

Algorithms and datastructures II

Lecture 1: text searching (1/2)

Jan Hubička

Department of Applied Mathematics
Charles University
Prague

Oct 5 2020

Algorithms and data-structures II

Lecturer: Jan Hubička, hubicka@kam.mff.cuni.cz.

- Consultations by appointment (email or zoom).
- Do we want a discussion forum, like discord or moodle?

Algorithms and data-structures II

Lecturer: Jan Hubička, hubicka@kam.mff.cuni.cz.

- Consultations by appointment (email or zoom).
- Do we want a discussion forum, like discord or moodle?

Lecture: Monday 9am over zoom.

- Lectures will be live: please come and do ask questions, tell me how you like the lecture!
- Recordings will be available on stream.cuni.cz (linked from webpage).
- Recordings will include all questions and communication. If you ask me after the lecture ends, I will edit out your question.
- I will try to keep lecture synchronized with Czech lecture by Jan Hric.
- If you spot mistake, please let me know. The slides are brand new.

Webpage: <https://iuuk.mff.cuni.cz/~hubicka/2020/adsII.html>

Syllabus

① String-searching

Syllabus

- ① String-searching
- ② Network flows

Syllabus

- ① String-searching
- ② Network flows
- ③ Algebraic algorithms (Fourier transformation)

Syllabus

- ① String-searching
- ② Network flows
- ③ Algebraic algorithms (Fourier transformation)
- ④ Parallel algorithms

Syllabus

- ① String-searching
- ② Network flows
- ③ Algebraic algorithms (Fourier transformation)
- ④ Parallel algorithms
- ⑤ Geometric algorithms

Syllabus

- ① String-searching
- ② Network flows
- ③ Algebraic algorithms (Fourier transformation)
- ④ Parallel algorithms
- ⑤ Geometric algorithms
- ⑥ Introduction to complexity

Syllabus

- ① String-searching
- ② Network flows
- ③ Algebraic algorithms (Fourier transformation)
- ④ Parallel algorithms
- ⑤ Geometric algorithms
- ⑥ Introduction to complexity
- ⑦ Approximation algorithm

Syllabus

- ① String-searching
- ② Network flows
- ③ Algebraic algorithms (Fourier transformation)
- ④ Parallel algorithms
- ⑤ Geometric algorithms
- ⑥ Introduction to complexity
- ⑦ Approximation algorithm
- ⑧ Probabilistic algorithms and cryptography



Today: 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 .

Today: 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)

Today: 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

Today: 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

Today: 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 .

Today: 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)

Today: 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

Today: 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)

Today: 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 $\alpha[i]\alpha[i+1]\dots\alpha[j-1]$

Today: 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

Today: 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

Today: 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

Today: 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 $\eta[i : i + |\nu|] = \nu$

Naive approach

string-searching

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

Search (ν, η) :

- ① For $i = 0, \dots, |\eta| - |\nu| - 1$:
- ② If $\eta[i : i + |\nu|] = \nu$: output i

$\nu = \text{coconut}$

η is some very long text about coconuts

Naive approach

string-searching

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

Search (ν, η) :

- ① For $i = 0, \dots, |\eta| - |\nu| - 1$:
- ② If $\eta[i : i + |\nu|] = \nu$: output i

$\nu = \text{coconut}$

η is some very long text about coconuts

Time complexity: $\Theta(|\nu| \cdot |\eta|)$.

Incremental algorithm

An incremental algorithm receives characters of η one by one and after receiving a new character it immediately outputs possible new occurrences of ν .

Incremental algorithm

An incremental algorithm receives characters of η one by one and after receiving a new character it immediately outputs possible new occurrences of ν .

Basic idea

We would like to remember an **state**. This is longest prefix of ν which is a suffix of η .

Observation: Whenever algorithm enters state ν it finds a new occurrence of ν in the input.

Assume that algorithm seen string η , is in state α and receives a new character c . How to update the state?

Incremental algorithm

An incremental algorithm receives characters of η one by one and after receiving a new character it immediately outputs possible new occurrences of ν .

Basic idea

We would like to remember an **state**. This is longest prefix of ν which is a suffix of η .

Observation: Whenever algorithm enters state ν it finds a new occurrence of ν in the input.

Assume that algorithm seen string η , is in state α and receives a new character c . How to update the state?

$$\text{State } (\alpha, c) = \begin{cases} \alpha c & \text{If } \alpha c \text{ is a prefix of } \nu \end{cases}$$

Incremental algorithm

An incremental algorithm receives characters of η one by one and after receiving a new character it immediately outputs possible new occurrences of ν .

