

Introduction to Unification Theory

Matching

Temur Kutsia

RISC, Johannes Kepler University of Linz, Austria
`kutsia@risc.uni-linz.ac.at`

Matching

- ▶ A special case of unification.
- ▶ An example of a matching equation:

$$f(f(a, X), Y) \doteq f(f(a, b), f(f(a, b), a)).$$

- ▶ $f(f(a, b), f(f(a, b), a))$ is ground.
- ▶ Any matching equation has at least one ground side.
- ▶ Solution (matcher): $\{X \mapsto b, Y \mapsto f(f(a, b), a)\}$.



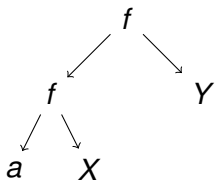
Matching

- ▶ Matching is needed in rewriting, functional programming, querying, etc.
- ▶ Often the following problem is required to be solved:
 - ▶ Given a ground term s (subject) and a term p (pattern)
 - ▶ Find all subterms in s to which p matches.
- ▶ In this lecture: An algorithm to solve this problem.
- ▶ Terms are represented as trees.

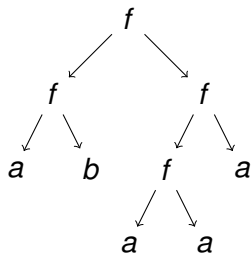


Tree Pattern Matching

Matching the pattern tree to the subject tree.



Pattern tree 1

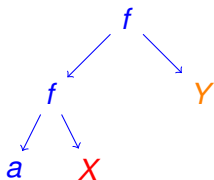


Subject tree

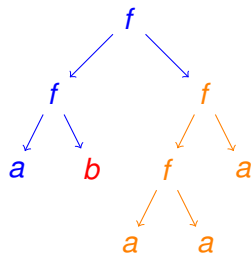
Tree Pattern Matching

Matching the pattern tree to the subject tree.

Pattern tree 1. First match:



Pattern tree 1

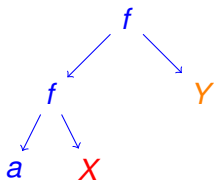


Subject tree

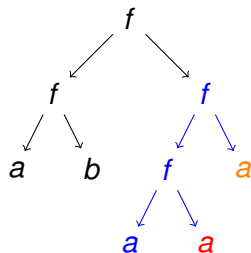
Tree Pattern Matching

Matching the pattern tree to the subject tree.

Pattern tree 1. Second match:



Pattern tree 1

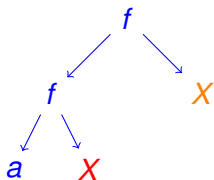


Subject tree

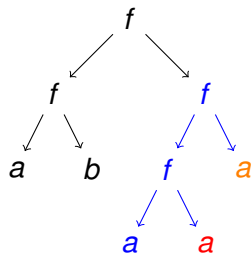
Tree Pattern Matching

Matching the pattern tree to the subject tree.

Pattern tree 2. Single match:



Pattern tree 2



Subject tree

Tree Pattern Matching

- ▶ Pattern tree 1 in the example is **linear**: Every variable occurs only once.
- ▶ Pattern tree 2 is **nonlinear**: X occurs twice.
- ▶ Two steps for nonlinear tree matching:
 1. Ignore multiplicity of variables (assume the pattern is linear) and do **linear tree pattern matching**.
 2. Verify that the substitutions computed for multiple occurrences of a variable are identical: **check consistency**.



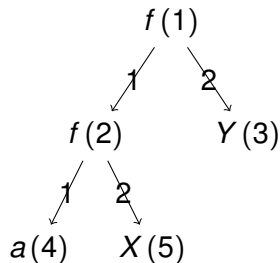
Terms

- ▶ \mathcal{V} : Set of variables.
- ▶ \mathcal{F} : Set of function symbols of fixed arity.
- ▶ $\mathcal{F} \cap \mathcal{V} = \emptyset$.
- ▶ Constants: 0-ary function symbols.
- ▶ **Terms:**
 - ▶ A variable or a constant is a term.
 - ▶ If $f \in \mathcal{F}$, f is n -ary, $n > 0$, and t_1, \dots, t_n are terms, then $f(t_1, \dots, t_n)$ is a term.

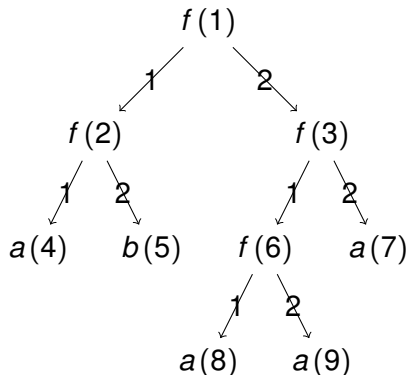


Term Trees, Nodes, Node Labels, Edges, Edge labels

Example



The tree for $f(f(a, X), Y)$

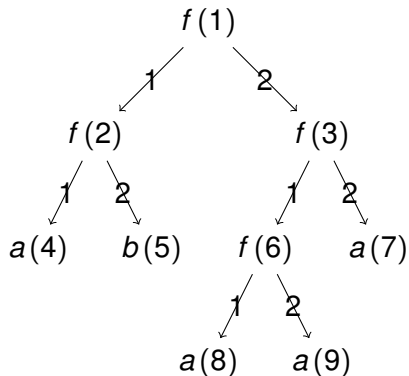
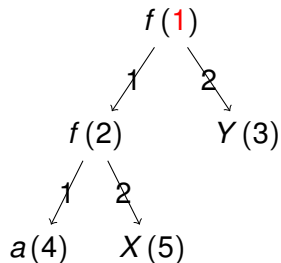


The tree for $f(f(a, b), f(f(a, a), a))$



Term Trees, Nodes, Node Labels, Edges, Edge labels

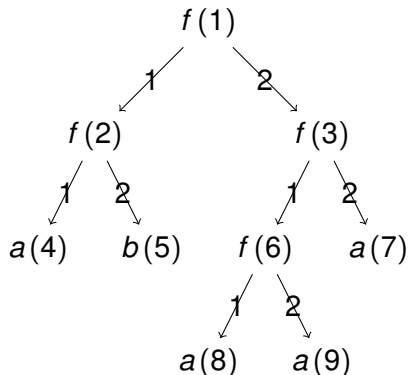
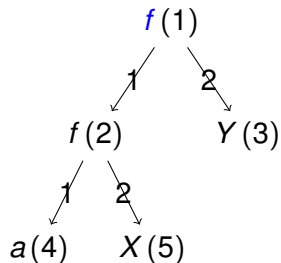
Example



Node

Term Trees, Nodes, Node Labels, Edges, Edge labels

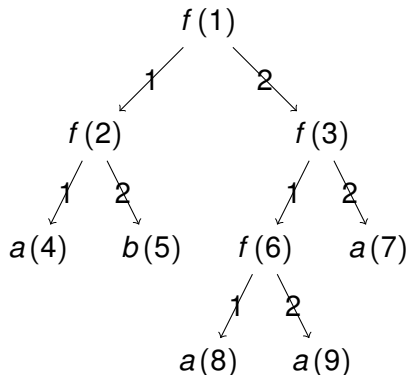
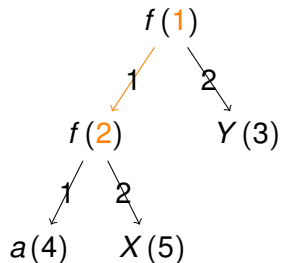
Example



Node label

Term Trees, Nodes, Node Labels, Edges, Edge labels

Example

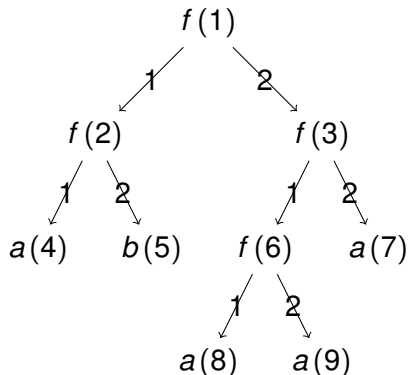
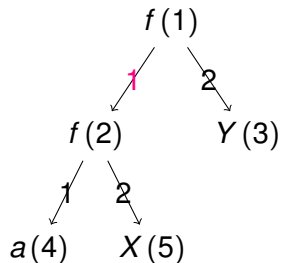


Edge



Term Trees, Nodes, Node Labels, Edges, Edge labels

Example



Edge label



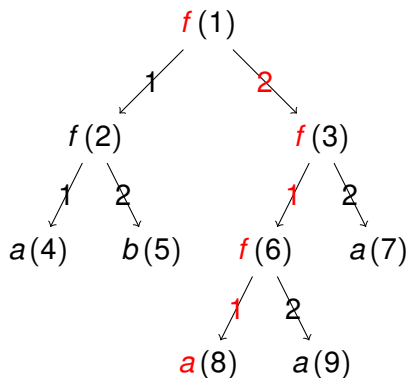
Labeled Path

- ▶ **Labeled path** $lp(n_1, n_q)$ in a term tree from the node n_1 to the node n_q :
A string formed by alternatively concatenating the node and edge labels from n_1 to n_q .



