

Automata approach to inexact tree pattern matching using 1-degree edit distance

Eliška Šestáková Ondřej Guth Jan Janoušek

Czech Technical University in Prague

PSC 2021

① The inexact tree pattern matching problem

② Automata-based solution

- Pushdown automaton
- Finite automaton

③ Dynamic programming

1-degree edit operations¹

Edit operations applied to a *labeled, ordered* tree $T = (V, E)$:

- **vertex relabel** change the label of a vertex $v \in V$,
- **leaf insert** insert a vertex v as a leaf of an existing vertex $u \in V$, and
- **leaf delete** delete a non-root leaf $v \in V$.

¹Stanley M Selkow. "The tree-to-tree editing problem". In: *Inf. Process. Lett.* 6.6 (1977), pp. 184–186.

Unit cost 1-degree edit distance

Definition (Unit cost 1-degree edit distance)

The *unit cost 1-degree edit distance* is a function $d : \text{TR}(\Sigma)^a \times \text{TR}(\Sigma) \rightarrow \mathbb{N}_0$. Given two trees T_1 and T_2 , the number $d(T_1, T_2)$ corresponds to the minimal number of 1-degree edit operations that transform T_1 into T_2 .

^a $\text{TR}(\Sigma)$ denotes the set of all labeled ordered trees over alphabet Σ

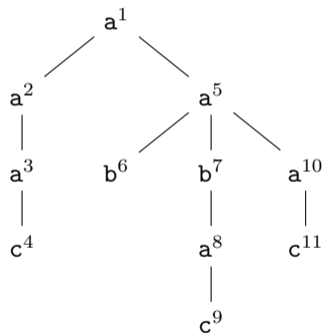
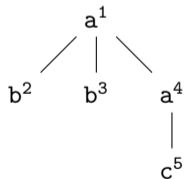
Inexact 1-degree tree pattern matching problem

Problem (Inexact 1-degree tree pattern matching problem)

Let Σ be an alphabet. Let $T = (V_T, E_T)$ be an input tree with n vertices over Σ . Let $P = (V_P, E_P)$ be a comparatively smaller tree pattern over Σ with m vertices. Let k be a non-negative integer representing the maximum number of errors allowed. Let d be the unit cost 1-degree edit distance function. Given T, P, k , and d , the inexact 1-degree tree pattern matching problem is to return a set $\{v : v \in V_T \wedge d(T_v, P) \leq k\}$.

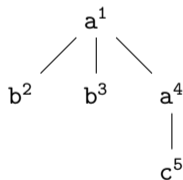
Inexact 1-degree tree pattern matching problem

Example, $k = 2$



Prefix bar notation of a tree²

$$\text{PREFIX-BAR}(P) = a \text{ PREFIX-BAR}(P_1) \text{ PREFIX-BAR}(P_2) \dots \text{ PREFIX-BAR}(P_s) |,$$



$$\text{PREFIX-BAR}(P) = ab|b|ac|||$$

²Jan Janoušek. *Arbology: Algorithms on trees and pushdown automata*. Habilitation thesis, Brno University of Technology, 2010.

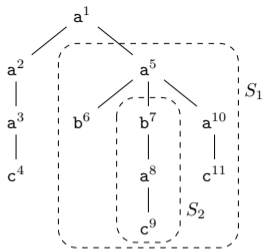
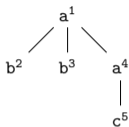
Prefix bar notation of a tree

ab|b|ac|||

(a(b)(b)(a(c))))

Unit cost 1-degree edit distance with prefix bar notation

- relabeling $R(i, b)$ that for $i \in \{1, \dots, r-1\}$, $b \in \Sigma$, and $\mathbf{x}[i] \in (\Sigma \setminus \{b\})$, change the symbol $\mathbf{x}[i]$ into symbol b ;
- insertion $I(i, a)$ that for $i \in \{2, \dots, r-1\}$ and $a \in \Sigma$ inserts the substring (leaf) “ $a|$ ” at position i ; and
- deletion $D(i)$ that for $i \in \{2, \dots, r-2\}$, $\mathbf{x}[i] \in \Sigma$, and $\mathbf{x}[i+1] = |$, deletes the substring (leaf) $\mathbf{x}[i]\mathbf{x}[i+1]$.



