

# Optimal Dynamic Strings

Paweł Gawrychowski<sup>2</sup>, **Tomasz Kociumaka**<sup>1</sup>,  
Adam Karczmarz<sup>1</sup>, Jakub Łącki<sup>3</sup>, Piotr Sankowski<sup>1</sup>

<sup>1</sup>University of Warsaw, Poland

<sup>2</sup>University of Wrocław, Poland

<sup>3</sup>Google Research NY, USA

## **SODA 2018**

New Orleans, LA, USA

January 9, 2018

# Dynamic Strings Problem

Maintain a multiset  $\mathcal{W}$  of non-empty strings subject to:

# Dynamic Strings Problem

Maintain a multiset  $\mathcal{W}$  of non-empty strings subject to:

- `make_string(w)`: insert  $w \in \Sigma^+$  to  $\mathcal{W}$ ;

$\mathcal{W}$ :    ab  
          1

`make_string(ab) = 1`

# Dynamic Strings Problem

Maintain a multiset  $\mathcal{W}$  of non-empty strings subject to:

- `make_string(w)`: insert  $w \in \Sigma^+$  to  $\mathcal{W}$ ;
- `concat(w1, w2)`: insert  $w_1 w_2$  to  $\mathcal{W}$  for  $w_1, w_2 \in \mathcal{W}$ ;

$\mathcal{W}$ :    ab    abab  
          1     2

`concat(1, 1) = 2`

# Dynamic Strings Problem

Maintain a multiset  $\mathcal{W}$  of non-empty strings subject to:

- `make_string(w)`: insert  $w \in \Sigma^+$  to  $\mathcal{W}$ ;
- `concat(w1, w2)`: insert  $w_1 w_2$  to  $\mathcal{W}$  for  $w_1, w_2 \in \mathcal{W}$ ;
- `split(w, k)`: insert  $w[1..k]$  and  $w[k + 1..|w|]$  for  $w \in \mathcal{W}$ ;

$\mathcal{W}$ :	ab	abab	a	bab
	1	2	3	4

`split(2, 1) = (3, 4)`

# Dynamic Strings Problem

Maintain a multiset  $\mathcal{W}$  of non-empty strings subject to:

- `make_string(w)`: insert  $w \in \Sigma^+$  to  $\mathcal{W}$ ;
- `concat(w1, w2)`: insert  $w_1 w_2$  to  $\mathcal{W}$  for  $w_1, w_2 \in \mathcal{W}$ ;
- `split(w, k)`: insert  $w[1..k]$  and  $w[k + 1..|w|]$  for  $w \in \mathcal{W}$ ;
- `equal(w1, w2)`: check whether  $w_1$  is the same string as  $w_2$ ;

$\mathcal{W}$ :    **ab**    abab    **a**    bab  
          1        2        3        4

`equal(1, 3) = false`

# Dynamic Strings Problem

Maintain a multiset  $\mathcal{W}$  of non-empty strings subject to:

- `make_string(w)`: insert  $w \in \Sigma^+$  to  $\mathcal{W}$ ;
- `concat(w1, w2)`: insert  $w_1 w_2$  to  $\mathcal{W}$  for  $w_1, w_2 \in \mathcal{W}$ ;
- `split(w, k)`: insert  $w[1..k]$  and  $w[k + 1..|w|]$  for  $w \in \mathcal{W}$ ;
- `equal(w1, w2)`: check whether  $w_1$  is the same string as  $w_2$ ;

$\mathcal{W}$ :	ab	abab	a	bab	ba	b
	1	2	3	4	5	6

$$\text{split}(4, 2) = (5, 6)$$

# Dynamic Strings Problem

Maintain a multiset  $\mathcal{W}$  of non-empty strings subject to:

- `make_string(w)`: insert  $w \in \Sigma^+$  to  $\mathcal{W}$ ;
- `concat(w1, w2)`: insert  $w_1 w_2$  to  $\mathcal{W}$  for  $w_1, w_2 \in \mathcal{W}$ ;
- `split(w, k)`: insert  $w[1..k]$  and  $w[k + 1..|w|]$  for  $w \in \mathcal{W}$ ;
- `equal(w1, w2)`: check whether  $w_1$  is the same string as  $w_2$ ;

$\mathcal{W}$ :	ab	abab	a	bab	ba	b	ba
	1	2	3	4	5	6	7

$$\text{concat}(6, 3) = 7$$



# Dynamic Strings Problem

Maintain a multiset  $\mathcal{W}$  of non-empty strings subject to:

- `make_string(w)`: insert  $w \in \Sigma^+$  to  $\mathcal{W}$ ;
- `concat(w1, w2)`: insert  $w_1 w_2$  to  $\mathcal{W}$  for  $w_1, w_2 \in \mathcal{W}$ ;
- `split(w, k)`: insert  $w[1..k]$  and  $w[k + 1..|w|]$  for  $w \in \mathcal{W}$ ;
- `equal(w1, w2)`: check whether  $w_1$  is the same string as  $w_2$ ;

$\mathcal{W}$ :	ab	abab	a	bab	ba	b	ba
	1	2	3	4	5	6	7

`equal(5, 7) = true`

# Dynamic Strings Problem

Maintain a multiset  $\mathcal{W}$  of non-empty strings subject to:

- `make_string(w)`: insert  $w \in \Sigma^+$  to  $\mathcal{W}$ ;
- `concat(w1, w2)`: insert  $w_1 w_2$  to  $\mathcal{W}$  for  $w_1, w_2 \in \mathcal{W}$ ;
- `split(w, k)`: insert  $w[1..k]$  and  $w[k + 1..|w|]$  for  $w \in \mathcal{W}$ ;
- `equal(w1, w2)`: check whether  $w_1$  is the same string as  $w_2$ ;

$\mathcal{W}$ :	ab	abab	a	bab	ba	b	ba	baab
	1	2	3	4	5	6	7	8

$$\text{concat}(7, 1) = 8$$

# Dynamic Strings Problem

Maintain a multiset  $\mathcal{W}$  of non-empty strings subject to:

- `make_string(w)`: insert  $w \in \Sigma^+$  to  $\mathcal{W}$ ;
- `concat(w1, w2)`: insert  $w_1 w_2$  to  $\mathcal{W}$  for  $w_1, w_2 \in \mathcal{W}$ ;
- `split(w, k)`: insert  $w[1..k]$  and  $w[k + 1..|w|]$  for  $w \in \mathcal{W}$ ;
- `equal(w1, w2)`: check whether  $w_1$  is the same string as  $w_2$ ;
- `LCP(w1, w2)`: return the length of the longest common prefix of  $w_1$  and  $w_2$ ;

$\mathcal{W}$ :	ab	abab	a	bab	ba	b	ba	baab
	1	2	3	4	5	6	7	8

$$\text{LCP}(4, 8) = 2$$

# Dynamic Strings Problem

Maintain a multiset  $\mathcal{W}$  of non-empty strings subject to:

- `make_string(w)`: insert  $w \in \Sigma^+$  to  $\mathcal{W}$ ;
- `concat(w1, w2)`: insert  $w_1 w_2$  to  $\mathcal{W}$  for  $w_1, w_2 \in \mathcal{W}$ ;
- `split(w, k)`: insert  $w[1..k]$  and  $w[k + 1..|w|]$  for  $w \in \mathcal{W}$ ;
- `equal(w1, w2)`: check whether  $w_1$  is the same string as  $w_2$ ;
- `LCP(w1, w2)`: return the length of the longest common prefix of  $w_1$  and  $w_2$ ;
- `compare(w1, w2)`: lexicographically compare  $w_1$  and  $w_2$ .

$\mathcal{W}$ :	ab	abab	a	bab	ba	b	ba	baab
	1	2	3	4	5	6	7	8

`compare(4, 8) = '>'`

## Setting

$n$  = total length of strings in the collection  
word RAM machine with word size  $\Omega(\log n)$

## Setting

$n$  = total length of strings in the collection  
word RAM machine with word size  $\Omega(\log n)$

	randomization	split, concat	equal	compare, LCP
„folklore”	Monte Carlo	$\mathcal{O}(\log n)$	$\mathcal{O}(1)$	$\mathcal{O}(\log^2 n)$

## Setting

$n$  = total length of strings in the collection  
word RAM machine with word size  $\Omega(\log n)$

	randomization	split, concat	equal	compare, LCP
„folklore”	Monte Carlo	$\mathcal{O}(\log n)$	$\mathcal{O}(1)$	$\mathcal{O}(\log^2 n)$
Mehlhorn et al. 1994	Las Vegas	$\mathcal{O}(\log^2 n)$ exp.	$\mathcal{O}(1)$	-
Mehlhorn et al. 1994	deterministic	$\mathcal{O}(\log^2 n \log^* n)$	$\mathcal{O}(1)$	-

## Setting

$n$  = total length of strings in the collection  
word RAM machine with word size  $\Omega(\log n)$

	randomization	split, concat	equal	compare, LCP
„folklore”	Monte Carlo	$\mathcal{O}(\log n)$	$\mathcal{O}(1)$	$\mathcal{O}(\log^2 n)$
Mehlhorn et al. 1994	Las Vegas	$\mathcal{O}(\log^2 n)$ exp.	$\mathcal{O}(1)$	-
Mehlhorn et al. 1994	deterministic	$\mathcal{O}(\log^2 n \log^* n)$	$\mathcal{O}(1)$	-
Alstrup et al. 2000	Las Vegas <sup>†</sup>	$\mathcal{O}(\log n \log^* n)$ w.h.p.	$\mathcal{O}(1)$	$\mathcal{O}(1)$

<sup>†</sup> due to dynamic dictionaries (hash tables) only.



## Setting

$n$  = total length of strings in the collection  
word RAM machine with word size  $\Omega(\log n)$

	randomization	split, concat	equal	compare, LCP
„folklore”	Monte Carlo	$\mathcal{O}(\log n)$	$\mathcal{O}(1)$	$\mathcal{O}(\log^2 n)$
Mehlhorn et al. 1994	Las Vegas	$\mathcal{O}(\log^2 n)$ exp.	$\mathcal{O}(1)$	-
Mehlhorn et al. 1994	deterministic	$\mathcal{O}(\log^2 n \log^* n)$	$\mathcal{O}(1)$	-
Alstrup et al. 2000	Las Vegas <sup>†</sup>	$\mathcal{O}(\log n \log^* n)$ w.h.p.	$\mathcal{O}(1)$	$\mathcal{O}(1)$
this work	Las Vegas	$\mathcal{O}(\log n)$ w.h.p.	$\mathcal{O}(1)$	$\mathcal{O}(1)$

<sup>†</sup> due to dynamic dictionaries (hash tables) only.

# Previous & Our Results

## Setting

$n$  = total length of strings in the collection  
word RAM machine with word size  $\Omega(\log n)$

	randomization	split, concat	equal	compare, LCP
„folklore”	Monte Carlo	$\mathcal{O}(\log n)$	$\mathcal{O}(1)$	$\mathcal{O}(\log^2 n)$
Mehlhorn et al. 1994	Las Vegas	$\mathcal{O}(\log^2 n)$ exp.	$\mathcal{O}(1)$	-
Mehlhorn et al. 1994	deterministic	$\mathcal{O}(\log^2 n \log^* n)$	$\mathcal{O}(1)$	-
Alstrup et al. 2000	Las Vegas <sup>†</sup>	$\mathcal{O}(\log n \log^* n)$ w.h.p.	$\mathcal{O}(1)$	$\mathcal{O}(1)$
this work	Las Vegas	$\mathcal{O}(\log n)$ w.h.p.	$\mathcal{O}(1)$	$\mathcal{O}(1)$

<sup>†</sup> due to dynamic dictionaries (hash tables) only.

## Unconditional Lower Bound (this work)

No Monte Carlo algorithm (correct with high probability) supports split, concat, and equal operations in  $o(\log n)$  amortized time. This is true even if split and concat invalidate their arguments.

## Karp–Rabin Fingerprints

- $H(w) = H(w') \implies w = w'$  w.h.p.
- $H(u), H(v) \rightsquigarrow H(uv)$  in  $\mathcal{O}(1)$  time.

## Karp–Rabin Fingerprints

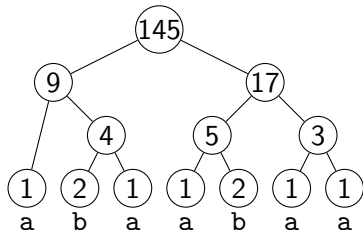
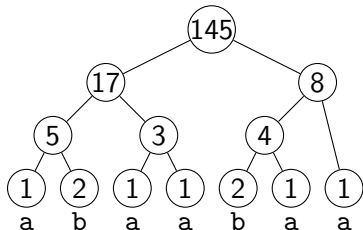
- $H(w) = H(w') \implies w = w'$  w.h.p.
- $H(u), H(v) \rightsquigarrow H(uv)$  in  $\mathcal{O}(1)$  time.

**Idea:** store each string in a persistent BST supporting split & join.

## Karp–Rabin Fingerprints

- $H(w) = H(w') \implies w = w'$  w.h.p.
- $H(u), H(v) \rightsquigarrow H(uv)$  in  $\mathcal{O}(1)$  time.

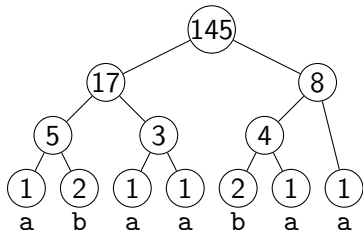
**Idea:** store each string in a persistent BST supporting split & join.



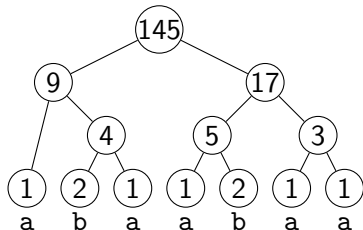
## Karp–Rabin Fingerprints

- $H(w) = H(w') \implies w = w'$  w.h.p.
- $H(u), H(v) \rightsquigarrow H(uv)$  in  $\mathcal{O}(1)$  time.