Labeled Path

Example

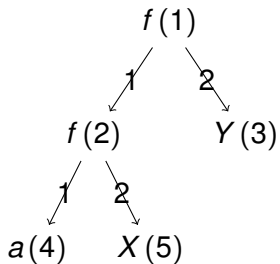


Labeled path from 1 to 8: $lp(1, 8) = f2f1f1a$



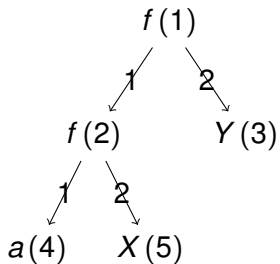
Euler Chains and Strings

- ▶ **Euler chain** for a term tree: a string of node labels obtained as follows:



Euler Chains and Strings

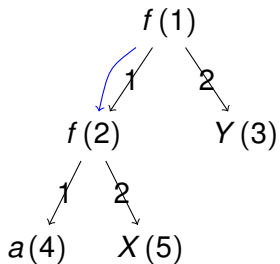
- ▶ **Euler chain** for a term tree: a string of node labels obtained as follows:



1

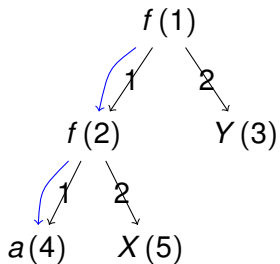
Euler Chains and Strings

- ▶ **Euler chain** for a term tree: a string of node labels obtained as follows:



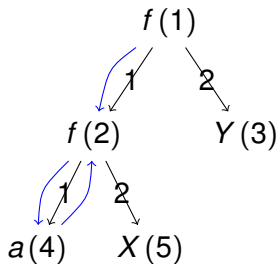
Euler Chains and Strings

- **Euler chain** for a term tree: a string of node labels obtained as follows:



Euler Chains and Strings

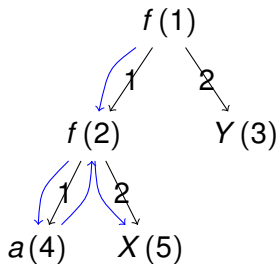
- **Euler chain** for a term tree: a string of node labels obtained as follows:



1242

Euler Chains and Strings

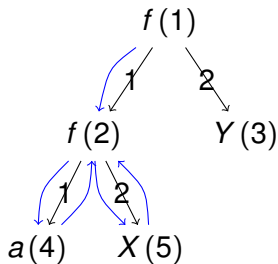
- **Euler chain** for a term tree: a string of node labels obtained as follows:



12425

Euler Chains and Strings

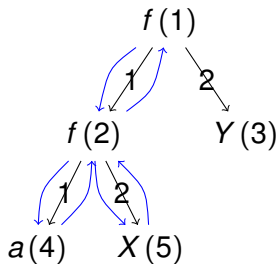
- ▶ **Euler chain** for a term tree: a string of node labels obtained as follows:



124252

Euler Chains and Strings

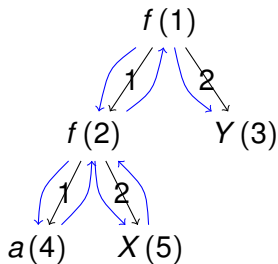
- **Euler chain** for a term tree: a string of node labels obtained as follows:



1242521

Euler Chains and Strings

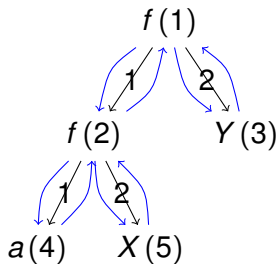
- **Euler chain** for a term tree: a string of node labels obtained as follows:



12425213

Euler Chains and Strings

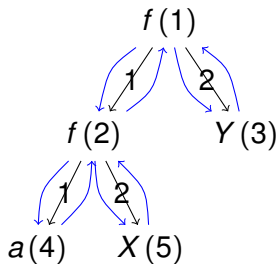
- **Euler chain** for a term tree: a string of node labels obtained as follows:



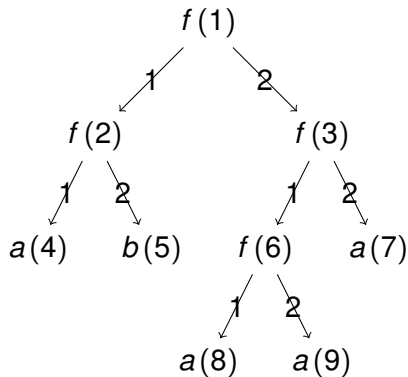
124252131

Euler Chains and Strings

- **Euler chain** for a term tree: a string of node labels obtained as follows:



124252131

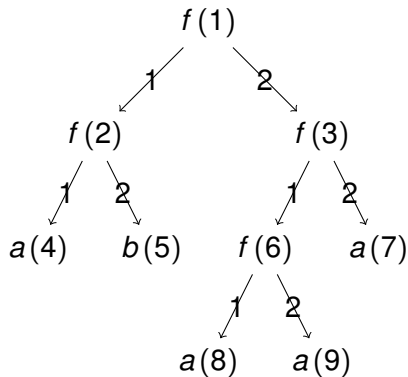
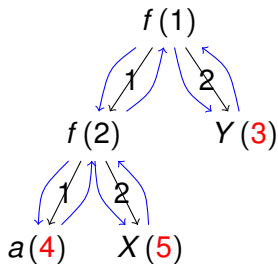


12425213686963731



Euler Chains and Strings

- Properties of Euler chains



The leaves occur only once:

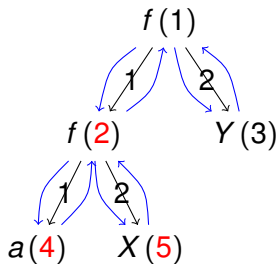
124252131

12425213686963731



Euler Chains and Strings

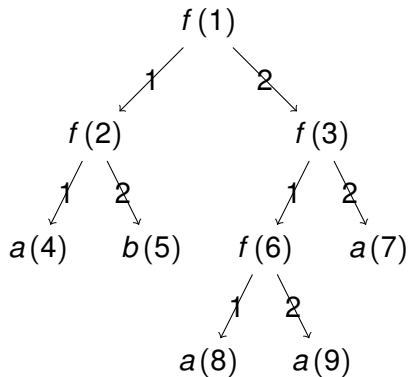
- Properties of Euler chains



The subchain between the first and last occurrence of a node:

The chain of the subtree rooted at that node:

1**24252**131

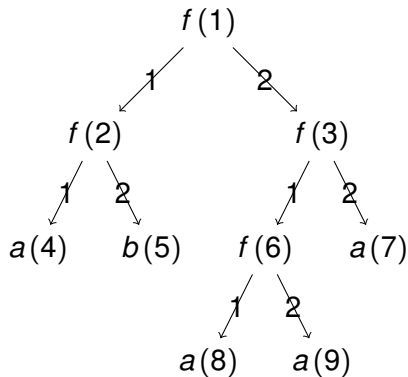
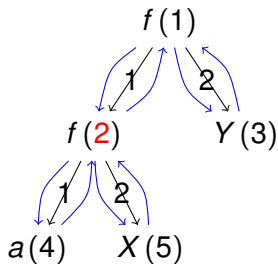


12425213686963731



Euler Chains and Strings

- Properties of Euler chains



A node with n children
occurs $n + 1$ times

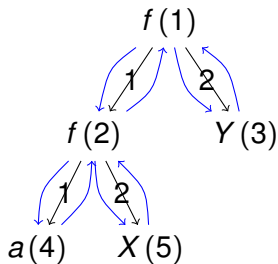
1**2**425**2**131

12425213686963731

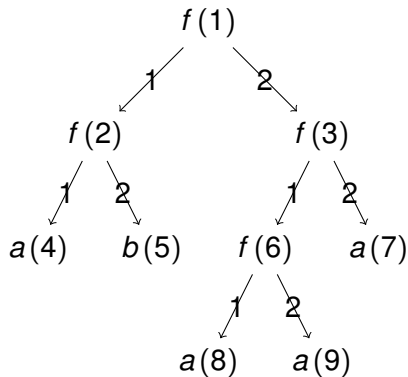


Euler Chains and Strings

- ▶ Euler strings: Replace nodes in Euler chains with node labels.



ffafXffYf



ffafbffffafaffaff

Tree Pattern Matching: Idea

- ▶ Instead of using the tree structure, the algorithm operates on Euler chains and Euler strings.