$$\text{PREF-BAR}(P) = ab|b|ac||| \quad \text{PREF-BAR}(S_1) = ab|bac|||ac||| \quad \text{PREF-BAR}(S_2) = bac|||$$

Example

$d(P, S_1) = 2$ and $d(P, S_2) = 3$ since

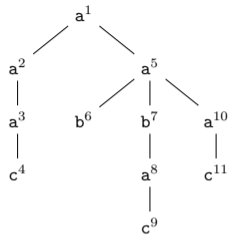
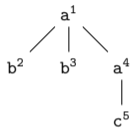
$$ab|b|ac||| \xrightarrow{I(5,a)} ab|ba||ac||| \xrightarrow{I(6,c)} ab|bac|||ac|||$$

$$ab|b|ac||| \xrightarrow{R(1,b)} bb|b|ac||| \xrightarrow{D(2)} bb|ac||| \xrightarrow{D(2)} bac|||$$

Problem statement with prefix bar notation

Problem (inexact 1-degree tree pattern matching)

The problem of inexact 1-degree tree pattern matching is finding all positions $j \in \{1, \dots, 2n\}$ in $\text{PREF-BAR}(T)$ such that $\text{PREF-BAR}(T)[j] = |$ and $d(\text{PREF-BAR}(P), \text{PREF-BAR}(T)[i \dots j]) \leq k$, where i is the start position of the subtree that ends at position j .



Example ($k = 2$)

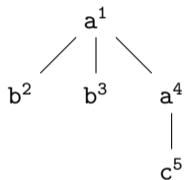
$$\text{PREF-BAR}(P) = ab|b|ac|||$$

$$\text{PREF-BAR}(T) = a \underbrace{aac}||| \overbrace{ab|bac}||| ac||| |$$

Algorithm outline

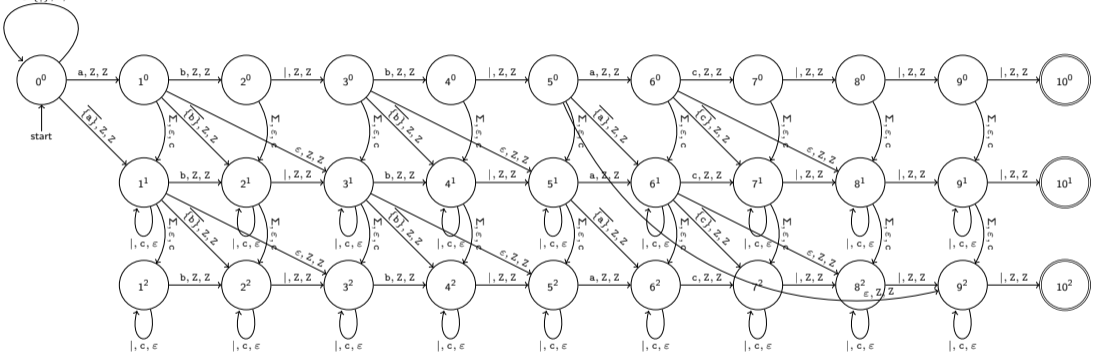
- ① Construct (string) pattern-matching automaton \mathcal{M}_{PDA} for $\text{PREF-BAR}(P)$ and k .
- ② Execute \mathcal{M}_{PDA} over $\text{PREF-BAR}(T)$.

