

Algorithms and datastructures II

Lecture 2: text searching (2/2)

Jan Hubička

Department of Applied Mathematics
Charles University
Prague

Oct 5 2020

String-searching (“searching needle in a haystack”)

string-searching

Given a string ν (“a needle”) and η (“a haystack”) find all occurrences of ν in η .

Some notation:

- ① Σ : an alphabet (finite set of characters)
- ② Σ^* : the set all finite words in alphabet Σ
- ③ α, β, \dots : words
- ④ $|\alpha|$: length of the word α .
- ⑤ ϵ : empty word (the only word of length 0)
- ⑥ $\alpha\beta$: concatenation of α and β
- ⑦ $\alpha[i]$: i -th character of α (starting from 0)
- ⑧ $\alpha[i : j]$: subword $\alpha[i]\alpha[i + 1]\dots\alpha[j - 1]$
- ⑨ $\alpha[: j]$: prefix of α of length j
- ⑩ $\alpha[i :]$: a suffix of α
- ⑪ $\alpha[:]$: whole word α

Occurrence of ν in η is any index i such that $\eta[i : i + |\nu|] = \nu$

Knuth–Morris–Pratt (KMP) algorithm (1974)

Searching automaton

- 1 **State 0, ..., $|\nu|$**
(state s corresponds to prefix $\nu[: s]$)
- 2 **Forward edges:** $s \rightarrow s + 1$
- 3 **Backward edges:** pointing from $s > 0$ to j such that $\nu[: j]$ is a proper suffix of $\nu[: s]$

Step (s, c) :

- 1 While $s \neq 0$ and $\nu[s] \neq c$:
- 2 $s \leftarrow b[s]$.
- 3 If $\nu[s] = c$: $s \leftarrow s + 1$.
- 4 Return s

KMPConstruction (ν) :

- 1 $b[0] \leftarrow$ undefined, $b[1] \leftarrow 0$, $s \leftarrow 0$.
- 2 For $i = 2, \dots, |\nu|$:
- 3 $s \leftarrow \text{Step } (s, \nu[i - 1])$.
- 4 $b[i] \leftarrow s$.

Search $(\eta, \text{automaton for } \nu)$:

- 1 $s \leftarrow 0$.
- 2 For $i = 0, \dots, |\eta| - 1$:
- 3 $s \leftarrow \text{Step } (s, \eta[i])$.
- 4 If $s = |\nu|$: report $i - |\nu| + 1$.

Theorem

Algorithm KMP will finish in time $\Theta(|\eta| + |\nu|)$.

Invariant: The state s corresponds to the longest suffix of $\eta[: i]$ that is a prefix of ν .

Today: Searching multiple needles at once

string-searching

Given a string ν (“a needle”) and η (“a haystack”) find all occurrences of ν in η .

Today: Searching multiple needles at once

string-searching

Given a string ν (“a needle”) and η (“a haystack”) find all occurrences of ν in η .

string-searching with multiple needles

Given strings $\nu_1, \nu_2, \dots, \nu_n$ (“a needles”) and η (“a haystack”) find all occurrences of $\nu_1, \nu_2, \dots, \nu_n$ in η .

We expect output in the form $S = \{(i, j) : \eta[i : i + |\nu_j|] = \nu_j\}$.

Today: Searching multiple needles at once

string-searching

Given a string ν (“a needle”) and η (“a haystack”) find all occurrences of ν in η .

string-searching with multiple needles

Given strings $\nu_1, \nu_2, \dots, \nu_n$ (“a needles”) and η (“a haystack”) find all occurrences of $\nu_1, \nu_2, \dots, \nu_n$ in η .

We expect output in the form $S = \{(i, j) : \eta[i : i + |\nu_j|] = \nu_j\}$.