Tree Pattern Matching: Idea

- ▶ Instead of using the tree structure, the algorithm operates on Euler chains and Euler strings.
- ▶ To declare a match of the pattern tree at a subtree of the subject tree, the algorithm
 - ▶ verifies whether their Euler strings are identical after replacing the variables in the pattern by Euler strings of appropriate terms.



Tree Pattern Matching: Idea

- ▶ Instead of using the tree structure, the algorithm operates on Euler chains and Euler strings.
- ▶ To declare a match of the pattern tree at a subtree of the subject tree, the algorithm
 - ▶ verifies whether their Euler strings are identical after replacing the variables in the pattern by Euler strings of appropriate terms.
- ▶ To justify this approach, Euler strings have to be related to the tree structures.



Tree Pattern Matching: Idea

- ▶ Instead of using the tree structure, the algorithm operates on Euler chains and Euler strings.
- ▶ To declare a match of the pattern tree at a subtree of the subject tree, the algorithm
 - ▶ verifies whether their Euler strings are identical after replacing the variables in the pattern by Euler strings of appropriate terms.
- ▶ To justify this approach, Euler strings have to be related to the tree structures.

Theorem

Two term trees are equivalent (i.e. they represent the same term) iff their corresponding Euler strings are identical.



Nonlinear Tree Pattern Matching: Ideas

Putting the ideas together:

1. Ignore multiplicity of variables (assume the pattern is linear) and do [linear tree pattern matching](#).
2. Verify that the substitutions computed for multiple occurrences of a variable are identical: [check consistency](#).
3. Instead of trees, operate on their [Euler strings](#).



Notation

- ▶ s : Subject tree.
- ▶ p : Pattern tree.
- ▶ C_s and E_s : Euler chain and Euler string for the subject tree.
- ▶ C_p and E_p : Euler chain and Euler string for the pattern tree.
- ▶ n : Size of s .
- ▶ m : Size of p .
- ▶ k : Number of variables in p .
- ▶ K : The set of all root-to-variable-leaf pathes in p .



Step 1. Linear Tree Pattern Matching

- ▶ Let v_1, \dots, v_k be the variables in p .
- ▶ v_1, \dots, v_k appear only once in E_p , because
 - ▶ only leaves are labeled with variables,
 - ▶ each leaf appears exactly once in the Euler string, and
 - ▶ each variable occurs exactly once in p (linearity).



Step 1. Linear Tree Pattern Matching

We start with a simple algorithm.

- ▶ E_s is stored in an array.

Step 1. Linear Tree Pattern Matching

We start with a simple algorithm.

- ▶ E_s is stored in an array.
- ▶ Split E_p into $k + 1$ strings, denoted $\sigma_1, \dots, \sigma_{k+1}$, by removing variables.



Step 1. Linear Tree Pattern Matching

We start with a simple algorithm.

- ▶ E_s is stored in an array.
- ▶ Split E_p into $k + 1$ strings, denoted $\sigma_1, \dots, \sigma_{k+1}$, by removing variables.
 - ▶ $ffafXffYf$ splits into $\sigma_1 = ffaf$, $\sigma_2 = ff$, and $\sigma_3 = f$.



Step 1. Linear Tree Pattern Matching

We start with a simple algorithm.

- ▶ E_s is stored in an array.
- ▶ Split E_p into $k + 1$ strings, denoted $\sigma_1, \dots, \sigma_{k+1}$, by removing variables.
 - ▶ $ffafXffYf$ splits into $\sigma_1 = ffaf$, $\sigma_2 = ff$, and $\sigma_3 = f$.
- ▶ Construct Boolean tables M_1, \dots, M_k , each having $|E_s|$ entries:

$$M_i[j] = \begin{cases} 1 & \text{if there is a match for } \sigma_i \text{ in } E_s \text{ starting at pos. } j \\ 0 & \text{otherwise.} \end{cases}$$



Step 1. Linear Tree Pattern Matching

Example

- ▶ $E_p = \text{ffafXffYf}$, $\sigma_1 = \text{ffaf}$, $\sigma_2 = \text{ff}$, $\sigma_3 = \text{f}$,
 $E_s = \text{ffafbffffafaffff}$.
- ▶ $M_1 = 10000001000010000$ (**f**fafb**ff**ffafa**f**ffaff).
- ▶ $M_2 = 10000111000010010$ (**f**fafb**fff**fafafa**f**ffaff).
- ▶ $M_3 = 11010111100011011$ (**ff**a**fb****ffff**a**f**a**ff**a**ff**).



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	12425213686963731
$E_p =$	ffaafXffYf	$E_s =$	ffaafbffffafaffaff
$\sigma_1 =$	ffaaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	110101111100011011

- ▶ We start from $M_1 = 10000001000010000$.



Step 1. Linear Tree Pattern Matching

$$\begin{array}{ll} p = f(f(a, X), Y) & s = f(f(a, b), f(f(a, a), a)) \\ C_p = 124252131 & C_s = 12425213686963731 \\ E_p = ffa f X f f Y f & E_s = ffa f b f f f f a f a f f a f f \\ \sigma_1 = ffa f & M_1 = 10000001000010000 \\ \sigma_2 = ff & M_2 = 10000111000010010 \\ \sigma_3 = f & M_3 = 11010111100011011 \end{array}$$

- ▶ We start from $M_1 = 10000001000010000$.
- ▶ The set of nodes where p matches s is a subset of the set of nodes with nonzero entries in M_1 .



Step 1. Linear Tree Pattern Matching

$$\begin{array}{ll} p = f(f(a, X), Y) & s = f(f(a, b), f(f(a, a), a)) \\ C_p = 124252131 & C_s = 12425213686963731 \\ E_p = ffa fXffYf & E_s = ffa fbf fffafa ffa fff \\ \sigma_1 = ffa f & M_1 = 10000001000010000 \\ \sigma_2 = ff & M_2 = 10000111000010010 \\ \sigma_3 = f & M_3 = 11010111100011011 \end{array}$$

- ▶ We start from $M_1 = 10000001000010000$.
- ▶ The set of nodes where p matches s is a subset of the set of nodes with nonzero entries in M_1 .
- ▶ Take a nonzero entry position i in M_1 that corresponds to the first occurrence of a node in the Euler chain, $i = 1$.



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	12425213686963731
$E_p =$	ffafXffYf	$E_s =$	ffafbffffafaffaff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	11010111100011011

- ▶ The replacement for X must be a string in E_s that starts at position $i + |\sigma_1| = 5$



Step 1. Linear Tree Pattern Matching

$$\begin{array}{ll} p = f(f(a, X), Y) & s = f(f(a, b), f(f(a, a), a)) \\ C_p = 124252131 & C_s = 12425213686963731 \\ E_p = ffa f X f f Y f & E_s = ffa f b f f f f a f a f f a f f \\ \sigma_1 = ffa f & M_1 = 10000001000010000 \\ \sigma_2 = ff & M_2 = 10000111000010010 \\ \sigma_3 = f & M_3 = 11010111100011011 \end{array}$$

- ▶ The replacement for X must be a string in E_s that starts at position $i + |\sigma_1| = 5$
- ▶ Moreover, this position must correspond to the first occurrence of a node in the Euler chain, because
 - ▶ variables can be substituted by subtrees only,
 - ▶ a subtree starts with the first occurrence of a node in the Euler chain.