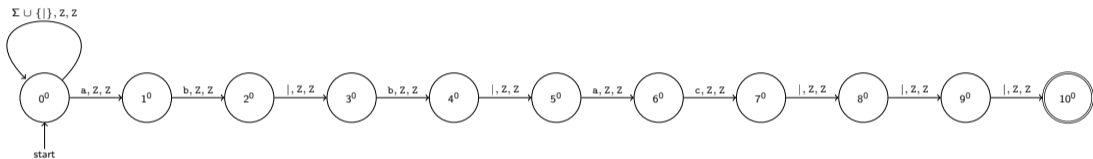
Example: Tree pattern P and $k = 2$



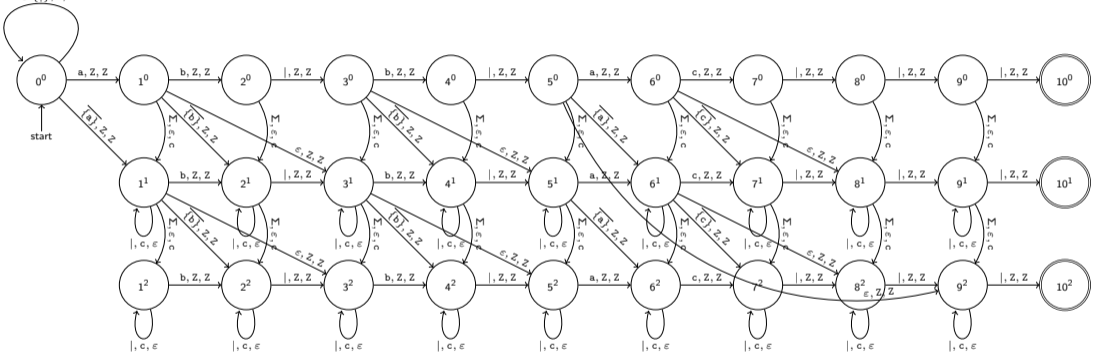
$ab|b|ac|||$

$\Sigma \cup \{\epsilon\}, Z, Z$





$\Sigma \cup \{\epsilon\}, Z, Z$

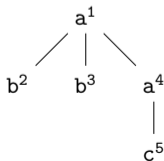


Subtree jump table

Definition (Subtree jump table^a)

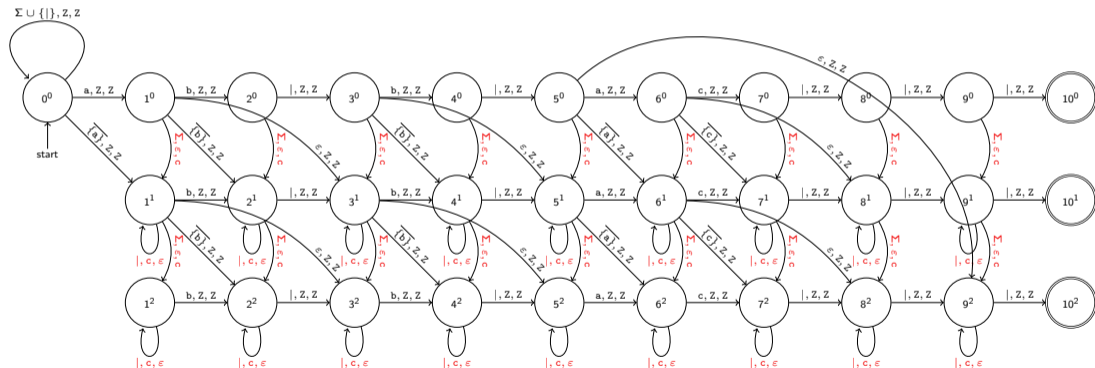
^aJan Trávníček. “(Nonlinear) Tree Pattern Indexing and Backward Matching”. PhD thesis. Faculty of Information Technology, Czech Technical University in Prague, 2018.

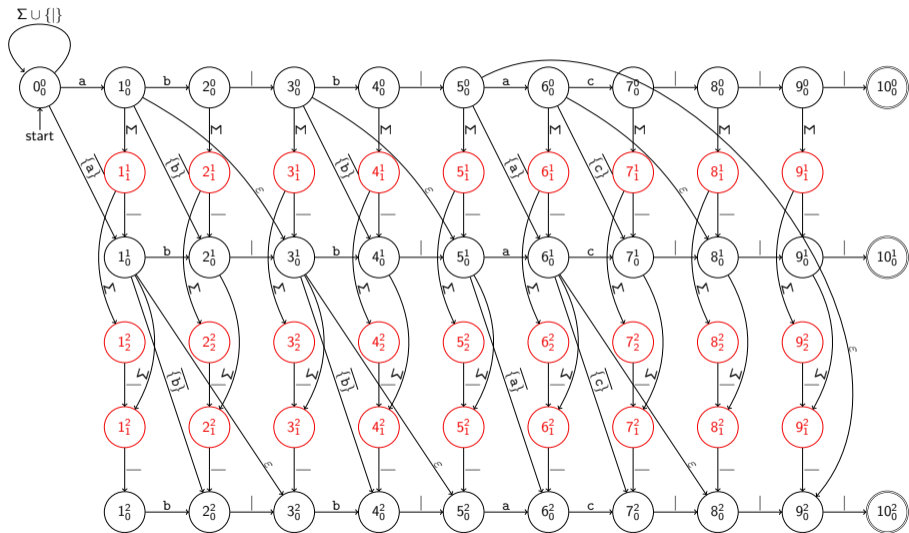
Given a tree T with n vertices and its prefix bar notation $\text{PREF-BAR}(T)$ with length $2n$, the *subtree jump table* \mathbf{S}_T for T is a mapping from a set of integers $\{1, \dots, 2n\}$ into a set of integers $\{0, \dots, 2n + 1\}$. If the substring $\mathbf{x}[i \dots j]$, where $1 \leq i < j$ is the prefix bar representation of a subtree of T , then $\mathbf{S}_T[i] = j + 1$ and $\mathbf{S}_T[j] = i - 1$.