**Idea:** store each string in a persistent BST supporting split & join.



$\text{split}(w, k) \mathcal{O}(\log |w|)$   
 $\text{concat}(w_1, w_2) \mathcal{O}(\log |w_1 w_2|)$



$\text{make\_string}(w) \mathcal{O}(|w|)$   
 $\text{equal}(w_1, w_2) \mathcal{O}(1)$

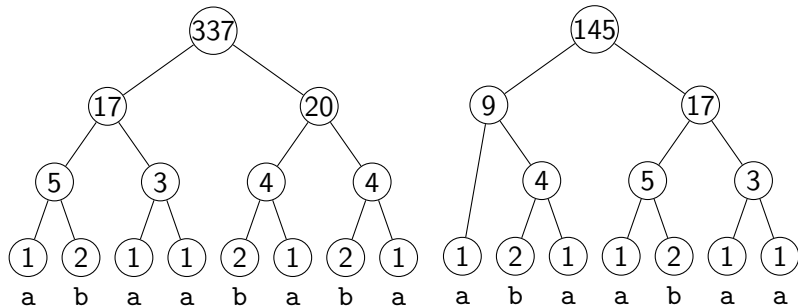
- 1 Monte Carlo: equal is correct w.h.p. only

- 1 Monte Carlo: `equal` is correct w.h.p. only
- 2 Answering `LCP(w1, w2)` and `compare(w1, w2)` is slow.



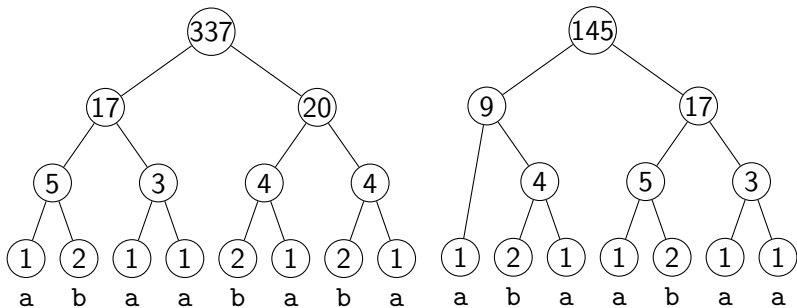
# „Folklore” Solution: Issues

- 1 Monte Carlo: equal is correct w.h.p. only
- 2 Answering  $LCP(w_1, w_2)$  and  $compare(w_1, w_2)$  is slow.



# „Folklore” Solution: Issues

- 1 Monte Carlo: equal is correct w.h.p. only
- 2 Answering  $LCP(w_1, w_2)$  and  $compare(w_1, w_2)$  is slow.



- The BST structure is not particularly helpful...
- Fall back to the naive solution:
  - LCP: binary search using split and equal;  $\mathcal{O}(\log^2 |w|)$  time;
  - compare: compare the characters following LCP.

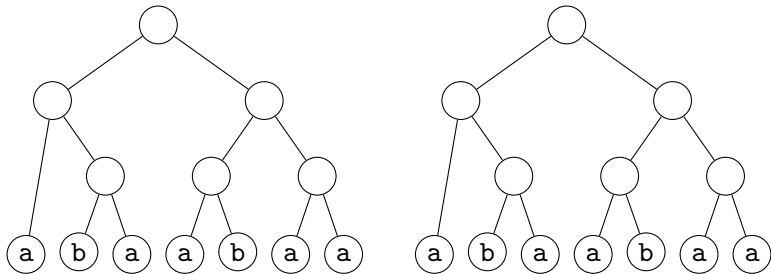
Consistent Parsing (Mehlhorn et al. SODA'94; Sahinalp–Vishkin STOC'94)

Shape of the **parse tree** is determined by the underlying string.

# Consistent Parsing

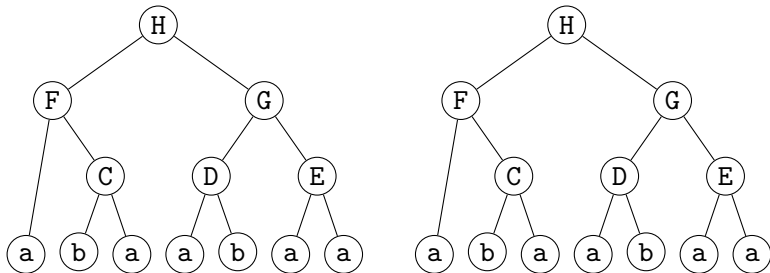
Consistent Parsing (Mehlhorn et al. SODA'94; Sahinalp–Vishkin STOC'94)

Shape of the **parse tree** is determined by the underlying string.



Consistent Parsing (Mehlhorn et al. SODA'94; Sahinalp–Vishkin STOC'94)

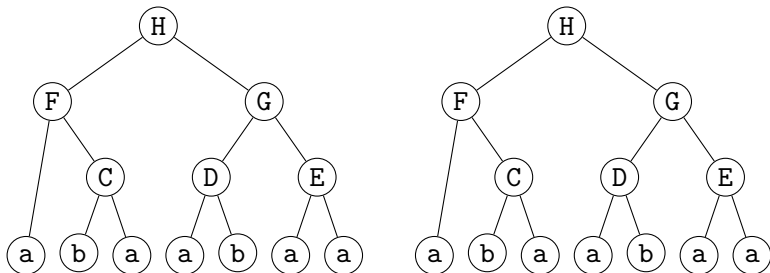
Shape of the **parse tree** is determined by the underlying string.



- **Name** each node based on the names of its children.

Consistent Parsing (Mehlhorn et al. SODA'94; Sahinalp–Vishkin STOC'94)

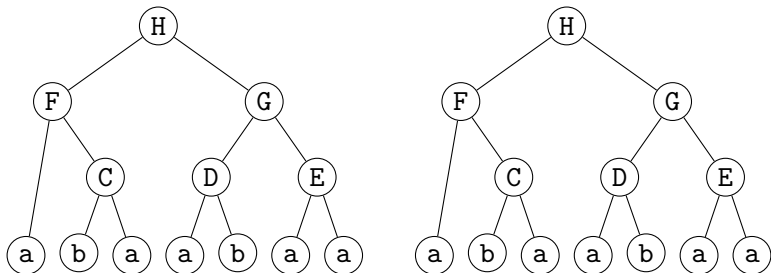
Shape of the **parse tree** is determined by the underlying string.



- **Name** each node based on the names of its children.
- Compare the roots' names for equal with no errors.

Consistent Parsing (Mehlhorn et al. SODA'94; Sahinalp–Vishkin STOC'94)

Shape of the **parse tree** is determined by the underlying string.



- **Name** each node based on the names of its children.
- Compare the roots' names for equal with no errors.
- Dictionaries needed for the naming function.

# Locally Consistent Parsing

Locally Consistent Parsing (Mehlhorn et al.; Sahinalp–Vishkin)

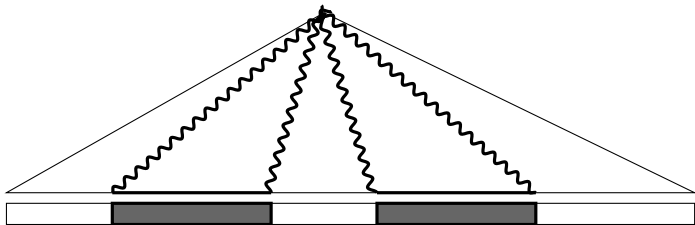
Equal **fragments** are parsed **almost** in the same way.



# Locally Consistent Parsing

Locally Consistent Parsing (Mehlhorn et al.; Sahinalp–Vishkin)

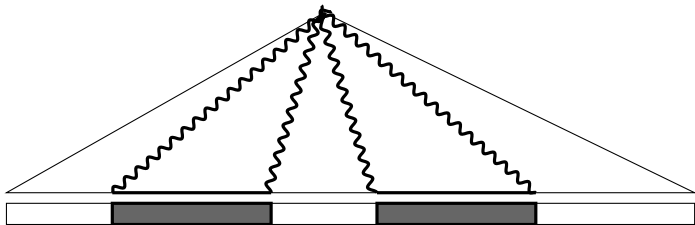
Equal **fragments** are parsed **almost** in the same way.



# Locally Consistent Parsing

Locally Consistent Parsing (Mehlhorn et al.; Sahinalp–Vishkin)

Equal **fragments** are parsed **almost** in the same way.

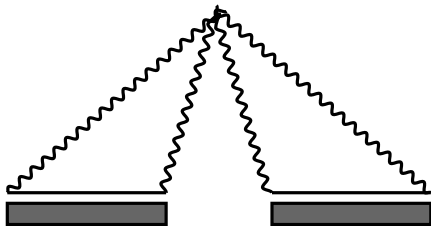


Why is local consistence useful?

# Locally Consistent Parsing

Locally Consistent Parsing (Mehlhorn et al.; Sahinalp–Vishkin)

Equal **fragments** are parsed **almost** in the same way.

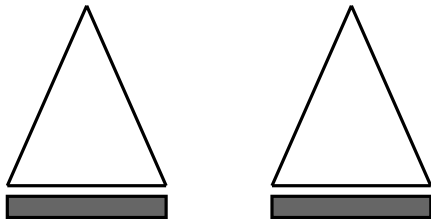


Why is local consistence useful?

- allows to maintain consistent parsing,

Locally Consistent Parsing (Mehlhorn et al.; Sahinalp–Vishkin)

Equal **fragments** are parsed **almost** in the same way.

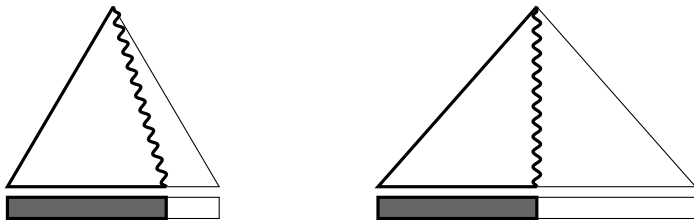


Why is local consistence useful?

- allows to maintain consistent parsing,

Locally Consistent Parsing (Mehlhorn et al.; Sahinalp–Vishkin)

Equal **fragments** are parsed **almost** in the same way.



Why is local consistence useful?

- allows to maintain consistent parsing,
- enables efficient LCP (and compare) queries.

How to parse highly-repetitive fragments?

# Parse Tree Construction: RLE

How to parse highly-repetitive fragments?

a a a a a a a a

# Parse Tree Construction: RLE

How to parse highly-repetitive fragments?

a a a a a a a a



# Parse Tree Construction: RLE

How to parse highly-repetitive fragments?

a a a a a a a a

# Parse Tree Construction: RLE

How to parse highly-repetitive fragments?

a a a **a a a a** a

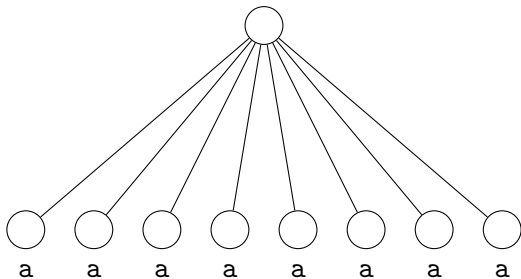
# Parse Tree Construction: RLE

How to parse highly-repetitive fragments?

a a a a a a a a

# Parse Tree Construction: RLE

How to parse highly-repetitive fragments?

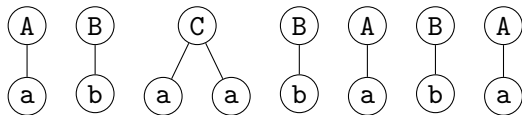


## Run-Length Encoding (RLE)

Replace each **run** with a new symbol:  $\underbrace{a \cdots a}_{k \text{ times}} \mapsto (a, k)$ .

$aaabbabaabbbaa \mapsto (a, 3)(b, 2)(a, 1)(b, 1)(a, 2)(b, 3)(a, 2)$

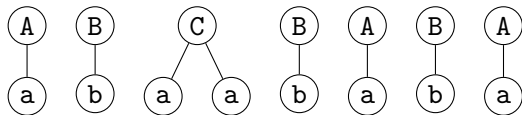
# Parse Tree Construction: COMPRESS



RLE

$A = (a, 1)$ ,  $B = (b, 1)$ ,  $C = (a, 2)$

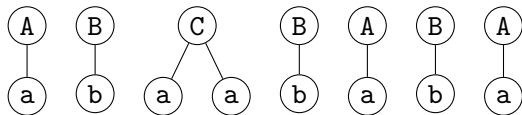
# Parse Tree Construction: COMPRESS



RLE  
 $A = (a, 1)$ ,  $B = (b, 1)$ ,  $C = (a, 2)$

How to **consistently** partition a string without non-trivial runs?

# Parse Tree Construction: COMPRESS



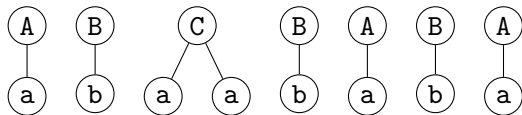
RLE  
 $A = (a, 1), B = (b, 1), C = (a, 2)$

How to **consistently** partition a string without non-trivial runs?

Mehlhorn et al. Local minima wrt. a **random order**:

- Block length:  $2 - \mathcal{O}(\log n)$  w.h.p.,  $\mathcal{O}(1)$  in expectation.
- Context size:  $\mathcal{O}(1)$ .

# Parse Tree Construction: COMPRESS



RLE  
 $A = (a, 1), B = (b, 1), C = (a, 2)$

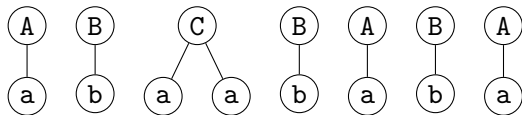
How to **consistently** partition a string without non-trivial runs?

Mehlhorn et al. Local minima wrt. a **random order**:

- Block length:  $2 - \mathcal{O}(\log n)$  w.h.p.,  $\mathcal{O}(1)$  in expectation.
- Context size:  $\mathcal{O}(1)$ .



# Parse Tree Construction: COMPRESS



RLE  
 $A = (a, 1), B = (b, 1), C = (a, 2)$

How to **consistently** partition a string without non-trivial runs?

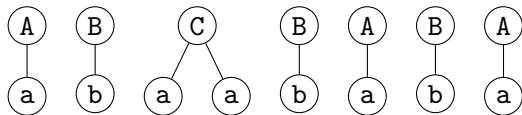
Mehlhorn et al. Local minima wrt. a **random order**:

- Block length:  $2 - \mathcal{O}(\log n)$  w.h.p.,  $\mathcal{O}(1)$  in expectation.
- Context size:  $\mathcal{O}(1)$ .