If this is not the case, take another nonzero entry position in M_1 .



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	12425213686963731
$E_p =$	ffafXffYf	$E_s =$	ffafbffffafaffff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	11010111100011011

- Replacement for X is a substring of E_s between the first and last occurrences of the node at position $i + |\sigma_1|$.



Step 1. Linear Tree Pattern Matching

$$\begin{array}{ll} p = f(f(a, X), Y) & s = f(f(a, b), f(f(a, a), a)) \\ C_p = 124252131 & C_s = 12425213686963731 \\ E_p = ffafXffYf & E_s = ffafbffffafaffff \\ \sigma_1 = ffaf & M_1 = 10000001000010000 \\ \sigma_2 = ff & M_2 = 10000111000010010 \\ \sigma_3 = f & M_3 = 11010111100011011 \end{array}$$

- ▶ Replacement for X is a substring of E_s between the first and last occurrences of the node at position $i + |\sigma_1|$.
- ▶ Let j be the position of the last occurrence from the previous item. Then $M_2[j + 1]$ should be 1: σ_2 should match E_s at this position.



Step 1. Linear Tree Pattern Matching

$$\begin{array}{ll} p = f(f(a, X), Y) & s = f(f(a, b), f(f(a, a), a)) \\ C_p = 124252131 & C_s = 12425213686963731 \\ E_p = ffafXffYf & E_s = ffafbffffafafff \\ \sigma_1 = ffaf & M_1 = 10000001000010000 \\ \sigma_2 = ff & M_2 = 10000111000010010 \\ \sigma_3 = f & M_3 = 11010111100011011 \end{array}$$

- ▶ Replacement for X is a substring of E_s between the first and last occurrences of the node at position $i + |\sigma_1|$.
- ▶ Let j be the position of the last occurrence from the previous item. Then $M_2[j + 1]$ should be 1: σ_2 should match E_s at this position.
- ▶ And proceed in the same way...



Step 1. Linear Tree Pattern Matching

$$\begin{array}{ll} p = f(f(a, X), Y) & s = f(f(a, b), f(f(a, a), a)) \\ C_p = 124252131 & C_s = 12425213686963731 \\ E_p = ffa f \mathbf{X} f f Y f & E_s = ffa f \mathbf{b} f f \mathbf{f} f a f a f f a f f \\ \sigma_1 = ffa f & M_1 = 10000001000010000 \\ \sigma_2 = f f & M_2 = 10000 \mathbf{1} 11000010010 \\ \sigma_3 = f & M_3 = 11010111100011011 \end{array}$$

- ▶ Replacement for X is a substring of E_s between the first and last occurrences of the node at position $i + |\sigma_1|$.
- ▶ Let j be the position of the last occurrence from the previous item. Then $M_2[j + 1]$ should be 1: σ_2 should match E_s at this position.
- ▶ And proceed in the same way...



Step 1. Linear Tree Pattern Matching

$$\begin{array}{ll} p = f(f(a, X), Y) & s = f(f(a, b), f(f(a, a), a)) \\ C_p = 124252131 & C_s = 12425213686963731 \\ E_p = ffafXffYf & E_s = ffafbfffafafff \\ \sigma_1 = ffaf & M_1 = 10000001000010000 \\ \sigma_2 = ff & M_2 = 10000111000010010 \\ \sigma_3 = f & M_3 = 11010111100011011 \end{array}$$

- ▶ Replacement for X is a substring of E_s between the first and last occurrences of the node at position $i + |\sigma_1|$.
- ▶ Let j be the position of the last occurrence from the previous item. Then $M_2[j + 1]$ should be 1: σ_2 should match E_s at this position.
- ▶ And proceed in the same way...



Step 1. Linear Tree Pattern Matching

$$\begin{array}{ll} p = f(f(a, X), Y) & s = f(f(a, b), f(f(a, a), a)) \\ C_p = 124252131 & C_s = 12425213686963731 \\ E_p = ffa f X f f Y f & E_s = ffa f b f f f f a f a f f a f f \\ \sigma_1 = f f a f & M_1 = 10000001000010000 \\ \sigma_2 = f f & M_2 = 10000111000010010 \\ \sigma_3 = f & M_3 = 11010111100011011 \end{array}$$

- ▶ Replacement for X is a substring of E_s between the first and last occurrences of the node at position $i + |\sigma_1|$.
- ▶ Let j be the position of the last occurrence from the previous item. Then $M_2[j + 1]$ should be 1: σ_2 should match E_s at this position.
- ▶ And proceed in the same way...



Step 1. Linear Tree Pattern Matching

$$\begin{array}{ll} p = f(f(a, X), Y) & s = f(f(a, b), f(f(a, a), a)) \\ C_p = 124252131 & C_s = 12425213686963731 \\ E_p = ffa f X f f Y f & E_s = ffa f b f f f f a f a f f a f f \\ \sigma_1 = ffa f & M_1 = 10000001000010000 \\ \sigma_2 = f f & M_2 = 10000111000010010 \\ \sigma_3 = f & M_3 = 11010111100011011 \end{array}$$

- ▶ Replacement for X is a substring of E_s between the first and last occurrences of the node at position $i + |\sigma_1|$.
- ▶ Let j be the position of the last occurrence from the previous item. Then $M_2[j + 1]$ should be 1: σ_2 should match E_s at this position.
- ▶ And proceed in the same way...



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	12425213686963731
$E_p =$	ffafXffYf	$E_s =$	ffafbfffaffaffaff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	110101111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	124252136886963731
$E_p =$	ffafXffYf	$E_s =$	ffafbffffafaffaff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	110101111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	124252136886963731
$E_p =$	ffafXffYf	$E_s =$	ffafbffffafaffaff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	11010111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	124252136886963731
$E_p =$	ffafXffYf	$E_s =$	ffafbffffaf a ffaff
$\sigma_1 =$	ffaf	$M_1 =$	1000000 1 000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	110101111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	12425213686963731
$E_p =$	ffaafXffYf	$E_s =$	ffaafbffffafaffaff
$\sigma_1 =$	ffaaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	11010111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	12425213686963731
$E_p =$	ffafXffYf	$E_s =$	ffafbfffffafaffaff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	11010111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	12425213686963731
$E_p =$	ffafXffYf	$E_s =$	ffafbfffffafaffaff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	11010111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	12425213686963731
$E_p =$	ffafXffYf	$E_s =$	ffafbfffffafaffaff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	11010111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	12425213686963731
$E_p =$	ffafXffYf	$E_s =$	ffafbfffffafaffaff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	11010111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	12425213686963731
$E_p =$	ffafXffYf	$E_s =$	ffafbfffffafaffaff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	11010111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	12425213686963731
$E_p =$	ffafXffYf	$E_s =$	ffafbfffffafaffaff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	11010111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	12425213686963731
$E_p =$	ffafXffYf	$E_s =$	ffafbfffffafaffaff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	11010111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt gives $X \rightarrow a, Y \rightarrow a$.



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	124252136886963731
$E_p =$	ffafXffYf	$E_s =$	ffafbffffafaffaff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	11010111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt gives $X \rightarrow a, Y \rightarrow a$.
- ▶ One more try...



Step 1. Linear Tree Pattern Matching

$p =$	$f(f(a, X), Y)$	$s =$	$f(f(a, b), f(f(a, a), a))$
$C_p =$	124252131	$C_s =$	124252136886963731
$E_p =$	ffafXffYf	$E_s =$	ffafbffffafaffaff
$\sigma_1 =$	ffaf	$M_1 =$	10000001000010000
$\sigma_2 =$	ff	$M_2 =$	10000111000010010
$\sigma_3 =$	f	$M_3 =$	11010111100011011

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt gives $X \rightarrow a, Y \rightarrow a$.
- ▶ One more try...



Step 1. Linear Tree Pattern Matching

$p = f(f(a, X), Y)$	$s = f(f(a, b), f(f(a, a), a))$
$C_p = 124252131$	$C_s = 12425213686963731$
$E_p = ffa fXffYf$	$E_s = ffa fbf fffafa ffa f f$
$\sigma_1 = ffa f$	$M_1 = 10000001000010000$
$\sigma_2 = ff$	$M_2 = 10000111000010010$
$\sigma_3 = f$	$M_3 = 11010111100011011$

- ▶ The first match found: $X \rightarrow b, Y \rightarrow f(f(a, a), a)$.
- ▶ The next attempt gives $X \rightarrow a, Y \rightarrow a$.
- ▶ One more try... fail.
- ▶ The last 1 in C_s is not the first occurrence of 1.