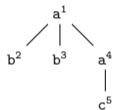


PREF-BAR(P)	a	b		b		a	c			
j	1	2	3	4	5	6	7	8	9	10
$\mathbf{S}_P[j]$	11	4	1	6	3	10	9	6	5	0

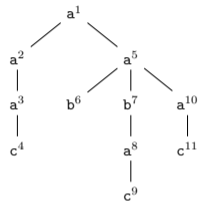
Use of pushdown store



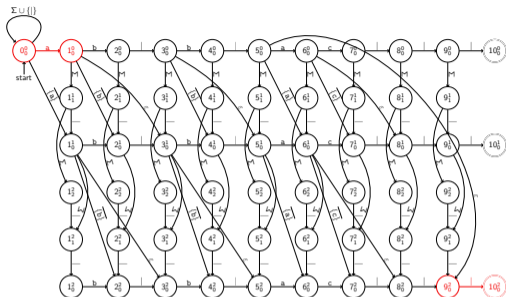
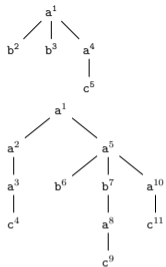




ab|b|ac||

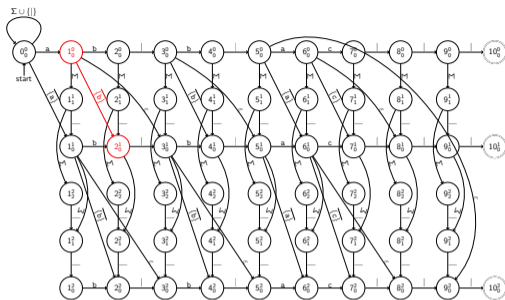
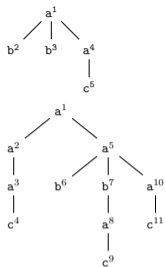


aaac|||ab|bac|||ac|||



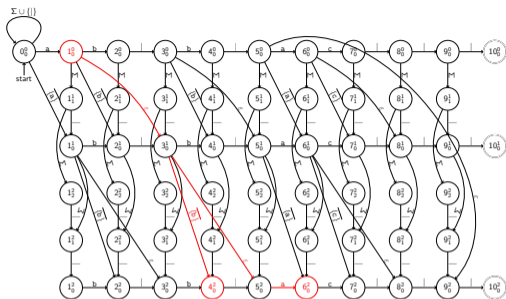
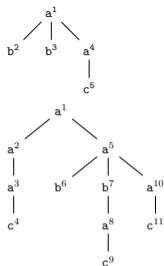
		1	2	3	4	5	6	7	8	...
		a	a	a	c				a	
a	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	
b	∞, ∞, ∞	∞, ∞, ∞	0, 1, ∞	0, 1, 2	1, 1, 2	1, 2, ∞	2, ∞, ∞	∞, ∞, ∞	0, ∞, ∞	
b	∞, ∞, ∞	∞, ∞, ∞	1, ∞, ∞	1, 2, ∞	1, 2, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
b	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
a	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
c	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
c	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
c	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
c	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
c	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	
c	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	
c	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	

$$D_{i,j,0} = D_{i-1,j-1,0} : t[i] = p[j] \wedge 1 \leq i \leq 2n \wedge 1 \leq j \leq 2m$$



