working space: $n \lg \sigma + O(\lg n)$ bits (including text)

in-place conversions



working space: $n \lg \sigma + O(\lg n)$ bits (including text)

any questions are welcome!