Time complexity using KMP:

$$O\left(|\eta| \cdot n + \sum_{i=1}^n |\nu_i|\right)$$

Today: Searching multiple needles at once

string-searching

Given a string ν (“a needle”) and η (“a haystack”) find all occurrences of ν in η .

string-searching with multiple needles

Given strings $\nu_1, \nu_2, \dots, \nu_n$ (“a needles”) and η (“a haystack”) find all occurrences of $\nu_1, \nu_2, \dots, \nu_n$ in η .

We expect output in the form $S = \{(i, j) : \eta[i : i + |\nu_j|] = \nu_j\}$.

Time complexity using KMP:

$$O\left(|\eta| \cdot n + \sum_{i=1}^n |\nu_i|\right)$$

We seek for:

$$O\left(|\eta| + \sum_{i=1}^n |\nu_i| + |S|\right)$$

Recall
oo

Aho–Corasick
○●○○○○

Network flows
ooo

Aho–Corasick (1975)



Aho–Corasick: The automaton

Searching automaton (Knuth–Morris–Pratt)

- ① **States**: $0, \dots, |\nu|$
- ② **Forward edges**: $s \rightarrow s + 1$
- ③ **Backward edges**: pointing from $s > 0$ to j such that $\nu[:j]$ is a proper suffix of $\nu[:s]$

Searching automaton (Aho–Corasick)

- **States**: $\{\alpha : \exists_i \alpha \text{ is a prefix of } \eta_i\}$.
- **Forward edges**: $\{(\alpha, \beta) : \exists_{x \in \Sigma} \beta = \alpha x\}$.
- **Backward edges**: $\{(\alpha, \beta) : \beta \text{ is the longest proper suffix of } \alpha \text{ that is a state}\}$.

Aho–Corasick: How to output the locations?

- 1 Output whenever a leaf state is reached

Aho–Corasick: How to output the locations?

- ① Output whenever a leaf state is reached (not working)

Aho–Corasick: How to output the locations?

- ① Output whenever a leaf state is reached (not working)
- ② Follow back edges and output all needles on the path

Aho–Corasick: How to output the locations?

- ① Output whenever a leaf state is reached (**not working**)
- ② Follow back edges and output all needles on the path (**slow**)

Aho–Corasick: How to output the locations?

- ① Output whenever a leaf state is reached (**not working**)
- ② Follow back edges and output all needles on the path (**slow**)
- ③ Precompute lists of needles to output when given state is reached

Aho–Corasick: How to output the locations?

- ① Output whenever a leaf state is reached (not working)
- ② Follow back edges and output all needles on the path (slow)
- ③ Precompute lists of needles to output when given state is reached (too large automaton)

Aho–Corasick: How to output the locations?

- ① Output whenever a leaf state is reached (not working)
- ② Follow back edges and output all needles on the path (slow)
- ③ Precompute lists of needles to output when given state is reached (too large automaton)
- ④ Output edges pointing from state α to nearest output state β reachable by back edges.

Representation of the searching automaton

- ① States: 1, 2, ..., N.

Aho–Corasick: How to output the locations?

- ① Output whenever a leaf state is reached (**not working**)
- ② Follow back edges and output all needles on the path (**slow**)
- ③ Precompute lists of needles to output when given state is reached (**too large automaton**)
- ④ Output edges pointing from state α to nearest output state β reachable by back edges.

Representation of the searching automaton

- ① **States:** 1, 2, ..., N.
- ② **Needles:** Needle[s].

Aho–Corasick: How to output the locations?

- ① Output whenever a leaf state is reached (not working)
- ② Follow back edges and output all needles on the path (slow)
- ③ Precompute lists of needles to output when given state is reached (too large automaton)
- ④ Output edges pointing from state α to nearest output state β reachable by back edges.

Representation of the searching automaton

- ① **States:** 1, 2, ..., N.
- ② **Needles:** Needle[s].
- ③ **Forward edges:** Forward[s][c].
(s is an index of a state, c is an incoming character.)

Aho–Corasick: How to output the locations?

- ① Output whenever a leaf state is reached (not working)
- ② Follow back edges and output all needles on the path (slow)
- ③ Precompute lists of needles to output when given state is reached (too large automaton)
- ④ Output edges pointing from state α to nearest output state β reachable by back edges.