Mehlhorn et al. **Deterministic coin tossing** (by Cole & Vishkin):

- Blocks length:  $2-4$ ; context size:  $\mathcal{O}(\log^* n)$ .

# Parse Tree Construction: COMPRESS



RLE  
 $A = (a, 1), B = (b, 1), C = (a, 2)$

How to **consistently** partition a string without non-trivial runs?

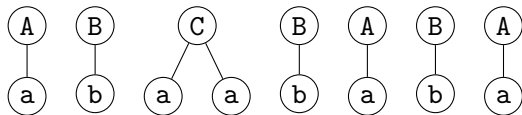
Mehlhorn et al. Local minima wrt. a **random order**:

- Block length:  $2 - \mathcal{O}(\log n)$  w.h.p.,  $\mathcal{O}(1)$  in expectation.
- Context size:  $\mathcal{O}(1)$ .

Mehlhorn et al. **Deterministic coin tossing** (by Cole & Vishkin):

- Blocks length:  $2-4$ ; context size:  $\mathcal{O}(\log^* n)$ .

# Parse Tree Construction: COMPRESS



RLE  
 $A = (a, 1), B = (b, 1), C = (a, 2)$

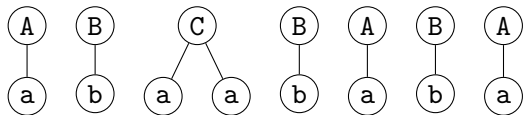
How to **consistently** partition a string without non-trivial runs?

Mehlhorn et al. Local minima wrt. a **random order**:

- Block length:  $2 - \mathcal{O}(\log n)$  w.h.p.,  $\mathcal{O}(1)$  in expectation.
- Context size:  $\mathcal{O}(1)$ .

Mehlhorn et al. **Deterministic coin tossing** (by Cole & Vishkin):

- Blocks length:  $2-4$ ; context size:  $\mathcal{O}(\log^* n)$ .
- Reused by Alstrup et al.



RLE  
 $A = (a, 1)$ ,  $B = (b, 1)$ ,  $C = (a, 2)$

How to **consistently** partition a string without non-trivial runs?

Mehlhorn et al. Local minima wrt. a **random order**:

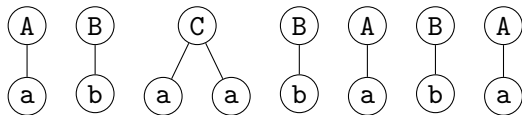
- Block length:  $2 - \mathcal{O}(\log n)$  w.h.p.,  $\mathcal{O}(1)$  in expectation.
- Context size:  $\mathcal{O}(1)$ .

Mehlhorn et al. **Deterministic coin tossing** (by Cole & Vishkin):

- Blocks length:  $2-4$ ; context size:  $\mathcal{O}(\log^* n)$ .
- Reused by Alstrup et al.

Jež (**Recompression**) Alphabet partitioning:

- Partition the alphabet into **left** and **right** symbols.



RLE  
 $A = (a, 1), B = (b, 1), C = (a, 2)$

How to **consistently** partition a string without non-trivial runs?

Mehlhorn et al. Local minima wrt. a **random order**:

- Block length:  $2 - \mathcal{O}(\log n)$  w.h.p.,  $\mathcal{O}(1)$  in expectation.
- Context size:  $\mathcal{O}(1)$ .

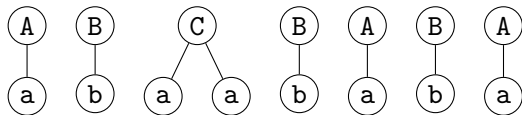
Mehlhorn et al. **Deterministic coin tossing** (by Cole & Vishkin):

- Blocks length:  $2-4$ ; context size:  $\mathcal{O}(\log^* n)$ .
- Reused by Alstrup et al.

Jež (**Recompression**) Alphabet partitioning:

- Partition the alphabet into **left** and **right** symbols.
- A block of length 2 when a right symbol follows a left symbol.
- Remaining characters in blocks of length 1.

# Parse Tree Construction: COMPRESS



RLE  
 $A = (a, 1), B = (b, 1), C = (a, 2)$

How to **consistently** partition a string without non-trivial runs?

Mehlhorn et al. Local minima wrt. a **random order**:

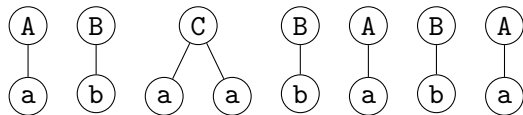
- Block length:  $2 - \mathcal{O}(\log n)$  w.h.p.,  $\mathcal{O}(1)$  in expectation.
- Context size:  $\mathcal{O}(1)$ .

Mehlhorn et al. **Deterministic coin tossing** (by Cole & Vishkin):

- Blocks length:  $2-4$ ; context size:  $\mathcal{O}(\log^* n)$ .
- Reused by Alstrup et al.

Jež (**Recompression**) Alphabet partitioning:

- Partition the alphabet into **left** and **right** symbols.
- A block of length 2 when a right symbol follows a left symbol.
- Remaining characters in blocks of length 1.
- Context size:  $\mathcal{O}(1)$ .



RLE  
 $A = (a, 1), B = (b, 1), C = (a, 2)$

How to **consistently** partition a string without non-trivial runs?

Mehlhorn et al. Local minima wrt. a **random order**:

- Block length:  $2 - \mathcal{O}(\log n)$  w.h.p.,  $\mathcal{O}(1)$  in expectation.
- Context size:  $\mathcal{O}(1)$ .

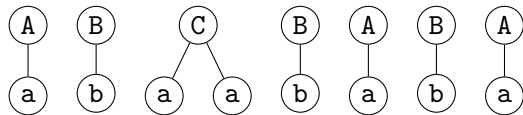
Mehlhorn et al. **Deterministic coin tossing** (by Cole & Vishkin):

- Blocks length:  $2-4$ ; context size:  $\mathcal{O}(\log^* n)$ .
- Reused by Alstrup et al.

Jež (**Recompression**) Alphabet partitioning:

- Partition the alphabet into **left** and **right** symbols.
- A block of length 2 when a right symbol follows a left symbol.
- Remaining characters in **blocks of length 1**.
- Context size:  $\mathcal{O}(1)$ .

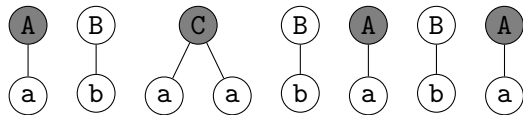
# Parse Tree Construction: Example



RLE  
 $A = (a, 1)$ ,  $B = (b, 1)$ ,  $C = (a, 2)$

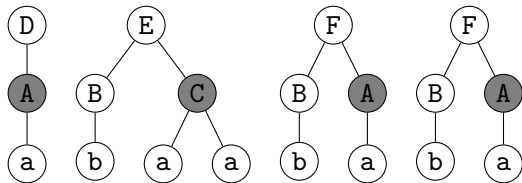


# Parse Tree Construction: Example



RLE  
 $A = (a, 1)$ ,  $B = (b, 1)$ ,  $C = (a, 2)$

# Parse Tree Construction: Example



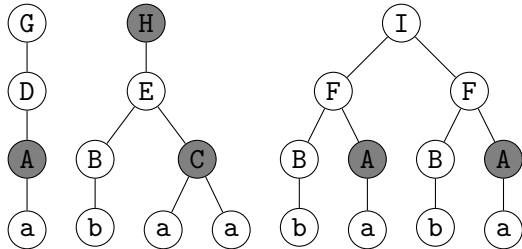
COMPRESS

$D = (A)$ ,  $E = (B, C)$ ,  $F = (B, A)$

RLE

$A = (a, 1)$ ,  $B = (b, 1)$ ,  $C = (a, 2)$

# Parse Tree Construction: Example



RLE

$G = (D, 1), H = (E, 1), I = (F, 2)$

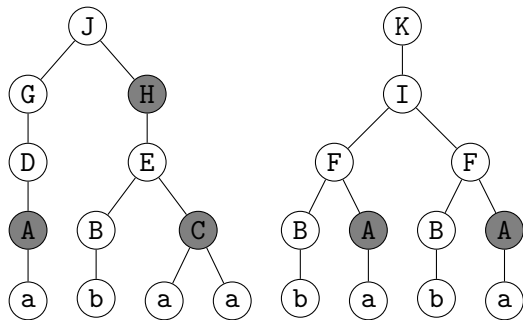
COMPRESS

$D = (A), E = (B, C), F = (B, A)$

RLE

$A = (a, 1), B = (b, 1), C = (a, 2)$

# Parse Tree Construction: Example



COMPRESS

$J = (G, H), K = (I)$

RLE

$G = (D, 1), H = (E, 1), I = (F, 2)$

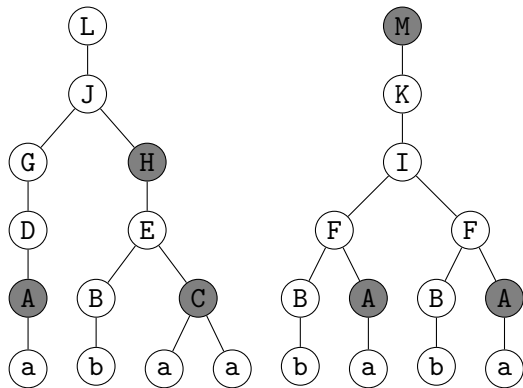
COMPRESS

$D = (A), E = (B, C), F = (B, A)$

RLE

$A = (a, 1), B = (b, 1), C = (a, 2)$

# Parse Tree Construction: Example



RLE

$L = (J, 1), M = (K, 1)$

COMPRESS

$J = (G, H), K = (I)$

RLE

$G = (D, 1), H = (E, 1), I = (F, 2)$

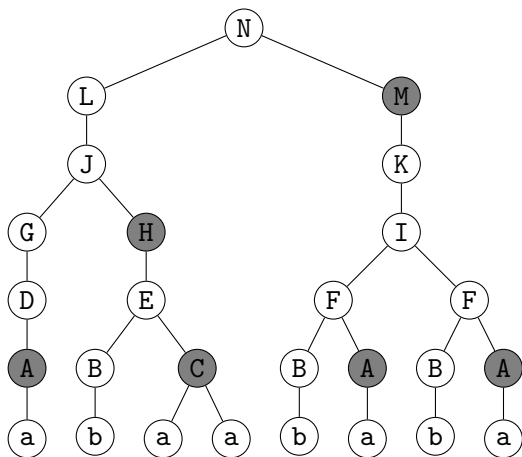
COMPRESS

$D = (A), E = (B, C), F = (B, A)$

RLE

$A = (a, 1), B = (b, 1), C = (a, 2)$

# Parse Tree Construction: Example



COMPRESS

$N = (L, M)$

RLE

$L = (J, 1), M = (K, 1)$

COMPRESS

$J = (G, H), K = (I)$

RLE

$G = (D, 1), H = (E, 1), I = (F, 2)$

COMPRESS

$D = (A), E = (B, C), F = (B, A)$

RLE

$A = (a, 1), B = (b, 1), C = (a, 2)$

How to partition the alphabet at each level?

# Parse Tree Construction: Alphabet Partitioning

How to partition the alphabet at each level?

**Jeż** To guarantee compression ratio  $r < 1$



# Parse Tree Construction: Alphabet Partitioning

How to partition the alphabet at each level?

Jeż To guarantee compression ratio  $r < 1$

- Works well in the static setting!

# Parse Tree Construction: Alphabet Partitioning

How to partition the alphabet at each level?

**Jež** To guarantee compression ratio  $r < 1$

- Works well in the static setting!
- Permanent decisions prone to be exploited by an adversary.

# Parse Tree Construction: Alphabet Partitioning

How to partition the alphabet at each level?

**Jeř** To guarantee compression ratio  $r < 1$

- Works well in the static setting!
- Permanent decisions prone to be exploited by an adversary.

**This work** Uniformly at random!

- Adjacent symbols form a block with probability  $\frac{1}{4}$ .

# Parse Tree Construction: Alphabet Partitioning

How to partition the alphabet at each level?

**Jež** To guarantee compression ratio  $r < 1$

- Works well in the static setting!
- Permanent decisions prone to be exploited by an adversary.

**This work** Uniformly at random!

- Adjacent symbols form a block with probability  $\frac{1}{4}$ .
- **Expected** compression ratio:  $\frac{3}{4}$ .

# Parse Tree Construction: Alphabet Partitioning

How to partition the alphabet at each level?

**Jeř** To guarantee compression ratio  $r < 1$

- Works well in the static setting!
- Permanent decisions prone to be exploited by an adversary.

**This work** Uniformly at random!

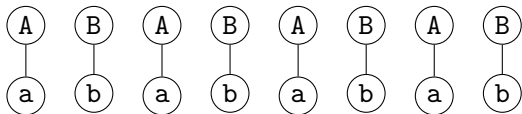
- Adjacent symbols form a block with probability  $\frac{1}{4}$ .
- **Expected** compression ratio:  $\frac{3}{4}$ .

## Lemma

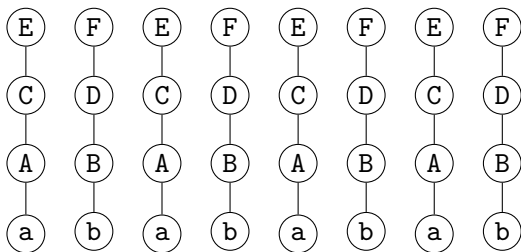
For any string  $w$  and any  $d \in \mathbb{R}_{\geq 0}$ :

$$\mathbb{P}[\text{DEPTH}(w) \leq 8(d + \ln |w|)] \geq 1 - e^{-d}$$

In short: depth  $\mathcal{O}(\log n)$  with high probability.