Basic idea

We would like to remember an **state**. This is longest prefix of ν which is a suffix of η .

Observation: Whenever algorithm enters state ν it finds a new occurrence of ν in the input.

Assume that algorithm seen string η , is in state α and receives a new character c . How to update the state?

$$\text{State } (\alpha, c) = \begin{cases} \alpha c & \text{If } \alpha c \text{ is a prefix of } \nu \\ \epsilon & \text{otherwise} \end{cases}$$

Incremental algorithm

An incremental algorithm receives characters of η one by one and after receiving a new character it immediately outputs possible new occurrences of ν .

Basic idea

We would like to remember an **state**. This is longest prefix of ν which is a suffix of η .

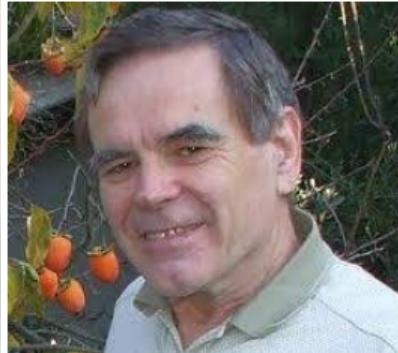
Observation: Whenever algorithm enters state ν it finds a new occurrence of ν in the input.

Assume that algorithm seen string η , is in state α and receives a new character c . How to update the state?

$$\text{State } (\alpha, c) = \begin{cases} \alpha c & \text{If } \alpha c \text{ is a prefix of } \nu \\ \epsilon & \\ \alpha' c & \alpha' c \text{ is a prefix of } \nu \text{ and } \alpha' \text{ is a suffix of } \alpha \end{cases}$$

We want to compute **backward function** b which tells for every prefix α of ν the longest proper suffix α' (of α) that is also a prefix of ν .

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



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

We want to compute **backward function b** which tells for every prefix α of ν the longest proper suffix α' (of α) that is also a prefix of ν .

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 $\nu[: s]$

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

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 $i - |\nu| + 1$.

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

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

Lemma

Search will run in time $O(|\eta|)$

Proof.

- 1 Search itself is clearly $O(|\eta|)$

□

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

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 $i - |\nu| + 1$.

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

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

Lemma

Search will run in time $O(|\eta|)$

Proof.

- ① Search itself is clearly $O(|\eta|)$
- ② Number of forward transitions in Step is $O(|\eta|)$

□

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

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 $i - |\nu| + 1$.

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

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

Lemma

Search will run in time $O(|\eta|)$

Proof.

- ① Search itself is clearly $O(|\eta|)$
- ② Number of forward transitions in Step is $O(|\eta|)$
- ③ Number of backward transitions in Step is $O(|\eta|)$

□

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

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 $i - |\nu| + 1$.

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

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

Lemma

Search will run in time $O(|\eta|)$

Proof.

- ① Search itself is clearly $O(|\eta|)$
- ② Number of forward transitions in Step is $O(|\eta|)$
- ③ Number of backward transitions in Step is $O(|\eta|)$

□

Overall runtime $\Theta(|\eta|)$.

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

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 $i - |\nu| + 1$.

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

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

How to obtain the automaton?

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

Search (η , automaton for ν):

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

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

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

How to obtain the automaton? We will steal it!

Imagine that someone has the automaton and we want to figure all backward edges

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

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 $i - |\nu| + 1$.

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

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

How to obtain the automaton? We will steal it!

Imagine that someone has the automaton and we want to figure all backward edges

To determine $b(s)$ we need to search $\nu[1 : s]$

Recall: $b(s)$ is j such that $\nu[: j]$ is the longest 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

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 $i - |\nu| + 1$.

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

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

How to obtain the automaton? We will steal it!

Imagine that someone has the automaton and we want to figure all backward edges

To determine $b(s)$ we need to search $\nu[1 : s]$

Recall: $b(s)$ is j such that $\nu[: j]$ is the longest proper suffix of $\nu[: s]$

KMPConstruction (ν):

- 1 $b[0] \leftarrow \text{undefined}$, $b[1] \leftarrow 0$, $s \leftarrow 0$.

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

Search (η , automaton for ν):

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

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

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

How to obtain the automaton? We will steal it!