Representation of the searching automaton

- ① **States:** 1, 2, ..., N.
- ② **Needles:** Needle[s].
- ③ **Forward edges:** Forward[s][c].
(s is an index of a state, c is an incoming character.)
- ④ **Backward edges:** Back[s].

Aho–Corasick: How to output the locations?

- ① Output whenever a leaf state is reached (not working)
- ② Follow back edges and output all needles on the path (slow)
- ③ Precompute lists of needles to output when given state is reached (too large automaton)
- ④ Output edges pointing from state α to nearest output state β reachable by back edges.

Representation of the searching automaton

- ① **States:** 1, 2, ..., N.
- ② **Needles:** Needle[s].
- ③ **Forward edges:** Forward[s][c].
(s is an index of a state, c is an incoming character.)
- ④ **Backward edges:** Back[s].
- ⑤ **Output edges:** Output[s].

Aho–Corasick: search loop

Representation of the searching automaton

- ① **States**: 1, 2, ..., N.
- ② **Needles**: Needle[s].
- ③ **Forward edges**: Forward[s][c].
(*s* is an index of a state, *c* is an incoming character.)
- ④ **Backward edges**: Back[s].
- ⑤ **Output edges**: Output[s].

Aho–Corasick: search loop

Representation of the searching automaton

- ① **States**: 1, 2, ..., N.
- ② **Needles**: Needle[s].
- ③ **Forward edges**: Forward[s][c].
(s is an index of a state, c is an incoming character.)
- ④ **Backward edges**: Back[s].
- ⑤ **Output edges**: Output[s].

Step (s,c):

- ① While $\text{Forward}[s][c] = \emptyset$ and $s \neq \text{root}$: $s \leftarrow \text{Back}[s]$.
- ② If $\text{Forward}[s][c] \neq \emptyset$: $s \leftarrow \text{Forward}[s][c]$.
- ③ Return s .

Aho–Corasick: search loop

Representation of the searching automaton

- ① **States:** 1, 2, ..., N.
- ② **Needles:** Needle[s].
- ③ **Forward edges:** Forward[s][c].
(s is an index of a state, c is an incoming character.)
- ④ **Backward edges:** Back[s].
- ⑤ **Output edges:** Output[s].

Search (η , automaton)

- ① $s \leftarrow \text{root}.$
- ② For characters c in η do:
 - ③ $s \leftarrow \text{Step}(s, c).$
 - ④ $j \leftarrow s.$
 - ⑤ While $j \neq \emptyset$:
 - ⑥ If Needle(j) $\neq \emptyset$: Report occurrence of needle.
 - ⑦ $j \leftarrow \text{Output}(j).$

Step (s, c):

- ① While $\text{Forward}[s][c] = \emptyset$ and $s \neq \text{root}$: $s \leftarrow \text{Back}[s].$
- ② If $\text{Forward}[s][c] \neq \emptyset$: $s \leftarrow \text{Forward}[s][c].$
- ③ Return $s.$

Aho–Corasick: search loop

Representation of the searching automaton

- ① **States:** $1, 2, \dots, N$.
- ② **Needles:** $\text{Needle}[s]$.
- ③ **Forward edges:** $\text{Forward}[s][c]$.
(s is an index of a state, c is an incoming character.)
- ④ **Backward edges:** $\text{Back}[s]$.
- ⑤ **Output edges:** $\text{Output}[s]$.

Step (s, c):

- ① While $\text{Forward}[s][c] = \emptyset$ and $s \neq \text{root}$: $s \leftarrow \text{Back}[s]$.
- ② If $\text{Forward}[s][c] \neq \emptyset$: $s \leftarrow \text{Forward}[s][c]$.
- ③ Return s .