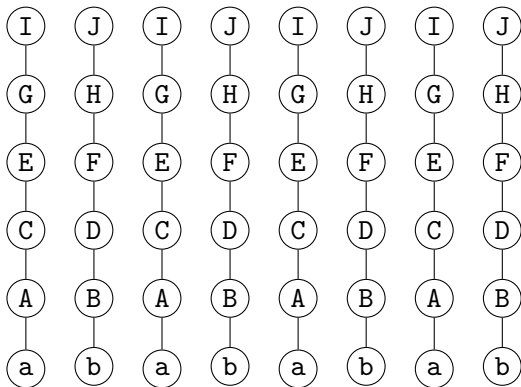
RLE



RLE

COMPRESS

RLE



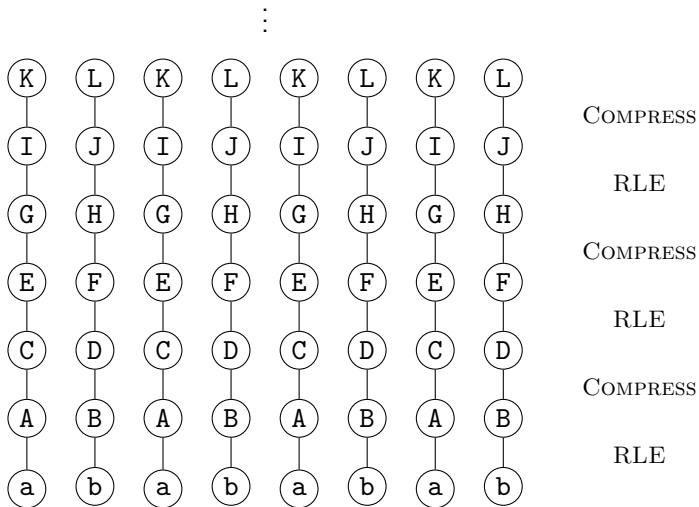
COMPRESS

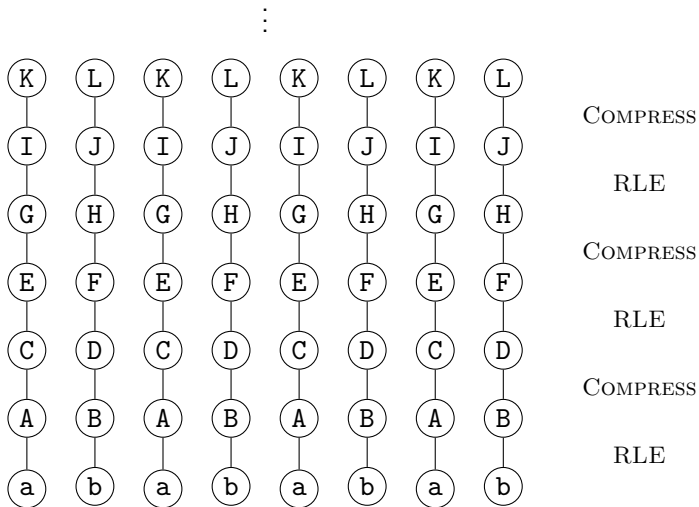
RLE

COMPRESS

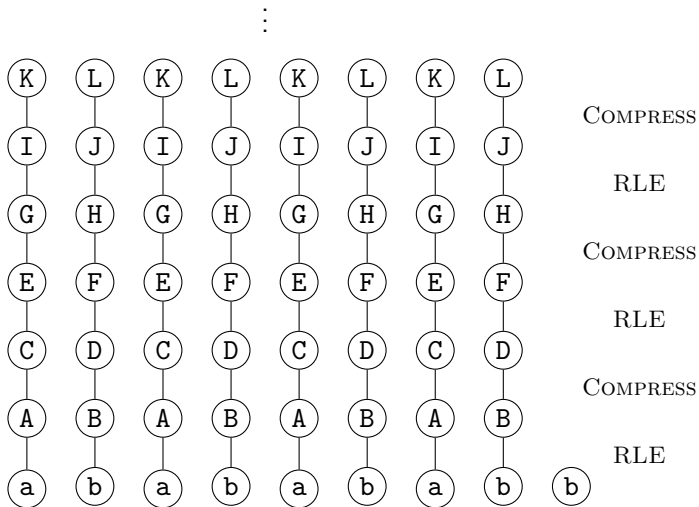
RLE



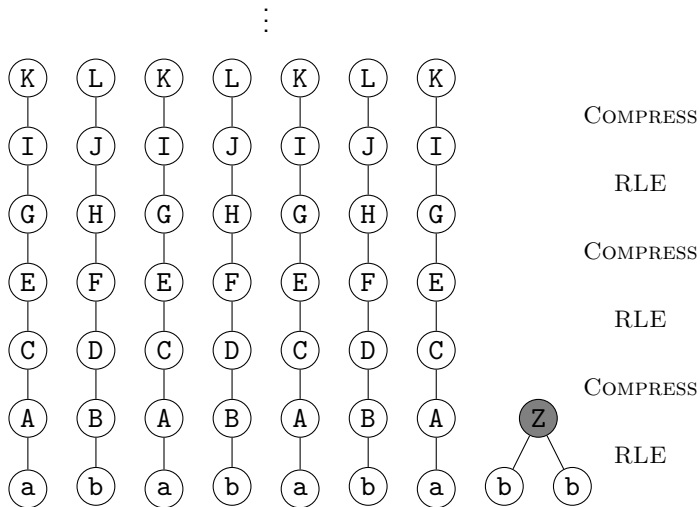




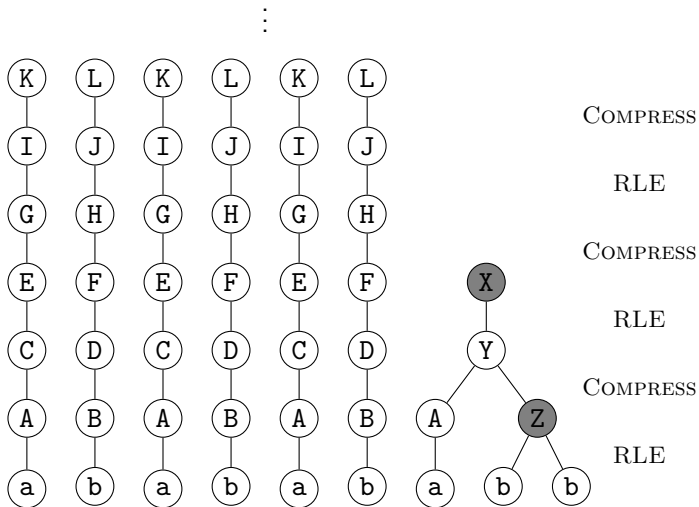
- $\Theta(n \log n)$  nodes with non-negligible probability  $\frac{1}{n^\epsilon}$ .



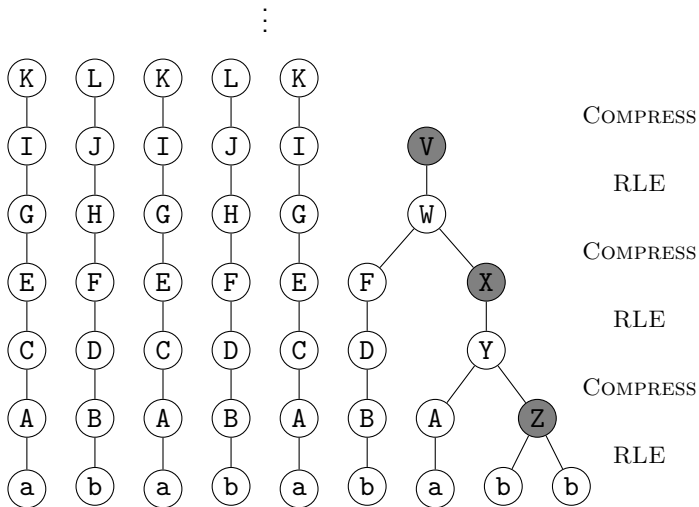
- $\Theta(n \log n)$  nodes with non-negligible probability  $\frac{1}{n^\epsilon}$ .



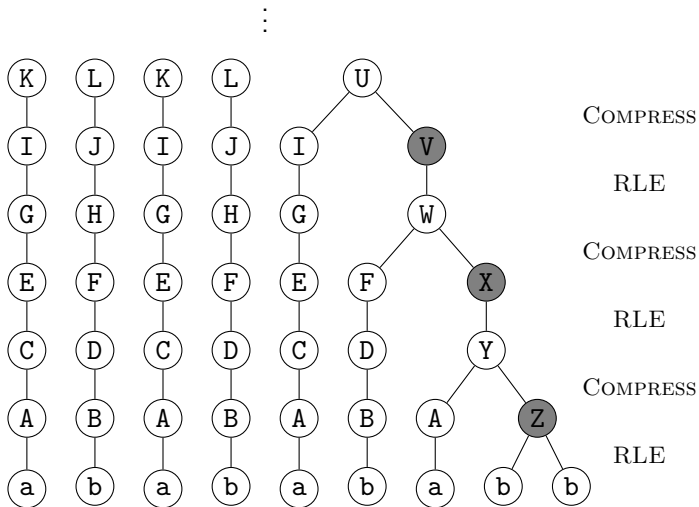
- $\Theta(n \log n)$  nodes with non-negligible probability  $\frac{1}{n^\epsilon}$ .



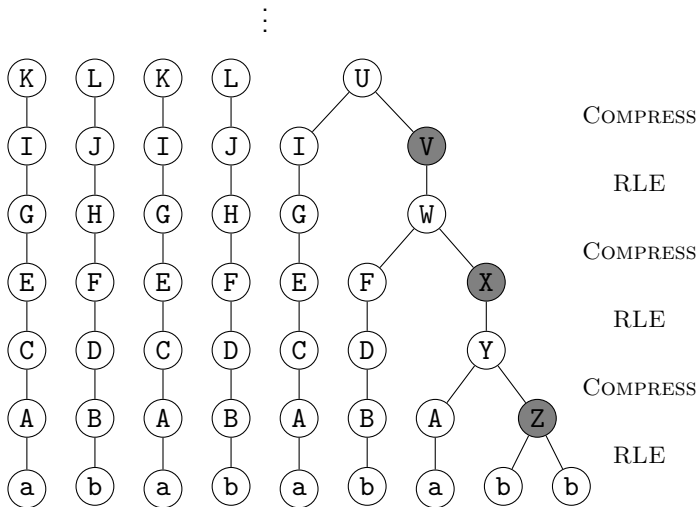
- $\Theta(n \log n)$  nodes with non-negligible probability  $\frac{1}{n^\epsilon}$ .



- $\Theta(n \log n)$  nodes with non-negligible probability  $\frac{1}{n^\epsilon}$ .



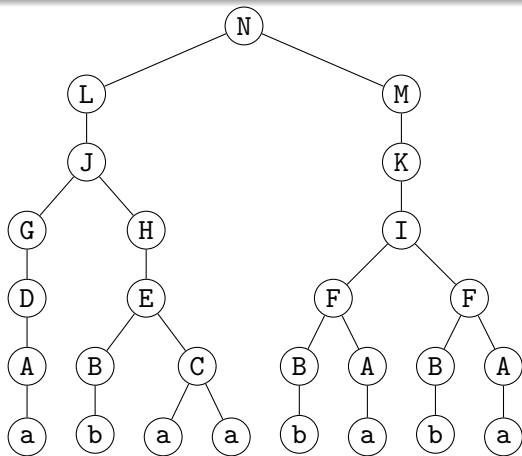
- $\Theta(n \log n)$  nodes with non-negligible probability  $\frac{1}{n^\epsilon}$ .



- $\Theta(n \log n)$  nodes with non-negligible probability  $\frac{1}{n^\epsilon}$ .
- Appending a letter might affect  $\Theta(\log^2 n)$  nodes.



# Workaround: Compress the Parse Tree



COMPRESS

$N = (L, M)$

RLE

$L = (J, 1), M = (K, 1)$

COMPRESS

$J = (G, H), K = (I)$

RLE

$G = (D, 1), H = (E, 1), I = (F, 2)$

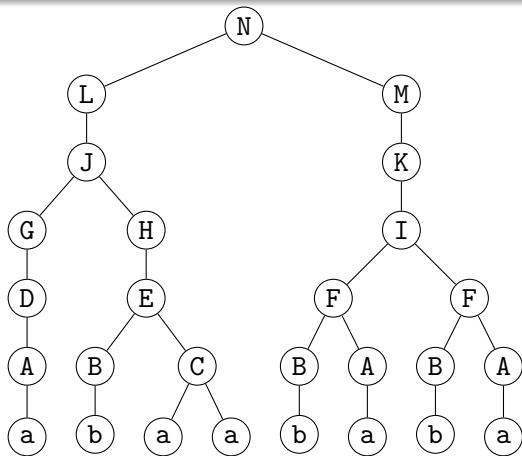
COMPRESS

$D = (A), E = (B, C), F = (B, A)$

RLE

$A = (a, 1), B = (b, 1), C = (a, 2)$

# Workaround: Compress the Parse Tree



COMPRESS

$N = (L, M)$

RLE

$L = (J, 1), M = (K, 1)$

COMPRESS

$J = (G, H), K = (I)$

RLE

$G = (D, 1), H = (E, 1), I = (F, 2)$

COMPRESS

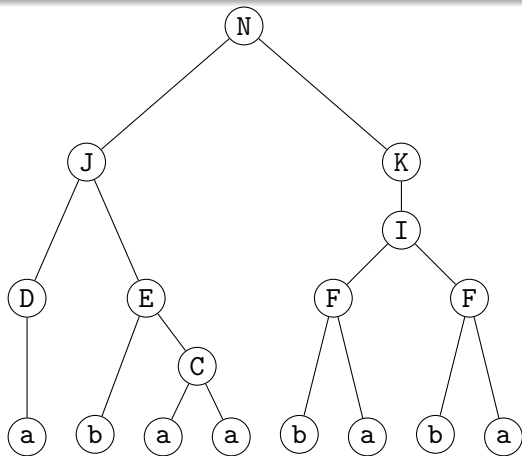
$D = (A), E = (B, C), F = (B, A)$

RLE

$A = (a, 1), B = (b, 1), C = (a, 2)$

- Do not replace  $S$  with  $(S, 1)$  in RLE.