		1	2	3	4	5	6	7	8	...
		a	a	a	c				a	
a	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	
b	∞, ∞, ∞	0, ∞, ∞	0, 1, ∞	0, 1, 2	1, 1, 2	1, 2, ∞	2, ∞, ∞	∞, ∞, ∞	0, ∞, ∞	
c	∞, ∞, ∞	∞, ∞, ∞	1, ∞, ∞	1, 2, ∞	1, 2, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
a	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	1, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
b	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	2, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
c	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
a	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
b	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
c	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
a	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
b	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	
c	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	

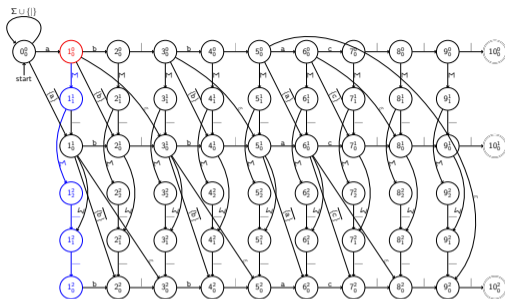
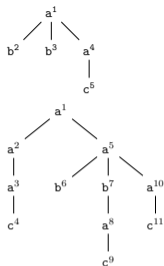
$$D_{i,j,0} = D_{i-1,j-1,0} + 1 : t[i], p[j] \in \Sigma \wedge 1 \leq i \leq 2n \wedge 1 \leq j \leq 2m$$



		1	2	3	4	5	6	7	8	...
		a	a	a	c	a	a	a	a	...
a	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	
b	∞, ∞, ∞	0, ∞, ∞	0, 1, ∞	0, 1, 2	1, 1, 2	1, 2, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
b	∞, ∞, ∞	∞, ∞, ∞	1, ∞, ∞	1, 2, ∞	1, 2, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
a	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	2, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
c	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
...	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	2, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	

$$D_{i,j,0} = D_{i-1, S_P[h], 0} + \frac{j - S_P[h] + 1}{2} : t[i] = p[j] \wedge p[j-1] = | \wedge 1 \leq i \leq 2n \wedge 2 \leq j \leq 2m \wedge 1 \leq h \leq 2m \wedge S_P[h] < j$$

$$D_{i,j,0} = D_{i-1, S_P[h], 0} + \frac{j - S_P[h] + 2}{2} : t[i], p[j] \in \Sigma \wedge t[i] \neq p[j] \wedge p[j-1] = | \wedge 1 \leq i \leq 2n \wedge 2 \leq j \leq 2m \wedge 1 \leq h \leq 2m \wedge S_P[h] < j$$



		1	2	3	4	5	6	7	8	...
		a	a	a	c				a	...
a		0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	0, ∞, ∞	
b		∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
b		∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
a		∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
c		∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
		∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
		∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
		∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
		∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
		∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	
		∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	∞, ∞, ∞	

$$D_{i,j,c} = D_{i-1,j,c-1} + 1 : \mathbf{t}[i] \in \Sigma \wedge 1 \leq i \leq 2n \wedge 1 \leq j \leq 2m \wedge 1 \leq c \leq k$$

$$D_{i,j,c} = D_{i-1,j,c+1} : \mathbf{t}[i] = \epsilon \wedge 1 \leq i \leq 2n \wedge 1 \leq j \leq 2m \wedge 0 \leq c < k$$

Complexity

Theorem (Space complexity)

Inexact 1-degree tree pattern matching can be solved using $\mathcal{O}(km)$ space.

Complexity

Theorem (Space complexity)

Inexact 1-degree tree pattern matching can be solved using $\mathcal{O}(km)$ space.

Theorem (Time complexity)

Inexact 1-degree tree pattern matching can be solved in $\mathcal{O}(kmn)$ time.

Conclusions

- solution for unit cost inexact 1-degree tree pattern matching
- string solution for unranked prefix bar notation
- pushdown automaton
- finite automaton
- dynamic programming
- all algorithms can be extended for non unit cost 1-degree distance