Complexity of Linear Tree Pattern Matching

- ▶ The simple algorithm computes $k + 1$ Boolean tables.
- ▶ Each table has $|E_s| = n$ size.
- ▶ In total, construction of the tables takes $O(nk)$ time.
- ▶ Room for improvement: Do not compute them explicitly.



Suffix Number, Suffix Index

Ψ : finite set of strings.

- ▶ **Suffix number** of a string λ in Ψ : The number of strings in Ψ which are suffixes of λ .
- ▶ **Suffix index** of Ψ (denoted Ψ^*): The maximum among all suffix numbers of strings in Ψ .
- ▶ If $|\Psi| = 0$ then $\Psi^* = 1$.

Example

- ▶ $\Psi = \{ffffX, fffb, ffb, fb\}$. $|\Psi| = 5$.
- ▶ Suffix number of $ffffX$ in Ψ is 1.
- ▶ Suffix number of ffb in Ψ is 3.
- ▶ Suffix number of $fffb$ in Ψ is 4.
- ▶ Suffix index of Ψ is 4.



Complexity of Linear Tree Pattern Matching

How many replacements at most are possible (independent of the algorithm)?



Complexity of Linear Tree Pattern Matching

How many replacements at most are possible (independent of the algorithm)?

- ▶ Assume p matches s at some node.

Complexity of Linear Tree Pattern Matching

How many replacements at most are possible (independent of the algorithm)?

- ▶ Assume p matches s at some node.
- ▶ If X at node i in p matches a subtree at node w in s then



Complexity of Linear Tree Pattern Matching

How many replacements at most are possible (independent of the algorithm)?

- ▶ Assume p matches s at some node.
- ▶ If X at node i in p matches a subtree at node w in s then
 - ▶ w is called a **legal replacement** for X ,
 - ▶ $path_1 \circ lp(r_p, i') = lp(r_s, w)$. (i' : i labeled with $(lab(w))$.)



Complexity of Linear Tree Pattern Matching

How many replacements at most are possible (independent of the algorithm)?

- ▶ Assume p matches s at some node.
- ▶ If X at node i in p matches a subtree at node w in s then
 - ▶ w is called a **legal replacement** for X ,
 - ▶ $path_1 \circ lp(r_p, i') = lp(r_s, w)$. (i' : i labeled with $(lab(w))$.)
- ▶ If another variable Y at node j in p matches w (in another match) then



Complexity of Linear Tree Pattern Matching

How many replacements at most are possible (independent of the algorithm)?

- ▶ Assume p matches s at some node.
- ▶ If X at node i in p matches a subtree at node w in s then
 - ▶ w is called a **legal replacement** for X ,
 - ▶ $path_1 \circ lp(r_p, i') = lp(r_s, w)$. (i' : i labeled with $(lab(w))$.)
- ▶ If another variable Y at node j in p matches w (in another match) then
 - ▶ $path_2 \circ lp(r_p, j') = lp(r_s, w)$.



Complexity of Linear Tree Pattern Matching

How many replacements at most are possible (independent of the algorithm)?

- ▶ Assume p matches s at some node.
- ▶ If X at node i in p matches a subtree at node w in s then
 - ▶ w is called a **legal replacement** for X ,
 - ▶ $path_1 \circ lp(r_p, i') = lp(r_s, w)$. (i' : i labeled with $(lab(w))$.)
- ▶ If another variable Y at node j in p matches w (in another match) then
 - ▶ $path_2 \circ lp(r_p, j') = lp(r_s, w)$.
- ▶ Therefore, $lp(r_p, j')$ is a suffix of $lp(r_p, i')$, or vice versa.



Complexity of Linear Tree Pattern Matching

How many replacements at most are possible (independent of the algorithm)?

- ▶ Assume p matches s at some node.
- ▶ If X at node i in p matches a subtree at node w in s then
 - ▶ w is called a **legal replacement** for X ,
 - ▶ $path_1 \circ lp(r_p, i') = lp(r_s, w)$. (i' : i labeled with $(lab(w))$.)
- ▶ If another variable Y at node j in p matches w (in another match) then
 - ▶ $path_2 \circ lp(r_p, j') = lp(r_s, w)$.
- ▶ Therefore, $lp(r_p, j')$ is a suffix of $lp(r_p, i')$, or vice versa.
- ▶ Hence, the subtree at w can be substituted at most K^* times over all matches and the number of all legal replacements that can be computed over all matches is $O(nK^*)$.



Complexity of Linear Tree Pattern Matching

Bound on the number of replacements computed by the simple algorithm:

- ▶ Assume e_1, \dots, e_w are Euler strings of subtrees in s rooted at nodes i_1, \dots, i_w .
- ▶ Assume the string $\sigma_1 \circ e_1 \circ \dots \circ e_w \circ \sigma_{w+1}$ matches a substring of E_s at position l .
- ▶ l is the position that corresponds to the first occurrence of a node j in s , i.e. $C_s[l] = j$ and $C_s[l'] \neq j$ for all $l' < l$.
- ▶ For each $1 < q < w$, $lp(r_p, v_q) = lp(j, i_q)$ (v 's are the corresponding variable nodes in p .)
- ▶ The strings $\sigma_1, \sigma_1 \circ e_1, \sigma_1 \circ e_1 \circ \sigma_2, \dots$ are computed incrementally.
- ▶ We have a match at j if we compute $\sigma_1 \circ e_1 \circ \dots \circ e_k \circ \sigma_{k+1}$, i.e. legal replacements for all variables in p .



Complexity of Linear Tree Pattern Matching

Bound on the number of replacements computed by the simple algorithm:

- ▶ In case of failed match attempt, we would have computed at most one illegal replacement.
- ▶ Hence, the total number of illegal replacements computed over match attempts at all nodes can be $O(n)$ at most.
- ▶ Therefore, the upper bound of the replacements computed by the algorithm is $O(nK^*)$.



Complexity of Linear Tree Pattern Matching

Bound on the number of replacements computed by the simple algorithm:

- ▶ In case of failed match attempt, we would have computed at most one illegal replacement.
- ▶ Hence, the total number of illegal replacements computed over match attempts at all nodes can be $O(n)$ at most.
- ▶ Therefore, the upper bound of the replacements computed by the algorithm is $O(nK^*)$.

That's fine, but how to keep the time-bound of the algorithm proportional to the number of replacements?



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ Do not spend more than $O(1)$ between replacements, **without** computing the tables explicitly.
- ▶ Do a replacement in $O(1)$.



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ Do not spend more than $O(1)$ between replacements, **without** computing the tables explicitly.
- ▶ Do a replacement in $O(1)$.
- ▶ Doing a replacement in $O(1)$ is easy:
 - ▶ Store an Euler string in an array along with a pointer from the first occurrence of a node to its last occurrence.
 - ▶ Check whether the replacement begins at the first occurrence of a node.
 - ▶ If yes, skip to its last occurrence.



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ Do not spend more than $O(1)$ between replacements, **without** computing the tables explicitly.
- ▶ Do a replacement in $O(1)$.
- ▶ Doing a replacement in $O(1)$ is easy:
 - ▶ Store an Euler string in an array along with a pointer from the first occurrence of a node to its last occurrence.
 - ▶ Check whether the replacement begins at the first occurrence of a node.
 - ▶ If yes, skip to its last occurrence.
- ▶ Not spending more than $O(1)$ between replacements, without computing the tables, needs more preprocessing.



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ What happens in the steps preceding the replacement for v_{i+1} , after computing a replacement for v_i ?



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ What happens in the steps preceding the replacement for v_{i+1} , after computing a replacement for v_i ?
- ▶ Determine whether pattern string σ_{i+1} matches E_s at the position following the replacement for v_i .



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ What happens in the steps preceding the replacement for v_{i+1} , after computing a replacement for v_i ?
- ▶ Determine whether pattern string σ_{i+1} matches E_s at the position following the replacement for v_i .
- ▶ Had we computed the tables, this can be done in $O(1)$, but how to achieve the same without the tables?



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ Problem: Given a position in E_s and a string σ_i , $1 \leq i \leq k + 1$, decide in $O(1)$ whether σ_i matches E_s in that position.