# Workaround: Compress the Parse Tree



COMPRESS

$N = (J, K)$

RLE

COMPRESS

$J = (D, E), K = (I)$

RLE

$I = (F, 2)$

COMPRESS

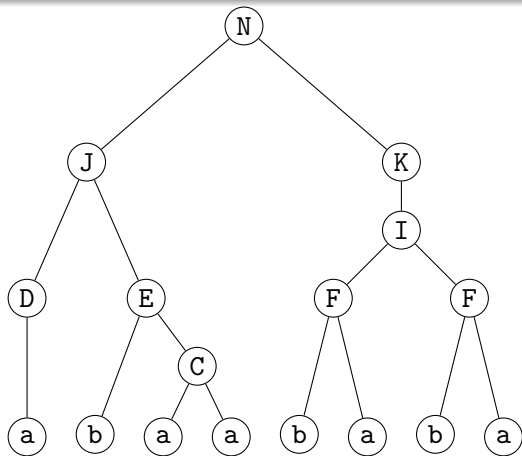
$D = (a), E = (b, C), F = (b, a)$

RLE

$C = (a, 2)$

- Do not replace  $S$  with  $(S, 1)$  in RLE.

# Workaround: Compress the Parse Tree



COMPRESS

$N = (J, K)$

RLE

COMPRESS

$J = (D, E), K = (I)$

RLE

$I = (F, 2)$

COMPRESS

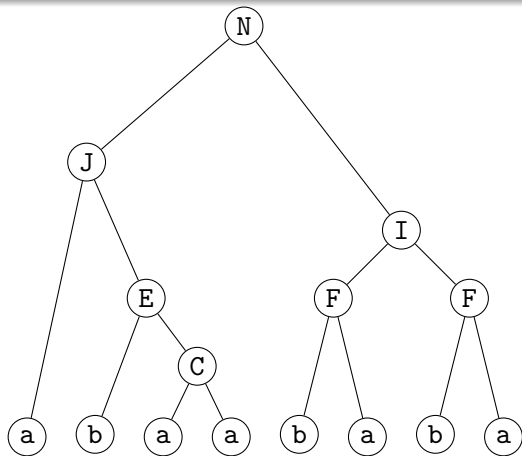
$D = (a), E = (b, C), F = (b, a)$

RLE

$C = (a, 2)$

- Do not replace  $S$  with  $(S, 1)$  in RLE.
- Do not introduce new symbols for unary blocks in COMPRESS.

# Workaround: Compress the Parse Tree



COMPRESS

$N = (J, I)$

RLE

COMPRESS

$J = (a, E)$

RLE

$I = (F, 2)$

COMPRESS

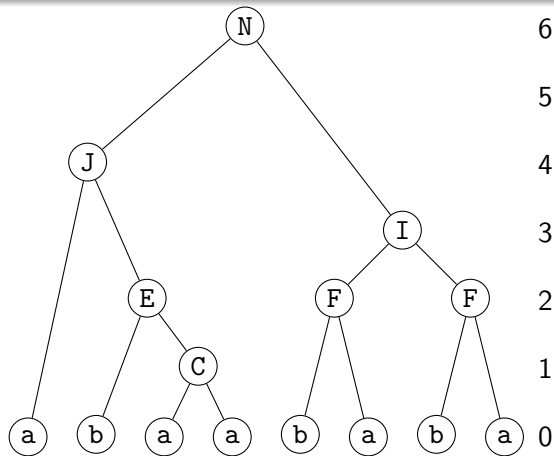
$E = (b, C), F = (b, a)$

RLE

$C = (a, 2)$

- Do not replace  $S$  with  $(S, 1)$  in RLE.
- Do not introduce new symbols for unary blocks in COMPRESS.

# Workaround: Compress the Parse Tree



6 COMPRESS  
 $N = (J, I)$   
RLE

5

4 COMPRESS  
 $J = (a, E)$   
RLE  
 $I = (F, 2)$

3

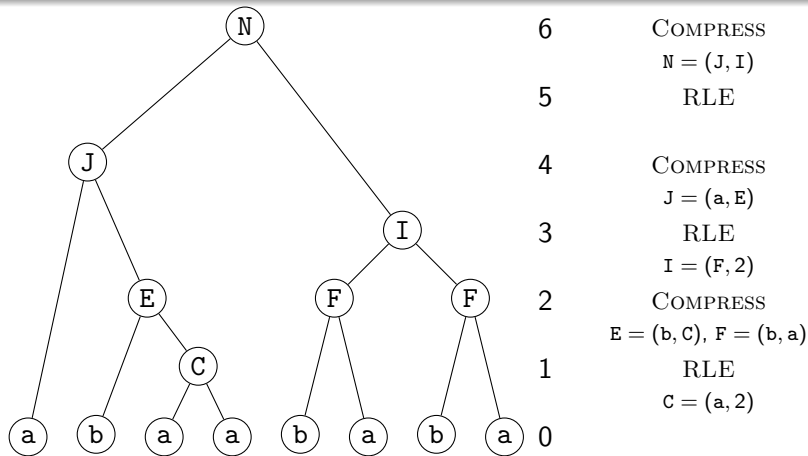
2 COMPRESS  
 $E = (b, C), F = (b, a)$

1 RLE  
 $C = (a, 2)$

0

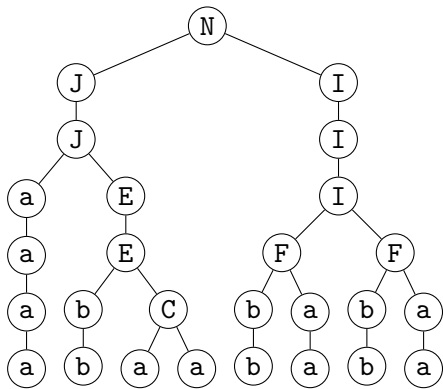
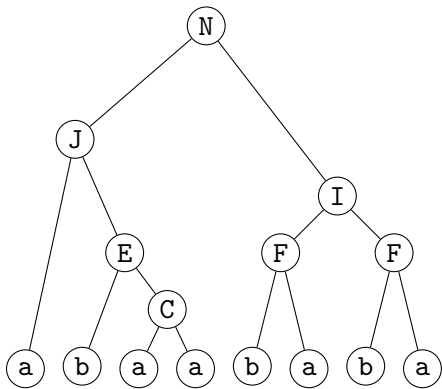
- Do not replace  $S$  with  $(S, 1)$  in RLE.
- Do not introduce new symbols for unary blocks in COMPRESS.
- Each symbol has a **level** where it appears ...

# Workaround: Compress the Parse Tree



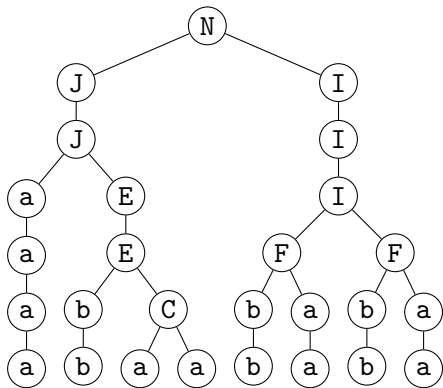
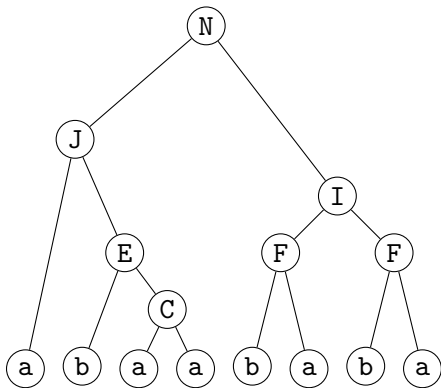
- Do not replace  $S$  with  $(S, 1)$  in RLE.
- Do not introduce new symbols for unary blocks in COMPRESS.
- Each symbol has a **level** where it appears ...
- and a random bit for each larger even level ( $\mathcal{O}(\log n)$  w.h.p.).

# Navigating Uncompressed Parse Trees





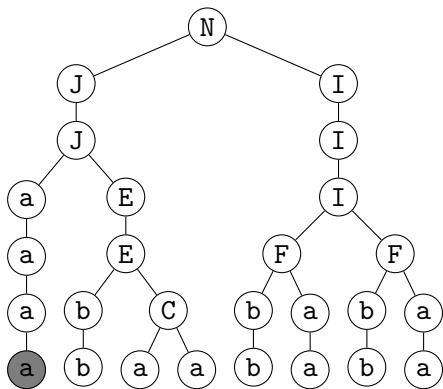
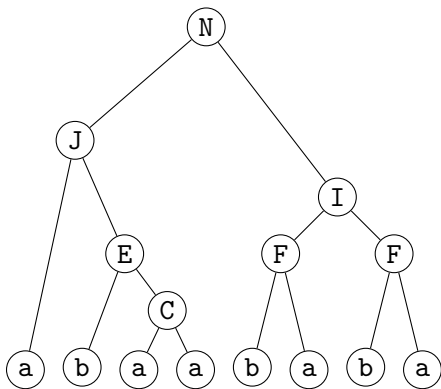
# Navigating Uncompressed Parse Trees



Navigate the **uncompressed parse trees** in  $\mathcal{O}(1)$  time:

- traverse edges: go to the parent or to the  $k$ -th child
- traverse levels: go left or right, perhaps skipping runs

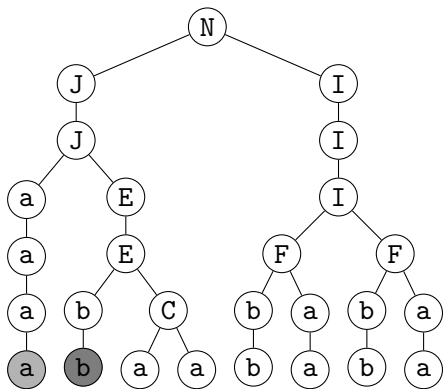
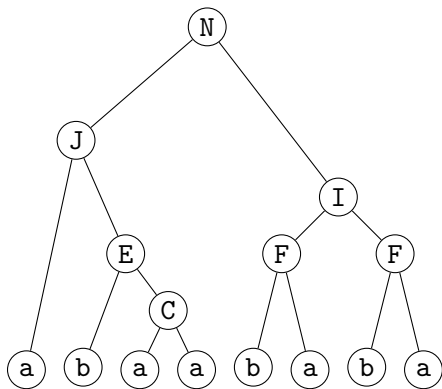
# Navigating Uncompressed Parse Trees



Navigate the **uncompressed parse trees** in  $\mathcal{O}(1)$  time:

- traverse edges: go to the parent or to the  $k$ -th child
- traverse levels: go left or right, perhaps skipping runs

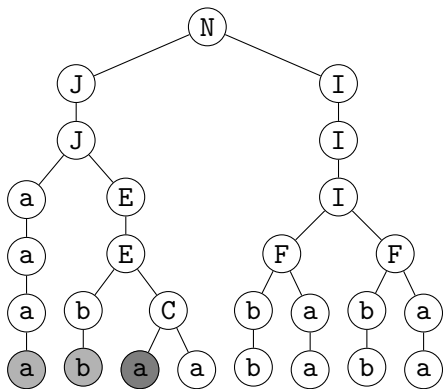
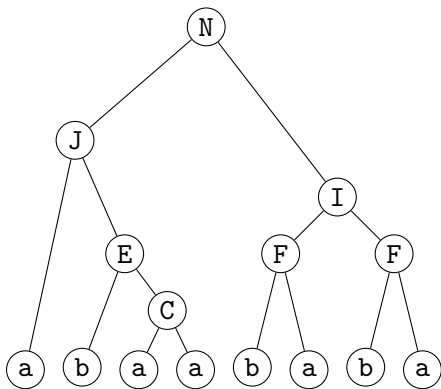
# Navigating Uncompressed Parse Trees



Navigate the **uncompressed parse trees** in  $\mathcal{O}(1)$  time:

- traverse edges: go to the parent or to the  $k$ -th child
- traverse levels: go left or right, perhaps skipping runs

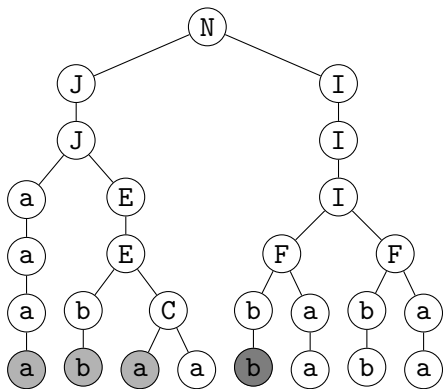
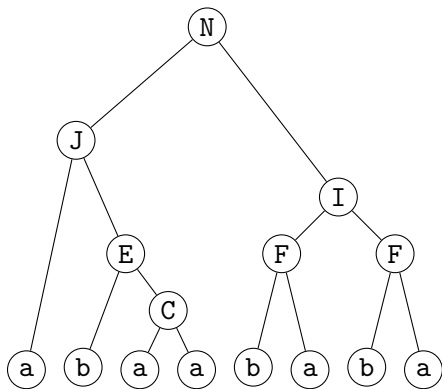
# Navigating Uncompressed Parse Trees



Navigate the **uncompressed parse trees** in  $\mathcal{O}(1)$  time:

- traverse edges: go to the parent or to the  $k$ -th child
- traverse levels: go left or right, perhaps skipping runs

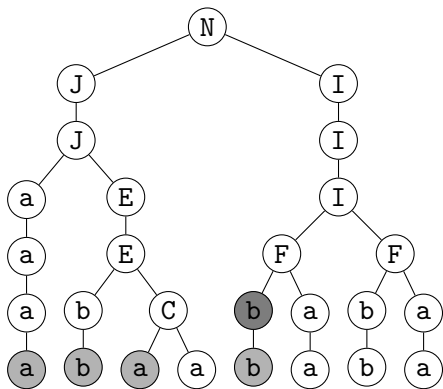
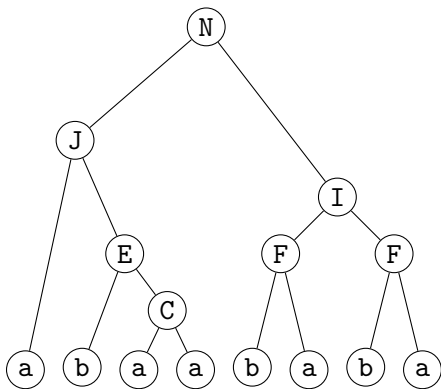
# Navigating Uncompressed Parse Trees



Navigate the **uncompressed parse trees** in  $\mathcal{O}(1)$  time:

- traverse edges: go to the parent or to the  $k$ -th child
- traverse levels: go left or right, perhaps skipping runs