Search (η , automaton)

- ① $s \leftarrow \text{root}$.
- ② For characters c in η do:
 - ③ $s \leftarrow \text{Step}(s, c)$.
 - ④ $j \leftarrow s$.
 - ⑤ While $j \neq \emptyset$:
 - ⑥ If $\text{Needle}(j) \neq \emptyset$: Report occurrence of needle.
 - ⑦ $j \leftarrow \text{Output}(j)$.

Invariant: The state s corresponds to the longest suffix of $\eta[: i]$ that is a prefix of some needle.

Lemma

Search runtime is $\Theta(|\eta| + |S|)$.

Recall
oo

Aho-Corasick
oooooo

Network flows
ooo



Aho–Corasick: automaton construction

Construct (ν_1, \dots, ν_n)

- ➊ Initialize a trie with root r .
(sets Forward and Needle arrays)
- ➋ Insert words ν_1, \dots, ν_n to the trie
- ➌ $\text{Back}(r) \leftarrow \emptyset$, $\text{Output}(r) \leftarrow \emptyset$.
- ➍ Create a queue Q and insert all sons of r .
- ➎ For every son s of r : $\text{Back}(s) \leftarrow r$, $\text{Output}(s) \leftarrow \emptyset$
- ➏ While $Q \neq \emptyset$:
 - ➐ Dequeue i from Q .
 - ➑ For every son s of i :
 - ➒ $b \leftarrow \text{Step}(\text{Back}[i], \text{letter on edge } (i, s))$.
 - ➓ $\text{Back}[s] \leftarrow b$.
 - ➔ If $\text{Needle}[b] \neq \emptyset$: $\text{Output}[s] \leftarrow b$
else $\text{Output}[s] \leftarrow \text{Output}[b]$.
 - ➕ Append s to Q .

Representation of the searching automaton

- ➊ **States**: $1, 2, \dots, N$.
- ➋ **Needles**: $\text{Needle}[s]$.
- ➌ **Forward edges**: $\text{Forward}[s][c]$.
(s is an index of a state, c is an incoming character.)
- ➍ **Backward edges**: $\text{Back}[s]$.
- ➎ **Output edges**: $\text{Output}[s]$.

Aho–Corasick: automaton construction

Construct (ν_1, \dots, ν_n)

- ➊ Initialize a trie with root r .
(sets Forward and Needle arrays)
- ➋ Insert words ν_1, \dots, ν_n to the trie
- ➌ $\text{Back}(r) \leftarrow \emptyset$, $\text{Output}(r) \leftarrow \emptyset$.
- ➍ Create a queue Q and insert all sons of r .
- ➎ For every son s of r : $\text{Back}(s) \leftarrow r$, $\text{Output}(s) \leftarrow \emptyset$
- ➏ While $Q \neq \emptyset$:
 - ➐ Dequeue i from Q .
 - ➑ For every son s of i :
 - ➒ $b \leftarrow \text{Step}(\text{Back}[i], \text{letter on edge } (i, s))$.
 - ➓ $\text{Back}[s] \leftarrow b$.
 - ➔ If $\text{Needle}[b] \neq \emptyset$: $\text{Output}[s] \leftarrow b$
else $\text{Output}[s] \leftarrow \text{Output}[b]$.
 - ➕ Append s to Q .

Representation of the searching automaton

- ➊ **States**: $1, 2, \dots, N$.
- ➋ **Needles**: $\text{Needle}[s]$.
- ➌ **Forward edges**: $\text{Forward}[s][c]$.
(s is an index of a state, c is an incoming character.)
- ➍ **Backward edges**: $\text{Back}[s]$.
- ➎ **Output edges**: $\text{Output}[s]$.

Lemma

Runtime of Construction is $O(\sum_{i=1}^n \nu_i)$.