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ Problem: Given a position in E_s and a string σ_i , $1 \leq i \leq k + 1$, decide in $O(1)$ whether σ_i matches E_s in that position.
- ▶ Idea:



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ Problem: Given a position in E_s and a string σ_i , $1 \leq i \leq k + 1$, decide in $O(1)$ whether σ_i matches E_s in that position.
- ▶ Idea:
 - ▶ Preprocess the pattern strings to produce an automaton that recognizes every instance of these $k + 1$ strings.



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ Problem: Given a position in E_s and a string σ_i , $1 \leq i \leq k + 1$, decide in $O(1)$ whether σ_i matches E_s in that position.
- ▶ Idea:
 - ▶ Preprocess the pattern strings to produce an automaton that recognizes every instance of these $k + 1$ strings.
 - ▶ Use the automaton to recognize these strings in E_s .



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ Problem: Given a position in E_s and a string σ_i , $1 \leq i \leq k + 1$, decide in $O(1)$ whether σ_i matches E_s in that position.
- ▶ Idea:
 - ▶ Preprocess the pattern strings to produce an automaton that recognizes every instance of these $k + 1$ strings.
 - ▶ Use the automaton to recognize these strings in E_s .
 - ▶ With every position in an array containing E_s , store the state of the automaton on reading the symbol in that position.



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ Problem: Given a position in E_s and a string σ_i , $1 \leq i \leq k + 1$, decide in $O(1)$ whether σ_i matches E_s in that position.
- ▶ Idea:
 - ▶ Preprocess the pattern strings to produce an automaton that recognizes every instance of these $k + 1$ strings.
 - ▶ Use the automaton to recognize these strings in E_s .
 - ▶ With every position in an array containing E_s , store the state of the automaton on reading the symbol in that position.
 - ▶ In order to decide whether a pattern string σ_i matches the substring of E_s at position j , look at the state of the automaton in position $j + |\sigma_i| - 1$.



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ Problem: Given a position in E_s and a string σ_i , $1 \leq i \leq k + 1$, decide in $O(1)$ whether σ_i matches E_s in that position.
- ▶ Idea:
 - ▶ Preprocess the pattern strings to produce an automaton that recognizes every instance of these $k + 1$ strings.
 - ▶ Use the automaton to recognize these strings in E_s .
 - ▶ With every position in an array containing E_s , store the state of the automaton on reading the symbol in that position.
 - ▶ In order to decide whether a pattern string σ_i matches the substring of E_s at position j , look at the state of the automaton in position $j + |\sigma_i| - 1$.
 - ▶ Lookup from the array in $O(1)$ time.



Complexity of Linear Tree Pattern Matching

Keep the time-bound of the algorithm proportional to the number of replacements:

- ▶ Problem: Given a position in E_s and a string σ_i , $1 \leq i \leq k + 1$, decide in $O(1)$ whether σ_i matches E_s in that position.
- ▶ Idea:
 - ▶ Preprocess the pattern strings to produce an automaton that recognizes every instance of these $k + 1$ strings.
 - ▶ Use the automaton to recognize these strings in E_s .
 - ▶ With every position in an array containing E_s , store the state of the automaton on reading the symbol in that position.
 - ▶ In order to decide whether a pattern string σ_i matches the substring of E_s at position j , look at the state of the automaton in position $j + |\sigma_i| - 1$.
 - ▶ Lookup from the array in $O(1)$ time.

How to preprocess the pattern strings?



Modifying the Linear Tree Pattern Matching

Pattern string preprocessing:

- ▶ The $k + 1$ pattern strings $\sigma_1, \dots, \sigma_{k+1}$ are preprocessed to produce an automation that recognizes every instance of these strings.
- ▶ Method: Aho-Corasick (AC) algorithm.
- ▶ The AC algorithm constructs the desired automaton in time proportional to the sum of the lengths of all pattern strings.



Modifying the Linear Tree Pattern Matching

What does a Aho-Corasick automaton for a set of pattern strings $\sigma_1, \dots, \sigma_{k+1}$ do?

- ▶ Takes the subject string E_S as input.
- ▶ Outputs the locations in E_S at which the σ 's appear as substrings, together with the corresponding σ 's.
- ▶ For example, a Aho-Corasick automaton for the strings *he*, *she*, *his*, *hers* returns on the input string *ushers* the locations 4 (match for *she* and *he*) and 6 (match for *hers*).



Modifying the Linear Tree Pattern Matching

Aho-Corasick automaton

- ▶ consists of a set of states, represented by numbers,
- ▶ processes the subject string by successively reading symbols in it, making state transitions and occasionally emitting output,
- ▶ is controlled by three functions:
 1. a goto function g ,
 2. a failure function f ,
 3. a output function *output*.



Modifying the Linear Tree Pattern Matching

Construction of the Aho-Corasick automation:

- ▶ Determine the states and the goto function.
- ▶ Compute the failure function.
- ▶ Computation of the output function begins on the first step and is completed on the second.



Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automation for the pattern strings *he, she, his, hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:

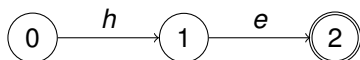


Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automation for the pattern strings *he*, *she*, *his*, *hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:



$$\text{output}(2) = \{he\}$$

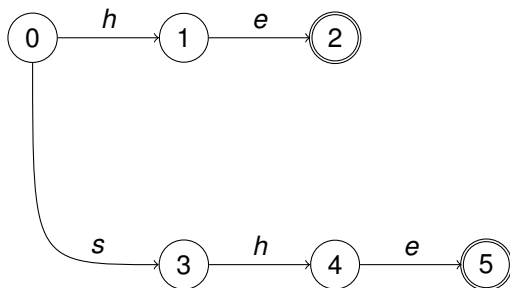


Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automation for the pattern strings *he*, *she*, *his*, *hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:



$output(2) = \{he\}$

$output(5) = \{she\}$

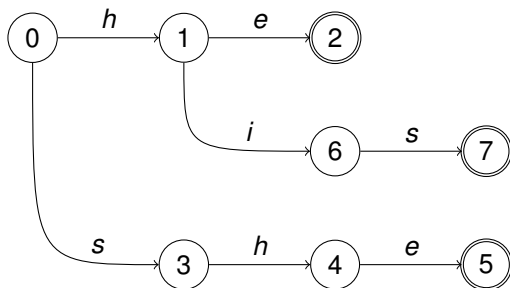


Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automaton for the pattern strings *he*, *she*, *his*, *hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:



$output(2) = \{he\}$

$output(5) = \{she\}$

$output(7) = \{his\}$

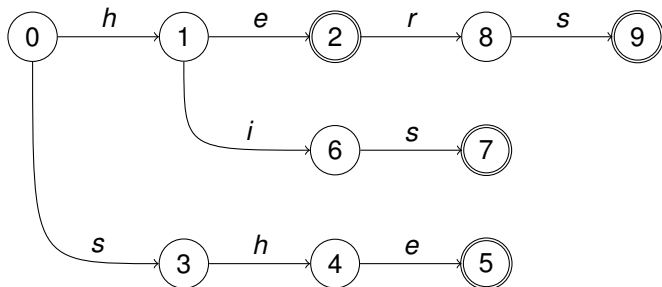


Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automation for the pattern strings *he*, *she*, *his*, *hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:



$output(2) = \{he\}$

$output(5) = \{she\}$

$output(7) = \{his\}$

$output(9) = \{hers\}$

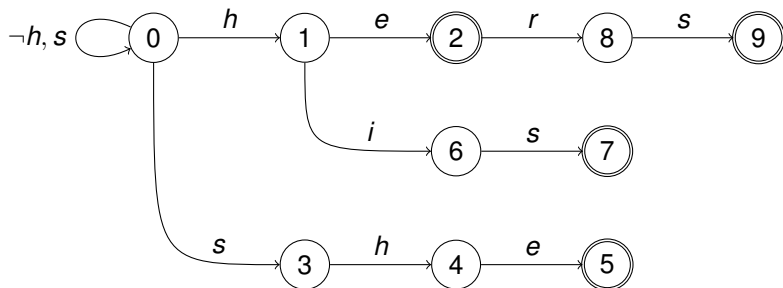


Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automation for the pattern strings *he*, *she*, *his*, *hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:



$output(2) = \{he\}$

$output(5) = \{she\}$

$output(7) = \{his\}$

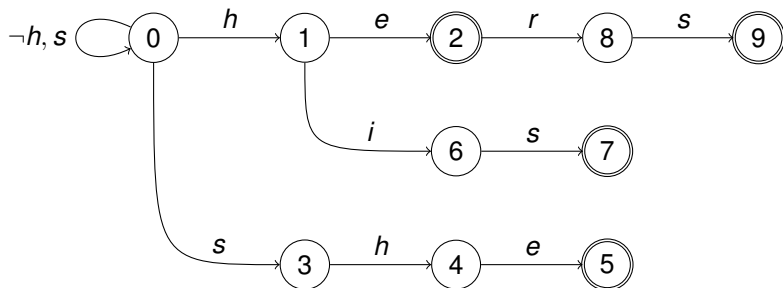
$output(9) = \{hers\}$

Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automation for the pattern strings *he*, *she*, *his*, *hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:



$output(2) = \{he\}$

$output(5) = \{she\}$

$output(7) = \{his\}$

$output(9) = \{hers\}$

$f(1) = f(3) = 0$

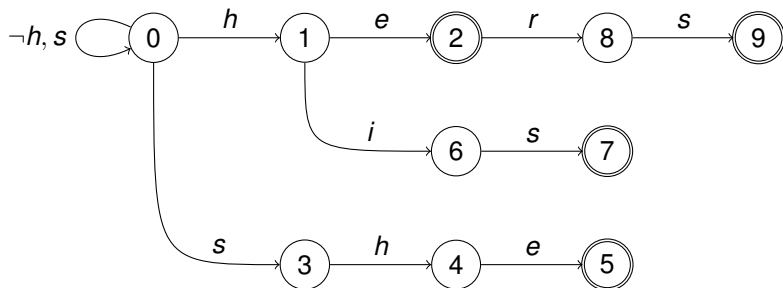


Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automation for the pattern strings *he*, *she*, *his*, *hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:



$output(2) = \{he\}$

$output(5) = \{she\}$

$output(7) = \{his\}$

$output(9) = \{hers\}$

$f(1) = f(3) = 0$

$f(2) = 0$

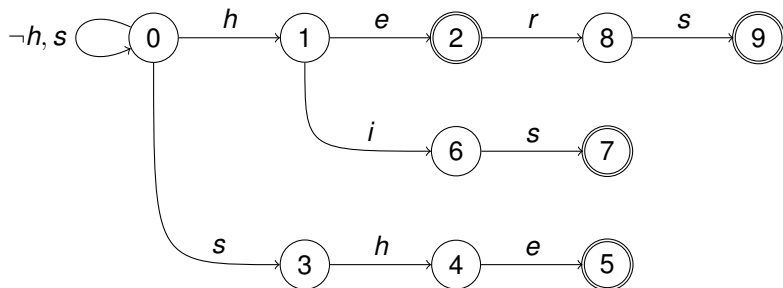


Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automaton for the pattern strings *he*, *she*, *his*, *hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:



$output(2) = \{he\}$

$output(5) = \{she\}$

$output(7) = \{his\}$

$output(9) = \{hers\}$

$f(1) = f(3) = 0$

$f(2) = 0$

$f(6) = 0$

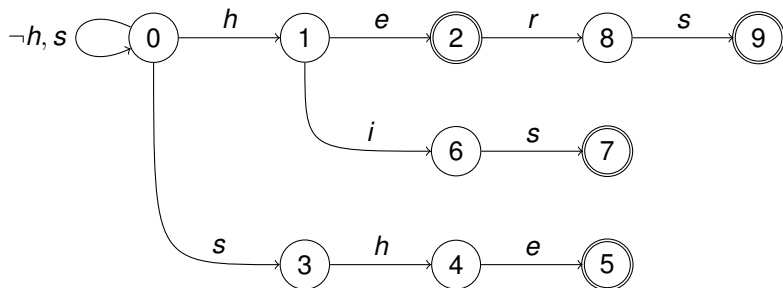


Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automation for the pattern strings *he*, *she*, *his*, *hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:



$output(2) = \{he\}$

$output(5) = \{she\}$

$output(7) = \{his\}$

$output(9) = \{hers\}$

$f(1) = f(3) = 0$

$f(2) = 0$

$f(6) = 0$

$f(4) = 1$

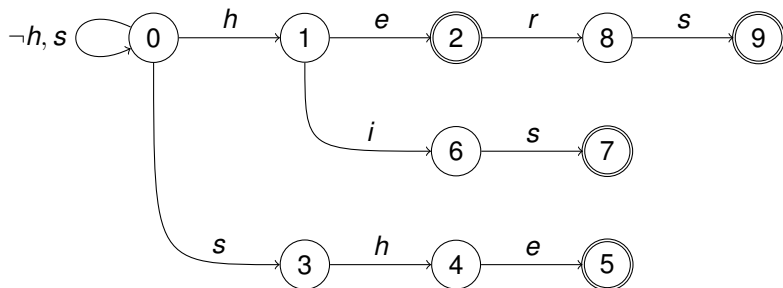


Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automation for the pattern strings *he*, *she*, *his*, *hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:



$output(2) = \{he\}$

$output(5) = \{she\}$

$output(7) = \{his\}$

$output(9) = \{hers\}$

$f(1) = f(3) = 0$ $f(8) = 0$

$f(2) = 0$

$f(6) = 0$

$f(4) = 1$

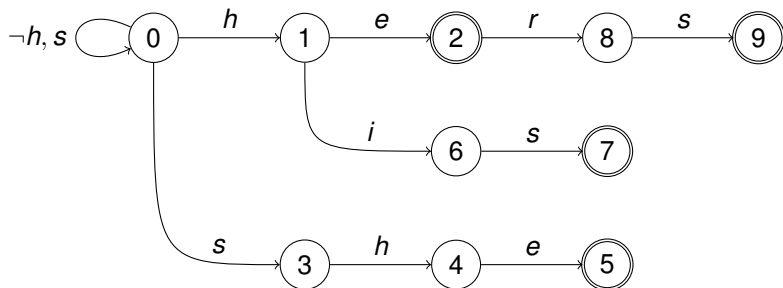


Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automation for the pattern strings *he*, *she*, *his*, *hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:



$output(2) = \{he\}$

$output(5) = \{she\}$

$output(7) = \{his\}$

$output(9) = \{hers\}$

$f(1) = f(3) = 0$

$f(2) = 0$

$f(6) = 0$

$f(4) = 1$

$f(8) = 0$

$f(7) = 3$

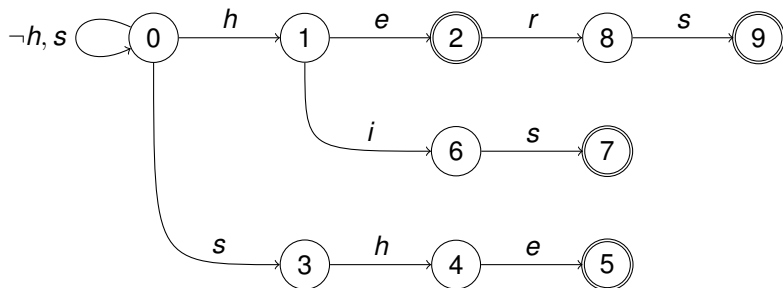


Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automation for the pattern strings *he*, *she*, *his*, *hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:



$output(2) = \{he\}$

$output(5) = \{she, he\}$

$output(7) = \{his\}$

$output(9) = \{hers\}$

$f(1) = f(3) = 0$

$f(2) = 0$

$f(6) = 0$

$f(4) = 1$

$f(8) = 0$

$f(7) = 3$

$f(5) = 2$

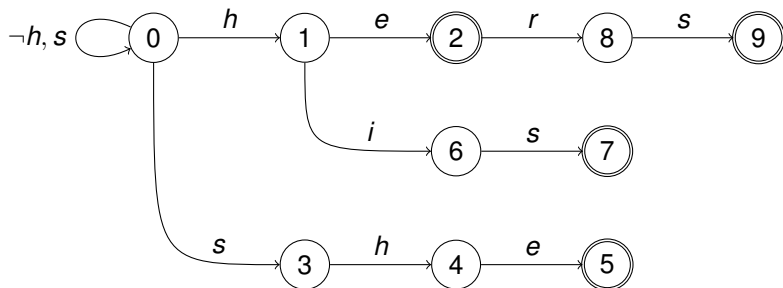


Modifying the Linear Tree Pattern Matching

Example

Construction of the Aho-Corasick automaton for the pattern strings *he*, *she*, *his*, *hers*.