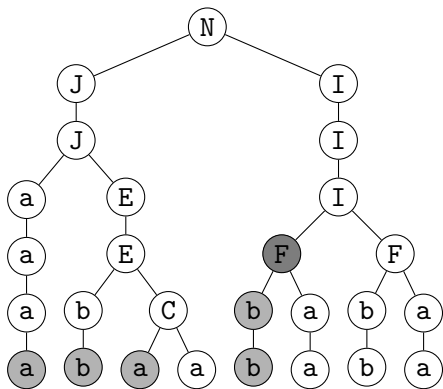
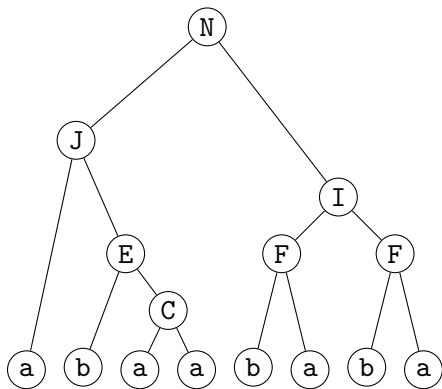
# Navigating Uncompressed Parse Trees



Navigate the **uncompressed parse trees** in  $\mathcal{O}(1)$  time:

- traverse edges: go to the parent or to the  $k$ -th child
- traverse levels: go left or right, perhaps skipping runs

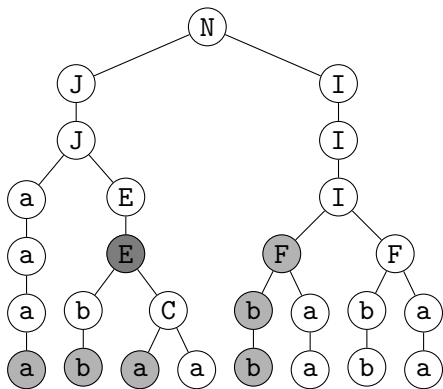
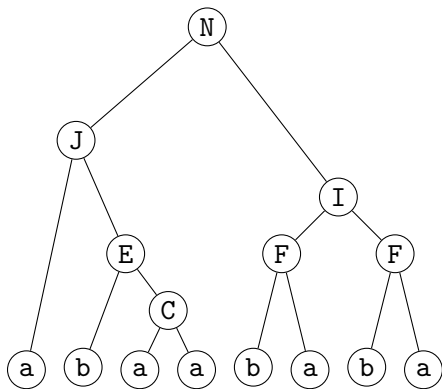
# Navigating Uncompressed Parse Trees



Navigate the **uncompressed parse trees** in  $\mathcal{O}(1)$  time:

- traverse edges: go to the parent or to the  $k$ -th child
- traverse levels: go left or right, perhaps skipping runs

# Navigating Uncompressed Parse Trees

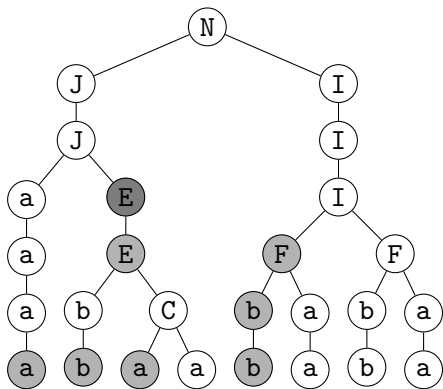
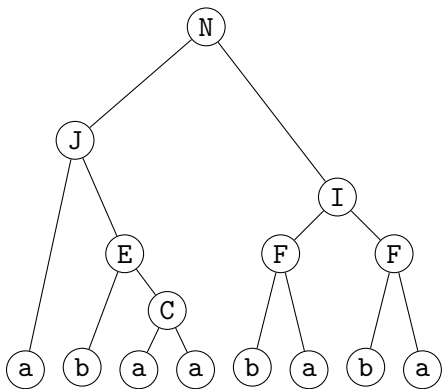


Navigate the **uncompressed parse trees** in  $\mathcal{O}(1)$  time:

- traverse edges: go to the parent or to the  $k$ -th child
- traverse levels: go left or right, perhaps skipping runs



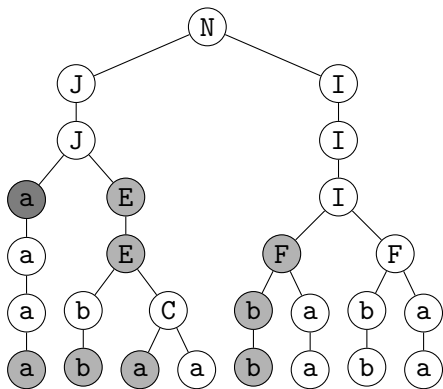
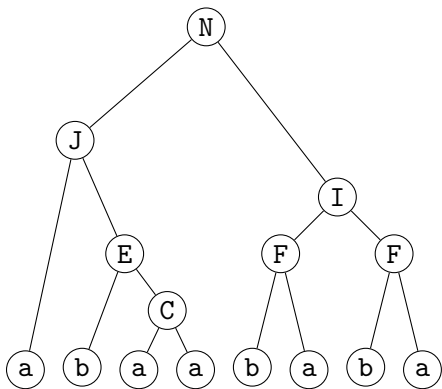
# Navigating Uncompressed Parse Trees



Navigate the **uncompressed parse trees** in  $\mathcal{O}(1)$  time:

- traverse edges: go to the parent or to the  $k$ -th child
- traverse levels: go left or right, perhaps skipping runs

# Navigating Uncompressed Parse Trees

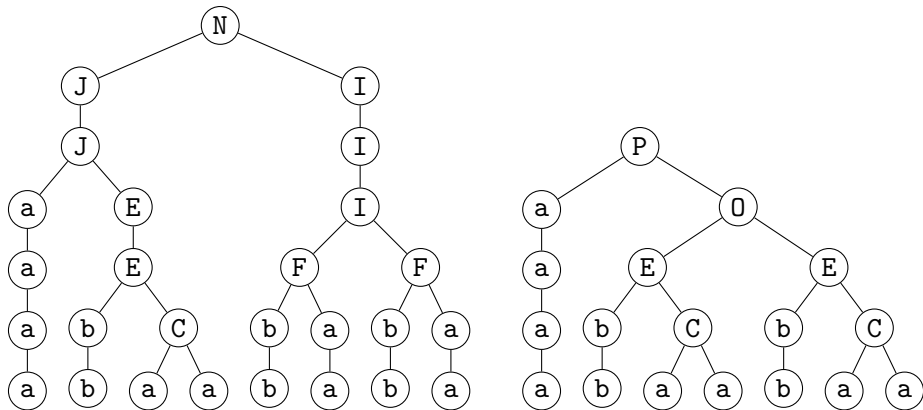


Navigate the **uncompressed parse trees** in  $\mathcal{O}(1)$  time:

- traverse edges: go to the parent or to the  $k$ -th child
- traverse levels: go left or right, perhaps skipping runs

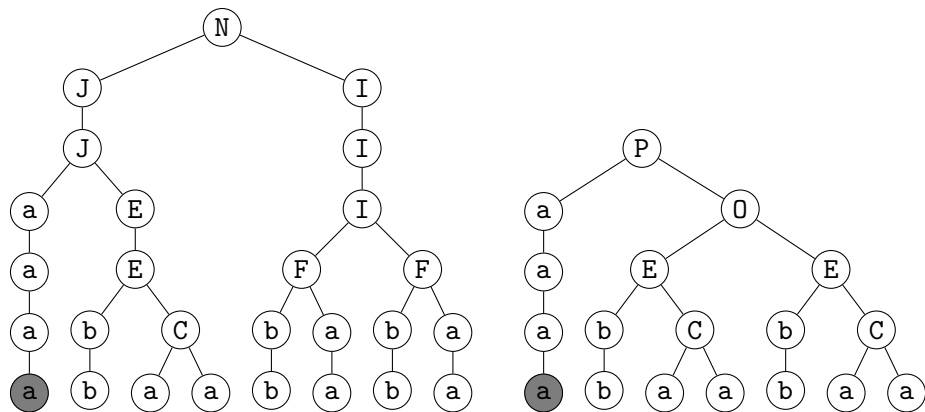


# Simple LCP and compare Implementation



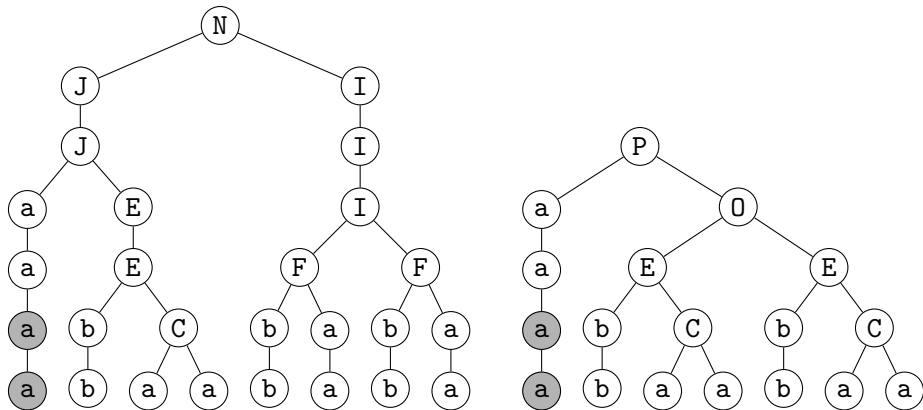
- 1 Start at the leftmost leaves and go up as far as possible.

# Simple LCP and compare Implementation



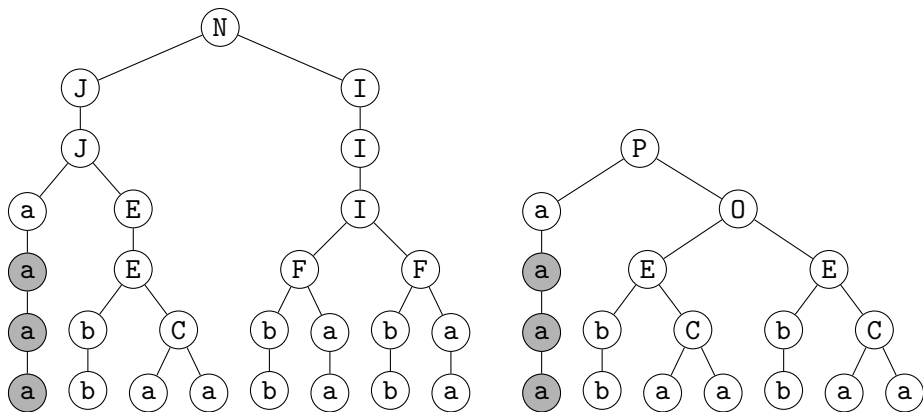
- 1 Start at the leftmost leaves and go up as far as possible.

# Simple LCP and compare Implementation



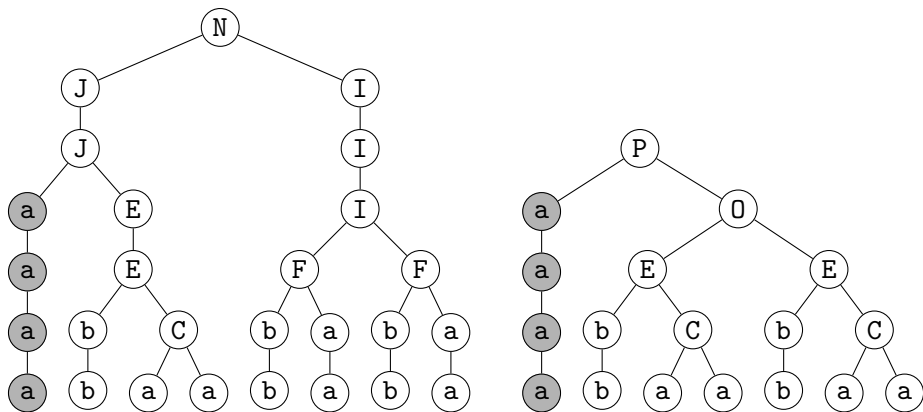
- 1 Start at the leftmost leaves and go up as far as possible.

# Simple LCP and compare Implementation



- 1 Start at the leftmost leaves and go up as far as possible.

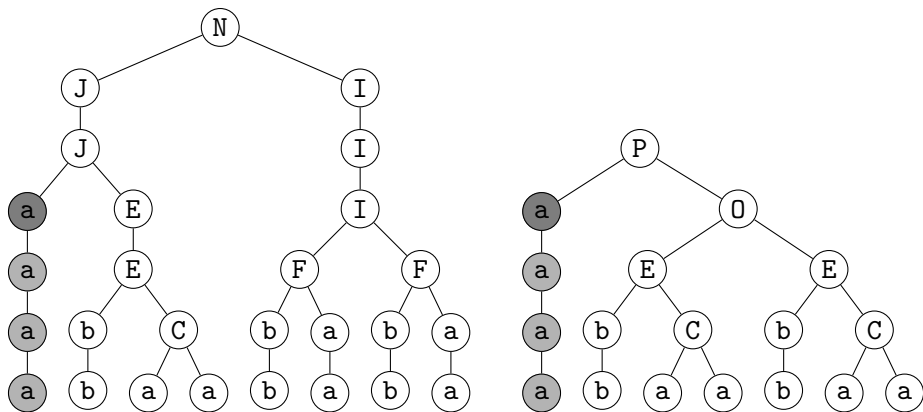
# Simple LCP and compare Implementation



- 1 Start at the leftmost leaves and go up as far as possible.