Imagine that someone has the automaton and we want to figure all backward edges

To determine $b(s)$ we need to search $\nu[1 : s]$

Recall: $b(s)$ is j such that $\nu[: j]$ is the longest proper suffix of $\nu[: s]$

KMPConstruction (ν):

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

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

Search (η , automaton for ν):

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

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

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

How to obtain the automaton? We will steal it!

Imagine that someone has the automaton and we want to figure all backward edges

To determine $b(s)$ we need to search $\nu[1 : s]$

Recall: $b(s)$ is j such that $\nu[: j]$ is the longest proper suffix of $\nu[: s]$

KMPConstruction (ν):

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

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

Search (η , automaton for ν):

- 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 .



Rabin–Karp algorithm (1987)

Recall the rotating/sliding hash function

$$H(x_1, x_2, \dots, x_K) = (x_1 P^{K-1} + x_2 P^{K-2} + \dots + x_{K-1} P^1 + x_K P^0) \bmod N$$

For $K = |\nu|$, prime number P and $N > 0$.

Rabin–Karp algorithm (1987)

Recall the rotating/sliding hash function

$$H(x_1, x_2, \dots, x_K) = (x_1 P^{K-1} + x_2 P^{K-2} + \dots + x_{K-1} P^1 + x_K P^0) \bmod N$$

For $K = |\nu|$, prime number P and $N > 0$. Observe that:

$$\begin{aligned} H(x_2, x_3, \dots, x_{K+1}) &= (x_2 P^{K-1} + x_3 P^{K-2} + \dots + x_K P^1 + x_{K+1} P^0) \bmod N \\ &= (P \cdot H(x_1, x_2, \dots, x_K) - x_1 P^K + x_{K+1}) \bmod N \end{aligned}$$

Rabin–Karp algorithm (1987)

Recall the rotating/sliding hash function

$$H(x_1, x_2, \dots, x_K) = (x_1 P^{K-1} + x_2 P^{K-2} + \dots + x_{K-1} P^1 + x_K P^0) \bmod N$$

For $K = |\nu|$, prime number P and $N > 0$. Observe that:

$$\begin{aligned} H(x_2, x_3, \dots, x_{K+1}) &= (x_2 P^{K-1} + x_3 P^{K-2} + \dots + x_K P^1 + x_{K+1} P^0) \bmod N \\ &= (P \cdot H(x_1, x_2, \dots, x_K) - x_1 P^K + x_{K+1}) \bmod N \end{aligned}$$

RabinKarpSearch (ν, η):

- ① Choose prime P and $N > 1$.
- ② Precompute $P^{|\nu|} \bmod N$.
- ③ $n \leftarrow H(\nu)$.
- ④ $h \leftarrow H(\eta[:|\nu|])$.
- ⑤ For $i = 0, \dots, |\eta| - |\nu| - 1$:
- ⑥ if $h = n$ and $\eta[i : i + |\nu|] = \nu$: report i .
- ⑦ $h \leftarrow (P \cdot h - \eta[i] \cdot P^{|\nu|} + \eta[i + |\nu|]) \bmod N$.

Rabin–Karp algorithm (1987)

Recall the rotating/sliding hash function

$$H(x_1, x_2, \dots, x_K) = (x_1 P^{K-1} + x_2 P^{K-2} + \dots + x_{K-1} P^1 + x_K P^0) \bmod N$$

For $K = |\nu|$, prime number P and $N > 0$. Observe that:

$$\begin{aligned} H(x_2, x_3, \dots, x_{K+1}) &= (x_2 P^{K-1} + x_3 P^{K-2} + \dots + x_K P^1 + x_{K+1} P^0) \bmod N \\ &= (P \cdot H(x_1, x_2, \dots, x_K) - x_1 P^K + x_{K+1}) \bmod N \end{aligned}$$

RabinKarpSearch (ν, η):

- ① Choose prime P and $N > 1$.
- ② Precompute $P^{|\nu|} \bmod N$.
- ③ $n \leftarrow H(\nu)$.
- ④ $h \leftarrow H(\eta[:|\nu|])$.
- ⑤ For $i = 0, \dots, |\eta| - |\nu| - 1$:
- ⑥ if $h = n$ and $\eta[i : i + |\nu|] = \nu$: report i .
- ⑦ $h \leftarrow (P \cdot h - \eta[i] \cdot P^{|\nu|} + \eta[i + |\nu|]) \bmod N$.

Expected time complexity: $O(|\nu| + |\eta| + C|\nu| + |\eta|/N \cdot |\nu|)$ where C is number of occurrences reported.