The goto function $g : \text{states} \times \text{letters} \rightarrow \text{states} \cup \{\text{fail}\}$:



$output(2) = \{he\}$

$output(5) = \{she, he\}$

$output(7) = \{his\}$

$output(9) = \{hers\}$

$f(1) = f(3) = 0$

$f(2) = 0$

$f(6) = 0$

$f(4) = 1$

$f(8) = 0$

$f(7) = 3$

$f(5) = 2$

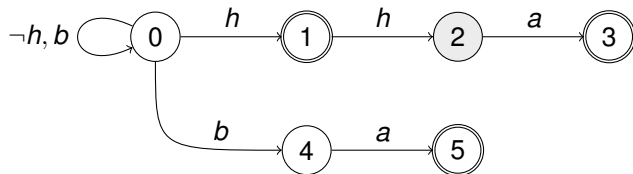
$f(9) = 3$



Modifying the Linear Tree Pattern Matching

Example

AC automation for the pattern strings hha , h , ba .



$$\begin{array}{llll} \text{output}'(1) = \{h\} & \text{output}'(2) = \{h\} & f(1) = 0 & f(5) = 0 \\ \text{output}'(3) = \{hha\} & & f(4) = 0 & f(3) = 0 \\ \text{output}'(5) = \{ba\} & & f(2) = 1 & \end{array}$$

1, 3, 5 : primary accepting states.

2 : secondary accepting state (h is a suffix of hh).



Modifying the Linear Tree Pattern Matching

- ▶ For each secondary accepting state there is a unique primary accepting state with exactly the same output set.
- ▶ Modify construction of AC automaton by maintaining pointers from secondary accepting states to the corresponding accepting states.



Modifying the Linear Tree Pattern Matching

- ▶ Construction of Aho-Corasick automaton takes $O(m)$ time.
- ▶ The output set is represented as a linked list, which is inappropriate for our purpose.
- ▶ Given an arbitrary string, we want to determine in constant time whether it is in the output set.



Modifying the Linear Tree Pattern Matching

- ▶ Construction of Aho-Corasick automaton takes $O(m)$ time.
- ▶ The output set is represented as a linked list, which is inappropriate for our purpose.
- ▶ Given an arbitrary string, we want to determine in constant time whether it is in the output set.
- ▶ Idea: Copy the output set into an array.



Modifying the Linear Tree Pattern Matching

- ▶ Construction of Aho-Corasick automaton takes $O(m)$ time.
- ▶ The output set is represented as a linked list, which is inappropriate for our purpose.
- ▶ Given an arbitrary string, we want to determine in constant time whether it is in the output set.
- ▶ Idea: Copy the output set into an array.
- ▶ Question: How many elements do we have to copy?



Modifying the Linear Tree Pattern Matching

- ▶ Construction of Aho-Corasick automaton takes $O(m)$ time.
- ▶ The output set is represented as a linked list, which is inappropriate for our purpose.
- ▶ Given an arbitrary string, we want to determine in constant time whether it is in the output set.
- ▶ Idea: Copy the output set into an array.
- ▶ Question: How many elements do we have to copy?
- ▶ Answer: As many as in the output sets of all primary accepting states, which is $O(m)$ (because any string in a primary accepting state is a suffix of the longest string in this state.)



Modifying the Linear Tree Pattern Matching

Linear Tree Pattern Matching:

1. Construct Aho-Corasick automaton for the pattern strings.
2. Visit each primary accepting state and copy its output set into a boolean array.
3. Scan the E_s with this automaton.
4. During this process, with each entry in E_s store the state of automaton upon reading the function symbol in that entry.
5. If this state is a secondary accepting state, instead of it store the corresponding primary accepting state.
6. To determine whether there is a match for a pattern string σ at the position i requires verifying the state associated with the $i + |\sigma| - 1$ 'th entry in E_s :
 - ▶ This should be a primary accepting state and
 - ▶ Its output set should contain σ .



Consistency Checking

For nonlinear patterns the computed replacements have to be checked for consistency.

- ▶ Idea: Assign integer codes (from 1 to n) to the nodes in the subject tree.
- ▶ Two nodes get the same encoding iff the subtrees rooted at them are identical.
- ▶ Such an encoding can be computed in $O(n)$.



Consistency Checking

Computing the encoding:

- ▶ Bottom up: First, sort the leaves with respect to their labels and take the ranks as the integers for encoding. Duplicates are assigned the same rank.
- ▶ Suppose the encoding for all nodes up to the height i is computed.
- ▶ Computing the encoding of the nodes at height $i + 1$:
 - ▶ Assign to each node v at the level $i + 1$ a vector $\langle f, j_1, \dots, j_n \rangle$.
 - ▶ f is the label of v and j_i is the encoding of its i 's child.
 - ▶ The vectors assigned to all nodes at $i + 1$ are radix sorted.
 - ▶ If the rank of v is α and the largest encoding among the nodes at level i is β , then the encoding for v is $\alpha + \beta$.



Consistency Checking

Checking consistency:

- ▶ Consistency of replacements is checked as they are computed.
- ▶ For each variable in the pattern, the encoding for the replacement of its first occurrence is computed and is entered into a table.
- ▶ For the next occurrence of the same variable, compare encoding of its replacement to the one in the table.
- ▶ If the check succeeds, proceed further. Otherwise report a failure and start matching procedure at another position in E_S .
- ▶ These steps do not increase the complexity of the algorithm.



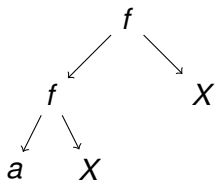
The Last Word

Nonlinear tree pattern matching can be done in $O(nK^*)$ time.



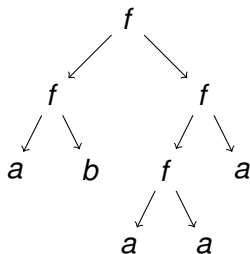
Example

$f(f(a, X), Y)$



Pattern tree p

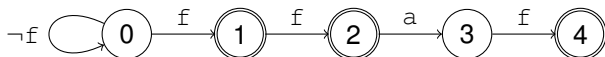
$f(f(a, b), f(f(a, a), a))$



Subject tree s

Example (Cont.)

- ▶ $E_p = \text{ffafXffXf}$
- ▶ AC automaton for the pattern strings ffaf , ff , f .



$output(1) = \{f\}$

$output(2) = \{ff, f\}$

$output(4) = \{ffaf, f\}$

$failure(1) = failure(3) = 0$

$failure(2) = failure(4) = 1$

	ffaf	ff	f
O_1	F	F	T
O_2	F	T	T
O_4	T	F	T



Example (Cont.)

$\sigma_1 = \text{ffaf}$, $\sigma_2 = \text{ff}$, $\sigma_3 = \text{f}$

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
C_s	1	2	4	2	5	2	1	3	6	8	6	9	6	3	7	3	1
E_s	f	f	a	f	b	f	f	f	f	a	f	a	f	f	a	f	f
$IsFirst$	T	T	T	F	T	F	F	T	T	T	F	T	F	F	T	F	F
$lastptr$	17	6	3	-	5	-	-	16	13	8	-	9	-	-	15	-	-
$IsLast$	F	F	T	F	T	T	F	F	F	T	F	T	T	F	T	T	T
$state$	1	2	3	4	0	1	2	2	2	3	4	0	2	2	3	4	2

- ▶ In the first row, the numbers from 1 to 17 - array indices.
- ▶ C_s and E_s - the Euler chain and the Euler string for s .
- ▶ For an index i ,
 - ▶ $IsFirst[i] = \text{T}$ iff $C_s[i]$ occurs first time in C_s .
 - ▶ if $IsFirst[i] = \text{T}$ then $lastptr[i] = j$ where j is the index of the last occurrence of the number $C_s[i]$ in C_s .
 - ▶ $IsLast[i] = \text{T}$ iff $C_s[i]$ occurs last time in C_s .
 - ▶ $state[i]$ is the state of the automaton after reading $E_s[i]$.



Reference



R. Ramesh and I. V. Ramakrishnan.
Nonlinear pattern matching in trees.
J. ACM, 39(2):295–316, 1992.