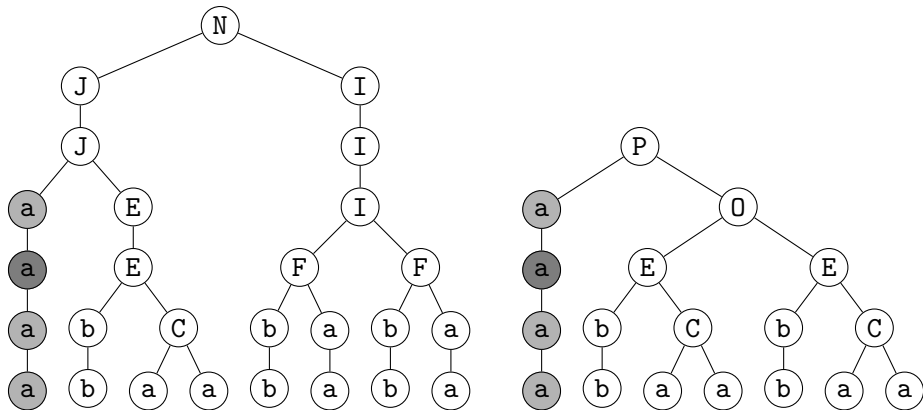


# Simple LCP and compare Implementation



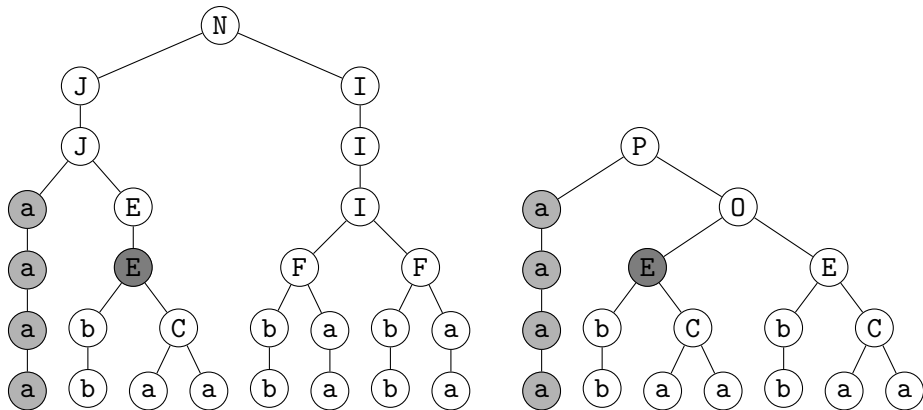
- 1 Start at the leftmost leaves and go up as far as possible.
- 2 Go right as far as possible, to the last child, and repeat.

# Simple LCP and compare Implementation



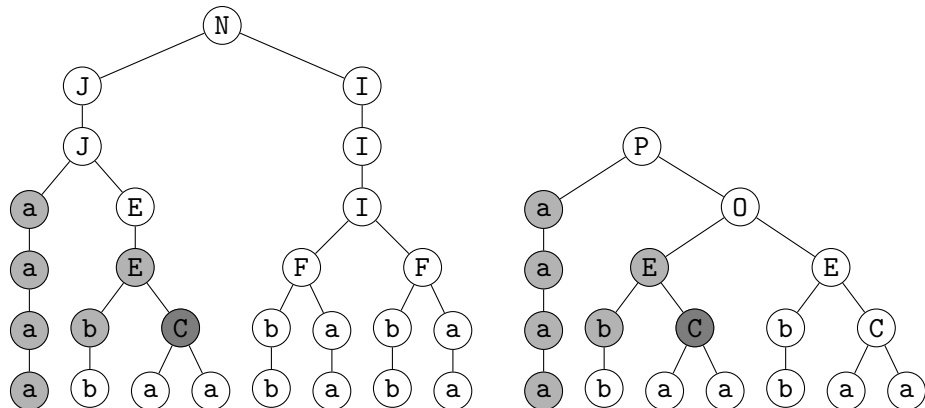
- 1 Start at the leftmost leaves and go up as far as possible.
- 2 Go right as far as possible, to the last child, and repeat.

# Simple LCP and compare Implementation



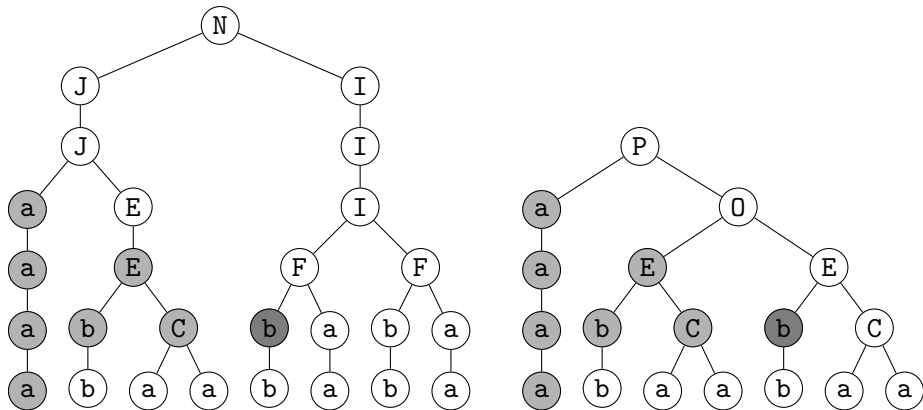
- 1 Start at the leftmost leaves and go up as far as possible.
- 2 Go right as far as possible, to the last child, and repeat.

# Simple LCP and compare Implementation



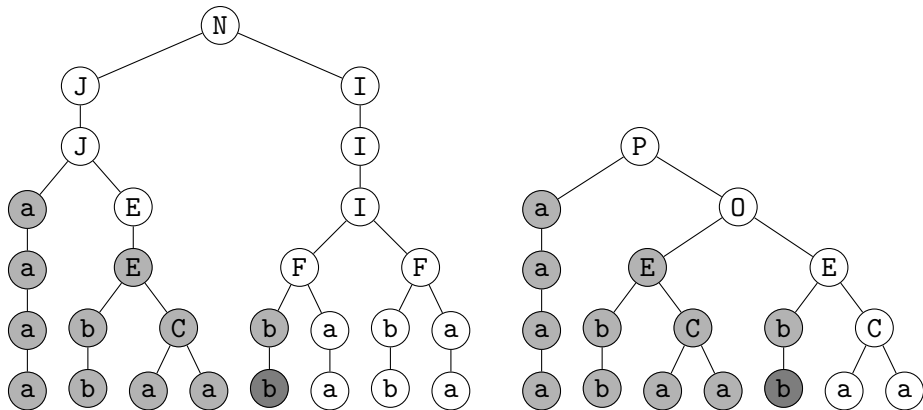
- 1 Start at the leftmost leaves and go up as far as possible.
- 2 Go right as far as possible, to the last child, and repeat.

# Simple LCP and compare Implementation



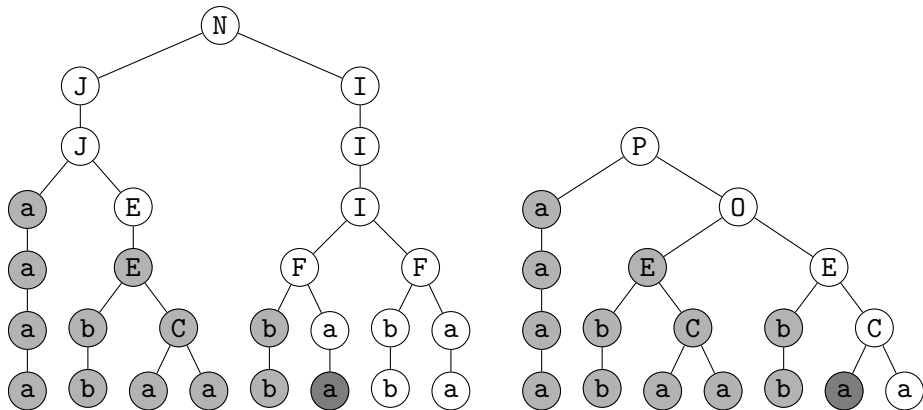
- 1 Start at the leftmost leaves and go up as far as possible.
- 2 Go right as far as possible, to the last child, and repeat.

# Simple LCP and compare Implementation



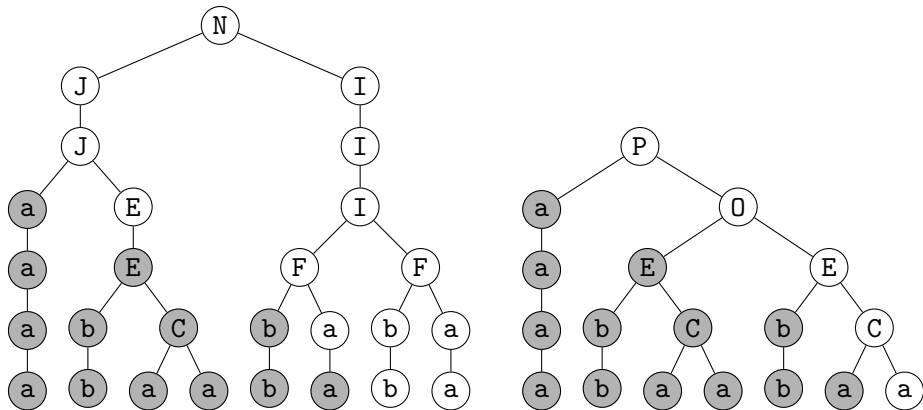
- 1 Start at the leftmost leaves and go up as far as possible.
- 2 Go right as far as possible, to the last child, and repeat.

# Simple LCP and compare Implementation



- 1 Start at the leftmost leaves and go up as far as possible.
- 2 Go right as far as possible, to the last child, and repeat.

# Simple LCP and compare Implementation

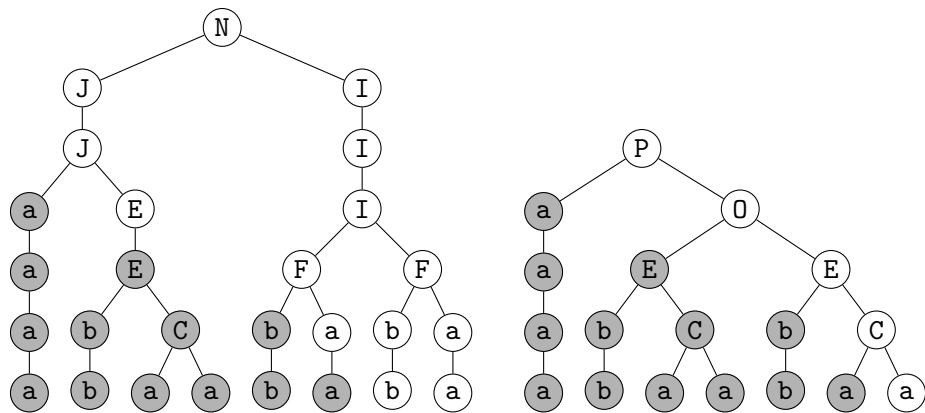


- 1 Start at the leftmost leaves and go up as far as possible.
- 2 Go right as far as possible, to the last child, and repeat.

**Fact.** While going to the right, at most run is traversed.



# Simple LCP and compare Implementation



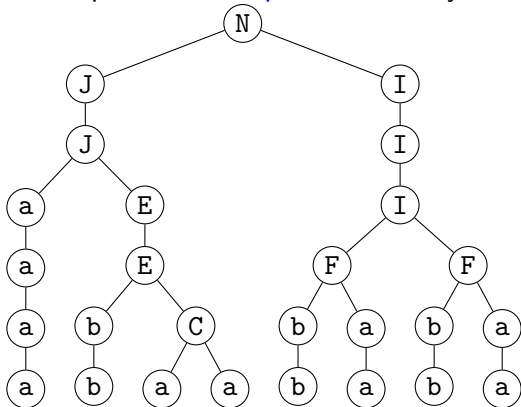
- 1 Start at the leftmost leaves and go up as far as possible.
- 2 Go right as far as possible, to the last child, and repeat.

**Fact.** While going to the right, at most run is traversed.

Total running time:  $\mathcal{O}(\min(\text{DEPTH}(w_1), \text{DEPTH}(w_2)))$ .

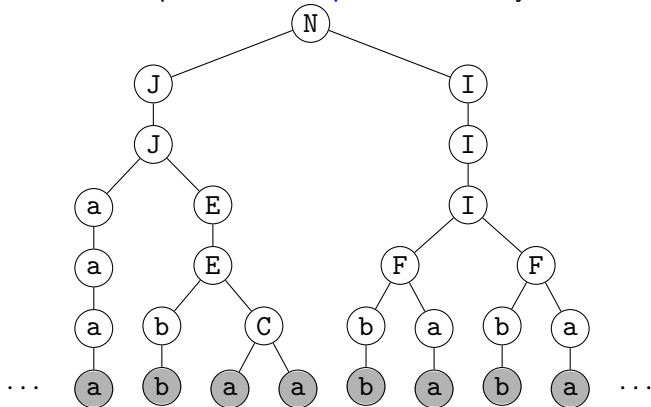
# Context-Insensitive Nodes

Some nodes in the parse tree are **preserved** in any **context**.



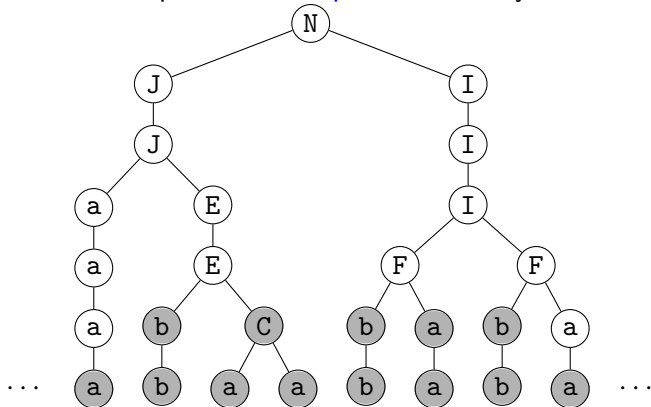
# Context-Insensitive Nodes

Some nodes in the parse tree are **preserved** in any **context**.



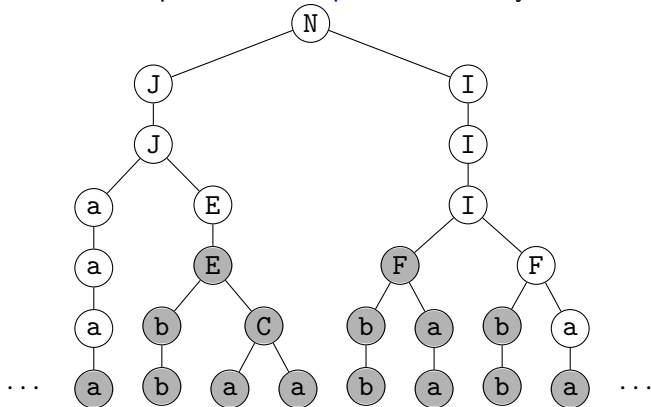
# Context-Insensitive Nodes

Some nodes in the parse tree are **preserved** in any **context**.