Aho–Corasick: automaton construction

Construct (ν_1, \dots, ν_n)

- ➊ Initialize a trie with root r .
(sets Forward and Needle arrays)
- ➋ Insert words ν_1, \dots, ν_n to the trie
- ➌ $\text{Back}(r) \leftarrow \emptyset$, $\text{Output}(r) \leftarrow \emptyset$.
- ➍ Create a queue Q and insert all sons of r .
- ➎ For every son s of r : $\text{Back}(s) \leftarrow r$, $\text{Output}(s) \leftarrow \emptyset$
- ➏ While $Q \neq \emptyset$:
 - ➐ Dequeue i from Q .
 - ➑ For every son s of i :
 - ➒ $b \leftarrow \text{Step}(\text{Back}[i], \text{letter on edge } (i, s))$.
 - ➓ $\text{Back}[s] \leftarrow b$.
 - ➔ If $\text{Needle}[b] \neq \emptyset$: $\text{Output}[s] \leftarrow b$
else $\text{Output}[s] \leftarrow \text{Output}[b]$.
 - ➕ Append s to Q .

Representation of the searching automaton

- ➊ **States**: $1, 2, \dots, N$.
- ➋ **Needles**: $\text{Needle}[s]$.
- ➌ **Forward edges**: $\text{Forward}[s][c]$.
(s is an index of a state, c is an incoming character.)
- ➍ **Backward edges**: $\text{Back}[s]$.
- ➎ **Output edges**: $\text{Output}[s]$.

Lemma

Runtime of Construction is $O(\sum_{i=1}^n \nu_i)$.

Theorem

Aho–Corasick algorithm will find all occurrences of needles in η and will terminate in time $O(|\eta| + \sum_{i=1}^n |\nu_i| + |S|)$.

Recall
oo

Aho–Corasick
ooooooo

Network flows
●oo

Network flow

Recall
oo

Aho–Corasick
ooooooo

Network flows
●oo

Network flow



Network flow

Definition (Network)

Network is an 4-tuple $N = (V, E, s, t, c)$ where

- ① (V, E) is a directed graph,
- ② $s \in V$ is a **source** vertex,
- ③ $t \in V$ is a **sink** vertex,
- ④ $c : E \rightarrow \mathbb{R}_0^+$ is a function assigning every edge a **capacity**.

- $f^+(v) = \sum_{u, (u,v) \in E} f(u, v)$ (**flow into a vertex**)
- $f^-(v) = \sum_{u, (v,u) \in E} f(v, u)$ (**flow out of a vertex**)
- $f^\Delta(v) = f^+(v) - f^-(v)$ (**surplus**)

Here $f : E \rightarrow \mathbb{R}_0^+$

Network flow

Definition (Network)

Network is an 4-tuple $N = (V, E, s, t, c)$ where

- ① (V, E) is a directed graph,
- ② $s \in V$ is a **source** vertex,
- ③ $t \in V$ is a **sink** vertex,
- ④ $c : E \rightarrow \mathbb{R}_0^+$ is a function assigning every edge a **capacity**.

- $f^+(v) = \sum_{u, (u,v) \in E} f(u, v)$ (**flow into a vertex**)
- $f^-(v) = \sum_{u, (v,u) \in E} f(v, u)$ (**flow out of a vertex**)
- $f^\Delta(v) = f^+(v) - f^-(v)$ (**surplus**)

Here $f : E \rightarrow \mathbb{R}_0^+$

Definition (Flow)

Function $f : E \rightarrow \mathbb{R}_0^+$ is **flow** if it satisfies

- ① **Capacity constraint:** $(\forall e \in E) : f(e) \leq c(e)$
- ② **Conservation of flows (Kirchoff's law):** $(\forall v \in V \setminus \{s, t\}) : f^\Delta = 0$

Value of the flow: $|f| = f^\Delta(t)$.

Network flow

Lemma

For every flow it holds that: $f^\Delta(s) + f^\Delta(t) = 0$.

Network flow

Lemma

For every flow it holds that: $f^\Delta(s) + f^\Delta(t) = 0$.

$$\sum_{v \in V} f^\Delta(v) = 0$$

Naive algorithm

Network flow problem

Given network $N = (V, E, s, t, c)$ find flow f maximizing $|f|$.

Naive approach: Start with 0 flow and keep improving as long as there is path from source to sink that can be improved.