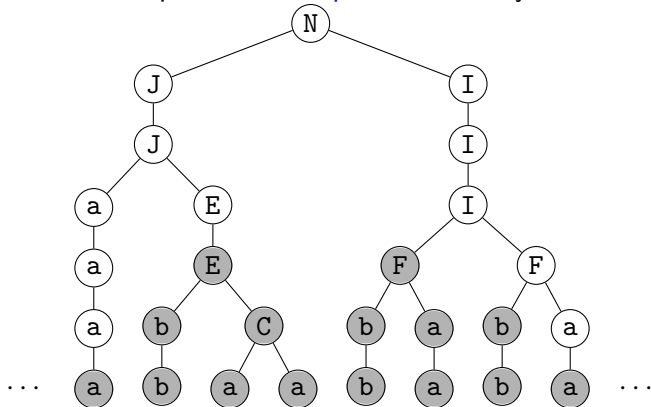
# Context-Insensitive Nodes

Some nodes in the parse tree are **preserved** in any **context**.



# Context-Insensitive Nodes

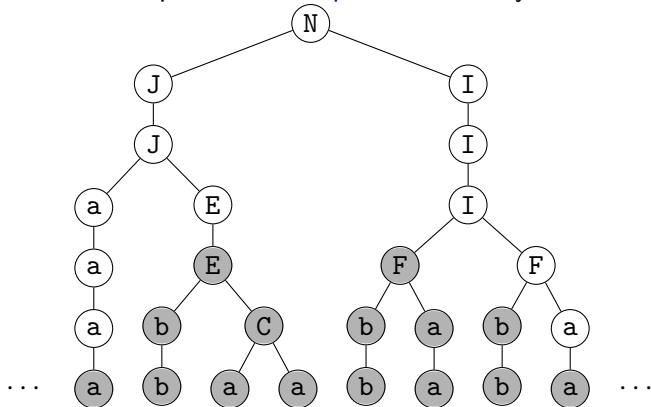
Some nodes in the parse tree are **preserved** in any **context**.



Context-insensitive decomposition  $D = aEFba$

# Context-Insensitive Nodes

Some nodes in the parse tree are **preserved** in any **context**.



Context-insensitive decomposition  $D = aEFba$

In general,  $|\text{RLE}(D)| \leq 2\text{DEPTH}(w)$ .

$\mathcal{O}(\text{DEPTH}(w))$  branching context-sensitive nodes.

## Lemma (Building over a decomposition)

*Given a run-length encoded decomposition  $D$  of  $w$ , we can add  $w$  to  $\mathcal{W}$  in  $\mathcal{O}(|\text{RLE}(D)| + \text{DEPTH}(w))$  time.*



## Lemma (Building over a decomposition)

*Given a run-length encoded decomposition  $D$  of  $w$ , we can add  $w$  to  $\mathcal{W}$  in  $\mathcal{O}(|\text{RLE}(D)| + \text{DEPTH}(w))$  time.*

- `make_string`: Build over the decomposition into letters.

## Lemma (Building over a decomposition)

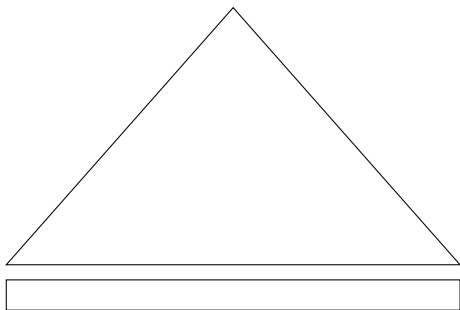
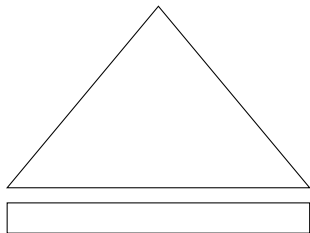
*Given a run-length encoded decomposition  $D$  of  $w$ , we can add  $w$  to  $\mathcal{W}$  in  $\mathcal{O}(|\text{RLE}(D)| + \text{DEPTH}(w))$  time.*

- `make_string`: Build over the decomposition into letters.
- `concat`, `split`: Extract the context-insensitive decompositions and build over them.

## Lemma (Building over a decomposition)

*Given a run-length encoded decomposition  $D$  of  $w$ , we can add  $w$  to  $\mathcal{W}$  in  $\mathcal{O}(|\text{RLE}(D)| + \text{DEPTH}(w))$  time.*

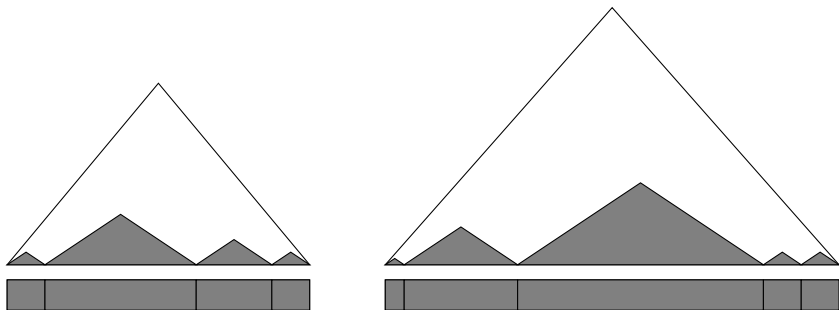
- `make_string`: Build over the decomposition into letters.
- `concat`, `split`: Extract the context-insensitive decompositions and build over them.



## Lemma (Building over a decomposition)

Given a run-length encoded decomposition  $D$  of  $w$ , we can add  $w$  to  $\mathcal{W}$  in  $\mathcal{O}(|\text{RLE}(D)| + \text{DEPTH}(w))$  time.

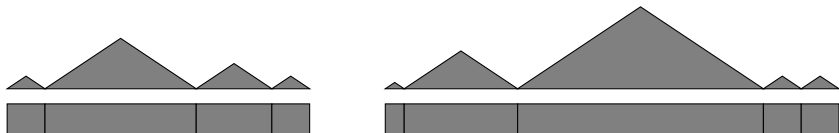
- `make_string`: Build over the decomposition into letters.
- `concat`, `split`: Extract the context-insensitive decompositions and build over them.



## Lemma (Building over a decomposition)

*Given a run-length encoded decomposition  $D$  of  $w$ , we can add  $w$  to  $\mathcal{W}$  in  $\mathcal{O}(|\text{RLE}(D)| + \text{DEPTH}(w))$  time.*

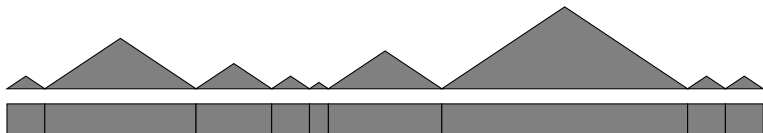
- `make_string`: Build over the decomposition into letters.
- `concat`, `split`: Extract the context-insensitive decompositions and build over them.



## Lemma (Building over a decomposition)

*Given a run-length encoded decomposition  $D$  of  $w$ , we can add  $w$  to  $\mathcal{W}$  in  $\mathcal{O}(|\text{RLE}(D)| + \text{DEPTH}(w))$  time.*

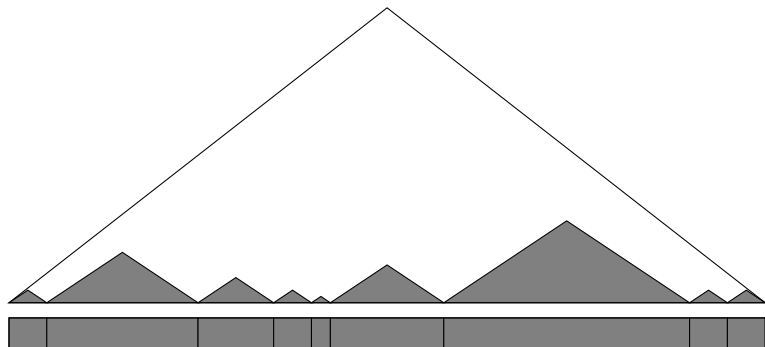
- `make_string`: Build over the decomposition into letters.
- `concat`, `split`: Extract the context-insensitive decompositions and build over them.



## Lemma (Building over a decomposition)

Given a run-length encoded decomposition  $D$  of  $w$ , we can add  $w$  to  $\mathcal{W}$  in  $\mathcal{O}(|\text{RLE}(D)| + \text{DEPTH}(w))$  time.

- `make_string`: Build over the decomposition into letters.
- `concat`, `split`: Extract the context-insensitive decompositions and build over them.



- 1 Indexing dynamic strings:
  - find all occurrences of a given pattern in the strings  $w \in \mathcal{W}$ ;



## ① Indexing dynamic strings:

- find all occurrences of a given pattern in the strings  $w \in \mathcal{W}$ ;
- improves upon Alstrup et al;
- described in the full ArXiv version.

- 1 Indexing dynamic strings:
  - find all occurrences of a given pattern in the strings  $w \in \mathcal{W}$ ;
  - improves upon Alstrup et al;
  - described in the full ArXiv version.
- 2 Introduce `drop_string(w)` to remove  $w \in \mathcal{W}$ :
  - garbage collection of unused symbols;
  - $\mathcal{O}(z \log n)$  symbols remain w.h.p., where  $z$  is the minimum number of operations to generate  $\mathcal{W}$ ;

- 1 Indexing dynamic strings:
  - find all occurrences of a given pattern in the strings  $w \in \mathcal{W}$ ;
  - improves upon Alstrup et al;
  - described in the full ArXiv version.
- 2 Introduce `drop_string(w)` to remove  $w \in \mathcal{W}$ :
  - garbage collection of unused symbols;
  - $\mathcal{O}(z \log n)$  symbols remain w.h.p., where  $z$  is the minimum number of operations to generate  $\mathcal{W}$ ;
  - tweaking hash tables needed to guarantee w.h.p. running time in  $\mathcal{O}(z \log n)$  space.

- 1 Indexing dynamic strings:
  - find all occurrences of a given pattern in the strings  $w \in \mathcal{W}$ ;
  - improves upon Alstrup et al;
  - described in the full ArXiv version.
- 2 Introduce `drop_string(w)` to remove  $w \in \mathcal{W}$ :
  - garbage collection of unused symbols;
  - $\mathcal{O}(z \log n)$  symbols remain w.h.p., where  $z$  is the minimum number of operations to generate  $\mathcal{W}$ ;
  - tweaking hash tables needed to guarantee w.h.p. running time in  $\mathcal{O}(z \log n)$  space.
- 3 Further  $\mathcal{O}(\log n)$ -time operations:

- 1 Indexing dynamic strings:
  - find all occurrences of a given pattern in the strings  $w \in \mathcal{W}$ ;
  - improves upon Alstrup et al;
  - described in the full ArXiv version.
- 2 Introduce `drop_string(w)` to remove  $w \in \mathcal{W}$ :
  - garbage collection of unused symbols;
  - $\mathcal{O}(z \log n)$  symbols remain w.h.p., where  $z$  is the minimum number of operations to generate  $\mathcal{W}$ ;
  - tweaking hash tables needed to guarantee w.h.p. running time in  $\mathcal{O}(z \log n)$  space.
- 3 Further  $\mathcal{O}(\log n)$ -time operations:
  - `power(w, k)`: insert the power  $w^k$  to  $\mathcal{W}$ ;

- 1 Indexing dynamic strings:
  - find all occurrences of a given pattern in the strings  $w \in \mathcal{W}$ ;
  - improves upon Alstrup et al;
  - described in the full ArXiv version.
- 2 Introduce `drop_string(w)` to remove  $w \in \mathcal{W}$ :
  - garbage collection of unused symbols;
  - $\mathcal{O}(z \log n)$  symbols remain w.h.p., where  $z$  is the minimum number of operations to generate  $\mathcal{W}$ ;
  - tweaking hash tables needed to guarantee w.h.p. running time in  $\mathcal{O}(z \log n)$  space.
- 3 Further  $\mathcal{O}(\log n)$ -time operations:
  - `power(w, k)`: insert the power  $w^k$  to  $\mathcal{W}$ ;
  - `reverse(w)`: insert the reverse  $w^R$  to  $\mathcal{W}$ ;

- 1 Indexing dynamic strings:
  - find all occurrences of a given pattern in the strings  $w \in \mathcal{W}$ ;
  - improves upon Alstrup et al;
  - described in the full ArXiv version.
- 2 Introduce `drop_string(w)` to remove  $w \in \mathcal{W}$ :
  - garbage collection of unused symbols;
  - $\mathcal{O}(z \log n)$  symbols remain w.h.p., where  $z$  is the minimum number of operations to generate  $\mathcal{W}$ ;
  - tweaking hash tables needed to guarantee w.h.p. running time in  $\mathcal{O}(z \log n)$  space.
- 3 Further  $\mathcal{O}(\log n)$ -time operations:
  - `power(w, k)`: insert the power  $w^k$  to  $\mathcal{W}$ ;
  - `reverse(w)`: insert the reverse  $w^R$  to  $\mathcal{W}$ ;
  - `cyclic(w)`: compute the **primitive root** and canonize it with respect to **cyclic equivalence**.

Data Structure	Ours	Impossible
make_string	$\mathcal{O}( w  + \log n)$	$o( w  \log n)$
concat	$\mathcal{O}(\log n)$	$o(\log n)$
split	$\mathcal{O}(\log n)$	$o(\log n)$
equal	$\mathcal{O}(1)$	$o(\log n)$
compare	$\mathcal{O}(1)$	-
LCP	$\mathcal{O}(1)$	-
persistent	YES	NO
	Las Vegas	Monte Carlo
	worst-case	amortized



Data Structure	Ours	Impossible
make_string	$\mathcal{O}( w  + \log n)$	$o( w  \log n)$
concat	$\mathcal{O}(\log n)$	$o(\log n)$
split	$\mathcal{O}(\log n)$	$o(\log n)$
equal	$\mathcal{O}(1)$	$o(\log n)$
compare	$\mathcal{O}(1)$	-
LCP	$\mathcal{O}(1)$	-
persistent	YES	NO
	Las Vegas	Monte Carlo
	worst-case	amortized

Thank you for your attention!